

OFFICIAL MICROSOFT LEARNING PRODUCT

20767C

Implementing a SQL Data Warehouse

MCT USE ONLY. STUDENT USE PROHIBITED

Information in this document, including URL and other Internet Web site references, is subject to change without notice. Unless otherwise noted, the example companies, organizations, products, domain names, e-mail addresses, logos, people, places, and events depicted herein are fictitious, and no association with any real company, organization, product, domain name, e-mail address, logo, person, place or event is intended or should be inferred. Complying with all applicable copyright laws is the responsibility of the user. Without limiting the rights under copyright, no part of this document may be reproduced, stored in or introduced into a retrieval system, or transmitted in any form or by any means (electronic, mechanical, photocopying, recording, or otherwise), or for any purpose, without the express written permission of Microsoft Corporation.

Microsoft may have patents, patent applications, trademarks, copyrights, or other intellectual property rights covering subject matter in this document. Except as expressly provided in any written license agreement from Microsoft, the furnishing of this document does not give you any license to these patents, trademarks, copyrights, or other intellectual property.

The names of manufacturers, products, or URLs are provided for informational purposes only and Microsoft makes no representations and warranties, either expressed, implied, or statutory, regarding these manufacturers or the use of the products with any Microsoft technologies. The inclusion of a manufacturer or product does not imply endorsement of Microsoft of the manufacturer or product. Links may be provided to third party sites. Such sites are not under the control of Microsoft and Microsoft is not responsible for the contents of any linked site or any link contained in a linked site, or any changes or updates to such sites. Microsoft is not responsible for webcasting or any other form of transmission received from any linked site. Microsoft is providing these links to you only as a convenience, and the inclusion of any link does not imply endorsement of Microsoft of the site or the products contained therein.

© 2018 Microsoft Corporation. All rights reserved.

Microsoft and the trademarks listed at <https://www.microsoft.com/en-us/legal/intellectualproperty/trademarks/en-us.aspx> are trademarks of the Microsoft group of companies. All other trademarks are property of their respective owners

Product Number: 20767C

Part Number (if applicable): X21-64458

Released: 01/2018

MICROSOFT LICENSE TERMS MICROSOFT INSTRUCTOR-LED COURSEWARE

These license terms are an agreement between Microsoft Corporation (or based on where you live, one of its affiliates) and you. Please read them. They apply to your use of the content accompanying this agreement which includes the media on which you received it, if any. These license terms also apply to Trainer Content and any updates and supplements for the Licensed Content unless other terms accompany those items. If so, those terms apply.

**BY ACCESSING, DOWNLOADING OR USING THE LICENSED CONTENT, YOU ACCEPT THESE TERMS.
IF YOU DO NOT ACCEPT THEM, DO NOT ACCESS, DOWNLOAD OR USE THE LICENSED CONTENT.**

If you comply with these license terms, you have the rights below for each license you acquire.

1. DEFINITIONS.

- a. "Authorized Learning Center" means a Microsoft IT Academy Program Member, Microsoft Learning Competency Member, or such other entity as Microsoft may designate from time to time.
- b. "Authorized Training Session" means the instructor-led training class using Microsoft Instructor-Led Courseware conducted by a Trainer at or through an Authorized Learning Center.
- c. "Classroom Device" means one (1) dedicated, secure computer that an Authorized Learning Center owns or controls that is located at an Authorized Learning Center's training facilities that meets or exceeds the hardware level specified for the particular Microsoft Instructor-Led Courseware.
- d. "End User" means an individual who is (i) duly enrolled in and attending an Authorized Training Session or Private Training Session, (ii) an employee of a MPN Member, or (iii) a Microsoft full-time employee.
- e. "Licensed Content" means the content accompanying this agreement which may include the Microsoft Instructor-Led Courseware or Trainer Content.
- f. "Microsoft Certified Trainer" or "MCT" means an individual who is (i) engaged to teach a training session to End Users on behalf of an Authorized Learning Center or MPN Member, and (ii) currently certified as a Microsoft Certified Trainer under the Microsoft Certification Program.
- g. "Microsoft Instructor-Led Courseware" means the Microsoft-branded instructor-led training course that educates IT professionals and developers on Microsoft technologies. A Microsoft Instructor-Led Courseware title may be branded as MOC, Microsoft Dynamics or Microsoft Business Group courseware.
- h. "Microsoft IT Academy Program Member" means an active member of the Microsoft IT Academy Program.
- i. "Microsoft Learning Competency Member" means an active member of the Microsoft Partner Network program in good standing that currently holds the Learning Competency status.
- j. "MOC" means the "Official Microsoft Learning Product" instructor-led courseware known as Microsoft Official Course that educates IT professionals and developers on Microsoft technologies.
- k. "MPN Member" means an active Microsoft Partner Network program member in good standing.

MCT USE ONLY. STUDENT USE PROHIBITED

- l. "Personal Device" means one (1) personal computer, device, workstation or other digital electronic device that you personally own or control that meets or exceeds the hardware level specified for the particular Microsoft Instructor-Led Courseware.
- m. "Private Training Session" means the instructor-led training classes provided by MPN Members for corporate customers to teach a predefined learning objective using Microsoft Instructor-Led Courseware. These classes are not advertised or promoted to the general public and class attendance is restricted to individuals employed by or contracted by the corporate customer.
- n. "Trainer" means (i) an academically accredited educator engaged by a Microsoft IT Academy Program Member to teach an Authorized Training Session, and/or (ii) a MCT.
- o. "Trainer Content" means the trainer version of the Microsoft Instructor-Led Courseware and additional supplemental content designated solely for Trainers' use to teach a training session using the Microsoft Instructor-Led Courseware. Trainer Content may include Microsoft PowerPoint presentations, trainer preparation guide, train the trainer materials, Microsoft One Note packs, classroom setup guide and Pre-release course feedback form. To clarify, Trainer Content does not include any software, virtual hard disks or virtual machines.

2. **USE RIGHTS.** The Licensed Content is licensed not sold. The Licensed Content is licensed on a **one copy per user basis**, such that you must acquire a license for each individual that accesses or uses the Licensed Content.

2.1 Below are five separate sets of use rights. Only one set of rights apply to you.

a. **If you are a Microsoft IT Academy Program Member:**

- i. Each license acquired on behalf of yourself may only be used to review one (1) copy of the Microsoft Instructor-Led Courseware in the form provided to you. If the Microsoft Instructor-Led Courseware is in digital format, you may install one (1) copy on up to three (3) Personal Devices. You may not install the Microsoft Instructor-Led Courseware on a device you do not own or control.
- ii. For each license you acquire on behalf of an End User or Trainer, you may either:
 - 1. distribute one (1) hard copy version of the Microsoft Instructor-Led Courseware to one (1) End User who is enrolled in the Authorized Training Session, and only immediately prior to the commencement of the Authorized Training Session that is the subject matter of the Microsoft Instructor-Led Courseware being provided, **or**
 - 2. provide one (1) End User with the unique redemption code and instructions on how they can access one (1) digital version of the Microsoft Instructor-Led Courseware, **or**
 - 3. provide one (1) Trainer with the unique redemption code and instructions on how they can access one (1) Trainer Content,**provided you comply with the following:**
- iii. you will only provide access to the Licensed Content to those individuals who have acquired a valid license to the Licensed Content,
- iv. you will ensure each End User attending an Authorized Training Session has their own valid licensed copy of the Microsoft Instructor-Led Courseware that is the subject of the Authorized Training Session,
- v. you will ensure that each End User provided with the hard-copy version of the Microsoft Instructor-Led Courseware will be presented with a copy of this agreement and each End User will agree that their use of the Microsoft Instructor-Led Courseware will be subject to the terms in this agreement prior to providing them with the Microsoft Instructor-Led Courseware. Each individual will be required to denote their acceptance of this agreement in a manner that is enforceable under local law prior to their accessing the Microsoft Instructor-Led Courseware,
- vi. you will ensure that each Trainer teaching an Authorized Training Session has their own valid licensed copy of the Trainer Content that is the subject of the Authorized Training Session,

- vii. you will only use qualified Trainers who have in-depth knowledge of and experience with the Microsoft technology that is the subject of the Microsoft Instructor-Led Courseware being taught for all your Authorized Training Sessions,
- viii. you will only deliver a maximum of 15 hours of training per week for each Authorized Training Session that uses a MOC title, and
- ix. you acknowledge that Trainers that are not MCTs will not have access to all of the trainer resources for the Microsoft Instructor-Led Courseware.

b. If you are a Microsoft Learning Competency Member:

- i. Each license acquired on behalf of yourself may only be used to review one (1) copy of the Microsoft Instructor-Led Courseware in the form provided to you. If the Microsoft Instructor-Led Courseware is in digital format, you may install one (1) copy on up to three (3) Personal Devices. You may not install the Microsoft Instructor-Led Courseware on a device you do not own or control.
- ii. For each license you acquire on behalf of an End User or Trainer, you may either:
 - 1. distribute one (1) hard copy version of the Microsoft Instructor-Led Courseware to one (1) End User attending the Authorized Training Session and only immediately prior to the commencement of the Authorized Training Session that is the subject matter of the Microsoft Instructor-Led Courseware provided, **or**
 - 2. provide one (1) End User attending the Authorized Training Session with the unique redemption code and instructions on how they can access one (1) digital version of the Microsoft Instructor-Led Courseware, **or**
 - 3. you will provide one (1) Trainer with the unique redemption code and instructions on how they can access one (1) Trainer Content,
provided you comply with the following:
- iii. you will only provide access to the Licensed Content to those individuals who have acquired a valid license to the Licensed Content,
- iv. you will ensure that each End User attending an Authorized Training Session has their own valid licensed copy of the Microsoft Instructor-Led Courseware that is the subject of the Authorized Training Session,
- v. you will ensure that each End User provided with a hard-copy version of the Microsoft Instructor-Led Courseware will be presented with a copy of this agreement and each End User will agree that their use of the Microsoft Instructor-Led Courseware will be subject to the terms in this agreement prior to providing them with the Microsoft Instructor-Led Courseware. Each individual will be required to denote their acceptance of this agreement in a manner that is enforceable under local law prior to their accessing the Microsoft Instructor-Led Courseware,
- vi. you will ensure that each Trainer teaching an Authorized Training Session has their own valid licensed copy of the Trainer Content that is the subject of the Authorized Training Session,
- vii. you will only use qualified Trainers who hold the applicable Microsoft Certification credential that is the subject of the Microsoft Instructor-Led Courseware being taught for your Authorized Training Sessions,
- viii. you will only use qualified MCTs who also hold the applicable Microsoft Certification credential that is the subject of the MOC title being taught for all your Authorized Training Sessions using MOC,
- ix. you will only provide access to the Microsoft Instructor-Led Courseware to End Users, and
- x. you will only provide access to the Trainer Content to Trainers.

c. If you are a MPN Member:

- i. Each license acquired on behalf of yourself may only be used to review one (1) copy of the Microsoft Instructor-Led Courseware in the form provided to you. If the Microsoft Instructor-Led Courseware is in digital format, you may install one (1) copy on up to three (3) Personal Devices. You may not install the Microsoft Instructor-Led Courseware on a device you do not own or control.
- ii. For each license you acquire on behalf of an End User or Trainer, you may either:
 1. distribute one (1) hard copy version of the Microsoft Instructor-Led Courseware to one (1) End User attending the Private Training Session, and only immediately prior to the commencement of the Private Training Session that is the subject matter of the Microsoft Instructor-Led Courseware being provided, **or**
 2. provide one (1) End User who is attending the Private Training Session with the unique redemption code and instructions on how they can access one (1) digital version of the Microsoft Instructor-Led Courseware, **or**
 3. you will provide one (1) Trainer who is teaching the Private Training Session with the unique redemption code and instructions on how they can access one (1) Trainer Content, **provided you comply with the following:**
- iii. you will only provide access to the Licensed Content to those individuals who have acquired a valid license to the Licensed Content,
- iv. you will ensure that each End User attending an Private Training Session has their own valid licensed copy of the Microsoft Instructor-Led Courseware that is the subject of the Private Training Session,
- v. you will ensure that each End User provided with a hard copy version of the Microsoft Instructor-Led Courseware will be presented with a copy of this agreement and each End User will agree that their use of the Microsoft Instructor-Led Courseware will be subject to the terms in this agreement prior to providing them with the Microsoft Instructor-Led Courseware. Each individual will be required to denote their acceptance of this agreement in a manner that is enforceable under local law prior to their accessing the Microsoft Instructor-Led Courseware,
- vi. you will ensure that each Trainer teaching an Private Training Session has their own valid licensed copy of the Trainer Content that is the subject of the Private Training Session,
- vii. you will only use qualified Trainers who hold the applicable Microsoft Certification credential that is the subject of the Microsoft Instructor-Led Courseware being taught for all your Private Training Sessions,
- viii. you will only use qualified MCTs who hold the applicable Microsoft Certification credential that is the subject of the MOC title being taught for all your Private Training Sessions using MOC,
- ix. you will only provide access to the Microsoft Instructor-Led Courseware to End Users, and
- x. you will only provide access to the Trainer Content to Trainers.

d. If you are an End User:

For each license you acquire, you may use the Microsoft Instructor-Led Courseware solely for your personal training use. If the Microsoft Instructor-Led Courseware is in digital format, you may access the Microsoft Instructor-Led Courseware online using the unique redemption code provided to you by the training provider and install and use one (1) copy of the Microsoft Instructor-Led Courseware on up to three (3) Personal Devices. You may also print one (1) copy of the Microsoft Instructor-Led Courseware. You may not install the Microsoft Instructor-Led Courseware on a device you do not own or control.

e. If you are a Trainer.

- i. For each license you acquire, you may install and use one (1) copy of the Trainer Content in the form provided to you on one (1) Personal Device solely to prepare and deliver an Authorized Training Session or Private Training Session, and install one (1) additional copy on another Personal Device as a backup copy, which may be used only to reinstall the Trainer Content. You may not install or use a copy of the Trainer Content on a device you do not own or control. You may also print one (1) copy of the Trainer Content solely to prepare for and deliver an Authorized Training Session or Private Training Session.

- ii. You may customize the written portions of the Trainer Content that are logically associated with instruction of a training session in accordance with the most recent version of the MCT agreement. If you elect to exercise the foregoing rights, you agree to comply with the following: (i) customizations may only be used for teaching Authorized Training Sessions and Private Training Sessions, and (ii) all customizations will comply with this agreement. For clarity, any use of “*customize*” refers only to changing the order of slides and content, and/or not using all the slides or content, it does not mean changing or modifying any slide or content.

2.2 **Separation of Components.** The Licensed Content is licensed as a single unit and you may not separate their components and install them on different devices.

2.3 **Redistribution of Licensed Content.** Except as expressly provided in the use rights above, you may not distribute any Licensed Content or any portion thereof (including any permitted modifications) to any third parties without the express written permission of Microsoft.

2.4 **Third Party Notices.** The Licensed Content may include third party code that Microsoft, not the third party, licenses to you under this agreement. Notices, if any, for the third party code are included for your information only.

2.5 **Additional Terms.** Some Licensed Content may contain components with additional terms, conditions, and licenses regarding its use. Any non-conflicting terms in those conditions and licenses also apply to your use of that respective component and supplements the terms described in this agreement.

3. **LICENSED CONTENT BASED ON PRE-RELEASE TECHNOLOGY.** If the Licensed Content’s subject matter is based on a pre-release version of Microsoft technology (“**Pre-release**”), then in addition to the other provisions in this agreement, these terms also apply:
- a. **Pre-Release Licensed Content.** This Licensed Content subject matter is on the Pre-release version of the Microsoft technology. The technology may not work the way a final version of the technology will and we may change the technology for the final version. We also may not release a final version. Licensed Content based on the final version of the technology may not contain the same information as the Licensed Content based on the Pre-release version. Microsoft is under no obligation to provide you with any further content, including any Licensed Content based on the final version of the technology.
 - b. **Feedback.** If you agree to give feedback about the Licensed Content to Microsoft, either directly or through its third party designee, you give to Microsoft without charge, the right to use, share and commercialize your feedback in any way and for any purpose. You also give to third parties, without charge, any patent rights needed for their products, technologies and services to use or interface with any specific parts of a Microsoft technology, Microsoft product, or service that includes the feedback. You will not give feedback that is subject to a license that requires Microsoft to license its technology, technologies, or products to third parties because we include your feedback in them. These rights survive this agreement.
 - c. **Pre-release Term.** If you are an Microsoft IT Academy Program Member, Microsoft Learning Competency Member, MPN Member or Trainer, you will cease using all copies of the Licensed Content on the Pre-release technology upon (i) the date which Microsoft informs you is the end date for using the Licensed Content on the Pre-release technology, or (ii) sixty (60) days after the commercial release of the technology that is the subject of the Licensed Content, whichever is earliest (“**Pre-release term**”). Upon expiration or termination of the Pre-release term, you will irretrievably delete and destroy all copies of the Licensed Content in your possession or under your control.

- 4. SCOPE OF LICENSE.** The Licensed Content is licensed, not sold. This agreement only gives you some rights to use the Licensed Content. Microsoft reserves all other rights. Unless applicable law gives you more rights despite this limitation, you may use the Licensed Content only as expressly permitted in this agreement. In doing so, you must comply with any technical limitations in the Licensed Content that only allows you to use it in certain ways. Except as expressly permitted in this agreement, you may not:
- access or allow any individual to access the Licensed Content if they have not acquired a valid license for the Licensed Content,
 - alter, remove or obscure any copyright or other protective notices (including watermarks), branding or identifications contained in the Licensed Content,
 - modify or create a derivative work of any Licensed Content,
 - publicly display, or make the Licensed Content available for others to access or use,
 - copy, print, install, sell, publish, transmit, lend, adapt, reuse, link to or post, make available or distribute the Licensed Content to any third party,
 - work around any technical limitations in the Licensed Content, or
 - reverse engineer, decompile, remove or otherwise thwart any protections or disassemble the Licensed Content except and only to the extent that applicable law expressly permits, despite this limitation.
- 5. RESERVATION OF RIGHTS AND OWNERSHIP.** Microsoft reserves all rights not expressly granted to you in this agreement. The Licensed Content is protected by copyright and other intellectual property laws and treaties. Microsoft or its suppliers own the title, copyright, and other intellectual property rights in the Licensed Content.
- 6. EXPORT RESTRICTIONS.** The Licensed Content is subject to United States export laws and regulations. You must comply with all domestic and international export laws and regulations that apply to the Licensed Content. These laws include restrictions on destinations, end users and end use. For additional information, see www.microsoft.com/exporting.
- 7. SUPPORT SERVICES.** Because the Licensed Content is “as is”, we may not provide support services for it.
- 8. TERMINATION.** Without prejudice to any other rights, Microsoft may terminate this agreement if you fail to comply with the terms and conditions of this agreement. Upon termination of this agreement for any reason, you will immediately stop all use of and delete and destroy all copies of the Licensed Content in your possession or under your control.
- 9. LINKS TO THIRD PARTY SITES.** You may link to third party sites through the use of the Licensed Content. The third party sites are not under the control of Microsoft, and Microsoft is not responsible for the contents of any third party sites, any links contained in third party sites, or any changes or updates to third party sites. Microsoft is not responsible for webcasting or any other form of transmission received from any third party sites. Microsoft is providing these links to third party sites to you only as a convenience, and the inclusion of any link does not imply an endorsement by Microsoft of the third party site.
- 10. ENTIRE AGREEMENT.** This agreement, and any additional terms for the Trainer Content, updates and supplements are the entire agreement for the Licensed Content, updates and supplements.
- 11. APPLICABLE LAW.**
- a. United States. If you acquired the Licensed Content in the United States, Washington state law governs the interpretation of this agreement and applies to claims for breach of it, regardless of conflict of laws principles. The laws of the state where you live govern all other claims, including claims under state consumer protection laws, unfair competition laws, and in tort.

b. Outside the United States. If you acquired the Licensed Content in any other country, the laws of that country apply.

- 12. LEGAL EFFECT.** This agreement describes certain legal rights. You may have other rights under the laws of your country. You may also have rights with respect to the party from whom you acquired the Licensed Content. This agreement does not change your rights under the laws of your country if the laws of your country do not permit it to do so.
- 13. DISCLAIMER OF WARRANTY. THE LICENSED CONTENT IS LICENSED "AS-IS" AND "AS AVAILABLE." YOU BEAR THE RISK OF USING IT. MICROSOFT AND ITS RESPECTIVE AFFILIATES GIVES NO EXPRESS WARRANTIES, GUARANTEES, OR CONDITIONS. YOU MAY HAVE ADDITIONAL CONSUMER RIGHTS UNDER YOUR LOCAL LAWS WHICH THIS AGREEMENT CANNOT CHANGE. TO THE EXTENT PERMITTED UNDER YOUR LOCAL LAWS, MICROSOFT AND ITS RESPECTIVE AFFILIATES EXCLUDES ANY IMPLIED WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NON-INFRINGEMENT.**
- 14. LIMITATION ON AND EXCLUSION OF REMEDIES AND DAMAGES. YOU CAN RECOVER FROM MICROSOFT, ITS RESPECTIVE AFFILIATES AND ITS SUPPLIERS ONLY DIRECT DAMAGES UP TO US\$5.00. YOU CANNOT RECOVER ANY OTHER DAMAGES, INCLUDING CONSEQUENTIAL, LOST PROFITS, SPECIAL, INDIRECT OR INCIDENTAL DAMAGES.**

This limitation applies to

- anything related to the Licensed Content, services, content (including code) on third party Internet sites or third-party programs; and
- claims for breach of contract, breach of warranty, guarantee or condition, strict liability, negligence, or other tort to the extent permitted by applicable law.

It also applies even if Microsoft knew or should have known about the possibility of the damages. The above limitation or exclusion may not apply to you because your country may not allow the exclusion or limitation of incidental, consequential or other damages.

Please note: As this Licensed Content is distributed in Quebec, Canada, some of the clauses in this agreement are provided below in French.

Remarque : Ce le contenu sous licence étant distribué au Québec, Canada, certaines des clauses dans ce contrat sont fournies ci-dessous en français.

EXONÉRATION DE GARANTIE. Le contenu sous licence visé par une licence est offert « tel quel ». Toute utilisation de ce contenu sous licence est à votre seule risque et péril. Microsoft n'accorde aucune autre garantie expresse. Vous pouvez bénéficier de droits additionnels en vertu du droit local sur la protection des consommateurs, que ce contrat ne peut modifier. La ou elles sont permises par le droit local, les garanties implicites de qualité marchande, d'adéquation à un usage particulier et d'absence de contrefaçon sont exclues.

LIMITATION DES DOMMAGES-INTÉRÊTS ET EXCLUSION DE RESPONSABILITÉ POUR LES DOMMAGES. Vous pouvez obtenir de Microsoft et de ses fournisseurs une indemnisation en cas de dommages directs uniquement à hauteur de 5,00 \$ US. Vous ne pouvez prétendre à aucune indemnisation pour les autres dommages, y compris les dommages spéciaux, indirects ou accessoires et pertes de bénéfices.

Cette limitation concerne:

- tout ce qui est relié au le contenu sous licence, aux services ou au contenu (y compris le code) figurant sur des sites Internet tiers ou dans des programmes tiers; et.
- les réclamations au titre de violation de contrat ou de garantie, ou au titre de responsabilité stricte, de négligence ou d'une autre faute dans la limite autorisée par la loi en vigueur.

Elle s'applique également, même si Microsoft connaissait ou devrait connaître l'éventualité d'un tel dommage. Si votre pays n'autorise pas l'exclusion ou la limitation de responsabilité pour les dommages indirects, accessoires ou de quelque nature que ce soit, il se peut que la limitation ou l'exclusion ci-dessus ne s'appliquera pas à votre égard.

EFFET JURIDIQUE. Le présent contrat décrit certains droits juridiques. Vous pourriez avoir d'autres droits prévus par les lois de votre pays. Le présent contrat ne modifie pas les droits que vous confèrent les lois de votre pays si celles-ci ne le permettent pas.

Revised July 2013

Welcome!

Thank you for taking our training! We've worked together with our Microsoft Certified Partners for Learning Solutions and our Microsoft IT Academies to bring you a world-class learning experience—whether you're a professional looking to advance your skills or a student preparing for a career in IT.

- **Microsoft Certified Trainers and Instructors**—Your instructor is a technical and instructional expert who meets ongoing certification requirements. And, if instructors are delivering training at one of our Certified Partners for Learning Solutions, they are also evaluated throughout the year by students and by Microsoft.
- **Certification Exam Benefits**—After training, consider taking a Microsoft Certification exam. Microsoft Certifications validate your skills on Microsoft technologies and can help differentiate you when finding a job or boosting your career. In fact, independent research by IDC concluded that 75% of managers believe certifications are important to team performance¹. Ask your instructor about Microsoft Certification exam promotions and discounts that may be available to you.
- **Customer Satisfaction Guarantee**—Our Certified Partners for Learning Solutions offer a satisfaction guarantee and we hold them accountable for it. At the end of class, please complete an evaluation of today's experience. We value your feedback!

We wish you a great learning experience and ongoing success in your career!

Sincerely,

Microsoft Learning
www.microsoft.com/learning

Microsoft | Learning

¹ IDC, Value of Certification: Team Certification and Organizational Performance, November 2006

Acknowledgements

Microsoft Learning would like to acknowledge and thank the following for their contribution towards developing this title. Their effort at various stages in the development has ensured that you have a good classroom experience.

Jamie Newman – Content Developer

Jamie Newman became an IT trainer in 1997, first for an IT training company and later for a university, where he became involved in developing courses as well as training them. He began to specialize in databases and eventually moved into database consultancy. In recent years he has specialized in SQL Server and has set up multi user systems that are accessed nationwide. Despite now being more involved with development work, Jamie still likes to deliver IT training courses when the opportunity arises!

John Daisley – Content Developer

John Daisley is a mixed vendor Business Intelligence and Data Warehousing contractor with a wealth of data warehousing and Microsoft SQL Server database administration experience. Having worked in the Business Intelligence arena for over a decade, John has extensive experience of implementing business intelligence and database solutions using a wide range of technologies and across a wide range of industries including airlines, engineering, financial services, and manufacturing.

Nick Anderson – Content Developer

Nick Anderson MBCS MISTC has been a freelance Technical Writer since 1987 and Trainer since 1999. Nick has written internal and external facing content in many business and technical areas including development, infrastructure and finance projects involving SQL Server, Visual Studio and similar tools. Nick provides services for both new and existing document processes from knowledge capture to publishing.

Phil Stollery – Content Developer

Phil has been providing IT consultancy to South West England since graduating in Computer Science. He has worked with small and large organizations to improve their use of SQL Server, predominantly focusing on business information and surrounding technologies such as SharePoint. Most recently, Phil worked with the National Health Service in Gloucestershire on a custom intranet built on SharePoint. A trusted partner, he can communicate at all levels, from technical staff to senior management. Phil brings a wealth of experience that enhances any project.

Geoff Allix – Technical Reviewer

Geoff Allix is a Microsoft SQL Server subject matter expert and professional content developer at Content Master—a division of CM Group Ltd. As a Microsoft Certified Trainer, Geoff has delivered training courses on SQL Server since version 6.5. Geoff is a Microsoft Certified IT Professional for SQL Server and has extensive experience in designing and implementing database and BI solutions on SQL Server technologies, and has provided consultancy services to organizations seeking to implement and optimize database solutions.

Lin Joyner – Technical Reviewer

Lin is an experienced Microsoft SQL Server developer and administrator. She has worked with SQL Server since version 6.0 and previously as a Microsoft Certified Trainer, delivered training courses across the UK. Lin has a wide breadth of knowledge across SQL Server technologies, including BI and Reporting Services. Lin also designs and authors SQL Server and .NET development training materials. She has been writing instructional content for Microsoft for over 15 years.

Contents

Module 1: Introduction to Data Warehousing	
Module Overview	1-1
Lesson 1: Overview of Data Warehousing	1-2
Lesson 2: Considerations for a Data Warehouse Solution	1-9
Lab: Exploring a Data Warehousing Solution	1-17
Module Review and Takeaways	1-22
Module 2: Planning Data Warehouse Infrastructure	
Module Overview	2-1
Lesson 1: Considerations for Data Warehouse Infrastructure	2-2
Lesson 2: Planning Data Warehouse Hardware	2-10
Lab: Planning Data Warehouse Infrastructure	2-19
Module Review and Takeaways	2-21
Module 3: Designing and Implementing a Data Warehouse	
Module Overview	3-1
Lesson 1: Data Warehouse Design Overview	3-2
Lesson 2: Designing Dimension Tables	3-8
Lesson 3: Designing Fact Tables	3-16
Lesson 4: Physical Design for a Data Warehouse	3-19
Lab: Implementing a Data Warehouse	3-34
Module Review and Takeaways	3-40
Module 4: Columnstore Indexes	
Module Overview	4-1
Lesson 1: Introduction to Columnstore Indexes	4-2
Lesson 2: Creating Columnstore Indexes	4-7
Lesson 3: Working with Columnstore Indexes	4-12
Lab: Using Columnstore Indexes	4-17
Module Review and Takeaways	4-21

Module 5: Implementing an Azure SQL Data Warehouse	
Module Overview	5-1
Lesson 1: Advantages of Azure SQL Data Warehouse	5-2
Lesson 2: Implementing an Azure SQL Data Warehouse Database	5-6
Lesson 3: Developing an Azure SQL Data Warehouse	5-11
Lesson 4: Migrating to an Azure SQL Data Warehouse	5-18
Lesson 5: Copying Data with the Azure Data Factory	5-26
Lab: Implement an Azure SQL Data Warehouse	5-33
Module Review and Takeaways	5-39
Module 6: Creating an ETL Solution	
Module Overview	6-1
Lesson 1: Introduction to ETL with SSIS	6-2
Lesson 2: Exploring Source Data	6-7
Lesson 3: Implementing Data Flow	6-14
Lab: Implementing Data Flow in an SSIS Package	6-28
Module Review and Takeaways	6-34
Module 7: Implementing Control Flow in an SSIS Package	
Module Overview	7-1
Lesson 1: Introduction to Control Flow	7-2
Lesson 2: Creating Dynamic Packages	7-11
Lesson 3: Using Containers	7-16
Lab A: Implementing Control Flow in an SSIS Package	7-22
Lesson 4: Managing Consistency	7-27
Lab B: Using Transactions and Checkpoints	7-33
Module Review and Takeaways	7-37
Module 8: Debugging and Troubleshooting SSIS Packages	
Module Overview	8-1
Lesson 1: Debugging an SSIS Package	8-2
Lesson 2: Logging SSIS Package Events	8-8
Lesson 3: Handling Errors in an SSIS Package	8-13
Lab: Debugging and Troubleshooting an SSIS Package	8-19
Module Review and Takeaways	8-24

Module 9: Implementing a Data Extraction Solution

Module Overview	9-1
Lesson 1: Introduction to Incremental ETL	9-2
Lesson 2: Extracting Modified Data	9-12
Lab A: Extracting Modified Data	9-26
Lesson 3: Loading Modified Data	9-38
Lesson 4: Temporal Tables	9-51
Lab B: Loading a Data Warehouse	9-62
Module Review and Takeaways	9-72

Module 10: Enforcing Data Quality

Module Overview	10-1
Lesson 1: Introduction to Data Quality	10-2
Lesson 2: Using Data Quality Services to Cleanse Data	10-11
Lab A: Cleansing Data	10-15
Lesson 3: Using Data Quality Services to Match Data	10-21
Lab B: Deduplicating Data	10-27
Module Review and Takeaways	10-31

Module 11: Master Data Services

Module Overview	11-1
Lesson 1: Introduction to Master Data Services	11-2
Lesson 2: Implementing a Master Data Services Model	11-6
Lesson 3: Hierarchies and Collections	11-18
Lesson 4: Creating a Master Data Hub	11-27
Lab: Implementing Master Data Services Model	11-34
Module Review and Takeaways	11-42

Module 12: Extending SQL Server Integration Services

Module Overview	12-1
Lesson 1: Using Scripts in SSIS	12-2
Lesson 2: Using Custom Components in SSIS	12-10
Lab: Using Custom Scripts	12-14
Module Review and Takeaways	12-16

Module 13: Deploying and Configuring SSIS Packages	
Module Overview	13-1
Lesson 1: Overview of SSIS Development	13-2
Lesson 2: Deploying SSIS Projects	13-5
Lesson 3: Planning SSIS Package Execution	13-14
Lab: Deploying and Configuring SSIS Packages	13-20
Module Review and Takeaways	13-24
Module 14: Consuming Data in a Data Warehouse	
Module Overview	14-1
Lesson 1: Introduction to Business Intelligence	14-2
Lesson 2: Introduction to Data Analysis	14-6
Lesson 3: Introduction to Reporting	14-9
Lesson 4: Analyzing Data with Azure SQL Data Warehouse	14-12
Lab: Using a Data Warehouse	14-15
Module Review and Takeaways	14-18
Lab Answer Keys	
Module 1 Lab: Exploring a Data Warehousing Solution	L01-1
Module 2 Lab: Planning Data Warehouse Infrastructure	L02-1
Module 3 Lab: Implementing a Data Warehouse	L03-1
Module 4 Lab: Using Columnstore Indexes	L04-1
Module 5 Lab: Implement an Azure SQL Data Warehouse	L05-1
Module 6 Lab: Implementing Data Flow in an SSIS Package	L06-1
Module 7 Lab A: Implementing Control Flow in an SSIS Package	L07-1
Module 7 Lab B: Using Transactions and Checkpoints	L07-8
Module 8 Lab: Debugging and Troubleshooting an SSIS Package	L08-1
Module 9 Lab A: Extracting Modified Data	L09-1
Module 9 Lab B: Loading a Data Warehouse	L09-19
Module 10 Lab A: Cleansing Data	L10-1
Module 10 Lab B: Deduplicating Data	L10-8
Module 11 Lab: Implementing Master Data Services Model	L11-1
Module 12 Lab: Using Custom Scripts	L12-1
Module 13 Lab: Deploying and Configuring SSIS Packages	L13-1
Module 14 Lab: Using a Data Warehouse	L14-1

About This Course

This section provides a brief description of the course, audience, suggested prerequisites, and course objectives.

Course Description

 **Note:** This first release ('C') MOC version of course 20767 has been developed on the RTM version of the SQL Server 2017 and supersedes any earlier releases.

This 5-day instructor led course describes how to implement a data warehouse platform to support a BI solution. Students will learn how to create a data warehouse with Microsoft® SQL Server® and with Azure SQL Data Warehouse, to implement ETL with SQL Server Integration Services, and to validate and cleanse data with SQL Server Data Quality Services and SQL Server Master Data Services.

Audience

The primary audience for this course are database professionals who need to fulfil a Business Intelligence Developer role. They will need to focus on hands-on work creating BI solutions including Data Warehouse implementation, ETL, and data cleansing.

Student Prerequisites

In addition to their professional experience, students who attend this training should already have the following technical knowledge:

- At least 2 years' experience of working with relational databases, including:
 - Designing a normalized database.
 - Creating tables and relationships.
 - Querying with Transact-SQL.
 - Some exposure to basic programming constructs (such as looping and branching).
 - An awareness of key business priorities such as revenue, profitability, and financial accounting is desirable.

Course Objectives

After completing this course, students will be able to:

- Describe the key elements of a data warehousing solution
- Describe the main hardware considerations for building a data warehouse
- Implement a logical design for a data warehouse
- Implement a physical design for a data warehouse
- Create columnstore indexes
- Implementing an Azure SQL Data Warehouse
- Describe the key features of SSIS
- Implement a data flow by using SSIS
- Implement control flow by using tasks and precedence constraints
- Create dynamic packages that include variables and parameters

- Debug SSIS packages
- Describe the considerations for implement an ETL solution
- Implement Data Quality Services
- Implement a Master Data Services model
- Describe how you can use custom components to extend SSIS
- Deploy SSIS projects
- Describe BI and common BI scenarios

Course Outline

The course outline is as follows:

Module 1, "Introduction to data warehousing"

Module 2, "Planning Data Warehouse Infrastructure"

Module 3, "Designing and Implementing a Data Warehouse "

Module 4, "Columnstore Indexes"

Module 5, "Implementing an Azure SQL Data Warehouse"

Module 6, "Creating an ETL Solution"

Module 7, "Implementing Control Flow in an SSIS Package"

Module 8, "Debugging and Troubleshooting SSIS Packages"

Module 9, "Implementing a Data Extraction Solution"

Module 10, "Enforcing Data Quality"

Module 11, "Master Data Services"

Module 12, "Extending SQL Server Integration Services"

Module 13, "Deploying and Configuring SSIS Packages"

Module 14, "Consuming Data in a Data Warehouse"

Course Materials

The following materials are included with your kit:

- **Course Handbook:** a succinct classroom learning guide that provides the critical technical information in a crisp, tightly-focused format, which is essential for an effective in-class learning experience.
 - **Lessons:** guide you through the learning objectives and provide the key points that are critical to the success of the in-class learning experience.
 - **Labs:** provide a real-world, hands-on platform for you to apply the knowledge and skills learned in the module.
 - **Module Reviews and Takeaways:** provide on-the-job reference material to boost knowledge and skills retention.
 - **Lab Answer Keys:** provide step-by-step lab solution guidance.



Additional Reading: Course Companion Content on the <http://www.microsoft.com/learning/en/us/companion-moc.aspx> **Site:** searchable, easy-to-browse digital content with integrated premium online resources that supplement the Course Handbook.

- **Modules:** include companion content, such as questions and answers, detailed demo steps and additional reading links, for each lesson. Additionally, they include Lab Review questions and answers and Module Reviews and Takeaways sections, which contain the review questions and answers, best practices, common issues and troubleshooting tips with answers, and real-world issues and scenarios with answers.
- **Resources:** include well-categorized additional resources that give you immediate access to the most current premium content on TechNet, MSDN®, or Microsoft® Press®.



Additional Reading: Student Course files on the <http://www.microsoft.com/learning/en/us/companion-moc.aspx> **Site:** includes the Allfiles.exe, a self-extracting executable file that contains all required files for the labs and demonstrations.

- **Course evaluation:** at the end of the course, you will have the opportunity to complete an online evaluation to provide feedback on the course, training facility, and instructor.
- To provide additional comments or feedback on the course, send email to mcspprt@microsoft.com. To inquire about the Microsoft Certification Program, send an email to mcphelp@microsoft.com.

Virtual Machine Environment

This section provides the information for setting up the classroom environment to support the business scenario of the course.

Virtual Machine Configuration

In this course, you will use Microsoft® Hyper-V® to perform the labs.



Note: At the end of each lab, you must revert the virtual machines to a snapshot. You can find the instructions for this procedure at the end of each lab

The following table shows the role of each virtual machine that is used in this course:

Virtual machine	Role
20767C-MIA-DC	Domain controller for the ADVENTUREWORKS domain.
20767C-MIA-SQL	SQL Server and SharePoint Server

Software Configuration

The following software is installed:

- Microsoft Windows Server 2016
- Microsoft SQL Server 2017
- Microsoft Office 2016
- Microsoft SharePoint Server 2013
- Microsoft Visual Studio 2015
- Microsoft Visio 2013

Course Files

The files associated with the labs in this course are located in the D:\Labfiles folder on the 20767C-MIA-SQL virtual machine.

Classroom Setup

Each classroom computer will have the same virtual machine configured in the same way.

Course Hardware Level

To ensure a satisfactory student experience, Microsoft Learning requires a minimum equipment configuration for trainer and student computers in all Microsoft Certified Partner for Learning Solutions (CPLS) classrooms in which Official Microsoft Learning Product courseware is taught.

Hardware Level 6+

- Intel Virtualization Technology (Intel VT) or AMD Virtualization (AMD-V) processor
- Dual 120 GB hard disks 7200 RPM SATA or better*
- 8GB or higher
- DVD drive

- Network adapter with Internet connectivity
- Super VGA (SVGA) 17-inch monitor
- Microsoft Mouse or compatible pointing device
- Sound card with amplified speakers

*Striped

In addition, the instructor computer must be connected to a projection display device that supports SVGA 1024 x 768 pixels, 16 bit colors.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 1

Introduction to Data Warehousing

Contents:

Module Overview	1-1
Lesson 1: Overview of Data Warehousing	1-2
Lesson 2: Considerations for a Data Warehouse Solution	1-9
Lab: Exploring a Data Warehousing Solution	1-17
Module Review and Takeaways	1-22

Module Overview

Data warehousing is a solution that organizations can use to centralize business data for reporting and analysis. Implementing a data warehouse solution can provide a business or other organization with significant benefits, including:

- Comprehensive and accurate reporting of key business information.
- A centralized source of business data for analysis and decision-making.
- The foundation for an enterprise business intelligence (BI) solution.

This module provides an introduction to the key components of a data warehousing solution and the high level considerations you must take into account when you embark on a data warehousing project.

Objectives

After completing this module, you will be able to:

- Describe the key elements of a data warehousing solution.
- Describe the key considerations for a data warehousing project.

Lesson 1

Overview of Data Warehousing

Data warehousing is a well-established technique for centralizing business data for reporting and analysis. Although the specific details of individual solutions can vary, there are some common elements in most data warehousing implementations. Familiarity with these elements will help you to better plan and build an effective data warehousing solution.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the business problem addressed by data warehouses.
- Define a data warehouse.
- Describe the commonly-used data warehouse architectures.
- Identify the components of a data warehousing solution.
- Describe a high level approach to implementing a data warehousing project.
- Identify the roles involved in a data warehousing project.
- Describe the components and features of Microsoft® SQL Server® and other Microsoft products that you can use in a data warehousing solution.

The Business Problem

Running a business effectively can present a significant challenge, particularly as it grows or is affected by trends in its target market or the global economy. To be successful, a business must be able to adapt to changing conditions. That requires individuals within the organization to make good strategic and tactical decisions. However, the following problems can often make effective business decision-making difficult:

- Key business data is distributed across multiple systems. This makes it hard to collate all the information necessary for a particular business decision.
- Finding the information required for business decision-making is time-consuming and error-prone. The need to gather and reconcile data from multiple sources results in slow, inefficient decision-making processes that can be further undermined by inconsistencies between duplicate, contradictory sources of the same information.
- Most business decisions require the answer to fundamental questions, such as “How many customers do we have?” or “Which products do we sell most often?” Although these questions may seem simple, the distribution of data throughout multiple systems in a typical organization can make them difficult, or even impossible, to answer.

By resolving these problems, you can make effective decisions that will help the business to be more successful, both at a strategic, executive level, and during day-to-day operations.

A successful business needs to be able to adapt—the following problems make that difficult:

1. Business data is spread across many systems
2. Data can be inconsistent, duplicated, and contradictory
3. Fundamental questions can't be easily answered

What Is a Data Warehouse?

A data warehouse provides a solution to the problem of distributed data that prevents effective business decision-making. There are many definitions for the term "data warehouse," and disagreements over specific implementation details. It is generally agreed, however, that a data warehouse is a centralized store of business data that can be used for reporting and analysis to inform key business decisions.

Typically, a data warehouse:

- Contains a large volume of data that relates to historical business transactions.
- Is optimized for read operations that support querying the data. This is in contrast to a typical Online Transaction Processing (OLTP) database that is designed to support data insert, in addition to update and delete operations.
- Is loaded with new or updated data at regular intervals.
- Provides the basis for enterprise BI applications.

A centralized store of business data for reporting and analysis that typically:

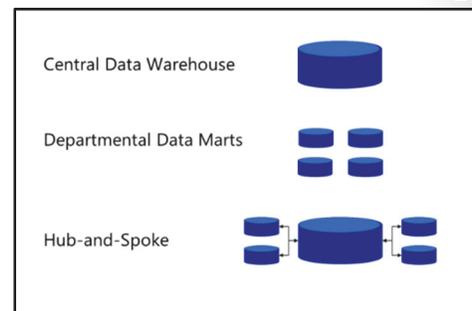
- Contains large volumes of historical data
- Is optimized for querying, as opposed to inserting or updating data
- Is incrementally loaded with new business data at regular intervals
- Provides the basis for enterprise BI solutions

Data Warehouse Architectures

There are many ways that you can implement a data warehouse solution in an organization. Some common approaches include:

- Creating a single, central enterprise data warehouse for all business units.
- Creating small, departmental data marts for individual business units.
- Creating a hub-and-spoke architecture that synchronizes a central enterprise data warehouse with departmental data marts containing a subset of the data warehouse data.

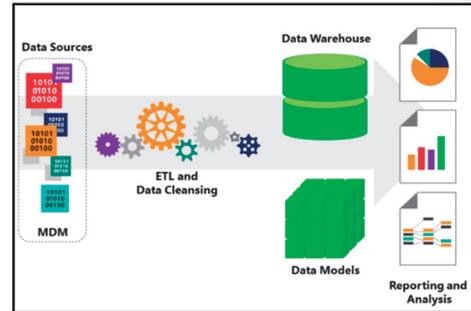
The correct architecture for a given business might be one of these, or a combination of various elements from all three approaches.



Components of a Data Warehousing Solution

A data warehousing solution usually consists of the following elements:

- **Data Sources.** Sources of business data for the data warehouse, often including OLTP application databases and data exported from proprietary systems, such as accounting applications.
- **An Extract, Transform, and Load (ETL) process.** A workflow for accessing data in the data sources, modifying it to conform to the data model for the data warehouse, and loading it into the data warehouse.
- **Data Staging Areas.** Intermediary locations where the data to be transferred to the data warehouse is stored. It is prepared here for import and to synchronize loading into the data warehouse.
- **Data Warehouse.** A relational database designed to provide high-performance querying of historical business data for reporting and analysis.
- **Data Models.** Based upon the data within the data warehouse, data models will be used by a business analyst to produce analytics and predictions for the business.



Many data warehousing solutions also include:

- **Data Cleansing and Deduplication.** A solution for resolving data quality issues before it is loaded into the data warehouse.
- **Master Data Management (MDM).** A solution that provides an authoritative data definition for business entities used by multiple systems across the organization.

Data Warehousing Projects

A data warehousing project has much in common with any other IT implementation, so you can apply most commonly-used methodologies, such as Agile, Waterfall, or Microsoft Solutions Framework (MSF). However, a data warehousing project often requires a deeper understanding of the key business objectives and metrics used to drive decision-making.

A high level approach to implementing a data warehousing project usually includes the following steps:

1. **Work** with business stakeholders and information workers to determine the questions the data warehouse must answer. They might include:
 - What was the total sales revenue for each geographic territory in a given month?
 - What are our most profitable products or services?

1. What are the business questions that need to be answered?
2. What data is required to answer them?
3. Where is this data and how easy is it to obtain?
4. Knowing the above, assess the importance of each question against the ability to answer it from existing data

- Are business costs growing or reducing over time?
 - Which sales employees are meeting their sales targets?
2. **Determine** the data required to answer these questions. It is normal to think of this data in terms of “facts” and “dimensions.” Facts contain the numerical measures required to aggregate data, so you can answer the business questions identified in step 1. For example, to determine sales revenue, you might need the sales amount for each individual transaction. Dimensions represent the different aspects of the business by which you want to aggregate the measures. For example, to determine sales revenue for each territory in a given month, you might need a geographic dimension, so you can aggregate sales by territory, and a time dimension where you can aggregate sales by month. Fact and dimensional modeling is covered in more detail in Module 3: *Designing and Implementing a Data Warehouse*.
 3. **Identify** data sources containing the data required to answer the business questions. These are commonly relational databases used by existing line-of-business applications, but they can also include:
 - Flat files or XML documents that have been extracted from proprietary systems.
 - Data in Microsoft SharePoint® lists.
 - Commercially available data that has been purchased from a data supplier such as the Microsoft Windows Azure® Marketplace.
 4. **Determine** the priority of each business question based on:
 - The importance of answering the question in relation to driving key business objectives.
 - The feasibility of answering the question from the available data.

A common approach to prioritizing the business questions you will address in the data warehousing solution is to work with key business stakeholders and plot each question on a quadrant-based matrix like the one shown below. The position of the questions in the matrix helps you to agree the scope and clarify the requirements of the data warehousing project.

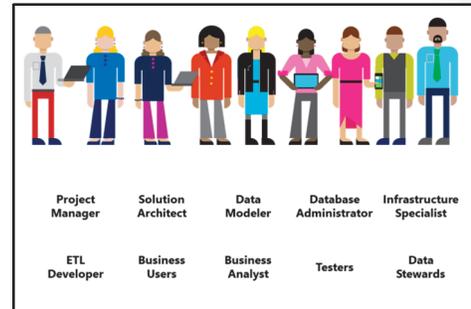
Business importance of the question	High importance, low feasibility	High importance, high feasibility
	Low importance, low feasibility	Low importance, high feasibility
Feasibility of answering the question		

If a large number of questions fall into the high importance, high feasibility category, you may want to consider breaking down the project into subprojects. Each subproject tackles the problem of implementing the data warehouse schema, ETL solution, and data quality procedures for a specific area of the business, starting with the highest priority questions. If you choose to adopt this incremental approach, take care to create an overall design for dimension and fact tables in early iterations of the solution, so that subsequent additions can reuse them.

Data Warehousing Project Roles

A data warehousing project typically involves several roles, including:

- A **project manager**. Coordinates project tasks and schedules, responsible for ensuring that the project is completed on time and within budget.
- A **solution architect**. Has overall responsibility for the technical design of the data warehousing solution.
- A **data modeler**. Designs the data warehouse schema.
- A **database administrator**. Designs the physical architecture and configuration of the data warehouse database. In addition, database administrators with responsibility for data sources used in the data warehousing solution must be involved in the project to provide access to the data sources required by the ETL process.
- An **infrastructure specialist**. Implements the server and network infrastructure for the data warehousing solution.
- An **ETL developer**. Builds the ETL workflow for the data warehousing solution.
- **Business users**. Provide requirements and help to prioritize the questions that the data warehousing solution will answer. Often, the team includes a business analyst as a full-time member, helping to interpret the questions and ensuring that the solution design meets the users' needs. They will also be the end users of the finished solution.
- **Testers**. Verify the business and operational functionality of the solution as it is developed.
- **Data stewards**. For each key subject area in the data warehousing solution. Determine data quality rules and validate data before it enters the data warehouse. Data stewards are sometimes referred to as data governors.



In addition to ensuring the appropriate assignment of these roles, you should consider the importance of executive-level sponsorship of the data warehousing project. It is significantly more likely to succeed if a high-profile executive sponsor is seen to actively support the creation of the data warehousing solution and its associated costs.

SQL Server As a Data Warehousing Platform

SQL Server includes components and features that you can use to implement various architectural elements of a data warehousing solution, including:

- **Database Engine.** The database engine component includes many features that are directly relevant to a data warehouse project, the main feature being a robust relational database. Other features include:
 - **In-Memory OLTP** and in-memory optimizations, which enable businesses to have a high-performance, real-time transactional database, while also layering column-based query processing on top to achieve up to 10 times the query performance of traditional row-based storage.
 - **Temporal Tables** that provide businesses with time context to data being captured by maintaining a full history of all data changes.
- **SQL Server Integration Services (SSIS).** A component for building data integration and transformation solutions, commonly referred to as Extract, Transform, and Load (ETL). Improvements in SQL Server include:
 - The ability to connect to Hadoop (Big Data) and Teradata sources.
 - Incremental package deployments, which enable individual packages to be upgraded without having to deploy the project as a whole.
- **Data Quality Services (DQS).** Consisting of Data Quality Server and Data Quality Client, businesses can use DQS to build a knowledge base around their data. You can then standardize, deduplicate, and enrich this data.
- **Master Data Services (MDS).** A component that gives businesses the tools to discover and define nontransactional lists of data. The goal is to compile a maintainable master list of business rules and data types.
- **SQL Server Analysis Services (SSAS).** An online analytical data engine that provides analytical data upon which business reports can be built (see the following description of Reporting Services). SSAS can create two different kinds of semantic models—Tabular (relational modeling) or Multidimensional (OLAP cubes).
- **SQL Server Reporting Services (SSRS).** The main goal of developing a data warehouse is to provide answers to business critical questions. These answers normally take the form of reports. SSRS is a server-based reporting platform that includes a comprehensive set of tools to create, manage, and deliver business reports.

Core Data Warehousing

- Database Engine
- Integration Services
- Master Data Services
- Data Quality Services

Business Intelligence

- Analysis Services
- Reporting Services

For more information about the other components as they relate more closely to BI, you should attend the following courses—20766A: *Developing SQL Data models* and 10988A: *Managing SQL Business Intelligence Operations*.

Sequencing Activity

Put the following steps in a data warehouse project in order by numbering each to indicate the correct order.

	Steps
	Work with business stakeholders to identify key questions.
	Identify the data that might be able to answer the identified questions.
	Identify where data can be sourced.
	Prioritize questions that need answers.
	Identify the people or groups who require access to the produced data.

Lesson 2

Considerations for a Data Warehouse Solution

Before starting a data warehousing project, you need to understand the considerations that will help you create a data warehousing solution to address your businesses-specific needs and constraints.

This lesson describes some of the key considerations for planning a data warehousing solution.

Lesson Objectives

After completing this lesson, you will be able to describe considerations for:

- Designing a data warehouse database.
- Columnstore indexes.
- Data sources.
- Designing an ETL process.
- Implementing data quality and master data management.

Data Warehouse Database and Storage

A data warehouse is a relational database that is optimized for reading data for analysis and reporting. When you are planning a data warehouse, you should take the following considerations into account:

Database Schema

The logical schema of a data warehouse is typically designed to denormalize data into a structure that minimizes the number of join operations required in the queries used to retrieve and aggregate data.

A common approach is to design a star schema in which numerical measures are stored in fact tables.

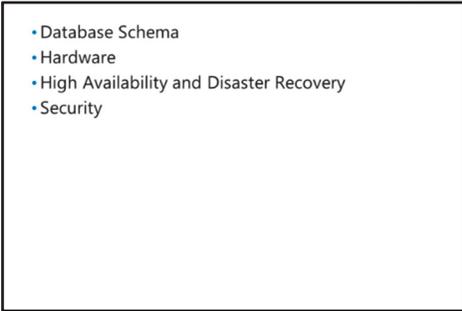
These have foreign keys to multiple dimension tables containing the business entities by which the measures can be aggregated.

Before designing your data warehouse, you must know:

- Which dimensions your business users employ when aggregating data.
- Which measures need to be analyzed and at what granularity.
- Which facts include those measures?

You must also carefully plan the keys that will be used to link facts to dimensions, and consider whether your data warehouse needs to support the use of dimensions that change over time. For example, handling dimension records for customers who change their address.

You must also consider the physical implementation of the database, because this will affect the performance and manageability of the data warehouse. It is common to use table partitioning to distribute large fact data across multiple file groups, each on a different physical disk. This can increase query performance and makes it easy for you to implement a file group-based backup strategy that can

- 
- Database Schema
 - Hardware
 - High Availability and Disaster Recovery
 - Security

help reduce downtime in the event of a single disk failure. You should also consider the appropriate indexing strategy for your data, and whether to use data compression when storing it.



Note: Designing a data warehouse schema is covered in more detail in Module 3: *Designing and Implementing a Data Warehouse*.

Hardware

The choice of hardware for your data warehouse solution can make a significant difference to the performance, manageability, and cost of the data warehouse. The hardware considerations include:

- Query processing requirements, including anticipated peak memory and CPU utilization.
- Storage volume and disk input/output requirements.
- Network connectivity and bandwidth.
- Component redundancy for high availability.
- In-memory requirements.

You can build your own data warehouse solution by purchasing and assembling individual components, using a pretested reference architecture, or purchasing a hardware appliance that includes preconfigured components in a ready-to-use package. Factors that will influence your choice of hardware include:

- How much budget is available?
- Existing enterprise agreements with hardware vendors.
- Time and expertise to construct the solution.
- Hardware assembly and configuration expertise.



Note: Hardware for data warehousing solutions is discussed in more detail in Module 2: *Planning Data Warehouse Infrastructure*.

High Availability and Disaster Recovery

A data warehouse can very quickly become a business-critical part of your overall application infrastructure, so it is essential to consider how you will ensure its availability. SQL Server includes support for several high-availability techniques, including database mirroring and server clustering. You must assess these technologies and choose the best one for your individual solution, based on:

- Failover time requirements.
- Hardware requirements and cost.
- Configuration and management complexity.

In addition to a server-level, high-availability solution, you must consider redundancy at the individual component level for network interfaces and storage arrays.

The most robust high-availability solution cannot protect your data warehouse from every eventuality, so you must also plan a suitable disaster recovery solution that includes a comprehensive backup strategy. This should take into account:

- The volume of data in the data warehouse.
- The frequency of changes to data in the data warehouse.

- The effect of the backup process on data warehouse performance.
- The time to recover the database in the event of a failure.

Security

Your data warehouse contains a huge volume of data that is typically commercially sensitive. You may want to provide access to some data by all users, but restrict access to other data for a subset of users.

Considerations for securing your data warehouse include:

- The authentication mechanisms that you must support to provide access to the data warehouse.
- The permissions that the various users who access the data warehouse will require.
- The connections over which data is accessed.
- The physical security of the database and backup media.

SQL Server includes a new option, Always Encrypted, which enables you to perform database operations on encrypted data. The keys are stored on the client, not the server, so the data owners (business users) no longer need to worry about data managers (DBAs and developers) being able to view sensitive data.

Columnstore Indexes

Microsoft introduced columnstore indexes in SQL Server 2012. This new index type references data in a column fashion and uses compression aggressively to reduce the disk space used and the I/O operations needed to respond to a query.

Microsoft has increased the capabilities of the database engine to incorporate the following new features for columnstore indexes:

- Clustered columnstore indexes with nonclustered row indexes (b-tree index) on top facilitates high performance reporting to sit on top of OLTP tables.
- Primary and foreign key constraints are supported through nonclustered row indexes.
- Now supported on in-memory tables.
- Batch mode enabled on more aggregate functions.

For details of the differences between the SQL Server 2012 and the latest implementations of columnstore indexes, see the SQL Server Technical Documentation:

 **Columnstore Indexes Versioned Feature Summary**

<http://aka.ms/t7wvb0>

Row Index or Columnstore Index?

When you are seeking specific data, row indexes generally work best. For example, an OLTP database will be updating, inserting, and deleting specific rows to perform changes, such as updating a customer's address or last name.

Row-based index	Column-based index
<ul style="list-style-type: none"> • Can be clustered and nonclustered • Improve performance on row level queries, inserts and updates • Best used in OLTP databases • All data in a row is processed 	<ul style="list-style-type: none"> • Can be clustered and nonclustered • Improve performance on queries that scan a table, aggregation and analytical queries • Best used in data warehouses • Only columns needed are processed
<p>A clustered columnstore index can be combined with multiple nonclustered row indexes to have the benefits of both types</p>	

When you are scanning data in a table, columnstore indexes are often the better choice. For example, in a data warehouse using aggregate functions to return the sum of all sales for a region.

In SQL Server 2012 and SQL Server 2014, database architects had to choose to use one type of clustered index. With the latest build of SQL Server, you can have either type of clustered index, in addition to some supporting nonclustered indexes of the other type.

Data Sources

You must identify the sources that provide the data for your data warehouse, and consider the following factors when planning your solution:

Data Source Connection Types

Your data warehouse might require data from a variety of sources. For each source, you must consider how the ETL process can connect and extract the data you want to use. In many cases, your data sources will be relational databases for which you can use an OLEDB or ODBC provider. However, some data sources might use proprietary storage that needs a bespoke provider, or for which no provider exists. In these cases, you must develop a custom provider or determine whether you can export data from the data source in a format that the ETL process can easily consume—such as XML, Excel files, or comma-delimited text.

- Data Source Connection Types
- Credentials and Permissions
- Data Formats
- Data Acquisition Windows

Credentials and Permissions

Most data sources require secure access with user authentication and, potentially, individual permissions on the data. You must work with the owners of the data sources used in your data warehousing solution to establish:

- The credentials that your ETL process can use to access the data source.
- The required permissions to access the data used by the data warehouse.

Data Formats

A data source might store data in a different format to the desired output. Your solution must take into account issues arising from this, including:

- Conversion of data from one data type to another—for example, extracting numeric values from a text file.
- Truncation of data when copying to a destination with a limited data length.
- Date and time formats used in data sources.
- Numeric formats, scales, and precisions.
- Support for Unicode characters.

Data Acquisition Windows

Depending on the workload patterns of the business, there might be times when individual data sources are not available or the usage level is such that the additional overhead of a data extraction is undesirable.

When planning a data warehousing solution, you must work with each data source owner to determine appropriate data acquisition windows based on:

- The workload pattern of the data source, and its resource utilization and capacity levels.
- The volume of data to be extracted, and the time that process takes.
- The frequency with which you need to update the data warehouse with fresh data.
- If applicable, the time zones in which business users are accessing the data.
- How the data is produced and whether it resides in a real-time or batch-based system.

Extract, Transform, and Load Processes

A significant part in the creation of a data warehouse solution is the implementation of one or more ETL processes. When you design an ETL process for a data warehousing solution, you must consider the following factors:

Staging

In some solutions, you can transfer data directly from data sources to the data warehouse without any intermediary staging. In many cases, however, you should consider staging data to:

- Synchronize a data warehouse refresh that includes source data extracted during multiple data acquisition windows.
- Perform data validation, cleansing, and deduplication operations on the data before loading it into the data warehouse.
- Perform transformations on the data that you cannot perform during the data extraction or data flow processes.

If you require a staging area in your solution, you must decide on a format for the staged data. Possible formats include:

- A relational database.
- Text or XML files.
- Raw files (binary files in a proprietary format of the ETL platform being used).

Your choice of format should be based on several factors including:

- The need to access and modify the staged data.
- The time taken to store and read the staging data.

Finally, if using a relational database as the staging area, you must decide where this database will reside.

Possible choices include:

- A dedicated staging server.
- A dedicated SQL Server instance on the data warehouse server.
- A dedicated staging database in the same instance of SQL Server as the data warehouse.
- A collection of staging tables (perhaps in a dedicated schema) in the data warehouse database.

Staging:

- What data must be staged?
- Staging data format

Required transformations:

- Transformations during extraction versus data flow transformations

Incremental ETL:

- Identifying data changes for extraction
- Inserting or updating when loading

Factors you should consider when deciding the location of the staging database include:

- Server hardware requirements and cost.
- The time taken to transfer data across network connections.
- The use of Transact-SQL loading techniques that perform better when the staging data and data warehouse are co-located on the same SQL Server instance.
- The server resource overheads associated with the staging and data warehouse load processes.

Required Transformations

Most ETL processes require that the data being extracted from data sources is modified to match the data warehouse schema. When planning an ETL process for a data warehousing solution, you must examine the source data and destination schema, and identify what transformations are required. You must then determine the optimal place in the ETL process to perform these transformations. Choices for implementing data transformations include:

- **During the data extraction.** For example, by concatenating two fields in a SQL Server data source into a single field in the Transact-SQL query used to extract the data.
- **In the data flow.** For example, by using a Derived Column data transformation task in a SQL Server Integration Services data flow.
- **In the staging area.** For example, by using a Transact-SQL query to apply default values to null fields in a staging table.

Factors affecting the choice of data transformation technique include:

- **The performance overhead of the transformation.** Typically, it is best to use the approach with the least performance overhead. Set-based operations that are performed in Transact-SQL queries usually function better than row-based transformations applied in a data flow.
- **The level of support for querying and updating in the data source or staging area.** Cases where you extract data from a comma-delimited file and stage it in a raw file limit your options to performing row-by-row transformations in the data flow.
- **Dependencies on data required for the transformation.** For example, you might need to look up a value in one data source to obtain additional data from another. In this case, you must perform the data transformation in a location where both data sources are accessible.
- **The complexity of the logic involved in the transformation.** In some cases, a transformation may require multiple steps and branches depending on the presence or value of specific data fields. In these cases, it is often easier to apply the transformation by combining several steps in a data flow, than it is to create a Transact-SQL statement to perform the transformation.

Incremental ETL

After the initial load of the data warehouse, you will usually need to incrementally add new or updated source data. When you plan your data warehousing solution, you must consider the following factors that relate to incremental ETL:

- How will you identify new or modified records in the data sources?
- Do you need to delete records in the data warehouse when corresponding records in the data sources are removed? If so, will you physically delete the records, or simply mark them as inactive (often referred to as a logical delete)?
- How will you determine whether a record that is to be loaded into the data warehouse should be new or an update to an existing one?

- Are there records in the data warehouse for which you need to preserve historical values by creating a new version of the row instead of updating the existing one?



Note: Managing data changes in an incremental ETL process is discussed in more detail in Module 9: *Implementing an Incremental ETL Process*.

Data Quality and Master Data Management

The value of a data warehouse is largely determined by the quality of the data it contains. When planning a data warehousing project, therefore, you should determine how you will ensure data quality. You should consider using a master data management solution.

Data Quality

To validate and enforce the quality of data in your data warehouse, it is recommended that business users who have knowledge of each subject area addressed by the data warehouse become data stewards. A data steward's responsibilities include:

- Building and maintaining a knowledge base that identifies common data errors and corrections.
- Validating data against the knowledge base.
- Ensuring that consistent values are used for data attributes where multiple forms of the value might be considered valid. For example, ensuring that a Country field always uses the value "United States" when referring to America, even though "USA," "The U.S." and "America" are also valid values.
- Identifying and correcting missing data values.
- Identifying and consolidating duplicate data entities, such as records for "Robert Smith" and "Bob Smith", which both refer to the same physical customer.

You can use SQL Server Data Quality Services to provide a data quality solution that helps the data steward to perform these tasks.

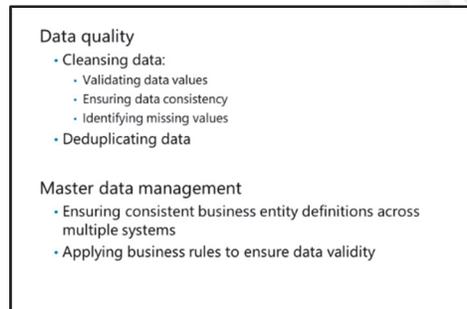


Note: SQL Server Data Quality Services is discussed in more detail in Module 10: *Enforcing Data Quality*.

Master Data Management

It is common for large organizations to have multiple business applications, and in many cases, these systems perform tasks that are related to the same business entities. For example, an organization may have an e-commerce application enabling customers to purchase products, and a separate inventory management system that also stores product data. A record representing a particular product might exist in both systems. In this scenario, it can be useful to implement a master data management system that provides an authoritative definition of each business entity (in this example, a particular product) that you can use across multiple applications to ensure consistency.

Master data management is especially important in a data warehousing scenario because it ensures that data conforms to the agreed definition for the business entities to be included in any analysis and



reporting solutions that it must support. You can use SQL Server Master Data Services to implement a master data management solution.



Note: SQL Server Master Data Services is discussed in more detail in Module 11: *Using Master Data Services*.

Categorize Activity

Which of the following are optimal places to perform data transformations in an ETL solution?

Items	
1	Perform simple transformations in the query that is used to extract the data.
2	All the data transformations should be completed by the derived reports.
3	Implement transformations in a data flow task, within an ETL process, when complex row-by-row processing is required.
4	There is no optimal place to perform data transformations.
5	Perform complex transformations in a staging database, where it is possible to combine data from multiple sources.

Category 1	Category 2
Optimal solution	Suboptimal solution

Lab: Exploring a Data Warehousing Solution

Scenario

The labs in this course are based on a fictional company called Adventure Works Cycles that manufactures and sells cycles and cycling accessories to customers all over the world. Adventure Works sells direct to customers through an e-commerce website and also through an international network of resellers.

Throughout this course, you will develop a data warehousing solution for Adventure Works Cycles, including: a data warehouse; an ETL process to extract data from source systems and populate the data warehouse; a data quality solution; and a master data management solution.

The lab for this module provides a high level overview of the solution that you will create in later labs. You can use this lab to become familiar with the various elements of the data warehousing solution that you will learn to build in later modules.

Objectives

After completing this lab, you will be able to:

- Describe the data sources in the Adventure Works data warehousing scenario.
- Describe the ETL process in the Adventure Works data warehousing scenario.
- Describe the data warehouse in the Adventure Works data warehousing scenario.

Estimated Time: 30 minutes

Virtual Machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Exploring Data Sources

Scenario

Adventure Works uses various software applications to manage different aspects of the business, and each has its own data store. Specifically, these are:

- Internet sales processed through an e-commerce web application.
- Reseller sales processed by sales representatives, who use a specific application. Sales employee details are stored in a separate human resources system.
- Reseller payments are processed by an accounting application.
- Products are managed in a product catalog and inventory system.

This distribution of data has made it difficult for users to answer key questions about the overall performance of the business.

In this exercise, you will examine some of the Adventure Works data sources that will be used in the data warehousing solution.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. View the Solution Architecture
3. View the Internet Sales Data Source
4. View the Reseller Sales Data Source

5. View the Products Data Source
6. View the Human Resources Data Source
7. View the Accounts Data Source
8. View the Staging Database

▶ **Task 1: Prepare the Lab Environment**

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run Setup.cmd in the **D:\Labfiles\Lab01\Starter** folder as Administrator. This script can take several minutes to complete.

▶ **Task 2: View the Solution Architecture**

1. View the **Adventure Works DW Solution.png** image in the **D:\Labfiles\Lab01\Starter** folder.
2. Note the data sources in the solution architecture.



Note: In addition to the data sources that you will examine in this lab, the diagram includes a Microsoft SQL Server Master Data Services model for product data and a SQL Server Data Quality Services task to cleanse data as it is staged. These elements form part of the complete solution for the lab scenario in this course, but they are not present in this lab.

▶ **Task 3: View the Internet Sales Data Source**

1. Use Microsoft SQL Server Management Studio (SSMS) to open the **View Internet Sales.sql** query file in the **D:\Labfiles\Lab01\Starter** folder. Use Windows® authentication to connect to the **MIA-SQL** instance of SQL Server.
2. Execute the query and examine the results. Note that this data source contains data about customers and the orders they have placed through the e-commerce web application.

▶ **Task 4: View the Reseller Sales Data Source**

1. Use SSMS to open the **View Reseller Sales.sql** query file.
2. Execute the query and examine the results. Note that this data source contains data about resellers and the orders they have placed through Adventure Works reseller account managers.

▶ **Task 5: View the Products Data Source**

1. Use SSMS to open the **View Products.sql** query.
2. Execute the query and examine the results. Note that this database contains data about products that Adventure Works sells, organized into categories and subcategories.

▶ **Task 6: View the Human Resources Data Source**

1. Use SSMS to open the **View Employees.sql** query file.
2. Execute the query and examine the results. Note that this database contains data about employees, including sales representatives.

► Task 7: View the Accounts Data Source

1. Examine the comma-delimited text files in the **D:\Accounts** folder by opening them using Notepad. Note that they contain details of payments made by resellers.
2. Close all files when you have finished reviewing them, without saving any changes.

► Task 8: View the Staging Database

- In SSMS, in the Object Explorer pane, examine the tables in the **Staging** database in the **MIA-SQL** instance of SQL Server. Ensure you examine the **Staging** database, not the **DQS_STAGING_DATA** database.

Note: All tables in this database are empty, except **dbo.ExtractLog**.

Results: After this exercise, you should have viewed data in the InternetSales, ResellerSales, Human Resources, and Products SQL Server databases; viewed payments data in comma-delimited files; and viewed an empty staging database.

Exercise 2: Exploring an ETL Process

Scenario

Now that you are familiar with the data sources in the Adventure Works data warehousing solution, you will examine the ETL process used to stage the data, and then load it into the data warehouse.

Adventure Works uses a solution based on SQL Server Integration Services to perform this ETL process.

The main tasks for this exercise are as follows:

1. View the Solution Architecture
2. Run the ETL Staging Process
3. View the Staged Data
4. Run the ETL Data Warehouse Load Process

► Task 1: View the Solution Architecture

1. Open the **Adventure Works DW Solution.png** file again.
2. Note that the ETL solution consists of two main phases: a process to **extract** data from the various data sources and load it into a **Staging** database, and another to **load** the data in the **Staging** database into the data warehouse. In this exercise, you will observe these ETL processes as they run.

► Task 2: Run the ETL Staging Process

1. Use Visual Studio® to open the **AdventureWorksETL.sln** solution file in the **D:\Labfiles\Lab01\Starter** folder.
2. In Solution Explorer, view the SSIS packages that this solution contains, and then double-click **Stage Data.dtsx** to open it in the designer.
3. View the control flow of the **Stage Data.dtsx** package, and then run the package by clicking **Start Debugging** on the **Debug** menu. The package will run other packages to perform the tasks in the control flow. This may take several minutes.

4. When the package has finished running, a message box is displayed, although this might be hidden by the Visual Studio window. Look for a new icon on the taskbar, and then click it to bring the message box to the front. After viewing this message box, stop the package by clicking **Stop Debugging** on the **Debug** menu.

► **Task 3: View the Staged Data**

1. Use SQL Server Management Studio to view the **Staging** database in the MIA-SQL instance of SQL Server. Take care to view the **Staging** database, not the **DQS_STAGING_DATA** database.
2. Note the following tables now contain data:
 - dbo.Customers
 - dbo.EmployeeInserts
 - dbo.InternetSales
 - dbo.Payments
 - dbo.Resellers
 - dbo.ResellerSales

► **Task 4: Run the ETL Data Warehouse Load Process**

1. In Visual Studio, in Solution Explorer, view the SSIS packages that the AdventureWorksETL.sln solution contains, and then double-click **Load DW.dtsx** to open it in the designer.
2. View the control flow of the Load DW.dtsx package, and then run the package by clicking **Start Debugging** on the **Debug** menu. The package will run other packages to perform the tasks in the control flow. This may take several minutes.
3. When the package has finished running, a message box will be displayed, although it might be hidden by the Visual Studio window. Look for a new icon on the taskbar, and then click it to bring the message box to the front. After viewing this message box, stop the package by clicking **Stop Debugging** on the **Debug** menu.
4. Close Visual Studio when you have finished.

Results: After this exercise, you should have viewed and run the SQL Server Integration Services packages that perform the ETL process for the Adventure Works data warehousing solution.

Exercise 3: Exploring a Data Warehouse

Scenario

Now that you have explored the ETL process that is used to populate the Adventure Works data warehouse, you can explore the data warehouse to see how business users can view key information.

The main tasks for this exercise are as follows:

1. View the Solution Architecture
2. Query the Data Warehouse

► Task 1: View the Solution Architecture

1. Open the **Adventure Works DW Solution.png** file again.
2. Note that the data warehouse provides a central data source for business reporting and analysis, then close Paint.

► Task 2: Query the Data Warehouse

1. Use SSMS to open the **Query DW.sql** query file.
2. Execute the query and examine the results. Note that the data warehouse contains the data necessary to view key metrics across multiple aspects of the business.

Results: After this exercise, you should have successfully retrieved business information from the data warehouse.

Module Review and Takeaways

In this module, you have learned about the key elements in a data warehousing solution.

Review Question(s)

Question: Why might you consider including a staging area in your ETL solution?

Check Your Knowledge

Question	
Which project role would be best assigned to a business user?	
Select the correct answer.	
<input type="checkbox"/>	Solution Architect
<input type="checkbox"/>	Data Modeler
<input type="checkbox"/>	Data Steward
<input type="checkbox"/>	Tester

Module 2

Planning Data Warehouse Infrastructure

Contents:

Module Overview	2-1
Lesson 1: Considerations for Data Warehouse Infrastructure	2-2
Lesson 2: Planning Data Warehouse Hardware	2-10
Lab: Planning Data Warehouse Infrastructure	2-19
Module Review and Takeaways	2-21

Module Overview

The server and hardware infrastructure for a business intelligence (BI) solution is a key consideration in any BI project. You must balance the performance and scalability gains that you can achieve by maximizing hardware specifications and distributing the elements of your BI solution across multiple servers against hardware and software licensing costs, and implementation complexity.

This module discusses considerations for selecting hardware and distributing SQL Server facilities across servers.

Objectives

After completing this module, you will be able to:

- Describe key considerations for data warehouse infrastructure.
- Plan data warehouse hardware.

Lesson 1

Considerations for Data Warehouse Infrastructure

Planning the infrastructure for a data warehousing solution that is based on Microsoft SQL Server requires an understanding of how the various SQL Server components work together, and how their typical workloads use hardware resources. This lesson considers all of the components in a BI solution, including the data warehouse, and servers for Extract, Transform and Load (ETL) workloads, SQL Server Analysis Services (SSAS), and SQL Server Reporting Services (SSRS).

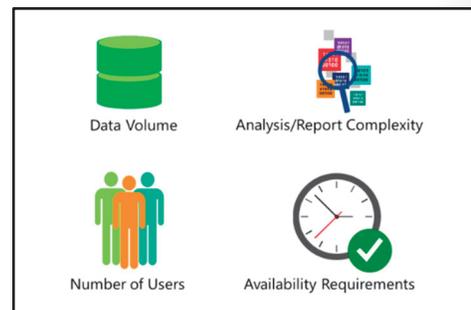
Lesson Objectives

After completing this lesson, you will be able to:

- Describe the considerations for characterizing the size of a BI solution.
- List the key features of the workloads that are associated with SQL Server components.
- Describe high level topology designs for BI solutions.
- Identify scale-out options for services in a BI solution.
- Describe the high availability options for services in a BI solution.

System Sizing Considerations

Each data warehouse solution has unique characteristics that depend on the business requirements that it is designed to address. There is no such thing as a “standard” data warehouse solution. However, when planning the infrastructure for a data warehouse, it is useful to start by classifying the type of solution to be implemented into one of three generic sizes—small, medium, or large. The factors that determine the size of a BI solution are:



- **Data Volume.** This is the amount of data that the data warehouse must store, and the size and frequency of incremental loads of new data. The primary consideration is the number of rows in fact tables, but you must also allow for dimension data, indexes, and data models that are stored on disk.
- **Analysis and Reporting Complexity.** This includes the number, complexity, and predictability of the queries that will be used to analyze the data or produce reports. Typically, BI solutions must support a mix of the following query types:
 - **Simple.** Relatively straightforward SELECT statements.
 - **Medium.** Repeatedly executed queries that include aggregations or many joins.
 - **Complex.** Unpredictable queries that have complex aggregations, joins, and calculations.
- **Number of Users.** This is the total number of information workers who will access the system, and how many of them will do so concurrently.

- **Availability Requirements.** These include when the system will need to be used, and what planned or unplanned downtime the business can tolerate.

Although it is difficult to be precise when categorizing a solution, the following table suggests some typical examples of the characteristics of small, medium, and large BI systems:

	Small	Medium	Large
Data Volume	Hundreds of GB to 1 TB	1 to 10 TB	10 TB to hundreds of TB
Analysis and Reporting Complexity	<ul style="list-style-type: none"> • More than 50% simple • 30% medium • Less than 10% complex 	<ul style="list-style-type: none"> • 50% simple • 30-35% medium • 10-15% complex 	<ul style="list-style-type: none"> • 0-5% simple • 40% medium • 20-25% complex
Number of Users	100 in total 10 to 20 concurrent	1,000 in total 100 to 200 concurrent	Thousands of concurrent users
Availability Requirements	Business hours	One hour of downtime per night	24/7 operations

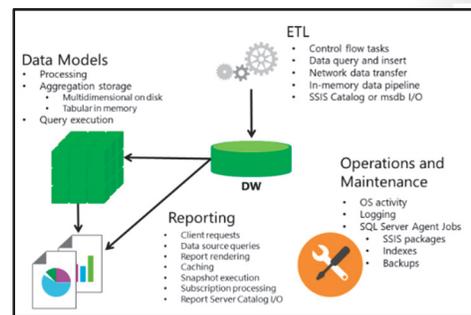
Data Warehouse Workloads

In addition to the size categorization of the solution that your infrastructure needs to support, it's useful to understand the workload types that might occur in a BI solution. It is important to assess the total impact of all workloads on hardware resource utilization and identify any potential for contention between workloads with similar requirements.

ETL Workloads

The ETL subsystem of the BI solution performs the following workloads:

- **Control Flow Tasks.** SQL Server Integration Services (SSIS) packages often include control flow tasks that require CPU processing time, memory, disk I/O, and network I/O. The specific resource requirements for control flow tasks can vary significantly, so if your ETL solution includes a substantial number of them, you should monitor resource utilization on a test system to better understand the workload profile.



MCT USE ONLY. STUDENT USE PROHIBITED

- **Data Query and Insert.** Fundamentally, ETL involves querying data sources, and inserting and updating data into staging and data warehouse tables. This incurs I/O and query processing on data sources, the staging databases, and the data warehouse, especially if data loads require rebuilding indexes or partition management.
- **Network Data Transfer.** Typically, ETL processes transfer large volumes of data from one server to another. This incurs network I/O and can require significant network bandwidth.
- **In-Memory Data Pipeline.** Data flow tasks in an SSIS package use an in-memory pipeline architecture to process transformations, and that places a demand on system memory resources. On systems where there is contention for memory resources between the SQL Server database engine and SSIS, data flow buffers might need to spill to disk, which reduces data flow performance and incurs disk I/O.
- **SSIS Catalog or msdb Database Activity.** SSIS packages that are deployed in project mode are stored in an SSIS Catalog database. Alternatively, packages that are deployed in package mode are stored either in the **msdb** database or on the file system. Whichever deployment model is used, the ETL process must access the package storage to load packages and their configuration. If the SSIS Catalog or **msdb** database is used, and it is located on a SQL Server instance that is hosting other databases used in the BI solution (such as the data warehouse, staging database, or Report Server Catalog), there might be contention for database server resources.

Data Model Workloads

The SSAS data models in a BI system require the following workloads:

- **Processing.** Data models contain aggregated values for the measures in the data warehouse. Values must be processed to load the required data from the data warehouse tables into the model and perform the necessary aggregation calculations. When data in the data warehouse is refreshed, the data models must be partially or fully processed again to reflect the new and updated data.
- **Aggregation Storage.** Data models must store the aggregated data in a structure that is optimized for analytical queries. Typically, multidimensional models are stored on disk and have some data cached in memory for performance reasons. Tabular models are usually stored completely in memory, so sufficient resources must be available.
- **Query Execution.** When users perform analytical queries against a data model, it must process the query and generate results. This requires CPU processing time, memory, and potentially disk I/O.

Reporting Workloads

Although reporting is often performed by client applications directly against the data warehouse or data models, many BI solutions include SSRS. An SSRS reporting solution typically involves the following workloads:

- **Handling Client Requests.** The report server must monitor for client requests that are submitted to Reporting Services over HTTP and process them.
- **Data Source Queries.** Reports are based on datasets that must be retrieved by querying data sources. Typically, data sources for the reports in a BI solution are SSAS data models or the data warehouse, so report execution incurs query processing overheads.
- **Report Rendering.** After the report data has been retrieved, SSRS must render it into the required format by using the Report Definition Language (RDL) and the specific rendering extension. Depending on the report's size, format, and complexity, rendering can incur substantial CPU and memory resources.

- **Caching.** To reduce query processing and rendering overheads, SSRS can cache datasets and reports in the report server temporary database. Datasets and reports can be cached on first use or you can use scheduled jobs to precache objects at a regular interval.
- **Snapshot Execution.** In addition to caching reports, you can create persisted report snapshots at a scheduled interval and store them in the report server database.
- **Subscription Processing.** You can configure SSRS to execute and deliver reports to file shares or email addresses on a scheduled basis.
- **Report Server Catalog I/O.** Report definitions and resources, such as images, are stored in the Report Server Catalog, and database I/O tasks are required to retrieve these when needed. Database I/O is required to retrieve cached reports and datasets from the temporary database, and to recover report snapshots from the catalog database. This database access may compete for resources with other databases that are hosted in the same SQL Server instance.

Operations and Maintenance

In addition to the fundamental BI activity, a BI solution must support ongoing system operations and maintenance activity, such as:

- Operating system activity.
- Logging activity.
- SQL Server Agent jobs, for example:
 - Execution of SSIS packages for ETL processes.
 - Index and statistics maintenance in databases, including the data warehouse.
 - Database backup tasks.

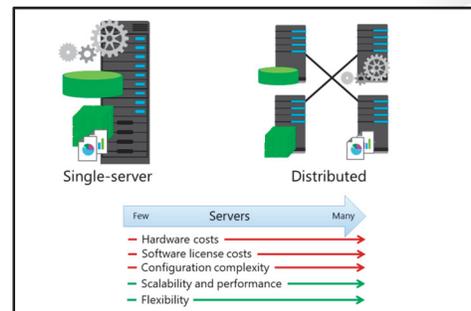
Typical Server Topologies for a BI Solution

SQL Server provides a versatile service architecture that enables you to choose a topology that best fits your needs. Selecting the right topology for your solution usually involves balancing cost and complexity against the need for high scalability and flexibility.

Single-Server BI Architecture

The simplest architecture for a BI solution that is based on SQL Server is to install all components on a single server. Depending on the business requirements of the BI solution, the components on the server typically include:

- **SQL Server Database Engine.** This component is used to store the data warehouse, staging database, Reporting Services databases, and SSIS Catalog database. In addition, the SQL Server Agent might be used to automate SSIS package execution and other operations by creating jobs, and schedules that are stored in the **msdb** database.
- **SQL Server Integration Services.** This component is used to execute packages that encapsulate ETL tasks and data flows to extract data from source systems into the staging database, and then load it into the data warehouse.



- **SQL Server Analysis Services.** This component is used to provide analytical data models and data mining functionality. Depending on business requirements, two instances of Analysis Services may be required—one for multidimensional models and data mining, and the other for tabular models.
- **SQL Server Reporting Services.** This component is used to provide report publishing and execution services. Reporting Services in native mode consists of a web service application, a web-based management user interface, and two SQL Server databases.

Depending on business requirements, you might also choose to install SQL Server Data Quality Services on the server to support data cleansing and matching capabilities for staged data before it is loaded into the data warehouse.



Note: If Microsoft SharePoint® Server is required, you can deploy the SharePoint farm and SQL Server integration components for Reporting Services and PowerPivot on the BI server. This architecture is not recommended for BI solutions that require significant scalability or performance.

A single-server architecture is suitable for test and development environments, and can be used in production locations that have minimal data volumes and scalability requirements.

Distributed BI Architecture

If your BI solution requires even moderate levels of scalability, it will benefit from expanding beyond a single-server architecture and distributing the BI workload across multiple servers. Typical approaches to distributing SQL Server components in a BI solution include:

- **Creating a Dedicated Report Server.** In many BI solutions, Analysis Services provides a data model that contains most (or even all) of the data in the data warehouse, and all reporting is performed against the data model. In these scenarios, there is little database activity in the data warehouse other than during ETL loading and data model processing (loading data from the data warehouse tables into the model and aggregating it). The server workloads that compete for resources most of the time are Analysis Services and Reporting Services, so you can increase scalability and performance by moving the reporting workload to a separate server. You can install the Reporting Services database on either server. Leaving it on the data warehouse server keeps all of the data engine elements on a single server, but can result in I/O workloads competing for disk resources. To install the Reporting Services database on the report server, the database engine must be installed on both servers, but the result means a cleaner separation of workloads. In extremely large enterprise solutions, which have intensive reporting requirements, you can add a third server as a dedicated host for the Reporting Services database.
- **Separating Data Warehouse and ETL Workloads from Analytical and Reporting Workloads.** If the data warehouse will be refreshed with new data frequently, or if it must support direct query access in addition to processing data models, the database I/O activity might compete with Analysis Services for disk, CPU, and memory resources. To prevent this, you can deploy Analysis Services on a separate server. Depending on analytical and reporting workloads, you might choose to co-locate Analysis Services and Reporting Services on the same server, or use a dedicated one for each. If you need to support tabular and multidimensional or data mining models, and a single Analysis Services server is not adequate to support both workloads, you could consider using separate servers for the different types of Analysis Services model.

- Using a Dedicated ETL Server.** If your ETL process requires frequent data extractions and loads, or involves particularly resource-intensive transformations, overall performance might benefit from moving Integration Services and the SSIS Catalog database to a dedicated server. Depending on the specific transformation and load operations that your ETL process requires, you can choose to locate the staging database on the ETL server or on the data warehouse server. Alternatively, you could use a two-phase staging approach where extracted data is staged on the ETL server for transformation and cleansing, and then loaded into staging tables on the data warehouse server before being inserted into the data warehouse tables.

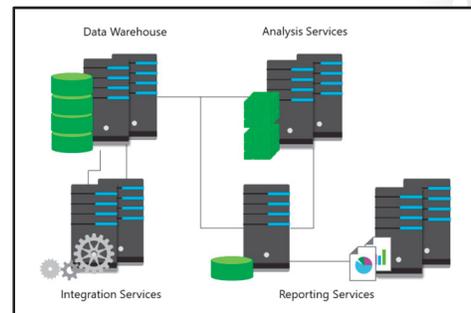
When designing a distributed architecture, the key goal is to eliminate contention for hardware resources. You then gain the greatest benefits by identifying workloads that have similar hardware utilization profiles and separating them.

Scaling Out a BI Solution

A distributed architecture increases scalability by separating the discrete workloads in a BI solution across multiple servers. To support the highest level of scalability, you can use a scale-out architecture to share the same workload across multiple servers.

You can use a scale-out architecture for the following components:

- The Data Warehouse.** You can scale out an extremely large data warehouse in several ways. Typically, this is done by partitioning the data across multiple database servers and using middle-tier logic to direct queries to the appropriate server instance. The Microsoft Analytics Platform System is available—it uses a massively parallel processing (MPP) integrated system that supports hybrid data warehouse scenarios to scale out data warehouse queries across multiple independent compute and storage nodes.
- SQL Server Analysis Services.** You can create a read-only copy of a multidimensional database and connect to it from multiple Analysis Services query servers. To accomplish this, an SSAS server processes the cubes in the database. The database is then detached, copied to a standby location, and reattached so that it can be used by applications that require write-back capabilities. The copied database is then attached in read-only mode to multiple SSAS instances, providing query services to clients. Client requests can be distributed across query servers by using a load-balancing middle tier, or you can assign specific subsets of the client population to dedicated query servers.
- SQL Server Integration Services.** Although it is not a pure scale-out technique, you can use multiple SSIS servers to perform a subset of the ETL processes in parallel. This approach requires extensive custom logic to ensure that all tasks are completed, and should be considered only if it is not possible to meet your ETL requirements through a scale-up approach in which you add hardware resources to a single SSIS server.
- SQL Server Reporting Services.** You can install the Reporting Services database on a single database server, and then install the Reporting Services report server service on multiple servers that all connect to the same Reporting Services database. This approach separates the report execution and rendering workloads from the database workloads that are required to store and retrieve report definitions, cached datasets and reports, and snapshots.





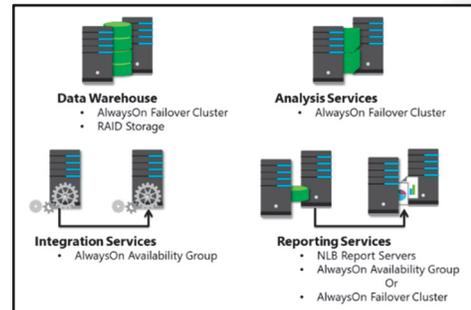
Microsoft Analytics Platform System

<http://aka.ms/yacqli>

Planning for High Availability

If one or more aspects of the BI solution require high availability, you can design a solution that uses redundant servers that have failover capabilities. SQL Server includes several high-availability technologies that you can use to protect critical services. These include:

- **AlwaysOn Failover Clustering.** This is server-level protection that uses Windows Server® cluster services to create a virtual server that has multiple redundant failover nodes.
- **AlwaysOn Availability Groups.** This is database-level protection that combines Windows Server cluster services with synchronization of database transactions across multiple SQL Server instances.
- **Log Shipping.** This is database-level protection that copies transaction logs from a primary server to a secondary server, where they can be applied to a secondary copy of the database.



Planning High Availability for the Data Warehouse

The data warehouse is fundamentally a SQL Server database and, theoretically, it can be protected by using any of the high-availability technologies that SQL Server provides. However, considering the large volume of data that is typical of data warehouses, and the fact that most activity consists of data reads or non-logged bulk load operations, AlwaysOn Failover Clustering is the most appropriate technology to consider.

Also, considering the importance of the data warehouse data, the database files should be stored on a redundant array of independent disks (RAID) that provides protection in the case of a disk failure.

Planning High Availability for Analysis Services

Analysis Services supports AlwaysOn Failover Clustering, so you should use this technology if you require high availability for an Analysis Services instance. If Analysis Services is to be installed on the same server as the database engine instance for the data warehouse, you can install both components in a single cluster operation.

Planning High Availability for Integration Services

Integration Services is not cluster-aware, and does not support failover from one cluster node to another. When SSIS packages are deployed in project mode, you can use an AlwaysOn availability group to protect the SSIS Catalog database. However, there are some additional considerations for using the AlwaysOn Availability Groups feature with the SSIS catalog:

- **Re-Encrypt Master Database Key.** On failover, you must re-encrypt the master database key in the new primary server before any packages that include encrypted sensitive data can be executed.

- **Recover ETL Process.** Your SSIS packages must be able to recover the ETL process as a whole to a consistent state if unplanned failover occurs. This means that you must include cleanup and failover logic in your packages.
- **SSIS Catalog and the Availability Group.** If you need to apply an update that modifies the SSIS catalog schema, you should remove the SSIS Catalog database from the availability group, update it, and then reestablish the availability group.



AlwaysOn for SSIS Catalog (SSISDB)

<http://aka.ms/vpndwp>

Planning High Availability for Reporting Services

As described earlier in this lesson, you can distribute Reporting Services across a multitier architecture in which a report server processes requests and executes reports. A separate database tier hosts the Report Server Catalog and temporary database. In this configuration, you can use Network Load Balancing (NLB) to distribute requests across report servers, and remove failed report servers from the NLB cluster to maintain a highly available solution for managing report requests. To protect the database tier, you can use either AlwaysOn Failover Clustering or an AlwaysOn availability group.

Check Your Knowledge

Question	
Which of these options is not a consideration for system sizing?	
Select the correct answer.	
<input type="checkbox"/>	Data volume
<input type="checkbox"/>	Number of users
<input type="checkbox"/>	Data center location
<input type="checkbox"/>	Analysis and reporting complexity
<input type="checkbox"/>	Availability requirements

Question: Which four factors determine the size of a BI solution?

Question: Discuss the considerations for deciding between single-server and distributed architecture.

Lesson 2

Planning Data Warehouse Hardware

The data warehouse is the foundation for a BI solution. You should consider several recommendations from Microsoft and its hardware partners when planning a data warehouse system. Data warehousing is substantially different from other database workloads, and the conventional database design approach for optimizing hardware to support the highest possible number of I/O operations per second (IOPS) is not always appropriate.

This lesson introduces Microsoft SQL Server Fast Track Data Warehouse reference architectures, and explains how some of the Fast Track design principles can be applied when planning data warehouse hardware.

Lesson Objectives

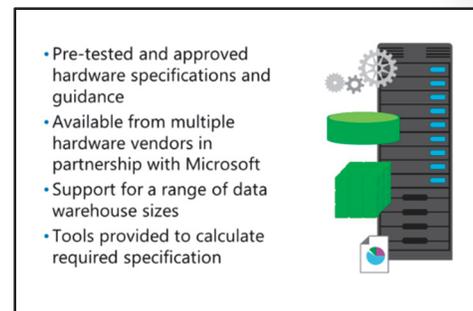
After completing this lesson, you will be able to:

- Describe the key features of SQL Server Fast Track Data Warehouse reference architectures.
- Explain how a core-balanced system architecture supports data warehousing workloads.
- Determine processor and memory requirements for a data warehouse.
- Determine storage requirements for a data warehouse.
- Describe considerations for choosing a storage solution.
- Describe and understand SQL Server data warehouse appliances.
- Understand SQL Server Parallel Data Warehouse.

SQL Server Fast Track Data Warehouse Reference Architectures

Data warehouse design is a specialist skill that many organizations may not have in-house, so the BI team must work with consultants and hardware vendors to determine the required hardware and its configuration. This adds cost and time to the project, and does not always guarantee that the resulting infrastructure is appropriate for the data warehousing workload.

To help organizations overcome these challenges, Microsoft has partnered with multiple hardware vendors to create Fast Track Data Warehouse reference architectures that reduce the time and effort that it takes to specify a data warehouse system. These reference architectures, together with hardware appliances that can be purchased as out-of-the-box solutions, make it easier and quicker to implement an effective data warehousing solution.



Pretested Specifications and Guidance

Fast Track Data Warehouse reference architectures are designed by SQL Server specialists and consultants from hardware partners. The reference architectures include specific hardware components that can be combined in a prescribed configuration to create a data warehouse system that is based on SQL Server. The architectures have been extensively tested with real-world data warehouse workloads that are based on thousands of customer case studies. They are certified to deliver specific performance levels based on the size of your data warehouse solution.

Multivendor Support

Microsoft has partnered with multiple vendors to create pretested system designs that use commodity hardware. If your organization has an existing supplier relationship with one of the Fast Track hardware vendors, you can easily specify a system that is based on components from that supplier. You can use the support and consulting services that the hardware vendor offers to create a data warehouse system that is based on a proven design.

Support for a Range of Data Warehouse Sizes

The Fast Track program is not a “one size fits all” solution. It includes multiple reference architectures for small to large data warehouse solutions, and there is flexibility within each design to support your specific requirements. You can evaluate the size of your required solution by using the principles that were described in the previous lesson, and then identify the reference architecture that best suits your needs.

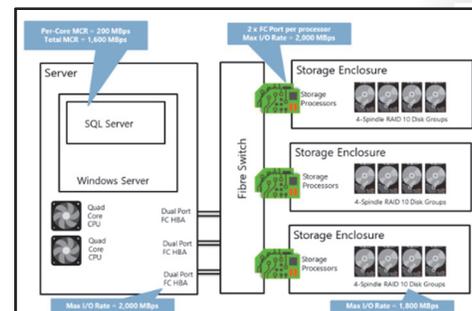
Core-Balanced System Architecture

SQL Server Fast Track Data Warehouse reference architectures are based on a core-balanced system design that reflects the performance characteristics of a typical data warehouse workload. Even if you don't plan on using one of the published Fast Track Data Warehouse reference architectures, you will benefit from understanding the principles on which they are designed and applying them to your own system specification.

The core-balanced system architecture is based on the fact that most data warehouse workloads need to transfer large amounts of data (usually accessed by sequential read operations) across multiple system components, from where the data is stored to the requesting applications. Each component through which the data is transferred is a potential bottleneck that will limit the overall performance of the system. The data can only flow to the requesting application at the rate of the slowest component. Any components that can operate at a higher rate are underutilized, which unbalances the system and can represent significant wasted cost.

The diagram on the slide shows a balanced system in which the I/O rates of each component are similar. This diagram represents a large-scale data warehousing system in which data is kept in a storage area network that has fiber channel connectivity and multiple storage enclosures, each containing multiple disk arrays. The same principles apply to a smaller architecture.

You can establish the I/O rate of hardware components such as hard disks, array storage processors, and fiber channel host bus adapters (HBAs), through careful testing and monitoring by using tools like SQLIO (see the reference that follows). Many manufacturers, particularly those who participate in the Fast Track program, publish the maximum rates. However, the initial figure that you need to start designing a data



warehouse system is the maximum consumption rate (MCR) of a single processor core, combined with the SQL Server database engine. After the MCR for the CPU core architecture that you intend to use has been established, you can determine the number of processors that are required to support your workload and the storage architecture that is needed to balance the system.

 **Note:** MCR is specific to a combination of a CPU and motherboard, and SQL Server. It is not a measure of pure processing speed or an indication of the performance that you can expect for all solution queries. Instead, MCR is a system-specific benchmark measure of maximum throughput per core for a data warehouse query workload. Calculating MCR requires you to execute a query that can be satisfied from cache while limiting execution to a single core—and reviewing the execution statistics to determine the number of megabytes of data processed per second.

 **Using SQLIO blog**

<http://aka.ms/h15oc8>

Demonstration: Calculating Maximum Consumption Rate (MCR)

In this demonstration, you will see how to:

- Create tables for benchmark queries.
- Execute a query to retrieve I/O statistics.
- Calculate MCR from the I/O statistics.

Demonstration Steps

Create Tables for Benchmark Queries

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the `D:\Demofiles\Mod02` folder, run **Setup.cmd** as Administrator.
3. In the User Account Control window, click **Yes**.
4. Start SQL Server Management Studio, and then connect to the **MIA-SQL** database engine by using Windows authentication.
5. In SQL Server Management Studio, open the **Create BenchmarkDB.sql** query file in the **D:\Demofiles\Mod02** folder.
6. Click **Execute**, and then wait for query execution to complete. This query creates a database that contains two tables, one with a clustered index and one without. Both tables contain a substantial number of rows.

Execute a Query to Retrieve I/O Statistics

1. In SQL Server Management Studio, open the **Measure MCR.sql** query file in the **D:\Demofiles\Mod02** folder.
2. Click **Execute**, and then wait for query execution to complete. The queries retrieve an aggregated value from each table, and are performed twice. This ensures that on the second execution (for which statistics are shown), the data is in cache, so the I/O statistics do not include disk reads. Note that the `MAXDOP 1` clause ensures that only a single core is used to process the query.

Calculate MCR from the I/O Statistics

1. In the results pane, click the **Messages** tab. This shows the statistics for the queries.
2. Add together the **logical reads** value for the two queries, and then divide the result by two to find the average.
3. Add the **CPU time** value for the two queries together, and then divide the result by two to find the average. Divide the result by 1000 to convert it to seconds.
4. Calculate MCR by using the following formula:

$$(\text{average logical reads} / \text{average CPU time}) * 8 / 1024$$

Determining Processor and Memory Requirements

After determining the MCR for the CPU cores that you intend to use, you can start to estimate the number of cores that will be required to support your anticipated query workload. You can also make an initial assessment of memory requirements.

Estimating CPU Requirements

MCR indicates the amount of data that can be processed by a single core in one second. To determine the number of cores that are required, you must apply this rate to the following factors:

- The amount of data that is returned by an average query.
- The number of concurrent users whom you need to support.
- The target response time for a query.

The specific formula to apply is:

$$((\text{Average query size in MB} \div \text{MCR}) \times \text{Concurrent users}) \div \text{Target response time}$$

For example, suppose the MCR of the CPU core that you intend to use is 200 megabytes per second (MBps). If an average query is expected to return 18,000 MB, the anticipated number of concurrent users is 10, and each query must respond within 60 seconds, the calculation to find the number of cores that are required is:

$$((18000 \div 200) \times 10) \div 60 = 15$$

CPU architecture does not include exactly 15 cores, so you round up to a requirement of 16.

Now that you know the number of cores that are required, you can make an initial determination of the number of processors.

For example, to provide 16 cores using quad-core processors, you would need four processors.

Alternatively, if dual-core processors are used, eight CPUs would be required. Remember that you need to balance the number of CPUs to closely match the number of storage arrays that will be used, which in turn may depend on the volume of data that your data warehouse must support.

Estimating CPU Requirements: 

- Determine core MCR
- Apply formula to estimate required number of cores:
 $((\text{Average query size in MB} + \text{MCR}) \times \text{Concurrent users}) \div \text{Target response time}$
- Spread cores across CPUs based on the number of storage arrays

Estimating RAM Requirements: 

- Use a minimum of 4 GB per core (or 64–128 GB per socket)
- Target 20% of data volume

Estimating RAM Requirements

It is difficult to calculate the amount of RAM that is required because many workloads can utilize memory to increase overall performance. You should generally consider a minimum figure for a small to medium-sized data warehouse system to be 4 GB per core, or 64 to 128 GB per CPU socket. If you intend to use columnstore indexes or support tabular data models on the data warehouse server, you should favor the higher end of these estimates.

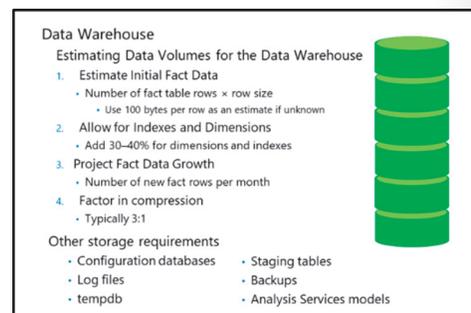
Another way to estimate memory requirements is to consider that, in an average data warehouse workload, users regularly need to access approximately 20 percent of the data that is stored. For example, in a data warehouse that stores five years of sales records, users most frequently query the data for the most recent year. Having enough memory to maintain approximately 20 percent of the data in cache will significantly enhance performance.

Determining Storage Requirements

Before you can fully determine CPU, memory, and storage hardware requirements, you must assess the volume of data that the system must support.

Estimating Data Volumes for the Data Warehouse

Most data warehouses consist predominantly of fact data. Determining the volume of fact data that the data warehouse must store is the most significant factor in assessing overall storage requirements.



1. Estimate Initial Fact Data

To start estimating data warehouse data volumes, determine the number of fact rows that will be initially loaded into the data warehouse and multiply that by the average size of a fact row. If you don't know the average fact row size at this stage, use a conservative estimate such as 100 bytes per row. For example, a data warehouse that will contain 200,000,000 fact rows, each 100 bytes in length, will have an initial fact data volume of approximately 20 GB.

2. Allow for Indexes and Dimensions

After estimating the initial fact data, add approximately 30 to 40 percent to allow for indexes and dimensions. To continue the example with 20 GB of fact data, you would add approximately 8 GB (40 percent of 20 GB), giving an initial data volume of approximately 28 GB.

3. Project Fact Data Growth

Most data warehouses are regularly refreshed with new data. To be sure that your storage solution will support the data warehouse in the future (say, three years from now), you must factor in the anticipated incremental data that will be loaded. For example, suppose the fact data in our data warehouse represents individual items that have been ordered in sales transactions, and the company typically sells 5,000,000 items a month—you can expect to load 5,000,000 rows (each containing 100 bytes of data), or approximately 500 MB each month. That equates to a data growth rate of 6 GB per year, so in three years, the example data warehouse would need to support the initial 28 GB of data plus another 18 GB (6 GB per year multiplied by three years), giving a total of 46 GB.

4. Factor in Compression

Finally, you should plan to compress the data in your data warehouse. Typically, SQL Server provides a compression factor of approximately 3:1, so the 46 GB of data should compress to approximately 15.5 GB on disk.

Other Storage Requirements

In addition to the data warehouse, you must include other data in your storage estimation. Additional storage is required for:

- **Configuration Databases.** If databases that are used by other BI services, including the SSIS Catalog and Reporting Services databases, are to be installed on the data warehouse server, you should include them in your storage estimate. In addition, the SQL Server instance includes system databases, although in practice, these are usually stored separately from the data warehouse data files.
- **Transaction Log Files.** Each database requires a transaction log. Typically, data warehouses are configured to use the simple recovery model and few transactions are actually logged.
- **tempdb.** Many data warehouse queries require temporary storage space. It is generally recommended to locate **tempdb** on a suitable storage column and assign an appropriate initial size to avoid the system having it grow automatically as needed.
- **Staging Tables.** Whether data is staged in a dedicated staging database, in tables within the data warehouse database itself, or in a combination of both, you must allocate enough space to allow for data staging during ETL processes.
- **Backups.** If you intend to back up the data warehouse and other databases to disk, you must ensure that the storage design provides space for backup files.
- **Analysis Services Models.** If you intend to host multidimensional Analysis Services data models on the data warehouse server, you must allocate sufficient disk space for them.

Considerations for Storage Hardware

The optimal storage hardware solution for a data warehouse depends on several factors, including the volume of data and the system MCR that data throughput from the storage system must support. When planning a storage solution, consider the following guidelines:

- **Disk Size.** Use more, smaller disks instead of fewer, larger disks. Although you can create a data warehouse that stores all of its data on a single, large hard disk, it is usually possible to achieve a better balance of throughput (and therefore overall system performance) by distributing the data across multiple small disks. This enables multiple disk reads to be performed in parallel and reduces wait times for I/O operations.
- **Disk Speed.** Use the fastest disks that you can afford. Disk technologies have advanced dramatically in recent years, with the speed of mechanical disks increasing and the advent of solid state disks with no moving parts. Most data warehouse read operations are sequential scans instead of the random I/O patterns of online transaction processing (OLTP) systems, so seek times are minimized. Regardless of this advantage, however, a faster disk means greater throughput when reading data. Solid state

- Use more smaller disks instead of fewer larger disks
- Use the fastest disks you can afford
 - Consider solid state disks—especially for random I/O
- Use RAID 10, or minimally RAID 5
- Consider a dedicated storage area network for manageability and extensibility
 - Balance I/O across enclosures, storage processors, and disk groups



disks are typically more expensive than mechanical disks, but if disk performance is critical, you might decide that the additional cost is worth paying. The lack of moving parts makes them particularly effective for random I/O data access, which is typical of queries against multidimensional data models.

- **RAID.** Use RAID 10, or minimally RAID 5. RAID 10 (in which data is both mirrored and striped) provides the best balance of read performance and protection from disk failure, and this should usually be the first choice for a data warehouse. However, the requirement for a complete set of redundant disks per array can make this an expensive option. As an alternative, you can use RAID 5, which provides striping for high read performance and parity-based data redundancy to protect against disk failure.
- **DAS or SAN?** Consider a dedicated storage area network (SAN). Although you can build a data warehouse that has direct-attached storage (DAS), using a SAN generally makes it easier to manage disk array configuration and to add future storage as the data warehouse grows. If you decide to use a SAN, it is best to have one that is dedicated to the BI solution and not shared with other business applications. In addition, try to balance the number of enclosures, storage processors per enclosure, and disk groups to achieve a consistent I/O rate that takes advantage of parallel core processing and matches the MCR of the system.

SQL Server Data Warehouse Appliances

Although Fast Track Data Warehouse reference architectures can reduce the time and effort that is taken to implement a data warehouse, organizations still require technical expertise to assemble the solution. To reduce the technical burden on organizations, and reduce the time that it takes to implement a solution, Microsoft has partnered with hardware vendors to create preconfigured data warehouse appliances that you can procure with a single purchase.

- Pre-built hardware and software solutions, based on tested configurations
- Part of a range of appliances that are based on SQL Server
- Available from multiple hardware vendors

The data warehouse appliances that are available from Microsoft and its hardware partners are based on tested configurations, including Fast Track reference architectures, and can significantly reduce the time that it takes to design, install, and optimize a data warehouse system.

Data warehouse appliances that are based on Fast Track Data Warehouse reference architectures are available for organizations or departments that need to deploy a solution quickly and with minimal installation and configuration effort. In addition, large organizations that require an enterprise data warehouse can purchase an appliance that is based on Microsoft SQL Server Parallel Data Warehouse for extreme scalability and performance.

Data warehouse appliances form part of a range of appliances, based on SQL Server, that Microsoft and its hardware partners have developed for common database workloads. Other types include business decision appliances that provide self-service BI capabilities, and database server consolidation appliances that use virtualization technologies to create a private cloud infrastructure for database servers. Appliances are available from multiple hardware vendors, and include technical support for the entire appliance, including software and hardware.

SQL Server Parallel Data Warehouse

Fast Track Data Warehouse systems, and appliances that are based on them, use a symmetric multiprocessing (SMP) architecture. For SMP, the system bus is the limiting component that prevents scaling up beyond a certain level. As the number of processors and the data load increases, the bus can become overloaded and a bottleneck occurs. For data warehouses that require greater scalability than an SMP system can provide, you can use an enterprise data warehouse appliance that is based on SQL Server Parallel Data Warehouse.

- A special SQL Server edition only available in hardware appliances
- Shared-nothing architecture
- Massively parallel processing
- Dedicated control nodes, compute nodes, and storage nodes

SQL Server Parallel Data Warehouse is an edition of SQL Server that is only available as a preinstalled and configured solution in enterprise data warehouse appliances from Microsoft and its hardware partners. Parallel Data Warehouse is designed specifically for extremely large-scale data warehouses that need to store and query hundreds of terabytes of data.

Microsoft Analytics Platform System

The Microsoft Analytics Platform System (APS) is designed to meet the requirements of data warehouse growth. APS is a network-attached storage (NAS) massively parallel processing integrated system that is suitable for scaling out. APS supports a variety of hybrid data warehouse scenarios, and uses PolyBase and other technologies for managing big data.

Microsoft Analytics Platform System

<http://aka.ms/yacqli>

Shared-Nothing Architecture

Systems that use shared components, such as memory or disk storage, can suffer from performance issues because of contention for those shared components. Contention occurs when multiple nodes attempt to access a component at the same time. It usually results in degraded performance as nodes queue to access resources. Shared-nothing architectures eliminate contention because each node has its own dedicated set of hardware, which the others do not use. Removing contention from a system results in improved performance, and enables the system to handle larger workloads.

Massively Parallel Processing

Parallel Data Warehouse uses a shared-nothing, massively parallel processing architecture, which delivers improved scalability and performance over SMP systems. MPP systems deliver much better performance than SMP servers for large data loads. MPP systems use multiple servers, called nodes, which process queries independently in parallel. Parallel processing involves distributing queries across the nodes so that each processes only a part of the query. The results of the partial queries are combined after processing completes to create a single result set.

Control Nodes, Compute Nodes, and Storage Nodes

A Parallel Data Warehouse appliance consists of a server acting as the control node, and multiple servers representing compute and storage nodes. Each compute node has its own dedicated processors and memory, and is associated with a dedicated storage node. A dual InfiniBand network connects the nodes, and dual fiber channels link the compute nodes to the storage nodes. The control node intercepts incoming queries and divides each query into multiple smaller operations, which it passes on to the

compute nodes to process. Each compute node returns the results of its processing to the control node. The control node integrates the data to create a result set, which it then returns to the client.

Control nodes are housed in a control rack. Several other types of node share this rack:

- **Management Nodes.** These are the nodes that administrators use to manage the appliance.
- **Landing Zone Nodes.** These nodes act as staging areas for data that you load into the data warehouse by using an ETL tool.
- **Backup Nodes.** These nodes are used to back up the data warehouse.
- **Compute Nodes and Storage Nodes.** These nodes are housed separately in a data rack. To scale the application, you can add more racks as required. Hardware components, including control and compute nodes, are duplicated to provide redundancy.

You can use a Parallel Data Warehouse appliance as the hub in a hub-and-spoke configuration, and populate data marts directly from the data warehouse. Using a hub-and-spoke configuration enables you to integrate the appliance with existing data marts or to create local data marts. If you use Fast Track Data Warehouse systems to build the data marts, you can achieve very fast transfers of data between the hub and the spokes.



Starting point for Parallel Data Warehouse Research

<http://aka.ms/Jtj6h2>

Check Your Knowledge

Question	
Which of these equations defines maximum consumption rate?	
Select the correct answer.	
<input type="checkbox"/>	$(\text{average logical writes} \div \text{average CPU time}) \times 4 \div 1024$
<input type="checkbox"/>	$(\text{average logical reads} \div \text{maximum CPU time}) \times 12 \div 1024$
<input type="checkbox"/>	$(\text{average physical reads} \div \text{minimum CPU time}) \times 8 \div 1024$
<input type="checkbox"/>	$(\text{average logical reads} \div \text{average CPU time}) \times 8 \div 1024$
<input type="checkbox"/>	$(\text{average physical reads} \div \text{average CPU time}) \times 8 \div 1024$

Question: What factors are used in estimating CPU requirements?

Question: What factors are involved when planning a storage solution?

Lab: Planning Data Warehouse Infrastructure

Scenario

You are planning a data warehouse solution for Adventure Works Cycles, and have been asked to specify the hardware that is required. You must design a solution that is based on SQL Server that provides the right balance of functionality, performance, and cost.

Objectives

After completing this lab, you will be able to:

- Plan data warehouse hardware for a BI solution that is based on SQL Server.

Estimated Time: 30 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Planning Data Warehouse Hardware

Scenario

Now that you have planned the server infrastructure, you must create a hardware specification for the data warehouse server. You will begin by calculating the MCR of the system that you are currently using, and then complete a planning worksheet for a new system with a published MCR figure.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Measure Maximum Consumption Rate
3. Estimate Server Hardware Requirements

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC**, **20767C-MIA-SQL**, and **20767C-MIA-CLI** virtual machines are both running, and then log on to the **20767C-MIA-SQL** virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab02\Starter** folder as Administrator.

► Task 2: Measure Maximum Consumption Rate

1. Use the **Create BenchmarkDB.sql** script in the **D:\Labfiles\Lab02\Starter** folder to create a benchmark database on the MIA-SQL instance of SQL Server.
2. Use the **Measure MCR.sql** script in the **D:\Labfiles\Lab02\Starter** folder to generate query performance statistics.
3. Use the statistics to calculate MCR for the database server.
4. Use the calculated MCR figure to estimate the number of cores that are required to support the following workload:
 - Average data per query: 500 MB
 - Concurrent users: 10
 - Target response time: 20 seconds

► Task 3: Estimate Server Hardware Requirements

1. Log on to the **20767C-MIA-CLI** virtual machine as
2. Use Microsoft Excel to open the **DWHardwareSpec.xlsx** workbook in the **D:\Labfiles\Lab02\Starter** folder.
3. Use the information in the workbook to:
 - Calculate the number of cores that are required in the data warehouse server.
 - Recommend the number and type (dual-core or quad-core) of processors to include in the data warehouse server.
 - Calculate the estimated volume of data in the data warehouse.
 - Suggest a suitable amount of memory for the data warehouse server.
 - Calculate the storage requirements for the data warehouse server, assuming a compression ratio of 3:1.
4. Use Microsoft Word to open the **Storage options.docx** document in the **D:\Labfiles\Lab02\Starter** folder, and then review the available storage options.
5. Based on the storage requirements that you have identified, select a suitable storage option, and then record your selection in the **DWHardwareSpec.xlsx** workbook.
6. Save your changes to **DWHardwareSpec.xlsx**, and then close Excel and Word.

Results: After this exercise, you should have a completed worksheet that specifies the required hardware for your data warehouse server.

Question: Review **DWHardwareSpec.xlsx** in the D:\Labfiles\Lab02\Solution folder. How does the hardware specification in this workbook compare to the one that you created in the lab?

Module Review and Takeaways

This module has described some key considerations for planning the hardware infrastructure for a BI solution that is based on SQL Server. You should use this as a starting point, and use the “more information” references to learn more about the supported, distributed, and scale-out architectures for SQL Server components, and about the design principles that are used in Fast Track Data Warehouse reference architectures.

Review Question(s)

Question: In a growing number of organizations, virtualization has become a core platform for infrastructure. Microsoft Hyper-V®, which is included as a feature from Windows 10, together with enterprise operations and management software such as Microsoft System Center, have enabled IT departments to benefit from simpler provisioning, management, mobility, and recoverability of services.

What components of a BI infrastructure would you consider virtualizing, and why?

MCT USE ONLY. STUDENT USE PROHIBITED

Module 3

Designing and Implementing a Data Warehouse

Contents:

Module Overview	3-1
Lesson 1: Data Warehouse Design Overview	3-2
Lesson 2: Designing Dimension Tables	3-8
Lesson 3: Designing Fact Tables	3-16
Lesson 4: Physical Design for a Data Warehouse	3-19
Lab: Implementing a Data Warehouse	3-34
Module Review and Takeaways	3-40

Module Overview

The data warehouse is at the heart of most business intelligence (BI) solutions. Designing the logical and physical implementation of the data warehouse is crucial to the success of the BI project. Although a data warehouse is fundamentally a database, there are some significant differences between the design process and best practices for an online transaction processing (OLTP) database and those for a data warehouse supporting online analytical processing (OLAP) and reporting workloads.

This module describes key considerations for the logical design of a data warehouse, and then discusses best practices for its physical implementation.

Objectives

After completing this module, you will be able to:

- Describe a process for designing a dimensional model for a data warehouse.
- Design dimension tables for a data warehouse.
- Design fact tables for a data warehouse.
- Design and implement effective physical data structures for a data warehouse.

Lesson 1

Data Warehouse Design Overview

Before designing individual database tables and relationships, it is important to understand the key concepts and design principles for a data warehouse. This lesson describes the dimensional model used in most data warehouse designs and the process used to translate business requirements into a data warehouse schema.

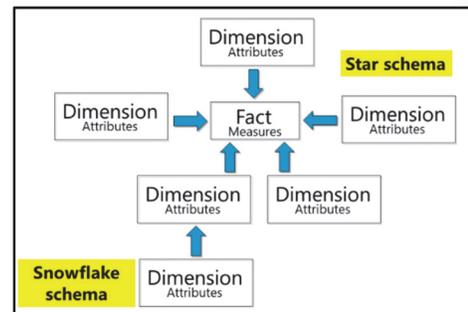
Lesson Objectives

After completing this lesson, you will be able to:

- Describe the dimensional model used for most data warehouses.
- Apply an effective process for data warehouse design.
- Use a business process-based approach to dimensional modeling.
- Document dimensional model designs.

The Dimensional Model

Although data warehouses can be implemented using normalized, relational database schemas, most data warehouse designs are based on the dimensional model advocated by Ralph Kimball. In the dimensional model, the numeric business measures that are analyzed and reported are stored in fact tables that are related to multiple dimension tables—where the attributes by which the measures can be aggregated are stored. For example, a fact table might store sales order measures, such as revenue and profit, and be related to dimension tables representing business entities such as product and customer. These relationships make it possible to aggregate the sales order measures by the attributes of a product (for example, to find the total profit for a particular product model) or a customer (for example, to find the total sales revenue for customers who live in a particular country).



Ideally, a dimensional model can be implemented in a database as a “star” schema, in which each fact table is directly related to its relevant dimension tables. However, in some cases, one or more dimensions may be normalized into a collection of related tables to form a “snowflake” schema. Generally, you should avoid creating snowflake dimensions because, in a typical data warehouse workload, the performance benefits of a single join between fact and dimension tables outweigh the data redundancy reduction benefits of normalizing the dimension data.

The query optimizer in the Enterprise edition of Microsoft® SQL Server® includes logic that detects star schema joins in queries and optimizes the way these joins are processed accordingly. Based on the selectivity of the query (that is, the proportion of rows from the fact table that the query is likely to return), the query optimizer uses bitmap filters to quickly eliminate nonqualifying rows from the fact table (which generally accounts for the largest cost in a data warehouse query).

SQL Server has an additional feature, Data Query Store, which, when enabled, uses the history of data queries previously run on the database, to improve efficiency when running queries.

For more detailed information about star join query optimization, see *Introduction to New Data Warehouse Scalability Features in SQL Server 2008* and *Data Warehouse Query Performance*:

 **Introduction to New Data Warehouse Scalability Features in SQL Server 2008**

<http://aka.ms/c9n1tv>

 **Data Warehouse Query Performance**

<http://aka.ms/asyxc1>

The Data Warehouse Design Process

Although every project has its unique considerations, there is a commonly-used process for designing a dimensional data warehouse that many BI professionals have found effective. The method is largely based on the data warehouse design patterns identified and documented by Ralph Kimball and the Kimball Group, though some BI professionals may adopt a varied approach to each task.

1. Determine analytical and reporting requirements
2. Identify the business processes that generate the required data
3. Examine the source data for those business processes
4. Conform dimensions across business processes
5. Prioritize processes and create a dimensional model for each
6. Document and refine the models to determine the database logical schema
7. Design the physical data structures for the database

 **Additional Reading:** For a detailed exploration of how to apply the Kimball dimensional modeling methodology to a SQL Server-based data warehouse design, read *The Microsoft Data Warehouse Toolkit* (Wiley, 2011).

Determine analytical and reporting requirements

After gathering the business requirements for the BI solution, you should interpret them in relation to the analytical and reporting capabilities that the solution is required to provide. Typically, analytical and reporting requirements support business requirements, so you will probably need to spend time with the stakeholders to understand the information that they need and to discuss how to achieve the best results. For example, a sales executive might express a business requirement as: "We want to improve the performance of sales representatives in the most poorly performing sales territories." To meet this requirement, you need to understand how "sales performance" is measured (for example, revenue, profitability, number of orders, or a combination of all three) and against what aspects of the business it should be aggregated.

Typically, asking questions such as: "How will you be able to tell whether the business requirement is being met?" leads the discussion toward the analytical and reporting requirements. For example, to determine if sales performance is improving, the sales executive might need to be able to see "order volume by territory" or "sales revenue by salesperson." Requirements expressed like this make it easier to determine the measures and dimensions that the solution must include, because the requirement often takes the form "measure by dimension".

Additionally, most analytical and reporting requirements include a time-based aggregation. For example, the sales executive might want to compare sales revenue by month or quarter.

Identify the business processes that generate the required data

Typically, the data required in the dimensional model is generated by an existing business process. After determining the data that you need to support analysis and reports, you should identify the business processes that generate the source data.

For example, a business might include the following processes:

- Order processing
- Stock management
- Order fulfillment
- Manufacturing
- Marketing and promotions
- Financial accounting
- Human resources management

Each process generates data that includes numeric values and events that can be counted (these can be sources for measures in a dimensional model) and information about key business entities (these can be sources for dimension attributes).

Examine the source data for those business processes

In most organizations, each business process captures data in at least one system. For example, order processing might store details of orders, customers, and products. A financial accounting process typically stores details of accounts and balances.

A significant part of the data warehouse solution design process involves exploring the data in these source systems and interviewing the users, system administrators, and application developers who understand it best. Initial exploration might simply involve running Transact-SQL queries to determine distinct value counts, average numerical values, and row counts. You can use the basic information gathered from these initial queries and discussions with data specialists as a foundation for deeper data profiling using tools such as the Data Profiling task in SQL Server Integration Services. This can determine minimum and maximum field lengths, data sparseness and null counts, and the reliability of relational dependencies.

At this stage, you do not need to perform a full data audit and start planning the extract, transform, and load (ETL) solution. You do, however, need to identify if and where the measures and dimension attributes you need to meet the reporting requirements are stored; what range of values exist for each required data field; what data is missing or unknown; and at what granularity the data is available.

Conform dimensions across business processes

Identifying the business processes and exploring the data generated by each one will help you to identify some key business entities that are common across the various processes. For example, the manufacturing and stock management business processes might both deal with a “product” entity. Similarly, the financial accounting, order processing and order fulfillment processes might all deal with a “customer” entity. Identifying dimensions that can be shared across multiple business processes is an important part of data warehouse design because it means you can define “conformed” dimensions. These ensure that the same definition of a business entity can be used to aggregate multiple facts and produce meaningful comparisons. For example, you can use a conformed product dimension—a data warehouse based on the order processing, manufacturing, and stock management business processes—to analyze a specific product and compare the number ordered, the number manufactured, and the number held in stock. If the product is perishable and has a limited shelf life, this kind of comparison could provide significant information for production planning and help reduce losses from spoiled, unsold products.

Prioritize business processes and define a dimensional model for each

Based on the business value of the identified reporting and analytical requirements, you can prioritize the business processes and create a dimensional model for each one that is required to provide the necessary analytical and reporting data. To do this, you should perform the following steps:

1. Identify the grain.
2. Select the required dimensions.
3. Identify the facts.

The details of these steps are discussed in the next topic.

Document and refine the models to determine the database logical schema

After you create initial dimensional models for each required business process, you can document the models to show:

- The measures in the fact tables.
- The related dimension tables.
- The attributes and hierarchies required in the dimension tables.

You can then iteratively refine the model to design the fact and dimension tables that will be required in the data warehouse database. Considerations for fact and dimension tables are discussed later in this module.

Design the physical data structures for the database

After completing the logical database design, you can consider the physical implementation of the database, including data file placement, indexes, table partitions, and data compression. These topics are discussed in more depth later in this module.

Dimensional Modeling

After you identify the business processes and conformed dimensions, you can document them in a matrix, as shown on the slide. This approach is based on the bus matrix design technique promoted by the Kimball Group.

 **Additional Reading:** For more information about using a bus matrix as part of a data warehouse design project, read *The Microsoft Data Warehouse Toolkit* (Wiley, 2011).

Business Processes	Time	Product	Customer	Salesperson	Factory Line	Shipper	Account	Department	Warehouse
Manufacturing	x	x				x			
Order Processing	x	x	x	x					
Order Fulfillment	x		x			x			
Financial Accounting	x						x	x	
Inventory Management	x	x							x

• **Grain:** 1 row per order item
 • **Dimensions:** Time (order date and ship date), Product, Customer, Salesperson
 • **Facts:** Item Quantity, Unit Cost, Total Cost, Unit Price, Sales Amount, Shipping Cost

You can then use the matrix to select each business process based on priority, and design a dimensional model by performing the following steps:

1. **Identify the grain.** The grain of a dimensional model is the lowest level of detail at which you can aggregate the measures. It is important to choose the level of grain that will support the most granular of reporting and analytical requirements—typically, the lowest level possible from the source data is the best option. For example, an order processing system might record data at two levels. There might be order-level data, such as the order date, salesperson, customer, and shipping cost, in addition to line item level data, such as the products in the order and their individual quantities, unit costs, and selling prices. To support the most granular analysis and reporting, the grain should be declared at the line item level, so the fact table will contain one row per line item.
2. **Select the required dimensions.** Next, determine which of the dimensions related to the business process should be included in the model. The selection of dimensions depends on the reporting and analytical requirements, specifically on the business entities by which users need to aggregate the measures. Almost all dimensional models include a time-based dimension, and the others generally become obvious as you review the requirements. At this stage, you might also begin to identify specific attributes of the dimensions that will be needed, such as the country, state, and city of a customer, or the color and size of a product. In the example on the slide, the Time, Customer, Product, and Salesperson dimensions are selected.



Note: In this example, the Time dimension is used for both order and ship date. Although it would be possible to define an individual dimension for each date type, it is more common to create a single time dimension and use it for multiple roles. In an analytical model, these multiple usages of the same dimension table are known as “role-playing dimensions”. This technique is most commonly used for timetables, but it can be applied to any dimension that is used in multiple ways. For example, a dimensional model for an airline flight-scheduling business process might use a single Airport dimension to support Origin and Destination role-playing dimensions.

3. **Identify the facts.** Finally, identify the facts that you want to include as measures. These are numeric values that can be expressed at the level of the grain chosen earlier and aggregated across the selected dimensions. Some facts will be taken directly from source systems, and others might be derived from the base facts. For example, you might choose Quantity and Unit Price facts from an order processing source system, and then calculate a total Sales Amount. Additionally, depending on the grain you choose for the dimensional model and the grain of the source data, you might need to allocate measures from a higher level of grain across multiple fact rows. For example, if the source system for the order processing business process includes a Tax measure at the order level—but the facts are to be stored at the line item level—you must decide how to allocate the tax amount across the line items. Typically, tax is calculated as a percentage of selling price, so it should be straightforward to apply the appropriate tax rate to each line item, based on the sales amount.

In the example on the slide, the Item Quantity, Unit Cost, and Unit Price measures are taken from the source system at the line item level. From these, the Total Cost and Sales Amount measures for each line item can be calculated. Additionally, the Shipping Cost measure is defined at the order level in the source system, so it must be allocated across the line items. You do this by dividing it equally across each row or applying a calculation that distributes the shared cost, based on the quantity of each item ordered, total line item weight, or some other appropriate formula.

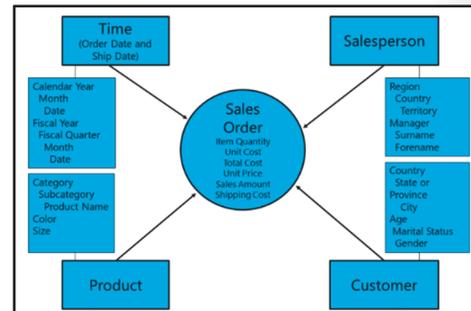
Documenting Dimensional Models

After you design the initial dimensional models for each business process, you can document them in a simple diagram. A common format for this documentation is a sun diagram, in which a fact table is shown at the center of the dimensions to which it is related.

As you refine the dimensional model, you can add more detail to the sun diagram, including the measures in the fact table and the attributes in the dimension tables. In most cases, some or all of the dimension attributes can be used to form a hierarchy for drill-down analysis—for example, enabling users to view aggregations of sales amount by year, month, and date or by country, state, and city. You can add these hierarchies to the sun diagram to help you communicate and validate model design with business stakeholders.

Eventually, the simple diagram will be refined to the point where it can be easily translated into a schema design for database tables. At this stage, you can use a diagramming tool such as Microsoft Visio® or a specialist database modeling tool to start designing the logical schema of your data warehouse.

Question: What is the difference between a snowflake schema and a star schema?



Lesson 2

Designing Dimension Tables

After designing the dimensional models for the data warehouse, you can translate the design into a logical schema for the database. However, before you design dimension tables, it is important to consider some common design patterns and apply them to your table specifications.

This lesson discusses some of the key considerations for designing dimension tables.

Lesson Objectives

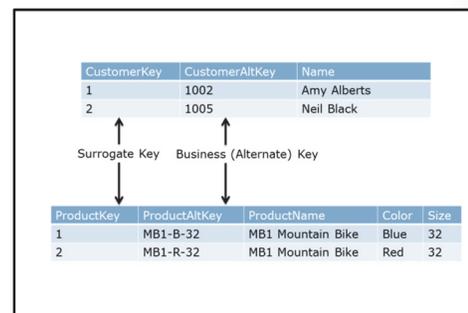
After completing this lesson, you will be able to:

- Describe considerations for dimension keys.
- Describe considerations for dimension attributes and hierarchies.
- Design dimensions that support values for “none” or “unknown”.
- Design appropriate slowly changing dimensions for your business requirements.
- Design time dimension tables.
- Design self-referencing dimension tables.
- Include junk dimensions in a data warehouse design where appropriate.

Considerations for Dimension Keys

Each row in a dimension table represents an instance of a business entity by which the measures in the fact table can be aggregated. Like other tables in a database, a key column uniquely identifies each row in the dimension table. In many scenarios, the dimension data is obtained from a source system in which a key is already assigned (sometimes referred to as the “business” key). When designing a data warehouse, however, it is standard practice to define a new “surrogate” key that uses an integer value to identify each row. A surrogate key is recommended for the following reasons:

- The data warehouse might use dimension data from multiple source systems, so it is possible that business keys are not unique.
- Some source systems use non-numeric keys, such as a globally unique identifier (GUID), or natural keys, such as an email address, to uniquely identify data entities. Integer keys are smaller and more efficient to use in joins from fact tables.
- Each row in a dimension table represents a specific version of a business entity instance. In some scenarios, the table might need to contain multiple rows that represent different versions of the same entity. These rows will have the same business key and won’t be uniquely identifiable without a surrogate key.



Typically, the business key is retained in the dimension table as an “alternate” key. Business keys that are based on natural keys can be familiar to users analyzing the data. For example, a **ProductCode** business key that users will recognize might be used as an alternate key in the **Product** dimension table. However, the main reason for retaining a business key is to make it easier to manage slowly changing dimensions when loading new data into the dimension table. The ETL process can use the alternate key as a lookup column to determine whether an instance of a business entity already exists in the dimension table.

Dimension Attributes and Hierarchies

In addition to the surrogate and alternate key columns, a dimension table includes a column for each attribute of the business entity that is needed to support reporting and analytical requirements. When designing a dimension table, you need to identify and include attributes that will be used in reports and analysis. Typically, dimension attributes are used in one of the following three ways:

CustKey	CustAltKey	Name	Country	State	City	Phone	Gender
1	1002	Amy Alberts	Canada	BC	Vancouver	555 123	F
2	1005	Neil Black	USA	CA	Irvine	555 321	M
3	1006	Ye Xu	USA	NY	New York	555 222	M

Hierarchies. Multiple attributes can be combined to form hierarchies that help users to drill down into deeper levels of detail. For example, the

Customer table in the slide includes **Country**, **State**, and **City** attributes that can be combined to form a natural geographical hierarchy. Business users can view aggregated fact data at each level—for example, to see sales order revenue by country. They can then access a specific country to see a breakdown by state, before drilling further into a specific state to see sales revenue by city.

Slicers. Attributes do not need to form hierarchies to be useful in analysis and reporting. Business users can group or filter data based on single-level hierarchies to create analytical subgroupings of data. For example, the **Gender** attribute in the **Customer** table can be used to compare sales revenue for male and female customers.

Drill-through detail. Some attributes have little value as slicers or members of a hierarchy. For example, it may be unlikely that a business user will need to analyze sales revenue by customer phone number. However, it can be useful to include entity-specific attributes to facilitate drill-through functionality in reports or analytical applications. For example, in a sales order report that helps users to drill down to the individual order level, users might want to double-click an order and drill through to see the name and phone number of the customer who placed it.

 **Note:** Terminology for interacting with report data can be confusing, and is sometimes used inconsistently. For clarity in this course, the term “drill down” means expanding a hierarchy to see the next level of aggregation, and “drill through” means viewing details outside the current hierarchy for a selected row. For example, while considering sales revenue by customer geography, you might view total revenue for a specific country in the hierarchy. You might then “drill down” to see a subtotal for each state within that country (the next level in the hierarchy), or “drill through” to see the country’s demographic details.

In the example on the slide, note that the **Name** column contains the full name of the customer. In a data warehouse table schema, it is not usually necessary to normalize the data to its most atomic level as is common in OLTP systems. In this example, it is unlikely that users will want to group or filter data by the customer’s first or last name; the data only has drill-through value at the full name level of detail.

Therefore, the **FirstName**, **MiddleName**, and **LastName** columns in the source system have been combined into a single **Name** field in the data warehouse.

Unknown and None

As a general rule, try to design your data warehouse in a way that eliminates, or at least minimizes, NULL values, particularly in fact table key columns that reference dimension tables. NULL values make it easy to accidentally eliminate rows from reports and produce misleading totals.

Identifying the semantic meaning of NULL

When you explore the source data for your BI solution, pay particular attention to how NULL values are used. The semantic meaning of NULL might be "None" or "Unknown," depending on the context. Only by examining the data and consulting the users, administrators, and developers who are familiar with the source system, will you be able to confidently determine which is relevant. In the example on the slide, the source data includes a column named **DiscountType**, in which two rows have a missing, or NULL, value. The fact that these rows include a non-zero Discount value indicates that NULL does not necessarily always mean "None", and is more likely to denote "Unknown." Additionally, on the rows where the **Discount** value is zero, the **DiscountType** value is consistently "N/A," implying that "N/A" is used in this system to mean "None."

To support these two cases, a row for each is added to the dimension table with appropriate surrogate keys, such as -1 for "Unknown" and 0 for "None". If the source systems were more ambiguous, you could add a single row to the dimension table to represent "None or Unknown."

NULL equality

Depending on the settings in a SQL Server database, you might not be able to compare NULL values for equality. In its strictest definition, NULL means unknown, so a "NULL = NULL" comparison is actually asking if one unknown value is the same as another; because both values are unknown, the answer is also unknown (and therefore NULL). You should not use NULL as the alternate key for the "Unknown" dimension row, because lookup queries during the ETL load process must compare this key to the data being loaded, to determine whether a dimension row already exists. Instead, use an appropriate key value that is unlikely to be the same as an existing business key; use the Transact-SQL ISNULL function to compare source rows with dimension rows, as shown in the following code sample:

```
SELECT d.DiscKey, s.DiscountType
FROM DimDiscount AS d
JOIN SourceData AS s ON ISNULL(s.DiscountType, 'Unknown') = d.DiscAltKey
```

OrderNo	Discount	DiscountType	Dimension Table		
			DiscKey	DiscAltKey	DiscountType
1000	1.20	Bulk Discount			
1001	0.00	N/A			
1002	2.00				
1003	0.50	Promotion	-1	Unknown	Unknown
1004	2.50	Other	0	N/A	None
1005	0.00	N/A	1	Bulk Discount	Bulk Discount
1006	1.50		2	Promotion	Promotion
		Source	3	Other	Other

- Identify the semantic meaning of NULL
 - Unknown or None?
- Do not assume NULL equality
 - Use ISNULL()

Designing Slowly Changing Dimensions

Slowly changing dimensions are a significant consideration in the design of dimension tables. You should try to identify requirements for maintaining historic dimension attribute values as early as possible in the design process.

There are three common techniques used to handle attribute value changes in slowly changing dimensions:

The diagram shows three types of slowly changing dimensions:

- Type 1:** A single table with columns: CustKey, CustAltKey, Name, Phone. An arrow points from a row (1, 1002, Amy Alberts, 555 123) to a new row (1, 1002, Amy Alberts, 555 222), indicating a direct update.
- Type 2:** A table with columns: CustKey, CustAltKey, Name, City, Current, Start, End. It shows a current row (1, 1002, Amy Alberts, Vancouver, Yes, 1/1/2000) and a historical row (4, 1002, Amy Alberts, Toronto, Yes, 1/1/2012).
- Type 3:** A table with columns: CustKey, CustAltKey, Name, Cars. It shows a current row (1, 1002, Amy Alberts, 0) and a historical row (1, 1002, Amy Alberts, 1) with a 'Prior Cars' column.

- Type 1.** These changes are the simplest type of slowly changing dimension to implement. Attribute values are updated directly in the existing dimension table row and no history is maintained. This makes Type 1 changes suitable for attributes that are used to provide drill-through details, but unsuitable for analytical slicers and hierarchy members where historic comparisons must reflect the attribute values as they were at the time of the fact event.
- Type 2.** These changes involve the creation of a fresh version of the dimension entity in the form of a new row. Typically, a bit column in the dimension table is used as a flag to indicate which version of the dimension row is the current one. Additionally, datetime columns are often used to indicate the start and end of the period for which a version of the row was (or is) current. Maintaining start and end dates makes it easier to assign the appropriate foreign key value to fact rows as they are loaded, so they are related to the version of the dimension entity that was current at the time the fact occurred.
- Type 3.** In a Type 3 change, the previous value (or sometimes a complete history of previous values) is maintained in the dimension table row. This requires modifying the dimension table schema to accommodate new values for each tracked attribute, and can result in a complex dimension table that is difficult to manage. These changes are rarely used.

After you define the dimensional model for the data warehouse and are evolving your design from a sun diagram to a database schema, it can be useful to annotate dimension attributes to indicate what kind of slowly changing dimension changes they must support. This will help you plan the metadata columns required for each dimension table.

Time Dimension Tables

Most analysis and reporting requirements include a need to aggregate values over time periods, so almost every data warehouse includes a time dimension table. When you design a time dimension table, you must take into account the following considerations:

Surrogate key. Although best practice for surrogate keys in dimension tables is normally to use a simple integer value with no semantic meaning, time dimension tables can benefit from an integer representation of the date or time that the row represents. Ideally, the values should be in

DateKey	DateAltKey	MonthDay	WeekDay	Day	MonthNo	Month	Year
00000000	01-01-1753	NULL	NULL	NULL	NULL	NULL	NULL
20130101	01-01-2016	1	3	Tue	01	Jan	2016
20130102	01-02-2016	2	4	Wed	01	Jan	2016
20130103	01-03-2016	3	5	Thu	01	Jan	2016
20130104	01-04-2016	4	6	Fri	01	Jan	2016

- Surrogate key
- Granularity
- Range
- Attributes and hierarchies
- Multiple calendars
- Unknown values

ascending order relative to the dates they represent; therefore, the best approach is to concatenate the integer values for each date part in descending order of scope. For example, using the YYYYMMDD pattern to represent dates, the value for January 31, 2016 would be 20160131. This ensures that the value used for the next sequential date of February 1, 2016, is a higher value of 20160201. Ascending values are recommended because data warehouse queries typically filter on a range of date or time values. The ascending numeric key means you can use indexes and partitions that store the fact data in chronological order; the query optimizer can then use sequential scans to read the data. Additionally, the actual datetime value for the row is generally used as the alternate key to support datetime functions or client applications that can apply datetime-specific logic.

Granularity. The level of granularity used for a time dimension table depends on business requirements. For many reporting and analysis scenarios, such as viewing details about sales orders, the lowest level of granularity likely to be required is a day. However, in some scenarios, users might need to aggregate facts by hours, minutes, or seconds, or even smaller increments. The lower the level of granularity used, the more rows will exist in the dimension table, and storing a row for increments of less than a day can result in extremely large tables. An alternative approach is to create a “date” dimension table that contains a row for each day, and a “time” dimension table that stores a row for each required time increment in a 24-hour period. Fact tables that are used for analysis of measures at the day level or higher can only be related to the date dimension table. Facts measured at smaller time increments can be related to both date and time dimension tables.

Range. Typically, a time dimension table stores a row for each increment between a start point and an end point—with no gaps. For example, a data warehouse time dimension used to analyze sales orders might have a row for each day between the first ever and most recent orders, even if no orders were placed during the intervening days. In reality, the start and end dates are typically based on key calendar dates. For example, the start date might be January 1 of the year the company started trading, or when its first fiscal year began. The end date is usually some future point, such as the end of the current year. To maintain a buffer of future dates, more rows are added automatically as the end date gets closer. If the data warehouse will be used to create and store projections or budget figures for future operations, you should choose an end date that is far enough into the future to accommodate any projections likely to be used in the data warehouse.

Attributes and hierarchies. You need to include attributes for each time period by which data will be aggregated—for example, year, quarter, month, week, and day. These attributes tend to form natural hierarchies. You can also add attributes to be used as slicers, such as “weekday” which, for example, would help users compare typical sales volumes for each day of the week. In addition to numeric values, you might want to include attributes for date element names, such as “month” and “day”. You will then have more accessible reports where, for example, users could compare sales in March and April instead of month 3 and 4. You should also include the numeric equivalents so that client applications can use them to sort data into the correct chronological order, such as sorting months by month number instead of by month name.

Multiple calendars. Many organizations need to support multiple calendars; for example, a normal year that runs from January to December, and a fiscal calendar, which might run from April to March. If this is the case in your data warehouse, you can either create a separate time dimension table for each calendar or, preferably, include attributes for all alternative calendar values in a single time dimension table. For example, a time dimension table might include a **Calendar Year** and a **Fiscal Year** attribute.

Unknown values. In common with other dimension tables, you might need to support facts for which a date or time value is unknown. Instead of requiring a NULL value in the fact table, consider creating a row for unknown values in the time dimension table. You can use an obvious surrogate key value, such as 00000000, for this row. However, because the alternate key must be a valid date, you should choose one outside the normal range of business operations, such as January 1, 1753, or December 31, 9999. These are the minimum and maximum values supported by the datetime data type.

Populating a time dimension table

Unlike most other tables in a data warehouse, time dimension tables are not usually populated with data that has been extracted from a source system. Generally, the data warehouse developer populates the time dimension table with rows at the appropriate granularity. These rows usually consist of a numeric primary key that is derived from the temporal value (such as 20160101 for January 1, 2016) and a column for each dimension attribute (such as the date, day of year, day name, month of year, month name, year). To generate the rows for the time dimension table, you can use one of the following techniques:

Create a Transact-SQL script. Transact-SQL includes many date and time functions that you can use in a loop construct to generate the required attribute values for a sequence of time intervals. The following Transact-SQL functions are commonly used to calculate date and time values:

DATEPART (datepart, date) returns the numerical part of a date, such as the weekday number, day of month, and month of year.

DATENAME (datepart, date) returns the string name of a part of the date, such as the weekday name and month name.

MONTH (date) returns the month number of the year for a given date.

YEAR (date) returns the year for a given date.

Use Microsoft® Excel®. Excel includes several functions that you can use to create formulas for date and time values. You can then use the autofill functionality in Excel to quickly create a large table of values for a sequence of time intervals.

Use a BI tool to autogenerate a time dimension table. Some BI tools include time dimension generation functionality that you can use to create a time dimension table quickly.

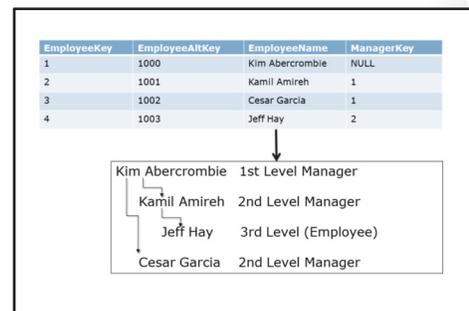
Self-Referencing Dimension Tables

A common requirement in a data warehouse is to support dimensions with parent-child hierarchies. For example, an employee dimension might consist of managers, each of whom has employees reporting to him or her, who in turn might have reports of their own.

Typically, parent-child hierarchies are implemented as self-referencing tables, in which a column in each row is used as a foreign key reference to a primary-key value in the same table. Some client applications, including SQL Server Analysis Services, are aware of self-joins and can automatically handle parent-child hierarchies in a dimension. For other applications, you might need to implement some custom, recursive logic, so that analysis and reporting of these hierarchies can take place.

When you implement a self-referencing dimension table in a data warehouse, you should think about the following considerations:

- Like all dimension load operations, when records are to be added to the dimension table, the ETL process must look up the alternate key to determine if a record already exists for the entity. However, the alternate key of the parent record must also be looked up to determine the correct surrogate key to use in the foreign key column.



- You may have to deal with a situation where you need to add a record for which the parent record has not yet been loaded.
- Supporting Type 2 slowly changing dimensions in a self-referencing dimension table can be complex. In a worst case scenario, you might perform a Type 2 change that results in a new row and, therefore, a new surrogate key. You might then need to cascade that Type 2 change to create new rows for all descendants of the entity, if the change has not altered the parent-child relationships.

Junk Dimensions

In some reporting and analytical requirements, there are useful attributes for grouping or filtering facts which do not belong in any of the dimensions defined in the dimensional model. If these attributes have low cardinality—when there are only a few discrete values—you can group them into a single dimension table containing miscellaneous analytical values. This kind of dimension table is generally referred to as a “junk dimension” and is used to avoid creating multiple, very small dimension tables.

For example, a sales order dimensional model might include “true or false” indicators for orders where goods were out of stock or where free shipping was provided. There may be a column that stores “credit” or “debit” to indicate the payment method. Instead of creating a dimension table for each of these attributes, you could merge them in every possible combination in a junk dimension table.

JunkKey	OutOfStockFlag	FreeShippingFlag	CreditOrDebit
1	1	1	Credit
2	1	1	Debit
3	1	0	Credit
4	1	0	Debit
5	0	1	Credit
6	0	1	Debit
7	0	0	Credit
8	0	0	Debit

- Combine low-cardinality attributes that don't belong in existing dimensions into a junk dimension
- Avoids creating many small dimension tables

Categorize Activity

Place each item into the appropriate category which is a dimension type. Indicate your answer by writing the category number to the right of each item.

Items	
1	Gender
2	Fiscal Year
3	Invoice Number
4	Month
5	Out of Stock Indicator

Category 1	Category 2	Category 3
Slicer	Time Dimension	Junk Dimension

Lesson 3

Designing Fact Tables

Fact tables contain the numeric measures that can be aggregated across the dimensions in your dimensional model, and can become extremely large. It is important to design them carefully with reporting and analytical requirements, performance, and manageability in mind.

This lesson discusses common considerations for fact table design.

Lesson Objectives

After completing this lesson, you will be able to:

- Determine appropriate columns for a fact table.
- Describe the types of measure that are stored in a fact table.
- Describe common types of fact table.

Fact Table Columns

A fact table usually consists of the following kinds of column:

Dimension keys. Fact tables reference dimension tables by storing the surrogate key for each related dimension. In this way, a row in a fact table is conceptually an intersection between the dimension tables to which it relates. For example, recording a sales order placed on a specific date, for a specific product, by a specific customer. You can add foreign key constraints to these columns, which will help the SQL Server query optimizer detect star joins. However, constraints can slow down data load operations, and because the surrogate keys are generated during a controlled ETL process, they do little to enforce referential integrity.

Measures. In most cases, a fact table is primarily used to store numeric measures that can be aggregated by the related dimensions. For example, a row in a fact table recording sales orders might include a column for the sales amount, which can then be aggregated by the dimensions to show sales amount by date, product, or customer.

In some cases, a fact table contains no measures and is simply used to indicate that an intersection occurred between the related dimensions. For example, a fact table in a manufacturing dimensional model might record a single row each time a product assembly is completed, indicating the product and date dimension keys. The fact table can then be used to calculate the number of times an assembly of each product was completed per time period by simply counting the distinct rows. A fact table with no numeric measure columns is sometimes referred to as a “factless fact table”.

Degenerate dimensions. Sometimes, a fact has associated attributes that are neither keys nor measures, but which can be useful to group or filter facts in a report or analysis. For example, a Purchase Order Number on a Purchase Order links all the order lines together for a specific Purchase Order, but the Purchase Order Number is not a dimension or a fact. You might include this column in the fact table where client applications can use it as a “degenerate dimension” by which the fact data can be

• Dimension Keys					
OrderDateKey	ProductKey	CustomerKey	OrderNo	Qty	SalesAmount
20160101	25	120	1000	1	350.99
20160101	99	120	1000	2	6.98
20160101	25	178	1001	2	701.98

• Measures					
OrderDateKey	ProductKey	CustomerKey	OrderNo	Qty	SalesAmount
20160101	25	120	1000	1	350.99
20160101	99	120	1000	2	6.98
20160101	25	178	1001	2	701.98

• Degenerate Dimensions					
OrderDateKey	ProductKey	CustomerKey	OrderNo	Qty	SalesAmount
20160101	25	120	1000	1	350.99
20160101	99	120	1000	2	6.98
20160101	25	178	1001	2	701.98

aggregated. In effect, including degenerate dimension columns in a fact table means it can be used additionally as a dimension table. Using degenerate dimensions can be a good alternative to using a junk dimension if the analytical attributes are specific to a single fact table.



Note: Unlike most database tables, a fact table does not necessarily require a primary key. Unless you have a business requirement to uniquely identify each row in the fact table, you should avoid creating a unique key column for the fact table and defining a primary key constraint. Facts are generally aggregated, and queries rarely need to identify an individual fact row.

In some cases, the combination of dimension keys can uniquely identify a fact row, but this is not guaranteed. For example, a customer could purchase the same product twice in one day.

Types of Measure

Fact tables can contain the following three kinds of measure:

Additive measures. These can be summed across all dimensions. For example, you could use a SalesAmount measure in a fact table with OrderDateKey, ProductKey, and CustomerKey dimension keys to calculate total sales amount by time period (such as month), product, or customer.

Semi-additive measures. These can be summed by some dimensions, but not others. Commonly, semi-additive measures cannot be summed by time dimensions. For example, the number of items in stock at the end of each day might be recorded as a StockCount measure in a fact table with DateKey and ProductKey dimension keys that can be used to calculate a total stock count for all products. However, summing the stock count across all the days in a month does not result in a total stock count value for that period. To find out how many products are in stock at the end of the month, you must use only the StockCount value for the final day.

Non-additive measures. These cannot be summed by any dimension. For example, a fact table for sales orders might include a ProfitMargin measure that records the profit margin for each order. However, you cannot calculate the overall margin for any dimension by summing the individual profit margins.

Generally, semi-additive and non-additive measures can be aggregated by using other functions. For example, you could find the minimum stock count for a month or the average profit margin for a product. Understanding the ways in which the measures can be meaningfully aggregated is useful when testing and troubleshooting data warehouse queries and reports.

• Additive			
OrderDateKey	ProductKey	CustomerKey	SalesAmount
20160101	25	120	350.99
20160101	99	120	6.98
20160102	25	178	701.98

• Semi-Additive			
DateKey	ProductKey	CustomerKey	StockCount
20160101	25	23	
20160101	99	118	
20160102	25	22	

• Non-Additive			
OrderDateKey	ProductKey	CustomerKey	ProfitMargin
20160101	25	120	25
20160101	99	120	22
20160102	25	178	27

Types of Fact Table

Generally, data warehouses include fact tables that are one of the following three types:

Transaction fact tables. The most common kind of fact table is a “transaction” fact table, in which each row records a transaction or event at an appropriate grain. For example, a fact table might record sales orders at the line item grain, in which each row records the purchase of a specific item. Transaction fact table measures are usually additive.

Periodic snapshot tables. These record measure values at a specific point in time. For example, a fact table might record the stock movement for each day, including the opening and closing stock count figures. Measures in a periodic snapshot fact table are often semi-additive.

Accumulating snapshot fact tables. These can be used in scenarios where you might want to use a fact table to track the progress of a business process through multiple stages. For example, a fact table might track an order from initial purchase through to delivery by including a date dimension-key field for the order date, shipping date, and delivery date. The ShipDate and DeliveryDate columns for orders that have been placed but not yet shipped will contain the dimension key for an “Unknown” or “None” row in the time dimension table. These will be updated to reflect the appropriate dimension key as the order is shipped and delivered.

OrderDateKey	ProductKey	CustomerKey	OrderNo	Qty	Cost	SalesAmount
20160101	25	120	1000	1	125.00	350.99
20160101	99	120	1000	2	2.50	6.98
20160101	25	178	1001	2	250.00	701.98

DateKey	ProductKey	OpeningStock	UnitsIn	UnitsOut	ClosingStock
20160101	25	25	1	3	23
20160101	99	120	0	2	118

OrderNo	OrderDateKey	ShipDateKey	DeliveryDateKey
1000	20160101	20160102	20160105
1001	20160101	20160102	00000000
1002	20160102	00000000	00000000

Check Your Knowledge

Question	
What kind of measure is a stock count?	
Select the correct answer.	
<input type="checkbox"/>	Additive
<input type="checkbox"/>	Semi-additive
<input type="checkbox"/>	Non-additive

Lesson 4

Physical Design for a Data Warehouse

After designing the logical schema for the data warehouse, you need to implement it as a physical database. This requires careful planning for file placement, data structures such as partitions and indexes, and compression. This lesson discusses considerations for all these aspects of the physical database design.

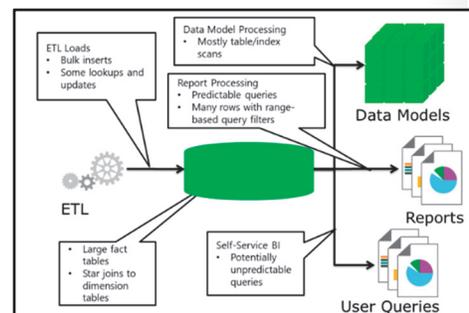
Lesson Objectives

After completing this lesson, you will be able to:

- Describe typical data warehouse I/O activity.
- Plan file placement for a data warehouse.
- Plan partitioning for data warehouse tables.
- Design effective indexes for data warehouse queries.
- Design the model to manage many-to-many relationships.
- Plan data compression for a data warehouse.
- Design views in a data warehouse.

Data Warehouse I/O Activity

Before designing the physical database for the data warehouse, it is useful to consider the types of workload it must support and the data it has to store. The database must store potentially very large fact tables with millions of rows, and dimension tables that are often related to fact tables by a single join. Typically, the I/O activity in the database is generated by one of the workloads described in this section or caused by maintenance operations, such as backups.



ETL

ETL processes affect the data warehouse when they load new or updated data into the data warehouse tables. In most cases, the inserts are performed as bulk load operations to minimize logging and constraint checking. The load process may involve some lookup operations to find alternate keys for slowly changing dimensions, and some update operations for Type 1 dimension changes or data modifications in fact tables where appropriate. Depending on the design of the data structures, ETL load operations might also involve dropping and rebuilding indexes and splitting partitions.

Data models—tabular or multidimensional

After each new load, any data models based on the data warehouse must be processed. The data model might be multidimensional, such as an OLAP cube, or it could be tabular. Loading the data into the data model involves reading data from the data warehouse tables into the data model and pre-aggregating measures to optimize analysis queries. Depending on the size of the data warehouse and the time window

for the processing operation, the entire data model can be completely processed after each load; or an incremental approach can be used, in which only new or modified data is handled.

Because of the volume of data being loaded into the model, the I/O activity typically involves sequential table scans to read entire tables, particularly when performing a full process of the data model.

Reports

In some scenarios, all reporting is performed against the data models, so it does not affect the data warehouse tables. However, it is common for some reports to query the data warehouse directly. In scenarios where IT-provided reports are supported, the queries are generally predictable and retrieve many rows with range-based query filters, often on a date field.

User queries

If self-service reporting is supported, users might be able to execute queries in the data warehouse or use tools that generate queries on their behalf. Depending on the query expertise of the users, this can result in complex, unpredictable queries.

Considerations for Database Files

In Module 2: *Planning Data Warehouse Infrastructure*, some key considerations for data warehouse storage hardware were discussed and it was recommended that storage be provided by multiple disks configured as RAID 10 or RAID 5 arrays. This storage is presented to the data warehouse server as multiple logical disks that are sometimes referred to as logical unit numbers (LUNs)—though technically a LUN is used to identify a unit of SCSI-based storage. When designing the file placement for your data warehouse, you must decide how best to use these logical disks.

- Data files and filegroups
- Staging tables
- tempdb
- Transaction logs
- Backup files

Data files and filegroups

Data files are used to pre-allocate disk storage for database objects. When planning files for a data warehouse, consider the following guidelines:

- Create files with an initial size, based on the eventual size of the objects that will be stored on them. This pre-allocates sequential disk blocks and helps avoid fragmentation.
- Disable autogrowth. If you begin to run out of space in a data file, it is more efficient to explicitly increase the file size by a large amount rather than rely on incremental autogrowth.

Because the logical disks for the database files are typically already configured as RAID 10 or RAID 5 arrays, you do not need to use filegroups to distribute tables across physical disk platters to improve I/O performance. However, you should consider the following guidance for using filegroups in a data warehouse:

- Create at least one filegroup in addition to the primary one, and then set it as the default filegroup so that you can separate data tables from system tables.
- Consider creating dedicated filegroups for extremely large fact tables and using them to place those fact tables on their own logical disks.

- If some tables in the data warehouse are loaded on a different schedule to others, consider using filegroups to separate the tables into groups that can be backed up independently.
- If you intend to partition a large fact table, create a filegroup for each one so that older, stable rows can be backed up, and then set as read-only.

Staging tables

Most data warehouses require staging tables to support incremental data loads from the ETL process. In some cases, you might use a separate staging database in addition to staging tables in the data warehouse itself. Consider the following recommendations for staging tables:

- If a separate staging database is to be used, create it on a logical disk, distinct from the data warehouse files.
- If the data warehouse will include staging tables, create a file and filegroup for them on a logical disk, separate from the fact and dimension tables.
- An exception to the previous guideline is made for staging tables that will be switched with partitions to perform fast loads. These must be created on the same filegroup as the partition with which they will be switched.

tempdb

The **tempdb** database is used for temporary objects required for query processing. To avoid fragmentation of data files, place it on a dedicated logical disk and set its initial size based on how much it is likely to be used. You can leave **autogrowth** enabled, but set the growth increment to be quite large to ensure that performance is not interrupted by frequent growth of the **tempdb** database. Additionally, consider creating multiple files for the **tempdb** database to help minimize contention during page free space (PFS) scans, as temporary objects are created and dropped.

Transaction logs

Generally, the transaction mode of the data warehouse staging database and the **tempdb** database should be set to **Simple** to avoid having to truncate transaction logs. Additionally, most of the inserts in a data warehouse are typically performed as bulk load operations which are not logged. To avoid disk resource conflicts between data warehouse I/O and logging, place the transaction log files for all databases on a dedicated logical disk.

Backup files

You will need to implement a backup routine for the data warehouse and, potentially, for a staging database. In most cases, you will back up these databases to disk, so allocate a logical disk for this purpose. You could allocate multiple logical disks and perform a mirrored backup, but because the disks are already configured as RAID 5 or RAID 10 arrays, this would be of little benefit from a performance perspective. Note that the backup files should be copied to offsite storage to provide protection in the case of a complete storage hardware failure or natural disaster.

Table Partitioning

Partitioning a table distributes data based on a function that defines a range of values for each partition. A partition scheme maps partitions to filegroups, and the table is partitioned by applying the partition scheme to the values in a specified column.



Additional Reading: For information about how to implement partitioning, see *Partitioned Tables and Indexes* in the SQL Server Technical Documentation.

OrderDateKey	ProductKey	CustomerKey	OrderNo	Qty	Cost	SalesAmount
20160101	25	120	1000	1	125.00	350.99
20160101	99	120	1000	2	2.50	6.98
20160101	25	178	1001	2	250.00	701.98
...						
20160201	23	76	2124	1	95.00	125.00
20160201	89	6	2125	1	45.00	76.99

Pre-2014	2014	2015	Jan 2016	Feb 2016

A partition scheme, grouping older records by year and current records by month.

Why use partitioning?

Partitioning a large table can produce the following benefits:

- Improved query performance.** By partitioning a table across filegroups, you can place specific ranges of data on different disk spindles, which can improve I/O performance. In most data warehouses, the disk storage is already configured as a RAID 10 or RAID 5 array, so partitioning usually has little benefit. However, when using a mix of fast solid state storage for recent, frequently-accessed data, and mechanical disks for older, less queried rows, you can use partitioning to balance disk performance against storage costs. The biggest performance gain from partitioning in a data warehouse is realized when queries return a range of rows that are filtered on the partitioning key. In this case, the query optimizer can eliminate partitions that are not within the filter range, and dramatically reduce the number of rows that need to be read.
- More granular manageability.** When you partition a large table, you can perform some maintenance operations at the partition level instead of on the whole table. For example, indexes can be created and rebuilt on a per-partition basis, and compression can be applied to individual partitions. You can also back up and restore partitions independently by mapping them to filegroups; you can then back up older data once, and configure the backed-up partitions as read-only. Future backups can be limited to the partitions that contain new or updated data.
- Improved data load performance.** The biggest benefit of using partitioning in a data warehouse is that it helps you to load many rows quickly by switching a staging table with a partition. This technique dramatically reduces the time taken by ETL data loads, and with the right planning, it can be achieved with minimal requirements to drop and rebuild indexes.

Best practices for partitioning in a data warehouse

When planning a data warehouse, consider the following best practices for partitioning:

- Partition large fact tables.** Fact tables of around 50 GB or more should usually be partitioned for the reasons described earlier. In general, fact tables benefit from partitioning more than dimension tables.
- Partition on an incrementing date key.** When defining a partition scheme for a fact table, use a date key that reflects the age of the data as it is incrementally loaded by the ETL process. For example, if a fact table contains sales order data, partitioning on the order date ensures that the most recent orders are in the last partition and the earliest ones are in the first.

- **Design the partition scheme for ETL and manageability.** In a data warehouse, the query performance gains realized by partitioning are small compared to the manageability and data load performance benefits. Ideally, your partitions should reflect the ETL load frequency (for example, monthly, weekly, daily) because this simplifies the load process. However, you may want to merge partitions periodically to reduce their overall number. For example, at the start of each year, you could merge the monthly partitions for the previous year into a single partition for the whole year.
- **Maintain an empty partition at the start and end of the table.** You can use an empty partition at the end of the table to simplify the loading of new rows. When a new set of fact table rows must be loaded, you can place them in a staging table, split the empty partition (to create two empty partitions), and then switch the staged data with the first empty partition. This loads the data into the table and leaves the second empty partition you created at the end of the table ready for the next load. You can use a similar technique to archive or delete obsolete data at the beginning of the table.

 **Note:** Partitioning is available only in SQL Server Enterprise edition. In SQL Server 2012 and later releases, the number of partitions per table is limited to 15,000 on a 64-bit system and 1,000 on a 32-bit system.

Demonstration: Partitioning a Fact Table

In this demonstration, you will see how to:

- Create a partitioned table
- View partition metadata
- Split a partition
- Merge partitions

Demonstration Steps

Create a Partitioned Table

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod03** folder, run **Setup.cmd** as Administrator.
3. In the **User Account Control** dialog box, click **Yes**, and then wait for the script to finish.
4. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the database engine by using Windows authentication.
5. On the **File** menu, point to **Open**, and then click **File**. In the **Open File** dialog box, navigate to the **D:\Demofiles\Mod03** folder, click **Partitions.sql**, and then click **Open**.
6. Select the code under the comment **Create partition function and scheme**, and then click **Execute**.
This creates a partition function that defines four ranges of values (less than 20140101, 20140101 to 20150100, 20150101 to 20160100, and 20160101 and higher), and a partition scheme that maps these ranges to the FG0000, FG2014, FG2015, and FG2016 filegroups.
7. Select the code under the comment **Create a partitioned table**, and then click **Execute**. This creates a partitioned table on the partition scheme you created previously.
8. Select the code under the comment **Insert data into the partitioned table**, and then click **Execute**. This inserts four records into the table.

View Partition Metadata

1. Select the code under the comment **Query the table**, and then click **Execute**. This retrieves rows from the table and uses the **\$PARTITION** function to show which partition the **datekey** value in each row is assigned to. This function is useful for determining which partition of a partition function a specific value belongs in.
2. Select the code under the comment **View filegroups, partitions, and rows**, and then click **Execute**. This code uses system tables to show the partitioned storage and the number of rows in each partition. Note that there are two empty partitions, one at the beginning of the table, and one at the end.

Split a Partition

1. Select the code under the comment **Add a new filegroup and make it the next used**, and then click **Execute**. This creates a new filegroup named FG2017 and adds it to the partition scheme as the next used partition.
2. Select the code under the comment **Split the empty partition at the end**, and then click **Execute**. This creates a new partition for values of 20170101 and higher and assigns it to the next used filegroup (**FG2017**), leaving an empty partition for values between 20160101 and 20170100.
3. Select the code under the comment **Insert new data**, and then click **Execute**. This inserts two new rows into the partitioned table.
4. Select the code under the comment **View partition metadata**, and then click **Execute**. This shows that the two rows inserted in the previous step are in partition 4, and that partition 5 (on **FG2017**) is empty.

Merge Partitions

1. Select the code under the comment **Merge the 2014 and 2015 partitions**, and then click **Execute**. This merges the partition that contains the value 20140101 into the previous partition.
2. Select the code under the comment **View partition metadata**, and then click **Execute**. This shows that partition 2 (on **FG2014**) now contains four rows, and that the partition previously on FG2015 has been removed.
3. Close **Partitions.sql** but keep SQL Server Management Studio open for the next demonstration.

Considerations for Indexes

Most databases use indexes to maximize query performance, and planning them is an important part of the database design process. Before deciding which indexes to create, you need to understand the workloads that the database must support, and balance the need for improved query performance against the effect that indexes will have on data inserts and updates—in addition to the overhead of maintaining indexes.

At first glance, a data warehouse seems to support mostly read operations, so some inexperienced BI professionals are tempted to create many indexes on all tables to support queries. However, another significant workload in most data warehouses is the

- Dimension table indexes
 - Clustered index on surrogate key column
 - Nonclustered index on business key and SCD columns
 - Nonclustered indexes on frequently searched columns
- Fact table indexes
 - Clustered index on most commonly searched date key
 - Nonclustered indexes on other dimension keys
 - Or
 - Columnstore index on all columns
 - Composite index key comprises up to 16 columns

regular ETL data load, which can often involve many inserts and updates. Too many indexes can slow down the ETL load process, and the need to periodically reorganize or rebuild indexes can create a significant maintenance overhead.

The first consideration is to determine whether any indexes are required in your data warehouse. It may seem unconventional to consider not creating indexes, but if the volume of fact data is relatively small, and all user access to the data is through a data model that is fully processed after each data load, there might be little performance benefit in maintaining indexes in the data warehouse. However, if your data warehouse does not match this restrictive description, you will probably need to consider creating some indexes. As with any database, the indexes you create depend on the specific queries your data warehouse must support and the need to balance the performance of those queries against data inserts and updates, and index maintenance.

However, in most data warehouse scenarios, you should consider the guidelines in this topic as a starting point for index design.

Dimension table indexes

When designing indexes for dimension tables, consider the following guidelines:

- Create a clustered index on the surrogate key column. This column is used to join the dimension table to fact tables, and a clustered index will help the query optimizer minimize the number of reads required to filter fact rows.
- Create a nonclustered index on the alternate key column and include the slowly changing dimension current flag, start date, and end date columns. This index will improve the performance of lookup operations during ETL data loads that need to handle slowly changing dimensions.
- Create nonclustered indexes on frequently searched attributes, and consider including all members of a hierarchy in a single index.

Fact table indexes

When designing indexes for a fact table, consider the following guidelines:

- Create a clustered index on the most commonly-searched date key. Date ranges are the most common filtering criteria in most data warehouse workloads, so a clustered index on this key should be particularly effective in improving overall query performance.
- Create nonclustered indexes on other, frequently-searched dimension keys.

Composite keys

A composite key is a primary key that includes more than one column. In a data warehouse, the fact table comprises the primary keys of the dimension tables as foreign key columns, and a combination of these columns is often a suitable candidate for a composite key, because a surrogate is not always a necessity on a fact table. When adding a composite key, consider the following guidelines:

- You can have up to 16 columns, which in combination must not exceed more than 900 bytes.
- All columns must be from the same table.
- Use a composite index when you have more than one column that would be searched simultaneously.
- When defining the list of columns, specify the most unique first.

When joining tables, you should join on all columns that comprise the composite key, with your query referencing the first column of the composite key for it to be considered for use by the query optimizer.

For more information on adding a surrogate key to a fact table, see:



Design Tip #81 Fact Table Surrogate Key

<https://aka.ms/a5fv47>

Columnstore indexes

SQL Server supports columnstore indexes, an in-memory solution that uses xVelocity compression technology to organize index data in a column-based format instead of the traditional row-based format. Columnstore indexes are specifically designed to enhance the performance of queries against large fact tables joined to smaller dimension tables in a star schema—they can significantly improve the speed of most data warehouse queries. In many cases, you can achieve the same performance improvements or better by replacing the recommended fact table indexes described previously with a single columnstore index that includes all the fact table columns. Some queries do not benefit from columnstore indexes. For example, queries that return an individual row from a dimension table will generally perform better by using a conventional clustered or nonclustered index. However, for most typical data warehouse queries that aggregate many fact rows by one or more dimension attributes, columnstore indexes can be very effective.

Columnstore indexes can be clustered or nonclustered.

Clustered columnstore indexes

A clustered columnstore index has the following characteristics:

- It can be created only in the Enterprise, Developer, and Evaluation editions of SQL Server.
- It includes all the columns in the table.
- It is the only index on the table.
- It does not store the columns in a sorted order, but rather optimizes storage for compression and performance.
- It can be updated.

Clustered columnstore indexes can be updated, and you can bulk load, insert, update, and delete data in a clustered columnstore indexed table using standard Transact-SQL statements.

Clustered columnstore indexes store the data in compressed columnstore segments, but some data is stored in a rowstore table, referred to as the *deltastore*, as an intermediary location until it can be compressed and moved into a columnstore segment. The following rules are used to manage data modifications:

- When you use an INSERT statement to insert a new row, it remains in the *deltastore* until there are enough rows to meet the minimum size for a *rowgroup*, which is then compressed and moved into the columnstore segments.
- When you execute a DELETE statement, affected rows in the *deltastore* are physically deleted, while affected data in the columnstore segments is marked as deleted; the physical storage is only reclaimed when the index is rebuilt.
- When you execute an UPDATE statement, affected rows in the *deltastore* are updated, while affected rows in the columnstore are marked as deleted and a new row is inserted into the *deltastore*.

Nonclustered columnstore indexes

A nonclustered columnstore index has the following characteristics:

- It can include some or all the columns in the table.
- It can be combined with other indexes on the same table.
- It cannot be updated; tables containing a nonclustered columnstore index are read-only.

Nonclustered columnstore indexes are read-only, but given that a typical data warehouse is a static database that is updated periodically through an ETL process, this is less of a limitation than might at first appear. However, administrators do need to plan how to handle updates to data in tables that have nonclustered columnstore indexes.

There are two ways to update nonclustered columnstore indexes:

- **Periodically drop the index, perform updates to the table, and then recreate the index.** This approach is the simplest way of handling updates, and fits in with the way that many organizations already perform loads into their data warehouses. The disadvantage of this approach is that creating a columnstore index can be time-consuming when the base table is very large; this can be problematic when the time window for performing a data load is relatively small.
- **Use table partitioning.** When you create an index on a partitioned table, SQL Server automatically aligns the index with the table, meaning both are divided in the same way. When you switch a partition out of the table, the aligned index partition follows. You can use partition switching to perform inserts, updates, merges, and deletes:
 - To perform a bulk insert, partition the table, load new data into a staging table, build a columnstore index on the staging table, and then use partition switching to load the data into the partitioned data warehouse table.
 - For other types of update, you can switch a partition out of the data warehouse table into a staging table; drop or disable the columnstore index on the staging table; perform the updates; recreate or rebuild the columnstore index on the staging table; and then switch the staging table back into the data warehouse table.

Demonstration: Creating Indexes

In this demonstration, you will see how to:

- Create an index on a dimension table.
- View index usage and execution statistics.
- Create an index on a fact table.
- Create a columnstore index.

Demonstration Steps

Create Indexes on Dimension Tables

1. Ensure that you have completed the previous demonstration in this module.
2. In SQL Server Management, on the **File** menu, point to **Open**, and then click **File**. In the **Open File** dialog box, navigate to the **D:\Demofiles\Mod03** folder, click **Indexes.sql**, and then click **Open**.

3. Select the code under the comment **Create indexes on the DimDate table**, and then click **Execute**. This creates a clustered index on the surrogate key column, and nonclustered indexes on commonly queried attribute columns.
4. Select the code under the comment **Create indexes on the DimCustomer table**, and then click **Execute**. This creates a clustered index on the surrogate key column, and nonclustered indexes on commonly queried attribute columns.
5. Select the code under the comment **Create indexes on the DimProduct table**, and then click **Execute**. This creates a clustered index on the surrogate key column, and nonclustered indexes on a commonly queried attribute column.

View Index Usage and Execution Statistics

1. Select the code under the comment **Create a fact table**, and then click **Execute**. This creates a fact table named **FactOrder** that contains more than 7.5 million rows from the existing data in the dimension tables.
2. On the toolbar, click **Include Actual Execution Plan**.
3. Select the code under the comment **View index usage and execution statistics**, and then click **Execute**.
This enables statistics messages, and then queries the tables in the data warehouse to view orders for the previous six months.
4. After query execution completes, in the Results pane, click the **Messages** tab. Note the CPU time and elapsed time for the query. Note the logical reads from each table. The number from the **FactOrder** table should be considerably higher than the number from the dimension tables.
5. Click the **Execution plan** tab, which shows a visualization of the steps the query optimizer used to execute the query. Scroll to the right-hand side and to the bottom, and note that a table scan was used to read data from the **FactOrder** table. Hold the mouse pointer over each of the **Index Scan** icons for the dimension tables to see which indexes were used.

Create Indexes on a Fact Table

1. Select the code under the comment **Create traditional indexes on the fact table**, and then click **Execute**. This creates a clustered index on the date dimension key, and nonclustered indexes on the other dimension keys (the operation can take a long time).
2. Select the code under the comment **Empty the cache**, and then click **Execute**. This clears any cached data.
3. Select the code under the comment **Test the traditional indexes**, and then click **Execute**. This executes the same query as earlier.
4. Click the **Messages** tab and compare the number of logical reads for the **FactOrders** table and the CPU and elapsed time values with the previous execution. They should all be lower.
5. Click the **Execution Plan** tab and note that the clustered index on the date key in the fact table was used.

Create a Columnstore Index

1. Select the code under the comment **Create a copy of the fact table with no indexes**, and then click **Execute**. This creates an unindexed copy of the **FactOrder** table named **FactOrderCS**.
2. Select the code under the comment **Create a columnstore index on the copied table**, and then click **Execute**. This creates a columnstore index on all columns in the **FactOrderCS** table.

3. Select the code under the comment **Empty the cache again**, and then click **Execute**. This clears any cached data.
4. Select the code under the comment **Test the columnstore index**, and then click **Execute**. This executes the same query as earlier.
5. Click the **Messages** tab and compare the number of logical reads for the **FactOrdersCS** table and the CPU and elapsed time values with the previous execution. They should all be lower.
6. Click the **Execution Plan** tab and note that the columnstore index on the fact table was used.
7. Close Indexes.sql but keep SQL Server Management Studio open for the next demonstration.

Managing Many-to-Many Relationships

Relationships in a data warehouse are usually many-to-one, with one row in the dimension table related to many rows in the fact table. For example, a single customer may have many orders. However, it is possible that a row in a fact table can be related to more than one record in a dimension table. A bank has joint accounts that belong to more than one customer, and customers can hold more than one account. Another example would be that of two realtor agents working together on the sale of one house. This is also known as a multivalued dimension.

There are several available methods for dealing with such relationships, and it might be that you have a combination in your data warehouse to manage the different types of data and relationship. No single method is regarded as better than another—rather, the method you choose will depend on the data you have and how this will be queried. The following solutions are available for managing many-to-many relationships:

Add a bridge table

Bridge tables are frequently used for handling many-to-many relationships, sitting between the fact and dimension table. In the example of two people making one sale, the bridge table would hold the key related to both salespeople, such as SalesOrderID, and the SalesPersonID. There would be one row in the bridge table for each salesperson. The bridge table can include a weighting column, enabling you to distribute the value of the sale between the two salespeople. For example, Salesperson A may take 75 percent of the commission, with the remaining 25 percent allocated to Salesperson B, because Salesperson A did the bulk of the work to achieve the sale. In this case, a weighting column can be added with the respect values of 75 and 25. When calculating sales, this means you can view the proportion of commission by salesperson.

Elect a primary value

A common solution is to select one value to be the main value, and ignore the rest. While this eliminates the problem of multivalues, the model may not be useful without the extra values. Furthermore, determining the business rules to decide which value will be the primary might be challenging, or even impossible.

- Data warehouse relationships are usually one-to-many, with one row in the dimension table related to many rows in the fact table
- Known as multivalued dimensions
- Many-to-many relationships can arise, such as two salespeople working on the same sales order
- A number of solutions are available:
 - Add a bridge table between the fact and dimension tables
 - Elect one of the values to be the primary value
 - Include multiple named attributes in the dimension
 - Change the grain of the table to allow multiple rows for a fact
 - Concatenate multiple values into one, and split during querying

Include multiple named attributes

The dimension table can be expanded to include multiple named attributes, but this is only suitable when you have a constant number of options.

Alter the grain of the fact table

You can alter the grain of the fact table if this fits with your data. Returning to the example of two salespeople selling one house, you could have a row in the fact table for each salesperson, with the amount of the sale split proportionally between the two rows.

Concatenate multiple values

Multiple attributes can be concatenated into a single value with a delimiter, such as the pipe character or a backslash, before and after each value. The values are easily displayed in the end reporting, although the ability to query such a string is not straightforward. To find matching values in your queries, you would need to perform wildcard searches using LIKE or CONTAINS predicates, which are unlikely to return precise results and also suffer from slow performance. This method also makes aggregations difficult to calculate on one of the concatenated values.

Data Compression

SQL Server Enterprise edition supports data compression at both page and row level. Row compression stores all fields in a variable width format and, if possible, reduces the number of bytes used to store each field. Page compression applies the same technique to rows on a page and also identifies redundant values, storing them only once per page. You can apply compression to a table, an index, or a partition.

Data compression in a data warehouse brings the following benefits:

- **Reduced storage requirements.** Although results vary, on average, most data warehouses can be compressed at a ratio of 3.5:1, reducing the amount of disk space required to host the data files by more than two thirds.
- **Improved query performance.** Compression can improve query performance in two ways: fewer pages must be read from disk, so I/O is reduced, and more data can be stored on a page and cached.

When page or row compression is used, data must be compressed and decompressed by the CPU. Performance gains resulting from compression must be balanced by the increase in CPU workload. However, in most adequately specified data warehouse servers, the additional workload on CPU is minimal compared to the benefits gained by compressing the data.



Best Practice: When planning tables, partitions, and indexes in a data warehouse, consider the following best practices for data compression:

- Use page compression on all dimension tables and fact table partitions.
- If performance is CPU-bound, revert to row compression on frequently-accessed partitions.

- Apply page compression on all dimension tables, indexes, and fact table partitions
- If performance becomes CPU-bound, fall back to row compression on the most queried partitions

Demonstration: Implementing Data Compression

In this demonstration, you will see how to:

- Create uncompressed tables and indexes.
- Estimate compression savings.
- Create compressed tables and indexes.
- Compare query performance.

Demonstration Steps

Create Uncompressed Tables and Indexes

1. Ensure that you have completed the previous demonstrations in this module.
2. Use Windows Explorer to view the contents of the **D:\Demofiles\Mod03** folder, set the folder window to **Details** view and resize it if necessary, so that you can see the **Size** column.
3. In SQL Server Management, on the **File** menu, point to **Open**, and then click **File**. In the **Open File** dialog box, navigate to the **D:\Demofiles\Mod03** folder, click **Compression.sql**, and then click **Open**.
4. Select the code under the comment **Create the data warehouse** (from line 2 to line 113 in the script), and then click **Execute**. This creates a database with uncompressed tables.
5. While the script is still executing, view the contents of the **D:\Demofiles\Mod03** folder and note the increasing size of **DemoDW.mdf**. This is the data file for the database.
6. When execution is complete (after approximately three minutes), note the final size of **DemoDW.mdf**.

Estimate Compression Savings

1. Return to SQL Server Management Studio.
2. Select the code under the comment **Estimate size saving** (line 119 in the script), and then click **Execute**. This uses the **sp_estimate_data_compression_savings** system stored procedure to compress a sample of the **FactOrder** table (which consists of one clustered and two nonclustered indexes).
3. View the results returned by the stored procedure, noting the current size and estimated compressed size of each index.

Create Compressed Tables and Indexes

1. Select the code under the comment **Create a compressed version of the data warehouse** (from line 125 to line 250 in the script), and then click **Execute**. This creates a database with compressed tables and indexes.
2. While the script is still executing, view the contents of the **D:\Demofiles\Mod03** folder and note the increasing size of **CompressedDemoDW.mdf**. This is the data file for the database.
3. When execution is complete (after approximately one minute), compare the final size of **CompressedDemoDW.mdf** with **DemoDW.mdf** (the file for the compressed database should be significantly smaller).

Compare Query Performance

1. Return to SQL Server Management Studio.
2. Select the code under the comment **Compare query performance** (from line 255 to line 277 in the script), and then click **Execute**. This executes an identical query in the compressed and uncompressed databases and displays execution statistics.
3. When execution is complete, click the **Messages** tab and compare the statistics for the two queries. The execution time statistics (the second and third set of figures labeled "SQL Server Execution Time") should be similar, but the second query (in the compressed database) should have used considerably fewer logical reads for each table than the first query.
4. Close SQL Server Management Studio.

Using Views to Abstract Base Tables

You can create views in a data warehouse to abstract the underlying fact and dimension tables. Although views are not always necessary, you should consider the following guidelines when planning a data warehouse:

- **Create a view for each dimension and fact table, and use the NOLOCK query hint in the view definition.** You can then use these views instead of the base tables for all data access from clients, which will eliminate locking overhead and optimize concurrency.
- **Create views with user-friendly view and column names.** Often, a naming convention (such as prefixing dimension tables with "dim" and fact tables with "fact") is used when creating the tables in a data warehouse. Such naming conventions are useful for database designers and administrators, but they can confuse business users. Creating a layer of views with accessible names makes it easier for users to create their own data models and reports from the data warehouse.
- **Do not include metadata columns in views.** Some columns are used for ETL operations or other administrative tasks, and can be omitted from views that will be consumed by business users. For example, slowly changing dimension current flag, start date, and end date columns might not be required for end user reporting or data models, so you can create views that do not include them.
- **Create views to combine snowflake dimension tables.** If you have included snowflake dimensions in your dimensional model, create a view for each set of related dimension tables to produce a single logical dimension table.
- **Partition-align indexed views.** SQL Server supports indexed views, which can be partitioned using the same partition scheme as the underlying table. If you use indexed views, you should partition-align them to support partition switching that does not require indexes on the views to be dropped and recreated.
- **Use the SCHEMABINDING option.** This ensures that the underlying tables cannot be dropped or modified in such a way as to invalidate the view, unless the view itself is dropped first. The SCHEMABINDING option is a requirement for index views.

```
CREATE VIEW dw_views.SalesOrder
WITH SCHEMABINDING
AS
SELECT [OrderDateKey]
      ,[ProductKey]
      ,[ShipDateKey]
      ,[CustomerKey]
      ,[OrderNumber]
      ,[OrderQuantity]
      ,[UnitPrice]
      ,[SalesAmount]
FROM [dbo].[FactSalesOrder]
WITH (NOLOCK)
```

Question: List three reasons for partitioning a large table.

Lab: Implementing a Data Warehouse

Scenario

You have gathered analytical and reporting requirements from stakeholders at Adventure Works Cycles. Now you must implement a data warehouse schema to support them.

Objectives

After completing this lab, you will be able to:

- Implement a dimensional star schema.
- Implement a snowflake schema.
- Implement a time dimension.

Estimated Time: 45 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Implementing a Star Schema

Scenario

Adventure Works Cycles requires a data warehouse that helps information workers and executives to create reports and perform analysis of key business measures. The company has identified two sets of related measures that it wants to include in fact tables. These are separate sales order measures relating to sales to resellers, and Internet sales. The measures will be aggregated by product, reseller (in the case of reseller sales), and customer (for Internet sales) dimensions.

The data warehouse has been partially completed; you must now add the necessary dimension and fact tables to complete a star schema.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. View a Data Warehouse Schema
3. Create a Dimension Table
4. Create a Fact Table
5. View the Revised Data Warehouse Schema

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab03\Starter** folder as Administrator.

► **Task 2: View a Data Warehouse Schema**

1. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the SQL Server database engine using Windows authentication.
2. Create a new database diagram in the **AWDataWarehouse** database (creating the required objects to support database diagrams if prompted). The diagram should include all of the tables in the database.
3. In the database diagram, modify the tables so they are shown in standard view, and arrange them so you can view the partially complete data warehouse schema.
4. Save the database diagram as **AWDataWarehouse Schema**.

► **Task 3: Create a Dimension Table**

1. Review the Transact-SQL code in the **DimCustomer.sql** query file in the **D:\Labfiles\Lab03\Starter** folder. Note that it creates a table named **DimCustomer** in the **AWDataWarehouse** database.
2. Execute the query to create the **DimCustomer** dimension table.

► **Task 4: Create a Fact Table**

1. Review the Transact-SQL code in the **FactInternetSales.sql** query file in the **D:\Labfiles\Lab03\Starter** folder. Note that it creates a table named **FactInternetSales** in the **AWDataWarehouse** database, and that this table is related to the **DimCustomer** and **DimProduct** tables.
2. Execute the query to create the **FactInternetSales** dimension table.

► **Task 5: View the Revised Data Warehouse Schema**

1. Add the tables that you have created in this exercise to the database diagram that you created.
Note: when adding tables to a diagram, you need to click **Refresh** in the **Add Table** dialog box to see tables that you have created or modified since the diagram was initially created.
2. Save the database diagram.
3. Keep SQL Server Management Studio open for the next exercise.

Results: After this exercise, you should have a database diagram in the **AWDataWarehouse** database that shows a star schema consisting of two fact tables (**FactResellerSales** and **FactInternetSales**) and four dimension tables (**DimReseller**, **DimEmployee**, **DimProduct**, and **DimCustomer**).

Exercise 2: Implementing a Snowflake Schema

Scenario

Having created a star schema, you have identified two dimensions that would benefit from being normalized to create a snowflake schema. Specifically, you want to create a hierarchy of related tables for product category, product subcategory, and product. You also want to create a separate geography dimension table that can be shared between the reseller and customer dimensions.

The main tasks for this exercise are as follows:

1. Create Dimension Tables That Form a Hierarchy
2. Create a Shared Dimension Table
3. View the Data Warehouse Schema

► Task 1: Create Dimension Tables That Form a Hierarchy

1. Review the Transact-SQL code in the **DimProductCategory.sql** query file in the **D:\Labfiles\Lab03\Starter** folder. Note that it:
 - Creates a table named **DimProductCategory**.
 - Creates a table named **DimProductSubcategory** that has a foreign key relationship to the **DimProductCategory** table.
 - Drops the **ProductSubcategoryName** and **ProductCategoryName** columns from the **DimProduct** table.
 - Adds a **ProductSubcategoryKey** column to the **DimProduct** table that has a foreign key relationship to the **DimProductSubcategory** table.
2. Execute the query.

► Task 2: Create a Shared Dimension Table

1. Review the Transact-SQL code in the **DimGeography.sql** query file in the **D:\Labfiles\Lab03\Starter** folder. Note that it:
 - Creates a table named **DimGeography**.
 - Drops the **City**, **StateProvinceName**, **CountryRegionCode**, **CountryRegionName**, and **PostalCode** columns from the **DimReseller** table.
 - Adds a **GeographyKey** column to the **DimReseller** table that has a foreign key relationship to the **DimGeography** table.
 - Drops the **City**, **StateProvinceName**, **CountryRegionCode**, **CountryRegionName**, and **PostalCode** columns from the **DimCustomer** table.
 - Adds a **GeographyKey** column to the **DimCustomer** table that has a foreign key relationship to the **DimGeography** table.
2. Execute the query.

► Task 3: View the Data Warehouse Schema

1. Delete the tables that you modified in the previous two tasks from the **AWDataWarehouse** Schema diagram (**DimProduct**, **DimReseller**, and **DimCustomer**).
2. Add the new and modified tables that you created in this exercise to the **AWDataWarehouse Schema** diagram and view the revised data warehouse schema, which now includes some snowflake dimensions. You will need to refresh the list of tables when adding tables, and you might be prompted to update the diagram to reflect new foreign key relationships.
3. Save the database diagram.

Results: After this exercise, you should have a database diagram in the **AWDataWarehouse** database showing a snowflake schema that contains a dimension consisting of a **DimProduct**, **DimProductSubcategory**, and **DimProductCategory** hierarchy of tables, in addition to a **DimGeography** dimension table that is referenced by the **DimCustomer** and **DimReseller** dimension tables.

Exercise 3: Implementing a Time Dimension Table

Scenario

The schema for the Adventure Works data warehouse now contains two fact tables and several dimension tables. However, users need to be able to analyze the fact table measures across consistent time periods. To make this happen, you must create a time dimension table.

Users should be able to aggregate measures across calendar years (which run from January to December) and fiscal years (which run from July to June). Your time dimension must include the following attributes:

- Date (this should be the business key).
- Day number of week (for example, 1 for Sunday, 2 for Monday).
- Day name of week (for example, Sunday, Monday, Tuesday).
- Day number of month.
- Day number of year.
- Week number of year.
- Month name (for example, January, February).
- Month number of year (for example, 1 for January, 2 for February).
- Calendar quarter (for example, 1 for dates in January, February, and March).
- Calendar year.
- Calendar semester (for example, 1 for dates between January and June).
- Fiscal quarter (for example, 1 for dates in July, August, and September).
- Fiscal year.
- Fiscal semester (for example, 1 for dates between July and December).

The main tasks for this exercise are as follows:

1. Create a Time Dimension Table
2. View the Database Schema
3. Populate the Time Dimension Table

► **Task 1: Create a Time Dimension Table**

1. Review the Transact-SQL code in the **DimDate.sql** query file in the D:\Labfiles\Lab03\Starter folder. Note that it:
 - Creates a table named **DimDate**.
 - Adds **OrderDateKey** and **ShipDateKey** columns to the **FactInternetSales** and **FactResellerSales** tables that have foreign key relationships to the **DimDate** table.
 - Creates indexes on the **OrderDateKey** and **ShipDateKey** foreign key fields in the **FactInternetSales** and **FactResellerSales** tables.
2. Execute the query.

► **Task 2: View the Database Schema**

1. Delete the tables that you modified in the previous two tasks from the **AWDataWarehouse Schema** diagram (**FactResellerSales** and **FactInternetSales**).
2. Add the new and modified tables that you created in this exercise to the **AWDataWarehouse Schema** diagram and view the revised data warehouse schema, which now includes a time dimension table named **DimDate**. You will need to refresh the list of tables when adding tables.
3. In Object Explorer, expand the list of tables in the **AWDataWarehouse** database, and refresh the list.
4. Expand the list of columns in the **dbo.FactInternetSales** table and refresh the list. Note that each of the columns, **OrderDateKey** and **ShipDateKey**, has a silver key symbol against it, denoting that each has a foreign key relationship. Note that, whereas the primary keys are also indicated by a key symbol on the table (pointing in the opposite direction) in **Object Explorer** and in the **Database Diagram**, the foreign key relationships are *not* indicated by a symbol on the database diagram. However, a line between the **FactResellerSales** table and the **DimDate** table indicates the presence of a foreign key relationship between the two tables, but does not show which columns are the foreign keys, and does not show that there are two foreign key relationships: **OrderDateKey** and **ShipDateKey** in the table **FactResellerSales** linked to **DateKey** on the table **DimDate**.
5. Save the database diagram.

► **Task 3: Populate the Time Dimension Table**

1. Review the Transact-SQL code in the **GenerateDates.sql** query file in the D:\Labfiles\Lab03\Starter folder. Note that it:
 - Declares a variable named **@StartDate** with the value 1/1/2000, and a variable named **@EndDate** with the value of the current date.
 - Performs a loop to insert appropriate values for each date between **@StartDate** and **@EndDate** into the **DimDate** table.
2. Execute the script.
3. When the query has completed, query the **DimDate** table to verify that it now contains time values.
4. Close Visual Studio®, saving your work if prompted.

Results: After this exercise, you should have a database that contains a **DimDate** dimension table that is populated with date values from January 1, 2000, to the current date.

Module Review and Takeaways

This module has described considerations for translating business requirements and information about business processes into a dimensional model, and then implementing that model as a data warehouse. Every business is different, and each has its unique challenges and processes. You should use the techniques and guidance in this module as a starting point, but be prepared to adapt typical data warehouse schema elements to match particular business requirements.

Review Question(s)

Question: When designing a data warehouse, is it better or worse to have a strong background in transactional database design?

Module 4

Columnstore Indexes

Contents:

Module Overview	4-1
Lesson 1: Introduction to Columnstore Indexes	4-2
Lesson 2: Creating Columnstore Indexes	4-7
Lesson 3: Working with Columnstore Indexes	4-12
Lab: Using Columnstore Indexes	4-17
Module Review and Takeaways	4-21

Module Overview

Columnstore indexes were introduced in Microsoft® SQL Server® 2012, and have been used by many organizations in large data warehouse solutions. This module highlights the benefits of using these indexes in the context of data warehousing, the improvements made in the latest build of SQL Server, and the considerations needed to use columnstore indexes effectively in your solutions.

Objectives

After completing this module, you will be able to:

- Describe columnstore indexes and identify suitable scenarios for their use.
- Create clustered and nonclustered columnstore indexes.
- Describe considerations for using columnstore indexes.

Lesson 1

Introduction to Columnstore Indexes

This lesson gives an overview of the types of columnstore indexes available in SQL Server; the advantages they have over their similar row based indexes; and under what circumstances they should be considered. By the end of this lesson, you will see the potential cost savings to your business of using clustered columnstore indexes, purely in terms of the gigabytes of disk storage available.

Lesson Objectives

After completing this lesson, you will be able to:

- Explain the differences between rowstore and columnstore indexes.
- Describe the properties of a nonclustered columnstore index.
- Describe the properties of a clustered columnstore index and how it differs from a nonclustered columnstore index.

What Are Columnstore Indexes?

Microsoft introduced columnstore indexes in SQL Server 2012. These indexes reference data in a columnar fashion and use compression aggressively to reduce the disk I/O when responding to queries.

Traditional rowstore tables are stored on disk in pages; each page contains a number of rows and includes all the associated columns with each row.

Columnstore indexes also store data in pages, but they store all the column values in a page—so the page consists of the same column of data from multiple rows.

Consider a data warehouse containing fact tables that are used to calculate aggregated data across multiple dimensions. These fact tables might consist of many rows, perhaps numbering 10s of millions.

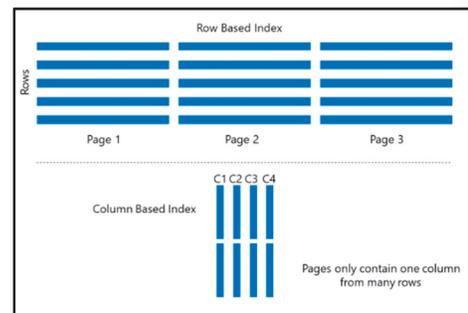
Using a code example:

Totaling Sales Orders by Product

```
SELECT ProductID
       ,SUM(LineTotal) AS 'ProductTotalSales'
FROM Sales.OrderDetail
GROUP BY ProductID
ORDER BY ProductID
```

The previous example, using a row based index, will need to load into memory all the rows and columns in all the pages, for all the products. With a column based index, the query only needs to load the pages associated with the two referenced columns, ProductID and LineTotal. This makes columnstore indexes a good choice for large data sets.

Using a columnstore index can improve the performance for a typical data warehouse query by up to 10 times. There are two key characteristics of columnstore indexes that impact this gain.



- **Storage.** Columnstore indexes store data in a compressed columnar data format instead of by row. This makes it possible to achieve compression ratios of seven times greater than a standard rowstore table.
- **Batch mode execution.** Columnstore indexes process data in batches (of 1,000-row blocks) instead of row by row. Depending on filtering and other factors, a query might also benefit from segment elimination, which involves bypassing million-row chunks (segments) of data and further reducing I/O.

Columnstore indexes perform well because:

- Columns often store matching data—for example, a set of States—enabling the database engine to compress the data better. This compression can reduce or eliminate any I/O bottlenecks in your system, while also reducing the memory footprint.
- High compression rates improve overall query performance because the results have a smaller in-memory footprint.
- Instead of processing individual rows, batch execution also improves query performance. This can typically be a performance improvement of around two to four times because processing takes place on multiple rows simultaneously.
- Aggregation queries often select only a few columns from a table, which also reduces the total I/O required from the physical media.



Note: Nonclustered and clustered indexes are supported in Azure® SQL Database V12 Premium Edition. For a full list of the columnstore features available in different versions of SQL server, see:



Columnstore Indexes Versioned Feature Summary

<http://aka.ms/J0nknq>

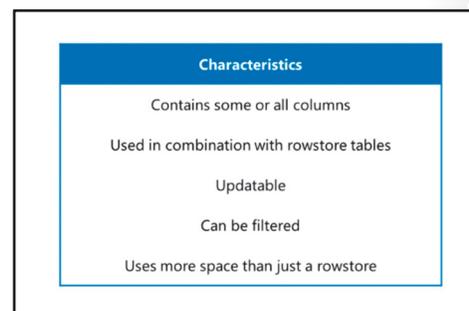
There are two types of columnstore indexes—nonclustered and clustered columnstore indexes—that both function in the same way. The difference is that a nonclustered index will normally be a secondary index created on top of a rowstore table; a clustered columnstore index will be the primary storage for a table.

Nonclustered Columnstore Indexes

Nonclustered columnstore indexes contain a copy of part or all of the columns in an underlying table. This kind of index represents a copy of the data; therefore, one of the disadvantages is that it will take up more space than using a rowstore table alone.

A nonclustered columnstore index has the following characteristics:

- It can include some or all of the columns in the table.
- It can be combined with other rowstore indexes on the same table.



- It is a full or partial copy of the data and takes more disk space than a rowstore table.

New Features for SQL Server

In SQL Server 2014, if a table has a nonclustered columnstore index, you cannot edit its contents without dropping and recreating the index—in subsequent builds of SQL Server, this restriction has been removed.

Also new in later releases of SQL Server is support for filtered nonclustered columnstore indexes. You can use a predicate condition to filter the rows that are included in the index; use this feature to create an index on only the cold data of an operational workload. This will greatly reduce the performance impact of having a columnstore index on an online transaction processing (OLTP) table.

Clustered Columnstore Indexes

Clustered columnstore indexes, like their rowstore alternatives, arrange the physical data on disk or in memory optimized for the index. In a columnstore index, this will store all the columns next to each other on disk; the structure of the index is used to store the data on disk.

To reduce the impact of fragmentation and to improve performance, a clustered columnstore index might use a deltastore. You can think of a deltastore as a temporary b-tree table with rows that a tuple-mover process moves into the clustered columnstore index at an appropriate time. This moving of row data is performed in the background. When querying the index, it will automatically combine results from the columnstore and deltastore to ensure that the query receives the correct results.

A clustered columnstore index has the following characteristics:

- It includes all of the columns in the table.
- It does not store the columns in a sorted order, but optimizes storage for compression and performance.
- It can be updated.

New Features for SQL Server

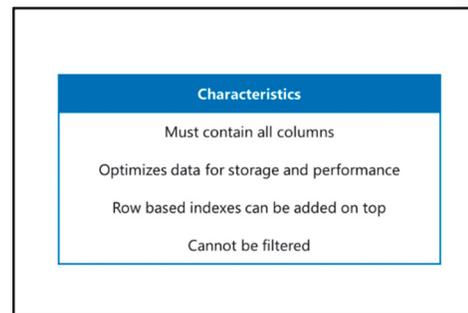
You can now have a nonclustered row index on top of a clustered columnstore index, so it's possible to have efficient table seeks on an underlying columnstore. Using a unique rowstore index, you can also enforce a primary key constraint.

In the latest builds of SQL Server, you can use columnstore indexes on memory optimized tables—the most relevant use being for real-time operational analytical processing.

For further information, see:

 **Columnstore Indexes for Real-Time Operational Analytics**

<https://aka.ms/mtq5h7>



Demonstration: The Benefits of Using Columnstore Indexes

In this demonstration, you will see how to create a columnstore index.

Demonstration Steps

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **D:\Demofiles\Mod04\Setup.cmd** as an administrator to revert any changes.
3. On the taskbar, click **SQL Server Management Studio**.
4. In the Connect to Server window, in the **Server name** box, type **MIA-SQL**. Ensure **Windows Authentication** is selected in the **Authentication** box.
5. Click **Connect**.
6. On the **File** menu, point to **Open**, click **File**, navigate to the **D:\Demofiles\Mod04\Demo** folder, click the **Demo.sql** script file, and then click **Open**.
7. Select the code and review the comments under **Step 1**, and then click **Execute**.
8. Select the code and review the comments under **Step 2**, and then click **Execute**.
9. Select the code under **Step 1** again, and then click **Execute**.
10. Select the code and review the comments under **Step 3**, and then click **Execute**.
11. Select the code and review the comments under the comment **Data space used**, and then click **Execute**.
12. Close SQL Server Management Studio without saving changes.

Categorize Activity

Categorize each index property into the appropriate index type. Indicate your answer by writing the category number to the right of each property.

Items	
1	Perform the best when seeking for specific data
2	A high degree of compression is possible, due to data being of the same category
3	Can improve the performance of database queries
4	Implemented as a b-tree index structure
5	Perform best when aggregating data
6	Can be stored in memory optimized tables

Category 1	Category 2	Category 3
Rowstore Index	Columnstore Index	Applies to both types of index

Lesson 2

Creating Columnstore Indexes

This lesson shows you the techniques required to create both types of columnstore indexes on a table in a SQL database. You will see how to quickly create indexes by using Transact-SQL or the SQL Server Management Studio user interface.

Lesson Objectives

After completing this lesson, you will be able to:

- Create nonclustered columnstore indexes.
- Create clustered columnstore indexes.
- Create tables that utilize both columnstore and rowstore indexes.

Creating a Nonclustered Columnstore Index

You can create a nonclustered columnstore index by using either a Transact-SQL statement or SQL Server Management Studio (SSMS).

Transact-SQL

To create a nonclustered columnstore index, use the **CREATE NONCLUSTERED COLUMNSTORE INDEX** statement as shown in the following code example:

Creating a Nonclustered Columnstore Index

```
CREATE NONCLUSTERED COLUMNSTORE INDEX
NCCSIX_FactInternetSales
ON FactInternetSales (
    Quantity
    ,Cost
    ,Discount);
```

A nonclustered index does not need to include all the columns from the underlying table. In the preceding example, only three columns are included in the index.

You can also restrict nonclustered indexes to a subset of the rows contained in a table:

Example of a Filtered Nonclustered Columnstore Index

```
CREATE NONCLUSTERED COLUMNSTORE INDEX NCCSIX_FactInternetSalesShippedOrders
ON FactInternetSales (
    Quantity
    ,Cost
    ,Discount
)
WHERE ShipDate < '2013-01-01';
```



The business reason for wanting to limit a columnstore index to a subset of rows is that you can use a single table for both OLTP and analytical processing. In the preceding example, the index supports analytical processing on historical orders that shipped before 2013.

SQL Server Management Studio

You can also use Object Explorer in SSMS to create columnstore indexes:

1. In **Object Explorer**, expand **Databases**.
2. Expand the required database; for example **AdventureWorksDW**.
3. Expand **Tables**, and then expand the required table; for example **FactFinance**.
4. Right-click **Indexes**, point to **New Index**, and then click **Non-Clustered Columnstore Index**.
5. Add at least one column to the index, and then click **OK**.

Creating a Clustered Columnstore Index

Similar to nonclustered columnstore indexes, a clustered columnstore index is created by using Transact-SQL or SSMS. The main difference between the declaration of a clustered index and nonclustered index is that the former must contain all the columns in the table being indexed.

Transact-SQL

To create a clustered columnstore index, use the `CREATE CLUSTERED COLUMNSTORE INDEX` statement as shown in the following code example:

Creating a Clustered Columnstore Index

```
CREATE CLUSTERED COLUMNSTORE INDEX CCSIX_FactSalesOrderDetails
ON FactSalesOrderDetails;
```

An optional parameter on a `CREATE` statement for a clustered index is `DROP_EXISTING`. You can use this to rebuild an existing clustered columnstore index or to convert an existing rowstore table into a columnstore table.

 **Note:** To use the `DROP_EXISTING` option, the new columnstore index must have the same name as the index it is replacing.



The following example creates a clustered rowstore table, and then converts it into a clustered columnstore table:

Converting a Rowstore Table to a Columnstore Table

```
CREATE TABLE ExampleFactTable (
    ProductKey [int] NOT NULL,
    OrderDateKey [int] NOT NULL,
    DueDateKey [int] NOT NULL,
    ShipDateKey [int],
    CostPrice [money] NOT NULL);
GO
CREATE CLUSTERED INDEX CCI_ExampleFactTable ON ExampleFactTable (ProductKey);
GO
CREATE CLUSTERED COLUMNSTORE INDEX CCI_ExampleFactTable ON ExampleFactTable
WITH (DROP_EXISTING = ON);
GO
```

SQL Server Management Studio

You can also create a clustered columnstore index by using SSMS:

1. In **Object Explorer**, expand **Databases**.
2. Expand the required **Database**, for example **AdventureWorksDW**.
3. Expand **Tables**, and then expand the required table, for example **FactFinance**.
4. Right-click **Indexes**, point to **New Index**, and then click **Clustered Columnstore Index**.
5. Click **OK**.

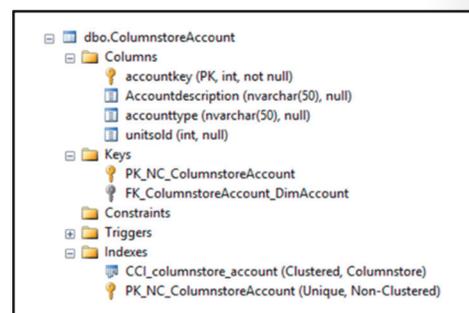


Note: You don't need to select columns to create a clustered columnstore index, because all the columns of a table must be included.

Creating a Clustered Columnstore Table with Primary and Foreign Keys

In SQL Server 2014, clustered columnstore indexes are limited in several ways, one of the most significant being that you cannot have any other index on the clustered columnstore table. This essentially means that, in SQL Server 2014, you cannot have a primary key, foreign keys, or unique value constraints, on a clustered columnstore table.

These limitations have since been removed— SQL Server now supports all these features.



This is an example of creating a table with both a primary key and a clustered columnstore index:

Create a Table with a Primary Key and Columnstore Index

```
CREATE TABLE ColumnstoreAccount (
    accountkey int NOT NULL,
    Accountdescription nvarchar (50),
    accounttype nvarchar(50),
    unitsold int,
    CONSTRAINT [PK_NC_ColumnstoreAccount] PRIMARY KEY NONCLUSTERED([accountkey] ASC),
    INDEX CCI_columnstore_account CLUSTERED COLUMNSTORE)
```

After you create the table, you can add a foreign key constraint:

Add a Foreign Key Constraint

```
ALTER TABLE ColumnstoreAccount WITH CHECK ADD CONSTRAINT FK_ColumnstoreAccount_DimAccount
FOREIGN KEY(AccountKey)
REFERENCES DimAccount (AccountKey)

ALTER TABLE ColumnstoreAccount CHECK CONSTRAINT FK_ColumnstoreAccount_DimAccount
GO
```

Checking the table shows that two indexes and two keys exist:

- **CCI_columnstore_account:** a clustered columnstore index.
- **PK_NC_ColumnstoreAccount:** a unique nonclustered rowstore index.
- **FK_ColumnstoreAccount_DimAccount:** a foreign key to the DimAccount table.
- **PK_NC_ColumnstoreAccount:** the primary key.

The previous Transact-SQL results in a columnstore index with a nonclustered index that enforces a primary key constraint on both indexes.

Demonstration: Creating Columnstore Indexes Using SQL Server Management Studio

In this demonstration, you will see how to:

- Create a nonclustered columnstore index using SSMS.
- Create a clustered columnstore index using SSMS.

Demonstration Steps

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. On the taskbar, click **SQL Server Management Studio**.
3. In the Connect to Server window, in the **Server name** box, type **MIA-SQL**. Ensure **Windows Authentication** is selected in the **Authentication** box.
4. Click **Connect**.
5. In Object Explorer, expand **Databases**, expand **AdventureWorksDW**, expand **Tables**, and then expand **dbo.AdventureWorksDWBuildVersion**.

6. Right-click **Indexes**, point to **New Index**, and then click **Clustered Columnstore Index**.
7. In the **New Index** dialog box, click **OK** to create the index.
8. In Object Explorer, expand **Indexes** to show the new clustered index.
9. In Object Explorer, expand **dbo.FactResellersSales**.
10. Right-click **Indexes**, point to **New Index**, and then click **Non-Clustered Columnstore Index**.
11. In the **Columnstore Columns** section of the **New Index** dialog box, click **Add**.
12. Select the following columns: **SalesOrderNumber**, **UnitPrice**, and **ExtendedAmount**, and then click **OK**.
13. In the **New Index** dialog box, click **OK**.
14. In Object Explorer, expand **Indexes** to show the new nonclustered index.
15. Close SQL Server Management Studio without saving changes.

Question: How will you create your indexes in a database—with SSMS or Transact-SQL?

Lesson 3

Working with Columnstore Indexes

When working with columnstore indexes, you should consider fragmentation and how SQL Server processes the insertion of data into the index. A new feature in SQL Server is to store columnstore tables in memory, meaning you can perform real-time operational analytics.

Lesson Objectives

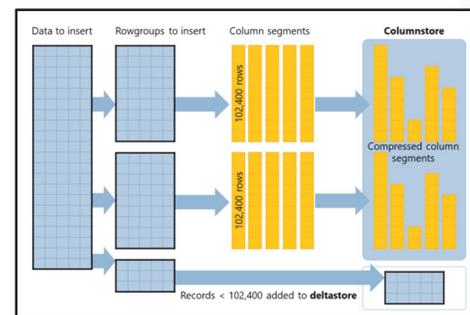
After completing this lesson, you will be able to:

- Efficiently add data into a columnstore table.
- Check the fragmentation of an index and choose the best approach to resolving the fragmentation.
- Create a memory optimized table to support real-time operational analytics.

Managing Columnstore Indexes

Management considerations for columnstore indexes are similar to those for rowstore indexes—however, special consideration must be given to DML operations.

SQL Server manages a deltastore so that data can be manipulated and changed in a columnstore table. The deltastore collects up to 1,048,576 rows before compressing them into a compressed rowgroup—and then marking that rowgroup as closed. The tuple-mover background process then adds the closed rowgroup back into the columnstore index.



Bypassing the Deltastore

To ensure that data is inserted directly into the columnstore, you should load the data in batches of between 102,400 and 1,048,576 rows. You can perform this bulk data loading using normal insertion methods, including the bcp utility, SQL Server Integration Services, and Transact-SQL insert statements from a staging table.

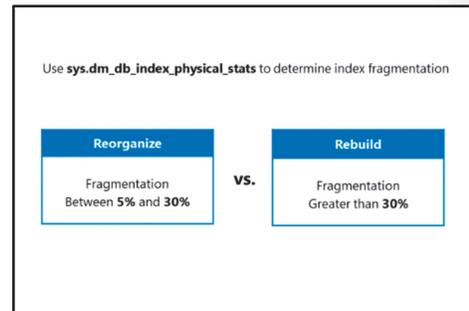
When bulk loading data, you have the following options for optimizations:

- **Parallel Load:** perform multiple concurrent bulk imports (bcp or bulk insert), each loading separate data.
- **Log Optimization:** the bulk load will be minimally logged when the data is loaded into a compressed rowgroup. Minimal logging is not available when loading data with a batch size of less than 102,400 rows.

Index Fragmentation

When it comes to managing, columnstore indexes are no different to rowstore indexes. The utilities and techniques used to keep an index healthy are the same.

Over time, and after numerous DML operations, indexes can become fragmented and their performance degrades. SSMS provides Transact-SQL commands and user interface elements to help you manage this fragmentation.



SQL Server Management Studio

You can examine the level of fragmentation by following these steps:

1. In **Object Explorer**, expand **Databases**.
2. Expand the required **Database**; for example, **AdventureWorksDW**.
3. Expand **Tables**, and then expand the required table; for example, **FactFinance**.
4. Expand **Indexes**, right-click the desired index; in the context menu, click **Properties**.
5. In the **Select a page** panel, click **Fragmentation**.

Transact-SQL

SQL Server provides dynamic management views and functions that help a database administrator to inspect and review the health of indexes.

One of these functions is **sys.dm_db_index_physical_stats** that can be run against all the databases on a server, a specific table in a database, or even a specific index.

The following code sample shows a useful query that joins the results from the **sys.dm_db_index_physical_stats** view with the system index table, and then returns the fragmentation and names of the indexes for a specific database:

Show Indexes with Fragmentation Greater Than Five Percent for a Specific Database

```
SELECT DB_NAME(database_id) AS 'Database'
      , OBJECT_NAME(dm.object_id) AS 'Table'
      , si.name AS 'Index'
      , dm.index_type_desc
      , dm.avg_fragmentation_in_percent
FROM sys.dm_db_index_physical_stats
     (DB_ID(N'AdventureworksDW'), NULL, NULL, NULL, 'DETAILED') dm
JOIN sys.indexes si ON si.object_id = dm.object_id AND si.index_id = dm.index_id
WHERE avg_fragmentation_in_percent > 5
ORDER BY avg_fragmentation_in_percent DESC;
```

When you identify that an index requires maintenance, two options are available in SQL Server: you can either rebuild or reorganize it. Previous guidance was:

- If the fragmentation is between five percent and 30 percent: **Reorganize**.
- If the fragmentation is greater than 30 percent: **Rebuild**.

The reorganizing of columnstore indexes has been enhanced in SQL Server, and it is now rarely necessary to rebuild an index.

Use the following Transact-SQL to reorganize an index:

Transact-SQL to Reorganize an Index Online

```
ALTER INDEX CCI_columnstore_account ON ColumnstoreAccount
REORGANIZE WITH (COMPRESS_ALL_ROW_GROUPS = ON)
GO

ALTER INDEX CCI_columnstore_account ON ColumnstoreAccount REORGANIZE
GO
```

The first statement in this code sample adds deltastore rowgroups into the columnstore index. Using the COMPRESS_ALL_ROW_GROUPS option forces all open and closed rowgroups into the index, in a similar way to rebuilding an index. After the query adds these deltastore rowgroups to the columnstore, the second statement then combines these, possibly smaller, rowgroups into one or more larger rowgroups. With a large number of smaller rowgroups, performing the reorganization a second time will improve the performance of queries against the index. Using these statements in SQL Server means that, in most situations, you no longer need to rebuild a columnstore index.



Note: Rebuilding an index will mean SQL Server can move the data in the index between segments to achieve better overall compression. If a large number of rows have been deleted, and the index fragmentation is more than 30 percent, then rebuilding the index may be the best option, rather than reorganizing.



For more information on all the available views and functions, see:

<http://aka.ms/V4odeb>

Columnstore Indexes and Memory Optimized Tables

The latest build of SQL Server has extended the in-memory capabilities of SQL Server 2014. You can now have both a clustered columnstore index and rowstore index on an in-memory table. Having the combination of these indexes on a table means you can perform real-time operational analytics and all of its associated benefits, including:

- **Simplicity.** No need to implement an ETL process to move data into analytical tables. Analytics can run directly on the operational data.
- **Reduced costs.** Removes the need to develop and support ETL processes. The associated hardware and infrastructure to support the ETL is no longer necessary.
- **Zero data latency.** Data is analyzed in real time. There are no background or schedule processes for moving data, so real-time analytics can be completed.

In-memory columnstore tables

- The index has to be declared at runtime
- Tables can be up to 2 TB in size
- Can be combined with rowstore index
- Enable real-time operational analytics

Memory Optimization Advisor can be used to support moving a table from being disk-based to being memory-optimized.

The combination of indexes means that analytical queries can run against the columnstore index and OLTP operations can run against the OLTP b-tree indexes. The OLTP workloads will continue to perform, but you may incur some additional overhead when maintaining the columnstore index.

 **For more information on Real-Time Operational Analytics, see:**

<http://aka.ms/l3nctd>

As with other similar in-memory tables, you must declare the indexes on memory optimized columnstore tables at creation. To support larger datasets—for example, those used in data warehouses—the size of in-memory tables has been increased in SQL Server, from a previous limit of 256 GB to 2 TB.

The Transact-SQL to create an in-memory table is simple. Add WITH (MEMORY_OPTIMIZED = ON) at the end of a table declaration.

Transact-SQL to Create a Columnstore In-Memory Table.

```
CREATE TABLE InMemoryAccount (
    accountkey int NOT NULL,
    Accountdescription nvarchar (50),
    accounttype nvarchar(50),
    unitsoId int,
    CONSTRAINT [PK_NC_InMemoryAccount] PRIMARY KEY NONCLUSTERED([accountkey] ASC),
    INDEX CCI_InMemoryAccount CLUSTERED COLUMNSTORE)
    WITH (MEMORY_OPTIMIZED = ON, DURABILITY = SCHEMA_AND_DATA)
GO
```

 **Note:** The above Transact-SQL will not work on a database without a memory optimized file group. Before creating any memory optimized tables, the database must have a memory optimized file group associated with it.

The following is an example of the code required to create a memory optimized filegroup:

Creating a Memory Optimized Filegroup

```
ALTER DATABASE AdventureWorksDW
ADD FILEGROUP AdventureWorksDW_Memory_Optimized_Data CONTAINS MEMORY_OPTIMIZED_DATA
GO

ALTER DATABASE AdventureWorksDW ADD
FILE (name='AdventureworksDW_MOD', filename='D:\AdventureworksDW_MOD')
TO FILEGROUP AdventureWorksDW_Memory_Optimized_Data
GO
```

You can alter and join this new in-memory table in the same way as its disk-based counterpart. However, you should use caution when altering an in-memory table, because this is an offline task. There should be twice the memory available to store the current table—a temporary table is used before being switched over when it has been rebuilt.

Depending on the performance requirements, you can control the durability of the table by using:

- **SCHEMA_ONLY:** creates a nondurable table.
- **SCHEMA_AND_DATA:** creates a durable table; this is the default if no option is supplied.

If SQL Server restarts, a schema-only table loses all its data. Temporary tables, or transient data, are examples of where you might use a nondurable table. Because SCHEMA_ONLY durability avoids both transaction logging and checkpoints, I/O operations can be reduced significantly.

Introduced in SQL Server 2014, the Memory Optimization Advisor is a GUI tool inside SSMS. The tool can analyze an existing disk-based table and warn if there are any features of that table—for example, an unsupported type of index—that aren't possible on a memory optimized table. It can then migrate the data contained in the disk-based table to a new memory optimized table. The Memory Optimization Advisor is available on the context menu of any table in Management Studio.

Check Your Knowledge

Question	
When would you consider converting a rowstore table, containing dimension data in a data warehouse, to a columnstore table?	
Select the correct answer.	
<input type="checkbox"/>	When mission critical analytical queries join one or more fact tables to the dimension table—and those fact tables are columnstore tables.
<input type="checkbox"/>	When the data contained in the dimension table has a high degree of randomness and uniqueness.
<input type="checkbox"/>	When the dimension table has very few rows.
<input type="checkbox"/>	When the dimension table has many millions of rows, with columns containing small variations in data.
<input type="checkbox"/>	It is never appropriate to convert a dimension table to a columnstore table.

Lab: Using Columnstore Indexes

Scenario

Adventure Works has created a data warehouse for analytics processing of its current online sales business. Due to large business growth, existing analytical queries are no longer performing as required. There is also a growing issue of disk space.

You have been tasked with optimizing the existing database workloads and, if possible, reducing the amount of disk space being used by the data warehouse.

Objectives

After completing this lab, you will be able to:

- Create clustered and nonclustered columnstore indexes.
- Examine an execution plan to check the performance of queries.
- Convert disk-based tables into memory optimized tables.

Lab Setup

Estimated Time: 45 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Dropping and recreating indexes can take time, depending on the performance of the lab machines.

Exercise 1: Create a Columnstore Index on the FactProductInventory Table

Scenario

You plan to improve the performance of the **AdventureWorksDW** data warehouse by using columnstore indexes. You need to improve the performance of queries that use the **FactProductInventory** tables without causing any database downtime, or dropping any existing indexes. Disk usage for this table is not an issue.

You must retain the existing indexes on the **FactProductInventory** table, and ensure you do not impact current applications by any alterations you make to the table.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Examine the Existing Size of the FactProductInventory Table and Query Performance
3. Create a Columnstore Index on the FactProductInventory Table

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab04\Starter** folder, run **Setup.cmd** as Administrator.

► Task 2: Examine the Existing Size of the FactProductInventory Table and Query Performance

1. In SQL Server Management Studio, open the **Query FactProductInventory.sql** script file in the **D:\Labfiles\Lab04\Starter** folder.
2. Configure SQL Server Management Studio to include the actual execution plan.
3. Execute the script against the **AdventureWorksDW** database. Review the execution plan, making a note of the indexes used, the execution time, and data space used.

► Task 3: Create a Columnstore Index on the FactProductInventory Table

1. Based on the scenario for this exercise, decide whether a clustered or nonclustered columnstore index is appropriate for the **FactProductInventory** table.
2. Create the required columnstore index. Re-execute the query to verify that the new columnstore index is used, along with existing indexes.
3. What, if any, are the disk space and query performance improvements?

Results: After completing this exercise, you will have created a columnstore index and improved the performance of an analytical query. This will have been done in real time without impacting transactional processing.

Exercise 2: Create a Columnstore Index on the FactInternetSales Table

Scenario

You need to improve the performance of queries that use the **FactInternetSales** table. The table has also become large and there are concerns over the disk space being used. You can use scheduled downtime to amend the table and its existing indexes.

Due to existing processing requirements, you must retain the foreign keys on the **FactInternetSales** table, but you can add any number of new indexes to the table.

The main tasks for this exercise are as follows:

1. Examine the Existing Size of the FactInternetSales Table and Query Performance
2. Create a Columnstore Index on the FactInternetSales Table

► Task 1: Examine the Existing Size of the FactInternetSales Table and Query Performance

1. In SQL Server Management Studio, in the **D:\Labfiles\Lab04\Starter** folder, open the **Query FactInternetSales.sql** script file.
2. Configure SQL Server Management Studio to include the actual execution plan.
3. Execute the script against the **AdventureWorksDW** database. Review the execution plan, making a note of the indexes that were used, the execution time, and data space used.

► Task 2: Create a Columnstore Index on the FactInternetSales Table

1. Based on the scenario for this exercise, decide whether a clustered or nonclustered columnstore index is appropriate for the **FactInternetSales** table.
2. Create the required columnstore index. Depending on your chosen index, you may need to drop and recreate existing indexes on the table.

3. Re-execute the query to verify that the new columnstore index is used, along with the existing indexes.
4. What, if any, are the disk space and query performance improvements?

Results: After completing this exercise, you will have greatly reduced the disk space taken up by the **FactInternetSales** table, and improved the performance of analytical queries against the table.

Exercise 3: Create a Memory Optimized Columnstore Table

Scenario

Due to the improved performance and reduced disk space that columnstore indexes provide, you have been tasked with taking the **FactInternetSales** table from disk and into memory.

The main tasks for this exercise are as follows:

1. Use the Memory Optimization Advisor
2. Enable the Memory Optimization Advisor to Create a Memory Optimized FactInternetSales Table
3. Examine the Performance of the Memory Optimized Table

► Task 1: Use the Memory Optimization Advisor

- In SQL Server Management Studio, run the **Memory Optimization Advisor** on the **FactInternetSales** table. Note that there are several issues to be resolved before the Memory Optimization Advisor can automatically convert the table.

► Task 2: Enable the Memory Optimization Advisor to Create a Memory Optimized FactInternetSales Table

1. Using either SQL Server Management Studio, or Transact-SQL statements, drop all the foreign keys.
2. Memory optimized tables cannot have more than eight indexes. Choose another three indexes to drop.



Note: Hint: consider the rows being used in the **Query FactInternetSales.sql** to guide your decision.

3. Run the **Memory Optimization Advisor** on the **FactInternetSales** table again. Create the clustered columnstore index when prompted, and then migrate the table.

► Task 3: Examine the Performance of the Memory Optimized Table

1. In SQL Server Management Studio, in the D:\Labfiles\Lab04\Starter folder, open the **Query FactProductInventory.sql** script file.
2. Configure SQL Server Management Studio to include the actual execution plan.
3. Execute the script against the **AdventureWorksDW** database, and then review the disk space used and the execution plan.

Results: After completing this exercise, you will have created a memory optimized version of the **FactInternetSales** disk-based table, using the Memory Optimization Advisor.

Question: Why do you think the disk space savings were so large?

Module Review and Takeaways

Columnstore indexes use disk compression to dramatically reduce disk I/O and, therefore, improve the speed of queries, especially in large data warehouse solutions. This module highlighted the benefits of using these indexes in the context of data warehousing, the improvements made in SQL Server, and the considerations needed to use columnstore indexes effectively in your solutions.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 5

Implementing an Azure SQL Data Warehouse

Contents:

Module Overview	5-1
Lesson 1: Advantages of Azure SQL Data Warehouse	5-2
Lesson 2: Implementing an Azure SQL Data Warehouse Database	5-6
Lesson 3: Developing an Azure SQL Data Warehouse	5-11
Lesson 4: Migrating to an Azure SQL Data Warehouse	5-18
Lesson 5: Copying Data with the Azure Data Factory	5-26
Lab: Implement an Azure SQL Data Warehouse	5-33
Module Review and Takeaways	5-39

Module Overview

If you are planning to implement an on-premises data warehouse, one of the more difficult considerations is the specification and cost of the hardware. Microsoft Azure SQL Data Warehouse provides a cloud-based alternative. By implementing an Azure SQL Data Warehouse, you can avoid hardware costs. If your data warehouse grows, you can easily scale the cloud-based resources—and therefore the costs—to grow with it.

In this module, you will learn about the many advantages of using Azure SQL Data Warehouse and how to implement an Azure SQL Data Warehouse database.

Objectives

After completing this module, you will be able to:

- Describe the advantages of Azure SQL Data Warehouse.
- Implement an Azure SQL Data Warehouse.
- Describe the considerations for developing an Azure SQL Data Warehouse.
- Migrate data to an Azure SQL Data Warehouse.
- Use Microsoft Azure Data Factory to copy data to an Azure SQL Data Warehouse.

Lesson 1

Advantages of Azure SQL Data Warehouse

One of the problems with traditional data warehouses is that they require expensive hardware that is not utilized all of the time. Azure SQL Data Warehouse offers a managed cloud-based solution capable of handling massive volumes of data. You can change the resource allocated to your data warehouse within seconds, so that you incur costs for resources only when you need them. Azure SQL Data Warehouse benefits from many features, including Microsoft's massively parallel processing architecture, columnstore indexes, and the hybrid cloud, which can be used to access structured and unstructured data within a familiar environment. In this lesson, you will learn about these features.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the features of an Azure SQL Data Warehouse.
- Understand the resource types, how they rescale, and their associated costs.
- Describe the security and availability features in Azure SQL Data Warehouse.
- Explain how you can query unstructured data with PolyBase.
- Describe what is possible in the hybrid cloud environment.

What is Azure SQL Data Warehouse?

Azure SQL Data Warehouse is a cloud-based database that can store and process relational and nonrelational data. It can be scaled to operate with workloads of different sizes up to enterprise level. Being fully integrated with Azure, it has access to many other useful Azure features, such as Microsoft Azure Blob Storage, Microsoft Azure Stream Analytics, and Azure Data Factory. The service is fully managed, including all maintenance, software patches, and backups.

Azure SQL Data Warehouse benefits from using many cutting edge technologies and features, including:

- **Massive Parallel Processing:** for complex queries, the workload can be spread across a number of nodes operating simultaneously, significantly improving performance.
- **Advanced Query Optimization:** by calculating complex statistics across the data, Azure SQL Data Warehouse can calculate the optimum query distribution across the nodes.
- **Columnstore Indexes:** tables in Azure SQL Data Warehouse use columnstore indexes. Compared to tables with rowstore indexes, tables with columnstore indexes can achieve up to five times more compression on storage and up to 10 times better performance on queries.
- **PolyBase:** with PolyBase, you can interact with data from different sources using Transact-SQL.
- **Auditing:** you can use auditing to keep a log of database events and use that information for security purposes.

- Cloud-based database
 - Relational and nonrelational
 - Enterprise workloads
 - Integrated with Azure
 - Fully managed service
- Benefits include:
 - Massive parallel processing
 - Advanced query optimization
 - Columnstore indexes
 - PolyBase integration
 - Auditing
 - Scalability

- **Scalability:** you can easily adjust the resources and cost allocated to your data warehouse—that gives you a great deal of flexibility. If your database grows and needs more resources, you can adjust the Azure SQL Data Warehouse resource allocation accordingly.

You can implement an Azure SQL Data Warehouse database in the Azure portal. You can then interact with it using Transact-SQL on a client computer.

Scalability and Cost

One of the most notable features of Azure SQL Data Warehouse is its scalability. There is no upfront cost and the amount of resource you wish to allocate to your Azure SQL Data Warehouse database is completely flexible. The cost is determined by the amount of resource you use. The two resource usage types in Azure are storage and compute—they operate independently.

Storage

The storage space and its associated cost automatically rescales when you add or remove data. Part of the Azure SQL Data Warehouse usage cost is based on the storage your Azure SQL Data Warehouse database uses.

For details of Azure SQL Data Warehouse storage costs, see:

 **Azure Storage Pricing**

<http://aka.ms/cn4u11>

Compute

The value of compute determines the performance of executions in Azure SQL Data Warehouse. It is measured in Data Warehouse Units (DWU). You can change the amount of DWUs used for a database as appropriate, depending on the expected workload requirements of your database. The DWU values are at discrete intervals, ranging from 100 DWU to 2000 DWU—you can change the DWU value at any time. If there is a period where you anticipate a high workload, you can increase the DWU value, and then decrease it at the end of that period.

Part of the cost of using Azure SQL Data Warehouse is based on the DWUs you use.

For more information on Azure SQL Data Warehouse compute usage costs, see:

 **SQL Data Warehouse Pricing**

<http://aka.ms/o0hvxg>

If you are not going to use your database during a certain time period, you can pause the database. This will release all compute resources and ensure that your database is not accruing any cost for compute. When you need to use the database again, you can restart it. You will find out how to change the DWU performance later in this module.

- No upfront cost
- Storage
 - Adjusts automatically
 - Cost based on storage used
- Compute
 - Determines execution performance
 - Data Warehouse Unit (DWU)
 - Increase or decrease DWU
 - Cost based on DWU used
 - Pause and start

Security and Availability

Security

Azure SQL Data Warehouse databases are protected by the server firewall—this will block any computers attempting to access the server, unless you give permission by adding the IP address of the computer to the firewall settings.

When you set up a server for an Azure SQL Data Warehouse database, you must initially specify a login for the server administrator. As with a SQL Server® database, you can then add other logins with certain levels of authorization as appropriate. You can also use transparent data encryption.

Auditing is a new security feature in Azure SQL Data Warehouse. You can use auditing to track any suspicious activity and write events to an audit log in your Azure storage account.

For further information on the auditing feature, see:



Get started with SQL database auditing

<http://aka.ms/fwggq18>

- Security
 - Firewall
 - Add logins
 - Set authorisation
 - Auditing
- Availability
 - Can restore in different region
 - Can choose restore point in last seven days

Availability

Azure replicates your data several times, including in a secondary region that is 100s of miles from the primary region. If a disruptive event in a region causes a database to become inaccessible, you can contact Microsoft support and arrange to restore the database in another region.

Each database has snapshots recorded at least every eight hours, and Azure retains the snapshots for seven days. If there is a user error, a database can be restored to the most suitable snapshot point.

PolyBase

Many data warehouses now have a requirement to use unstructured, nonrelational data sources. This can result in a demand to understand the different systems that store the unstructured data. You can use PolyBase to query such data sources without needing to directly use the systems where they are stored. You can use PolyBase with Azure SQL Data Warehouse to set up an external table, which is a reference to the data source and the system it uses. You can then query the external table as you would a local SQL Server table.

For information about how to create an external table in PolyBase, see:



CREATE EXTERNAL TABLE (Transact-SQL)

<http://aka.ms/hxdhod>

- Can access unstructured data in other systems
- Set up external table to link to data source
- Query external table as normal table

Hybrid Cloud

Azure SQL Data Warehouse provides a hybrid cloud environment. You can integrate seamlessly between on-premises databases, cloud databases, and unstructured data sources.

You can use PolyBase to query data in the different types of sources and to copy data from one type of source to another. Regardless of the type of data source, you can interact with the data in the same way, using Transact-SQL.

You can also use Azure Data Factory to set up pipeline activities to schedule copying data from one type of source to another—this is explained later in the module.

Question: You are working on a data warehouse that sometimes requires executions with very high workloads, but also goes for long periods of time with a very low workload requirement. If you migrate this data warehouse to Azure SQL Data Warehouse, how can you improve the efficiency of resource allocation?

- Can integrate between on-premises, cloud, and unstructured data sources
- Use PolyBase to query and copy data with Transact-SQL
- Schedule data copy using Azure Data Factory

Lesson 2

Implementing an Azure SQL Data Warehouse Database

In this lesson, you will learn how to implement an Azure SQL Data Warehouse database. As part of the process, you must also create a logical server for the database. To connect to the database, you must configure the firewall settings of the server to allow the IP address of your computer.

Lesson Objectives

After completing this lesson, you will be able to:

- Create a logical server.
- Create a database.
- Configure the server firewall.
- Connect to an Azure SQL Data Warehouse database with Microsoft SQL Server Management Studio.

Creating a Server

To implement an Azure SQL Data Warehouse database, you must specify a logical server to which it is assigned. A logical server is not a physical server—it is part of the overall service provided by Azure SQL Data Warehouse. When creating a logical server, you must specify:

- **Server name:** this must be a name that is globally unique.
- **Server admin logon:** this is the username for the server administrator. You can specify other users after the server is created.
- **Password:** this is the password for the server administrator. The password must be at least eight characters long, and include at least one character from the following categories:
 - Uppercase letters
 - Lowercase letters
 - Numbers
 - Non alpha numeric characters
- **Location:** the location determines where the resources are applied to your logical server. You should choose the location nearest to where you are.

- Logical server
- Specify:
 - Server name that has not been used
 - Server admin logon
 - Password
 - Location nearest to you
- Create database in same process

To create a server in the Azure portal, you must also create a database in the same process. If you do not yet need to use a database, you can create a temporary database, and then delete it after the server is created. The process of creating a database, and the server it is assigned to, is described in the next topic.

Creating a Database

To create an Azure SQL Data Warehouse database, in the Azure portal, click **New**, click **Data + Storage**, and then click **SQL Data Warehouse**. You must then configure the parameters in the following list, and then click **Create**:

- **Name:** the name of your database.
- **Performance:** drag the slider to choose the amount of DWUs you wish to use.
- **Server:** you can either create a new server or use an existing one.
- **Select source:** you should use a blank database unless you wish to practice with the sample database. The sample database is a version of the Microsoft AdventureWorksDW database.
- **Resource group:** resource groups provide a means of organizing different resources that you may use in Azure. When using resources with other services, such as Azure SQL Database, and Azure Data Factory, you can also specify resource groups. If, for example, you have several projects that use a number of different resources, you can set up a resource group for each project, so that each resource can be categorized according to project. You can either create a new resource group or select an existing one.

- Create database
 - Name of database
 - Drag slider to change DWU performance
 - Create a new server or use existing server
 - Source
 - Create a new resource group or use existing resource group
- DWU settings
 - Scale
 - Pause/start

After you have created the database, you can change the DWUs of the database as required. To view or change these settings, click **All resources**, and then click the database you wish to view. To change the DWUs, click **Scale**, and then adjust the slider as appropriate. To pause or start the database, click **Pause** or **Start**.

Configuring the Server Firewall

For security reasons, client computers are not allowed to connect to an Azure server until the firewall has been configured to permit connections from the client IP address. To configure the firewall to permit connections from the client machine you are using, in the Azure portal, click **All resources**, click the server you wish to configure, click **Show firewall settings**, click **Add client IP**, and then click **Save**.

If the IP address of your computer changes, you will need to repeat this process. To avoid this, you can manually specify a rule to allow for the range of IP addresses a client computer may use. You can also manually add IP addresses for client computers that will need to access the server.

For more information on how to manually specify a rule and how to add other IP addresses, see:

 **How to configure firewall settings on SQL Database using the Azure Portal**

<http://aka.ms/iddzy1>

Connecting to Azure Database Using SQL Server Management Studio

After you have created your Azure SQL Data Warehouse database and configured the firewall of the server, you can connect to the database using SQL Server Management Studio (SSMS) on your computer. You must use the fully qualified name of the logical server, which is **<Server Name>.database.windows.net**, where <Server Name> is the name you specified when you created the server. To retrieve this information from the Azure portal, click **All resources**, and then click the name of the database you wish to connect to. The fully qualified server name is shown below the **Server name**.

- Fully qualified server name
- Connect to server using SSMS
- USE statement not supported
- Right-click database, New Query
- Most Transact-SQL supported in Azure SQL Data Warehouse databases

To connect to the logical server in SSMS, click **Connect Object Explorer**. In the **Connect to Server** dialog box, specify:

- **Server name:** type the fully qualified server name.
- **Authentication:** use SQL Server Authentication.
- **User name:** use the Server Admin Logon specified when you created the server (you can set up other users after you are connected).
- **Password:** use the password specified when you created the server.

After you have connected to the server you created, you can navigate to your database using Object Explorer. The USE statement is not supported in Azure SQL Data Warehouse, so to run a query against your database from Object Explorer, right-click the database, and then click **New Query**. You can then type the query and run it as required.

Most Transact-SQL used in SQL Server is also supported in Azure SQL Data Warehouse. There are, however, some differences, which are described later in this module.

Demonstration: Creating and Configuring an Azure SQL Data Warehouse Database

In this demonstration, you will see how to:

- Create an Azure SQL Data Warehouse database and server.
- Change the performance settings.
- Configure the Azure firewall.
- Connect to the Azure server with SSMS.

Demonstration Steps

Create an Azure SQL Data Warehouse Database and Server

1. Ensure that the **MT17B-WS2016-NAT**, **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod05** folder, right-click **Setup.cmd**, and then click **Run as administrator**.

3. Click **Yes** when prompted to confirm that you want to run the command file, and then wait for the script to finish.
4. Open Microsoft Internet Explorer® and go to **https://portal.azure.com/**.
5. Sign in to the Azure portal with your Azure pass or Microsoft account credentials.
6. Click **New**, click **Databases**, and then click **SQL Data Warehouse**.
7. In the **SQL Data Warehouse** blade, in the **Database Name** box, type **BikeSalesDW**.
8. Under **Resource group**, click **Create new**, and then in the name box, type **BikeSalesRG**,
9. Under **Server**, click **Configure required settings**, and then click **Create a new server**.
10. On the **New server** blade, in the **Server name** box, type **20767c** followed by your initials, followed by today's date. For example, if your initials are CN and the date is 15 March 2018, type **20767ccn15mar18**. The server name must be unique; if the name is unique and valid, a green tick appears.
11. In the **Server admin login** box, type **BikeSalesadmin**.
12. In the **Password** and **Confirm password** boxes, type **Pa55w.rd**
13. Under **Location**, select a region nearest your current geographical location, and then click **Select**.
14. In the **SQL Data Warehouse** blade, click **Performance Tier**.
15. Move the slider to **DW200**, and then click **Apply**.
16. On the **SQL Data Warehouse** blade, verify that **Select source** has the value **Blank database**, and then click **Create**.
17. It will take some time for the new server and database to be created. The Azure portal will notify you when this step is finished.

Change the Performance Settings

1. On the **Microsoft Azure** blade, click **Resource groups**, then click **BikeSalesRG**.
2. On the **BikeSalesRG** blade, click **BikeSalesDW**.
3. On the **BikeSalesDW** blade, click **Scale**.
4. On the **Scale** blade, drag the slider to **DW100**, click **Save**, and then, in the **Save performance level** message box, click **Yes**.
5. Close the Scale blade.

Configure the Azure Firewall

1. On the **Microsoft Azure** blade, click **All resources**.
2. In the **All resources** blade, click **<Server Name>** (where **<Server Name>** is the name of the SQL server you created).
3. In the **<Server Name>** blade, under **Settings**, click **Firewall / Virtual Networks**.
4. In the firewall settings blade, click **Add client IP**, and then click **Save**.
5. When the firewall changes are complete, click **OK**.
6. Leave Internet Explorer open.

Connect to the Azure Server with SQL Server Management Studio

1. Start SQL Server Management Studio and connect to the **MIA-SQL** database engine instance using Windows® authentication.
2. On the **File** menu, click **Connect Object Explorer**.
3. In the **Connect to Server** dialog box, enter the following values, and then click **Connect**:
 - **Server name:** <Server Name>.database.windows.net (where <Server Name> is the name of the server you created)
 - **Authentication:** SQL Server Authentication
 - **Login:** BikeSalesadmin
 - **Password:** Pa55w.rd
4. In Object Explorer, under <Server Name>.database.windows.net, expand **Databases**. Note that the **BikeSalesDW** database is listed.
5. Leave SQL Server Management Studio open.

Check Your Knowledge

Question	
Which of the following statements about a logical server is correct?	
Select the correct answer.	
<input type="radio"/>	You must create a new logical server each time you create a new database.
<input type="radio"/>	A logical server is not a physical server.
<input type="radio"/>	You are charged for each logical server you create.
<input type="radio"/>	The more databases that a logical server hosts, the more the performance of the logical server will decrease.
<input type="radio"/>	You can change the DWU performance of the logical server.

Lesson 3

Developing an Azure SQL Data Warehouse

In Azure SQL Data Warehouse, many of the features that you can use in SQL Server are also supported. However, because of the new technologies used, there are some new features and considerations that you need to be aware of. This lesson describes these issues, in addition to showing how you can work around features that are not supported.

Lesson Objectives

After completing this lesson, you will be able to:

- Explain concurrency and memory allocation.
- Describe the different types of data distributions.
- Use a CREATE TABLE AS SELECT statement.
- Manage the GROUP BY limitations.
- Understand how you can apply temporary tables.
- Apply user defined schemas to organize your data.

Concurrency and Memory Allocation

In Azure SQL Data Warehouse, the query execution performance is determined by a number of concepts: the resource class, the concurrency slots, the concurrent queries, and the memory allocation.

Resource Class

A resource class is a role given to a user and determines the resources that will be allocated to queries run by that user. There are four different resource classes:

- smallrc
- mediumrc
- largerc
- xlargerc

By default, users are members of the smallrc resource class; if a user needs more resources to run their queries, they can be allocated a higher resource class.

Concurrency Slots

When a query is run in Azure SQL Data Warehouse, it can use one or more concurrency slots. The more slots that are used, the greater the performance of that query. The number of slots a query will use is a function of the resource class of the user running the query and the DWU service level of your data warehouse.

- Resource class
- Concurrency slots
 - Query may use more than one concurrency slot
 - Dependent on resource class and DWU service level
- Concurrent queries
 - Maximum of 32 queries
 - Maximum slots dependent on DWU service level
- Memory allocation
 - Dependent on resource class and DWU service level



Note: If a user with a higher resource class is running a query that does not have a high resource requirement, the query might not be allocated all of the concurrency slots that would normally be available for that resource class.

Concurrent Queries

The number of queries that can run concurrently in a data warehouse is determined by two factors:

- A maximum of 32 queries can run concurrently, regardless of how many concurrency slots each query uses.
- The maximum number of concurrency slots that an Azure SQL Data Warehouse database can use is a function of the DWU service level.

If more than 32 queries are started concurrently, or if the number of concurrency slots has been exceeded, any other query will be queued until the resources are available.

Memory Allocation

The performance of a query will improve if you allocate more memory to it. The memory that is available to a query is calculated from the resource class of the user running the query, and the DWU service level of your data warehouse.

For more information on concurrency and workload management in SQL Data Warehouse, see:



Concurrency and workload management in SQL Data Warehouse

<http://aka.ms/v9tpol>

Data Distribution

In Azure SQL Data Warehouse tables, the data in each table is spread across a number of different locations known as distributions. This results in a significant improvement in performance, because the workload is shared among the distributions.

There are two methods by which you can distribute data: round-robin distribution and hash distribution. You specify the distribution method when creating the table.

- Data in tables allocated to distributions
- Round-robin distribution
 - Random distribution allocation
- Hash distribution
 - Choose hashed column
 - Distribution determined by function of column value
 - Ensure hashed column has even spread of data

Round-Robin Distribution

When a table uses a round-robin distribution, buffers containing rows of data are each randomly allocated to different distributions until all of the available distributions have been used. Because the Azure SQL Data Warehouse engine does not have to determine which distribution to allocate the records to, the process of adding records is quicker. However, because the data is located at random distributions, it takes longer to find and retrieve it.

The following example creates a table using a round-robin distribution:

Create a Table Using Round-Robin Distribution

```
CREATE TABLE Countries
(CountryKey int NOT NULL
,Country varchar(50) NOT NULL
)
WITH
(CLUSTERED COLUMNSTORE INDEX
,DISTRIBUTION = ROUND_ROBIN
);
```



Note: When you create tables in an Azure SQL Data Warehouse database, it uses round-robin distribution by default. It is therefore not necessary to explicitly specify a round-robin distribution as in the above example, although it is good practice to do so.

Hash Distribution

When a table uses a hash distribution, each row is allocated to a distribution, based on a hash function applied to one of the columns in the table. The distribution used for each row will depend on the value of the hashed column, so that any rows with the same value are allocated to the same distribution. Because the Azure SQL Data Warehouse engine can quickly determine the distribution of each row, the process of finding and retrieving data from tables using the hash tag distribution is quicker. However, the process of adding records is slower.

The following example creates a table using hash distribution, where the hashed column is CountryKey:

Create a Table Using Hash Distribution

```
CREATE TABLE Countries
(CountryKey int NOT NULL
,Country varchar(50) NOT NULL
)
WITH
(CLUSTERED COLUMNSTORE INDEX
,DISTRIBUTION = HASH(CountryKey)
);
```

When choosing the hash column, avoid columns where nulls are allowed or one where many of the values are the same. Because the Azure SQL Data Warehouse engine places columns with the same value on the same distribution, this will result in an uneven spread of data across the distributions, which will reduce the performance. The ideal columns to use are those that have a unique value for each record. If there is not a suitable candidate column, it is better to use a round-robin distribution.

CREATE TABLE AS SELECT

CREATE TABLE AS SELECT (also known as CTAS) is an operation that is particular to Azure SQL Data Warehouse. You can use it to make a copy of a table, with the option to set properties on the new table. One of the most common applications is if you want to change the distribution type or the index properties of a table. For example, there might be a table that uses a round-robin distribution that you want to change so that it uses a hash distribution.

The following example makes a copy of the Countries table with the new table using a hash distribution on the CountryKey column:

CREATE TABLE AS SELECT

```
CREATE TABLE Countries_New
WITH
(
  CLUSTERED COLUMNSTORE INDEX,
  DISTRIBUTION = HASH(CountryKey)
)
AS SELECT * FROM Countries
;
```

Assuming that you do not wish to keep the original table, you can rename the tables and delete the original table as shown in the following example:

Rename Tables and Delete Original Table

```
RENAME OBJECT Countries TO Countries_old;
RENAME OBJECT Countries_new TO Countries;
DROP TABLE Countries_old;
```

In Azure SQL Data Warehouse, the SELECT INTO and MERGE statements are not supported. Also, joins are not supported in UPDATE and DELETE statements. In addition to the example shown above, CTAS can be used to provide a workaround for these unsupported features.



Note: Differences between SQL Server and Azure SQL Data Warehouse will be discussed in more detail later in this module.

- Makes copy of a table
- Can set index properties and distribution type

```
CREATE TABLE Countries_New
WITH
(
  CLUSTERED COLUMNSTORE INDEX,
  DISTRIBUTION = HASH(CountryKey)
)
AS SELECT * FROM Countries
;
```

- Use to work around unsupported features

GROUP BY Limitations

In SQL Server, you can use the GROUP BY clause in a SELECT statement to arrange rows in groups, typically to support aggregations. If you need to group by different attributes at the same time—for example, to report at different levels—you can use the GROUPING SETS, CUBE, and ROLLUP subclauses that enable multiple sets to be returned in the same query.

For more information on GROUPING SETS, CUBE, and ROLLUP subclauses, see:

Using GROUP BY with ROLLUP, CUBE, and GROUPING SETS

<http://aka.ms/g1bke3>

- GROUP BY clause is supported
- GROUPING SETS, CUBE and ROLLUP subclauses are not supported
- UNION ALL operator is supported
- When migrating to Azure SQL Data Warehouse, ensure queries containing unsupported clauses are amended

When querying data in an Azure SQL Data Warehouse database, you can use the GROUP BY clause, but the GROUPING SETS, CUBE, and ROLLUP subclauses are not supported. You can achieve the same results by using the UNION ALL operator, which is supported in Azure SQL Data Warehouse. If you are migrating a data warehouse from SQL Server to Azure SQL Data Warehouse, ensure that any queries containing GROUPING SETS, CUBE, or ROLLUP subclauses are amended to use the UNION ALL operator to return the same result set.

Temporary Tables

Temporary tables are very useful programming tools. In SQL Server, a temporary table can only be accessed while the procedure in which it is created is running. However, in Azure SQL Data Warehouse, when you create a local temporary table, you can access it from anywhere in the same session.

Consider the following code example that creates a stored procedure that copies the **Countries** table to the local temporary table **#tempcountrycopy** using a CTAS statement:

Stored Procedure to Create Local Temporary Table

```
CREATE PROCEDURE prctempcountrycopy
AS
CREATE TABLE #tempcountrycopy
WITH
(CLUSTERED COLUMNSTORE INDEX
,DISTRIBUTION = ROUND_ROBIN
)
AS SELECT * FROM Countries
;
```

- Local temporary tables can be accessed anywhere within session
- Global temporary tables are not supported

The following code example then executes the stored procedure:

Execute the Stored Procedure

```
EXECUTE prctempcountrycopy
```

After running the stored procedure, the temporary table **#tempcountrycopy** can still be referred to, as shown in the following code example:

Code That Refers to Temporary Table

```
SELECT * FROM #tempcountrycopy
```

Global temporary tables are not supported in Azure SQL Data Warehouse.

User Defined Schemas

In Azure SQL Data Warehouse, joins between databases are not supported. Therefore, you must have all of your data in one database. If you are migrating from a data warehouse that used more than one database, into a single database on Azure SQL Data Warehouse, you are advised to restructure the schemas and tables as follows:

- Use user-defined schemas to identify the databases that were in the legacy system.
- If the legacy databases used schemas that you wish to maintain a reference to, you can prefix the name of each new table with a reference to the schema where the legacy table is situated.

- All data in one database
- Use schemas to identify legacy databases

The following tables illustrate an example of how a data warehouse from a legacy system could be restructured in an Azure SQL Data Warehouse. The legacy data warehouse uses a staging database and a data warehouse database.

Legacy Data Warehouse			
Databases	VintageHatsSalesStaging		VintageHatsDataWarehouse
User Defined Schemas	Sales	Production	Not Used
Tables	OrderHeader OrderDetail Customers	Products ProductCategory ProductSubCategory	FactSales DimProducts DimCustomers

Azure SQL Data Warehouse		
Databases	VintageHats	
User Defined Schemas	Staging	DataWarehouse
Tables	SalesOrderHeader SalesOrderDetail SalesCustomers ProductionProducts ProductionProductCategory ProductionProductSubCategory	FactSales DimProducts DimCustomers

 **Note:** Any security settings applied at database level in the legacy data warehouse should be applied to the equivalent schemas in the Azure SQL Data Warehouse.

Verify the correctness of the statement by placing a mark in the column to the right.

Statement	Answer
A gender column in a table would be a suitable hashed column. True or false?	

Lesson 4

Migrating to an Azure SQL Data Warehouse

In this lesson, you will learn how to use the Microsoft Data Warehouse Migration Utility to migrate data. You will also see the table features and data types that are not supported in Azure SQL Data Warehouse, and learn how to check if a database contains any of these features.

Lesson Objectives

After completing this lesson, you will be able to:

- Install the Data Warehouse Migration Utility.
- Migrate a database to Azure SQL Data Warehouse.
- Consider other available options for migrating data to an Azure SQL Data Warehouse.
- Describe differences between schemas in SQL Server and Azure SQL Data Warehouse.
- List the Transact-SQL features that are not supported in Azure SQL Data Warehouse.

The Data Warehouse Migration Utility

You can migrate data from an on-premises database to an Azure SQL Data Warehouse database in a number of ways. One of the best methods is by using the Data Warehouse Migration Utility. The Data Warehouse Migration Utility has the following benefits:

- It is relatively straightforward to use.
- You can copy multiple tables.
- You can specify the distribution type to use for each table.
- It can notify you of features in the database that are incompatible with Azure SQL Data Warehouse before starting the migration.

- Advantages
 - Straightforward
 - Multiple tables
 - Specify distribution type
 - Notification of incompatibility
- Download from Internet
- Must have BCP and Excel installed

Installing the Data Warehouse Migration Utility

You can download the Data Warehouse Migration Utility from the Internet.

To access the page containing the download link, see:

 **Data Warehouse Migration Utility (Preview)**

<http://aka.ms/taqsbg>

Scroll down the page to the link where you can download the Data Warehouse Migration Utility—and then click the link. After the download is complete, double-click the executable file you have downloaded to install the Data Warehouse Migration Utility. After the installation, a link to the Data Warehouse Migration Utility will be available on your desktop.

For the Data Warehouse Migration Utility to function fully, you must have the BCP command-line utility and Microsoft Excel® installed.

Migrating Data with the Data Warehouse Migration Utility

You can use the Data Warehouse Migration Utility to migrate the schema and data from an existing data warehouse to Azure SQL Data Warehouse. Before migrating the data, you should check if your database schema has any features that are incompatible with Azure SQL Data Warehouse.

Check Compatibility

In a previous topic, you saw the features in SQL Server schemas that are not compatible with Azure SQL Data Warehouse. You can check if your schema contains any incompatible features by performing the following steps:

1. Open the Data Warehouse Migration Utility, specifying the server where your database is situated as the source.
2. In the Data Warehouse Migration Utility window, select your database.
3. Activate the Check Compatibility process, specifying the file name and location of the compatibility report file.
4. Open the compatibility report in Excel.
5. The compatibility report will list any types of incompatibilities that were found in the schema of the database, and how many instances there were of each incompatibility.

If your schema contains any incompatible features, you should consider the impact of the loss of these features, and make any necessary changes to the tables affected before migrating.

 **Note:** When you begin the migration process, the migration will proceed whether or not incompatible features are found. The features will be removed when the database is migrated.

Migrate Schema

After you have resolved any incompatibilities, you can migrate your schema by performing the following steps:

1. If you have not done so already, create the database you are migrating to, using the Azure portal.
2. Open the **Data Warehouse Migration Utility**, specifying the server where your database is contained as the source.
3. In the **Data Warehouse Migration Utility** window, select your database.
4. Select all of the tables you wish to migrate, specifying the distribution type, and the distributed column if applicable.
5. Activate the **Migrate Schema** process. You will then see the Transact-SQL script that will be used to create the tables in Azure SQL Data Warehouse.
6. Activate the **Run Script** process then, when prompted for details of the destination database, specify the details of the Azure SQL Data Warehouse database to which you are migrating.

- Check compatibility
- Migrate schema
- Migrate data
 - bcp commands to export and import



Note: Instead of running the script, you can copy it to a file or the clipboard, and then run it from SSMS later.

Migrate Data

After you have migrated the schema, you can migrate the data in the tables to the Azure SQL Data Warehouse database. The Data Warehouse Migration Utility uses the parameters you specify to generate bcp commands contained in two bat files. One file exports the legacy data to a flat file on the server, and the other imports the data into the Azure SQL Data Warehouse database. This procedure is described in the following steps:

1. Open the **Data Warehouse Migration Utility**, specifying the server where your database is contained as the source.
2. In the Data Warehouse Migration Utility window, select your database.
3. Select all of the tables you wish to migrate and activate the **Migrate Data** process, specifying the folder location of the bat files, and the details of the destination database.
4. In the folder you specified, run the bat file to export the data, followed by the bat file to import the data.



Note: You are advised to only use the Migrate Data process for small amounts of data, as the Data Warehouse Migration Utility does not have built-in retries.

Other Migration Tools

In addition to the Data Warehouse Migration Utility, there are several other options for loading data into your Azure SQL Data Warehouse. The method you choose will be determined by the amount of data you want to migrate, and the data source.

Azure Feature Pack for Integration Services (SSIS)

The Azure Feature pack is a downloadable extension for SSIS that facilitates the transfer of data between on-premises data and data stored in Azure. The pack includes Azure Storage and Azure Subscription connection managers, and you can use the Azure Blob source and destinations in your data flows tasks. Task components include the Azure Blob upload task, Azure Blob download task, Azure HDInsight® Hive task, Azure HDInsight Pig task, Azure HDInsight Create Cluster task and the Azure HDInsight Delete Cluster task. The pack also includes an Azure Blob enumerator component.



Azure Feature Pack for Integration Services (SSIS)

<https://aka.ms/yvglgt>

- Options for loading data into an Azure SQL Data Warehouse include:
 - **Azure Feature Pack for Integration Services (SSIS)**
 - Downloadable extension for SSIS that facilitates the movement of data between on-premises and cloud
 - **SSIS**
 - Add Azure SQL Data Warehouse connection in data flows
 - Use SQL Agent to schedule regular transfer of data
 - **Bulk Copy Program (bcp)**
 - Useful for small data, use bcp to copy data to flat files and load into the data warehouse destination
 - **AZCopy**
 - Copy data from flat files into Blob storage, and use PolyBase to load into data warehouse
 - **Import/Export**
 - For data larger than 10TB, bcp data to files, copy to disks and ship to Microsoft
 - **PolyBase and T-SQL**
 - Move UTF-8 formatted data in text files to Azure Blob storage or HDInsight, then use T-SQL command to load into the data warehouse
 - PolyBase uses the massively parallel processing (MPP) architecture for fast loading

Other options for loading data include:

- **Microsoft SQL Server Integration Services (SSIS):** in your SSIS packages, you can add a connection manager to your Azure SQL Data Warehouse and use this as the destination in your data flow tasks. SSIS is a good option if you have data residing in other vendor databases that requires copying to the cloud. Ensure you check the data types of column mappings, and other unsupported features before executing the packages (restrictions are covered in the next topic) if you are amending an existing package from an on-premises SQL Server destination to an Azure destination. Furthermore, you can create a SQL Agent job to schedule the migration to run on a regular basis.
- **Bulk Copy Program (bcp):** bcp is ideal for loading small amounts of data from SQL Server. You can use the bcp utility to copy data from SQL Server to flat files, and then use bcp to load the data in the flat files into your Azure SQL Data Warehouse.
- **AZCopy:** you can use the AZCopy command-line tool in combination with bcp, though this is recommended for data under a terabyte in size. Firstly, use bcp to copy the data from SQL Server into flat files. You can then use AZCopy to copy the data from the files to Azure Blob storage, and use PolyBase to load the data into your Azure SQL Data Warehouse.
- **Import/Export:** if your data is greater than 10TB, Microsoft recommends that you copy the data to disk and ship it to them. Use bcp to copy the data from SQL Server to flat files, save the files to transferrable disk, and send to Microsoft for loading into your Azure SQL Data Warehouse.
- **PolyBase and T-SQL:** you can move your data to Azure Blob storage or HDInsight, and use T-SQL to load the data into a new database table. Firstly, format your data as UTF-8 for PolyBase, and then move it to Azure Blob storage or HDInsight, storing it in text files. For HDInsight, you can also use ORC or Parquet format. You configure the external objects in your Azure SQL Data Warehouse, defining the location and format of the data. You then use a T-SQL command to load the data into a new table. PolyBase has the advantage of being faster than SSIS because it uses the massively parallel processing (MPP) architecture in SQL Data Warehouse.

For detailed instructions and examples of how to perform each of the above methods, see the following article:



Load data into Azure SQL Data Warehouse

<https://aka.ms/vnqc26>

Differences Between SQL Server and Azure SQL Data Warehouse Schemas

Although most features in Azure SQL Data Warehouse and SQL Server run in the same way, there are some differences in the schemas. It is important to be aware of these differences, especially if you are migrating your data warehouse to Azure.

Table Features

Tables in Azure SQL Data Warehouse do not support the following features:

- Primary Keys
- Foreign Keys
- Unique Indexes

- Some table features not supported
 - Primary Keys
 - Foreign Keys
 - Unique Indexes
 - Constraints
- Some data types not supported
 - numeric
 - nvarchar(max)
 - varchar(max)

- Constraints

Data Types

In Azure SQL Data Warehouse, you can use most of the data types that you use in SQL Server. However, some data types are not supported, including:

- **Numeric:** use decimal instead.
- **nvarchar(max):** use nvarchar(4000) or smaller.
- **varchar(max):** use varchar(8000) or smaller.

For a more detailed list of table features and data types that are not supported in Azure SQL Data Warehouse, see:

 **Migrate your schema to SQL Data Warehouse**

<http://aka.ms/xbdb1k>

Updating Transact-SQL

Most of the Transact-SQL that you can use in SQL Server can also be used in Azure SQL Data Warehouse. However, there are some Transact-SQL features that are not supported, including:

- SELECT INTO
- Joins in UPDATE statements
- Joins in DELETE statements
- GROUPING SETS, CUBE, and ROLLUP subclauses in GROUP BY clauses
- MERGE statement
- Cross database joins

For more information on all of the main TRANSACT-SQL features not supported, see:

 **Migrate your SQL code to SQL Data Warehouse**

<http://aka.ms/w681ix>

If you have some Transact-SQL in your legacy database that you are intending to use in a migrated database, you should check the Transact-SQL for features that are not supported. In most cases, you can rewrite the code using features that are supported to achieve the same result.

- Some Transact-SQL not supported
- Rewrite to achieve same result

Demonstration: Migrating a Database to Azure SQL Data Warehouse

In this demonstration, you will see how to:

- Install the Data Warehouse Migration Utility.
- Check compatibility of the legacy database.
- Migrate the schema.
- Migrate the data.

Demonstration Steps

Install the Data Warehouse Migration Utility

1. In the **D:\Demofiles\Mod05\Installdwmigrationutility** folder, double-click **Data Warehouse Migration Utility.msi**.
2. In the **Data Warehouse Migration Utility** dialog box, on the **Select Installation Folder** page, click **Next**.
3. On the **License Agreement** page, review the License Agreement, click **I Agree**, and then click **Next**.
4. In the **User Account Control** dialog box, click **Yes**.
5. On the **Installation Complete** page, click **Close**.
6. Minimize all open applications and note the **Data Warehouse Migration Utility** icon on the desktop.

Check Compatibility of the Legacy Database

1. In SQL Server Management Studio, in Object Explorer, under **MIA-SQL**, expand **Databases**, expand **BikeSalesDW**, and then expand **Tables**.
2. Right-click each of the following tables, click **Select Top 1000 Rows**, and then review the data they contain:
 - a. **dbo.DimCustomers**
 - b. **dbo.DimProducts**
 - c. **dbo.DimResellers**
 - d. **dbo.FactInternetSales**
 - e. **dbo.FactResellerSales**
3. Minimize all windows, and then, on the desktop, double-click the **Data Warehouse Migration Utility** icon.
4. In the **DATA WAREHOUSE MIGRATION UTILITY** dialog box, note the following settings, and then click **Next**.
 - **Source Type:** SQL Server
 - **Destination Type:** Azure SQL Data Warehouse
5. In the **DATA WAREHOUSE MIGRATION UTILITY MIGRATION SOURCE** dialog box, in the **Server Name** box, type **MIA-SQL**, and then click **Connect**.
6. In the **Database List** pane, click **BikeSalesDW**, and then click **Check Compatibility**.
7. In the **Save As** dialog box, save the **BikeSalesDW – DatabaseCompatibilityReport.xlsx** file to the **D:\Demofiles\Mod05** folder.

8. In the **File Saved Successfully** dialog box, click **No**.
9. Switch to the **20767C-MIA-CLI** virtual machine, and log in as **Student** with the password **Pa55w.rd**
10. In File Explorer, in the address bar, type `\\MIA-SQL\D$\Demofiles\Mod05`, and then press Enter
11. Double-click the **BikeSalesDW – DatabaseCompatibilityReport.xlsx** file, and review the data in the Excel spreadsheet, and then close Excel without saving changes.
12. Switch to the **20767C-MIA-SQL** virtual machine.

Migrate the Schema

1. In the **DATA WAREHOUSE MIGRATION UTILITY MIGRATION SOURCE** dialog box, in the **Database List** pane, click **BikeSalesDW**, and then click **Migrate Selected**.
2. In the **Migrate Database** pane, select the **Table Name** check box. Note that the check boxes for all tables in the database are now selected, and then click **Migrate Schema**.
3. In the **Migrate Schema** pane, note the script that has been generated for each table, and then click **Run Script**.
4. In the **DATA WAREHOUSE MIGRATION UTILITY MIGRATION DESTINATION** dialog box, enter the following values, click **OK**, and wait for the script run:
 - **Database:** BikeSalesDW
 - **Server Name:** <Server Name>.database.windows.net (where <Server Name> is the name of the server you created previously)
 - **Authentication:** Standard
 - **User Name:** BikeSalesAdmin
 - **Password:** Pa55w.rd
5. In the **Script Applied Successfully** message box, click **OK**.

Migrate the Data

1. In the **Migrate Schema** pane, click **Migrate Data**.
2. In the **Migrate Data** pane, click **Run Migration**.
3. In the **DATA WAREHOUSE MIGRATION UTILITY** dialog box, in the **Package Path** box, type `D:\Demofiles\Mod05`, and then click **Next**.
4. In the **DATA WAREHOUSE MIGRATION UTILITY MIGRATION DESTINATION** dialog box, click **Generate**.
5. In the **Package(s) Generated Successfully** message box, click **OK**.
6. In File Explorer navigate to `D:\Demofiles\Mod05` folder, and run the **BikeSalesDW_Export.bat** script as Administrator.
7. Run the **BikeSalesDW_Import.bat** script, and wait for the script to complete.
8. Close the **DATA WAREHOUSE MIGRATION UTILITY** window.

Check the Data Has Been Migrated Successfully

1. In SQL Server Management Studio, in **Object Explorer**, under **<Server Name>.database.windows.net**, expand **Databases**, right-click **BikeSalesDW**, and then click **New Query**.
2. In the **Microsoft SQL Server Management Studio** message box, click **Cancel**.

- In the query pane, type the following SQL:

```
SELECT * FROM DimCustomers
SELECT * FROM DimProducts
SELECT * FROM DimResellers
SELECT * FROM FactInternetSales
SELECT * FROM FactResellerSales
```

- To check the data has been migrated to each of the five tables, highlight the first line, click **Execute**, review the data, and then repeat the process for the four other lines.
- Close SQL Server Management Studio without saving any changes. Close any other windows that are open.

Verify the correctness of the statement by placing a mark in the column to the right.

Statement	Answer
If you are migrating a schema using the Data Warehouse Migration Utility, and some of the data in the schema has features that are incompatible with Azure SQL Data Warehouse, the migration will continue.	

Lesson 5

Copying Data with the Azure Data Factory

This lesson provides a brief introduction to the Azure Data Factory and how you can use it to set up a scheduled activity to copy data from one table to another. You will see that the activities can reference data that is on-premises or in the cloud.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the entities in a data factory that work together in a data flow.
- Create a data factory.
- Set up a data gateway to access on-premises data.
- Set up a linked service to access databases.
- Set up a dataset to access tables in databases.
- Understand how you can schedule an activity to copy data from an on-premises database to a cloud database.
- Use a data factory diagram to monitor the data flow.

The Azure Data Factory

The Azure Data Factory is a separate service to Azure SQL Data Warehouse. However, you can integrate it with your data warehouse to set up scheduled processes to copy, transform and analyze data, whether the data is cloud-based, on-premises, structured or unstructured.

This lesson introduces the data factory and provides an overview of how you can use Azure Data Factory with your Azure SQL Data Warehouse, in particular to set up a schedule to copy data from an on-premises database to the cloud.

For further information on the capabilities of the data factory, see:

 **Data factory documentation**

<http://aka.ms/p1h4ha>

Pricing

Because the Azure Data Factory is a separate service to Azure SQL Data Warehouse, it is priced separately.

For price details of Azure Data Factory, see:

 **Data Factory Pricing**

<http://aka.ms/wbsjpo>

- Capabilities and application to Azure SQL Data Warehouse
- Entities
 - Activity
 - Pipeline
 - Dataset
 - Linked service
- Scheduling
- JSON templates
 - Edit parameters in script
 - Replace \ with \\

Entities

A data factory contains key entities that work together in a data flow:

- **Activity:** an activity is a certain process that is performed, such as copying data from one location to another. Activities run repeatedly in defined time intervals referred to as data slices—for example, you can schedule an activity to run every hour.
- **Pipeline:** a pipeline is a grouping of activities that work together to perform a task. A pipeline must contain at least one activity.
- **Dataset:** a dataset references the source or destination data structure, such as tables, documents, files, or folders. Datasets also contain information about the data slice interval.
- **Linked Service:** a linked service defines the system where the dataset is contained, such as an Azure SQL Data Warehouse, an on-premises SQL database, or an Azure Blob storage account.

Scheduling

Activities in pipelines run repeatedly at defined time intervals, which must be specified when a pipeline is created. Time intervals are also specified when defining the availability of a dataset. The unit of data consumed by an activity run is known as a data slice. A full explanation of scheduling in the Azure Data Factory is beyond the scope of this course.

For more information about scheduling, see:



Scheduling and Execution with Data Factory

<http://aka.ms/q3c72w>

JSON Templates

As you will see later in this lesson, when you create an entity, a JavaScript Object Notation (JSON) template script will be created. In this script, you can change the parameters to apply to your process as necessary.

When changing parameters in JSON, it is important to be aware that the \ character is treated as a special character. If any of your parameters, such as a connection string or a user name, contain a \ character, you must replace it with \\. A more detailed instruction on using JSON is beyond the scope of this course.

For further information about using JSON, see:



JSON Data (SQL Server)

<http://aka.ms/q4hitr>

Creating a Data Factory

The entities required for an activity, such as copying data, must all be contained in the same data factory. A data factory does not store any data—instead, the entities reference the systems where the data is stored, such as an on-premises server or an Azure SQL Data Warehouse server.

To create a data factory in the Azure portal, click **New**, click **Data + Analytics**, and then click **Data Factory**. You must then specify the following parameters, and then click **Create**:

- **Name:** the data factory name must be globally unique and be at least three characters long.
- **Resource Group Name:** you can either create a new resource group or select an existing one.
- **Region:** choose the location nearest to you.

- Factory contains entities for activities
- Specify
 - Name
 - Resource group name
 - Region

Setting Up a Data Gateway for the On-Premises Server

For the data factory to access data in an on-premises server, you must set up a gateway to that server. Use the following procedure to set up a data gateway:

1. In the Azure portal, using the server where you are adding the data gateway, click **All resources**, click the data factory where you are adding the gateway, click **Author and Deploy**, click **More commands**, and then click **New integration runtime (gateway)**.
2. Specify the name of the gateway, click **OK**, and then click **Install directly on this computer**.
3. When the download completes, select the default options on the dialog boxes to install the gateway on your computer.

- Access data factory from on-premises server
- Create new gateway
- Install on computer

 **Note:** If you need to install the gateway on multiple computers, you can download and install the integration runtime package manually. When you do this, you must provide the security key generated for your gateway when you configure the integration runtime.

Setting up a Linked Service

A linked service defines the system where the source or destination data is situated, such as an on-premises SQL Server database, an Azure SQL Data Warehouse database, or an Azure Blob storage account. To copy data from an on-premises database to the cloud, you must create a linked service that references the on-premises database and a linked service that references the cloud database. A linked service is set up by performing the following steps:

1. In the Azure portal, click **All resources**, click the data factory where you are adding the linked service, click **Author and Deploy**, click **New data set**, and then click the type of the service to which you wish to link.
2. A JSON template script will be generated. Change the parameters of the script as follows, and then click **Deploy**:
 - o **name**: the name you wish to give to your linked service.
 - o **connectionString**: the connection string of your linked service. If the linked service refers to a SQL server using Windows security, in the connection string, set **Integrated Security to True**, and delete the **User ID** and **Password** information from the connection string.
 - o **gatewayName**: if you are creating a linked service for an on-premises database, you must specify the name of the data gateway for the server where the database is contained.
 - o **userName**: you must specify a username if you are using Windows authentication for a SQL Server database. Otherwise, delete the **userName** line.
 - o **password**: you must specify a password if you are using Windows authentication for a SQL Server database. Otherwise, delete the **password** line.

- New data store
- Edit parameters in JSON script
 - name
 - connectionString
 - Integrated Security
 - User ID
 - Password
 - gatewayName
 - userName: Use for Windows authentication
 - password: Use for Windows authentication

Setting Up a Dataset

A dataset references the structure of the source or destination data. Types of structures that you can reference include tables, documents, files, and folders. Datasets also contain information about the availability frequency of the data slices. To copy data from an on-premises database to the cloud, you must create a dataset that references a table in the on-premises database and a dataset that references a table in the cloud database. A dataset is set up by performing the following steps:

1. In the Azure portal, click **All resources**, click the data factory where you are adding the dataset, click **Author and Deploy**, click **New data set**, and then click the type of dataset you are creating (note that this is the same blade that you use to create a linked service).

- New dataset
- Edit parameters in JSON script
 - name
 - linkedServiceName
 - tableName
 - frequency
 - interval

2. A JSON template script will be generated. Change the parameters of the script as follows, and then click **Deploy**:
 - **name**: the name you wish to give to your dataset.
 - **linkedServiceName**: the name of the linked service that refers to the system where the data is contained.
 - **tableName**: the name of the table if your dataset refers to a database table.
 - **frequency**: the time interval for the data slice availability. You can use Minute, Hour, Day, Week, or Month.
 - **interval**: the interval for the time frequency specified. For example, to set an availability of 15 minutes, type **Minute** for the frequency, and type **15** for the interval.

Setting Up a Pipeline Activity to Copy Data

A pipeline contains at least one activity to achieve a particular task. To copy data from an on-premises database to the cloud, you must create a pipeline with a copy activity that specifies a dataset referencing a table in the on-premises database as the input, and specifies a dataset referencing a table in the cloud database as the output. To do this, perform the following steps:

1. In the Azure portal, click **All resources**, click the data factory you are adding the gateway to, click **Author and Deploy**, click **More commands**, and then click **New pipeline**. A JSON template script will be generated as shown here:

- New pipeline
- Edit parameters in JSON script
 - name
 - start
 - end
- Add activity script

```
{
  "name": "PipelineTemplate",
  "properties": {
    "description": "<Enter the Pipeline description here>",
    "activities": [],
    "start": "<The start date-time of the duration in which data processing will occur or the data slices will be processed. Example : 2014-05-01T00:00:00Z>",
    "end": "<The end date-time of the duration in which data processing will occur or the data slices will be processed. Example: 2014-05-05T00:00:00Z>"
  }
}
```

2. Change the parameters of the script as specified here:
 - **name**: the name you wish to give to your pipeline.
 - **start**: the date and time your pipeline is to start running.
 - **end**: the date and time your pipeline is to stop running.

- For the activities parameter, you must define any activities the pipeline is to contain. The following activity will copy data from an on-premises SQL dataset to an Azure SQL Data Warehouse dataset:

```
"activities": [
  {
    "type": "Copy",
    "typeProperties": {
      "source": {
        "type": "SqlSource",
        "sqlReaderQuery": "select * from <Enter Source Table>"
      },
      "sink": {
        "type": "SqlSink",
        "writeBatchSize": 0,
        "writeBatchTimeout": "00:00:00"
      }
    },
    "inputs": [
      {
        "name": "<Enter Source Dataset>"
      }
    ],
    "outputs": [
      {
        "name": "<Enter Destination Dataset>"
      }
    ],
    "policy": {
      "timeout": "01:00:00",
      "concurrency": 1,
      "executionPriorityOrder": "NewestFirst",
      "style": "StartOfInterval"
    },
    "scheduler": {
      "frequency": "<Enter Activity Frequency>",
      "interval": <Enter Activity Interval>
    },
    "name": "<Enter Activity Name>",
    "description": "Copy data from on-prem SQL server to AzureSQLDW"
  }
]
```

Change *<Enter Source Table>* to the name of the table in the on-premises database.

Change *<Enter Source Dataset>* to the name of the dataset that references the table in the on-premises database.

Change *<Enter Destination Dataset>* to the name of the dataset that references the table in the cloud database.

Change *<Enter Activity Frequency>* to the time interval for the activity. Minute, Hour, Day, Week, or Month can be used.

Change *<Enter Activity Interval>* to the interval for the time frequency specified. For example, to set the activity to run every 15 minutes, type **Minute** for the frequency, and type **15** for the interval.

Change *<Enter Activity Name>* to the name you wish to give to the activity.

- After you have entered the activity details, click **Deploy**.

Data Factory Diagram

The data factory diagram shows your pipelines and datasets, and helps you to monitor their progress. To view the diagram, in the Azure portal, click **All resources**, click the data factory you are viewing, and then click **Diagram**. The diagram will show the data flow from the source datasets to the pipeline input, and the dataflow from the pipeline output to the destination datasets.

Double-click the pipeline from the diagram to view the properties of the pipeline, including the activities of the pipeline, and the datasets consumed and produced by the pipeline.

Double-click the datasets on the diagram to view the data slices scheduled, and monitor their progress.

Question: In this lesson, the pipeline activity that copies data from the on-premises database to the Azure SQL Data Warehouse database does not contain any connection strings for the databases. Why is it not necessary to specify the connection strings?

- Shows data flow
- Pipeline properties
 - Activities
 - Datasets
- Dataset properties
 - Data slices

Lab: Implement an Azure SQL Data Warehouse

Scenario

A data warehouse containing food orders might need to rapidly expand; however, this is not definite, so the executive board have decided not to purchase the hardware to support the expanded database and instead wish to implement an Azure SQL Data Warehouse. You have been asked to implement a preliminary test system that uses a cutdown version of the data warehouse.

In this lab, you will create an Azure SQL Data Warehouse database on a new Azure logical server. You will then use the Data Warehouse Migration Utility to migrate the data in the **FoodOrdersDW** database on the MIA-SQL server to the new Azure SQL Data Warehouse database. Finally, you will test the capabilities of the Azure Data Factory by setting up a scheduled pipeline activity that copies data from another database on the MIA-SQL server to a table in the new Azure SQL Data Warehouse database.

Objectives

After completing this lab, you will be able to:

- Implement an Azure SQL Data Warehouse.
- Migrate data from a SQL Server database to an Azure SQL Data Warehouse database.
- Use the Azure Data Factory to schedule activities to copy data.

Estimated Time: 60 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Create an Azure SQL Data Warehouse Database

Scenario

So that you can migrate your **FoodOrdersDW** database to Azure SQL Data Warehouse, you must first create a new Azure SQL Data Warehouse database using the Azure portal. As part of the process, you must also create a logical server for the database and configure the firewall settings so that you can connect to the server using SSMS on your computer. You will decide the name of the logical server, which must be globally unique.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Create a Data Warehouse and Server
3. Configure the Azure Firewall
4. Connect to Azure SQL Data Warehouse

► Task 1: Prepare the Lab Environment

1. Ensure that the **MT17B-WS2016-NAT**, **20767C-MIA-DC**, **20767C-MIA-CLI**, and **20767C-MIA-SQL** virtual machines are running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab05\Starter** folder as Administrator.

► Task 2: Create a Data Warehouse and Server

1. Open the Azure portal <https://portal.azure.com/> with your Azure pass or Microsoft account.
2. Create a new SQL Data Warehouse named **FoodOrdersDW** in a resource group named **FoodOrdersRG**. Add the data warehouse to a new server with a name of your choosing. This must be unique throughout the whole Azure service, so cannot be specified here. A suggested format is **20767c**, followed by your initials, followed by today's date. For example, **20767ccn15mar18**.
3. The new server should use the admin login name **foodordersadmin** with the password **Pa55w.rd**
4. Set the performance tier to **DW100**.

► Task 3: Configure the Azure Firewall

- Amend the firewall configuration for the server created in the previous task to allow connections from your client IP address.

► Task 4: Connect to Azure SQL Data Warehouse

- Connect to the Azure server you created with SQL Server Management Studio.

Results: After this exercise, you should have an Azure SQL Data Warehouse database called **FoodOrdersDW** on a logical server with a name you specified. The firewall settings of the server should be set so that you can connect to the server using SSMS on your computer.

Exercise 2: Migrate to an Azure SQL Data Warehouse Database

Scenario

Now that you have implemented the **FoodOrdersDW** database on Azure SQL Data Warehouse, you are ready to migrate the data in the **FoodOrdersDW** database on the MIA-SQL server to the Azure SQL Data Warehouse database by using the Data Warehouse Migration Utility. You must first check the compatibility of the legacy data with Azure SQL Data Warehouse, migrate the schema, and then migrate the data in the tables.

The main tasks for this exercise are as follows:

1. Install the Data Warehouse Migration Utility
2. Check Compatibility of the Legacy Database
3. Migrate the Schema
4. Migrate the Data
5. Verify That the Data Has Been Migrated Successfully

► Task 1: Install the Data Warehouse Migration Utility

- The **Data Warehouse Migration Utility.msi** file in the **D:\Labfiles\Lab05\Starter\Installdwmigrationutility** folder has been downloaded from the Internet. Use this file to install the Data Warehouse Migration Utility.

► **Task 2: Check Compatibility of the Legacy Database**

1. Use SQL Server Management Studio to review the data in all of the tables in the **FoodOrdersDW** database in the MIA-SQL database engine.
2. Use the Data Warehouse Migration Utility to check the compatibility of the **FoodOrdersDW** database on the MIA-SQL server with Azure SQL Data Warehouse.
3. Switch to the **20767C-MIA-CLI** virtual machine, and log in as **Student** with the password **Pa55w.rd**.
4. Open the **FoodOrdersDW - DatabaseCompatibilityReport.xlsx**, in the **\\MIA-SQL\D\$\Labfiles\Lab05\Starter\Ex2** folder, review the data in the Excel spreadsheet (there should be two items assessed to have low impact), and then close Excel without saving changes.
5. Switch to the **20767C-MIA-SQL** virtual machine.

► **Task 3: Migrate the Schema**

- Use the Data Warehouse Migration Utility to migrate the schema for all the tables in the **FoodOrdersDW** database on the **MIA-SQL** server to the **FoodOrdersDW** data warehouse on the Azure server you created previously.

► **Task 4: Migrate the Data**

- Use the Data Warehouse Migration Utility to migrate the data for all the tables in the **FoodOrdersDW** database on the **MIA-SQL** server to the **FoodOrdersDW** data warehouse on the Azure server you created previously.

► **Task 5: Verify That the Data Has Been Migrated Successfully**

- Use SQL Server Management Studio to verify that the data has been successfully migrated to the **FoodOrdersDW** data warehouse on the Azure server you created previously.

Results: After this exercise, all of the tables in the **FoodOrdersDW** database on the **MIA-SQL** server should be copied to the **FoodOrdersDW** database that you created in Azure SQL Data Warehouse, in the previous exercise. All of the tables should contain data, apart from the **FactCallCentreOrders** table.

Exercise 3: Copy Data with the Azure Data Factory

Scenario

In the legacy data warehouse, the latest data for the call centre orders is currently entered into a table in a separate database called **RecentCallCentreOrders**. You are considering a scenario where you have a process that regularly copies the data in this table to a table in the Azure SQL Data Warehouse. As a prototype, you want to schedule a pipeline activity in the Azure Data Factory that will copy the data in the **CallCentreOrders** table, in the **RecentCallCentreOrders** database on the **MIA-SQL** server, to the **FactCallCentreOrders** table in the Azure SQL Data Warehouse database.

The main tasks for this exercise are as follows:

1. Create an Azure Data Factory
2. Create a Data Management Gateway
3. Create a Linked Server for the On-Premises SQL Server Database
4. Create a Linked Server for the Azure SQL Data Warehouse Database
5. Create a Dataset for a Table in the On-Premises SQL Server Database
6. Create a Dataset for a Table in the Azure SQL Data Warehouse Database
7. Create a Pipeline Activity to Copy the Data
8. Check the Progress of the Pipeline

► Task 1: Create an Azure Data Factory

- Create a new data factory in the **FoodOrdersRG** resource group. The name of the data factory must be unique throughout the whole Azure service, so cannot be specified here. A suggested format is **20767c<your initials> <today's date>DF**.

► Task 2: Create a Data Management Gateway

1. In the data factory you have created, create a Data Management Gateway named **MIA-SQLGW**.
2. Configure the Data Management Gateway so that it is installed on your computer; download and install the Integration Runtime manually, but accept the default options.
3. When the Microsoft Integration Runtime Configuration Manager runs, provide the **key1** key of the **MIA-SQLGW** gateway from the Azure portal, and verify that the gateway connects without any errors.

► Task 3: Create a Linked Server for the On-Premises SQL Server Database

- In the data factory you have created, create a **SQL Server table** linked service named **RecentCallCentreOrdersSQLServerLS** that uses the **MIA-SQLGW** gateway to access the **RecentCallCentreOrders** database in the **MIA-SQL** server. For **userName** and **password**, use the details you use to log on to 20767C-MIA-SQL. Remember that, in JSON, you must replace the **** character with ****.

► Task 4: Create a Linked Server for the Azure SQL Data Warehouse Database

- In the data factory you have created, create an **Azure SQL Data Warehouse** linked service named **FoodOrdersAzureSqlDWLS** that connects to the **FoodOrdersDW** Azure SQL Data Warehouse database you created earlier.

► **Task 5: Create a Dataset for a Table in the On-Premises SQL Server Database**

- In the data factory you have created, create a **SQL Server table** dataset named **CallCentreOrdersSQLServerDS** that references the **CallCentreOrders** table using the **RecentCallCentreOrdersSqlServerLS** linked service. Set the availability to 15 minutes.

► **Task 6: Create a Dataset for a Table in the Azure SQL Data Warehouse Database**

- In the data factory you have created, create an **Azure SQL Data Warehouse** dataset named **FactCallCentreOrdersAzureSQLDWDS** that references the **FactCallCentreOrders** table using the **FoodOrdersAzureSqlDWLS** linked service. Set the availability to 15 minutes.

► **Task 7: Create a Pipeline Activity to Copy the Data**

1. In the data factory you have created, create a pipeline named **FoodOrdersPL**. Set the pipeline to start at the most recent exact quarter past the hour, and to finish a quarter of an hour after it starts.
2. For the pipeline activities, add the following JSON script:

```
"activities": [
  {
    "type": "Copy",
    "typeProperties": {
      "source": {
        "type": "SqlSource",
        "sqlReaderQuery": "select * from <Enter Source Table>"
      },
      "sink": {
        "type": "SqlSink",
        "writeBatchSize": 0,
        "writeBatchTimeout": "00:00:00"
      }
    },
    "inputs": [
      {
        "name": "<Enter Source Dataset>"
      }
    ],
    "outputs": [
      {
        "name": "<Enter Destination Dataset>"
      }
    ],
    "policy": {
      "timeout": "01:00:00",
      "concurrency": 1,
      "executionPriorityOrder": "NewestFirst",
      "style": "StartOfInterval"
    },
    "scheduler": {
      "frequency": "<Enter Activity Frequency>",
      "interval": <Enter Activity Interval>
    },
    "name": "<Enter Activity Name>",
    "description": "Copy data from on-prem SQL server to AzureSQLDW"
  }
]
```

You can copy this script from the **ActivityText.txt** file in the **D:\Labfiles\Lab05\Starter\Ex3** folder.

3. Edit the script so that the activity is named **CopyFromSQLtoAzureSQLDW**, the SQL statement references the **CallCentreOrders** table, and the activity will copy data from the **CallCentreOrdersSQLServerDS** dataset to the **FactCallCentreOrdersAzureSQLDWDS** dataset every 15 minutes, and then deploy the pipeline.

► **Task 8: Check the Progress of the Pipeline**

1. In the data factory you have created, view the diagram, and then view the data slices of the two datasets.
2. In the **FactCallCentreOrdersAzureSQLDWDS** dataset, monitor the status of the data slice that has the slice start time and slice end time that you specified for the pipeline activity in the previous task. When the pipeline activity has completed for that data slice, the status will show as **Ready**.
3. Use SQL Server Management Studio to verify that the data has been successfully copied to the **FactCallCentreOrders** table in the **FoodOrdersDW** database on the Azure server you created earlier.
4. Pause the **FoodOrdersDW** database and log out of Azure.

Results: After this exercise, the data in the **CallCentreOrders** table, in the **RecentCallCentreOrders** database on the **MIA-SQL** server, should be copied to the **FactCallCentreOrders** table in the Azure SQL Data Warehouse database.

Module Review and Takeaways

In this module, you have learned how to use Azure SQL Data Warehouse and about the issues to consider when migrating to Azure SQL Data Warehouse.

Review Question(s)

Question: How do you add other logins to an Azure SQL Data Warehouse logical server?

MCT USE ONLY. STUDENT USE PROHIBITED

Module 6

Creating an ETL Solution

Contents:

Module Overview	6-1
Lesson 1: Introduction to ETL with SSIS	6-2
Lesson 2: Exploring Source Data	6-7
Lesson 3: Implementing Data Flow	6-14
Lab: Implementing Data Flow in an SSIS Package	6-28
Module Review and Takeaways	6-34

Module Overview

Successful data warehousing solutions rely on the efficient and accurate transfer of data from the various data sources in the business. This is referred to as an extract, transform, and load (ETL) process. ETL is a core process that is required in any data warehousing project.

This module discusses considerations for implementing an ETL process, and then focuses on Microsoft® SQL Server® Integration Services (SSIS) as a platform for building ETL solutions.

Objectives

After completing this module, you will be able to:

- Describe the key features of SSIS.
- Explore source data for an ETL solution.
- Implement a data flow using SSIS.

Lesson 1

Introduction to ETL with SSIS

There are several ways to implement an ETL solution, but SSIS is the primary ETL tool for SQL Server. Before using SSIS to implement an ETL solution, it is important to understand some of the key features and components that SSIS provides.

This lesson describes possible ETL solution options, and then introduces SSIS.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the options for ETL.
- Describe the key features of SSIS.
- Describe the high level architecture of an SSIS project.
- Identify key elements of the SSIS design environment.
- Describe strategies for upgrading SSIS solutions from previous versions of SQL Server.

Options for ETL

An ETL solution generally involves the transfer of data from one or more data sources to a destination, often transforming the data structure and values in the process. There are several tools and technologies you can use to accomplish this task, each having specific strengths and weaknesses that you should take into account when choosing the approach for your ETL solution.

- SQL Server Integration Services
- The Import and Export Data Wizard
- Transact-SQL
- The bcp utility
- Replication

The following list describes some common techniques:

- **SQL Server Integration Services.** This is the primary platform for ETL solutions that are provided with SQL Server, and generally offers the most flexible way to implement an enterprise ETL solution.
- **Import and Export Data Wizard.** This wizard is included with the SQL Server management tools, and provides a simple way to create an SSIS-based data transfer solution. You should consider using the Import and Export Data Wizard when your ETL solution requires only a few, simple data transfers that do not include any complex transformations in the data flow.
- **Transact-SQL.** Transact-SQL is a powerful language that you can use to implement complex data transformations when extracting, inserting, or modifying data. Most ETL solutions include some Transact-SQL logic combined with other technologies. In some scenarios, such as when data sources and destinations are co-located, you can implement a complete ETL solution by using only Transact-SQL queries.

- **Bulk copy program (bcp) utility.** This utility provides an interface based on the command line for extracting data from, and inserting data into, SQL Server. The bcp utility is a versatile tool you can use to create scheduled data extractions and insertions, but its relatively complex syntax and console-based operation make it difficult to create a manageable, enterprise-scale ETL solution on its own.
- **Replication.** SQL Server includes built-in replication functionality that you can use to synchronize data across SQL Server instances. You can also include other relational data sources, such as Oracle databases, in a replication solution. Replication is a suitable technology for ETL when all data sources are supported in a replication topology, and the data requires minimal transformations.

What Is SSIS?

SSIS is an extensible platform for building complex ETL solutions. It is included with SQL Server and consists of a Microsoft Windows® service that manages the execution of ETL workflows, along with tools and components for developing them. The SQL Server setup program installs the SSIS service when you select **Integration Services** on the **Feature Selection** page of the wizard. After you have installed SSIS, you must use the **DCOM Configuration** tool (Dcomcnfg.exe) to grant specific permission to users who need to access the SSIS service.

- A platform for ETL operations
- Installed as a feature of SQL Server
- Control flow engine:
 - Runtime resources and operational support for data flow
- Data flow engine:
 - Pipeline architecture for buffer-oriented rowset processing

For instructions on how to grant users access to the SSIS service, see the SQL Server Technical Documentation:

 **Grant Permissions to Integration Services Service**

<http://aka.ms/m1b8cb>

The SSIS Windows service is primarily a control flow engine that manages the execution of task workflows, which are defined in packages, and can be performed on demand or at scheduled times. When you are developing an SSIS package, the task workflow is referred to as the control flow.

The control flow can include a special type of task to perform data flow operations. SSIS executes these tasks using a data flow engine that encapsulates the data flow in a pipeline architecture. Each step in the Data Flow task operates in sequence on a rowset of data as it passes through the pipeline. The data flow engine uses buffers to optimize the data flow rate, resulting in a high-performance ETL solution.

In addition to the SSIS Windows service, SSIS includes:

- **SSIS Designer.** A graphical design interface for developing SSIS solutions in the Microsoft Visual Studio® development environment. Typically, you start the SQL Server Data Tools (SSDT) application to access this.
- **Wizards.** Graphical utilities you can use to quickly create, configure, and deploy SSIS solutions.
- **Command-line tools.** Utilities you can use to manage and execute SSIS packages.

SSIS Projects and Packages

An SSIS solution usually consists of one or more SSIS projects, each containing at least one SSIS package.

SSIS Projects

In SQL Server, a project is the unit of deployment for SSIS solutions. You can define project-level parameters so that users can specify run-time settings, and project-level connection managers that reference data sources and destinations used in package data flows. You can then deploy projects to an SSIS catalog in a SQL Server instance, and configure project-level parameter values and connections as appropriate for execution environments. You can use SSDT to create, debug, and deploy SSIS projects.

SSIS Packages

A project contains one or more packages, each defining a workflow of tasks to be executed. The workflow of tasks is referred to as its control flow. A package control flow can include one or more Data Flow tasks, each of which encapsulates its own pipeline. You can include package-level parameters so that the package receives dynamic values at run time.

In previous SSIS releases, deployment was managed at the package level. In SQL Server, you can still deploy individual packages in a package deployment model.

 **Note:** Deployment of SSIS solutions is discussed in more detail in Module 13: *Deploying and Configuring SSIS Packages*.

The SSIS Design Environment

You can use SSDT to develop SSIS projects and packages. SSDT is an add-in for Visual Studio and provides a graphical development environment for business intelligence (BI) solutions. When you create an Integration Services project, the design environment includes the following elements:

- **Solution Explorer.** A pane in the SSDT user interface you can use to create and view project-level resources, including parameters, packages, data connection managers, and other shared objects. A solution can contain multiple projects, where each project is shown in Solution Explorer.
- **The Properties pane.** A pane in the SSDT user interface you can use to view and edit the properties of the currently-selected object.
- **The Control Flow design surface.** A graphical design surface in SSIS Designer you can use to define the workflow of tasks for a package.

- **Project Deployment Model:**
 - Multiple packages are deployed in a single project
- **Package Deployment Model:**
 - SSIS packages are deployed and managed individually

- Solution Explorer
- Properties pane
- Control Flow design surface
- Data Flow design surface
- Parameters tab
- Event Handlers design surface
- Package Explorer
- Connection Managers pane
- Variables pane
- SSIS Toolbox

- **The Data Flow design surface.** A graphical design surface in SSIS Designer you can use to define the pipeline for a Data Flow task within a package.
- **The Parameters Tab.** A tab you can use to define the package level parameters.
- **The Event Handlers design surface.** A graphical design surface in SSIS Designer that you can use to define the workflow for an event handler within a package.
- **Package Explorer.** A tree view providing access to the components in a package.
- **The Connection Managers pane.** A list of the connection managers used in a package.
- **The Variables pane.** A list of variables used in a package. You can display this pane by clicking the **Variables** button at the upper right of the design surface.
- **The SSIS Toolbox.** A collection of components you can add to a package control flow or data flow. You can display this pane by clicking the **SSIS Toolbox** button at the upper right of the design surface or by clicking **SSIS Toolbox** on the **SSIS** menu. Note that this pane is distinct from the standard SSDT Toolbox pane.

Upgrading from Previous Versions

If you have developed ETL solutions by using SSIS in SQL Server 2005, 2008, 2012, or 2014, you should consider how you will include them in an SSIS solution for the latest build of SQL Server.

Upgrading Packages in a Project

To upgrade SSIS packages in an SSIS project, you can open the project file using SSDT. If the project contains any packages that were built using a previous version of SQL Server, the SSIS Package Upgrade Wizard will automatically start. You can use this wizard to select the packages that you want to upgrade; it will then upgrade them. Any package upgrades made using the SSIS Package Upgrade Wizard are permanent whether or not you subsequently save the project file. You are advised to back up any packages before you upgrade them.

- Upgrading Packages in a Project:
 - Open project file in SQL Server Data Tools
 - SSIS Package Upgrade Wizard will automatically be activated
- Upgrading a Single Package:
 - Open package file in SQL Server Data Tools
 - Package will automatically be upgraded
- Scripts:
 - Migrated VSA scripts are automatically updated to VSTA
 - Microsoft ActiveX scripts are no longer supported and must be replaced

Upgrading a Single Package

To upgrade a single SSIS package, you can open the package file using SSDT. If the package was built using a previous version of SQL Server, SSDT will automatically upgrade the file to the latest version. If you subsequently save the package file, the upgrade will be made permanent. If you close the file without saving, the file will revert to the previous version. You are advised to back up the package before upgrading it.

Scripts

SSIS packages can include script tasks to perform custom actions. In previous releases of SSIS, you could implement scripted actions by including a Microsoft ActiveX® script task, written in Microsoft Visual Basic® Scripting Edition (VBScript). You could also employ a script task written for the .NET Visual Studio for Applications, or VSA runtime, in a control flow. In the latest build of SQL Server, the ActiveX script task is no longer supported, and any VBScript-based custom logic must be replaced. In addition, the SQL Server script task now uses the Visual Studio Tools for Applications (VSTA) runtime, which differs from the

VSA runtime used in previous releases. When you upgrade a package that includes a script component, the script is automatically updated for the VSTA runtime.

Verify the correctness of the statement by placing a mark in the column to the right.

Statement	Answer
The Connection Managers pane in SSDT shows all of the connection managers in the current project. True or false?	

Lesson 2

Exploring Source Data

Now that you understand the basic architecture of SSIS, you can start planning the data flows in your ETL solution. However, before you start implementing an ETL process, you should explore the existing data in the sources that your solution will use. By gaining a thorough knowledge of the source data on which your ETL solution will be based, you can design the most effective SSIS data flows for transferring the data, and anticipate any quality issues you may need to resolve in your SSIS packages.

This lesson discusses the value of exploring source data, in addition to describing techniques for examining and profiling it.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the value of exploring source data.
- Examine an extract of data from a data source.
- Profile source data by using the Data Profiling SSIS task.

Why Explore Source Data?

The design and integrity of your data warehouse ultimately rely on the data it contains. Before you can design an appropriate ETL process to populate the data warehouse, you must have a thorough knowledge of the source data that your solution will consume.

Specifically, you need to understand:

- The business entities that are represented by the source data, and their attributes. For example, the specific attributes that fully describe a product or a customer entity might be stored in multiple columns, tables, or even databases across the organization.
- How to interpret data values and codes. For example, does a value of **1** in an **InStock** column in a **Products** table mean that the company has a single unit in stock—or does **1** indicate the value **true**, meaning that there is an unspecified quantity of units in stock?
- The relationships between business entities, and how those relationships are modeled in the data sources.

In addition to understanding the data modeling of the business entities, you should also examine source data to help identify:

- Column data types and lengths for specific attributes that will be included in data flows. For example, what maximum lengths exist for string values? What formats are used to indicate date, time, and numeric values?
- Data volume and sparseness. For example, how many rows of sales transactions are typically recorded in a single trading day? Are there any attributes that frequently contain null values?

- Understand business data:
 - What business entities are represented
 - How to interpret values and codes
 - Relationships between business entities
- Examine data for:
 - Column data types and lengths
 - Data volume and sparseness
 - Data quality issues

- Data quality issues. For example, are there any obvious data entry errors? Are there commonly-used values that are synonyms for one another?

Finding the answers to questions like these, before you implement the ETL solution, can help you anticipate data flow problems and proactively design effective solutions for them.

Examining Source Data

To understand the data that your solution will consume, you must first extract it from its source, and then examine it using an appropriate application, such as Microsoft Excel®.

The following list describes some of the approaches you can use to extract the data:

- Run queries against data sources in Microsoft SQL Server Management Studio (SSMS). To view the data in another application, you can copy and paste the results as required.
- Create an SSIS package with a data flow that extracts a sampling of data or a row count for a specific data source.
- Use the Import and Export Data Wizard to extract a data sample. You must choose a format for the sample that is compatible with the application you are intending to use to view the data. For example, if you are intending to use Excel, you should extract the data in a comma-delimited text format.

- Extract a sample of data:
 - Run queries in SSMS
 - Create an SSIS package that extracts a sample of data
 - Use the Import and Export Data Wizard
- Examine the data using an appropriate application such as Excel

After extracting the sample data, you should examine it. One of the most effective tools for examining data is Excel. By using Excel, you can:

- Sort the data by columns.
- Apply column filters to help identify the range of values used in a particular column.
- Use formulas to calculate minimum, maximum, and average values for numerical columns.
- Search the data for specific string values.

Demonstration: Exploring Source Data

In this demonstration, you will see how to:

- Extract data by using the Import and Export Data Wizard.
- Explore the data.

Demonstration Steps

Extract Data with the Import and Export Data Wizard

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod06** folder, right-click **Setup.cmd**, and then click **Run as administrator**.
3. If you are prompted to confirm, click **Yes**, and then wait for the batch file to complete.

4. Click **Search Windows**, type **Import and Export**, and then click **SQL Server Import and Export Data (64-bit)**.
5. On the **Welcome to SQL Server Import and Export Wizard** page, click **Next**.
6. On the **Choose a Data Source** page, set the following options, and then click **Next**:
 - o **Data source:** SQL Server Native Client 11.0
 - o **Server name:** MIA-SQL
 - o **Authentication:** Windows Authentication
 - o **Database:** ResellerSales
7. On the **Choose a Destination** page, select the following options, and then click **Next**:
 - o **Destination:** Flat File Destination
 - o **File name:** D:\Demofiles\Mod06\Top 500 Resellers.csv
 - o **Text qualifier:** " (this is used to enclose exported text values in quotation marks. This is required because some European address formats include a comma, and these must be distinguished from the commas that are used to separate each column value in the exported text file)
8. On the **Specify Table Copy or Query** page, select **Write a query to specify the data to transfer**, and then click **Next**.
9. On the **Provide a Source Query** page, enter the following Transact-SQL code, and then click **Next**:

```
SELECT TOP 500 * FROM Resellers
```
10. On the **Configure Flat File Destination** page, click **Next**.
11. On the **Save and Run Package** page, select only **Run immediately**, and then click **Next**.
12. On the **Complete the Wizard** page, click **Finish**.
13. When the data extraction has completed successfully, click **Close**.

Explore Source Data

1. Start Notepad, and open the file **D:\Demofiles\Mod06\Top 500 Resellers.csv**.
2. Scroll through the data to examine it. Note that many of the addresses for resellers in France include a comma. If no text qualifier had been selected in the Import and Export Data Wizard, these commas would have created additional columns in these rows, making the data difficult to examine as a table.
3. Close Notepad without saving any changes.

Profiling Source Data

In addition to examining samples of source data, you can use the Data Profiling task in an SSIS package to obtain statistics about the data. This can help you to understand the structure of the data that you will extract and identify columns where null or missing values are likely. Profiling source data can help you plan effective data flows for your ETL process.

You can include multiple profile requests in a single instance of the Data Profiling task. The following types of profile request are available:

- **Candidate Key** determines whether you can use a column (or set of columns) as a key for the selected table.
- **Column Length Distribution** reports the range of lengths for string values in a column.
- **Column Null Ratio** reports the percentage of null values in a column.
- **Column Pattern** identifies regular expressions that are applicable to the values in a column.
- **Column Statistics** reports statistics such as minimum, maximum, and average values for a column.
- **Column Value Distribution** reports the groupings of distinct values in a column.
- **Functional Dependency** determines if the value of a column is dependent on the value of other columns in the same table.
- **Value Inclusion** reports the percentage of time that a column value in one table matches a column in another.

- Use the Data Profiling task in SSIS to report data statistics:
 - Candidate key
 - Column length distribution
 - Column null ratio
 - Column pattern
 - Column statistics
 - Column value distribution
 - Functional dependency
 - Value inclusion
- View the profile in the Data Profile Viewer

The Data Profiling task gathers the requested profile statistics and writes them to an XML document. You can save this as a file for later analysis or write it to a variable for programmatic analysis in the control flow.

You can also use the Data Profile Viewer to view the profile statistics. This is available as a stand-alone tool in which you can open the XML file that the Data Profiling task generates. Alternatively, you can open the Data Profile Viewer window in SSDT, from the **Properties** dialog box of the Data Profiling task, while the package is running in the development environment.

Use the following procedure to collect and view data profile statistics:

1. Create an SSIS project that includes a package.
2. Add an ADO.NET connection manager for each data source that you want to profile.
3. Add the Data Profiling task to the control flow of the package.
4. Configure the Data Profiling task to specify:
 - The file or variable to which the resulting profile statistic should be written.
 - The individual profile requests that should be included in the report.
5. Run the package.
6. View the resulting profile statistics in the Data Profile Viewer.

Demonstration: Using the Data Profiling Task

In this demonstration, you will see how to:

- Use the Data Profiling task.
- View a Data Profiling report.

Demonstration Steps

Use the Data Profiling Task

1. Ensure you have completed the previous demonstration in this module.
2. Start SQL Server Data Tools.
3. On the **File** menu, point to **New**, and then click **Project**.
4. In the **New Project** dialog box, select the following values, and then click **OK**.
 - **Project Template:** Integration Services Project
 - **Name:** ProfilingDemo
 - **Location:** D:\Demofiles\Mod06
 - **Solution name:** ProfilingDemo
5. In the Solution Explorer pane, right-click **Connection Managers**, and then click **New Connection Manager**.
6. In the **Add SSIS Connection Manager** dialog box, click **ADO.NET**, and then click **Add**.
7. In the **Configure ADO.NET Connection Manager** dialog box, click **New**.
8. In the **Connection Manager** dialog box, enter the following values, and then click **OK**:
 - **Server name:** localhost
 - **Log on to the server:** Windows Authentication
 - **Select or enter a database name:** ResellerSales
9. In the **Configure ADO.NET Connection Manager** dialog box, verify that a data connection named **localhost.ResellerSales** is listed, and then click **OK**.
10. If the SSIS Toolbox pane is not visible, click the **Package.dtsx [Design]** pane, and then on the **SSIS** menu, click **SSIS Toolbox**.
11. In the SSIS Toolbox pane, in the **Common** section, double-click **Data Profiling Task** to add it to the Control Flow surface.
12. Double-click the **Data Profiling Task** icon on the Control Flow surface.
13. In the **Data Profiling Task Editor** dialog box, on the **General** page, in the **Destination** property drop-down list, click **<New File connection...>**.
14. In the **File Connection Manager Editor** dialog box, in the **Usage type** drop-down list, click **Create file**. In the **File** box, type **D:\Demofiles\Mod06\Reseller Sales Data Profile.xml**, and then click **OK**.
15. In the **Data Profiling Task Editor** dialog box, on the **General** page, set **OverwriteDestination** to **True**.
16. In the **Data Profiling Task Editor** dialog box, on the **Profile Requests** page, in the **Profile Type** drop-down list, click **Column Statistics Profile Request**, and then click the **RequestID** column.

17. In the Request Properties pane, set the following property values:
 - **ConnectionManager:** localhost.ResellerSales
 - **TableOrView:** [dbo].[SalesOrderHeader]
 - **Column:** OrderDate
18. In the **Data Profiling Task Editor** dialog box, click **OK**.
19. On the **File** menu, point to **Open**, and then click **Project/ Solution**.
20. In the **Open Project** dialog box, open the **Explore Reseller Sales.sln** file in the **D:\DemoFiles\Mod06\Profilerequests** folder, saving your changes to the previous solution if you are prompted.
21. In Solution Explorer, expand **SSIS Packages** and double-click **Package.dtsx**.
22. Double-click the **Data Profiling Task** icon on the Control Flow surface.
23. In the **Data Profiling Task Editor** dialog box, click the **Profile Requests** page, and then note the profile requests.
24. Click the **Column Length Distribution Profile Request** row and note the following settings in the Request Properties pane:
 - **ConnectionManager:** localhost.ResellerSales
 - **TableOrView:** [dbo].[Resellers]
 - **Column:** AddressLine1
25. Click the **Column Null Ratio Profile Request** row and note the following settings in the Request Properties pane:
 - **ConnectionManager:** localhost.ResellerSales
 - **TableOrView:** [dbo].[Resellers]
 - **Column:** AddressLine2
26. Click the **Value Inclusion Profile Request** row and note the following settings in the Request Properties pane:
 - **ConnectionManager:** localhost.ResellerSales
 - **SubsetTableOrView:** [dbo].[SalesOrderHeader]
 - **SupersetTableOrView:** [dbo].[PaymentTypes]
 - **InclusionColumns:**
 - **Subset side Columns:** PaymentType
 - **Superset side Columns:** PaymentTypeKey
 - **InclusionThresholdSetting:** None
 - **SupersetColumnsKeyThresholdSetting:** None
27. In the **Data Profiling Task Editor** dialog box, click **Cancel**.
28. On the **Debug** menu, click **Start Debugging**.

View a Data Profiling Report

1. When the Data Profiling task has completed, with the package still running, double-click **Data Profiling Task**, and then in the **Data Profiling Task Editor**, click **Open Profile Viewer**.
2. Maximize the Data Profile Viewer window and under the **[dbo].[SalesOrderHeader]** table, click **Column Statistics Profiles**. Review the minimum and maximum values for the **OrderDate** column.
3. Under the **[dbo].[Resellers]** table, click **Column Length Distribution Profiles**. Click the bar chart for any of the column lengths, and then click the **Drill Down** button to view the source data that matches the selected column length.
4. Close the Data Profile Viewer window, and then in the **Data Profiling Task Editor** dialog box, click **Cancel**.
5. On the **Debug** menu, click **Stop Debugging**, and then close SSDT, saving your changes if prompted.
6. On the Start screen, type **Data Profile**, and then click **SQL Server Data Profile Viewer**. When the **Data Profile Viewer** window appears, maximize it.
7. On the toolbar, click **Open**, and open **Reseller Sales Data Profile.xml** in the **D:\Demofiles\Mod06** folder.
8. Under the **[dbo].[Resellers]** table, click **Column Null Ratio Profiles** and view the null statistics for the **AddressLine2** column. Select the **AddressLine2** column, and then click the **Drill Down** button to view the source data.
9. Under the **[dbo].[SalesOrderHeader]** table, click **Inclusion Profiles** and review the inclusion statistics for the **PaymentType** column.
10. Select the inclusion violation for the payment type value of **0**, and then click the **Drill Down** button to view the source data. **Note:** The **PaymentTypes** table includes two payment types, using the value 1 for invoice-based payments and the value 2 for credit account payments. The Data Profiling task has revealed that for some sales, the value 0 is used, which might indicate an invalid data entry or may be used to indicate some other kind of payment that does not exist in the **PaymentTypes** table.
11. Close the Data Profile Viewer window.

Check Your Knowledge

Question	
<p>You are examining a database that has a Sales table containing a column named Product. It appears that this column references a column named ProductKey in a Product table. You want to use the Data Profiling task to check whether all values of the Product column in the Sales table have a matching value in the ProductKey column of the Product table. What type of profiling request can you use to achieve this?</p>	
Select the correct answer.	
<input type="checkbox"/>	Candidate Key
<input type="checkbox"/>	Column Null Ratio
<input type="checkbox"/>	Value Inclusion
<input type="checkbox"/>	Column Statistics

Lesson 3

Implementing Data Flow

After you have thoroughly explored the data sources for your data warehousing solution, you can start to implement an ETL process by using SSIS. This process will consist of one or more SSIS packages, with each containing one or more Data Flow tasks. Data flow is at the core of any SSIS-based ETL solution, so it is important to understand how you can use the components of an SSIS data flow pipeline to extract, transform, and load data.

This lesson describes the various components you can use to implement a data flow, and gives guidance for optimizing data flow performance.

Lesson Objectives

After completing this lesson, you will be able to:

- Create a connection manager.
- Add a Data Flow task to a package control flow.
- Add a Source component to a data flow.
- Add a Destination component to a data flow.
- Add Transformation components to a data flow.
- Optimize data flow performance.

Connection Managers

To extract or load data, an SSIS package must be able to connect to the data source and destination. In an SSIS solution, you define data connections by creating a connection manager for each data source component and each destination component used in the workflow. A connection manager encapsulates the following information, which it uses to make a connection to the data source or destination:

- The data provider to be used. For example, you can create a connection manager for a relational database by using an OLE DB or ADO.NET provider. Alternatively, you can create a connection manager for a text file by using a flat file provider.
- The connection string used to locate the data source. For a relational database, the connection string includes the network name of the database server and the name of the database. For a file, the file name and path must be specified.
- The credentials used to access the data source.

- A connection to a data source or destination:
 - Provider (for example, ADO.NET, OLE DB, or flat file)
 - Connection string
 - Credentials
- Project or package level:
 - Project-level connection managers can be shared across multiple packages
 - Package-level connection managers exist only in that package

You can create a connection manager at the project level or at the package level:

- Project-level connection managers are listed in Solution Explorer and can be shared across multiple packages in the same project. Use project-level connection managers when multiple packages need to access the same data source. To create a project-level connection manager, on the **Project** menu, click **New Connection Manager**.
- Package-level connection managers exist only in the package in which they are defined. Both project-level and package-level connection managers used by a package are shown in its Connection Managers pane in the SSIS Designer. To create a package-level connection manager, right-click in the Connection Managers pane and choose the type you want to create. Alternatively, create a new connection manager in the dialog box of a source, destination, or transformation component.

 **Note:** To create a new connection manager, you can use SSDT to select connection details that you have created previously, even if they relate to connection managers that do not exist in the current project or have been deleted.

The Data Flow Task

A package defines a control flow for actions that the SSIS run-time engine is to perform. A package control flow can contain several different tasks, and include complex branching and iteration, but the core of any ETL control flow is the Data Flow task.

To include a data flow in a package control flow, drag the Data Flow task from the SSIS Toolbox pane into the Control Flow surface. Alternatively, you can double-click the Data Flow task icon in the SSIS Toolbox pane to add the task to the design surface. After you add a Data Flow task to the control flow, you can rename it and set its properties in the Properties pane.

- The core control flow task in most SSIS packages
- Encapsulates a data flow pipeline
- Define the pipeline for the task on the Data Flow tab

 **Note:** This module focuses on the Data Flow task. Other control flow tasks will be discussed in detail in *Module 7: Implementing Control Flow in an SSIS Package*.

To define the pipeline for the Data Flow task, double-click the task. SSIS Designer will display a design surface where you can add data flow components. Alternatively, click the Data Flow tab in SSIS Designer, and then select the Data Flow task that you want to edit in the drop-down list displayed at the top of the design surface.

A typical data flow pipeline includes one or more data source components—transformation components that operate on the data as it flows through the pipeline—and one or more destination components for the data. The pipeline flow is defined by connecting the output from one component to the input of the next.

Data Source Components

The starting point for a data flow is a data source component. When you create a data source component, you must specify:

- The connection manager it will use to connect to the source data.
- The table, view, or query it will use to extract the data (if it is a database source).
- The columns to include in the output and pass to the next component in the data flow pipeline.

- The source of data for a data flow:
 - Connection manager
 - Table, view, or query
 - Columns that are included
- Many sources supported:
 - Database
 - File
 - Custom

The following tables describe the types of data source components that SSIS supports:

Databases	
ADO.NET	Any database for which an ADO.NET data provider is installed.
OLE DB	Any database for which an OLE DB provider is installed.
SQL Server	A SQL Server database.
Oracle	An Oracle database.
CDC Source	A SQL Server or Oracle database in which change data capture (CDC) has been enabled. CDC is discussed in <i>Module 9: Implementing an Incremental ETL Process</i> .

Files	
Excel	An Excel workbook.
Flat file	Data in a text file, such as comma-delimited text.
XML	A file that contains data in XML format.
Raw file	An SSIS-specific binary format file.

Other sources	
Script component	A custom source that is implemented as a script.
HDFS File Source	A Hadoop Distributed File System.
OData Source	An OData data feed.

In addition to those listed in the table, you can install connection managers to access SAP BI and Teradata data sources.

For further information on how to install these connection managers, see the SQL Server Technical Documentation:

Integration Services (SSIS) Connections

<http://aka.ms/a9dyg8>

You can also use DB2 data sources by installing Microsoft OLE DB Provider for DB2, and then using an OLE DB data source component within SSDT.

For more information on how to install Microsoft OLE DB Provider for DB2, see the SQL Server Technical Documentation:

Installing Data Provider

<http://aka.ms/ypbw43>

To add a data source component for SQL Server, Excel, a flat file, or Oracle to a data flow, drag the **Source Assistant** icon from the **Favorites** section of the SSIS Toolbox pane to the design surface and use the wizard to select or create a connection manager for the source component. For other data source components, drag the appropriate icon from the **Common** or **Other Sources** section of the SSIS Toolbox pane to the design surface, and then double-click the data source component on the design surface to define the connection, data, and output columns.

By default, the output from a data source component is represented as an arrow at the bottom of the data source component icon on the design surface. To create a data flow, you drag this arrow and connect it to the next component in the data flow pipeline, which might be a destination or a transformation component.

Data Destination Components

A destination component is an endpoint for a data flow. It has input columns, which are determined by the connection from the previous component in the data flow pipeline, but no output.

A destination component includes:

- The connection manager for the data store where it will insert data.
- The table or view into which it will insert the data (if it is a database destination).

- Endpoint for a data flow:
 - Connection manager
 - Table or view
 - Column mapping
- Multiple destination types:
 - Database
 - File
 - SQL Server Analysis Services
 - Rowset
 - Other

The types of destination components that SSIS supports are described in the following tables:

Databases	
ADO.NET	Any database for which an ADO.NET data provider is installed.
OLE DB	Any database for which an OLE DB provider is installed.
SQL Server	A SQL Server database.
SQL Server Compact	An instance of SQL Server Compact.
Oracle	An Oracle database.

Files	
Excel	An Excel workbook.
Flat file	A text file.
Raw file	An SSIS-specific binary format file.

SQL Server Analysis Services	
Data mining model training	Used to build data mining models for data analysis.
Dimension processing	Used to populate a dimension in an online analytical processing (OLAP) cube.
Partition processing	Used to populate a partition in an OLAP cube.

Rowsets	
DataReader	An ADO.NET DataReader interface that can be read by another application.
Recordset	An ADO Recordset interface that can be read by another application.

Other destinations	
Script component	A custom destination that is implemented as a script.
HDFS File Destination	A Hadoop Distributed File System.
Data Streaming Destination	Allows a linked server to be created with a tabular result set from the package that can be queried using SQL.

To add a SQL Server, Excel, or Oracle destination component to a data flow, drag the **Destination Assistant** icon from the **Favorites section** of the SSIS Toolbox pane to the design surface, and then use the wizard to select or create a connection manager. For other types of destination components, drag the

appropriate icon from the **Common** or **Other Destinations** section of the SSIS Toolbox pane to the design surface.

After you add a destination component to the data flow, connect the output from the previous component in the data flow to the destination component, double-click it, and then edit it to define:

- The connection manager and destination table (if it is a database destination) to use when loading the data.
- The column mappings between the input columns and the columns in the destination.

Mapping Columns

In the destination component, SSIS makes a best guess to automatically map columns with similar names in the source and destination. You need to manually map unmatched columns, and correct any wrong guesses by SSIS. The destination component shows a red error icon when columns are unmapped. In the source component, you can change the names of the columns passing into the data flow, so it might be helpful to alter names to match the destination columns. You can also use the Resolve References Editor to align the input and output columns between components. The OLE DB data source enables you to cancel columns that are not needed. This is not available in other data sources, so ensure you only include essential columns to keep the performance of the data flow as fast as possible—there is no point in bringing in data that will not be copied to the destination.

When the source data enters the data flow, it is automatically converted to an SSIS data type. Data types that cannot be converted will raise an error. In the source editor, you can change the data type of each column, and set type specific properties such as precision, scale, and width. Some data type mapping may be necessary between the source and destinations. For example, SSIS raises an error when trying to insert a Unicode column into a non-Unicode column destination. Furthermore, you might need to adjust column widths of string types to avoid truncation. Keeping widths as narrow as possible helps performance when the data is moved from the source to destination. You can also use the Data Conversion and Derived Column transformations, along with expressions, and the cast operator to manage types.

For a full list of Integration Services data types, see the following:

Integration Services Data Types

<https://aka.ms/k7bnin>

Data Transformation Components

You can use data transformation components to perform operations on rows of data as they pass through the data flow pipeline. Transformation components have both inputs and outputs. These transformations use memory in one of three different ways:

- **Non-blocking:** each row passes through the transformation and on to the next task without waiting.
- **Partial-blocking:** after sufficient rows have been read, they will then pass on to the next task.

- Perform operations on rows of data
- Use inputs and outputs
- Multiple transformation types:
 - Row Transformations
 - Rowset Transformations
 - Split and Join Transformations
 - Auditing Transformations
 - BI Transformations
 - Custom Transformations
- Blocking types:
 - Non-blocking
 - Partial-blocking
 - Blocking

- **Blocking:** the transformation reads all rows before passing them on.

The way in which a transformation uses memory can considerably affect the performance of your package.

The Balanced Data Distributor transformation can help alleviate the loss of performance in blocking transformations by splitting the flow of data so it can be processed in parallel.

The following tables list the transformation components in SSIS. The right-hand column denotes the type of blocking used by the transformation: non-blocking (N), partial-blocking (P), and blocking (B):

Row transformations – update column values or create new columns for each row in the data flow		
Character Map	Applies string functions to column values, such as conversion from lowercase to uppercase.	N
Copy Column	Creates a copy of a column and adds it to the data flow.	N
Data Conversion	Converts data of one type to another, for example, numerical values to strings.	N
Derived Column	Adds a new column based on an expression. For example, you could use an expression to multiply a Quantity column by a UnitPrice column to create a new TotalPrice column.	N
Export Column	Saves the contents of a column as a file.	N
Import Column	Reads data from a file and adds it as a column in the data flow.	N
OLE DB Command	Runs a SQL command for each row in the data flow.	N

Rowset transformations – create new rowsets		
Aggregate	Creates a new rowset by applying aggregate functions such as SUM.	B
Sort	Creates a new sorted rowset.	B
Percentage Sampling	Creates a rowset by randomly selecting a specified percentage of rows.	N
Row Sampling	Creates a rowset by randomly selecting a specified number of rows.	B
Pivot	Creates a rowset by condensing multiple records with a single column into a single record with multiple columns.	P
Unpivot	Creates a rowset by expanding a single record with multiple columns into multiple records with a single column.	P

Split and Join transformations – merge or branch data flows		
Conditional Split	Splits a single-input rowset into multiple-output rowsets based on conditional logic.	N
Multicast	Distributes all input rows to multiple outputs.	N
Union All	Adds multiple inputs into a single output.	P
Merge	Merges two sorted inputs into a single output.	P
Merge Join	Joins two sorted inputs to create a single output based on a FULL, LEFT, or INNER join operation.	P
Lookup	Looks up columns in a data source by matching key values in the input. It creates an output for matched rows and a second output for rows with no matching value in the lookup data source.	N
Cache Transform	Caches data from a data source to be used by a lookup transformation.	N
CDC Splitter	Splits inserts, updates, and deletes from a CDC source into separate data flows. CDC is discussed in Module 9: <i>Implementing an Incremental ETL Process</i> .	N
Balanced Data Distributor	Improves the performance of the package by separating the output into different threads.	N/A

Auditing transformations – add audit information or count rows		
Audit	Provides execution environment information that can be added to the data flow.	N
RowCount	Counts the rows in the data flow and writes the result to a variable.	N

BI transformations – perform BI tasks		
Slowly Changing Dimension	Redirects rows when loading a data warehouse to preserve historical dimension values.	N
Fuzzy Grouping	Uses fuzzy logic to deduplicate rows in the data flow.	B
Fuzzy Lookup	Looks up columns in a data source by finding approximate matches for values in the input.	B
Term Extraction	Extracts nouns or noun phrases from text for statistical analysis.	P
Term Lookup	Matches terms extracted from text with terms in a reference table.	P

BI transformations – perform BI tasks		
Data Mining Query	Runs a data mining prediction query against the input to predict unknown column values.	P
DQS Cleansing	Applies a Data Quality Services knowledge base to data as it flows through the pipeline.	P

Custom transformations – perform custom operations		
Script Component	Runs custom script code for each row in the input.	P

Row transformations generally perform well because they work at the row level, without needing other rows from the data source. The Merge and Merge Join transformations need to sort the data from the two inputs to define the sort order, preserving the sort output, or for matching across the sorted rows. They use partial blocking, so rows are delayed in moving to the output path. Multirow transformations, such as the Aggregate, Row-Sampling, and Sort transformations tend to be memory intensive: they need to read all input rows before moving them to the output path. The more data you are passing through your data flow transformations, the higher the potential performance cost of using blocking operations. If you have large datasets, it can be helpful to look for ways of avoiding blocking transformations. For example, if you can avoid using the Merge Join transformation by performing the operation using Transact-SQL instead, this might help the performance of the package. If not, you can sort the data in the source query before it moves to the Merge Join transformation, and in the Advanced Editor, set the IsSorted property to True. In some cases you might find that a non-blocking RowCount transformation delivers the same result as using a blocking Aggregate transformation, or that the pivoting or unpivoting of data can be done in your staging tables instead.

To add a transformation component to a workflow, drag it from the **Common** or **Other Transforms** section of the SSIS Toolbox pane to the design surface, and then connect the required inputs to the transformation component. Double-click the transformation component to configure the specific operation that it will perform, and then define the columns to include in the outputs from the transformation component.

For more information on SSIS Transformations, see the SQL Server Technical Documentation:

 **SSIS Transformations**

<http://aka.ms/i&t9cs>

Optimizing Data Flow Performance

There are several techniques that you can apply to optimize the performance of a data flow. When you are implementing a data flow, consider the following guidelines:

- **Optimize queries.** To reduce the overall volume of data in the data flow, select only the rows and columns you need.
- **Avoid unnecessary sorting.** If you require sorted data from a single data source, sort it during the extraction by using a query with an **ORDER BY** clause. If subsequent transformation components in your data flow rely on sorted data, use the **IsSorted** property of the output to indicate that the data is already sorted. To set the **IsSorted** property for the output from a source component, right-click the source component and click **Show Advanced Editor**. In the **Advanced Editor** dialog box, click the **Input and Output Properties** tab. In the **Inputs and outputs** folder list, click the **<source type> Output** folder (where **<source type>** is the type of source you are setting the property for). This will display a list of output properties, including the **IsSorted** property. To indicate that the data is already sorted, set the **IsSorted** property to **True**.
- **Optimize performance.** Configure Data Flow task properties using the following properties:
 - **DefaultBufferSize** and **DefaultBufferMaxRows.** Configuring the size of the buffers the data flow uses can significantly improve performance. When there is sufficient memory available, you should try to achieve a small number of large buffers without incurring any disk paging. The default values for these properties are 10 MB and 10,000 rows respectively.
 - **BufferTempStoragePath** and **BLOBTempStoragePath.** Using these properties to locate temporary objects that the data flow creates on a fast disk, or spreading them across multiple storage devices, can improve performance.
 - **EngineThreads.** Setting the number of threads available to the Data Flow task can improve execution performance, particularly in packages where the **MaxConcurrentExecutables** property is set to enable parallel execution of the package's tasks across multiple processors.
 - **RunInOptimizedMode.** Setting a Data Flow task to run in optimized mode increases performance by removing any columns or components that are not required further downstream in the data flow.

- Optimize queries:
 - Select only the rows and columns that you need
- Avoid unnecessary sorting:
 - Use pre-sorted data where possible
 - Set the **IsSorted** property where applicable
- Configure Data Flow task properties:
 - Buffer size
 - Temporary storage location
 - Parallelism
 - Optimized mode

For more information on data flow performance, see the SQL Server Technical Documentation:

 **Data Flow Performance Features**

<http://aka.ms/kqi50q>

Demonstration: Implementing a Data Flow

In this demonstration, you will see how to:

- Configure a data source component.
- Use a derived column transformation.
- Use a lookup transformation.
- Configure a destination.

Demonstration Steps

Configure a Data Source Component

1. Ensure you have completed the previous demonstrations in this module.
2. Start SQL Server Management Studio and connect to the **MIA-SQL** database engine instance using Windows authentication.
3. In Object Explorer, expand **Databases**, expand **Products**, and then expand **Tables**. Right-click each of the following tables, click **Select Top 1000 Rows**, and then review the data they contain.
 - **dbo.Product**
 - **dbo.ProductCategory**
 - **dbo.ProductSubcategory**
4. In Object Explorer, under **Databases**, expand **DemoDW**, and then expand **Tables**. Right-click **dbo.DimProduct**, and then click **Select Top 1000 Rows** to verify that this table is empty.
5. Start SQL Server Data Tools and create a new Integration Services project named **DataFlowDemo** in the **D:\Demofiles\Mod06** folder.
6. In Solution Explorer, expand **SSIS Packages**, right-click **Package.dtsx**, and then click **Rename**. Change the package name to **ExtractProducts.dtsx**.
7. In Solution Explorer, right-click **Connection Managers**, and then click **New Connection Manager**. Add a new **OLEDB** connection manager with the following settings:
 - **Server name:** localhost
 - **Log on to the server:** Windows Authentication
 - **Select or enter a database name:** Products
8. In the **SSIS Toolbox**, in the **Favorites** section, double-click **Data Flow Task** to add it to the Control Flow surface.
9. Right-click **Data Flow Task**, and then click **Rename**. Change its name to **Extract Products**.
10. Double-click **Extract Products** to switch to the **Data Flow** tab.
11. In the **SSIS Toolbox**, in the **Favorites** section, double-click **Source Assistant** to add a source component to the Data Flow surface.
12. In the **Source Assistant - Add New Source** dialog box, in the list of types, click **SQL Server**. In the list of connection managers, click **localhost.Products**, and then click **OK**.
13. Rename the **OLE DB Source** to **Products**, and then double-click it to edit its settings.
14. In the **OLE DB Source Editor** dialog box, on the **Connection Manager** page, view the list of available tables and views in the drop-down list.

15. Change the data access mode to **SQL Command**, and then enter the following Transact-SQL code:

```
SELECT ProductKey, ProductName FROM Product
```

16. Click **Build Query** to open the **Query Builder** dialog box.
17. In the **Product** table, select the **ProductSubcategoryKey**, **StandardCost**, and **ListPrice** columns, and then click **OK**.
18. In the **OLE DB Source Editor** dialog box, click **Preview** to see a data preview, and then click **Close** to close the **Preview Query Results** dialog box.
19. In the **OLE DB Source Editor** dialog box, on the **Columns** page, view the list of external columns that the query has returned and the output columns generated by the data source, and then click **OK**.

Use a Derived Column Transformation

1. In the SSIS Toolbox pane, in the **Common** section, double-click **Derived Column** to add a derived column transformation component to the Data Flow surface, and then position it below the **Products** source component.
2. Rename the **Derived Column** transformation component to **Calculate Profit**.
3. Select the **Products** source component, and then drag the blue output arrow to the **Calculate Profit** transformation component.
4. Double-click the **Calculate Profit** transformation component to edit its settings, and then in the **Derived Column Name** box, type **Profit**.
5. Ensure that **<add as new column>** is selected in the **Derived Column** box.
6. Expand the **Columns** folder, and then drag the **ListPrice** column to the Expression box.
7. In the Expression box, after [ListPrice], type a minus sign (-), and then drag the **StandardCost** column to the Expression box to create the following expression:

```
[ListPrice]-[StandardCost]
```

8. Click the **Data Type** box, ensure that it is set to **Currency [DT_CY]**, and then click **OK**.

Use a Lookup Transformation

1. In the SSIS Toolbox pane, in the **Common** section, double-click **Lookup** to add a lookup transformation component to the Data Flow surface, and then position it below the **Calculate Profit** transformation component.
2. Rename the **Lookup** transformation component to **Lookup Category**.
3. Select the **Calculate Profit** transformation component, and then drag the blue output arrow to the **Lookup Category** transformation component.
4. Double-click the **Lookup Category** transformation component to edit its settings.
5. In the **Lookup Transformation Editor** dialog box, on the **General** page, in the **Specify how to handle rows with no matching entries** list, select **Redirect rows to no match output**.
6. In the **Lookup Transformation Editor** dialog box, on the **Connection** page, ensure that the **localhost.Products** connection manager is selected, and then click **Use results of an SQL query**.
7. Click **Browse**, move to the **D:\Demofiles\Mod06** folder, and then open the **LookupProductCategories.sql** query.

8. Click **Preview** to view the product category data, note that it includes a **ProductSubcategoryKey** column, and then click **Close** to close the preview.
9. In the **Lookup Transformation Editor** dialog box, on the **Columns** page, in the **Available Input Columns** list, drag **ProductSubcategoryKey** to **ProductSubCategoryKey** in the **Available Lookup Columns** list.
10. Select the **ProductSubcategoryName** and **ProductCategoryName** columns to add them as new columns to the data flow, and then click **OK**.

Configure a Destination Component

1. In Solution Explorer, create a new OLE DB connection manager with the following settings:
 - o **Server name:** localhost
 - o **Log on to the server:** Windows Authentication
 - o **Select or enter a database name:** DemoDW
2. In the SSIS Toolbox pane, in the **Favorites** section, double-click **Destination Assistant** to add a destination component to the Data Flow surface.
3. In the **Destination Assistant - Add New Destination** dialog box, in the list of types, click **SQL Server**. In the list of connection managers, click **localhost.DemoDW**, and then click **OK**.
4. Rename the OLE DB destination component to **DemoDW** and position it below the **Lookup Category** transformation component.
5. Select the **Lookup Category** transformation component, and then drag the blue output arrow to the **DemoDW** destination component.
6. In the **Input Output Selection** dialog box, in the **Output** list, click **Lookup Match Output**, and then click **OK**.
7. Double-click the **DemoDW** destination component to edit its settings, and then in the **Name of the table or the view list**, click **[dbo].[DimProduct]**.
8. In the **OLE DB Destination Editor** dialog box, on the **Mappings** page, note that input columns are automatically mapped to destination columns with the same name.
9. In the **Available Input Columns** list, drag the **ProductKey** column to the **ProductID** column in the **Available Destination Columns** list, and then click **OK**.
10. In the SSIS Toolbox pane, in the **Other Destinations** section, double-click **Flat File Destination** to add a destination component to the Data Flow surface, and then position it to the right of the **Lookup Category** transformation component.
11. Rename the flat file destination component **Uncategorized Products**.
12. Select the **Lookup Category** transformation component, and then drag the blue output arrow to the **Uncategorized Products** destination component. The **Lookup No Match Output** is automatically selected.
13. Double-click the **Uncategorized Products** destination component to edit its settings, and then click **New**. In the **Flat File Format** dialog box, select **Delimited**, and then click **OK**.
14. In the **Flat File Connection Manager Editor** dialog box, name the new connection manager **Unmatched Products**, specify the file name **D:\Demofiles\Mod06\UnmatchedProducts.csv**, and then click **OK**.

15. In the **Flat File Destination Editor** dialog box, click the **Mappings** page and note that the input columns are mapped to destination columns with the same names, and then click **OK**.
16. On the **Debug** menu, click **Start Debugging**, and observe the data flow as it runs. Note that the number of rows transferred along each path is shown in front of the arrows connecting the components in the data flow.
17. When the data flow has completed, on the **Debug** menu, click **Stop Debugging**.
18. Close SSDT, saving your changes if you are prompted.
19. Use Notepad to examine the **Unmatched Products.csv** flat file in the **D:\Demofiles\Mod06** folder. Note that there were no unmatched products.
20. Use SQL Server Management Studio to view the contents of the **DimProduct** table in the **DemoDW** database, and note that the product data has been transferred.
21. Close SQL Server Management Studio without saving any files.

Question: In a lookup transformation component, the source of the lookup columns is a table, view, or SQL statement that you must specify. Why do you not need to specify the source of the input columns?

Lab: Implementing Data Flow in an SSIS Package

Scenario

In this lab, you will focus on the extraction of customer and sales order data from the **InternetSales** database used by the company's e-commerce site, which you must load into the **Staging** database. This database contains customer data (in a table named **Customers**), and sales order data (in tables named **SalesOrderHeader** and **SalesOrderDetail**). You will extract sales order data at the line item level of granularity. The total sales amount for each sales order line item is then calculated by multiplying the unit price of the product purchased by the quantity ordered. Additionally, the sales order data includes only the ID of the product purchased, so your data flow must look up the details of each product in a separate **Products** database.

Objectives

After completing this lab, you will be able to:

- Extract and profile source data.
- Implement a data flow.
- Use transformation components in a data flow.

Estimated Time: 60 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Exploring Source Data

Scenario

You have designed a data warehouse schema for Adventure Works Cycles. You must now design an ETL process to populate it with data from various source systems. Before creating the ETL solution, you have decided to examine the source data so you can understand it better.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Extract and View Sample Source Data
3. Profile Source Data

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **20767C-MIA-SQL** virtual machine, run **Setup.cmd** in the **D:\Labfiles\Lab06\Starter** folder as Administrator.

► Task 2: Extract and View Sample Source Data

1. Use the SQL Server Import and Export Data Wizard to extract a sample of customer data from the **InternetSales** database on the **localhost** instance of SQL Server to a comma-delimited flat file.
 - Your sample should consist of the first 1,000 records in the **Customers** table.
 - You should use a text qualifier because some string values in the table may contain commas.

2. After you have extracted the sample data, use Notepad to view it.

Note: You may observe some anomalies in the data, such as invalid gender codes and multiple values for the same country or region. The purpose of examining the source data is to identify as many of these problems as possible, so that you can resolve them in the development of the ETL solution. You will address the problems in this data in later labs.

► Task 3: Profile Source Data

1. Create an Integration Services project named **Explore Internet Sales** in the **D:\Labfiles\Lab06\Starter** folder.
2. Add an ADO.NET connection manager that uses Windows authentication to connect to the **InternetSales** database on the **localhost** instance of SQL Server.
3. Use a Data Profiling task to generate the following profile requests for data in the **InternetSales** database:
 - Column statistics for the **OrderDate** column in the **SalesOrderHeader** table. You will use this data to find the earliest and latest dates on which orders have been placed.
 - Column length distribution for the **AddressLine1** column in the **Customers** table. You will use this data to determine the appropriate column length to allow for address data.
 - Column null ratio for the **AddressLine2** column in the **Customers** table. You will use this data to determine how often the second line of an address is null.
 - Value inclusion for matches between the **PaymentType** column in the **SalesOrderHeader** table and the **PaymentTypeKey** column in the **PaymentTypes** table. Do not apply an inclusion threshold and set a maximum limit of 100 violations. You will use this data to find out if any orders have payment types that are not present in the table of known payment types.
4. Run the SSIS package and view the report that the Data Profiling task generates in the Data Profile Viewer.

Results: After this exercise, you should have a comma-separated text file that contains a sample of customer data, and a data profile report that shows statistics for data in the **InternetSales** database.

Exercise 2: Transferring Data by Using a Data Flow Task

Scenario

Now you have explored the source data in the **InternetSales** database, you are ready to start implementing data flows for the ETL process. A colleague has already implemented data flows for reseller sales data, and you plan to model your Internet sales data flows on those.

The main tasks for this exercise are as follows:

1. Examine an Existing Data Flow
2. Create a Data Flow Task
3. Add a Data Source Component to a Data Flow
4. Add a Data Destination Component to a Data Flow
5. Test the Data Flow Task

► Task 1: Examine an Existing Data Flow

1. Open the **D:\Labfiles\Lab06\Starter\Ex2\AdventureWorksETL.sln** solution in SSDT.
2. Open the **Extract Reseller Data.dtsx** package and examine its control flow. Note that it contains two Data Flow tasks.
3. On the **Data Flow** tab, view the **Extract Resellers** task and note that it contains a source component named **Resellers** and a destination component named **Staging DB**.
4. Examine the **Resellers** source component, noting the connection manager that it uses, the source of the data, and the columns that its output contains.
5. Examine the **Staging DB** destination component, noting the connection manager that it uses, the destination table for the data, and the mapping of input columns to destination columns.
6. Right-click anywhere on the Data Flow design surface, click **Execute Task**, and then observe the data flow as it runs, noting the number of rows transferred.
7. When the data flow has completed, stop the debugging session.

► Task 2: Create a Data Flow Task

1. Add a new package to the project and name it **Extract Internet Sales Data.dtsx**.
2. Add a Data Flow task named **Extract Customers** to the new package's control flow.

► Task 3: Add a Data Source Component to a Data Flow

1. Create a new project-level OLE DB connection manager that uses Windows authentication to connect to the **InternetSales** database on the localhost instance of SQL Server.
2. In the **Extract Customers** data flow, add a source component that uses the connection manager you created for the **InternetSales** database, and name it **Customers**.
3. Configure the **Customers** source component to extract all columns from the **Customers** table in the **InternetSales** database.

► Task 4: Add a Data Destination Component to a Data Flow

1. Add a destination component that uses the existing **localhost.Staging** connection manager to the **Extract Customers** data flow, and then name it **Staging DB**.
2. Connect the output from the **Customers** source component to the input of the **Staging DB** destination component.
3. Configure the **Staging DB** destination component to load data into the **Customers** table in the **Staging** database.
4. Ensure that all columns are mapped, and in particular that the **CustomerKey** input column is mapped to the **CustomerBusinessKey** destination column.

► Task 5: Test the Data Flow Task

1. Right-click anywhere on the Data Flow design surface, click **Execute Task**, and then observe the data flow as it runs, noting the number of rows transferred.
2. When the data flow has completed, stop the debugging session.

Results: After this exercise, you should have an SSIS package that contains a single Data Flow task, which extracts customer records from the **InternetSales** database and inserts them into the **Staging** database.

Exercise 3: Using Transformation Components in a Data Flow

Scenario

You have implemented a simple data flow to transfer customer data to the staging database. Now you must implement a data flow for Internet sales records. The new data flow adds a new column containing the total sales amount for each line item (which is derived by multiplying the list price by the quantity of units purchased), and uses a product key value to find additional data in a separate Products database. Once again, you will model your solution on a data flow that a colleague has already implemented for reseller sales data.

The main tasks for this exercise are as follows:

1. Examine an Existing Data Flow
2. Create a Data Flow Task
3. Add a Data Source Component to a Data Flow
4. Add a Derived Column Transformation Component to a Data Flow
5. Add a Lookup Transformation Component to a Data Flow
6. Add a Data Destination Component to a Data Flow
7. Test the Data Flow Task

► Task 1: Examine an Existing Data Flow

1. Open the **D:\Labfiles\Lab06\Starter\Ex3\AdventureWorksETL.sln** solution in SSDT.
2. Open the **Extract Reseller Data.dtsx** package and examine its control flow. Note that it contains two Data Flow tasks.
3. On the **Data Flow** tab, view the **Extract Reseller Sales** task.
4. Examine the **Reseller Sales** source component, noting the connection manager that it uses, the source of the data, and the columns that its output contains.
5. Examine the **Calculate Sales Amount** transformation component, noting the expression that it uses to create a new derived column.
6. Examine the **Lookup Product Details** transformation component, noting the connection manager and query that it uses to look up product data, and the column mappings used to match data and add columns to the data flow.
7. Examine the arrows from the **Lookup Product Details** transformation component to the **Staging DB** destination component and the **Orphaned Sales** destination component.
8. Examine the **Staging DB** destination component, noting the connection manager that it uses, the destination table for the data, and the mapping of input columns to destination columns.
9. Examine the **Orphaned Sales** destination component, noting the connection manager.
10. Right-click anywhere on the Data Flow design surface, click **Execute Task**, and then observe the data flow as it runs, noting the number of rows transferred.
11. When the data flow has completed, stop the debugging session.

► Task 2: Create a Data Flow Task

1. Open the **Extract Internet Sales Data.dtsx** package, and then add a new Data Flow task named **Extract Internet Sales** to its control flow.
2. Connect the pre-existing **Extract Customers** Data Flow task to the new **Extract Internet Sales** task.

► Task 3: Add a Data Source Component to a Data Flow

1. Add a source component that uses the existing **localhost.InternetSales** connection manager to the **Extract Internet Sales** data flow, and then name it **Internet Sales**.
2. Configure the **Internet Sales** source component to use the Transact-SQL code in the **D:\Labfiles\Lab06\Starter\Ex3\InternetSales.sql** query file to extract Internet sales records.

► Task 4: Add a Derived Column Transformation Component to a Data Flow

1. Add a derived column transformation component named **Calculate Sales Amount** to the **Extract Internet Sales** data flow.
2. Connect the output from the **InternetSales** source component to the input of the **Calculate Sales Amount** transformation component.
3. Configure the **Calculate Sales Amount** transformation component to create a new column named **SalesAmount**, containing the **UnitPrice** column value multiplied by the **OrderQuantity** column value.

► Task 5: Add a Lookup Transformation Component to a Data Flow

1. Add a lookup transformation component named **Lookup Product Details** to the **Extract Internet Sales** data flow.
2. Connect the output from the **Calculate Sales Amount** transformation component to the input of the **Lookup Product Details** transformation component.
3. Configure the **Lookup Product Details** transformation component to:
 - Redirect unmatched rows to the no match output.
 - Use the **localhost.Products** connection manager and the **Products.sql** query in the **D:\Labfiles\Lab06\Starter\Ex3** folder to retrieve product data.
 - Match the **ProductKey** input column to the **ProductKey** lookup column.
 - Add all lookup columns other than **ProductKey** to the data flow.
4. Add a flat file destination component named **Orphaned Sales** to the **Extract Internet Sales** data flow. Redirect nonmatching rows from the **Lookup Product Details** transformation component to the **Orphaned Sales** destination component, which should save any orphaned records in a comma-delimited file named **Orphaned Internet Sales.csv** in the **D:\Labfiles\Lab06\Starter\Ex3** folder.

► Task 6: Add a Data Destination Component to a Data Flow

1. Add a destination component that uses the **localhost.Staging** connection manager to the **Extract Customers** data flow, and name it **Staging DB**.
2. Connect the match output from the **Lookup Product Details** transformation component to the input of the **Staging DB** destination component.
3. Configure the **Staging DB** destination component to load data into the **InternetSales** table in the **Staging** database. Ensure that all columns are mapped. In particular, ensure that all Key input columns are mapped to the equivalent BusinessKey destination columns.

► **Task 7: Test the Data Flow Task**

1. Right-click anywhere on the Data Flow design surface, click **Execute Task**, and then observe the data flow as it runs, noting the number of rows.
2. When the data flow has completed, stop the debugging session.

Results: After this exercise, you should have a package that contains a Data Flow task including derived column and lookup transformation components.

Module Review and Takeaways

In this module, you have learned how to explore source data and how to use SQL Server Integration Services to implement a data flow.

Review Question(s)

Question: How could you determine the range of OrderDate values in a data source to plan a time dimension table in a data warehouse?

Module 7

Implementing Control Flow in an SSIS Package

Contents:

Module Overview	7-1
Lesson 1: Introduction to Control Flow	7-2
Lesson 2: Creating Dynamic Packages	7-11
Lesson 3: Using Containers	7-16
Lab A: Implementing Control Flow in an SSIS Package	7-22
Lesson 4: Managing Consistency	7-27
Lab B: Using Transactions and Checkpoints	7-33
Module Review and Takeaways	7-37

Module Overview

You can use control flow in SQL Server Integration Services (SSIS) packages to implement complex extract, transform, and load (ETL) solutions that combine multiple tasks and workflow logic. By learning how to implement control flow, you can design robust ETL processes for a data warehousing solution that coordinate data flow operations with other automated tasks.

Objectives

After completing this module, you will be able to:

- Implement control flow with tasks and precedence constraints.
- Create dynamic packages that include variables and parameters.
- Use containers in a package control flow.
- Enforce consistency with transactions and checkpoints.

Lesson 1

Introduction to Control Flow

Control flow in an SSIS package consists of one or more tasks, usually executed as a sequence, based on precedence constraints that define a workflow. Before you can implement a control flow, you need to know what tasks are available and how to define a workflow sequence using precedence constraints. You should also understand how to use multiple packages to create complex ETL solutions.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the control flow tasks provided by SSIS.
- Define a workflow for tasks by using precedence constraints.
- Use design time features of SSIS to help you develop control flow efficiently.
- Use multiple packages in an SSIS solution.
- Create reusable package templates.

Control Flow Tasks

A control flow consists of one or more tasks. SSIS includes the following control flow tasks that you can use in a package:

- Data Flow tasks
- Database tasks
- File and Internet tasks
- Process Execution tasks
- WMI tasks
- Custom Logic tasks
- Database Transfer tasks
- Analysis Services tasks
- SQL Server Maintenance tasks

Data Flow Tasks	
Data Flow	Encapsulates a data flow that transfers data from a source to a destination.
Database Tasks	
Data Profiling	Generates statistical reports based on a data source.
Bulk Insert	Inserts data into a data destination in a bulk load operation.
Execute SQL	Runs a structured query language (SQL) query in a database.
Execute T-SQL Statement	Runs a Transact-SQL query in a Microsoft® SQL Server® database.
CDC Control	Performs a change data capture (CDC) status management operation.

Data Flow Tasks	
File and Internet Tasks	
File System	Performs file system operations, such as creating folders or deleting files.
FTP	Performs file transfer protocol (FTP) operations, such as copying files.
XML	Performs XML processing operations, such as applying a style sheet.
Web Service	Calls a method on a specific web service.
Send Mail	Sends an email message.
Message Queue	Sends and receives MSMQ messages between an SSIS package and an Application Queue.
Process Execution Tasks	
Execute Package	Runs a specified SSIS package.
Execute Process	Runs a specified app.
WMI Tasks	
WMI Data Reader	Runs a Windows Management Instrumentation (WMI) query.
WMI Event Watcher	Monitors a specific WMI event.
Custom Logic Tasks	
Script	Runs Microsoft Visual Studio® Tools for Applications (VSTA) script written in Microsoft Visual Basic® or Microsoft Visual C#®.
Custom Task	A custom task implemented as a .NET assembly.
Database Transfer Tasks	
Transfer Database	Transfers a database from one SQL Server instance to another.
Transfer Error Messages	Transfers custom error messages from one SQL Server instance to another.
Transfer Jobs	Transfers SQL Agent jobs from one SQL Server instance to another.
Transfer Logins	Transfers logins from one SQL Server instance to another.
Transfer Master Stored Procedures	Transfers stored procedures in the master database from one SQL Server instance to another.

Data Flow Tasks	
Transfer SQL Server Objects	Transfers database objects, such as tables and views, from one SQL Server instance to another.
Analysis Services Tasks	
Analysis Services Execute DDL	Runs a data definition language (DDL) statement in an Analysis Services instance—for example, to create a cube.
Analysis Services Processing	Processes an Analysis Services object, such as a cube or data mining model.
Data Mining Query	Runs a prediction query using a data mining model.
Big Data Tasks	
Hadoop File System	Runs a Hadoop File System task.
Hadoop Hive	Runs a Hive Script on a Hadoop Cluster.
Hadoop Pig	Runs a Pig Script on a Hadoop Cluster.
SQL Server Maintenance Tasks	
Backup Database	Backs up a SQL Server database.
Check Database Integrity	Checks the integrity of a SQL Server database.
History Cleanup	Deletes out-of-date history data for SQL Server maintenance operations.
Maintenance Cleanup	Deletes files left by maintenance operations.
Notify Operator	Sends a notification by email message, pager message, or network alert to a SQL Agent operator.
Rebuild Index	Rebuilds a specified index on a SQL Server table or view.
Reorganize Index	Reorganizes a specified index on a SQL Server table or view.
Shrink Database	Reduces the size of the specified SQL Server database.
Update Statistics	Updates value distribution statistics for tables and views in a SQL Server database.

 **Note:** The Execute T-SQL Statement task is only used to query a Microsoft SQL Server Database and lacks the flexibility of the Execute SQL task. If you need to query a database produced by another vendor, run parameterized queries, or assign the results of a query to variables, you should use the Execute SQL task instead.

Precedence Constraints

A control flow usually defines a sequence of tasks to be executed. You define the sequence by connecting tasks with precedence constraints. These constraints evaluate the outcome of a task to determine the flow of execution.

Control Flow Conditions

You can define precedence constraints for one of the following three conditions:

- **Success.** The execution flow to follow when a task completes successfully. In the control flow designer, success constraints are shown as green arrows.
- **Failure.** The execution flow to follow when a task fails. In the control flow designer, failure constraints are shown as red arrows.
- **Completion.** The execution flow to follow when a task completes, regardless of whether it succeeds or fails. In the control flow designer, complete constraints are shown as black arrows.

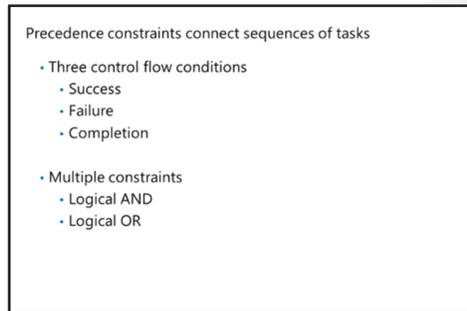
By using these conditional precedence constraints, you can define a control flow that executes tasks based on conditional logic. For example, you could create a control flow with the following steps:

1. An FTP task downloads a file of sales data to a local folder.
2. If the FTP download succeeds, a Data Flow task imports the downloaded data into a SQL Server database. However, if the FTP download fails, a Send Mail task notifies an administrator that there is a problem.
3. When the Data Flow task completes, regardless of whether it fails or succeeds, a File System task deletes the folder from where the customer data file was downloaded.

Using Multiple Constraints

You can connect multiple precedence constraints to a single task. For example, a control flow might include two Data Flow tasks, and a Send Mail task that you want to use to notify an administrator if something goes wrong. To accomplish this, you could connect a failure precedence constraint from each of the Data Flow tasks to the Send Mail task. However, you need to determine whether the notification should be sent if either one of the Data Flow tasks fails, or only if both fail.

By default, when multiple precedence constraints connect to a single task, a logical AND operation is applied to the precedence condition, meaning that all the precedence constraints must evaluate to True to execute the connected task. In the example above, this means that the Send Mail task would only execute if both Data Flow tasks fail. In the control flow designer, logical AND constraints are shown as solid arrows.



You can double-click a precedence constraint to edit and configure it to use a logical OR operation, in which case the connected task executes if any of the connections evaluates to True. Setting the constraints in the example above to use a logical OR operation would result in the Send Mail task executing if either (or both) of the Data Flow tasks fails. In the control flow designer, logical OR constraints are shown as dotted arrows.

Grouping and Annotations

As your control flows become more complex, it can be difficult to interpret the control flow surface. The SSIS Designer includes two features to help SSIS developers work more efficiently.

Grouping Tasks

You can group multiple tasks on the design surface to manage them as a single unit. A task grouping is a “design time only” feature and has no effect on run-time behavior. With a grouped set of tasks, you can:

- Move the tasks around the design surface as a single unit.
- Show or hide the individual tasks to make the best use of space on the screen.

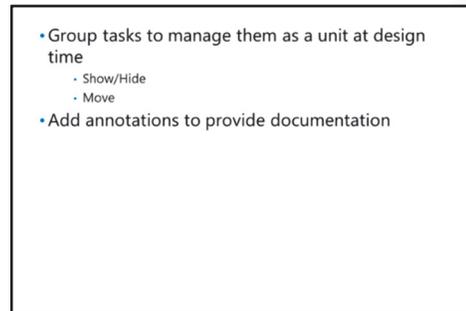
To create a group of tasks, select the ones you want by dragging around or clicking them while holding the CTRL key. Right-click any of the selected tasks, and then click **Group**.

Adding Annotations

You can add annotations to the design surface to document your workflow. An annotation is a text-based note that you use to describe important features of your package design. To add an annotation, right-click the design surface, click **Add Annotation**, and then type the annotation text.



Note: You can add annotations to the Control Flow design surface, the Data Flow design surface, and the Event Handler design surface.



Demonstration: Implementing Control Flow

In this demonstration, you will see how to:

- Add tasks to a control flow.
- Use precedence constraints to define a control flow.

Demonstration Steps

Add Tasks to a Control Flow

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**.
2. In the **D:\Demofiles\Mod07** folder, run **Setup.cmd** as Administrator.

3. Start **SQL Server Data Tools** and open **ControlFlowDemo.sln** from the **D:\Demofiles\Mod07** folder.
4. In Solution Explorer, double-click **Control Flow.dtsx**.
5. If the SSIS Toolbox is not visible, on the **SSIS** menu, click **SSIS Toolbox**.
6. From the SSIS Toolbox, drag a **File System Task** to the control flow surface.
7. Double-click the **File System Task**, configure the following settings, and then click **OK**:
 - o **Name**: Delete Files
 - o **Operation**: Delete directory content
 - o **SourceConnection**: A new connection with a **Usage type** of **Create folder**, and a **Folder** value of **D:\Demofiles\Mod07\Demo**
8. From the SSIS Toolbox, drag a second **File System Task** to the control flow surface.
9. Double-click the **File System Task**, configure the following settings, and then click **OK**.
 - o **Name**: Delete Folder
 - o **Operation**: Delete directory
 - o **SourceConnection**: Demo
10. From the SSIS Toolbox, drag a third **File System Task** to the control flow surface.
11. Double-click the **File System Task**, configure the following settings, and then click **OK**:
 - o **Name**: Create Folder
 - o **Operation**: Create directory
 - o **UseDirectoryIfExists**: True
 - o **SourceConnection**: Demo
12. From the SSIS Toolbox, drag a fourth **File System Task** to the control flow surface.
13. Double-click the **File System Task**, configure the following settings, and then click **OK**:
 - o **Name**: Copy File
 - o **Operation**: Copy file
 - o **DestinationConnection**: Demo
 - o **OverwriteDestination**: True
 - o **SourceConnection**: A new connection with a **Usage type** of **Existing file**, and a **File** value of **D:\Demofiles\Mod07\Demo.txt**
14. From the SSIS Toolbox, drag a **Send Mail Task** to the control flow surface.
15. Double-click the **Send Mail Task**, configure the following settings, and then click **OK**:
 - o **Name** (on the **General** tab): Send Failure Notification
 - o **SmtpConnection** (on the **Mail** tab): Create a new SMTP connection manager with a **Name** property of **Local SMTP Server** and an **SMTP Server** property of **localhost**. Use the default values for all other settings
 - o **From** (on the **Mail** tab): demo@adventureworks.msft
 - o **To** (on the **Mail** tab): student@adventureworks.msft

- **Subject** (on the **Mail** tab): Control Flow Failure
- **MessageSource** (on the **Mail** tab): A task failed

Use Precedence Constraints to Define a Control Flow

1. Select the **Delete Files** task and drag its green arrow to the **Delete Folder** task.
2. Connect the **Delete Folder** task to the **Create Folder** task.
3. Connect the **Create Folder** task to the **Copy File** task.
4. Connect each of the File System tasks to the **Send Failure Notification** task.
5. Right-click the connection between **Delete Files** and **Delete Folder**, and then click **Completion**.
6. Right-click the connection between **Delete Folder** and **Create Folder** and click **Completion**.
7. Click the connection between the **Delete Files** task and the **Send Failure Notification** task to select it. Hold the **Ctrl** key and click each connection between the remaining File System tasks and the **Send Failure Notification** task, to select them all.
8. Press F4 and in the Properties pane, set the **Value** property to **Failure**.
9. Click anywhere on the control flow surface to clear the current selection, and then double-click any of the red constraints connected to the **Send Failure Notification** task.
10. In the **Precedence Constraint Editor** dialog box, in the **Multiple constraints** section, click **Logical OR. One constraint must evaluate to True**, and then click **OK**. Note that all connections to the **Send Failure Notification** task are now dotted to indicate that a logical OR operation is applied.
11. Right-click the control flow surface next to the **Send Failure Notification** task and click **Add Annotation**. Then type **Send an email message if a task fails**.
12. Select the **Delete Files** and **Delete Folder** tasks, then right-click either of them and click **Group**. Drag the group to rearrange the control flow so you can see that the **Delete Folder** task is still connected to the **Create Folder** task.
13. On the **Debug** menu, click **Start Debugging** to run the package, and note that the **Delete Files** and **Delete Folder** tasks failed because the specified folder did not previously exist. This caused the **Send Failure Notification** task to be executed.
14. You can view the email message that was sent by the **Send Failure Notification** task in the **C:\inetpub\mailroot\Drop** folder. Use Notepad to examine the .eml file.
15. In SQL Server Data Tools, on the **Debug** menu, click **Stop Debugging**, and then run the package again. This time, all the File System tasks should succeed because the folder was created during the previous execution. Consequently, the **Send Failure Notification** task is not executed.
16. Stop debugging and close **SQL Server Data Tools**. Save the solution files if prompted.

Using Multiple Packages

Although you can implement an SSIS solution that includes only one package, most enterprise solutions include multiple packages. By dividing your solution into multiple packages, you can:

- Create reusable units of workflow that can be used multiple times in a single ETL process.
- Run multiple control flows in parallel, taking advantage of multiprocessing computers and improving the overall throughput of your ETL processes.
- Separate data extraction workflows to suit data acquisition windows.

You can execute each package independently, in addition to using the Execute Package task to run one package from another.

With multiple packages you can:

- Create reusable units of workflow
- Run multiple control flows in parallel
- Separate ETL workflows to fit data acquisition windows

Create a Package Template

SSIS developers often need to create multiple similar packages. To make the development process more efficient, you can adopt the following procedure to create a package template—this can be reused to create multiple packages with predefined objects and settings:

1. Create a package that includes the elements you want to reuse. These elements can include:
 - Connection Managers
 - Tasks
 - Event Handlers
 - Parameters and Variables

Save the package to the **DataTransformationItems** folder on your development workstation. By default, this folder is located at C:\Program Files (x86)\Microsoft Visual Studio <version>\Common7\IDE\PrivateAssemblies\ProjectItems\DataTransformationProject.

2. When you want to reuse the package, add a new item to the project and select the package in the **Add New Item** dialog box.
3. Change the Name and ID properties of the new package to avoid naming conflicts.

1. Create a package that contains elements you want to reuse
 - Connection Managers
 - Tasks
 - Event Handlers
 - Parameters and Variables
2. Save the package to the **DataTransformationItems** folder
3. Add the package to a project from the **Add New Item** dialog box
4. Change the package name and ID

Check Your Knowledge

Question	
Which of the following is one of the key advantages of developing multiple packages within an SSIS solution?	
Select the correct answer.	
<input type="checkbox"/>	It isolates individual units of workflow in the solution.
<input type="checkbox"/>	It creates reusable units of workflow that you can use multiple times in an ETL process.
<input type="checkbox"/>	Multiple packages result in a less cluttered workspace.
<input type="checkbox"/>	Multiple small packages use fewer resources than a single large package.

Lesson 2

Creating Dynamic Packages

You can use variables, parameters, and expressions to make your SSIS packages more dynamic. For example, rather than hard coding a database connection string or file path in a data source, you can create a package that sets the value dynamically at run time. This produces a more flexible and reusable solution, helping to mitigate differences between the development and production environments.

This lesson describes how you can create variables and parameters, and then use them in expressions.

Lesson Objectives

After completing this lesson, you will be able to:

- Create variables in an SSIS solution.
- Create parameters in an SSIS solution.
- Use expressions in an SSIS solution.

Variables

You can use variables to store values that a control flow uses at run time. Variable values can change when you execute the package to reflect run-time conditions. For example, a variable storing a file path might change, depending on the specific server on which the package is running. You can use variables to:

- Set property values for tasks and other objects.
- Store an iterator or enumerator value for a loop.
- Set input and output parameters for a SQL query.
- Store results from a SQL query.
- Implement conditional logic in an expression.

SSIS packages can contain user and system variables.

User Variables

You can define user variables to store dynamic values that your control flow uses. To create a variable, view the Variables pane in SSIS Designer and click **Add Variable**. For each user variable, you can specify the following properties:

- **Name.** A name for the variable. The combination of name and namespace must be unique within the package. Note that variable names are case-sensitive.
- **Scope.** The scope of the variable. Variables can be accessible throughout the whole package, or scoped to a particular container or task. You cannot set the scope in the Variable pane; it is determined by the object selected when you create the variable.
- **Data Type.** The type of data the variable will hold; for example string, datetime, or decimal.

- **User Variables:**
 - Variables created by an SSIS developer to hold dynamic values
 - Defined in the **User** namespace by default
 - Defined at a specified scope
- **System Variables**
 - Built-in variables with dynamic system values
 - Defined in the **System** namespace

- **Value.** The initial value of the variable.
- **Namespace.** The namespace in which the variable name is unique. By default, user variables are defined in the User namespace, but you can create additional namespaces as required.
- **Raise Change Event.** A true/false value specifying whether to raise an event when the variable value changes. You can then implement an event handler to perform some custom logic.
- **IncludeInDebugDump.** A true/false value specifying whether to include the variable value in debug dump files.



Note: If you inadvertently create a variable with the wrong object selected, resulting in an incorrect scope, you can change the scope of the variable by selecting it in the Variables window and clicking **Move Variable**.

System Variables

System variables store information about the running package and its objects, and are defined in the System namespace. Some useful system variables include:

- **MachineName.** The computer on which the package is running.
- **PackageName.** The name of the package that is running.
- **StartTime.** The time that the package started running.
- **UserName.** The account name of the user who started the package.



Note: For a full list of system variables, see the SQL Server Integration Services documentation in the SQL Server Technical Documentation.

System Variables



<http://aka.ms/scqkpo>

Parameters

You can use parameters to pass values to a project or package at run time. When you define a parameter, you can set a default value, which can be overridden when the package is executed in a production environment. For example, you could use a parameter to specify a database connection string for a data source, using one value during development and a different value when the project is deployed to a production environment.

Parameters have three kinds of value:

- **Design default value.** A default value assigned to the parameter in the design environment.

- **Project parameters**
 - Accessible from any package in the project
- **Package parameters**
 - Exist only at the package level

- **Server default value.** A default value assigned to the parameter during deployment. This value overrides the design default value.
- **Execution value.** A value for a specific execution of a package. This value overrides both the server and design default values.

When you deploy the project to an SSIS Catalog, you can define multiple environments and specify server default parameter values for each environment.

SSIS supports two kinds of parameter:

- **Project parameters.** These are defined at the project level and can be used in any packages within the project.
- **Package parameters.** These are scoped at the package level and are only available within the package for which they are defined.

 **Note:** Parameters are only supported in the project deployment model. When using the legacy deployment model, you can set dynamic package properties by using package configurations. Deployment is discussed in Module 12: *Deploying and Configuring SSIS Packages*.

Expressions

SSIS provides a rich expression language that you can use to set values for numerous elements in an SSIS package, including:

- Properties.
- Conditional split transformation criteria.
- Derived column transformation values.
- Precedence constraint conditions.

- Expressions set values dynamically:
 - Properties
 - Conditional split criteria
 - Derived column values
 - Precedence constraints
- Based on SSIS expression syntax
 - Can include variables and parameters
- Type expressions or create using Expression Builder

Expressions are based on SSIS expression syntax, which uses similar functions and keywords to common programming languages like Microsoft

C#. Expressions can include variables and parameters, meaning you can set values dynamically, based on specific run-time conditions.

For example, you could use an expression in a Data Flow task to specify the location of a file to be used as a data source.

The following sample code shows an expression that concatenates a parameter containing a folder path and a variable containing a file name to produce a full file path:

An SSIS Expression

```
@[$Project::folderPath]+@[User::fName]
```

Note that you must prefix variable names with an @ symbol and, to support identifiers with names containing spaces, you use square brackets to enclose identifier names. Also, note that the expression uses fully qualified parameter and variable names, including the namespace, and that you prefix the parameter name with a \$ symbol.

You can type expressions or, in many cases, create them using the Expression Builder. You can use this graphical tool to drag variables, parameters, constants, and functions to build up the expression that you require. The Expression Builder automatically adds prefixes and text qualifiers for variables and parameters, simplifying the task of creating complex expressions.

Demonstration: Using Variables and Parameters

In this demonstration, you will see how to:

- Create a variable.
- Create a parameter.
- Use variables and parameters in an expression.

Demonstration Steps

Create a Variable

1. Ensure you have completed the previous demonstration in this module.
2. Start **SQL Server Data Tools** and open the **VariablesAndParameters.sln** solution in the **D:\Demofiles\Mod07** folder.
3. In Solution Explorer, double-click **Control Flow.dtsx**.
4. On the **View** menu, point to **Other Windows**, and click **Variables**.
5. In the Variables pane, click **Add Variable** and add a variable with the following properties:
 - **Name:** fName
 - **Scope:** Control Flow
 - **Data type:** String
 - **Value:** Demo1.txt

Create a Parameter

1. In Solution Explorer, double-click **Project.params**.
2. In the Project.params [Design] pane, click **Add Parameter** and add a parameter with the following properties:
 - **Name:** FolderPath
 - **Data type:** String
 - **Value:** D:\Demofiles\Mod07\Files\
 - **Sensitive:** False
 - **Required:** True
 - **Description:** Folder containing text files

Note: Be sure to include the trailing “\” in the Value property.

3. Save all files and close the Project.params [Design] window.

Use a Variable and a Parameter in an Expression

1. On the **Control Flow.dtsx** package design surface, in the Connection Managers pane, click the **Demo.txt** connection manager, and then press F4.
2. In the Properties pane, in the **Expressions** property box, click the ellipsis (...) button.
3. In the **Property Expressions Editor** dialog box, in the **Property** box, select **ConnectionString** and in the **Expression box**, click the ellipsis (...) button.
4. In the **Expression Builder** dialog box, expand the **Variables and Parameters** folder, and then drag the **\$Project::folderPath** parameters to the **Expression** box.
5. In the **Expression** box, type a plus (+) symbol and then drag the **User::fName** variable to the **Expression** box to create the following expression:

```
@[$Project::folderPath]+@[User::fName]
```

6. In the **Expression Builder** dialog box, click **Evaluate Expression** and verify that the expression produces the result **D:\Demofiles\Mod07\Files\Demo1.txt**.
7. Click **OK** to close the **Expression Builder** dialog box, and then in the **Property Expressions Editor** dialog box, click **OK**.
8. Run the project, and when it has completed, stop debugging and close **SQL Server Data Tools**.
9. View the contents of the **D:\Demofiles\Mod07\Demo** folder and verify that **Demo1.txt** has been copied.

Lesson 3

Using Containers

You can create containers in SSIS packages to group related tasks together or define iterative processes. Using containers in packages helps you create complex workflows and a hierarchy of execution scopes that you can use to manage package behavior.

This lesson describes the types of containers that are available and how to use them in an SSIS package control flow.

Lesson Objectives

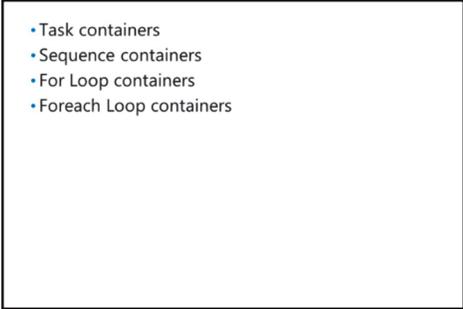
After completing this lesson, you will be able to:

- Describe the types of container available in an SSIS package.
- Use a Sequence container to group related tasks.
- Use a For Loop container to repeat a process until a specific condition is met.
- Use a Foreach Loop container to process items in an enumerated collection.

Introduction to Containers

SSIS packages can have the following kinds of containers:

- **Task containers.** Each control flow task has its own implicit container.
- **Sequence containers.** You can group tasks and other containers into a Sequence container. This creates an execution hierarchy and means you can set properties at the container level that apply to all elements within the container.
- **For Loop containers.** You can use a For Loop container to perform an iterative process until a specified condition is met. For example, you could use a For Loop container to execute the same task a specific number of times.
- **Foreach Loop containers.** You can use a Foreach Loop container to perform an iterative task that processes each element in an enumerated collection. For example, you could use a Foreach Loop container to execute a Data Flow task that imports data from each file in a specified folder into a database.

- 
- Task containers
 - Sequence containers
 - For Loop containers
 - Foreach Loop containers

Containers can be start or endpoints for precedence constraints and you can nest containers within other containers.

Sequence Containers

You can use a Sequence container to group tasks and other containers together, and define a subset of the package control flow. By using a Sequence container, you can:

- Manage properties for multiple tasks as a unit.
- Disable a logical subset of the package for debugging purposes.
- Create a scope for variables.
- Manage transactions at a granular level.

To create a Sequence container, drag the Sequence container icon from the SSIS Toolbox pane to the design surface. You then drag the tasks and other containers you want to include into the Sequence container.

Sequence containers

- Define a control flow subset
- Help you to manage properties for multiple tasks
- Create a scope for variables, transactions, and precedence

Demonstration: Using a Sequence Container

In this demonstration, you will see how to use a Sequence container.

Demonstration Steps

Use a Sequence Container

1. Ensure you have completed the previous demonstrations in this module.
2. Start **SQL Server Data Tools** and open the **SequenceContainer.sln** solution in the **D:\Demofiles\Mod07** folder.
3. In Solution Explorer, double-click **Control Flow.dtsx**.
4. Right-click the **Group** indicator around the **Delete Files** and **Delete Folder** tasks, and then click **Ungroup** to remove it.
5. Drag a **Sequence container** from the SSIS Toolbox to the control flow design surface.
6. Right-click the precedence constraint that connects **Delete Files** to **Send Failure Notification**, and click **Delete**.
7. Delete the precedence constraints connecting the **Delete Folder** to **Send Failure Notification** and **Delete Folder** to **Create Folder**.
8. Click the **Delete Files** and **Delete Folder** tasks to select them both, and then drag into the Sequence container.
9. Drag a precedence constraint from the Sequence container to **Create Folder**, then right-click the precedence constraint and click **Completion**.
10. Drag a precedence constraint from the Sequence container to **Send Failure Notification**. Then right-click the precedence constraint and click **Failure**.
11. Run the package and view the results, then click **stop debugging**.
12. Click the **Sequence container** and press F4, then in the Properties pane, set the **Disable** property to **True**.
13. Run the package again and note that neither of the tasks in the Sequence container is executed.

14. Stop debugging and close **SQL Server Data Tools**.

For Loop Containers

You can use a For Loop container to repeat a portion of the control flow until a specific condition is met. For example, you could run a task a specified number of times.

Conceptually, a For Loop container behaves similarly to a “for construct”, in common programming languages such as Microsoft C#. A For Loop container uses the following expression-based properties to determine the number of iterations it performs:

```
For Loop containers
• Implement iterative control flow
• Similar to a C# For Loop
  - Initialization expression
    @Count = 0
  - Evaluation expression
    @Count < 4
  - Iteration expression
    @Count = @Count + 1
```

- An optional initialization expression, which sets a counter variable to an initial value.
- An evaluation expression that typically evaluates a counter variable, to exit the loop when it matches a specific value.
- An iteration expression that typically modifies the value of a counter variable.

To use a For Loop container in a control flow, drag the For Loop container icon from the SSIS Toolbox to the control flow surface, and then double-click it to set the expression properties required to control the number of loop iterations. Drag the tasks and containers you want to repeat into the For Loop container on the control flow surface.

Demonstration: Using a For Loop Container

In this demonstration, you will see how to use a For Loop container.

Demonstration Steps

Use a For Loop Container

1. Ensure you have completed the previous demonstrations in this module.
2. Start **SQL Server Data Tools** and open the **ForLoopContainer.sln** solution in the **D:\Demofiles\Mod07** folder.
3. In Solution Explorer, double-click **Control Flow.dtsx**.
4. If the Variables window is not open, on the **View** menu, point to **Other Windows**, and then click **Variables**.
5. Add a variable with the following properties:
 - **Name:** counter
 - **Scope:** Control Flow
 - **Data type:** Int32
 - **Value:** 0
6. From the SSIS Toolbox, drag a **For Loop Container** to the control flow design surface.

7. Double-click the **For Loop Container**, set the following properties, and then click **OK**:
 - **InitExpression**: @counter = 1
 - **EvalExpression**: @counter < 4
 - **AssignExpression**: @counter = @counter + 1
8. From the SSIS Toolbox, drag an **Execute Process Task** into the For Loop container.
9. Double-click the **Execute Process Task**, set the following properties, and then click **OK**:
 - **Name** (on the **General** tab): Open File
 - **Executable** (on the **Process** tab): Notepad.exe
 - **Expressions** (on the **Expressions** tab): Use the **Property Expressions Editor** to set the following expression for the **Arguments** property:


```
@[$Project::folderPath] + "Demo" + (DT_WSTR,1)[User::counter] + ".txt"
```
10. Drag a precedence constraint from the **For Loop Container** to the **Sequence Container**.
11. Run the package, and note that the For Loop starts Notepad three times, opening the text file with the counter variable value in its name (Demo1.txt, Demo2.txt, and Demo3.txt). Close Notepad each time it opens, and when the execution is complete, stop debugging.
12. Close **SQL Server Data Tools**, saving the solution files if prompted.

Foreach Loop Containers

You can use a Foreach Loop container to perform an iterative process on each item in an enumerated collection. SSIS supports various Foreach Loop enumerators including:

- **ADO**. You can use this enumerator to loop through elements of an ADO object; for example, records in a Recordset.
- **ADO.NET Schema Rowset**. You can use this enumerator to iterate through objects in an ADO.NET schema; for example, tables in a dataset or rows in a table.
- **File**. You can use this enumerator to iterate through files in a folder.
- **From Variable**. You can use this enumerator to iterate through elements, in a variable that contains an array.
- **Item**. You can use this enumerator to iterate through a property collection for an SSIS object.
- **Nodelist**. You can use this enumerator to iterate through elements and attributes in an XML document.
- **SMO**. You can use this enumerator to iterate through a collection of SQL Server Management Objects.



To use a Foreach Loop container in a control flow:

1. Drag the Foreach Loop container icon from the SSIS Toolbox to the control flow surface.
2. Double-click the Foreach Loop container and select the enumerator you want to use. Each enumerator has specific properties you need to set; for example, the File enumerator requires the path to the folder containing the files through which you want to iterate.
3. Specify the variable in which you want to store the enumerated collection value during each iteration.
4. Drag the tasks you want to perform during each iteration into the Foreach Loop container and configure their properties appropriately to reference the collection value variable.

Demonstration: Using a Foreach Loop Container

In this demonstration, you will see how to use a Foreach Loop container.

Demonstration Steps

Use a Foreach Loop Container

1. Ensure you have completed the previous demonstrations in this module.
2. Start **SQL Server Data Tools** and open the **ForeachLoopContainer.sln** solution in the **D:\Demofiles\Mod07** folder.
3. In Solution Explorer, double-click **Control Flow.dtsx**.
4. From the SSIS Toolbox, drag a **Foreach Loop Container** to the control flow design surface.
5. Double-click the **Foreach loop Container** to view the **Foreach Loop Editor** dialog box.
6. On the **Collection** tab, in the **Enumerator** list, click **Foreach File Enumerator**. In the **Expressions** box, click the ellipsis (...) button.
7. In the **Property Expressions Editor** dialog box, in the **Property** list, click **Directory** and in the **Expression** box click the ellipsis (...) button.
8. In the **Expression Builder** dialog box, expand the **Variables and Parameters** folder and drag the **\$Project::folderPath** parameter to the **Expression** box to specify that the loop should iterate through files in the folder referenced by the **folderPath** project parameter. Click **OK** to close the **Expression Builder**, and then in the **Property Expressions Editor** dialog box, click **OK**.
9. In the **Foreach Loop Editor** dialog box, on the **Collection** tab, in the **Retrieve file name** section, select **Name and extension** to return the file name and extension for each file the loop finds in the folder.
10. In the **Foreach Loop Editor** dialog box, on the **Variable Mappings** tab, in the **Variable** list, click **User::fName** and in the **Index** column, select **0** to assign the file name of each file found in the folder to the **fName** variable. Click **OK**.
11. Remove the precedence constraints that are connected to and from the **Copy File** task, and then drag the **Copy File** task into the **Foreach Loop Container**.
12. Create a precedence constraint from the **Create Folder** task to the **Foreach Loop Container**, and a precedence constraint from the **Foreach Loop Container** to the **Send Failure Notification** task.
13. Right-click the constraint between the **Foreach Loop Container** and the **Send Failure Notification** task, and click **Failure**.

14. Run the package, closing each instance of Notepad as it opens. When the package execution has completed, stop debugging and close SQL Server Data Tools, saving the solution files if prompted.
15. Verify that the **D:\Demofiles\Mod07\Demo** folder contains a copy of each of the files from the **D:\Demofiles\Mod07\Files** folder.

Lab A: Implementing Control Flow in an SSIS Package

Scenario

You are implementing an ETL solution for Adventure Works Cycles and must ensure that the data flows you have already defined are executed as a workflow that notifies operators of success or failure by sending an email message. You must also implement an ETL solution that transfers data from text files generated by the company's financial accounting package to the data warehouse.

Objectives

After completing this lab, you will be able to:

- Use tasks and precedence constraints.
- Use variables and parameters.
- Use containers.

Estimated Time: 60 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Using Tasks and Precedence in a Control Flow

Scenario

You have implemented data flows to extract data and load it into a staging database as part of the ETL process for your data warehousing solution. Now you want to coordinate these data flows by implementing a control flow that notifies an operator of the outcome of the process.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. View the Control Flow
3. Add Tasks to a Control Flow
4. Test the Control Flow

► Task 1: Prepare the Lab Environment

1. Ensure the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab07A\Starter** folder as Administrator.

► Task 2: View the Control Flow

1. Use **SQL Server Data Tools** to open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab07A\Starter\Ex1** folder.
2. Open the **Extract Reseller Data.dtsx** package and examine its control flow. Note that it contains two **Send Mail** tasks—one that runs when either the **Extract Resellers** or **Extract Reseller Sales** tasks fail, and one that runs when the **Extract Reseller Sales** task succeeds.
3. Examine the settings for the precedence constraint connecting the **Extract Resellers** task to the **Send Failure Notification** task to determine the conditions under which this task will be executed.

4. Examine the settings for the **Send Mail** tasks, noting that they both use the **Local SMTP Server** connection manager.
5. Examine the settings of the **Local SMTP Server** connection manager.
6. Run the package, and observe the control flow as the task executes.
7. When package execution is complete, stop debugging, and verify that the email message has been delivered to the **C:\inetpub\mailroot\Drop** folder. You can examine the email file by using Notepad, although the text of the message will be encoded.

► Task 3: Add Tasks to a Control Flow

1. Open the **Extract Internet Sales Data.dtsx** package and examine its control flow.
2. Add a **Send Mail** task to the control flow, configure it with the following settings, and create a precedence constraint that runs this task if the **Extract Internet Sales** task succeeds:
 - **Name:** Send Success Notification
 - **SmtpConnection:** A new SMTP Connection Manager named **Local SMTP Server** that connects to the **localhost** SMTP server
 - **From:** ETL@adventureworks.msft
 - **To:** Student@adventureworks.msft
 - **Subject:** Data Extraction Notification - Success
 - **MessageSourceType:** Direct Input
 - **MessageSource:** The Internet Sales data was successfully extracted
 - **Priority:** Normal
3. Add a second **Send Mail** task to the control flow, configure it with the following settings, and create a precedence constraint that runs this task if either the **Extract Customers** or **Extract Internet Sales** task fails:
 - **Name:** Send Failure Notification
 - **SmtpConnection:** The Local SMTP Server connection manager you created previously
 - **From:** ETL@adventureworks.msft
 - **To:** Student@adventureworks.msft
 - **Subject:** Data Extraction Notification - Failure
 - **MessageSourceType:** Direct Input
 - **MessageSource:** The Internet Sales data extraction process failed
 - **Priority:** High

► Task 4: Test the Control Flow

1. Set the **ForceExecutionResult** property of the **Extract Customers** task to **Failure**. Then run the package and observe the control flow.
2. When package execution is complete, stop debugging, and verify that the failure notification email message has been delivered to the **C:\inetpub\mailroot\Drop** folder. You can use Notepad to examine the email message and verify the subject line.

3. Set the **ForceExecutionResult** property of the **Extract Customers** task to **None**. Then run the package and observe the control flow.
4. When package execution is complete, stop debugging, and verify that the success notification email message has been delivered to the **C:\inetpub\mailroot\Drop** folder.
5. Close **SQL Server Data Tools** when you have completed the exercise.

Results: After this exercise, you should have a control flow that sends an email message if the **Extract Internet Sales** task succeeds, or sends an email message if either the **Extract Customers** or **Extract Internet Sales** tasks fail.

Exercise 2: Using Variables and Parameters

Scenario

You need to enhance your ETL solution to include the staging of payments data that is generated in comma-separated value (CSV) format from a financial accounts system. You have implemented a simple data flow that reads data from a CSV file and loads it into the staging database. You must now modify the package to construct the folder path and file name for the CSV file dynamically at run time instead of relying on a hard-coded name in the Data Flow task settings.

The main tasks for this exercise are as follows:

1. View the Control Flow
2. Create a Variable
3. Create a Parameter
4. Use a Variable and a Parameter in an Expression

► Task 1: View the Control Flow

1. View the contents of the **D:\Accounts** folder and note the files it contains. In this exercise, you will modify an existing package to create a dynamic reference to one of these files.
2. Open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab07A\Starter\Ex2** folder.
3. Open the **Extract Payment Data.dtsx** package and examine its control flow. Note that it contains a single Data Flow task named **Extract Payments**.
4. View the **Extract Payments** data flow and note that it contains a flat file source named **Payments File**, and an OLE DB destination named **localhost.Staging**. This destination connects to the Staging database.
5. View the settings of the **Payments File** source and note that it uses a connection manager named **Payments File**.
6. In the Connection Managers pane, double-click **Payments File**, and note that it references the **Payments.csv** file in the **D:\Labfiles\Lab07A\Starter\Ex2** folder. This file has the same data structure as the **payments** file in the **D:\Accounts** folder.
7. Run the package, and stop debugging when it has completed.
8. On the **Execution Results** tab, find the following line in the package execution log:

[Payments File [2]: The processing of the file "D:\Labfiles\Lab07A\Starter\Ex2\Payments.csv" has started

► Task 2: Create a Variable

- Add a variable with the following properties to the package:
 - **Name:** fName
 - **Scope:** Extract Payment Data
 - **Data type:** String
 - **Value:** Payments - US.csv

Note that the Value property includes a space on either side of the "-" character.

► Task 3: Create a Parameter

- Add a project parameter with the following settings:
 - **Name:** AccountsFolderPath
 - **Data type:** String
 - **Value:** D:\Accounts\
 - **Sensitive:** False
 - **Required:** True
 - **Description:** Path to accounts files



Note: Be sure to include the trailing "\" in the Value property.

► Task 4: Use a Variable and a Parameter in an Expression

1. Set the **Expressions** property of the **Payments File** connection manager in the **Extract Payment Data** package so that the **ConnectionString** property uses the following expression:

```
@[$Project::AccountsFolderPath]+ @[User::fName]
```
2. Run the package and view the execution results to verify that the data in the **D:\Accounts\Payments - US.csv** file was loaded.
3. Close **SQL Server Data Tools** when you have completed the exercise.

Results: After this exercise, you should have a package that loads data from a text file based on a parameter that specifies the folder path where the file is stored, and a variable that specifies the file name.

Exercise 3: Using Containers

Scenario

You have created a control flow that loads Internet sales data and sends a notification email message to indicate whether the process succeeded or failed. You now want to encapsulate the Data Flow tasks for this control flow in a Sequence container so you can manage them as a single unit.

You have also successfully created a package that loads payments data from a single CSV file, based on a dynamically derived folder path and file name. Now you must extend this solution to iterate through all the files in the folder and import data from each one.

The main tasks for this exercise are as follows:

1. Add a Sequence Container to a Control Flow
2. Add a Foreach Loop Container to a Control Flow

► **Task 1: Add a Sequence Container to a Control Flow**

1. Open the **AdventureWorksETL** solution in the **D:\Labfiles\Lab07A\Starter\Ex3** folder.
2. Open the **Extract Internet Sales Data.dtsx** package and modify its control flow so that:
 - The **Extract Customers** and **Extract Internet Sales** tasks are contained in a Sequence container named **Extract Customer Sales Data**.
 - The **Send Failure Notification** task is executed if the **Extract Customer Sales Data** container fails.
 - The **Send Success Notification** task is executed if the **Extract Customer Sales Data** container succeeds.
3. Run the package to verify that it successfully completes both Data Flow tasks in the sequence, and then executes the **Send Success Notification** task.

► **Task 2: Add a Foreach Loop Container to a Control Flow**

1. In the **AdventureWorksETL** solution, open the **Extract Payment Data.dtsx** package.
2. Move the existing **Extract Payment** Data Flow task into a new **Foreach Loop Container**.
3. Configure the **Foreach Loop Container** so that it loops through the files in the folder referenced by the **AccountsFolderPath** parameter, adding each file to the **fName** variable.
4. Run the package and count the number of times the **Foreach Loop** is executed.
5. When execution has completed, stop debugging and view the results to verify that all files in the **D:\Accounts** folder were processed.
6. Close **SQL Server Data Tools** when you have completed the exercise.

Results: After this exercise, you should have one package that encapsulates two Data Flow tasks in a Sequence container, and another that uses a Foreach Loop to iterate through the files in a folder specified in a parameter—and uses a Data Flow task to load their contents into a database.

Lesson 4

Managing Consistency

SSIS solutions are generally used to transfer data from one location to another. Often, the overall SSIS solution can include multiple data flows and operations; it might be important to ensure that the process always results in data that is in a consistent state, even if some parts of the process fail.

This lesson discusses techniques for ensuring data consistency when packages fail.

Lesson Objectives

After completing this lesson, you will be able to:

- Configure failure behavior.
- Use transactions.
- Use checkpoints.

Configuring Failure Behavior

An SSIS package control flow can contain nested hierarchies of containers and tasks. You can use the following properties to control how a failure in one element of the control flow determines the overall package outcome:

- **FailPackageOnFailure.** When set to **True**, the failure of the task or container results in the failure of the package in which it is defined. The default value for this property is **False**.
- **FailParentOnFailure.** When set to **True**, the failure of the task or container results in the failure of its container. If the item with this property is not in a container, then its parent is the package, in which case this property has the same effect as the **FailPackageOnFailure** property. When setting this property on a package executed by an **Execute Package** task in another package, a value of **True** causes the calling package to fail if this package fails. The default value for this property is **False**.
- **MaximumErrorCount.** This property specifies the maximum number of errors that can occur before the item fails. The default value for this property is **1**.

You can use these properties to achieve fine-grained control of package behavior in the event of an error that causes a task to fail.

Use properties to control failure propagation:

- FailPackageOnFailure
- FailParentOnFailure
- MaximumErrorCount

Using Transactions

Transactions ensure that all data changes in a control flow either succeed or fail as a single, atomic unit of work. When tasks are enlisted in a transaction, a failure of any single task causes all tasks to fail, ensuring that the data affected by the control flow remains in a consistent state, with no partial data modifications.

A task, container, or package's participation in a transaction is determined by the

TransactionOption property, which you can set to one of three possible values:

- **Required.** This object requires a transaction and will create a new one if none exists.
- **Supported.** This object will enlist in a transaction if its parent is participating in one.
- **NotSupported.** This object does not support transactions and will not enlist in an existing transaction.

SSIS transactions rely on the Microsoft Distributed Transaction Coordinator (MSDTC), a system component that coordinates transactions across multiple data sources. An error will occur if an SSIS package attempts to start a transaction when the MSDTC service is not running.

SSIS supports multiple concurrent transactions in a single hierarchy of packages, containers, and tasks, but does not support nested transactions. To understand how multiple transactions behave in a hierarchy, consider the following facts:

- If a container with a **TransactionOption** value of **Required** includes a container with a **TransactionOption** of **NotSupported**, the child container will not participate in the parent transaction.
- If the child container includes a task with a **TransactionOption** value of **Supported**, the task will not start a new transaction, but will join any transaction started by its parent.
- If the child container contains a task with a **TransactionOption** value of **Required**, the task will start a new transaction. However, the new transaction is unrelated to the existing transaction, and the outcome of one transaction will have no effect on the other.

For more detailed information on transactions in SSIS, see the SQL Server Technical Documentation:



Integration Services Transactions

<http://aka.ms/tt04jf>

Set the **TransactionOption** property of a task, container, or package:

- Required
- Supported
- NotSupported

Demonstration: Using a Transaction

In this demonstration, you will see how to use a transaction.

Demonstration Steps

Use a Transaction

1. If you did not complete the previous demonstrations in this module, ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod07** folder, run **Setup.cmd** as administrator.
3. Start **SQL Server Management Studio** and connect to the **MIA-SQL** database engine instance using Windows authentication.
4. In Object Explorer, expand **Databases**, expand **DemoDW**, and then expand **Tables**.
5. Right-click **dbo.StagingTable** and click **Select Top 1000 Rows** to verify that it contains product data.
6. Right-click **dbo.ProductionTable** and click **Select Top 1000 Rows** to verify that it is empty.
7. Start **SQL Server Data Tools** and open the **Transactions.sln** solution in the **D:\Demofiles\Mod07** folder.
8. In Solution Explorer, double-click **Move Products.dtsx**. Note that the control flow consists of a Data Flow task named **Copy Products** that moves products from a staging table to a production table, and a **SQL Command** task named **Update Prices** that sets the product price.
9. On the **Debug** menu, click **Start Debugging** to run the package and note that the **Update Prices** task fails. Then on the **Debug** menu, click **Stop Debugging**.
10. In **SQL Server Management Studio**, select the top 1,000 rows from the **dbo.ProductionTable** table, noting that it now contains product data but the prices are all set to 0.00. You want to avoid having products with invalid prices in the production table, so you need to modify the SSIS package to ensure that, when the price update task fails, the production table remains empty.
11. On the **File** menu, point to **New**, and then click **Query with Current Connection**, type the following Transact-SQL code, and then click **Execute**. This deletes all rows in the **dbo.ProductionTable** table:

```
TRUNCATE TABLE DemoDW.dbo.ProductionTable;
```
12. In **SQL Server Data Tools**, click anywhere on the Control Flow surface and press F4.
13. In the Properties pane, set the **TransactionOption** property to **Required**.
14. Click the **Copy Products** task, and in the Properties pane, set the **FailPackageOnFailure** property to **True** and ensure the **TransactionOption** property is set to **Supported**.
15. Repeat the previous step for the **Update Prices** task.
16. Run the package and note that the **Update Prices** task fails again. Stop debugging.
17. In **SQL Server Management Studio**, select the top 1,000 rows from the **dbo.ProductionTable** table, noting that it is empty, even though the **Copy Products** task succeeded. The transaction has rolled back the changes to the production table because the **Update Prices** task failed.
18. In **SQL Server Data Tools**, double-click the **Update Prices** task and change the **SQLStatement** property to **UPDATE ProductionTable SET Price = 100**. Click **OK**.
19. Run the package and note that all tasks succeed.

20. Stop debugging and close SQL Server Data Tools.
21. In **SQL Server Management Studio**, select the top 1,000 rows from the **dbo.ProductionTable** table, noting that it now contains products with valid prices.
22. Close SQL Server Management Studio.

Using Checkpoints

Another way you can manage data consistency is to use checkpoints. You can use checkpoints to restart a failed package after the issue that caused it to fail has been resolved. Any tasks that were previously completed successfully are ignored, and the execution resumes at the point in the control flow where the package failed. While checkpoints do not offer the same level of atomic consistency as a transaction, they can provide a useful solution when a control flow includes a long-running or resource-intensive task that you do not wish to repeat unnecessarily—such as downloading a large file from an FTP server.

- Checkpoints allow failed packages to be restarted without repeating previously successful tasks
- Enable checkpoints by setting package properties:
 - CheckpointFileName
 - CheckpointUsage
 - SaveCheckpoints

Checkpoints work by saving information about work in progress to a checkpoint file. When a failed package is restarted, the checkpoint file is used to identify where to resume execution in the control flow. To enable a package to use checkpoints, you must set the following properties of the package:

- **CheckpointFileName.** The full file path where you want to save the checkpoint file.
- **SaveCheckpoints.** A Boolean value used to specify whether the package should save checkpoint information to the checkpoint file.
- **CheckpointUsage.** An enumeration with one of the following values:
 - **Always:** The package will always look for a checkpoint file when starting. If none exists, the package will fail with an error.
 - **Never:** The package will never use a checkpoint file to resume execution and will always begin execution with the first task in the control flow.
 - **IfExists:** If a checkpoint file exists, the package will use it to resume where it failed previously. If no checkpoint file exists, the package will begin execution with the first task in the control flow.

Demonstration: Using a Checkpoint

In this demonstration, you will see how to use a checkpoint.

Demonstration Steps

Use a Checkpoint

1. Ensure you have completed the previous demonstrations in this module.
2. Start **SQL Server Management Studio** and connect to the **MIA-SQL** database engine instance using Windows authentication.

3. In Object Explorer, expand **Databases**, expand **DemoDW**, and expand **Tables**.
4. Right-click **dbo.StagingTable** and click **Select Top 1000 Rows** to verify that it contains product data.
5. Using Notepad, open **Products.csv** in the **D:\Demofiles\Mod07** folder. Note that it contains details for three more products. Then close Notepad.
6. Start **SQL Server Data Tools** and open the **Checkpoints.sln** solution in the **D:\Demofiles\Mod07** folder.
7. In Solution Explorer, double-click **Load Data.dtsx**. Note that the control flow consists of a File System task to create a folder, a second File System task to copy the products file to the new folder, and a Data Flow task that loads the data in the products file into the staging table.
8. Click anywhere on the Control Flow surface to select the package, and press F4.
9. In the Properties pane, set the following properties:
 - o **CheckpointFileName**: D:\Demofiles\Mod07\Checkpoint.chk
 - o **CheckpointUsage**: IfExists
 - o **SaveCheckpoints**: True
10. Set the **FailPackageOnFailure** property for all three tasks in the control flow to **True**.
11. On the **Debug** menu, click **Start Debugging** to run the package and note that the **Load to Staging Table** task fails.
12. On the **Debug** menu, click **Stop Debugging**.
13. In the **D:\Demofiles\Mod07** folder, note that a file named **Checkpoint.chk** has been created, and that the File System tasks that succeeded have created a folder named **Data** and copied the **Products.csv** file into it.
14. In **SQL Server Data Tools**, view the **Data Flow** tab for the **Load to Staging Table** task, and double-click the **Derive Columns** transformation. Change the expression for the **NewPrice** column to **100**, and click **OK**.
15. View the **Control Flow** tab, and then on the **Debug** menu, click **Start Debugging**. Note that the **Create Folder** and **Copy File** tasks, which succeeded previously, are not re-executed. Only the **Load to Staging Table** task is executed.
16. Stop debugging, and verify that the **Checkpoint.chk** file has been deleted now that the package has been executed successfully.
17. In SQL Server Management Studio, select the top 1,000 rows from the **dbo.StagingTable**, and note that it now contains data six products.
18. Close SQL Server Management Studio and SQL Server Data Tools.

Check Your Knowledge

Question	
<p>You have developed an SSIS package containing two Execute SQL Task tasks. The first task extracts data and the second loads a table. A checkpoint is configured for the package. During execution, the second task fails after loading 50 percent of the table data. What will happen when the package is next executed (assuming the error that caused the failure has been rectified)?</p>	
Select the correct answer.	
<input type="checkbox"/>	The package will rerun all tasks from the beginning.
<input type="checkbox"/>	The package will fail.
<input type="checkbox"/>	The package will rerun the second task, loading the remaining 50 percent of the table data.
<input type="checkbox"/>	The package will rerun the second task from the beginning.
<input type="checkbox"/>	The package will rerun the first task but skip the second task to avoid another failure.

Lab B: Using Transactions and Checkpoints

Scenario

You are concerned that, if the Adventure Works ETL data flow fails, it will leave you with a partially loaded staging database. To avoid this and ensure data integrity, you intend to use transactions and checkpoints.

Objectives

After completing this lab, you will be able to:

- Use transactions.
- Use checkpoints.

Estimated Time: 30 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Using Transactions

Scenario

You have created an SSIS package that uses two data flows to extract, transform, and load Internet sales data. You now want to ensure that package execution always results in a consistent data state, so that if any of the data flows fail, no data is loaded.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. View the Data in the Database
3. Run a Package to Extract Data
4. Implement a Transaction
5. Observe Transaction Behavior

► Task 1: Prepare the Lab Environment

1. Ensure the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab07B\Starter** folder as Administrator.

► Task 2: View the Data in the Database

1. Start **SQL Server Management Studio** and connect to the **MIA-SQL** database engine instance by using Windows authentication.
2. In the **Staging** database, view the contents of the **dbo.Customers** and **dbo.InternetSales** tables to verify that they are both empty.

► Task 3: Run a Package to Extract Data

1. Use **SQL Server Data Tools** to open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab07B\Starter\Ex1** folder.
2. Open the **Extract Internet Sales Data.dtsx** package and examine its control flow.

3. Run the package, noting that the **Extract Customers** task succeeds, but the **Extract Internet Sales** task fails. When execution is complete, stop debugging.
4. In **SQL Server Management Studio**, verify that the **dbo.InternetSales** table is still empty, but the **dbo.Customers** table now contains customer records.
5. In **SQL Server Management Studio**, execute the following Transact-SQL query to reset the staging tables:

```
TRUNCATE TABLE Staging.dbo.Customers;
```

► Task 4: Implement a Transaction

1. Configure the **Extract Customer Sales Data** sequence container in the **Extract Internet Sales Data.dtsx** package so that it requires a transaction.
2. Ensure that the **Extract Customers** and **Extract Internet Sales** tasks both support transactions, and configure them so that, if they fail, their parent also fails.

► Task 5: Observe Transaction Behavior

1. Run the **Extract Internet Sales Data.dtsx** package, noting once again that the **Extract Customers** task succeeds, but the **Extract Internet Sales** task fails. Note also that the **Extract Customer Sales Data** sequence container fails. When execution is complete, stop debugging.
2. In **SQL Server Management Studio**, verify that both the **dbo.InternetSales** and **dbo.Customers** tables are empty.
3. View the data flow for the **Extract Internet Sales** task, and modify the expression in the **Calculate Sales Amount** derived column transformation to resemble that shown below:

```
UnitPrice * OrderQuantity
```

4. Run the **Extract Internet Sales Data.dtsx** package, noting that the **Extract Customers** and **Extract Internet Sales** tasks both succeed. When execution is complete, stop debugging.
5. In **SQL Server Management Studio**, verify that both the **dbo.InternetSales** and **dbo.Customers** tables contain data.
6. Close **SQL Server Data Tools** when you have completed the exercise.

Results: After this exercise, you should have a package that uses a transaction to ensure that all Data Flow tasks succeed or fail as an atomic unit of work.

Exercise 2: Using Checkpoints

Scenario

You have created an SSIS package that uses two data flows to extract, transform, and load reseller sales data. You now want to ensure that, if any task in the package fails, it can be restarted without re-executing the tasks that had previously succeeded.

The main tasks for this exercise are as follows:

1. View the Data in the Database
2. Run a Package to Extract Data
3. Implement Checkpoints
4. Observe Checkpoint Behavior

► Task 1: View the Data in the Database

1. Use **SQL Server Management Studio** to view the contents of the **dbo.Resellers** and **dbo.ResellerSales** tables in the **Staging** database on the **MIA-SQL** database engine instance.
2. Verify that both of these tables are empty.

► Task 2: Run a Package to Extract Data

1. Open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab07B\Starter\Ex2** folder.
2. Open the **Extract Reseller Data.dtsx** package and examine its control flow.
3. Run the package, noting that the **Extract Resellers** task succeeds, but the **Extract Reseller Sales** task fails. When execution is complete, stop debugging.
4. In **SQL Server Management Studio**, verify that the **dbo.ResellerSales** table is still empty, but the **dbo.Resellers** table now contains reseller records.
5. In **SQL Server Management Studio**, execute the following Transact-SQL query to reset the staging tables:

```
TRUNCATE TABLE Staging.dbo.Resellers;
```

► Task 3: Implement Checkpoints

1. Set the following properties of the **Extract Reseller Data** package:
 - **CheckpointFileName**: D:\ETL\CheckPoint.chk
 - **CheckpointUsage**: IfExists
 - **SaveCheckpoints**: True
2. Configure the properties of the **Extract Resellers** and **Extract Reseller Sales** tasks so that, if they fail, the package also fails.

► Task 4: Observe Checkpoint Behavior

1. View the contents of the **D:\ETL** folder and verify that no file named **CheckPoint.chk** exists.
2. Run the **Extract Reseller Data.dtsx** package, noting once again that the **Extract Resellers** task succeeds, but the **Extract Reseller Sales** task fails. When execution is complete, stop debugging.
3. View the contents of the **D:\ETL** folder and verify that a file named **CheckPoint.chk** has been created.

4. In **SQL Server Management Studio**, verify that the **dbo.ResellerSales** table is still empty, but the **dbo.Resellers** table now contains reseller records.
5. View the data flow for the **Extract Reseller Sales** task, and modify the expression to match the following code sample:

```
UnitPrice * OrderQuantity
```

6. Run the **Extract Reseller Sales Data.dtsx** package, noting the **Extract Resellers** task is not re-executed, and package execution starts with the **Extract Reseller Sales** task, which failed on the last attempt. When execution is complete, stop debugging.
7. In **SQL Server Management Studio**, verify that the **dbo.ResellerSales** table now contains data.
8. Close **SQL Server Data Tools** when you have completed the exercise.

Results: After this exercise, you should have a package that uses checkpoints, so that execution can be restarted at the point of failure of the previous execution.

Module Review and Takeaways

In this module, you have learned how to implement control flow in an SSIS package, and how to use transactions and checkpoints to ensure data integrity when a package fails.

Review Question(s)

Question: You have an existing SSIS package containing three tasks. You want Task 3 to run if Task 1 or Task 2 fails. How can you accomplish this?

Question: Which container should you use to perform the same task once for each file in a folder?

Question: Your package includes an FTP task that downloads a large file from an FTP folder and a Data Flow task that inserts data from the file into a database. The Data Flow task might fail if the database is unavailable, in which case you plan to run the package again, after bringing the database online. How can you avoid downloading the file again when the package is re-executed?

MCT USE ONLY. STUDENT USE PROHIBITED

Module 8

Debugging and Troubleshooting SSIS Packages

Contents:

Module Overview	8-1
Lesson 1: Debugging an SSIS Package	8-2
Lesson 2: Logging SSIS Package Events	8-8
Lesson 3: Handling Errors in an SSIS Package	8-13
Lab: Debugging and Troubleshooting an SSIS Package	8-19
Module Review and Takeaways	8-24

Module Overview

As you develop more complex SQL Server® Integration Services (SSIS) packages, it's important to be familiar with the tools and techniques you can use to debug package execution and handle errors. This module describes how you can debug packages to find the cause of errors that occur during execution. It then discusses the logging functionality built into SSIS, which you can use to log events for troubleshooting purposes. Finally, the module describes common approaches for handling errors in control flow and data flow.

Objectives

After completing this module, you will be able to:

- Debug an SSIS package.
- Implement logging for an SSIS package.
- Handle errors in an SSIS package.

Lesson 1

Debugging an SSIS Package

When you are developing an application, misconfiguration of tasks or data flow components, or errors in variable definitions or expressions, can lead to unexpected behavior. Even if you develop your package perfectly, many potential problems might arise during execution—for example, there may be a missing or misnamed file, or an invalid data value. It is therefore important to be able to use debugging techniques to find the cause of these problems, and formulate a solution.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the tools and techniques for debugging SSIS packages.
- View package execution events.
- Use breakpoints to pause package execution.
- View variable values and status while debugging.
- Use data viewers to view data flow values while debugging.

Overview of SSIS Debugging

Debugging is the process of finding the source of problems that occur during package execution, either during development or in a package that has been deployed to a production environment.

Debugging During Development

At design time, SSIS developers can use a variety of Microsoft® Visual Studio® debugging techniques to find problems in control flow and data flow processes. These techniques include:

- Observing row counts and task outcome indicators when running packages in the debugging environment.
- Viewing events that are recorded during package execution. These events are shown in the Progress tab during execution, and in the Execution Results tab after execution. Events are also shown in the Output window during and after each execution.
- Stepping through package execution by setting breakpoints that pause execution at specific points in the control flow.
- Viewing variable values while debugging.
- Viewing the rows that pass through the data flow pipeline by attaching data viewers to data flow paths.

- Debugging During Development
 - Observe row counts and task outcome
 - View events in the Output window and Progress/Execution Results tab
 - Step through package execution
 - Track variable values
 - View data in the data flow
- Debugging in the Production Environment
 - View package execution logs
 - Create a dump file

 **Note:** Microsoft Visual Studio includes a number of debugging windows and tools that are primarily designed for debugging software solutions with programming languages, such as Microsoft Visual C#®. This lesson focuses on the debugging tools in the Visual Studio environment, which has the SQL Server Data Tools for BI add-in and is designed for debugging SSIS packages.

Debugging in the Production Environment

It is common for problems to occur during execution after a package has been deployed to the production environment. In this scenario, if the package source project is available, you can use the techniques previously described. However, you can also debug the package by examining any log files that it is configured to generate, or by using the **dtexec** or **dtutil** utilities to generate a dump file. These files contain information about system variable values and settings that you can use to diagnose a problem with a package.

Viewing Package Execution Events

You can think of a package execution as a sequence of events generated by the tasks and containers in the package control flow. When you run a package in debug mode in the development environment, these events are recorded and displayed in two locations. To run a package in debug mode, you can use any of the following techniques:

- On the **Debug** menu, click **Start Debugging**.
- Click the **Start Debugging** button on the toolbar.
- Press F5.

- Package execution is a sequence of events generated by tasks and containers
- During debugging, events are shown:
 - Progress/Execution Results tab
 - Output window

The Progress/Execution Results Tab

During execution, the Progress tab of the SSIS package designer shows a hierarchical view of the package and its containers and tasks, displaying information about events that occur during execution. When execution is complete, the tab's name changes to Execution Results and shows the entire event tree for the completed execution.

You can enable or disable the display of messages on the Progress tab by toggling Debug Progress Reporting on the SSIS menu. Disabling progress reporting can help improve performance when debugging complex packages.

The Output Window

The Output window shows the list of events that occur during execution. After execution completes, you can review the Output window to find details of the events that occurred.

The Output window and the Progress/Execution Results tab are useful resources for troubleshooting errors during package execution. As an SSIS developer, you should habitually review the events in these windows when debugging packages.

Breakpoints

To troubleshoot a problem with an event in a package, you can use breakpoints to pause execution at the stage in the control flow where the error occurs.

You can create a breakpoint for events raised by any container or task in a control flow. The simplest way to create a breakpoint is to select the task or container where you want to pause execution and, on the Debug menu, click **Toggle Breakpoint**. This adds a breakpoint at the OnPreExecute event of the selected task or container, which is the first event a task or container raises during package execution.

For greater control of when a breakpoint pauses execution, you can right-click any task or container and click **Edit Breakpoints** to display the Set Breakpoints dialog box for the task or container. In this dialog box you can:

- Enable breakpoints for any event supported by the task or container.
- Specify a Hit Count Type and Hit Count Value to control how often the event is ignored before the breakpoint pauses execution. You can set the Hit Count Type to one of the following settings:
 - **Always**. The Hit Count Value is ignored and execution is always paused at this event.
 - **Hit count equals**. Execution is paused when the event has been raised the number of times specified in the Hit Count property.
 - **Hit greater or equal**. Execution is paused when the event has been raised the number of times specified in the Hit Count property or more.
 - **Hit count multiple**. Execution is paused when the event has been raised a number of times that is a multiple of the Hit Count property or more.

You can view and manage all the breakpoints that are set in a package in the Breakpoints window. You can display this window by clicking the **Debug** menu, clicking **Windows**, and clicking **Breakpoints**.

- Add breakpoints to halt execution when debugging
- Specify breakpoint conditions:
 - Event
 - Hit Count
- Manage breakpoints in the Breakpoints window

Variable and Status Windows

When you have used a breakpoint to pause package execution, it can be useful to view the current values assigned to variables, parameters, and other system settings. SQL Server Data Tools provides two windows you can use to observe these values while debugging.

The Locals Window

The Locals window is a pane in the SQL Server Data Tools environment that lists all the system settings, variables, and parameters that are currently in scope. You can use this window to find current values for these settings, variables, and parameters in the execution context.

- Locals window: shows in-scope variables and status
- Watch windows: show selected variables

To view the Locals window when package execution is paused by a breakpoint, click **Windows** on the **Debug** menu, and then click **Locals**.

Watch Windows

If you want to track specific variable or parameter values while debugging, you can add a watch for each value you want to track. Watched values are shown in watch windows named Watch 1, Watch 2, Watch 3, and Watch 4. However, in most SSIS debugging scenarios, only Watch 1 is used.

To display a watch window while debugging, on the **Debug** menu, click **Windows**, click **Watch**, and then click the watch window you want to display.

To add a value to the Watch 1 window, right-click the variable or parameter you want to track in the Locals window and click **Add Watch**.

To add a variable or parameter to another watch window, drag it from the Locals window to the watch window in which you want it to display.

Data Viewers

Most SSIS packages are designed primarily to transfer data. When debugging a package, it can be useful to examine the data as it passes through the data flow. Data viewers provide a way to view the data rows as they pass along data flow paths between sources, transformations, and destinations.

Enabling a Data Viewer

To enable a data viewer, right-click a data flow path on the **Data Flow** tab and click **Enable Data Viewer**. Alternatively, you can double-click a data flow path on the **Data Viewer** tab of the **Data Flow Path Editor** dialog box. The dialog box also enables you to select specific columns to include in the data viewer.

Viewing Data in the Data Flow

A data viewer behaves like a breakpoint and pauses execution at the data flow path on which it is defined. When a data viewer pauses execution, a window containing the data in the data flow path is displayed, which enables you to examine the data at various stages. After you examine the data, you can resume execution by clicking the green **Continue** arrow button in the data viewer window. If you no longer require the data viewer, you can remove it by clicking the **Detach** button in the data viewer window.

Copying Data from a Data Viewer

The data viewer window includes a Copy button, which you can use to copy the contents of the data viewer to the Microsoft Windows® clipboard. When a data flow contains a large number of rows, it can be useful to copy the contents of a data viewer and paste the data into a tool such as Notepad or Microsoft Excel® for further examination.

- Enable data viewers on data flow paths
- View data as it passes through the data flow
- Copy data for further investigation

Demonstration: Debugging a Package

In this demonstration, you will see how to:

- Add a breakpoint.
- View variables while debugging.
- Enable a data viewer.

Demonstration Steps

Add a Breakpoint

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are started, and log onto **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod08** folder, run **Setup.cmd** as Administrator. Click **Yes** when prompted.
3. Start Visual Studio and open the **Debugging.sln** solution in the **D:\Demofiles\Mod08** folder.
4. In Solution Explorer, double-click **Debugging Demo.dtsx**. This package includes a control flow that performs the following tasks:
 - Copies a text file using a variable named **User::sourceFile** to determine the source path and a variable named **User::copiedFile** to determine the destination path.
 - Uses a data flow to extract the data from the text file, convert columns to appropriate data types, and load the resulting data into a database table.
 - Deletes the copied file.
5. On the **Debug** menu, click **Start Debugging**, and note that the first task fails.
6. On the **Debug** menu, click **Stop Debugging**.
7. Click the **Copy Source File** task and on the **Debug** menu, click **Toggle Breakpoint**.
8. Right-click the **Copy Source File** task and click **Edit Breakpoints**. Note that you can use this dialog box to control the events and conditions for breakpoints in your package. When you toggle a breakpoint, by default it is enabled for the **OnPreExecute** event with a **Hit Count Type** value of **Always**. Click **OK**.
9. On the **Debug** menu, click **Start Debugging**. Note that execution stops at the breakpoint.

View Variables While Debugging

1. With execution stopped at the breakpoint, on the **Debug** menu, point to **Windows**, and then click **Locals**.
2. In the Locals pane, expand **Variables** and find the **User::copiedFile** variable. Right-click it, and then click **Add Watch**. The Watch 1 pane appears with the **User::copiedFile** variable displayed.
3. Click the Locals pane, right-click the **User::sourceFile** variable, and then click **Add Watch**. The Watch 1 pane displays the **User::copiedFile** and **User::sourceFile** variables.
Note that the value of the **User::sourceFile** variable is **D:\\Demofiles\\Mod08\\Products.txt** (“\” is used as an escape character, so “\\” is used to indicate a literal “\” string). In the **D:\Demofiles\Mod08** folder, note that the file is actually named **Products.csv**.
4. On the **Debug** menu, click **Stop Debugging**.
5. On the **SSIS** menu, click **Variables**.
6. In the Variables window, change the value for the **sourceFile** variable to **D:\Demofiles\Mod08\Products.csv**.

7. Close the Variables window.
8. On the **Debug** menu, click **Start Debugging**.
9. Observe the variable values in the Watch 1 pane.
Note that the sourceFile variable now refers to the correct file.
10. On the **Debug** menu, click **Continue** and note that the **Load Data** task now fails.
11. On the **Debug** menu, click **Stop Debugging**.

Enable a Data Viewer

1. Double-click the **Load Data** task to view the data flow design surface.
2. Right-click the data flow path between **Products File** and **Data Conversion**, and click **Enable Data Viewer**.
3. Double-click the data flow path between **Products File** and **Data Conversion**.
4. In the **Data Flow Path Editor** dialog box, click the **Data Viewer** tab. Note that you can use this tab to enable the data viewer and specify which columns should be included, and that by default, all columns are included. Click **OK**.
5. Click the **Control Flow** tab and verify that a breakpoint is still enabled on the **Copy Source File** task. Then, on the **Debug** menu, click **Start Debugging**.
6. When execution stops at the breakpoint, on the **Debug** menu, click **Continue**.
7. When the data viewer window is displayed, resize it so you can see the data it contains, and note that the Price column for the second row contains a "-" character instead of a number.
8. In the data viewer window, click **Copy Data** and then click the green **continue** button in the data viewer window. Close the window.
9. When execution stops because the data flow task has failed, on the **Debug** menu, click **Stop Debugging**.
10. Start Notepad, and on the **Edit** menu, click **Paste**.
11. Review the data you have pasted from the data viewer.
12. Close Notepad without saving the file, and then close Visual Studio.

Question: You have executed a package in Visual Studio and a task failed unexpectedly. Where can you review information about the package execution to help determine the cause of the problem?

Lesson 2

Logging SSIS Package Events

Visual Studio debugging tools can be extremely useful when developing a package. However, after a package is in production, it can be easier to diagnose a problem if the package provides details in a log of the events that occurred during execution. In addition to using a log for troubleshooting, you might want to record details of package execution for auditing or performance benchmarking purposes. Planning and implementing a suitable logging solution is an important part of developing a package. SSIS includes built-in functionality to help you accomplish this.

Lesson Objectives

After completing this lesson, you will be able to:

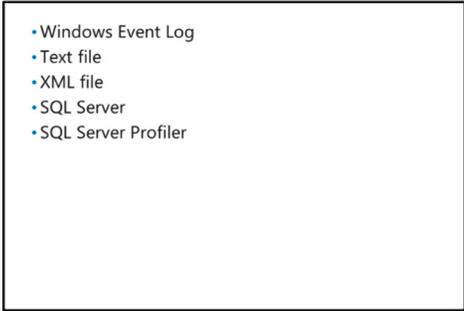
- Describe the log providers available in SSIS.
- Describe the events that can be logged and the schema for logging information.
- Implement logging in an SSIS package.
- View logged events.

SSIS Log Providers

The SSIS logging architecture supports the recording of event information to one or more logs. Each log is accessed through a log provider that determines its type and the connection details used to access it.

SSIS includes the following log providers:

- **Windows Event Log.** This logs event information in the Application Windows event log. No connection manager is required for this log provider.
- **Text File.** This logs event information to a text file specified in a file connection manager.
- **XML File.** This logs event information to an XML file specified in a file connection manager.
- **SQL Server.** This logs event information in the **syssislog** system table in a SQL Server database, which is specified in an OLE DB connection manager.
- **SQL Server Profiler.** This logs event information in a .trc file that can be examined in SQL Server profiler. The location of the .trc file is specified in a file connection manager. This log provider is only available in 32-bit execution environments.

- 
- Windows Event Log
 - Text file
 - XML file
 - SQL Server
 - SQL Server Profiler

Additionally, software developers can use the Microsoft .NET Framework to develop custom log providers.

When deciding which log providers to include in your logging solution, you should generally try to comply with standard logging procedures in the existing IT infrastructure environment. For example, if administrators in the organization typically use the Windows Event Log as the primary source of troubleshooting information, you should consider using it for your SSIS packages. When using files or SQL Server tables for logging, you should also consider the security of the log, which might contain sensitive information.

Log Events and Schema

When you have determined the log providers you want to use, you can select the events for which you want to create log entries and the details you wish to include.

Log Events

SSIS logging supports the following events:

- **OnError.** This event is raised when an error occurs.
- **OnExecStatusChanged.** This event is raised when a task is paused or resumed.
- **OnInformation.** This event is raised during validation and execution to report information.
- **OnPipelinePostComponentCall.** This event is raised when the data flow engine finishes a call to a component method.
- **OnPipelinePostEndOfRowset.** This event is raised when a component has processed all rows.
- **OnPipelinePreComponentCall.** This event is raised when the data flow engine will call a component method.
- **OnPipelinePreEndOfRowset.** This event is raised when a component will be sent the end of rowset signal.
- **OnPipelinePrePrimeOutput.** This event is raised when the PrimeOutput method will be called on a component.
- **OnPipelineRowsSent.** This event is raised when rows are sent to a data flow component.
- **OnPostExecute.** This event is raised when execution of an executable has completed.
- **OnPreExecute.** This event is raised before an executable starts running.
- **OnPreValidate.** This event is raised when validation of an executable begins.
- **OnProgress.** This event is raised to indicate execution progress for an executable.
- **OnQueryCancelled.** This event is raised when execution is cancelled.
- **OnTaskFailed.** This event is raised when a task fails.
- **OnVariableChangedValue.** This event is raised when a variable has its value changed.
- **OnWarning.** This event is raised when a warning occurs.
- **Diagnostic.** This event is raised to provide diagnostic information.
- **DiagnosticEX.** This event logs data flow lineage to an XML log file.
- **Executable-specific events.** Some containers and tasks provide events that are specific to the executable. For example, a Foreach Loop container provides an event that is raised at the start of each loop iteration.

Log Events	Log Schema
- OnError	- StartTime
- OnExecStatusChanged	- EndTime
- OnInformation	- DataCode
- OnPipelinePostComponentCall	- Computer
- OnPipelinePostEndOfRowset	- Operator
- OnPipelinePostPrimeOutput	- MessageText
- OnPipelinePreComponentCall	- DataBytes
- OnPipelinePreEndOfRowset	- SourceName
- OnPipelinePrePrimeOutput	- SourceID
- OnPipelineRowsSent	- ExecutionID
- OnPostExecute	
- OnPreExecute	
- OnPreValidate	
- OnProgress	
- OnQueryCancelled	
- OnTaskFailed	
- OnVariableChangedValue	
- OnWarning	
- Diagnostic	
- DiagnosticEX	

You should consider the performance overhead incurred by the logging process before you log every event. The choice of events to log depends on the purposes of the logging solution. For example, if your goal is primarily to provide troubleshooting information when exceptions occur, you should consider logging the OnError, OnWarning, and OnTaskFailed events. If your log is to be used for audit purposes, you might want to log the OnInformation event; if you want to use your log to measure package performance, you might consider logging the OnProgress and PipelineComponentTime events.

Log Schema

The specific details that can be logged for each event are defined in the SSIS log schema. This schema includes the following values:

- **StartTime**: when the executable started running.
- **EndTime**: when the executable finished.
- **DataCode**: an integer value indicating the execution result:
 - **0**: Success
 - **1**: Failure
 - **2**: Completed
 - **3**: Cancelled
- **Computer**: the name of the computer on which the package was executed.
- **Operator**: the Windows account that initiated package execution.
- **MessageText**: a message associated with the event.
- **DataBytes**: a byte array specific to the log entry.
- **SourceName**: the name of the executable.
- **SourceID**: the unique identifier for the executable.
- **ExecutionID**: a unique identifier for the running instance of the package.

You can choose to include all elements of the schema in your log, or select individual values to reduce log size and performance overhead. However, the **StartTime**, **EndTime**, and **DataCode** values are always included in the log.

Implementing SSIS Logging

Packages can be thought of as a hierarchy of containers and tasks, with the package itself as the root of the hierarchy. You can configure logging at the package level in the hierarchy, and by default, all child containers, and tasks inherit the same logging settings. If required, you can override inherited logging settings for any container or task in the hierarchy. For example, you might choose to log only OnError events to the Windows Event Log provider at the package level—and inherit these settings for most child containers and tasks— but configure a data flow task within the package to log OnInformation and Diagnostic events to the SQL Server log provider.

1. Add and configure log providers
2. Select containers and tasks to include
3. Select events and details to log
4. Override log settings for child executables if required

To implement logging for an SSIS package, in SQL Server Data Tools, with the package open in the designer, on the SSIS menu, click **Logging** to display the **Configure SSIS Logs** dialog box. Then perform the following steps:

1. **Add and configure log providers.** On the **Providers and Logs** tab of the dialog box, add the log providers you want to use. For providers other than the one for Windows Event Log, you must also specify a connection manager that defines the file or SQL Server instance where you want to write the log information.

2. **Select the containers and tasks to include.** Choose the package container which, by default, selects all child containers and tasks with inherited log settings. You can then clear individual containers and tasks that you do not want to include.
3. **Select events and details to log.** On the **Details** tab of the dialog box, select the events you want to include in the log. By default, all schema fields are logged for the selected events, but you can click the **Advanced** button to choose individual fields.
4. **Override log settings for child executables if required.** If you want to specify individual logging settings for a specific child executable, select the executable in the Containers tree and specify the log provider, events, and details you want to use for that executable.

Viewing Logged Events

You can view logged events in Visual Studio by displaying the Log Events window. When logging is configured for a package, this window shows the selected log event details, even when no log provider is specified.

The Log Events window is a useful tool for troubleshooting packages during development, and also for testing and debugging logging configuration.

To display the Log Events window, on the **SSIS** menu, click **Log Events**.

- Logged events are displayed in the Log Events window
 - Even if no log provider is specified
- Useful for:
 - Troubleshooting
 - Testing a logging strategy

Demonstration: Logging Package Execution

In this demonstration, you will see how to:

- Configure SSIS logging.
- View logged events.

Demonstration Steps

Configure SSIS Logging

1. Ensure you have completed the previous demonstration in this module.
2. Start Visual Studio and open the **Logging.sln** solution in the **D:\Demofiles\Mod08** folder.
3. In Solution Explorer, double-click **Logging Demo.dtsx**.
4. On the **SSIS** menu, click **Logging**.
5. In the **Configure SSIS Logs: Logging Demo** dialog box, in the **Provider type** list, select **SSIS log provider for Windows Event Log**, and then click **Add**. Then select **SSIS log provider for SQL Server** and click **Add**.
6. In the **Configuration** column for the **SSIS log provider for SQL Server**, click the drop-down arrow, and then click the **(local).DemoDW** connection manager.

Note that the Windows Event Log provider requires no configuration.

7. In the **Containers** tree, select the **Logging Demo** check box, and then, on the **Providers and Logs** tab, select the **SSIS log provider for Windows Event Log** check box.
8. On the **Details** tab, select the **OnError** and **OnInformation** events.
9. In the **Containers** tree, clear the **Load Data** check box.
10. In the **Containers** tree, clear the **Load Data** check box. This can override the inherited logging settings for the **Load Data** task.
11. With **Load Data** selected, on the **Providers and Logs** tab, select the **SSIS log provider for SQL Server** check box.
12. On the **Details** tab, select the **OnError** and **OnInformation** events.
13. Click **Advanced** and clear the **Operator** column for the two selected events.
14. Click **OK**.

View Logged Events

1. On the **Debug** menu, click **Start Debugging**.
2. When the **Load Data** task fails, on the **Debug** menu click **Stop Debugging**.
3. On the **SSIS** menu, click **Log Events**. This shows the events that have been logged during the debugging session (if the log is empty, rerun the package, and then view the Log Events window again).
4. In the Windows taskbar, click **Search Windows**, type **Event Viewer**, and then press Enter.
5. Expand **Windows Logs**, and then click **Application**. Note the log entries with a source of **SQLISPackage140**. These are the logged events for the package.
6. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the database engine by using Windows authentication.
7. In Object Explorer, expand **Databases**, expand **DemoDW**, expand **Tables**, and then expand **System Tables**.
8. Right-click **dbo.sysssislog**, and then click **Select Top 1000 Rows**.
9. Review the contents of the table.
10. Close SQL Server Management Studio without saving any files, then close Event Viewer and Visual Studio.

Question: Can you think of some advantages and disadvantages of the different SSIS log providers?

Lesson 3

Handling Errors in an SSIS Package

No matter how much debugging you perform or how much information you log during package execution, exceptions that cause errors can occur in any ETL process. For example, servers can become unavailable, files can be renamed or deleted, and data sources can include invalid entries. A good SSIS solution includes functionality to handle errors by performing compensating tasks and continuing with execution wherever possible, or by ensuring that temporary resources are cleaned up and operators notified where execution cannot be continued.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe approaches for handling errors in an SSIS package.
- Implement event handlers.
- Handle errors in data flows.

Introduction to Error Handling

Errors can occur at any stage in the execution of a package, and SSIS provides several ways to handle them and take corrective action if possible.

Handling Errors in Control Flow

You can use the following techniques to handle errors in package control flow:

- Use Failure precedence constraints to redirect control flow when a task fails. For example, you can use a Failure precedence constraint to execute another task that performs a compensating alternative action to allow the control flow to continue, or to delete any temporary files and send an email notification to an operator.
- Implement event handlers to execute a specific set of tasks when an event occurs in the control flow. For example, you could implement an event handler for the OnError event of the package, and include tasks to delete files and send email notifications in the OnError event handler.

- Handling errors in control flow
 - Failure precedence constraints
 - Event handlers
- Handling errors in data flow
 - Ignore or redirect failed rows

 **Note:** Precedence constraints are discussed in Module 7: *Implementing Control Flow in an SSIS Package*. The remainder of this lesson focuses on using event handlers to handle errors in control flow.

Handling Errors in Data Flow

Errors in data flow can often be caused by invalid or unexpected data values in rows being processed by the data flow pipeline. SSIS data flow components provide the following configuration options for handling rows that cause errors:

- Fail the task if any rows cause an error.

- Ignore errors and continue the data flow.
- Redirect rows that cause an error to the error output of the data flow component.

Implementing Event Handlers

You can add an event handler for the events supported by each executable in the package. To add an event handler to a package, click the **Event Handlers** tab, select the executable and event for which you want to implement a handler, and click the hyperlink on the design surface. This creates a new control flow surface on which you can define the control flow for the event handler. Event handlers can be used for all kinds of control flow tasks, and are not specific to handling errors. However, the OnError and OnTaskFailed events are commonly used for handling error conditions.

- Add an event handler on the Event Handler tab
- Each event handler has its own control flow
- Use contextualized system variables to implement custom logging or notifications

The system variables and configuration values available to tasks in your event handler are specific to the context of the event. For example, the **System::ErrorDescription** variable is populated with an error message during the OnError event.

Because a package is a hierarchy of containers and tasks, each with their own events, you need to consider where best to handle each possible error condition. For example, you can handle a task-specific error in the task's own OnError event or, depending on the MaxErrors property of containers and the package itself, the error could trigger OnError events further up the package hierarchy, where you could also handle the error condition. In general, if you anticipate specific errors that you can resolve or compensate for and continue execution, you should use the OnError event handler for the task or container where the error is likely to occur. To catch errors that cannot be resolved, perform clean-up tasks and notify operators, you should use the OnError event handler of the package.

Handling Data Flow Errors

Data flow components participate in a data flow pipeline through which rows of data are passed along data flow paths. Errors can occur in data flow for a number of reasons, including:

- Rows that contain data of a type that is incompatible with a transformation or destination, such as a decimal field in a text file that is mapped to an integer column in a destination.
- Rows that contain invalid data values, such as a text file that contains a date field with an invalid date value.
- Rows that contain data that will be truncated by a transformation or destination, such as a 50-character text field that is loaded into a table where the mapped column has a maximum length of 40 characters.

- Configure Error Output for data flow components:
 - Fail component
 - Ignore failure
 - Redirect row
- Redirect failed rows with error output path
- DiagnosticEX event for more detailed error information
- Use custom ETL data lineage to trace data back from the destination to the source, retaining all history of transformations and deduplication

- Rows that contain data values that will cause an exception during a transformation, such as a numerical field with a value of zero that is used as a divisor in a derived column transformation.

By default, rows that cause an error result in the failure of the data flow component. However, you can configure many data flow components to ignore rows that contain errors, or to redirect them to the error output data flow path of the component. Ignoring or redirecting failed rows means that the data flow can complete for all other rows; if you have chosen to redirect failed rows, you can use transformations to attempt to correct the invalid data values, or save the failed rows to a file or table for later analysis.

When you configure error output for a data flow component, you can specify different actions for truncations and errors. For example, you could choose to ignore truncations, but redirect rows that contain other errors. Additionally, with some components, you can specify different actions for each column in the data flow. You could ignore errors in one column while redirecting rows that have an invalid value in another.

Redirected rows include all the input columns for the data flow component, and two additional columns:

- **ErrorCode:** the numeric code for the error.
- **ErrorColumn:** the original number of the column that caused the error.



Note: Occasionally, ErrorColumn might be populated with a zero. This means the error relates to the entire row of data and not a specific column.

The information contained in the columns ErrorCode and ErrorColumn of the error output may not be sufficient for you to easily identify the column within the data flow that is responsible for the error.

To assist in these situations, you can enable logging for the **DiagnosticEX** event. The **DiagnosticEX** event produces lineage of the data flow in an XML document that can be used to look up column names from the ID found in the **ErrorColumn** column of an error output.

Adding Custom ETL Data Lineage

When you push data through an ETL process and into a data warehouse to be used for reporting, you must be confident that the data you present is entirely accurate. Knowing where data has come from, and how it has changed in the process, is known as data lineage. You can implement custom data lineage within your ETL processes. Data lineage is vital if you are to discover any issues with the data in your dimension tables. You will need to trace the data back from the dimension table to the source, determining any applied transformations along the way, in addition to considering rows removed by deduplication.

Data might be extracted into the initial staging tables from several different sources, including databases, files, and cloud sources. This data might be loaded in parallel, and additional data may be generated during the load. Logging the source in a table and giving it a unique identifier enables you to store each data row with a corresponding source key. Each row can then have its own unique identifier. As data is cleaned, shaped, and transformed, it can be moved or copied to other tables: if the source keys are stored with the data, the lineage can be traced. These relationships are particularly helpful during deduplication. Discarded rows can be redirected in the data flow to an archive table, along with their source keys, retaining a full history of every row, whether it reaches the destination data warehouse or not.

Demonstration: Handling Errors

In this demonstration, you will see how to:

- Implement an event handler.
- Redirect failed rows.
- Add DiagnosticEX logging.

Demonstration Steps

Implement an Event Handler

1. Ensure you have completed the previous demonstration in this module.
2. Start Visual Studio and open the **ErrorHandling.sln** solution in the **D:\Demofiles\Mod08** folder.
3. In Solution Explorer, double-click **Error Handling Demo.dtsx**, and then click the **Event Handlers** tab.
4. On the **Event Handlers** tab, in the **Executable** list, ensure **Error Handling Demo** is selected, and, in the **Event handler** list, ensure **OnError** is selected. Click the hyperlink in the middle of the design surface.
5. On the **SSIS** menu, click **SSIS Toolbox**.
6. In the SSIS Toolbox, double-click the **Send Mail** task.
7. On the design surface, right-click **Send Mail** task, click **Rename**, and change the name to **Notify Administrator**.
8. Double-click the **Notify Administrator** task.
9. In the **Send Mail Task Editor** dialog box, on the **Mail** tab, configure the following properties:
 - **SmtpConnection**: A new connection to the mia-sql.adventureworks.msft SMTP server with default settings
 - **From**: etl@adventureworks.msft
 - **To**: administrator@adventureworks.msft
 - **Subject**: An error has occurred
10. On the **Expressions** tab, in the **Expressions** box, click the ellipsis (...).
11. In the **Property Expressions Editor** dialog box, in the **Property** list, click **MessageSource**, and in the **Expression** box, click the ellipsis (...).
12. In the **Expression Builder** dialog box, expand **Variables and Parameters**, expand **System Variables**, and drag **System::ErrorDescription** to the **Expression** box. Click **OK**.
13. In the **Property Expressions Editor** dialog box, in the row under the **MessageSource** property, in the **Property** list, click **FileAttachments**.
14. In the **Expression** box, click the ellipsis (...).
15. In the **Expression Builder** dialog box, expand **Variables and Parameters**, and drag **User::sourceFile** to the **Expression** box. Click **OK**.
16. In the **Property Expressions Editor** dialog box, click **OK**.
17. In the **Send Mail Task Editor** dialog box, click **OK**.
18. Click the **Control Flow** tab and then, on the **Debug** menu, click **Start Debugging**.

19. When the **Load Data** task fails, click the **Event Handlers** tab to verify that the **OnError** event handler has been executed and then, on the **Debug** menu, click **Stop Debugging**.
20. In the **C:\inetpub\mailroot\Drop** folder, open the most recent email message using Notepad and review its contents. The message body will be encoded, but you should be able to verify that the subject line is "An error has occurred"; this is the subject line you specified when you configured the Send Mail task. Close Notepad.

Redirect Failed Rows

1. In Visual Studio, click the **Data Flow** tab.
2. In the **Data Flow Task** drop-down list, ensure **Load Data** is selected.
3. Double-click the **Data Conversion** task and then, in the **Data Conversion Transformation Editor** dialog box, click **Configure Error Output**.
4. In the **Configure Error Output** dialog box, click the **Error** cell for the Numeric **ProductID** column, then hold the **Ctrl** key and click the **Error** cell for the **Numeric Price** column so that both cells are selected.
5. In the **Set this value to the selected cells** list, click **Redirect row** and click **Apply**.
6. In the **Configure Error Output** dialog box, click **OK**.
7. In the **Data Conversion Transformation Editor** dialog box, click **OK**.
8. In the SSIS Toolbox, in the **Other Destinations** section, double-click **Flat File Destination**.
9. On the design surface, right-click the **Flat File Destination** task, click **Rename**, and change the name to **Invalid Rows**.
10. Move **Invalid Rows** to the right of the **Data Conversion** task then click the **Data Conversion** task and drag the red data path to the **Invalid Rows** task.
11. In the **Configure Error Output** dialog box, verify that the **Numeric ProductID** and **Numeric Price** columns both have an **Error** value of **Redirect row**, and click **OK**.
12. Double-click the **Invalid Rows** task and then click **New**.
13. In the **Flat File Format** dialog box, ensure **Delimited** is selected and click **OK**.
14. In the **Flat File Connection Manager Editor** dialog box, in the **Connection manager name** box, type **Invalid Rows CSV File**.
15. In the **File name** box, type **D:\Demofiles\Mod08\InvalidRows.csv**, and click **OK**.
16. In the **Flat File Destination Editor** dialog box, click the **Mappings** tab and note that the input columns include the columns from the data flow, an **ErrorCode** column, and an **ErrorColumn** column. Click **OK**.
17. Click the **Control Flow** tab, on the **Debug** menu, click **Start Debugging**. Note that all tasks succeed, and on the **Data Flow** tab note the row counts that pass through each data flow path.
18. On the **Debug** menu, click **Stop Debugging**.
19. Using Notepad, open **InvalidRows.csv** in the **D:\Demofiles\Mod08** folder.
20. Review the rows that were redirected, and then close Notepad without saving any changes.
21. Return to Visual Studio.

Enable DiagnosticEX Event Logging

1. On the **SSIS** menu, click **Logging**.
2. In the **Configure SSIS Logs: Error Handling Demo**, in the **Provider type** list, select **SSIS Log Provider for XML Files** and click **Add**.
3. In the **Configuration** column for the **SSIS log provider for XML Files**, click the drop-down arrow, and then click **<New connection>**.
4. In the **File Connection Manager Editor** dialog box, in the **Usage type** list, click **Create file**, in the **File** box, type **D:\DemoFiles\Mod08\DiagnosticEX.xml**, and then click **OK**.
5. In the **Containers** tree, select the **Error Handling Demo** check box, and then, on the **Providers and Logs** tab, select the **SSIS log provider for XML files** check box.
6. On the **Details** tab, select the **DiagnosticEX** check box, and then click **OK**.
7. Click the **Control Flow** tab and on the **Debug** menu, click **Start Debugging**. Note that all tasks succeed.
8. Start Notepad and open **InvalidRows.csv** in the **D:\Demofiles\Mod08** folder. Note the value 10 at the end of the row—this is the **ErrorColumn** value.
9. Still using Notepad, open **DiagnosticEX.xml** in the **D:\Demofiles\Mod08** folder.
10. In the file, find the text **DTS:ID="10"**. This will be followed by the text **DTS:IdentificationString="Data Conversion.Outputs[Data Conversion Output].Columns[Numeric Price]"** which details the column details for the error row.
11. Close Notepad.
12. In Visual Studio, stop debugging and close without saving any changes.

Question: In what situations might you use the package level OnError event handler?

Lab: Debugging and Troubleshooting an SSIS Package

Scenario

The ETL process for Adventure Works Cycles occasionally fails when extracting data from text files generated by the company's financial accounts system. You plan to debug the ETL process to identify the source of the problem and implement a solution to handle any errors.

Objectives

After completing this lab, you will be able to:

- Debug an SSIS package.
- Log SSIS package execution.
- Implement an event handler in an SSIS package.
- Handle errors in a data flow.

Estimated Time: 60 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Debugging an SSIS Package

Scenario

You have developed an SSIS package to extract data from text files exported from a financial accounts system and load the data into a staging database. However, while developing the package, you have encountered some errors and you need to debug it to identify the cause.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Run an SSIS Package
3. Add a Breakpoint
4. Add a Data Viewer
5. View Breakpoints
6. Observe Variables while Debugging
7. View Data Copied from a Data Viewer

► Task 1: Prepare the Lab Environment

1. Ensure the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab08\Starter** folder as Administrator.

► Task 2: Run an SSIS Package

1. Open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex1** folder.
2. Open the **Extract Payment Data.dtsx** package and examine its control flow.
3. Run the package, noting that it fails. When execution is complete, stop debugging.

► Task 3: Add a Breakpoint

1. Add a breakpoint to the **Extract Payments** data flow task in the Extract Payment Data.dtsx package.
2. Select the appropriate event to ensure that execution is always paused at the beginning of every iteration of the loop.

► Task 4: Add a Data Viewer

- In the data flow for the **Extract Payments** task, add a data viewer to the data flow path between the **Payments File** source and the **Staging DB** destination.

► Task 5: View Breakpoints

- View the Breakpoints window, and note that it contains the breakpoint you defined on the Foreach Loop container and the data viewer you added to the data flow.

► Task 6: Observe Variables while Debugging

1. Start debugging the package, and note that execution pauses at the first breakpoint.
2. View the Locals window and view the configuration values and variables it contains.
3. Add the following variables to the Watch 1 window:
 - User::fName
 - \$Project::AccountsFolderPath
4. Continue execution and note that it pauses at the next breakpoint, which is the data viewer you defined in the data flow. Also note that the data viewer window contains the contents of the file indicated by the **User::fName** variable in the Watch 1 window.
5. Continue execution, observing the value of the **User::fName** variable and the contents of the data viewer window during each iteration of the loop.
6. When the **Extract Payments** task fails, note the value of the **User::fName** variable, and copy the contents of the data viewer window to the clipboard. Then stop debugging and close Visual Studio.

► Task 7: View Data Copied from a Data Viewer

1. Start Notepad and paste the data you copied from the data viewer window into a worksheet.
2. Examine the data and try to determine the errors it contains (Hint: look at payments 4074, 4102, and 4124).

Results: After this exercise, you should have observed the variable values and data flows for each iteration of the loop in the Extract Payment Data.dtsx package. You should also have identified the file that caused the data flow to fail and examined its contents to find the data errors that triggered the failure.

Exercise 2: Logging SSIS Package Execution

Scenario

You have debugged the Extract Payment Data.dtsx package and found some errors in the source data. Now you want to implement logging to assist in diagnosing future errors when the package is deployed in a production environment.

The main tasks for this exercise are as follows:

1. Configure SSIS Logs
2. View Logged Events

► Task 1: Configure SSIS Logs

1. Open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex2** folder.
2. Open the **Extract Payment Data.dtsx** package and configure logging with the following settings:
 - Enable logging for the package, and inherit logging settings for all child executables.
 - Log package execution events to the Windows Event Log.
 - Log all available details for the OnError and OnTaskFailed events.

► Task 2: View Logged Events

1. Run the package in Visual Studio, and note that it fails. When execution is complete, stop debugging.
2. View the Log Events window and note the logged events it contains. If no events are logged, rerun the package and look again, then close Visual Studio.
3. View the Application window event log in the Event Viewer administrative tool.

Results: After this exercise, you should have a package that logs event details to the Windows Event Log.

Exercise 3: Implementing an Event Handler

Scenario

You have debugged the Extract Payment Data.dtsx package and observed how errors in the source data can cause it to fail. You now want to implement an error handler that copies the invalid source data file to a folder for later examination and notifies an administrator.

The main tasks for this exercise are as follows:

1. Create an Event Handler
2. Add a File System Task
3. Add a Send Mail Task
4. Test the Event Handler

► Task 1: Create an Event Handler

1. Open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex3** folder.
2. In the **Extract Payment Data.dtsx** package, create an event handler for the **OnError** event of the **Extract Payment Data** executable.

► Task 2: Add a File System Task

1. Add a file system task to the control flow for the OnError event handler, and name it **Copy Failed File**.
2. Configure the Copy Failed File task as follows:
 - The task should perform a **Copy file** operation.
 - Use the **Payments File** connection manager as the source connection.
 - Create a new connection manager to create a file named **D:\ETL\FailedPayments.csv** for the destination.
 - The task should overwrite the destination file if it already exists.
3. Configure the connection manager you created for the destination file to use the following expression for its **ConnectionString** property. Instead of the name configured in the connection manager, the copied file is named using a combination of the unique package execution ID and the name of the source file:

```
"D:\\ETL\\" + @[System::ExecutionInstanceGUID] + @[User::fName]
```

► Task 3: Add a Send Mail Task

1. Add a send mail task to the control flow for the OnError event handler, and name it **Send Notification**.
2. Connect the **Copy Failed File** task to the **Send Notification** task with a **Completion precedence** constraint.
3. Configure the **Notification** task as follows:
 - Use the **Local SMTP Server** connection manager.
 - Send a high-priority email message from etl@adventureworks.msft to student@adventureworks.msft with the subject "An error occurred".
 - Use the following expression to set the **MessageSource** property:

```
@[User::fName] + " failed to load. " +  
@[System::ErrorDescription]
```

► Task 4: Test the Event Handler

1. Run the package in Visual Studio and verify that the event handler is executed. Close Visual Studio, saving your changes, if prompted.
2. Verify that the source file containing invalid data is copied to the **D:\ETL** folder with a name similar to {1234ABCD-1234-ABCD-1234-ABCD1234}Payments - EU.csv.
3. View the contents of the **C:\inetpub\mailroot\Drop** folder, and verify that an email was sent for each error that occurred. You can examine each email file by using Notepad.

Results: After this exercise, you should have a package that includes an event handler for the OnError event. The event handler should create a copy of files that contain invalid data and send an email message.

Exercise 4: Handling Errors in a Data Flow

Scenario

You have implemented an error handler that notifies an operator when a data flow fails. However, you would like to handle errors in the data flow so that only rows containing invalid data are not loaded, and the rest of the data flow succeeds.

The main tasks for this exercise are as follows:

1. Redirect Data Flow Errors
2. View Invalid Data Flow Rows

► Task 1: Redirect Data Flow Errors

1. Open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex4** folder.
2. Open the **Extract Payment Data.dtsx** package and view the data flow for the **Extract Payments** task.
3. Configure the error output of the **Staging DB** destination to redirect rows that contain an error.
4. Add a flat file destination to the data flow and name it **Invalid Rows**. Then configure the **Invalid Rows** destination as follows:
 - Create a new connection manager named **Invalid Payment Records** for a delimited file named **D:\ETL\InvalidPaymentsLog.csv**.
 - Do not overwrite data in the text file if it already exists.
 - Map all columns, including **ErrorCode** and **ErrorColumn**, to fields in the text file.

► Task 2: View Invalid Data Flow Rows

1. Run the package in debug mode and note that it succeeds. When execution is complete, stop debugging and close Visual Studio.
2. Use Notepad to view the contents of the **InvalidPaymentsLog.csv** file in the **D:\ETL** folder and note the rows that contain invalid values.

Results: After this exercise, you should have a package that includes a data flow where rows containing errors are redirected to a text file.

Module Review and Takeaways

In this module, you have learned how to:

- Debug an SSIS package.
- Implement logging for an SSIS package.
- Handle errors in an SSIS package.

Review Question(s)

Question: You have configured logging with the SSIS log provider for SQL Server. Where can you view the logged event information?

Question: You suspect a data flow is failing because some values in a source text file are too long for the columns in the destination. How can you handle this problem?

Module 9

Implementing a Data Extraction Solution

Contents:

Module Overview	9-1
Lesson 1: Introduction to Incremental ETL	9-2
Lesson 2: Extracting Modified Data	9-12
Lab A: Extracting Modified Data	9-26
Lesson 3: Loading Modified Data	9-38
Lesson 4: Temporal Tables	9-51
Lab B: Loading a Data Warehouse	9-62
Module Review and Takeaways	9-72

Module Overview

A data warehousing solution generally needs to refresh the Data Warehouse (DW) regularly to reflect new and modified data in the source systems on which it is based. It is important to implement a refresh process that has a minimal impact on network and processing resources, and which helps you to retain historical data in the DW, while reflecting changes and additions to business entities in transactional systems.

This module describes the techniques you can use to implement a SQL Server® Integration Services (SSIS) solution that supports incremental DW loads and changing data using extract, transform and load (ETL).

Objectives

After completing this module, you will be able to:

- Plan and implement incremental ETL within SSIS packages.
- Understand temporal tables and Slowly Changing Dimensions (SCDs).
- Extract modified data, use Change Data Capture (CDC) and Change Tracking (CT).
- Transform and load incremental data within the data warehouse.
- Describe the considerations for planning data loads.
- Use SSIS to load new and modified data into a data warehouse.
- Use Transact-SQL techniques to load data into a data warehouse.

Lesson 1

Introduction to Incremental ETL

Most data warehousing solutions use an incremental ETL process to refresh the data warehouse with new and modified data from source systems. Implementing an effective incremental ETL process presents a number of challenges, for which common solution designs have been identified. By understanding some of the key features of an incremental ETL process, you can design an effective data warehouse refresh solution that meets your analytical and reporting needs, while maximizing performance and resource efficiency.

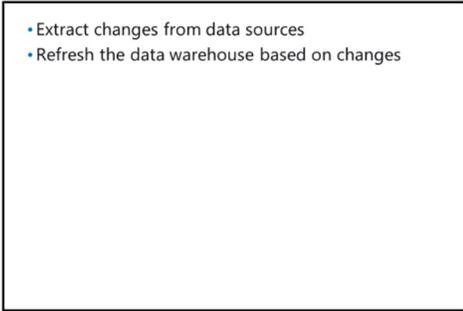
Lesson Objectives

After completing this lesson, you will be able to:

- Describe a typical data warehouse refresh scenario.
- Describe considerations for implementing an incremental ETL process.
- Describe common ETL architectures.
- Plan data extraction windows.
- Plan transformations for different types of data, such as dynamic or Slowly Changing Dimensions.
- Document data flows.
- Understand the basics of Slowly Changing Dimensions.

Overview of Data Warehouse Load Cycles

Data warehousing solutions will include regular refreshing of data in the data warehouse to reflect new and modified data in the source systems on which it is based. For each load cycle, data is extracted from the source systems, usually to a staging area, and then loaded into the data warehouse. The frequency of the refresh process depends on how up to date the analytical and reporting data in the data warehouse needs to be. In some cases, you might choose to implement a different refresh cycle for each group of related data sources.

- 
- Extract changes from data sources
 - Refresh the data warehouse based on changes

In rare cases, it can be appropriate to replace the data warehouse data with fresh data from the data sources during each load cycle. However, a more common approach is to use an incremental ETL process to extract only rows that have been inserted or modified in the source systems. Rows are then inserted or updated in the data warehouse to reflect the extracted data. This reduces the volume of data being transferred, minimizing the effect of the ETL process on network bandwidth and processing resources.

Considerations for Incremental ETL

When planning an incremental ETL process, there are a number of factors to consider:

Data Modifications to Be Tracked

One of the primary considerations for planning an incremental ETL process is to identify the kinds of data modifications you need to track in source systems. Specifically, you should consider the following kinds of modifications:

- **Inserts:** for example, new sales transactions or new customer registrations.
- **Updates:** for example, a change of a customer's telephone number or address.
- **Deletes:** for example, the removal of a discontinued item from a product catalog.

- Data modifications to be tracked
- Load order
- Dimension keys
- Updating dimension members
- Updating fact records

Most data warehousing solutions include inserted and updated records in refresh cycles. However, you must give special consideration to deleted records because propagating deletions to the data warehouse results in the loss of historical reporting data.

Load Order

A data warehouse can include dependencies between tables. For example, rows in a fact table generally include foreign key references to rows in dimension tables; some dimension tables include foreign key references to subdimension tables. For this reason, you should generally design your incremental ETL process to load subdimension tables first, then dimension tables, and finally fact tables. If this is not possible, you can load inferred members as minimal placeholder records for dimension members that are referenced by other tables, and which will be loaded later.

 **Note:** Inferred members are normally used to create a placeholder record for a missing dimension member referenced by a fact record. For example, the data to be loaded into a fact table for sales orders might include a reference to a product for which no dimension record has yet been loaded. In this case, you can create an inferred member for the product that contains the required key values but null columns for all other attributes. You can then update the inferred member record on a subsequent load of product data.

Dimension Keys

The keys used to identify rows in dimension tables are usually independent from the business keys used in source systems, and are referred to as surrogate keys. When loading data into a data warehouse, you need to consider how you will identify the appropriate dimension key value to use in the following scenarios:

- Determining whether or not a staged record represents a new dimension member or an update to an existing one, and if it is an update, applying the update to the appropriate dimension record.
- Determining the appropriate foreign key values to use in a fact table that references a dimension table, or in a dimension table that references a subdimension table.

In many data warehouse designs, the source business key for each dimension member is retained as an alternative key in the data warehouse, and can therefore be used to look up the corresponding dimension key. In other cases, dimension members must be found by matching a unique combination of multiple columns.

Updating Dimension Members

When refreshing dimension tables, you must consider whether changes to individual dimension attributes will have a material effect on historical reporting and analysis. Dimension attributes can be categorized as one of three kinds:

- **Fixed:** the attribute value cannot be changed. For example, you might enforce a rule that prevents changes to a product name after it has been loaded into the dimension table.
- **Changing:** the attribute value can change without affecting historical reporting and analytics. For example, a customer's telephone number might change, but it is unlikely that any historical business reporting or analytics will aggregate measures by telephone number. The change can therefore be made without having to retain the previous telephone number.
- **Historical:** the attribute value can change, but the previous value must be retained for historical reporting and analysis. For example, a customer might move from Edinburgh to New York, but reports and analysis must associate all sales that occurred to that customer before the move with Edinburgh; and all sales after the move with New York.

Updating Fact Records

When refreshing the data warehouse, you must consider whether you will accept updates to fact records. Often, your data warehouse design will only contain complete fact records, so incomplete records will not be loaded. In some cases, however, you might want to include an incomplete fact record in the data warehouse that will be updated during a later refresh cycle.

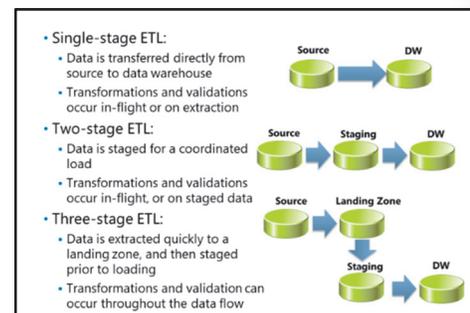
For example, you might choose to include a fact record for a sales order where the sale has been completed, but the item has not yet been delivered. If the record includes a column for the delivery date, you might initially store a null value in this column, and then update the record during a later refresh, after the order has been delivered.

While some data warehousing professionals allow updates to the existing record in the fact table, others prefer to support changes by replacing the existing fact record with a new one. In most cases, the delete operation is actually a logical deletion that is achieved by setting a bit value on a column that indicates whether the record is active or not—rather than actually deleting the record from the table.

Common ETL Data Flow Architectures

ETL is concerned with data flow from source systems to the data warehouse. The data flow process can be performed directly from source to target, or in stages. Factors that affect the choice of data flow architecture include:

- The number of data sources.
- The volume of data to be transferred.
- The complexity of validation and transformation operations to be applied to the data.



- How frequently data is generated in source systems, and how long it is retained.
- Suitable times to extract source data while minimizing the impact on performance for users.

Single-stage ETL

In a very small Business Intelligence (BI) solution with few data sources and simple requirements, it might be possible to copy data from data sources to the data warehouse in a single data flow. Basic data validations (such as checking for NULL fields or specific value ranges) and transformations (such as concatenating multiple fields into a single field, or looking up a value from a key), can either be performed during extraction (for example, in the Transact-SQL statement used to retrieve data from a source database) or in-flight (for example, by using transformation components in an SSIS data flow task).

Two-stage ETL

In many cases, a single-stage ETL solution is not suitable because of the complexity or volume of data being transferred. Additionally, if multiple data sources are used, it is common to synchronize loads into the data warehouse. This ensures consistency and integrity across fact and dimension data from different sources, and minimizes the performance impact of the load operations on data warehouse activity. If the data is not ready to extract from all systems at the same time, or if some sources are only available at specific times when others are not, a common approach is to stage the data in an interim location before loading into the data warehouse.

Typically, the structure of the data in the staging area is similar to the source tables, which minimizes the extraction query complexity and duration in the source systems. When all source data is staged, it can then be processed to conform to the data warehouse schema during the load operation—either as it is extracted from the staging tables or during the data flow to the data warehouse.

Staging the data also provides a recovery point for data load failures and means you can retain extracted data for audit and verification purposes.

Three-stage ETL

A two-stage data flow architecture can reduce the extraction overhead and source systems, enabling a coordinated load of data from multiple sources. However, performing validation and transformations during the data flow into the data warehouse can affect load performance, and cause the load to negatively affect data warehouse activity. When large volumes of data must be loaded into the data warehouse, it is important to minimize load times by preparing the data as much as possible before performing the operation.

For BI solutions that involve loading large volumes of data, a three-stage ETL process is recommended. In this data flow architecture, the data is initially extracted to tables that closely match the source system schemas (often referred to as a “landing zone”). From here, the data is validated and transformed as it is loaded into staging tables that more closely resemble the target data warehouse tables. Finally, the conformed and validated data can be loaded into the data warehouse tables.

Planning Extraction Windows

To help you determine when to perform the data extraction process, consider the following questions:

How frequently is new data generated in the source systems, and for how long is it retained?

Some business applications generate only a few transactions per day, and store the details permanently. Others generate transient feeds of data that must be captured in real time. The volume of changes and storage interval of the source data will determine the frequency of extraction required to support the business requirements.

What latency between changes in source systems and reporting is tolerable?

Another factor in planning data extraction timings is the requirement for the data warehouse to be kept up to date with changes in the source systems. If real-time (or near real-time) reporting must be supported, data should be extracted and loaded into the data warehouse as soon as possible after each change. Alternatively, if all reporting and analysis is historical, you might be able to leave a significant period (for example, a month) between data warehouse loads. However, you do not need to match data extractions one-to-one with data loads. If less overhead is created in the data source by a nightly extraction of the day's changes than a monthly extraction, you might choose to stage the data nightly, and then load it into the data warehouse in one operation at the end of the month.

How long does data extraction take?

Perform a test extraction and note the time taken to extract a specific number of rows. Then, based on how many new and modified rows are created in a particular period, estimate the time an extraction would take if performed hourly, daily, weekly, or at any other interval that makes sense, based on your answers to the first two questions.

During what time periods are source systems least heavily used?

Some data sources might be available only during specific periods; others might be too heavily used during business hours to support the additional overhead of an extraction process. You must work closely with the administrators and users of the data sources to identify the ideal data extraction periods for each source.

After you consider these questions for all source systems, you can start to plan extraction windows for the data. Note that it is common to have multiple sources with different extraction windows, so that the elapsed time to stage all the data might be several hours or even days.



- How frequently is new data generated in the source systems, and for how long is it retained?
- What latency between changes in source system and reporting is tolerable?
- How long does data extraction take?
- During what time periods are source systems least heavily used?

Planning Transformations

When planning an ETL solution, you must consider the transformations that need to be applied to the data to validate it and make it conform to the target table schemas.

Where to Perform Transformations

One consideration for transformations is where they should be applied during the ETL process.

Performing Transformations on Extraction

If the data sources support it, you can perform transformations in the queries used to extract data. For example, in a SQL Server data source, you can use joins, ISNULL expressions, CAST and CONVERT expressions, and concatenation expressions in the SELECT query used to extract the data. In an enterprise BI solution, this technique can be used during the following:

- Extraction from the source system.
- Extraction from the landing zone.
- Extraction from the staging area.

Performing Transformations in the Data Flow

You can use SSIS data flow transformations to transform data during the data flow. For example, you can use lookups, derived column transformations, and custom scripts to validate and modify rows in a data flow. You can also use merge and split transformations to combine or create multiple data flow paths. In an enterprise BI solution, this technique can be used during the following data flows:

- Source to landing zone.
- Landing zone to staging.
- Staging to data warehouse.

Performing Transformations In-Place

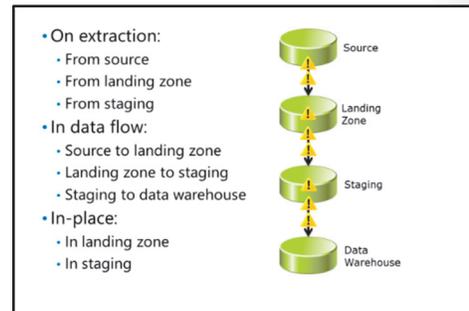
In some cases, it might make sense to transfer data from sources into one or more database tables, and then perform UPDATE operations to modify the data in-place before the next phase of the ETL data flow. For example, you might extract data from one source and stage it, and then update coded values, based on data from another source that is extracted during a later extraction window. In an enterprise BI solution, this technique can be used in the following locations:

- Landing zone tables.
- Staging tables.

Guidelines for Choosing Where to Perform Transformations

Although there is no single correct place in the data flow to perform transformations, consider the following guidelines for designing your solutions:

- Minimize the extraction workload on source systems. In this way, you can extract the data in the shortest time possible with minimal adverse effect on business processes and applications, using the data source.



- Perform validations and transformations in the data flow as soon as possible. This enables you to remove or redirect invalid rows and unnecessary columns early in the extraction process, reducing the amount of data being transferred across the network.
- Minimize the time it takes to load the data warehouse tables. This means you can get the new data into production as soon as possible and perform the load with minimal adverse effect on data warehouse users.

How to Perform Transformations

You can use Transact-SQL statements to transform or validate columns during extraction or in-place. Alternatively, you can use SSIS data flow transformations to modify the data during the data flow. The following table lists some typical validation and transformation scenarios, together with information about how to use Transact-SQL or data flow transformations to implement a solution:

Scenario	Transact-SQL	Data flow transformations
Data type conversion	Use the CAST or CONVERT function.	Use the Data Conversion transformation.
Concatenation	Concatenate fields in the SELECT clause of the query.	Use the Derived Column transformation.
Replacing NULL values	Use the ISNULL function.	Use the Derived Column transformation with an expression containing the ReplaceNull function.
Looking up related values where referential integrity is enforced	Use an INNER JOIN.	Use the Lookup transformation. Looking up related values where referential integrity is not enforced.
Looking up related values where referential integrity is NOT enforced	Use an OUTER JOIN, and optionally use ISNULL to replace null values where no matching rows exist.	Use the Lookup transformation with the Ignore failure option, and then add a transformation later in the data flow to handle null values (either by replacing them with a Derived Column or redirecting them with a Condition Split). Alternatively, use the Redirect rows to no match output option and handle the nulls before using a Merge transformation to return the fixed rows to the main data flow.

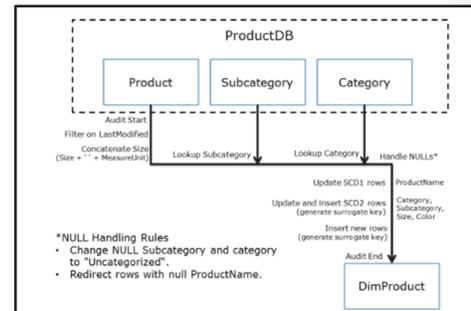
 **Note: SQL Server Integration Services (SSIS).** Some people who are unfamiliar with SSIS make the incorrect assumption that the data flow processes rows sequentially, and that data flow transformations are inherently slower than set-based transformations performed with Transact-SQL. However, the SSIS pipeline performs set-based operations on buffered batches of rows—it is designed to provide high-performance transformation in data flows.

Documenting Data Flows

An important part of designing an ETL solution is to document the data flows you need to implement. The diagrams and notes that document your data flows are commonly referred to as “source-to-target” documentation. Typically, this starts with a simple high level diagram for each table in the data warehouse.

The diagram will show:

- **Source Tables:** from which the data warehouse table fields originate.
- **Validation and Transformations:** that must be applied during the data flow.



Generally, you use a consistent diagramming approach for each table, and include as much detail about validation rules, transformations, and potential issues as you can. It is common for these high level diagrams to start simple and be refined as the ETL design evolves.

As your ETL design is refined, you will start to develop a clear idea of what fields will be extracted, generated, and validated at each stage of the data flow. To help document the lineage of the data as it flows from the source to the data warehouse tables, you can create source-to-target mappings that show detailed information for the fields at each stage.

 **Note: Spreadsheet.** A common way to create a source-to-target mapping is to use a spreadsheet divided into a set of columns for each stage in the data flow. Start with the fields in the target table, and then work backward to determine the required staging, landing zone, and source fields, along with any validation rules and transformations that must be applied. The goal is to create a single document in which the origins of a field in the target table can be traced back across a row to its source.

Like high level data flow diagrams, many BI professionals have adopted different variations of source-to-target mapping. The organization where you are working might not have a standard format for this kind of documentation. It’s important, therefore, to use a consistent format that is helpful during ETL design and easy to understand for anyone who needs to troubleshoot or maintain the ETL system in future.

Slowly Changing Dimensions

Dimensions contain data that can be relatively static. For example, in a **Customer Account** dimension table, data changes can be unpredictable in size and frequency. You would expect people to log on to a retail system at any time of the day to make large or small changes to dimension data, such as address details. Fact tables are linked to dimension tables by foreign keys and use the dimension data as a source of data, such as address details. Fact tables might change rapidly, such as when an order is placed; the orders list may be out of date if dimension data is not loaded quickly and efficiently.

When updating dimension tables, you need to apply the appropriate logic, depending on the kind of dimension attributes being modified. Changes to fixed attributes are not supported, so any updates to columns containing them must either be discarded or cause an error. However, modifications to changing and historical attributes must be supported, and a widely used set of techniques for handling these changes has been identified.

The management of loading SCD data ensures that the fact tables access the latest dimension data.



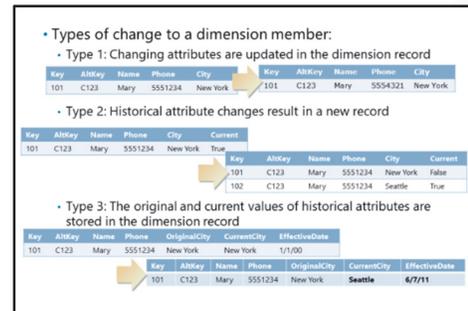
Note: Temporal or System-Versioned tables is a SQL Server feature that can be used to manage version control for SCD. This topic is dealt with later in the module.

SCDs change over time, while retaining historical attributes for reporting and analysis. Changes to dimension members are usually categorized as the following types:

- **Type 1:** Changing attributes are updated in the existing record and the previous value is lost.
- **Type 2:** Historical attribute changes result in a new record in the dimension table, representing a new version of the dimension member. A column is used to indicate which version of the dimension member is the current one, either with a flag value to indicate the current record, or by storing the date and time when each version becomes effective. You can use this technique to store a complete history of all versions of the dimension member.
- **Type 3:** Historical attribute changes are stored in the existing record, in which the original value is also retained and a column indicates the date on which the new value becomes effective. You can use this technique to store the original and latest versions of the dimension member.

Question: What are the missing words indicated by <xx xx xx> in the following statement about a data warehouse ETL flow?

“For BI solutions that involve loading large volumes of data, a <xx xx xx> process is recommended. In this data flow architecture, the data is initially extracted to tables that closely match the source system schemas (often referred to as a landing zone).”



Check Your Knowledge

Question	
Which of the following questions is NOT a consideration when planning extraction windows?	
Select the correct answer.	
<input type="checkbox"/>	How much impact on business is likely during the extraction window?
<input type="checkbox"/>	How frequently is new data generated in the source systems, and for how long is it retained?
<input type="checkbox"/>	During what time periods are source systems least heavily used?
<input type="checkbox"/>	How long does data extraction take?
<input type="checkbox"/>	What latency is tolerable between changes in source systems and reporting?

Question: Which three features does SQL Server offer for SCD control?

Lesson 2

Extracting Modified Data

An incremental ETL process starts by extracting data from source systems. To avoid including unnecessary rows of data in the extraction, the solution must be able to identify records that have been inserted or modified since the last refresh cycle, and limit the extraction to those records.

This lesson describes a number of techniques for identifying and extracting modified records.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the common options for extracting modified records.
- Implement an ETL solution that extracts modified rows based on a DateTime column.
- Configure the CDC feature in the SQL Server Enterprise Edition database engine.
- Implement an ETL solution that extracts modified rows by using the CDC feature.
- Use the CDC Control Task and data flow components to extract CDC records.
- Configure the Change Tracking feature in the Microsoft® SQL Server database engine.
- Implement an ETL solution that extracts modified rows by using Change Tracking.

Options for Extracting Modified Data

There are a number of commonly used techniques employed to extract data as part of a data warehouse refresh cycle.

Extract All Records

The simplest solution is to extract all source records and load them to a staging area, before using them to refresh the data warehouse. This technique works with all data sources and ensures that the refresh cycle includes all inserted, updated, and deleted source records. However, this technique can require the transfer and storage of large volumes of data, making it inefficient and impractical for many enterprise data warehousing solutions.

Store a Primary Key and Checksum

Another solution is to store the primary key of all previously extracted rows in a staging table, along with a checksum value that is calculated from source columns in which you want to detect changes. For each refresh cycle, your ETL process can extract source records for which the primary key is not recorded in the table of previous extractions—in addition to rows where the checksum value calculated from the columns in the source record does not match the checksum recorded during the previous extraction. Additionally, any primary keys recorded in the staging table that no longer exist in the source, represent deleted records.

- Extract all records
- Store a primary key and checksum
- Use a datetime column as a “high water mark”
- Use Change Data Capture
- Use Change Tracking

This technique limits the extracted records to those that have been inserted or modified since the previous refresh cycle. However, for large numbers of rows, the overhead of calculating a checksum to compare with each row can significantly increase processing requirements.

Use a Datetime Column As a “High Water Mark”

Tables in data sources often include a column to record the date and time of the initial creation and last modification to each record. If your data source includes such a column, you can log the date and time of each refresh cycle and compare that with the last modified value in the source data records—to identify records that have been inserted or modified. This technique is commonly referred to as using the date and time of each extraction as a “high water mark” because of its similarity to the way a tide or flood can leave an indication of the highest water level.

Use Change Data Capture

Change Data Capture (CDC) is a feature of SQL Server Enterprise Edition that uses Transaction Log Sequence Numbers (LSNs) to identify insert, update, and delete operations that have occurred within a specified time period. To use CDC, your ETL process must store the date and time or LSN of the last extraction as described for the high water mark technique. However, it is not necessary for tables in the source database to include a column that indicates the date and time of the last modification.

CDC is an appropriate technique when:

- The data source is a database in the Enterprise Edition of SQL Server 2008 or later.
- You need to extract a complete history that includes each version of a record that has been modified multiple times.

Use Change Tracking

Change Tracking (CT) is another SQL Server technology you can use to record the primary key of modified records that and extract records based on a version number that is incremented each time a row is inserted, updated, or deleted. To use CT, you must log the version that is extracted, and then compare the logged version number to the current version to identify modified records during the next extraction.

CT is an appropriate technique when:

- The data source is a SQL Server 2008 or later database.
- You need to extract the latest version of a row that has been modified since the previous extraction, but you do not need a full history of all interim versions of the record.

Considerations for Handling Deleted Records

If you need to propagate record deletions in source systems to the data warehouse, you should consider the following guidelines:

- You need to be able to identify which records have been deleted since the previous extraction. One way to accomplish this is to store the keys of all previously extracted records in the staging area and compare them to the values in the source database as part of the extraction process. Alternatively, CDC and CT both provide information about deletions, helping you to identify deleted records without maintaining the keys of previously extracted records.
- If the source database supports logical deletes by updating a Boolean column, to indicate that the record is removed, then deletions are conceptually just a special form of update. You can implement custom logic in the extraction process to treat data updates and logical deletions separately if necessary.

Slowly Changing Dimensions

SCDs are static data tables that change slowly and unpredictably. Changes may not be made to dependent fact tables. Your ETL solution must be able to optimize the transforming of SCD while maintaining historically accurate reporting. The demonstrations, lessons and labs in this module use static data of this type.



Note: SCD was described earlier in this module. How System-Versioned tables are used to assist in version control for SCD is described later in this module.

Extracting Rows Based on a Datetime Column

If your data source includes a column to indicate the date and time each record was inserted or modified, you can use the high water mark technique to extract modified records. The high level steps your ETL process must perform to use the high water mark technique are:

1. Note the current time.
2. Retrieve the date and time of the previous extraction from a log table.
3. Extract records where the modified date column is later than the last extraction time, but before or equal to the current time you noted in step 1. This disregards any insert or update operations that have occurred since the start of the extraction process.
4. In the log, update the last extraction date and time with the time you noted in step 1.

1. Note the current time
2. Retrieve the last extraction time from an extraction log
3. Extract and transfer records that were modified between the last extraction and the current time
4. Replace the stored last extraction value with the current time

Demonstration: Using a Datetime Column

In this demonstration, you will see how to use a datetime column to extract modified data.

Demonstration Steps

Use a Datetime Column to Extract Modified Data

1. Ensure 20767C-MIA-DC and 20767C-MIA-SQL are started, and log onto 20767C-MIA-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod09 folder**, run **Setup.cmd** as Administrator.
3. In the **User Account Control** dialog box, click **Yes**.
4. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the database engine using Windows authentication.
5. In Object Explorer, expand **Databases**, expand **DemoDW**, and then expand **Tables**. Note that the database includes tables in three schemas (**dw**, **src**, and **stg**) to represent the data sources, staging database, and data warehouse in an ETL solution.

6. Right-click each of the following tables and click **Select Top 1000 Rows**:
 - **stg.Products**: this table is used for staging product records during the ETL process, and is currently empty.
 - **stg.ExtractLog**: this table logs the last extraction date for each source system.
 - **src.Products**; this table contains the source data for products, including a **LastModified** column that records when each row was last modified.
7. Start Visual Studio® and open the **IncrementalETL.sln** solution in the **D:\Demofiles\Mod09** folder.
8. In Solution Explorer, double-click the **Extract Products.dtsx** SSIS package.
9. On the **SSIS** menu, click **Variables**, and note that the package contains two user variables named **CurrentTime** and **LastExtractTime**. Close the window.
10. On the control flow surface, double-click **Get Current Time**.
11. In the **Expression Builder** dialog box, note that the Expression pane of this task sets the **CurrentTime** user variable to the current date and time. Click **Cancel**.
12. Double-click **Get Last Extract Time**, note the following configuration settings, and then click **Cancel**:
 - On the **General** tab, the **ResultSet** property is set to return a single row, and the **SQLStatement** property contains a query to retrieve the maximum **LastExtractTime** value for the products source in the **stg.ExtractLog** table.
 - On the **Result Set** tab, the **LastExtractTime** value in the query results row is mapped to the **User::LastExtractTime** user variable.
13. Double-click **Stage Products** to view the data flow surface.
14. Double-click **Products Source** and note that the SQL command used to extract products data includes a WHERE clause that filters the query results.
15. Click **Parameters**, and note that the parameters in the Transact-SQL query are mapped to the **LastExtractTime** and **CurrentTime** variables.
16. Click **Cancel** in all dialog boxes, and then click the **Control Flow** tab.
17. On the control flow surface, double-click **Update Last Extract Time**, note the following configuration settings, and then click **Cancel**:
 - On the **General** tab, the **SQLStatement** property contains a Transact-SQL UPDATE statement that updates the **LastExtractTime** in the **stg.ExtractLog** table.
 - On the **Parameter Mapping** tab, the **CurrentTime** user variable is mapped to the parameter in the Transact-SQL statement.
18. In SQL Server Management Studio, open **Modify Products.sql** in the **D:\Demofiles\Mod09** folder.
19. Execute the script to modify some rows in the **src.Products** table.
20. In Visual Studio, on the **Debug** menu, click **Start Debugging**. When package execution is complete, on the **Debug** menu, click **Stop Debugging**.
21. In SQL Server Management Studio, right-click each of the following tables and click **Select Top 1000 Rows**:
 - **stg.ExtractLog**: note that the **LastExtractTime** for the **Products** data source has been updated.
 - **stg.Products**: note the rows that have been extracted from the **src.Products** table.
22. Minimize SQL Server Management Studio and Visual Studio.

Change Data Capture

The CDC feature in SQL Server Enterprise Edition provides a number of functions and stored procedures that you can use to identify modified rows. To use CDC, perform the following high level steps:

1. Enable CDC in the data source. You must enable CDC for the database and for each table in the database where you want to monitor changes. The following Transact-SQL code sample shows how to:
 - a. Check CDC is enabled on a database using the **is_cdc_enabled** system stored procedure.
 - b. Enable CDC on a database using the **sp_cdc_enable_db** system stored procedure.
 - c. Enable CDC on a table using the **sp_cdc_enable_table** system stored procedure.

```
IF (SELECT is_cdc_enabled FROM sys.databases WHERE name='Customers') = 'FALSE'
    BEGIN
        EXEC sys.sp_cdc_enable_db
        EXEC sys.sp_cdc_enable_table
            @source_schema      = N'dbo',
            @source_name        = N'Customers',
            @role_name          = NULL,
            @supports_net_changes = 1
    END
GO
```

2. In the ETL process used to extract the data, map start and end times (based on the logged date and time of the previous extraction and the current date and time) to log sequence numbers. The following Transact-SQL code sample shows how to use the **fn_cdc_map_time_to_lsn** system function to map Transact-SQL variables named **@StartDate** and **@EndDate** to log sequence numbers:

```
DECLARE @from_lsn binary(10), @to_lsn binary(10);
SET @from_lsn = sys.fn_cdc_map_time_to_lsn('smallest greater than', @StartDate)
SET @to_lsn = sys.fn_cdc_map_time_to_lsn('largest less than or equal', @EndDate)
```

3. Include logic to handle errors if either of the log sequence numbers is null. This can happen if no changes have occurred in the database during the specified time period. The following Transact-SQL code sample shows how to check for null log sequence numbers:

```
IF (@from_lsn IS NULL) OR (@to_lsn IS NULL)
    -- There may have been no transactions in the timeframe.
```

4. Extract records that have been modified between the log sequence numbers. When you enable CDC for a table, SQL Server generates table-specific system functions that you can use to extract data modifications to that table. The following Transact-SQL code sample shows how to use the **fn_cdc_get_net_changes_dbo_Customers** system function to retrieve rows that have been modified in the **dbo.Customers** table:

```
SELECT * FROM cdc.fn_cdc_get_net_changes_dbo_Customers(@from_lsn, @to_lsn, 'all')
```

1. Enable Change Data Capture:


```
EXEC sys.sp_cdc_enable_db
EXEC sys.sp_cdc_enable_table @source_schema = N'dbo', @source_name = N'Customers',
                             @role_name = NULL, @supports_net_changes = 1
```
2. Map start and end times to log sequence numbers:


```
DECLARE @from_lsn binary(10), @to_lsn binary(10);
SET @from_lsn = sys.fn_cdc_map_time_to_lsn('smallest greater than', @StartDate)
SET @to_lsn = sys.fn_cdc_map_time_to_lsn('largest less than or equal', @EndDate)
```
3. Handle null log sequence numbers:


```
IF (@from_lsn IS NULL) OR (@to_lsn IS NULL)
    -- There may have been no transactions in the timeframe
```
4. Extract changes between log sequence numbers:


```
SELECT * FROM cdc.fn_cdc_get_net_changes_dbo_Customers(@from_lsn, @to_lsn, 'all')
```

 **About Change Data Capture (SQL Server)**

<http://aka.ms/x8xtv2>

 **Change Data Capture Functions (Transact-SQL)**

<http://aka.ms/r4i219>

Demonstration: Using Change Data Capture

In this demonstration, you will see how to:

- Enable Change Data Capture.
- Use Change Data Capture to Extract Modified Data.

Demonstration Steps

Enable CDC on the Customers Table

1. Ensure you have completed the previous demonstration in this module.
2. Return to SQL Server Management Studio, and in Object Explorer, in the **DemoDW** database, right-click the **src.Customers** table and click **Select Top 1000 Rows**. This table contains source data for customers.
3. Open **Using CDC.sql** in the D:\Demofiles\Mod09 folder.
4. In the code window, select the Transact-SQL code under the comment **Enable CDC on src.Customers table**, and then click **Execute**. This enables CDC in the **DemoDW** database, and starts logging modifications to data in the **src.Customers** table.
5. Select the Transact-SQL code under the comment **Select all changed customer records since the last extraction**, and then click **Execute**. This code uses CDC functions to map dates to log sequence numbers, and retrieve records in the **src.Customers** table that have been modified between the last logged extraction in the **stg.ExtractLog** table, and the current time. There are no changed records because no modifications have been made since CDC was enabled.

Use CDC to Extract Modified Data

1. Select the Transact-SQL code under the comment **Insert a new customer**, and then click **Execute**. This code inserts a new customer record.
2. Select the Transact-SQL code under the comment **Make a change to a customer**, and then click **Execute**. This code updates a customer record.
3. Select the Transact-SQL code under the comment **Now see the net changes**, and then click **Execute**. This code uses CDC functions to map dates to log sequence numbers, and retrieve records in the **src.Customers** table that have been modified between the last logged extraction in the **stg.ExtractLog** table, and the current time. Two records are returned.
4. Wait 10 seconds then select the Transact-SQL code under the comment Check for changes in an interval with no database activity, and then click **Execute**. Because there has been no activity in the database during the specified time interval, the system returns There has been no logged database activity in the specified time interval. This demonstrates the importance of checking for a null log sequence number value when using CDC.
5. Minimize SQL Server Management Studio.

Extracting Data with Change Data Capture

To extract data from a CDC-enabled table in an SSIS-based ETL solution, you can create a custom control flow that uses the same principles as the high water mark technique described earlier. The general approach is to establish the range of records to be extracted, based on a minimum and maximum Log Sequence Number (LSN), extract those records, and then log the endpoint of the extracted range to be used as the starting point for the next extraction.

You can choose to log the high water mark as an LSN or a datetime value that can be mapped to an LSN by using the `fn_cdc_map_time_to_lsn` system function.

The following procedure describes one way to create an SSIS control flow for extracting CDC data:

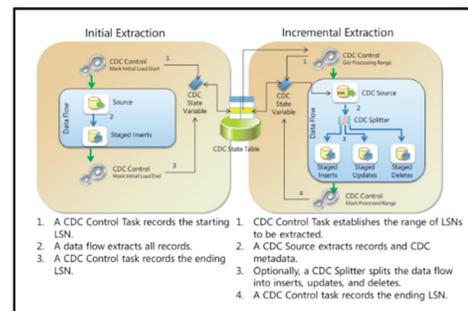
1. Use an Expression task to assign the current time to a datetime variable.
2. Use a SQL command task to retrieve the logged datetime value that was recorded after the previous extraction.
3. Use a data flow task in which a source employs the `fn_cdc_map_time_to_lsn` system function to map the current time and previously extracted time to the corresponding LSNs, and then uses the `cdc.fn_cdc_get_net_changes_capture_instance` function to extract the data that was modified between those LSNs. You can use the `$_operation` column in the resulting dataset to split the records into different data flow paths for inserts, updates, and deletes.
4. Use a SQL command task to update the logged datetime value to be used as the starting point for the next extraction.

1. Identify the endpoint for the extraction (LSN or DateTime)
2. Retrieve the last extraction endpoint from an extraction log
3. Extract and transfer records that were modified during the LSN range defined by the previous extraction endpoint and the current endpoint
4. Replace the logged endpoint value with the current endpoint

The CDC Control Task and Data Flow Components

To make it easier to implement packages that extract data from CDC-enabled sources, SSIS includes CDC components that abstract the underlying CDC functionality. The CDC components included in SSIS are:

- **CDC Control Task:** a control flow task that you can use to manage CDC state, providing a straightforward way to track CDC data extraction status.
- **CDC Source:** a data flow source that uses the CDC state logged by the CDC Control Task to extract a range of modified records from a CDC-enabled data source.
- **CDC Splitter:** A data flow transformation that splits output rows from a CDC Source into separate data flow paths for inserts, updates, and deletes.



All the CDC components in SSIS require the use of ADO.NET connection managers to the CDC-enabled data source and the database where the CDC state is to be stored.

Performing an Initial Extraction with the CDC Control Task

When using the CDC Control Task to manage extractions from a CDC-enabled data source, it is recommended practice to create a package that will be executed once to perform the initial extraction. This package should contain the following control flow:

1. **A First CDC Control Task:** configured to perform the **Mark initial load start** operation. This writes an encoded value, including the starting LSN to a package variable, and optionally persists it to a state tracking table in a database.
2. **A Data Flow:** this extracts all rows from the source and loads them into a destination—typically a staging table. This data flow does not require CDC-specific components.
3. **A Second CDC Control Task:** configured to perform the **Mark initial load end** operation. This writes an encoded value, including the ending LSN to a package variable, and optionally persists it to a state tracking table in a database.

Performing Incremental Extractions with the CDC Control Task

After the initial extraction has been performed, subsequent extractions should use an SSIS package with the following control flow:

1. **A First CDC Control Task:** configured to perform the **Get processing range** operation. This establishes the range of records to be extracted and writes an encoded representation to a package variable, which can also be persisted to a state tracking table in a database.
2. **A Data Flow:** that uses a CDC source, using the encoded value in the CDC state package variable to extract modified rows from the data source.
3. **Data Flow CDC Splitter Task (Optional):** if required, the data flow can include a CDC Splitter task, which uses the **\$_operation** column in the extracted rowset to redirect inserts, updates, and deletes to separate data flow paths. These can then be connected to appropriate destinations for staging tables.
4. **A Second CDC Control Task:** configured to perform the **Mark processed range** operation. This writes an encoded value, including the ending LSN to a package variable, and optionally persists it to a state tracking table in a database. This value is then used to establish the starting point for the next extraction.

Demonstration: Using CDC Components

In this demonstration, you will see how to use the CDC Control Task to:

- Perform an Initial Extraction.
- Extract Changes.

Demonstration Steps

Perform an Initial Extraction

1. Ensure you have completed the previous demonstrations in this module.
2. Return to SQL Server Management Studio.
3. Open the **CDC Components.sql** script file in the **D:\Demofiles\Mod09** folder. Note that the script enables CDC for the **src.Shippers** table, and then click **Execute**.

4. Open the **Fix CDC.sql** script file in the **D:\Demofiles\Mod09** folder. Note that this script corrects an issue that can arise with CDC and that will be fixed in a future release. Click **Execute** to run the script.
5. In Object Explorer, right-click each of the following tables in the **DemoDW** database, and click **Select Top 1000 Rows** to view their contents:
 - **src.Shippers**: this table should contain four records.
 - **stg.ShipperDeletes**: this table should be empty.
 - **stg.ShipperInserts**: this table should be empty.
 - **stg.ShipperUpdates**: this table should be empty.
6. Return to Visual Studio, in which the **IncrementalETL.sln** solution should be open, and in Solution Explorer, double-click the **Extract Initial Shippers.dtsx** SSIS package. Note that the CDC Control Tasks in the control flow contain errors, which you will resolve.
7. Double-click the **Mark Initial Load Start** CDC Control Task.
8. In the editor, set the following properties, and then click **OK**:
 - **SQL Server CDC database ADO.NET connection manager**: localhost DemoDW ADO NET.
 - **CDC control operation**: Mark initial load start.
 - **Variable containing the CDC state**: click **New** and create a new variable named **CDC_State**.
 - **Automatically store state in a database table**: selected.
 - **Connection manager for the database where the state is stored**: localhost DemoDW ADO NET.
 - **Table to use for storing state**: click **New**, and then click **Run** to create the **cdc_states** table.
 - **State name**: CDC_State.
9. Double-click the **Extract All Shippers** data flow task.
10. On the Data Flow surface, note that an ADO.NET source named **Shippers** is used to extract all rows from the **src.Shippers** table, and an ADO.NET destination named **Shipper Inserts** is used to load the extracted rows into the **stg.ShipperInserts** table.
11. On the **Control Flow** tab, double-click the **Mark Initial Load End** CDC Control Task.
12. Set the following properties, and then click **OK**:
 - **SQL Server CDC database ADO.NET connection manager**: localhost DemoDW ADO NET.
 - **CDC control operation**: Mark initial load end (make sure you do not select Mark initial load start).
 - **Variable containing the CDC state**: Select **User::CDC_State** from the list (this is the variable you created earlier).
 - **Automatically store state in a database table**: selected.
 - **Connection manager for the database where the state is stored**: localhost DemoDW ADO NET.
 - **Table to use for storing state**: [dbo].[cdc_states] (the table you created earlier).
 - **State name**: CDC_State.
13. On the **Debug** menu, click **Start Debugging** and wait for the package execution to complete.

14. On the **Debug** menu, click **Stop Debugging**.
15. In SQL Server Management Studio, right-click the **Tables** folder for the **DemoDW** database and click **Refresh**. Note that a table named **dbo.cdc_states** has been created.
16. Right-click **dbo.cdc_states** and click **Select Top 1000 Rows** to view the logged **CDC_State** value (which should begin "ILEND...").
17. Right-click **stg.ShipperInserts** and click **Select Top 1000 Rows** to verify that the initial set of shippers has been extracted.

Extract Changes

1. In SQL Server Management Studio, open the file **Update Shippers.sql** in the D:\Demofiles\Mod09 folder.
2. Select the code under the comment **Cleanup previous extraction**, noting that it truncates the **stg.ShipperInserts** table, and click **Execute**.
3. In Visual Studio, in Solution Explorer, double-click **Extract Changed Shippers.dtsx**.
4. On the **Control Flow** tab, double-click the **Get Processing Range** CDC Control Task. Note this task gets the processing range, storing it in the **CDC_State** variable and the **cdc_states** table. Click **Cancel**.
5. Double-click the **Extract Modified Shippers** data flow task.
6. On the **Data Flow** tab, view the properties of the **Shipper CDC Records** CDC Source component, noting that it extracts modified records, based on the range stored in the **CDC_State** variable. Click **Cancel**.
7. Note that the **CDC Splitter** transformation has three outputs, one each for inserts, updates, and deletes. Each of these outputs is connected to an ADO.NET destination that loads the records into the **stg.ShipperInserts**, **stg.ShipperUpdates**, and **stg.ShipperDeletes** tables respectively.
8. On the **Control Flow** tab, double-click the **Mark Processed Range** CDC Control Task. Note it updates the **CDC_State** variable and the **cdc_states** table when the extraction is complete. Click **Cancel**.
9. On the **Build** menu, click **Rebuild Solution**.
10. On the **Debug** menu, click **Start Debugging**. When execution is complete, stop debugging.
11. In SQL Server Management Studio, in Object Explorer, right-click each of the following tables in the **DemoDW** database, and click **Select Top 1000 Rows** to view their contents:
 - **stg.ShipperDeletes**: this table should still be empty.
 - **stg.ShipperInserts**: this table should still be empty.
 - **stg.ShipperUpdates**: this table should still be empty.

No rows are transferred, because the source data is unchanged since the initial extraction.

12. In the **Update Shippers.sql** script file, select the code under the comment **Modify source data**, noting that it performs an INSERT, an UPDATE, and a DELETE operation on the **src.Shippers** table, and click **Execute**.
13. In Visual Studio, click the **Control Flow** tab, and then on the **Debug** menu, click **Start Debugging**.
14. When execution is complete, view the **Extract Modified Shippers** data flow task and note the number of rows transferred (if no rows were transferred, stop debugging and re-run the package). When three rows have been transferred (one to each output of the CDC Splitter transformation), stop debugging and close Visual Studio.

15. In SQL Server Management Studio, right-click each of the following tables and click **Select Top 1000 Rows** to view their contents. Each table should contain a single row:
 - stg.ShipperDeletes
 - stg.ShipperInserts
 - stg.ShipperUpdates
16. Minimize SQL Server Management Studio.

Change Tracking

The CT feature in SQL Server provides a number of functions and stored procedures that you can use to identify modified rows. To use CT, perform the following high level steps:

1. Enable CT in the data source. You must enable CT for the database, and for each table in the database for which you want to monitor changes. The following Transact-SQL code sample shows how to enable CT in a database named **Sales** and monitor data modifications in the **Salespeople** table. Note that you can choose to track which columns were modified, but the change table only contains the primary key of each modified row—not the modified column values:

```
ALTER DATABASE Sales
SET CHANGE_TRACKING = ON (CHANGE_RETENTION = 7 DAYS, AUTO_CLEANUP = ON)

ALTER TABLE Salespeople
ENABLE CHANGE_TRACKING WITH (TRACK_COLUMNS_UPDATED = OFF)
```

2. For the initial data extraction, record the current version (which by default is 0), and extract all rows in the source table, then log the current version as the last extracted version. The following Transact-SQL code sample shows how to use the **Change_Tracking_Current_Version** system function to retrieve the current version, extract the initial data, and assign the current version to a variable so it can be stored as the last extracted version:

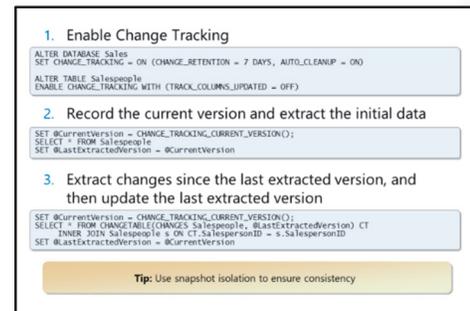
```
SET @CurrentVersion = CHANGE_TRACKING_CURRENT_VERSION();

SELECT * FROM Salespeople

SET @LastExtractedVersion = @CurrentVersion
```

3. For subsequent refresh cycles, extract changes that have occurred between the last extracted version and the current one. The following Transact-SQL code sample shows how to determine the current version; use the **Changetable** system function in a query that joins the primary key of records in the change table to records in the source table; and update the last extracted version:

```
SET @CurrentVersion = CHANGE_TRACKING_CURRENT_VERSION();
SELECT * FROM CHANGETABLE(CHANGES Salespeople, @LastExtractedVersion) CT
INNER JOIN Salespeople s ON CT.SalespersonID = s.SalespersonID
SET @LastExtractedVersion = @CurrentVersion
```



4. When using Change Tracking, a best practice is to enable snapshot isolation in the source database and use it to ensure that any modifications occurring during the extraction do not affect records that were modified between the version numbers that define the lower and upper bounds of your extraction range.



About Change Tracking (SQL Server)

<http://aka.ms/i10jx3>



Change Tracking Functions (Transact-SQL)

<http://aka.ms/ul08f4>

Demonstration: Using Change Tracking

In this demonstration, you will see how to:

- Enable Change Tracking.
- Use Change Tracking.

Demonstration Steps

Enable Change Tracking

1. Ensure you have completed the previous demonstrations in this module.
2. Return to SQL Server Management Studio.
3. In Object Explorer, in the **DemoDW** database, right-click the **src.Salespeople** table and click **Select Top 1000 Rows**. This table contains source data for sales employees.
4. Open the file **Using CT.sql** in the **D:\Demofiles\Mod09** folder.
5. Select the Transact-SQL code under the comment **Enable Change Tracking**, and click **Execute**. This code enables CT in the **DemoDW** database, and starts logging changes to data in the **src.Salespeople** table.
6. Select the Transact-SQL code under the comment **Obtain the initial data and log the current version number**, and then click **Execute**. This code uses the **CHANGE_TRACKING_CURRENT_VERSION** function to determine the current version, and retrieves all records in the **src.Salespeople** table.
7. Select the Transact-SQL code under the comment **Insert a new salesperson**, and then click **Execute**.
8. Select the Transact-SQL code under the comment **Update a salesperson**, and then click **Execute**.
9. Select the Transact-SQL code under the comment **Retrieve the changes between the last extracted and current versions**, and then click **Execute**.
10. In Object Explorer, in the **DemoDW** database, right-click the **src.Salespeople** table and click **Select Top 1000 Rows**.
11. Close SQL Server Management Studio without saving changes.

Extracting Data with Change Tracking

You can create an SSIS package that uses the CT feature in SQL Server in a similar way to the high water mark technique described earlier in this lesson. The key difference is that, rather than storing the date and time of the previous extraction, you must store the CT version number that was extracted, and update this with the current version during each extract operation.

A typical control flow for extracting data from a CT-enabled data source includes the following elements:

1. A SQL command that retrieves the previously extracted version from a log table and assigns it to a variable.
2. A data flow that contains a source to extract records that have been modified since the previously extracted version, and then return the current version.
3. A SQL command that updates the logged version number with the current version.

Question: What are the four steps used when extracting rows based on a datetime column?

1. Retrieve the last version number that was extracted from an extraction log
2. Extract and transfer records that were modified since the last version, retrieving the current version number
3. Replace the logged version number with the current version number

Check Your Knowledge

Question	
<p>In the following code to enable CDC in a table, two system stored procedures have been omitted (<XXXXXX> and <YYYYY>). Choose the correct stored procedures from the following options:</p> <pre>IF (SELECT is_cdc_enabled FROM sys.databases WHERE name='Customers') = 'FALSE' BEGIN EXEC <xxxxxx> EXEC <yyyyyy> @source_schema = N'dbo', @source_name = N'Customers', @role_name = NULL, @supports_net_changes = 1 END GO</pre>	
Select the correct answer.	
<input type="checkbox"/>	<pre><XXXXXX> = sys.sp_cdc_enable_table <YYYYYY> = sys.sp_cdc_enable_db</pre>
<input type="checkbox"/>	<pre><XXXXXX> = sys.sp_cdc_enable_db <YYYYYY> = sys.sp_cdc_create_table</pre>
<input type="checkbox"/>	<pre><XXXXXX> = sys.sp_cdc_enable_db <YYYYYY> = sys.sp_cdc_enable_table</pre>
<input type="checkbox"/>	<pre><XXXXXX> = sys.sp_cdc_enable_db <YYYYYY> = sys.sp_cdc_mark_table</pre>
<input type="checkbox"/>	<pre><XXXXXX> = sys.sp_cdc_start_db <YYYYYY> = sys.sp_cdc_add_table</pre>

Question: What three steps would you expect to find in a data flow for extracting data from a CT-enabled data source?

Lab A: Extracting Modified Data

Scenario

You have developed SSIS packages that extract data from various data sources and load it into a staging database. However, the current solution extracts all source records each time the ETL process is run. This results in unnecessary processing of records that have already been extracted and consumes a sizeable amount of network bandwidth to transfer a large volume of data. To resolve this problem, you must modify the SSIS packages to extract only data that has been added or modified since the previous extraction.

Objectives

After completing this lab, you will be able to:

- Use a datetime column to extract modified rows.
- Use Change Data Capture to extract modified rows.
- Use the CDC Control Task to extract modified rows.
- Use Change Tracking to extract modified rows.

Estimated Time: 60 minutes.

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Using a Datetime Column to Incrementally Extract Data

Scenario

The **InternetSales** and **ResellerSales** databases contain source data for your data warehouse. The sales order records in these databases include a **LastModified** date column that is updated with the current date and time when a row is inserted or updated. You have decided to use this column to implement an incremental extraction solution that compares record modification times to a logged extraction date and time in the staging database. This restricts data extractions to rows that have been modified since the previous refresh cycle.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. View Extraction Data
3. Examine an Incremental Data Extraction
4. Define Variables for Execution Times
5. Modify a Data Source to Filter Data
6. Add a Task to Update the Extraction Log
7. Test the Package

► Task 1: Prepare the Lab Environment

1. Ensure that the 20767C-MIA-DC and 20767C-MIA-SQL virtual machines are both running, and then log on to 20767C-MIA-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab09\Starter** folder as **Administrator**.

► Task 2: View Extraction Data

1. Start SQL Server Management Studio and connect to the **MIA-SQL** database instance by using Windows authentication.
2. In the **Staging** database, view the contents of the **dbo.ExtractLog** table, noting that it contains the date and time of previous extractions from the **InternetSales** and **ResellerSales** databases. This is initially set to the date and time "1900-01-01 00:00:00.000".
3. In the **InternetSales** database, view the contents of the **dbo.SalesOrderHeader** table and note that the **LastModified** column contains the date and time when each record was inserted or modified.

► Task 3: Examine an Incremental Data Extraction

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex1** folder. Open the **Extract Reseller Data.dtsx** SSIS package.
2. View the variables defined in the package, and note that they include two **DateTime** variables named **CurrentTime** and **ResellerSalesLastExtract**.
3. Examine the tasks in the package and note the following:
 - The **Get Current Time** task uses the GETDATE() function to assign the current date and time to the **CurrentTime** variable.
 - The **Get Last Extract Time** task uses a Transact-SQL command to return a single row containing the **LastExtract** value for the **ResellerSales** data source from the **dbo.ExtractLog** table in the **Staging** database, and assigns the **LastExtract** value to the **ResellerSalesLastExtract** variable.
 - The **Extract Reseller Sales** data flow task includes a data source named **Reseller Sales** that uses a WHERE clause to extract records with a **LastModified** value between two parameterized values. The parameters are mapped to the **ResellerSalesLastExtract** and **CurrentTime** variables.
 - The **Update Last Extract Time** task updates the **LastExtract** column for the **ResellerSales** data source in the **dbo.ExtractLog** table with the **CurrentTime** variable.
4. Run the **Extract Reseller Data.dtsx** package. When it has completed, stop debugging, and in SQL Server Management Studio, verify that the **ExtractLog** table in the **Staging** database has been updated to reflect the most recent extraction from the **ResellerSales** source.
5. Leave SQL Server Management Studio open.

► Task 4: Define Variables for Execution Times

1. In Visual Studio, open the **Extract Internet Sales Data.dtsx** SSIS package, and add **DateTime** variables named **CurrentTime** and **InternetSalesLastExtract**.
2. Add two **DateTime** variables called **CurrentTime** and **InternetSalesLastExtract**.
3. Add an **Expression** task named **Get Current Time** to the **Extract Customer Sales Data** sequence in the control flow of the package, and configure it to apply the following expression:

```
@[User::CurrentTime] = GETDATE()
```

4. Add an **Execute SQL** task named **Get Last Extract Time** to the **Extract Customer Sales Data** sequence in the control flow of the package, and set the following configuration properties:
5. On the **General** tab, set the **Connection** property to **localhost.Staging**.

- On the **General** tab, set the **SQLStatement** to the following Transact-SQL query:

```
SELECT MAX>LastExtract) LastExtract
FROM ExtractLog
WHERE DataSource = 'InternetSales'
```

- On the **General** tab, set the **ResultSet** property to **Single** row.
- On the **Result Set** tab, add a result that maps the **LastExtract** column in the result **User::InternetSalesLastExtract** variable.
- Connect the precedence constraints of the new tasks so that the **Get Current Time** task runs first, followed by the **Get Last Extract Time** task, and then the **Extract Customers** task.

► Task 5: Modify a Data Source to Filter Data

- On the data flow tab for the **Extract Internet Sales** data flow task, make the following changes to the **Internet Sales** source:

- Add the following WHERE clause to the query in the SQL command property:

```
WHERE LastModified > ?
AND LastModified <= ?
```

- Specify the following input parameter mappings:

- **Parameter0:**
User::InternetSalesLastExtract
- **Parameter1:**
User:CurrentTime

► Task 6: Add a Task to Update the Extraction Log

- On the control flow tab, add an **Execute SQL** task named **Update Last Extract Time** to the **Extract Customer Sales Data** sequence in the control flow of the **Extract Internet Sales Data.dtsx** package, and set the following configuration properties.
- On the **General** tab, set the **Connection** property to **localhost.Staging**.
- On the **General** tab, set the **SQLStatement** property to the following Transact-SQL query:

```
UPDATE ExtractLog
SET LastExtract = ?
WHERE DataSource = 'InternetSales'
```

- On the **Parameter Mapping** tab, add the following parameter mapping:
 - **Variable Name:**
User:CurrentTime
 - **Direction:**
Input
 - **Data Type:**
DATE
 - **Parameter Name:**
0
 - **Parameter Size:**
-1

5. Connect the precedence constraint of the **Extract Internet Sales** task to the **Update Last Extract Time** task.

► Task 7: Test the Package

1. View the **Extract Internet Sales** data flow and then start debugging the package. Note the number of rows transferred. When package execution is complete, stop debugging.
2. In SQL Server Management Studio, view the contents of the **dbo.ExtractLog** table in the **Staging** database and verify that the **LastExtract** column for the InternetSales data source has been updated.
3. View the contents of the **dbo.InternetSales** table and note the rows that have been extracted.
4. In Visual Studio, debug the package again and verify that this time no rows are transferred in the **Extract Internet Sales** data flow (because no data has been modified since the previous extraction). When package execution is complete, stop debugging.
5. Close Visual Studio.

Results: After this exercise, you should have an SSIS package that uses the high water mark technique to extract only records that have been modified since the previous extraction.

Exercise 2: Using Change Data Capture

Scenario

The Internet Sales database contains a **Customers** table that does not include a column to indicate when records were inserted or modified. You plan to use the CDC feature of SQL Server Enterprise Edition to identify records that have changed between data warehouse refresh cycles, and restrict data extractions to include only modified rows.

The main tasks for this exercise are as follows:

1. Enable Change Data Capture
2. Create a Stored Procedure to Retrieve Modified Rows
3. Use the Stored Procedure in a Data Flow
4. Test the Package

► Task 1: Enable Change Data Capture

1. In SQL Server Management Studio, execute Transact-SQL statements to enable CDC in the **InternetSales** database, and monitor net changes in the **Customers** table. You can use the **Enable CDC.sql** file in the **D:\Labfiles\Lab09\Starter\Ex2** folder to accomplish this.
2. Open the **Test CDC.sql** file in the **D:\Labfiles\Lab09\Starter\Ex2** folder and examine it.
3. Execute the code under the comment **Select all changed customer records between 1/1/1900 and today** and note that no rows are returned because no changes have occurred since CDC was enabled.
4. Execute the code under the comment **Make a change to all customers (to create CDC records)** to modify data in the **Customers** table.
5. Execute the code under the comment **Now see the net changes**, and note that all customer records are returned because they have all been modified within the specified time period while CDC was enabled.

► Task 2: Create a Stored Procedure to Retrieve Modified Rows

1. In SQL Server Management Studio, execute a Transact-SQL statement that creates a stored procedure named **GetChangedCustomers** in the **InternetSales** database. The stored procedure should perform the following tasks. You can execute the **Create SP.sql** file in the D:\Labfiles\Lab09\Starter\Ex2 folder to accomplish this:
 - Retrieve the log sequence numbers for the dates specified in **StartDate** and **EndDate** parameters.
 - If neither of the log sequence numbers is null, return all records that have changed in the **Customers** table.
 - If either of the log sequence numbers is null, return an empty rowset.
2. Test the stored procedure by running the following query:

```
USE InternetSales
GO
EXEC GetChangedCustomers '1/1/1900', '1/1/2099';
GO
```

► Task 3: Use the Stored Procedure in a Data Flow

1. Reset the staging database by running the Transact-SQL code in the **Reset Staging.sql** file in the **D:\Labfiles\Lab09\Starter\Ex2** folder.
2. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex2** folder. Open the **Extract Internet Sales Data.dtsx** SSIS package.
3. On the **Data Flow** tab for the **Extract Customers** task, modify the **Customers** source to execute the following Transact-SQL command, and map the **@StartDate** and **@EndDate** input parameters to the **User::InternetSalesLastExtract** and **User::CurrentTime** variables:

```
EXEC GetChangedCustomers ?, ?
```

► Task 4: Test the Package

1. View the **Extract Customers** data flow, and then start debugging the package and note the number of rows transferred. When package execution is complete, stop debugging.
2. In SQL Server Management Studio, view the contents of the **dbo.ExtractLog** table in the **Staging** database and verify that the **LastExtract** column for the **InternetSales** data source has been updated.
3. View the contents of the **dbo.Customers** table in the **Staging** database and note the rows that have been extracted.
4. Debug the package again and verify that no rows are transferred in the **Extract Customers** data flow this time. When package execution is complete, stop debugging and close Visual Studio.

Results: After this exercise, you should have a database in which CDC is enabled, and an SSIS package that uses a stored procedure to extract modified rows, based on changes monitored by CDC.

Exercise 3: Using the CDC Control Task

Scenario

The **HumanResources** database contains an **Employee** table in which employee data is stored. You plan to use the CDC feature of SQL Server Enterprise Edition to identify modified rows in this table. You will also use the CDC Control Task in SSIS to manage the extractions from this table by creating a package to perform the initial extraction of all rows, and a second package that uses the CDC data flow components to extract rows modified since the previous extraction.

The main tasks for this exercise are as follows:

1. Enable Change Data Capture
2. View Staging Tables
3. Create Connection Managers for CDC Components
4. Create a Package for Initial Data Extraction
5. Test Initial Extraction
6. Create a Package for Incremental Data Extraction
7. Test Incremental Extraction

► Task 1: Enable Change Data Capture

1. Enable CDC in the **HumanResources** database, and monitor net changes in the **Employee** table. You can use the **Enable CDC.sql** file in the D:\Labfiles\Lab09\Starter\Ex3 folder to accomplish this.
2. Open the **Fix CDC.sql** file in the D:\Labfiles\Lab09\Starter\Ex3 folder, and then click **Execute** to run this file. This script corrects an incompatibility between the CDC functions created when enabling CDC in the database and the CDC SSIS controls that you will use in Visual Studio.

► Task 2: View Staging Tables

- Verify that the following tables in the **Staging** database are empty:
 - dbo.EmployeeDeletes
 - dbo.EmployeeInserts
 - dbo.EmployeeUpdates

► Task 3: Create Connection Managers for CDC Components

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex3** folder.
2. Note that the project already contains an **OLE DB** connection manager for the **Staging** database, but the CDC components in SSIS require an **ADO.NET** connection manager. You will therefore need to create **ADO.NET** connection managers for the **HumanResources** and **Staging** databases.
3. Create a new, project-level connection manager that produces an ADO.NET connection to the **HumanResources** database on the localhost server by using Windows authentication. After the connection manager has been created, rename it to **localhost.HumanResources.ADO.NET.conmgr**.
4. Create another new, project-level connection manager that creates an ADO.NET connection to the **Staging** database on the **localhost** server by using Windows authentication. After the connection manager has been created, rename it to **localhost.Staging.ADO.NET.conmgr**.

- **Task 4: Create a Package for Initial Data Extraction**
1. Add a new SSIS package named **Extract Initial Employee Data.dtsx** to the project.
 2. Add a CDC Control Task to the control flow, and rename the task to **Mark Initial Load Start**. Configure the **Mark Initial Load Start** task as follows:
 - **SQL Server CDC database ADO.NET connection manager:** **localhost HumanResources ADO NET.**
 - **CDC control operation:** **Mark initial load start.**
 - **Variable containing the CDC state:** a new variable named **CDC_State** in the **Extract Initial Employee Data** container.
 - **Automatically store state in a database table:** selected.
 - **Connection manager for the database where the state is stored:** localhost Staging ADO NET.
 - **Table to use for storing state:** **a new table named [dbo].[cdc_states] in the Staging database.**
 - **State name:** CDC_State.
 3. Add a data flow task named **Extract Initial Employee Data** to the control flow, connecting the success precedence constraint from the **Mark Initial Load Start** task to the **Extract Initial Employee Data** task.
 4. In the **Extract Initial Employee Data** data flow, create an ADO.NET Source named **Employees**, with the following settings:
 - **ADO.NET connection manager:** **localhost HumanResources ADO NET.**
 - **Data access mode:** **table or view.**
 - **Name of the table or view:** **"dbo"."Employee".**
 5. Connect the data flow from the **Employees** source to a new ADO.NET Destination named **Employee Inserts** with the following settings:
 - **Connection manager:** **localhost Staging ADO NET.**
 - **Use a table or view:** **"dbo"."EmployeeInserts".**
 - **Mappings:** **map all available input columns to destination columns of the same name.**
 6. On the control flow, add a second CDC Control Task named **Mark Initial Load End** and connect the success precedence constraint from the **Extract Initial Employee Data** task to the **Mark Initial Load End** task. Configure the **Mark Initial Load End** task as follows:
 - **SQL Server CDC database ADO.NET connection manager:** localhost HumanResources ADO NET.
 - **CDC control operation:** mark initial load end*.
 - **Variable containing the CDC state:** User:: CDC_State.

- **Automatically store state in a database table:** selected.
- **Connection manager for the database where the state is stored:** localhost Staging ADO NET.
- **Table to use for storing state:** [dbo].[cdc_states].
- **State name:** CDC_State.

Note: Be careful not to select "Mark initial control start".

7. When you have completed the control flow, save the package.

► Task 5: Test Initial Extraction

1. Start debugging the **Extract Initial Employee Data.dtsx** package. When package execution is complete, stop debugging.
2. In SQL Server Management Studio, view the contents of the **dbo.EmployeeInserts** table in the **Staging** database to verify that the employee records have been transferred.
3. Refresh the view of the tables in the **Staging** database, and verify that a new table named **dbo.cdc_states** has been created. This table should contain an encoded string that indicates the CDC state.

► Task 6: Create a Package for Incremental Data Extraction

1. In Visual Studio, add a new SSIS package named **Extract Changed Employee Data.dtsx** to the **AdventureWorksETL** project.
2. Add a CDC Control Task to the control flow of the new package, and rename the task to **Get Processing Range**. Configure the **Get Processing Range** task as follows:
 - **SQL Server CDC database ADO.NET connection manager:** localhost HumanResources ADO NET.
 - **CDC control operation:** get processing range.
 - **Variable containing the CDC state:** a new variable named **CDC_State**.
 - **Automatically store state in a database table:** selected.
 - **Connection manager for the database where the state is stored:** localhost Staging ADO NET.
 - **Table to use for storing state:** [dbo].[cdc_states].
 - **State name:** CDC_State.
3. Add a data flow task named **Extract Changed Employee Data** to the control flow, connecting the success precedence constraint from the **Get Processing Range** task to the **Extract Changed Employee Data** task.
4. In the **Extract Changed Employee Data** data flow, add a **CDC Source** (in the **Other Sources** section of the SSIS Toolbox) named **Employee Changes**, with the following settings:
 - **ADO.NET connection manager:** localhost HumanResources ADO NET.
 - **CDC enabled table:** [dbo].[Employee].
 - **Capture instance:** dbo_Employee.
 - **CDC processing mode:** Net.
 - **Variable containing the CDC state:** User::CDC_State.

5. Connect the data flow from the **Employee Changes** source to a **CDC Splitter** transformation (in the **Other Transforms** section of the SSIS Toolbox).
6. Add an ADO.NET Destination named **Employee Inserts** below and to the left of the CDC Splitter transformation. Connect the **InsertOutput** data flow output from the CDC Splitter transformation to the **Employee Inserts** destination. Configure the **Employee Inserts** destination as follows:
 - **Connection manager:**
localhost Staging ADO NET.
 - **Use a table or view:**
"dbo"."EmployeeInserts".
 - **Mappings:** map all available input columns other than **__\$start_Isn**, **__\$operation**, and **__\$update_mask** to destination columns of the same name.
7. Add an ADO.NET Destination named **Employee Updates** directly below the CDC Splitter transformation, and connect the **UpdateOutput** data flow output from the CDC Splitter transformation to the **Employee Updates** destination. Configure the **Employee Updates** destination as follows:
 - **Connection manager:** localhost Staging ADO NET.
 - **Use a table or view:** "dbo"."EmployeeUpdates".
 - **Mappings:** map all available input columns other than **__\$start_Isn**, **__\$operation**, and **__\$update_mask** to destination columns of the same name.
8. Add an ADO.NET Destination named **Employee Deletes** below and to the right of the CDC Splitter transformation. Connect the **DeleteOutput** data flow output from the CDC Splitter transformation to the **Employee Deletes** destination. Configure the **Employee Deletes** destination as follows:
 - **Connection manager:**
localhost Staging ADO NET.
 - **Use a table or view:**
"dbo"."EmployeeDeletes".
 - **Mappings:** map all available input columns other than **__\$start_Isn**, **__\$operation**, **__\$update_mask**, and **__\$command_id** to destination columns of the same name.
9. On the control flow, add a second CDC Control Task named **Mark Processed Range**, and connect the success precedence constraint from the **Extract Changed Employee Data** task to the **Mark Processed Range** task. Configure the **Mark Processed Range** task as follows:
 - **SQL Server CDC database ADO.NET connection manager:** localhost HumanResources ADO NET.
 - **CDC control operation:** mark processed range.
 - **Variable containing the CDC state:** User:: CDC_State.
 - **Automatically store state in a database table:** selected.
 - **Connection manager for the database where the state is stored:** localhost Staging ADO NET.
 - **Table to use for storing state:**
[dbo].[cdc_states].
 - **State name:** CDC_State.
10. When you have completed the control flow, save the package.

► Task 7: Test Incremental Extraction

1. In Visual Studio, view the control flow for the **Extract Changed Employee Data.dtsx** package and start debugging.
2. When execution has completed, view the **Extract Changed Employee Data** data flow to verify that no rows were extracted (because the data is unchanged since the initial extraction).
3. Return to SQL Server Management Studio, open the **Change Employees.sql** Transact-SQL script file in the **D:\Labfiles\Lab09\Starter\Ex3** folder, and execute the script to:
 - Truncate the **dbo.EmployeeInserts**, **dbo.EmployeeUpdates**, and **dbo.EmployeeDeletes** tables in the **Staging** database.
 - Insert a new record in the **Employee** table in the **HumanResources** database.
 - Update **employee 281** to change the **Title** column value.
 - Delete **employee 273**.
4. In Visual Studio, start debugging again.
5. When execution has completed, stop debugging.
6. In SQL Server Management Studio, view the contents of the **dbo.EmployeeDeletes**, **dbo.EmployeeInserts**, and **dbo.EmployeeUpdates** tables in the **Staging** database to verify that they contain the inserted, updated, and deleted rows respectively.
7. Close Visual Studio and minimize SQL Server Management Studio when you are finished.

Results: After this exercise, you should have a **HumanResources** database in which CDC has been enabled, and an SSIS package that uses the CDC Control to extract the initial set of employee records. You should also have an SSIS package that uses the CDC Control and CDC data flow components to extract modified employee records, based on changes recorded by CDC.

Exercise 4: Using Change Tracking

Scenario

The **ResellerSales** database contains a **Resellers** table that does not include a column to indicate when records were inserted or modified. You plan to use the CT feature of SQL Server to identify records that have changed between data warehouse refresh cycles, and restrict data extractions to include only modified rows.

The main tasks for this exercise are as follows:

1. Enable Change Tracking
2. Create a Stored Procedure to Retrieve Modified Rows
3. Modify a Data Flow to Use the Stored Procedure
4. Test the Package

► Task 1: Enable Change Tracking

1. Enable Change Tracking in the **ResellerSales** database, and track changes in the **Resellers** table. You do not need to track which columns were modified. You can use the **Enable CT.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder to accomplish this.

2. Use the **Test CT.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder to perform the following tasks:
 - Get the current Change Tracking version number.
 - Retrieve all data from the **Resellers** table.
 - Store the current version number as the previously retrieved version.
 - Update **Resellers** table.
 - Get the new current version number.
 - Get all changes between the previous and current versions.
 - Store the current version number as the previously retrieved version.
 - Update the **Resellers** table again.
 - Get the new current version number.
 - Get all changes between the previous and current versions.
3. View the results returned by the script and verify that:
 - The first resultset shows all reseller records.
 - The second resultset indicates that the previously retrieved version was numbered 0, and the current version is numbered 1.
 - The third resultset indicates that the previously retrieved version was numbered 1, and the current version is numbered 2.

► **Task 2: Create a Stored Procedure to Retrieve Modified Rows**

1. In SQL Server Management Studio, execute a Transact-SQL statement that enables snapshot isolation and creates a stored procedure named **GetChangedResellers** in the **ResellerSales** database. The stored procedure should perform the following tasks (you can use the **Create SP.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder to accomplish this):
 - Set the isolation level to snapshot.
 - Retrieve the current Change Tracking version number.
 - If the **LastVersion** parameter is -1, assume that no previous versions have been retrieved, and return all records from the **Resellers** table.
 - If the **LastVersion** parameter is not -1, retrieve all changes between **LastVersion** and the current version.
 - Update the **LastVersion** parameter to the current version, so the calling application can store the last version retrieved for next time.
 - Set the isolation level back to read "committed".
2. Test the stored procedure by running the following query:

```
USE ResellerSales
GO
DECLARE @p BigInt = -1;
EXEC GetChangedResellers @p OUTPUT;
SELECT @p LastVersionRetrieved;
EXEC GetChangedResellers @p OUTPUT;
```

► Task 3: Modify a Data Flow to Use the Stored Procedure

1. In SQL Server Management Studio, open the **Reset Staging.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder and execute it.
2. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex4** folder. Open the **Extract Reseller Data.dtsx** SSIS package.
3. Add a decimal variable named **PreviousVersion** to the package.
4. Add an Execute SQL task named **Get Previously Extracted Version** to the control flow, and configure it to return a single row resultset by executing the following Transact-SQL statement in the **Staging** database, mapping the **LastVersion** column in the result set to the **User::PreviousVersion** variable:

```
SELECT MAX>LastVersion) LastVersion
FROM ExtractLog
WHERE DataSource = 'ResellerSales'
```

5. Add an Execute SQL task named **Update Previous Version** to the control flow, and configure it to execute the following Transact-SQL statement in the **Staging** database, mapping the **User:PreviousVersion** variable to parameter 0 in the query:

```
UPDATE ExtractLog
SET LastVersion = ?
WHERE DataSource = 'ResellerSales'
```

6. Make the necessary changes to the precedence constraint connections in the control flow so that:
 - The **Get Previously Extracted Version** task is executed after the **Get Last Extract Time** task, and before the **Extract Resellers** task.
 - The **Update Previous Version** task is executed after the **Update Last Extract Time** task, and before the **Send Success Notification** task.
7. Modify the **Resellers** source in the **Extract Resellers** data flow to retrieve reseller data by executing the following Transact-SQL statement, mapping the **@LastVersion** input/output parameter to the **User::PreviousVersion** variable:

```
EXEC GetChangedResellers ? OUTPUT
```

► Task 4: Test the Package

1. View the **Extract Resellers** data flow, and then start debugging the package and note the number of rows transferred. When package execution is complete, stop debugging.
2. In SQL Server Management Studio, view the contents of the **dbo.ExtractLog** table in the **Staging** database and verify that the **LastVersion** column for the **ResellerSales** data source has been updated.
3. View the contents of the **dbo.Resellers** table and note the rows that have been extracted. Close SQL Server Management Studio without saving any files.
4. In Visual Studio, debug the package again and verify that no rows are transferred in the **Extract Resellers** data flow. When package execution is complete, stop debugging.
5. Close Visual Studio.

Results: After this exercise, you should have a database in which CT has been enabled, and an SSIS package that uses a stored procedure to extract modified rows, based on changes recorded by CT.

Lesson 3

Loading Modified Data

A data warehousing solution loads data into the data warehouse at periodic intervals. Loading data into data warehouse tables presents some challenges, because you need to maintain historic versions of some data and minimize the performance impact on reporting workloads as the data warehouse is loaded.

This lesson describes the techniques you can use to implement a data warehouse load process.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the considerations for planning data loads.
- Use SQL Server Integration Services (SSIS) to load new and modified data into a data warehouse.
- Use Transact-SQL techniques to load data into a data warehouse.

Options for Incrementally Loading Data

There are a number of commonly used techniques for loading incremental changes to a data warehouse. The specific technique you should use to load a particular dimension or fact table depends on a number of factors. These include performance; having to update existing records in addition to inserting new ones; having to retain historical dimension attributes; and the location of the staging and data warehouse tables.

- Insert, Update, or Delete from CDC Output Tables
- Use a Lookup transformation
- Use the Slowly Changing Dimension transformation
- Use the MERGE statement
- Use a checksum
- Considerations for Deleting Data Warehouse records:
 - Use a logical deletion technique
 - Technique depends on how deleted records are staged

Insert, Update, or Delete Data Based on CDC Output Tables

If you use the CDC Splitter to stage modified source data into operation-specific tables, you can create an SSIS package that uses the business keys or unique column combinations in the staging tables to apply the appropriate changes to associated tables in the data warehouse. In most cases, delete operations in the source are applied as logical delete operations in the data warehouse—a deleted flag is set, instead of actually deleting the matching record.

Use a Lookup Transformation

You can use a Lookup transformation to determine whether, in the data warehouse, a matching record exists for a record that has been extracted from the data sources. You can then use the no match output of the Lookup transformation to create a data flow for new records that have to be inserted into the data warehouse. Optionally, you can also use the match output of the Lookup transformation to create a data flow that updates existing data warehouse records with new values from the extracted ones.

Use the Slowly Changing Dimension Transformation

The SCD transformation helps you to create a complex data flow that inserts new dimension members. It applies type 1 or type 2 changes to existing dimension members, depending on which attributes have been updated. In many data warehousing scenarios, the SCD transformation provides an easy-to-implement solution for refreshing dimension tables. However, its performance can be limited for ETL

processes with extremely large numbers of rows, for which you might need to create a custom solution for SCD.

Use the MERGE Statement

The MERGE statement is a Transact-SQL construct that you can use to perform insert, update, and delete operations in the same statement. It works by matching rows in a source rowset with rows in a target table, taking appropriate action to merge the source with the target.

The MERGE statement is appropriate when the source or staging tables and the data warehouse tables are implemented in SQL Server databases and the ETL process can execute the MERGE statement using a connection through which all source and target tables can be accessed. In practical terms, this requires that:

- The staging tables and data warehouse are co-located in the same SQL Server database.
- The staging and data warehouse tables are located in multiple databases in the same SQL Server instance, and the credentials used to execute the MERGE statement have appropriate user rights in both databases.
- The staging tables are located in a different SQL Server instance than the data warehouse, but a linked server has been defined, which means that the MERGE statement can access both databases. The performance of the MERGE statement over the linked server connection is acceptable.

Use a Checksum

You can use the columns in the staged dimension records to generate a checksum value, and then compare this with a checksum generated from the historic attributes in the corresponding dimension table, to identify rows that require a type 2 or type 3 change. When combined with a Lookup transformation to identify new or modified rows, this technique can form the basis of a custom solution for SCD.

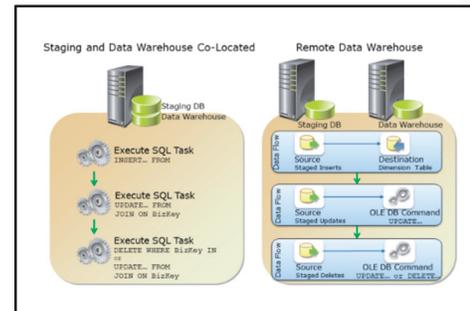
Considerations for Deleting Data Warehouse Records

If you need to propagate record deletions in source systems to the data warehouse, you should consider the following guidelines:

- In most cases, you should use a logical deletion technique where you indicate that a record is no longer valid, by setting a Boolean column value. It is not common practice to physically delete records from a data warehouse unless you have a compelling business reason to discard all historical information relating to them.
- The techniques you can use to delete records (or mark them as logically deleted) when loading data depend on how you have staged deleted records:
 - If the staging tables for a dimension or fact table contain all valid records (not just those that have been modified since the previous refresh cycle), you can delete any existing records in the data warehouse that do not exist in the staging tables.
 - If the staged data indicates logical deletions in the form of a Boolean column value, and you need to apply logical deletes in the data warehouse, you can treat these deletions as updates.
 - If the keys of records to be deleted are stored separately from new and updated records in the staging database, you might want to perform two distinct load operations for each dimension or fact table—one to load new and updated records, and another for deletions.

Using CDC Output Tables

As described earlier in this module, you can use the CDC Splitter or a custom data flow to stage CDC insert, update, or delete records in separate tables, based on the **Operation** value in the CDC output. If your data extraction process has taken this approach, you can apply the appropriate changes to the corresponding data warehouse tables. The specific control flow tasks you can use to perform this type of data load depend on the relative locations of the staging tables and the data warehouse.



Loading a Data Warehouse from a Co-located Staging Database

If the staging tables and the data warehouse are co-located in the same database or server instance—or can be connected by using a linked server—you can use SQL command control flow tasks to execute set-based Transact-SQL statements that load the data warehouse tables.

For example, you could use the following Transact-SQL statements to load data from CDC output tables into a co-located data warehouse table:

Inserting, Updating, and Deleting Data Based on Co-located CDC Output Tables

```
-- Insert new records
INSERT INTO dw.DimProduct (ProductAltKey, ProductName, Description, Price)
SELECT ProductBizKey, ProductName, Description, Price
FROM stg.ProductInserts;

-- Update modified records
UPDATE dw.DimProduct
SET    dw.DimProduct.ProductName = stg.ProductUpdates.ProductName,
       dw.DimProduct.Description = stg.ProductUpdates.Description,
       dw.DimProduct.Price = stg.ProductUpdates.Price
FROM dw.DimProduct JOIN stg.ProductUpdates
ON dw.DimProduct.ProductAltKey = stg.ProductUpdates.ProductBizKey;

-- Mark deleted records as deleted
UPDATE dw.DimProduct
SET dw.DimProduct.Deleted = 1
FROM dw.DimProduct JOIN stg.ProductDeletes
ON dw.DimProduct.ProductAltKey = stg.ProductDeletes.ProductBizKey;
```

You could use a DELETE statement to remove deleted records from the data warehouse table. However, deleting records in this way can result in loss of historical analytical and reporting data. This can be further complicated by the presence of foreign key constraints between fact and dimension tables. A more common approach is to perform a logical delete by setting a column to a Boolean value, to indicate that the record was deleted from the source system.

Loading a Remote Data Warehouse

If the data warehouse is stored on a different server from the staging database, and no linked server connection is available, you can apply the necessary updates in the staging tables to the data warehouse by creating a data flow for each operation.

To load records from an inserts staging table, create a data flow that includes a source component to extract records from the staging table, and a destination that maps the extracted column to the appropriate data warehouse table.

To apply the changes recorded in an updates staging table to a data warehouse table, create a data flow that includes a source component to extract records from the staging table, and an OLE DB Command transformation to execute an UPDATE statement that sets the changeable data warehouse table columns to the corresponding staging table values. This is based on a join between the business key in the staging table and the alternative key in the data warehouse table.

To apply deletes to data warehouse tables based on records in a deletes staging table, create a data flow that includes a source component to extract records from the staging table. You can also create an OLE DB Command transformation to execute a DELETE statement that matches the business key in the staging table to the alternative key in the data warehouse table. Alternatively, use the OLE DB command to perform a logical delete by executing an UPDATE statement that sets the deleted flag to 1 for records where the business key in the staging table matches the alternative key in the data warehouse table.

Demonstration: Using CDC Output Tables

In this demonstration, you will see how to load data from CDC output tables.

Demonstration Steps

Load Data from CDC Output Tables

1. Ensure you have completed the previous demonstrations in this module.
2. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the SQL Server database engine by using Windows authentication.
3. In Object Explorer, expand **Databases**, expand **DemoDW**, and expand **Tables**.
4. Right-click each of the following tables and click **Select Top 1000 Rows**:
 - o **dw.DimShipper**. This is the dimension table in the data warehouse.
 - o **stg.ShipperDeletes**. This is the table of records that have been deleted in the source system.
 - o **stg.ShipperInserts**. This is the table of new records in the source system.
 - o **stg.ShipperUpdates**. This is the table of rows that have been updated in the source system.
5. Start Visual Studio and open the **IncrementalETL.sln** solution in the **D:\Demofiles\Mod09** folder.
6. In Solution Explorer, double-click the **Load Shippers.dtsx** SSIS package.
7. On the control flow surface, double-click the **Load Inserted Shippers** Execute SQL task. Note that the SQL Statement inserts data into **dw.DimShippers** from the **stg.ShipperInserts** table. Click **Cancel**.
8. On the control flow surface, double-click the **Load Updated Shippers** Execute SQL task. Note that the SQL statement updates data in **dw.DimShippers** with new values from the **stg.ShipperUpdates** table. Click **Cancel**.
9. On the control flow surface, double-click the **Load Deleted Shippers** data flow task. On the data flow surface, note that the task extracts data from the **stg.ShipperDeletes** table, and then uses an OLE DB Command transformation to update the **Deleted** column in **dw.DimShippers** for the extracted rows.
10. On the **Debug** menu, click **Start Debugging**, and observe the control flow as it executes.

11. When execution is complete, on the **Debug** menu, click **Stop Debugging** and minimize Visual Studio.
12. In SQL Server Management Studio, right-click the **dw.DimShipper** table and click **Select Top 1000 Rows** to review the changes made.
13. Minimize SQL Server Management Studio.

The Lookup Transformation

To use a Lookup transformation when loading a data warehouse table, connect the source output that extracts the staged data to the Lookup transformation and apply the following configuration settings:

- Redirect nonmatched rows to the no match output.
- Look up the primary key column or columns in the dimension or fact table you want to refresh by matching it to one or more input columns from the staged data. If the staged data includes a business key column and the business key is stored as an alternative key in the data warehouse table, match one to the other. Alternatively, you can match a combination of columns that uniquely identifies a fact or dimension member.
- Connect the no match output from the Lookup transformation to a data flow that ultimately inserts new records into the data warehouse.
- To update existing data warehouse records with modified values in the staged data, connect the match output of the Lookup transformation to a data flow that uses an OLE DB Command transformation. This is based on the primary key you retrieved in the Lookup transformation.

- Redirect nonmatched rows to the *no match* output
- Look up extracted data in a dimension or fact table based on a business key or unique combination of keys
- If no match is found, insert a new record
- Optionally, if a match is found, update non-key columns to apply a type 1 change

 **Note:** The Lookup transformation uses an in-memory cache to optimize performance. If the same data set is used in multiple lookup operations, you can persist the cache to a file and use a Cache Connection Manager to reference it. This further improves performance by decreasing the time it takes to load the cache, but results in lookup operations against a data set that might not be as up to date as the data in the database. For more information about configuring caching for the Lookup transformation, see *Lookup Transformation* in the SQL Server Technical Documentation.

Demonstration: Using the Lookup Transformation

In this demonstration, you will see how to:

- Use a Lookup Transformation to Insert Rows.
- Use a Lookup Transformation to Insert and Update Rows.

Demonstration Steps

Use a Lookup Transformation to Insert Rows

1. Ensure you have completed the previous demonstrations in this module.
2. Return to Visual Studio, and in Solution Explorer, double-click the **Load Geography.dtsx** SSIS package.
3. On the control flow surface, double-click **Load Geography Dimension** to view the data flow surface.
4. On the data flow surface, double-click **Staged Geography Data**. Note that the SQL command used by the OLE DB source extracts geography data from the **stg.Customers** and **stg.Salespeople** tables, and then click **Cancel**.
5. On the data flow surface, double-click **Lookup Existing Geographies**.
6. Note the following configuration settings of the Lookup transformation, and then click **Cancel**:
 - On the **General** tab, unmatched rows are redirected to the no match output.
 - On the **Connection** tab, the data to be matched is retrieved from the **dw.DimGeography** table.
 - On the **Columns** tab, the **GeographyKey** column is retrieved for rows where the input columns are matched.
7. On the data flow surface, note that the data flow arrow connecting **Lookup Existing Geographies** to **New Geographies** represents the no match data flow.
8. Double-click **New Geographies**.
9. Note that the rows in the no match data flow are inserted into the **dw.DimGeography** table, and then click **Cancel**.
10. On the **Debug** menu, click **Start Debugging**, and observe the data flow as it executes. Note that, while four rows are extracted from the staging tables, only one does not match an existing record. The new record is loaded into the data warehouse, and the rows that match existing records are discarded.
11. When execution is complete, on the **Debug** menu, click **Stop Debugging**.

Use a Lookup Transformation to Insert and Update Rows

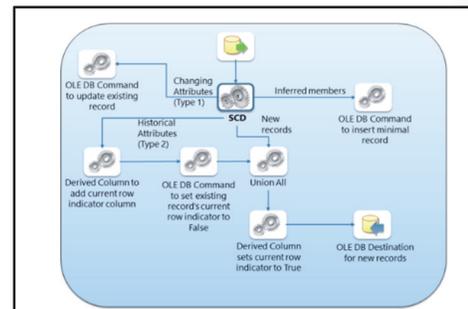
1. In Visual Studio, in Solution Explorer, double-click the **Load Products.dtsx** SSIS package.
2. On the control flow surface, double-click **Load Product Dimension** to view the data flow surface.
3. On the data flow surface, double-click **Staged Products**.
4. Note that the SQL command used by the OLE DB source extracts product data from the **stg.Products** table, and then click **Cancel**.
5. On the data flow surface, double-click **Lookup Existing Products**.

6. Note the following configuration settings of the Lookup transformation, and then click **Cancel**:
 - On the **General** tab, unmatched rows are redirected to the no match output.
 - On the **Connection** tab, the data to be matched is retrieved from the **dw.DimProduct** table.
 - On the **Columns** tab, the **ProductKey** column is retrieved for rows where the **ProductBusinessKey** column in the staging table matches the **ProductAltKey** column in the data warehouse dimension table.
7. On the data flow surface, note that the data flow arrow connecting **Lookup Existing Products** to **Insert New Products** represents the no match data flow. The data flow arrow connecting **Lookup Existing Products** to **Update Existing Products** represents the match data flow.
8. Double-click **Insert New Products**.
9. Note that the rows in the no match data flow are inserted into the **dw.DimProduct** table, and then click **Cancel**.
10. Double-click **Update Existing Products**.
11. Note the following configuration settings, and then click **Cancel**:
 - On the **Connection Managers** tab, the OLE DB Command transformation connects to the **DemoDW** database.
 - On the **Component Properties** tab, the **SQLCommand** property contains a parameterized Transact-SQL statement that updates the **ProductName**, **ProductDescription**, and **ProductCategoryName** columns for a given **ProductKey**.
 - On the **Column Mapping** tab, the **ProductName**, **ProductDescription**, **ProductCategoryName**, and **ProductKey** input columns from the match data flow are mapped to the parameters in the SQL command.
12. On the **Debug** menu, click **Start Debugging**, and observe the data flow as it executes. Note the number of rows extracted from the staging tables, and how the Lookup transformation splits these rows to insert new records and update existing ones.
13. When execution is complete, on the **Debug** menu, click **Stop Debugging**, and then minimize Visual Studio.

The Slowly Changing Dimension Transformation

The Slowly Changing Dimension (SCD) transformation provides a wizard you can use to generate a complex data flow to handle inserts and updates for a dimension table. Using the SCD wizard, you can specify:

- The columns containing keys that can be used to look up existing dimension records in the data warehouse.
- The non-key columns that are fixed, changing, or historic attributes.
- Whether or not changes to a fixed column should produce an error or be ignored.



- The column in the dimension table that should be used to indicate the current version of a dimension member for which historic attributes have changed over time.
- Whether or not the staged data includes inferred members for which a minimal record should be inserted. Inferred members are identified, based on the value of a specified column in the source.

After completing the wizard, the SCD transformation generates a data flow that includes the following:

- A path to insert new dimension records.
- A path to update dimension records where a changing attribute has been modified. This is an implementation of a type 1 change.
- A path to update the current record indicator and insert a new record for dimension members where a historic attribute has been modified. This is an implementation of a type 2 change.
- If specified, a path to insert minimal records for inferred members in the source data.

The MERGE Statement

The MERGE statement combines insert and update operations to merge changes from one table into another. In a data warehouse loading scenario, it can be used to load new and updated rows in dimension and fact tables. When used to load a data warehouse table, a MERGE statement includes the following components:

- An INTO clause containing the name of the target table being loaded.
- A USING clause containing a query that defines the source data to be loaded into the target table. This is usually executed against tables in a staging database, often using JOIN clauses to look up dimension keys in the data warehouse.
- An ON clause specifying the criteria used to match rows in the source rowset with rows in the target table.
- A WHEN MATCHED clause specifying the action to be taken for rows in the target table that match rows in the source rowset—often an UPDATE statement to apply a type 1 change.
- A WHEN NOT MATCHED clause specifying the action to be taken when no matches are found in the target table for rows in the source rowset. This is usually an INSERT statement to add a new record to the data warehouse.

- Matches source and target rows
- Performs insert, update, or delete operations based on row matching results

The following code example shows how the MERGE statement can be used to load new records and perform type 1 updates in a dimension table:

Applying Type 1 Updates with MERGE

```
MERGE INTO DimProduct WITH (TABLOCK) AS tgt
USING
-- Query to return staged data
(SELECT ProductAltKey, ProductName, Color
 FROM StagedProducts) AS src (ProductAltKey, ProductName, Color)
-- Match staged records to existing dimension records
ON (src.ProductAltKey = tgt.ProductAltKey)
-- If a record for this product already exists, update it
WHEN MATCHED THEN
UPDATE
SET ProductName = src.ProductName,
    Color = src.Color
-- If not, insert a new record
WHEN NOT MATCHED THEN
INSERT (ProductAltKey, ProductName, Color)
VALUES (src.ProductAltKey, src.ProductName, src.Color);
```

Additionally, when combined with the OUTPUT clause, you can use the **\$action** metadata column to detect updates, and implement type 2 changes.

The following code example shows how to perform type 2 updates with the MERGE statement:

Applying Type 2 Updates with the MERGE Statement

```
INSERT INTO DimCustomer WITH (TABLOCK)
(CustomerAltKey, Name, City, CurrentFlag, StartDate, EndDate)
SELECT CustomerAltKey, Name, City, 1, getdate(), NULL
FROM
(MERGE INTO DimCustomer AS tgt
USING
-- Query to return staged customer data
(SELECT CustomerAltKey, Name, City
 FROM [StagedCustomer]
 AS src (CustomerAltKey, Name, City)
-- Match staged customers to existing (current) dimension records
ON (src.CustomerAltKey = tgt.CustomerAltKey AND tgt.CurrentFlag = 1)
-- If a current record for this customer already exists, mark it as a type 2 change
WHEN MATCHED THEN
UPDATE
SET tgt.CurrentFlag = 0, tgt.EndDate = getdate()
-- If not, insert a new record
WHEN NOT MATCHED THEN
INSERT (CustomerAltKey, Name, City, CurrentFlag, StartDate, EndDate)
VALUES (CustomerAltKey, Name, City, 1, getdate(), NULL)
-- Now output the records you've inserted or updated
OUTPUT $action, CustomerAltKey, Name, City)
AS Type2Changes(MergeAction, CustomerAltKey, Name, City)
-- filter them so you insert new records for the type 2 updates.
WHERE MergeAction = 'UPDATE';
```

When you implement an incremental ETL process with SQL Server Integration Services, you can use a SQL command task in the control flow to execute a MERGE statement. However, you must ensure that the connection manager assigned to the SQL command task provides access to the source and target tables.

Demonstration: Using the Merge Statement

In this demonstration, you will see how to use the MERGE statement.

Demonstration Steps

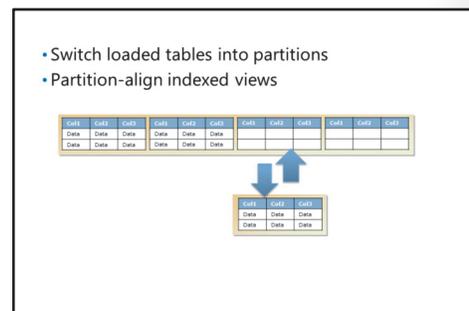
Use the MERGE Statement

1. In the **D:\Demofiles\Mod09\LoadModifiedData** folder, run **SetupB.cmd** as Administrator.
2. In the **User Account Control** dialog box, click **Yes**.
3. In SQL Server Management Studio, in Object Explorer, in the **DemoDW** database, right-click the **stg.SalesOrders** table, and click **Select Top 1000 Rows**. This table contains staged sales order data.
4. Right-click the **dw.FactSalesOrders** table and click **Select Top 1000 Rows**. This table contains sales order fact data. Note that the staged data includes three order records that do not exist in the data warehouse fact table (with **OrderNo** and **ItemNo** values of 1005 and 1; 1006 and 1; and 1006 and 2 respectively), and one record that does exist but for which the Cost value has been modified (OrderNo 1004, ItemNo 1).
5. Open the **Merge Sales Orders.sql** file in the **D:\Demofiles\Mod09\LoadModifiedData** folder.
6. Review the Transact-SQL code, noting the following details:
 - o The MERGE statement specifies the **DemoDW.dw.FactSalesOrders** table as the target.
 - o A query that returns staged sales orders and uses joins to look up dimension keys in the data warehouse is specified as the source.
 - o The target and source tables are matched on the **OrderNo** and **ItemNo** columns.
 - o Matched rows are updated in the target.
 - o Unmatched rows are inserted into the target.
7. Click **Execute** and note the number of rows affected.
8. Right-click the **dw.FactSalesOrders** table and click **Select Top 1000 Rows**. Then compare the contents of the table with the results of the previous query you performed in step 3.
9. Minimize SQL Server Management Studio.

Partition Switching

Fact tables are often partitioned to simplify management and data loads. When using partitioned fact tables, you should consider the following guidelines for your ETL data load processes:

- Switch loaded tables into partitions.
- Partition-align indexed views.



Switch Loaded Tables into Partitions

After you create partitioned fact tables, you can optimize data load operations by switching a loaded table into an empty partition. This technique can be used to load a partition from a table that:

- Has the same schema as the partition, including column names and data types.
- Has the same indexes as the partition, including columnstore indexes.
- Has the same compression setting as the partition.
- Has a check constraint that uses the same criteria as the partition function.
- Is stored on the same filegroup as the partition.

To use this technique to load new data into a partition, maintain an empty partition at the end of the table. The lower bound of the partition range for the empty partition should be the date key value for the next set of data to be loaded. The basic technique to load new data into a partition uses the following procedure:

1. If each partition is stored on its own filegroup, add a filegroup to the database and set it as the next used filegroup for the partition scheme.
2. Split the empty partition at the end of the table, specifying the key for the upper bound of the data to be loaded. This creates one empty partition for the new data, and another to be maintained at the end of the table for the next load cycle.
3. Create a table on the same filegroup as the second to last, empty partition, with the same columns and data types as the partitioned table. For fastest load performance, create this table as a heap (a table with no indexes).
4. Bulk insert the staged data into the load table you created in the previous step.
5. Add a constraint that checks that the partitioning key column values are within the range of the target partition to the load table.
6. Add indexes to the load table that match those on the partitioned table.
7. Switch the partition and the load table.
8. Drop the load table.

This technique works best when the table is partitioned on a date key that reflects the data warehouse load cycle, so each new load is performed into a new partition. However, it can also be used when partitions do not match load intervals.

- When partitions are based on more frequent intervals than load cycles (for example, each partition holds a week's worth of data, but the data is loaded monthly), you can switch multiple load tables into multiple partitions.
- When partitions are based on less frequent intervals than load cycles (for example, each partition holds a month's worth of data, but the data is loaded daily), you can:
 - Create a new partition for the load, and then merge it with the previous partition.
 - Switch out a partially loaded partition, drop the indexes on the partially populated load table, insert the new rows, recreate the indexes, and switch the partition back in. This technique can also be used for late arriving facts (rows that belong in partitions that have previously been loaded) and updates.



Top 10 Best Practices for Building a Large Scale Relational Data Warehouse

<http://aka.ms/twl2xa>

Partition-Align Indexed Views

If you plan to use indexed views in the data warehouse, align the indexes to the partitions on the underlying table. When indexed views are partition-aligned, you can switch partitions without having to drop and recreate the indexes on the views.

Demonstration: Partition Switching

In this demonstration, you will see how to:

- Split a partition.
- Create a load table.
- Switch a partition.

Demonstration Steps

Split a Partition

1. Ensure you have completed the previous demonstration in this module.
2. Return to SQL Server Management Studio and open the **Load Partition.sql** file in the **D:\Demofiles\LoadModifiedData** folder.
3. Select the code under the comment **Create a partitioned table**, and then click **Execute**. This code creates a database with a partitioned fact table, on which a columnstore index has been created.
4. Select the code under the comment **View partition metadata**, and then click **Execute**. This code shows the partitions in the table with their starting and ending range values, and the number of rows they contain. Note that the partitions are shown once for each index (or for the heap if no clustered index exists). The final partition (4) is for key values of 20020101 or higher and currently contains no rows.
5. Select the code under the comment **Add a new filegroup and make it the next used**, and then click **Execute**. This code creates a filegroup, and configures the partition scheme to use it for the next partition to be created.
6. Select the code under the comment **Split the empty partition at the end**, and then click **Execute**. This code splits the partition function to create a new partition for keys with the value 20030101 or higher.
7. Select the code under the comment **View partition metadata** again, and then click **Execute**. This time, the query is filtered to avoid including the same partition multiple times. Note that the table now has three empty partitions (1, 4 and 5).

Create a Load Table

1. Select the code under the comment **Create a load table**, and then click **Execute**. This code creates a table on the same filegroup as partition 4, with the same schema as the partitioned table.
2. Select the code under the comment **Bulk load new data**, and then click **Execute**. This code inserts the data to be loaded into the load table (in a real solution, this would typically be bulk loaded from staging tables).
3. Select the code under the comment **Add constraints and indexes**, and then click **Execute**. This code adds a check constraint to the table that matches the partition function criteria, and a columnstore index that matches the index on the partitioned table.

Switch a Partition

1. Select the code under the comment **Switch the partition**, and then click **Execute**. This code switches the load table with the partition on which the value 20020101 belongs. Note that the required partition number is returned by the \$PARTITION function.
2. Select the code under the comment **Clean up and view partition metadata**, and then click **Execute**. This code drops the load table and returns the metadata for the partitions. Note that partition 4 now contains two rows that were inserted into the load table.
3. Close SQL Server Management Studio.

Check Your Knowledge

Question	
Which of these options is NOT used in a MERGE statement?	
Select the correct answer.	
<input type="checkbox"/>	WHEN NOT MATCHED
<input type="checkbox"/>	USING
<input type="checkbox"/>	ON
<input type="checkbox"/>	WHEN MATCHED
<input type="checkbox"/>	OFF

Verify the correctness of the statement by placing a mark in the column to the right.

Statement	Answer
<p>Is the following statement true or false?</p> <p>"The Lookup transformation uses an in-memory cache to optimize performance."</p>	

Lesson 4

Temporal Tables

Temporal is a new feature, based on the ANSI SQL 2011 standard. This feature addresses the problem of capturing and storing data changes in *System-Versioned* tables. Data changes need to be captured and saved for auditing or reporting purposes on historical data, in addition to current data.



Note: The terms *Temporal tables* and *System-Versioned tables* are both used. In this lesson, we will use *temporal* to describe the SQL Server feature and *System-Versioned* to describe the tables that have been configured for the temporal feature.

The Slowly Changing Dimensions (SCD) feature in a data warehouse is a common way to capture changes from the OLTP system, so that data can be viewed and analyzed over time. SCD coordinates the changes and insertions of dimension data in data warehouses using the SCD wizard. SCD was covered in a previous topic in this module.



Note: The SCD wizard is specific to SQL Server connections.

Developers need to address the problem of storing information about changed data, including when it was changed, and who changed it. In addition to the SCD wizard used in data warehousing, it is common for triggers or custom code to be added to databases to extract the changes and store this for future reference. The introduction of System-Versioned tables means that SQL Server can capture this information automatically.



Note: System-Versioned tables are also described in another module of this course. This lesson describes the data warehouse specific use of System-Versioned tables.

The following link provides background information to Temporal Tables:



Temporal Tables

<http://aka.ms/fut3c9>

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the benefits of using System-Versioned tables in your database.
- Understand considerations, including the limitations, of using System-Versioned tables.
- Create System-Versioned tables.
- Query versioned tables.

About System-Versioned Tables

The temporal feature helps you to store all changes to your data. Temporal operates by creating a historical table linked by a primary key to your System-Versioned (or primary) table. Every record in a System-Versioned table has a full history of changes stored in this historical table. The historical table will have a row for each record change (insert, update or delete) that has occurred in the primary table—this includes the data itself, in addition to the time and date of the change. Only current records are stored on the System-Versioned table, so you have to query the historical table for noncurrent records.

- System-Versioned tables summary:
 - Enable a full history of data changes
 - Current and historical tables operate as a pair
 - Feature can be added to existing tables
 - Historical table can be named, or take system name
 - Current table must have a primary key
- System-versioned tables operation:
 - Versioning is automatic
 - SysStartTime and SysEndTime columns define the validity period for data versions
 - History specific queries are used to obtain version information

Temporal tables reduce the need to interrogate and restore backup records. You can use Transact-SQL to query historical records that are candidates for restoring, without having to access remote backup data.

In summary, you can use the temporal feature to:

- **View Historical Data.** You can use Transact-SQL to show how your data was at a specific point in time.
- **Compare Historical Data.** You can compare how your data was at any two or more times, using Transact-SQL.
- **Repair Damage.** If data corruption occurs on the primary table, you can use Transact-SQL to query for the last known “good” record on the history table.
- **Perform Audit.** Auditing is easier because you have a full record of all data as it changes.
- **Manage Slowly Changing Dimension (SCD) Tables.** The retention of temporal data history records can be long term with little extra overhead. This means that audit and point-in-time analysis for slowly changing data is easy to deploy using Transact-SQL. Data errors can also be flagged and the records be returned to an early period, again using Transact-SQL.

System-Versioned Tables Summary

To use the feature, you start by creating a System-Versioned table, by either creating a new table or modifying an existing table.

When you create a System-Versioned table, SQL Server creates two tables linked by a primary key:

- **Current Data Table.** The System-Versioned table, also known as the primary table.
- **Historical Data Table.** A table in which data history is stored.

 **Note:** You can name the historical table yourself; otherwise, SQL Server uses a default naming convention.

 **Note:** The current table requires a primary key so that SQL Server can establish a relationship to the historical table. You can also join back to this table to see the full history for any given record.

Two datetime2 columns, in both current and historical tables, are included when a System-Versioned table is created. These are:

- **SysStartTime Column** or **Period Start Column**: a record of the start time and date from which the data row is valid.
- **SysEndTime Column** or **Period End Column**: a record of the time and date at which the row was superseded by a change or delete operation.

Summary of System-Versioned Tables Operation

Data is always created on the historical table first. The **SysStartTime** and **SysEndTime** columns are updated to reflect the current validity of the data. The row created is then written to the current table.



Note: In the history tables, an “open” row contains data that is current; a “closed” row contains historic data.

The following list provides summaries of how SQL Server processes temporal data:

- **Inserts.** The row needs to be written as an open row to the historical table, and then copied to the current table:
 - **Create a Historical Row.** **SysStartTime** is set to the begin time for the transaction. **SysEndTime** is set to **9999-12-31**. The data is written to the historical table, and this new row is now marked as open.
 - **Create a Current Row.** The row data is then inserted into the current table.
- **Updates.** There are existing rows in both current and historic tables that need to be updated:
 - **Update the Existing Historical Row.** Update the **SysEndTime** for the existing row to the transaction begin time. This row is now marked as closed.
 - **Insert a New Historical Row.** A new row is inserted, based on the new data. The **SysStartTime** is set to the transaction begin time and the **SysEndTime** is set to **9999-12-31**. The new row is now marked as open and the old row closed.
 - **Update the Current Table.** The current table row is updated.
- **Deletes.** There are existing rows in both current and historic tables:
 - **Update the Existing Historical Row. Update the SysEndTime for the existing row to the transaction end time. The record is marked as closed. Only history specific queries will return the data records.**
 - **Delete the Current Row.** The row is removed from the current table.



Note: Only history data specific queries return data for a deleted row.

Considerations for System-Versioned Tables

There are a number of points that you should consider and tasks to carry out before you can use a System-Versioned table:

- The system date columns **SysStartTime** and **SysEndTime** must use the datetime2 data type.
- The System-Versioned table must have a primary key, but the history table cannot use constraints.
- To name the history table, you must specify the schema name, in addition to the table name.
- By default, the compression of the history table is set to PAGE.
- The history table must reside in the same database as the current table.
- System-Versioned tables are not compatible with FILETABLE or FILESTREAM features because SQL Server cannot track changes that happen outside of itself.
- Columns with a BLOB data type, such as varchar(max) or image, can result in high storage requirements because the history table will store the history values as the same type.
- INSERT and UPDATE statements cannot reference the **SysStartTime** or **SysEndTime** columns.
- You cannot modify data in the history table directly.
- You cannot truncate a System-Versioned table. Turn SYSTEM_VERSIONING OFF to truncate the table.
- Merge replication is not supported.



Temporal Table Considerations and Limitations

<http://aka.ms/vlejh6>

Creating System-Versioned Tables

To create a new System-Versioned table, use the standard CREATE TABLE code with three additional columns, and then specify to turn on system-versioning.

- Main considerations when using System-Versioned tables:
 - Current table must have a primary key
 - SysStartTime and SysEndTime must be datetime2
 - Data in the historical table cannot be directly modified
 - Current table cannot be truncated when SYSTEM_VERSIONING is ON
 - FILETABLE or FILESTREAM are not supported
 - History table is PAGE compressed by default

- Create a new System-Versioned table:

```
CREATE TABLE dbo.Employee
(
    EmployeeID int NOT NULL PRIMARY KEY CLUSTERED,
    ManagerID int NULL,
    FirstName varchar(50) NOT NULL,
    LastName varchar(50) NOT NULL,
    SysStartTime datetime2 GENERATED ALWAYS AS ROW START NOT NULL,
    SysEndTime datetime2 GENERATED ALWAYS AS ROW END NOT NULL,
    PERIOD FOR SYSTEM_TIME ('SysStartTime', 'SysEndTime')
)
WITH (SYSTEM_VERSIONING = ON (HISTORY_TABLE = dbo.EmployeeHistory));
```

Create a new table and set the SYSTEM_VERSIONING feature ON.

Create Temporal Table

```
CREATE TABLE dbo.Employee
(
    EmployeeID int NOT NULL PRIMARY KEY CLUSTERED,
    ManagerID int NULL,
    FirstName varchar(50) NOT NULL,
    LastName varchar(50) NOT NULL,
    SysStartTime datetime2 GENERATED ALWAYS AS ROW START NOT NULL,
    SysEndTime datetime2 GENERATED ALWAYS AS ROW END NOT NULL,
    PERIOD FOR SYSTEM_TIME (SysStartTime,SysEndTime)
)
WITH (SYSTEM_VERSIONING = ON (HISTORY_TABLE = dbo.EmployeeHistory));
```

So that SQL Server can name the historical table, this statement can be run without including the (HISTORY_TABLE = dbo.EmployeeHistory) clause.

You must include the **SysStartTime** and **SysEndTime** columns and the PERIOD FOR SYSTEM_TIME that references these columns. You can change the name of these columns and change the references to them in the PERIOD FOR SYSTEM_TIME parameters.

The Manager table in the following example has the date columns named as **DateFrom** and **DateTo**. These names are also used in the history table:

Change the Names of the Start and End System Columns

```
CREATE TABLE dbo.Manager
(
    ManagerID int NOT NULL PRIMARY KEY CLUSTERED,
    FirstName varchar(50) NOT NULL,
    LastName varchar(50) NOT NULL,
    DateFrom datetime2 GENERATED ALWAYS AS ROW START NOT NULL,
    DateTo datetime2 GENERATED ALWAYS AS ROW END NOT NULL,
    PERIOD FOR SYSTEM_TIME (DateFrom,DateTo)
)
WITH (SYSTEM_VERSIONING = ON (HISTORY_TABLE = dbo.ManagerHistory));
```

An existing table can be converted to a System-Versioned table by:

- Adding in the start and end date columns.
- Setting the system-versioning flag on.

The following code adds the system datetime2 columns to the **Sales** table. After these are created, the code alters the **Sales** table so that SYSTEM_VERSIONING is ON and changes are stored. In this example, the history table has been named:

Make an Existing Table System-Versioned

```
ALTER TABLE dbo.Sales
ADD
    SysStartTime datetime2(0) GENERATED ALWAYS AS ROW START HIDDEN CONSTRAINT DF_SysStartTime
    DEFAULT SYSUTCDATETIME(),
    SysEndTime datetime2(0) GENERATED ALWAYS AS ROW END HIDDEN CONSTRAINT DF_SysEndTime
    DEFAULT CONVERT(datetime2 (0), '9999-12-31 23:59:59'),
    PERIOD FOR SYSTEM_TIME (SysStartTime, SysEndTime);
GO

ALTER TABLE dbo.Sales
SET (SYSTEM_VERSIONING = ON (HISTORY_TABLE = dbo.SalesHistory));
```

Querying System-Versioned Tables

As you have already learned in this lesson, a System-Versioned table always has a primary key. This primary key links the current and historical tables so you can view history for your data over a period of time. However, you do not need to join the current and historical tables explicitly, because a new clause has been added to Transact-SQL to do this for you. You can use the FOR SYSTEM_TIME clause with one of the following four subclauses to return the data for a specific date or within a date range:

- System-Versioned tables can be queried using the FOR SYSTEM_TIME clause and one of the following four subclauses:
 - AS OF <date_time>
 - FROM <start_date_time> TO <end_date_time>
 - BETWEEN <start_date_time> AND <end_date_time>
 - CONTAINED IN (<start_date_time>, <end_date_time>)
- Or use ALL to return everything

- AS OF <date_time>
- FROM <start_date_time> TO <end_date_time>
- BETWEEN <start_date_time> AND <end_date_time>
- CONTAINED IN (<start_date_time>, <end_date_time>)

You can also use the ALL subclause to return data from the current and historical tables without any time restrictions.

The AS OF subclause returns the data at a given point in time. The following code uses the datetime2 format to return all employees that had a **ManagerID** of **3** on a specific date:

The AS OF Subclause

```
SELECT *
FROM dbo.Employee
FOR SYSTEM_TIME AS OF '2015-12-21 09:00:00'
WHERE ManagerID = 3
```

Using the FROM TO subclause returns both current and historical data between two dates. The following example returns all employees that existed in a six-month period, along with all changes to their records:

The FROM TO Subclause

```
SELECT *
FROM dbo.Employee
FOR SYSTEM_TIME FROM '2015-07-01' TO '2015-12-31'
ORDER BY EmployeeID, SysStartTime
```

The FROM TO subclause returns all current and historical rows that were active during the timespan, regardless of whether they were active before or after those times. The results will include rows that were active precisely on the lower boundary defined by the FROM date; however, this excludes rows that became inactive on the upper boundary defined by the TO date.

The BETWEEN AND and CONTAINED IN subclauses work in much the same way as FROM TO, but each has a subtle difference:

- The BETWEEN AND subclause is identical to FROM TO, except that it includes rows that were active on the upper time boundary.
- The CONTAINED IN subclause returns a table with the values for all rows that were opened and closed within the specified timespan. It includes rows that became active exactly on the lower boundary or became inactive exactly on the upper boundary.

If your queries are not returning the results you expect, check the subclause in case you are including or excluding rows because of the rules stated above.



Note: If you want to search historical data only, use the CONTAINED IN subclause because this gives the best performance.

Demonstration: Creating System-Versioned Tables

Demonstration Steps

Create System-Versioned Tables

1. In the **D:\Demofiles\Mod09\Temporal** folder, run **SetupC.cmd** as Administrator.
2. In the **User Account Control** dialog box, click **Yes**.
3. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the database engine using Windows authentication.
4. Open the **Demo.ssmssl** solution in the **D:\Demofiles\Mod09\Temporal\Demo** folder.
5. In Solution Explorer, in the **Queries** folder, open the **6 – Temporal Tables.sql** script file.
6. Select the Transact-SQL code under the comment **Step 1: Connect to AdventureWorks**, and then click **Execute**.
7. Select the Transact-SQL code under the comment **Step 2: Create System-Versioned Table**, and then click **Execute**.
8. In Object Explorer, expand **Databases**, expand **AdventureWorks**, expand **Tables**, point out that System-Versioned tables are tagged with the **System-Versioned** marker. Expand the **HumanResources.EmployeeBenefit** table to show the history table located under the main table node. Expand the columns to show the tables are the same, other than the PK.
9. Select the Transact-SQL code under the comment **Step 4 - Alter Existing Table**, and then click **Execute**.
10. In Object Explorer, right-click **Tables** and then click **Refresh**.
11. Expand the **Sales.SalesTaxRate** table to show the history table located under the main table node. Expand the columns to show the tables are the same, other than the PK.
12. Select the Transact-SQL code under the comment **Step 6 - Show Dates**, and then click **Execute**.
13. Close SQL Server Management Studio. If prompted, do not save changes.

Using System-Versioned Tables to Implement Slowly Changing Dimensions

Slowly Changing Dimensions (SCD) change unpredictably and slowly, as the name suggests.

SQL Server contains functionality for managing SCD.

SCD Functionality

SQL Server has the following features for handling SCD using Temporal Tables:

- **Changed Incoming Records Management.** Incoming rows with changes that require new rows and also where expired rows need refreshing.
- **Changed Historical Records Management.** Incoming rows with historical changes that require new rows and also where expired history rows need refreshing.
- **Matching Incoming Records.** Identify which rows in the target tables need updating and creating.
- **Identify Incoming Records with Disallowed Changes.** Identify which rows contain unpermitted changes.
- **Inferred Member Record Management.** Inferred members exist when dimension data is not fully up to date—that is, not loaded. These rows are identified and managed by the SCD feature.

SCD Change Types and Transformation Outputs

SCD transformation outputs are used to direct incoming rows to the correct mapping for the SCD change type. These outputs are configured in the SCD wizard as described below.

SCD supports four types of changes:

- **Changing Attribute (Type 1 Change).** Overwrite existing rows. Incoming data is directed to the **Changing Attributes Updates Output**.
- **Historical Attribute (Type 2 Change).** Create new rows. The only change allowed is that a row can be “expired”. Incoming data is directed to the **Historical Attribute Insert Output** and **New Output**.
- **Fixed Attribute.** The row value must not be changed. Incoming data is directed to the **Fixed Attribute Insert Output**.
- **Inferred Member.** Create a minimal row so that the relevant dimension data can be loaded later. The SCD transformation data is directed to the **Inferred Member Output**. An update is made when the data is loaded rather than inserted.



Slowly Changing Dimension Transformation

<http://aka.ms/it2upg>

Slowly Changing Dimensions (SCD):

- SCD manages incoming data:
 - Changed, historical, fixed, disallowed, inferred
- SCD change types and transformation outputs:
 - Changing attributes updates
 - Historical attribute insert output and new
 - Fixed attribute inserts
 - Inferred member
- The SCD wizard:
 - Choose data source and dimension table
 - Configure mapping
 - Set attribute options
 - Review, run and update

The Slowly Changing Dimension Wizard

Configuring and building SCD data flows can be quite difficult and time consuming. The SCD wizard simplifies the configuration and building of SCD data flows.

The following steps describe how to use the SCD wizard:

1. **Data Source.** Choose the data source that contains the dimension table to be updated.
2. **Dimension.** Select the dimension table or view.
3. **Mapping.** Configure the key attributes and columns to be mapped.
4. You choose the change type for each column:
 - **Changing Attribute.** Overwrites values.
 - **Historical Attribute.** Creates new rows.
 - **Fixed Attribute.** The column value must not be amended.
5. **Set Fixed and Changing Attributes.** You can set various rules, such as with changing attributes, whether or not all matching records (including expired) are to be updated.
6. **Set Historical Attribute Options.** You must choose how SCD will differentiate between current and expired records. The following list describes the options:
 - **Current.** Select **Current, True**.
 - **Expired.** Select **Expired, or False**.
7. **Inferred Members Support.** You can set SCD to create minimal records for inferred members. Two types of inferred members can be created. Rows in which:
 - All columns with change types are NULL.
 - A Boolean column can be set to mark the row as an inferred member.
8. **Review.** You can review the configurations before you run the SCD.
9. **Updates.** You can update the SCD transformation using the SCD wizard, the Advanced Editor, or by using Transact-SQL.

Using the SCD Dimension Wizard

<http://aka.ms/au4plu>

Change Data Capture Compared to System-Versioned Tables

The temporal feature does not replace CDC—they have different uses and are deployed and used in different ways:

- **Usage.** CDC is best used for managing changes for short periods—for example, where short-term retention is required to protect data in an ETL process. System-Versioned tables are used to retain historical data for much longer periods—for example, for audits, SCD, historical data comparison or fast recovery of corrupted records.

Parameter	CDC	System-Versioning
Usage	Rapid data change, short retention; ETL, data checks, fact tables, Transact-SQL options, SSMS, and Visual Studio	Slower data change, long retention, SCD, audit, fix corrupt data, period reporting, and comparison, Transact SQL and SSMS
Deployment	Database and table levels	Table level
Methodology	Transaction logs many tables	One historical table per source table
Scope	Down to row/column level	Source table tracked at row level
Performance	Handles large data amounts quickly	Large data amounts slower
Dependencies	SQL Server Agent	None
Operations	Can be quickly enabled/disabled at database level	Requires each table to be enabled/disabled
Operational	Easy to enable/disable for maintenance, and so on	More difficult
Primary Tables	No extra columns in primary	No extra columns in primary
Primary Keys	Not required	Required



Note: CDC Tasks are also part of SQL Server Visual Studio.

Deployment. CDC is enabled at both database and table level. System-versioning is deployed at table level.

- **Methodology.** CDC reads data from transaction logs and copies them to a CDC specific table. For each System-Versioned table, a historical table is created in which all changes are stored by start and end date. CDC requires CDC schema and users—this concept is not required for system-versioning.
- **Scope.** CDC tracks down data to table or column level and can handle changes to schema. System-versioning applies to entire tables.
- **Performance Issues.** CDC can produce large amounts of historical data very quickly—for example in an ETL process. However, this data might only be required for a short time. System-versioning is more likely to impact performance when many changes are being made quickly to many tables. However, this is offset if historical tables can be stored on a different device.
- **Dependencies.** CDC requires SQL Server Agent to be up and running. If it is off, CDC does not operate. System-versioning does not require SQL Server Agent, so will work as long as SQL Server is operational.
- **Operational Issues.** CDC is easy to switch on and off at database level using the **sys.sp_cdc_disable_db** and **sys.sp_cdc_enable_db** stored procedures. System-versioning can only be switched off one table at a time using Transact-SQL. System-versioning also requires Transact-SQL statements to disable a table for a temporary period or permanently.
- **Primary Tables.** Neither CDC nor System-Versioned tables require extra columns in the primary tables.
- **Primary Keys.** CDC enabled tables do not require a primary key. System-Versioned tables require a primary key to provide a relationship between the source and historical data.

Question: When you create a System-Versioned table, what is the minimum SQL Server creates in terms of tables, columns and keys?

Check Your Knowledge

Question	
Which of these options is not a subclause of FOR SYSTEM_TIME?	
Select the correct answer.	
<input type="checkbox"/>	AS OF <date_time>
<input type="checkbox"/>	BEFORE <date_time>
<input type="checkbox"/>	FROM <start_date_time> TO <end_date_time>
<input type="checkbox"/>	BETWEEN <start_date_time> AND <end_date_time>
<input type="checkbox"/>	CONTAINED IN (<start_date_time>, <end_date_time>)

Lab B: Loading a Data Warehouse

Scenario

You are ready to start developing the SSIS packages that load data from the staging database into the data warehouse.**Objectives**

After completing this lab, you will be able to:

- Load data from CDC output tables.
- Use a Lookup transformation to load data.
- Use the SCD transformation.
- Use the MERGE statement.

Estimated Time: 60 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Loading Data from CDC Output Tables

Scenario

The staging database in your ETL solution includes tables named as follows:

- **EmployeeInserts**: containing employee records that have been inserted in the employee source system.
- **EmployeeUpdates**: containing records modified in the employee source system.
- **EmployeeDeletes**: containing records that have been deleted in the employee source system.

You must use these tables to load and update the **DimEmployee** dimension table, which uses a **Deleted** flag to indicate records that have been deleted in the source system.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Create a Data Flow for Inserts
3. Create a Data Flow for Updates
4. Create a Data Flow for Deletes
5. Test the Package

► Task 1: Prepare the Lab Environment

1. Ensure that the 20767C-MIA-DC and 20767C-MIA-SQL virtual machines are both running, and then log on to MIA-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **SetupB.cmd** in the **D:\Labfiles\Lab09\Starter** folder as Administrator.

► Task 2: Create a Data Flow for Inserts

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex5** folder.
2. In Solution Explorer, note that a connection manager for the **AWDataWarehouse** database has been created.
3. Add a new SSIS package named **Load Employee Data.dtsx** to the project.
4. Add a data flow task named **Insert Employees** to the control flow.
5. In the **Insert Employees** data flow, create an OLE DB source named **Staged Employee Inserts** that extracts data from the **[dbo].[EmployeeInserts]** table in the **Staging** database.
6. Connect the **Staged Employees** source to a new OLE DB destination named **New Employees** that uses a fast load option to load data into the **DimEmployee** table in the **AWDataWarehouse** database.

► Task 3: Create a Data Flow for Updates

1. On the control flow surface of the **Load Employee Data.dtsx** package, connect the success precedence constraint of the **Insert Employees** data flow task to a new data flow task named **Update Employees**.
2. In the **Update Employees** data flow, create an OLE DB source named **Staged Employee Updates** that extracts data from the **[dbo].[EmployeeUpdates]** table in the **Staging** database.
3. Connect the data flow from the **Staged Employee Updates** source to a new OLE DB Command transformation named **Update Existing Employees** that executes the following Transact-SQL statement in the **AWDataWarehouse** database:

```
UPDATE dbo.DimEmployee
SET FirstName = ?, LastName = ?, EmailAddress = ?, Title = ?, HireDate = ?
WHERE EmployeeAlternateKey = ?
```

Use the following column mappings:

- **FirstName:** Param_0
- **LastName:** Param_1
- **EmailAddress:** Param_2
- **Title:** Param_3
- **HireDate:** Param_4
- **EmployeeID:** Param_5

► Task 4: Create a Data Flow for Deletes

1. On the control flow surface of the **Load Employee Data.dtsx** package, connect the success precedence constraint of the **Update Employees** data flow task to a new data flow task named **Delete Employees**.
2. In the **Delete Employees** data flow, create an OLE DB source named **Staged Employee Updates** that extracts data from the **[dbo].[EmployeeDeletes]** in the **Staging** database.

3. Connect the data flow from the **Staged Employee Deletes** source to a new OLE DB Command transformation named **Delete Existing Employees** that executes the following Transact-SQL statement in the **AWDataWarehouse** database, mapping the **EmployeeID** column to **Param_0**:

```
UPDATE dbo.DimEmployee
SET Deleted = 1
WHERE EmployeeAlternateKey = ?
```

► Task 5: Test the Package

1. In Visual Studio, start debugging the **Load Employee Data.dtsx** package. When execution is complete, view the data flow surface for each of the data flow tasks, noting the numbers of rows processed in each task.
2. When you have finished, stop debugging and close Visual Studio.

Results: After this exercise, you should have an SSIS package that uses data flows to apply inserts, updates, and logical deletes in the data warehouse, based on staging tables extracted by the CDC Control Task and data flow components.

Exercise 2: Using a Lookup Transformation to Insert or Update Dimension Data

Scenario

Another BI developer has partially implemented an SSIS package to load product data into a hierarchy of dimension tables. You must complete this package by creating a data flow that uses a Lookup transformation to determine whether a product dimension record already exists, and then insert or update a record in the dimension table accordingly.

The main tasks for this exercise are as follows:

1. View Data Flows
2. Create a Data Flow
3. Add a Lookup Transformation for Parent Keys
4. Configure a Lookup Transformation
5. Add a Destination for New Products
6. Add an OLE DB Command for Updated Product Records
7. Test the Package

► Task 1: View Data Flows

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex6** folder. Then open the **Load Products Data.dtsx** SSIS package.
2. View the data flow for the **Load Product Category Dimension** task, and note the following:
 - The **Staged Product Category Data** source extracts product category data from the **InternetSales** and **ResellerSales** tables in the **Staging** database.
 - The **Lookup Existing Product Categories** task retrieves the **ProductCategoryKey** value for product categories that exist in the **DimProductCategory** table in the **AWDataWarehouse**

database by matching the **Product Category** business key in the **Staging** database to the **Product Category** alternative key in the data warehouse.

- The **Lookup No Match Output** data flow path from the **Lookup Existing Product Categories** task connects to the **New Product Categories** destination, and the **Lookup Match Output** data flow path connects to the **Update Existing Product Categories** task.
 - The **New Product Categories** destination loads new product category records into the **DimProductCategory** table.
 - The **Update Existing Product Categories** task executes a Transact-SQL statement to update the **ProductCategoryName** column in the **DimProductCategory** table for an existing row based on the **ProductCategoryKey**.
3. View the data flow for the **Load Product Subcategory Dimension** task, and note that this data flow inserts or updates product category dimension data using a similar approach to the **Load Product Category Dimension** data flow. Additionally, it has a lookup task to retrieve the **ProductCategoryKey** in **AWDataWarehouse** for the parent category, which should have already been loaded.

► Task 2: Create a Data Flow

1. Add a data flow task named **Load Product Dimension** to the control flow of the **Load Products Data.dtsx** package, and connect the success precedence constraint from the **Load Product Subcategory Dimension** task to the **Load Product Dimension** task.
2. In the data flow for the **Load Product Dimension** data flow task, add an OLE DB source named **Staged Product Data** that uses the following Transact-SQL command to retrieve data from the **Staging** database:

```
SELECT DISTINCT ProductSubcategoryBusinessKey, ProductBusinessKey, ProductName,
StandardCost, Color, ListPrice, Size, Weight, Description
FROM dbo.InternetSales
UNION
SELECT DISTINCT ProductSubcategoryBusinessKey, ProductBusinessKey, ProductName,
StandardCost, Color, ListPrice, Size, Weight, Description FROM dbo.ResellerSales
```

► Task 3: Add a Lookup Transformation for Parent Keys

1. In the **Load Product Dimension** data flow, add a Lookup transformation named **Lookup Parent Subcategory**, connect the output data flow path from the **Staged Product Data** source to it, and configure it as follows:
 - The component should fail if there are rows with no matching entries.
 - The components should look up rows in the **[dbo].[DimProductSubcategory]** table in the **AWDataWarehouse** database by matching the **ProductSubcategoryBusinessKey** column to the **ProductSubcategoryAlternateKey** lookup column.
2. Each matching lookup should return the **ProductSubCategoryKey** lookup column and add it to the data flow.

► Task 4: Configure a Lookup Transformation

1. In the **Load Product Dimension** data flow, add another Lookup transformation named **Lookup Existing Products**, connect the **Lookup Match Output** data flow path from the **Lookup Parent Subcategory** transformation to it.

2. Configure the new Lookup transformation as follows:
 - Rows with no matching entries should be redirected to the no match output.
 - The components should look up rows in the **[dbo].[DimProduct]** table in the **AWDataWarehouse** database by matching the **ProductBusinessKey** column to the **ProductAlternateKey** lookup column.
 - Each matching lookup should return the **ProductKey** lookup column and add it to the data flow.

► **Task 5: Add a Destination for New Products**

1. In the **Load Product Dimension** data flow, connect the **Lookup No Match** output from the **Lookup Existing Products** transformation to a new OLE DB destination named **New Products** that loads unmatched product records into the **DimProduct** table in the **AWDataWarehouse** database.
2. Use the following column mappings when loading the data:
 - **<ignore>**: ProductKey
 - **ProductBusinessKey**: ProductAlternateKey
 - **ProductName**: ProductName
 - **StandardCost**: StandardCost
 - **Color**: Color
 - **ListPrice**: ListPrice
 - **Size**: Size
 - **Weight**: Weight
 - **Description**: Description
 - **<ignore>**: ProductSubcategoryKey

► **Task 6: Add an OLE DB Command for Updated Product Records**

1. In the **Load Product Dimension** data flow, connect the **Lookup Match Output** data flow path from the **Lookup Existing Products** transformation to a new OLE DB Command transformation named **Update Existing Products**.
2. Configure the **Update Existing Products** transformation to use the following Transact-SQL command to update the **DimProduct** table in the **AWDataWarehouse** database:

```
UPDATE dbo.DimProduct
SET ProductName = ?, StandardCost = ?, Color = ?, ListPrice = ?, Size = ?,
Weight = ?, Description = ?, ProductSubcategoryKey = ?
WHERE ProductKey = ?
```

3. Use the following column mappings for the query:
 - **ProductName**: Param_0
 - **StandardCost**: Param_1
 - **Color**: Param_2
 - **ListPrice**: Param_3
 - **Size**: Param_4
 - **Weight**: Param_5

- **Description:** Param_6
- **ProductSubcategoryKey:** Param_7
- **ProductKey:** Param_8

► Task 7: Test the Package

Note: When debugging, the rows do not always appear on the data flow path. An alternate check is to open the database and ensure the changes have been successful.

1. With the **Load Product Dimension** data flow visible, start debugging the package and verify that all rows flow to the **New Products** destination (because the data warehouse contained no existing product records). When package execution is complete, stop debugging.
2. Debug the package again and verify that all rows flow to the **Update Existing Products** transformation this time (because all staged product records were loaded to the data warehouse during the previous execution, so they all match existing records). When package execution is complete, stop debugging and close Visual Studio.

Results: After this exercise, you should have an SSIS package that uses a Lookup transformation to determine whether product records already exist, and updates them or inserts them as required.

Exercise 3: Implementing a Slowly Changing Dimension

Scenario

You have an existing SSIS package that uses an SCD transformation to load reseller dimension records into a data warehouse. You want to examine this package and then create a new one that uses an SCD transformation to load customer dimension records into the data warehouse.

The main tasks for this exercise are as follows:

1. Execute a Package to Load a Non-Changing Dimension
2. Observe a Data Flow for a Slowly Changing Dimension
3. Implement a Slowly Changing Dimension Transformation
4. Test the Package

► Task 1: Execute a Package to Load a Non-Changing Dimension

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex7** folder.
2. Open the **Load Geography Data.dtsx** package and review the control flow and data flow defined in it. This package includes a simple data flow to load staged geography data into the data warehouse. Only new rows are loaded, and rows that match existing data are discarded.
3. Start debugging and observe the package execution as it loads geography data into the data warehouse.
4. When package execution has completed, stop debugging.

► Task 2: Observe a Data Flow for a Slowly Changing Dimension

1. Open the **Load Reseller Data.dtsx** SSIS package.
2. Examine the data flow for the **Load Reseller Dimension** task, and note the following features:
 - The **Staged Reseller Data** source extracts data from the **Resellers** table in the **Staging** database.
 - The **Lookup Geography Key** transformation looks up the geography key for the reseller in the **DimGeography** table in the **AWDataWarehouse** database.
 - The **Reseller SCD** is an SCD transformation that has generated the remaining transformations and destinations. You can double-click the **Reseller SCD** transformation to view the wizard used to configure the SCD, and then click **Cancel** to avoid making any unintentional changes.
 - The **Reseller SCD** transformation maps the **ResellerBusinessKey** input column to the **ResellerAlternateKey** dimension column and uses it as a business key to identify existing records.
 - The **Reseller SCD** transformation treats **AddressLine1**, **AddressLine2**, **BusinessType**, **GeographyKey**, and **NumberEmployees** as historical attributes, **Phone** and **ResellerName** as changing attributes, and **YearOpened** as a fixed attribute.
3. Start debugging and observe the data flow as the dimension is loaded.
4. When package execution is complete, stop debugging.

► Task 3: Implement a Slowly Changing Dimension Transformation

1. Open the **Load Customer Data.dtsx** SSIS package and view the data flow for the **Load Customer Dimension** task. Note that the data flow already contains a source named **Staged Customer Data**, which extracts customer data from the **Staging** database, and a Lookup transformation named **Lookup Geography Key**, which retrieves a **GeographyKey** value from the **AWDataWarehouse** database.
2. Add an SCD transformation named **Customer SCD** to the data flow and connect the **Lookup Match Output** data flow path from the **Lookup Geography Key** transformation to the **Customer SCD** transformation.
3. Use the SCD wizard to configure the data flow to load the **DimCustomer** table in the **AWDataWarehouse** database.
 - Map input columns to dimension columns with the same name.
 - Map the **CustomerBusinessKey** input column to the **CustomerAlternateKey** dimension column, and use this column as the business key.
 - Do not map the **CurrentRecord** dimension column to any input column.
 - Specify the following slowly changing dimension columns:

Dimension Columns	Change Type
AddressLine1	Historical attribute
AddressLine2	Historical attribute
BirthDate	Changing attribute

Dimension Columns	Change Type
CommuteDistance	Historical attribute
EmailAddress	Changing attribute
FirstName	Changing attribute
Gender	Historical attribute
GeographyKey	Historical attribute
HouseOwnerFlag	Historical attribute
LastName	Changing attribute
MaritalStatus	Historical attribute
MiddleName	Changing attribute
NumberCarsOwned	Historical attribute
Occupation	Historical attribute
Phone	Changing attribute
Suffix	Changing attribute
Title	Changing attribute

- Use the **CurrentRecord** column to show current and expired records, using the value **True** for current records and **False** for expired records.
4. Do not enable inferred member support.

► Task 4: Test the Package

Note: When debugging, the rows do not always appear on the data flow path. An alternative check is to open the database and ensure the changes have been successful.

1. Debug the package and verify that all rows pass through the **New Output** data flow path. When package execution is complete, stop debugging.
2. Debug the package again and verify that no rows pass through the **New Output** data flow path, because they already exist and no changes have been made. When package execution is complete, stop debugging.
3. Use SQL Server Management Studio to execute the **Update Customers.sql** script in the localhost instance of the database engine. This script updates two records in the **Staging** database, changing one customer's phone number and another's marital status.

4. In Visual Studio, debug the package again and verify that one row passes through the **Historical Attribute Inserts Output** data flow path, and another passes through the **Changing Attributes Updates Output**. When package execution is complete, stop debugging.
5. Close Visual Studio.

Results: After this exercise, you should have an SSIS package that uses an SCD transformation to load data into a dimension table.

Exercise 4: Using the MERGE Statement

Scenario

Your **Staging** database is located on the same server as the data warehouse. You want to take advantage of this collocation of data and use the MERGE statement to insert and update staged data into the Internet sales fact table. An existing package already uses this technique to load data into the reseller sales fact table.

The main tasks for this exercise are as follows:

1. Examine a Control Flow That Uses the MERGE Statement
2. Create a Package That Uses the MERGE Statement
3. Test the Package

► Task 1: Examine a Control Flow That Uses the MERGE Statement

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the D:\Labfiles\Lab09\Starter\Ex8 folder then open the **Load Reseller Sales Data.dtsx** SSIS package.
2. Examine the configuration of the **Merge Reseller Sales** task and note the following details:
 - The task uses the **localhost.Staging** connection manager to connect to the **Staging** database.
 - The task executes a Transact-SQL MERGE statement that retrieves reseller sales and related dimension keys from the **Staging** and **AWDataWarehouse** databases. It then matches these records with the **FactResellerSales** table, based on the **SalesOrderNumber** and **SalesOrderLineNumber** columns, updates rows that match, and inserts new records for rows that do not.
3. Start debugging to run the package and load the reseller data. When package execution is complete, stop debugging.

► Task 2: Create a Package That Uses the MERGE Statement

1. Add a new SSIS package named **Load Internet Sales Data.dtsx**.
2. Add an Execute SQL task named **Merge Internet Sales Data** to the control flow of the **Load Internet Sales Data.dtsx** package.
3. Configure the **Merge Internet Sales Data** task using the **localhost.Staging** connection manager and execute a MERGE statement that retrieves Internet sales and related dimension keys from the **Staging** and **AWDataWarehouse** databases. It also matches these records with the **FactInternetSales** table, based on the **SalesOrderNumber** and **SalesOrderLineNumber** columns, updates rows that match, and inserts new records for rows that do not. To accomplish this, you can use the code in the **Merge Internet Sales.sql** script file in the D:\Labfiles\Lab09\Starter\Ex8 folder.

► **Task 3: Test the Package**

1. Start debugging the package, observing the execution of the **Merge Internet Sales Data** task. When execution is complete, stop debugging.
2. Close Visual Studio.

Results: After this exercise, you should have an SSIS package that uses an Execute SQL task to execute a MERGE statement that inserts or updates data in a fact table.

Module Review and Takeaways

In this module, you have learned how to plan and implement an ETL solution that extracts and loads data incrementally from data sources.

Review Question(s)

Question: What should you consider when choosing between Change Data Capture and Change Tracking?

Question: What should you consider when deciding whether to use the MERGE statement to load staging data into a data warehouse?

Module 10

Enforcing Data Quality

Contents:

Module Overview	10-1
Lesson 1: Introduction to Data Quality	10-2
Lesson 2: Using Data Quality Services to Cleanse Data	10-11
Lab A: Cleansing Data	10-15
Lesson 3: Using Data Quality Services to Match Data	10-21
Lab B: Deduplicating Data	10-27
Module Review and Takeaways	10-31

Module Overview

Ensuring the high quality of data is essential if the results of data analysis are to be trusted. SQL Server® includes Data Quality Services (DQS) to provide a computer-assisted process for cleansing data values, in addition to identifying and removing duplicate data entities. This process reduces the workload of the data steward to a minimum while maintaining human interaction to ensure accurate results.

Objectives

After completing this module, you will be able to:

- Describe how DQS can help you manage data quality.
- Use DQS to cleanse your data.
- Use DQS to match data.

Lesson 1

Introduction to Data Quality

Data quality is a major concern for anyone building a data warehousing solution. In this lesson, you will learn about the types of data quality issues that must be addressed in a data warehousing solution—and how SQL Server Data Quality Services (DQS) can help.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the need for data quality management.
- Describe the features and components of DQS.
- Describe the features of a knowledge base.
- Describe the features of a domain.
- Explain how reference data can be used in a knowledge base.
- Create a DQS knowledge base.

What Is Data Quality and Why Do You Need It?

As organizations consume larger volumes of data from more sources, the need for data quality management has become increasingly common. Data quality is especially important in a Business Intelligence (BI) solution, because the reports and analysis generated from data in the data warehouse can form the basis of important decisions. To make these decisions, business users must be able to trust the data they use.

- Business decisions should be made on trusted data
- Data quality issues in sources can be propagated into the data warehouse:
 - Invalid data values
 - Inconsistencies
 - Duplicate business entities

Data Quality Issues

Common data quality issues include the following categories:

- **Invalid Data Values:** for example, an organization might categorize its stores as “wholesale” or “retail”. However, a user might have an application that allows free-form data entry to create a store with a category of “reseller” instead of “retail”, or they might accidentally type “wholesale” instead of “wholesale”. Any analysis or reporting that aggregates data by store type will then produce inaccurate results because of the additional, invalid categories.
- **Inconsistencies:** for example, an organization might have one application for managing customer accounts in which US states are stored using two-letter codes (such as “WA” for Washington), and a second application that stores supplier addresses with a full state name (such as “California”). When data from both these systems is loaded, the data warehouse will contain inconsistent values for states.
- **Duplicate Business Entities:** for example, a customer relationship management system might contain records for Jim Corbin, Jimmy Corbin, James Corbin, and J Corbin. If the address and telephone number for these customers are all the same, it might be reasonable to assume that all the records relate to the same customer. Of course, it’s also possible that Jim Corbin has a wife named

Jennifer and a son called James—so you must be confident that you have matched the records appropriately before deduplicating the data.

Using the SSIS DQS Cleansing Component

<http://aka.ms/iajqdh>

Data Quality Services Overview

DQS is a knowledge-based solution for managing data quality. With DQS, you can perform the following kinds of data quality management:

- **Data Cleansing:** identify invalid or inconsistent data values and correct them.
- **Data Matching:** find duplicate data entities.

DQS is installed from the SQL Server installation media, and consists of the following components:

- **Data Quality Services Server:** a service that uses a knowledge base to apply data quality rules to data. The server must be installed on the same instance as the data that you wish to analyze. With two SQL Server catalogs installed, you can monitor, maintain, back up, and perform other administrative tasks from within SQL Server Management Studio. DQS_MAIN includes stored procedures, the DQS engine, and published knowledge bases. DQS_PROJECT includes data required for knowledge base management and data quality project activities.
- **Data Quality Client:** a wizard-based application that data stewards (typically business users) can use to create and manage data quality services knowledge bases and perform data quality services tasks. The client can either be installed on the same computer as the DQS server, or be used remotely.
- **Data Cleansing SSIS Transformation:** a data flow transformation for SQL Server Integration Services (SSIS) that you can use to cleanse data as it passes through a data flow pipeline.

- DQS is a knowledge-based solution for:
 - Data cleansing
 - Data matching
- DQS Components:
 - Server
 - Client
 - Data cleansing SSIS transformation

What Is a Knowledge Base?

DQS helps you improve data quality by creating a knowledge base for the data, then applying the rules it contains to perform data cleansing and matching. A knowledge base stores all the knowledge related to a specific aspect of the business. For example, you could maintain one knowledge base for a customer database and another for a product database.

Each knowledge base contains:

- Domains that define valid values and correction rules for data fields.
- Matching policies that define rules for identifying duplicate data entities.

- Repository of knowledge about data:
 - Domains define values and rules for each field
 - Matching policies define rules for identifying duplicate records
- Determine data for a DQS knowledge base:
 - Analyze source databases and data warehouses for inconsistencies, inaccuracies, and incompleteness
 - Audit website and software forms used for data entry to find free-form fields prone to creating low quality data
 - Look at dependent reporting systems and find incorrect results

Knowledge bases are usually created and maintained by data stewards—they are often business users who have particular expertise in a specific area.

DQS provides a basic knowledge base that includes domains for US address data, such as states and cities. You can use it to learn about data quality services and as a starting point for your own knowledge bases.

Determining the Knowledge Base Data

To determine which data requires a DQS knowledge base, you can use one or both approaches:

- Analyze the data in your source databases and data warehouses for inconsistencies, inaccuracies, and incompleteness.
- Audit the forms used on websites and internal software systems to identify free-form input fields that have no validation and are prone to accepting inaccurate data.
- Look at dependent reporting solutions to see where data is missing or inaccurate, thereby causing incorrect results.

When analyzing data, use queries to look for the following issues:

- **Missing data:** look for NULL values or empty strings in columns where there should be data. This can help highlight data that could be causing invisible issues in reports. Columns in dimensions that are missing a value might cause rows to be excluded from aggregations, thereby producing incorrect results.
- **Duplication:** identify where data has been duplicated inside the same table, across tables, or across databases. Has data been changed in one table or database, but not another?
- **Consistency:** find values that have the same meaning but are stored differently. For example, gender values may be stored as M or F in one database, and Male or Female in another.
- **Accuracy:** determine what accuracy means in your data. Do you want to include customers who have not bought a product or service for more than 12 or 24 months?
- **Validity:** identify incorrect data, such as dates falling outside of acceptable ranges. Are there customers in your database with a birthday less than 18 years ago, when all customers must be 18 or over to order goods or services?

The ETL packages used to load data into your data warehouse are a good starting point for sorting through your data. Look at the data that is brought into staging and identify any values that are missing, or data that might be inconsistent between different source systems, or needs cleaning, or deduplicating. Including a data steward in this process can help determine the business rules for your data that will be applied in your DQS knowledge base.

It is also useful to review the code behind the forms on your website or internal software systems that are used for data entry. Ensuring data is entered correctly at the beginning helps reduce cleansing efforts later. For example, providing select lists instead of free form entry fields prevents inconsistency. It is impossible to completely prevent all data entry mistakes, but reducing the potential delivers higher quality data.

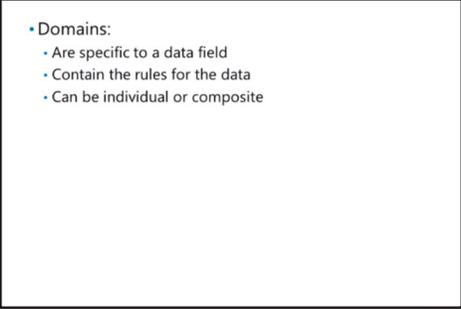
Users with a deep understanding of the business who know the results they expect to see in reports can help identify incorrect data. Perhaps they know that the sales figures for one of the retail stores is lower than it should be. However, by looking at the data, you might uncover the fact that some salespeople regularly fill in a form incorrectly, omitting data that is required to be included in aggregations. By working back through the queries that produce the results in a report, you can start to identify rows that have missing or inconsistent data. However, this could be a slower process if the queries run on very large datasets and take a long time to return a result.

After you have identified the data that requires some form of cleaning, you can begin to build your DQS knowledge base, working alongside a data steward to create rules. When the data loaded into your data warehouse is in a good shape, you can widen your scope of analysis to review other data in your systems. Customer data might be extracted from the main OLTP system and copied to a marketing database for mailings and promotions. Ensuring this data is clean and free of inaccuracies helps retain customer satisfaction by ensuring names and address are correct, and prevents duplicate correspondence being sent out.

What Is a Domain?

Domains are central to a DQS knowledge base. Each domain identifies the possible values and rules for a data field. The values for each domain are categorized as:

- **Valid:** for example, valid values for a US state domain may include "California" or "CA".
- **Invalid:** for example, invalid values for a US state domain may include "8".
- **Error:** for example, a common error for a US state domain may be "California" (with a missing "i").

- 
- Domains:
 - Are specific to a data field
 - Contain the rules for the data
 - Can be individual or composite

Values can be grouped as synonyms. For example, you might group "California", "CA", and "California" as synonyms for California. You can specify a leading value to which all synonyms are corrected. For example, you could configure the domain so that instances of "CA" and "California" are automatically corrected to "California".

In addition to defining the values for a domain, you can create domain rules that validate new data values. For example, you could create a rule to ensure that all values in an **Age** domain are numbers or that all values in an **Email Address** domain include a "@" character. You can also specify standardization settings for a text-based domain to enforce correct capitalization. This means you can ensure that cleansed text values have consistent formats.

Often, you can create domains to represent the most granular level of your data—for example, **First Name**—but the actual unit of storage consists of multiple domains—for example, **Full Name**. In this instance, you can combine the **First Name** and **Last Name** domains to form a **Full Name** composite domain. Composite domains are also used for address fields made up of a combination of address, city, state, postal code, and country data. Another use of composite domains is a rule that combines data from multiple domains. For example, you can verify that the string "98007" in a **Postal Code** domain corresponds to the string "Bellevue" in a city domain.

Matching can be performed on the individual domains that comprise the composite domain, but not on the composite domain itself.

What Is a Reference Data Service?

Many data quality problems are outside the core specialization of the organization in which they are used. For example, your organization might be an Internet retailer that ships goods based on incorrect address data, a core data problem that produces unnecessary costs. Your website could be as user-friendly as possible, but there might still be an unacceptably high number of incorrectly addressed orders.

- The Azure Marketplace hosts specialist data cleansing providers, where you can:
 - Set up an account
 - Subscribe to a reference service
 - Map your domain to the reference service

Reference Data Services

To cleanse data that is outside the knowledge of your organization, you can subscribe to third-party Reference Data Service (RDS) providers. By using the Windows® Azure® Data Market, it is straightforward to subscribe to an RDS service and use it to validate and cleanse your data.

To continue the example, by using the Windows Azure Data Market, you could purchase a subscription to an address verification service. You can then send data there for it to be verified and cleansed, to reduce incorrect address information and, therefore, cut your postage costs.

To use RDS to cleanse your data, you must follow these steps:

1. Create a free account key at the Windows Azure Marketplace.
2. Subscribe to a free or paid-for RDS provider's service at the Marketplace.
3. Configure the RDS details in DQS.
4. Map your domain to the RDS.
5. Finally, you can use the knowledge base containing the domain that maps to the RDS to cleanse the data.

One of the key advantages of using the Azure Data Market to provide DQS services is that, typically, the cost of the data service is based on the number of times you use it per month. This means you can scale up at busy times and reduce costs when the business is quieter.



Reference Data Services in DQS

<http://aka.ms/y7anf2>

Creating a Knowledge Base

Building a DQS knowledge base is an iterative process that involves the following steps:

1. **Knowledge Discovery:** use existing data to identify domain values.
2. **Domain Management:** categorize discovered values as valid, invalid, or errors; specify synonyms and leading values; define correction rules and other domain configuration tasks.

- Creating a knowledge base is an iterative process:
 1. Knowledge discovery
 2. Domain management

The data steward can create the initial knowledge base from scratch, center it on an existing one, or import one from a data file. The knowledge discovery process is then used to identify data fields that need to be managed, map these fields to domains in the knowledge base (which can be created during knowledge discovery if required), and identify values for these fields.

After the knowledge base is populated by the knowledge discovery process, the data steward manages the domains, to control how DQS validates and corrects data values. Additionally, domain management might include configuration of reference data services, or the setup of term-based or cross-field relationships.

Knowledge base development is an ongoing activity. A data steward will continually use the knowledge discovery and domain management processes to enhance the knowledge base and manage the quality of new data values and domains.



DQS Knowledge Bases and Domains

<http://aka.ms/o1reyr>

Demonstration: Creating a Knowledge Base

In this demonstration, you will see how to create a DQS knowledge base.

Demonstration Steps

Create a Knowledge Base

1. Ensure that the 20767C-MIA-DC20767C-MIA-CLI, and 20767C-MIA-SQL virtual machines are all running, and log into the 20767C-MIA-SQL virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod10** folder, right-click **Setup.cmd** and click **Run as administrator**.
3. When prompted to allow changes, click **Yes**.
4. Wait for the script to finish, and then press any key to close the command prompt.
5. Log into the 20767C-MIA-SQL virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
6. Click **Search Windows**, type **Data Quality Client**, and then press Enter
7. In the **Connect to Server** dialog box, type **MIA-SQL**, and then click **Connect**.

8. In SQL Server Data Quality Services, in the **Knowledge Base Management** section, click **New Knowledge Base**.
9. In the **Name** box, type **Demo KB**.
10. In the **Create Knowledge Base from** list, click **Existing Knowledge Base**.
11. In the **Select Knowledge Base** list, ensure **DQS Data** is selected.
12. In the **Select Activity** section, ensure **Domain Management** is selected, and then click **Next**.
13. In the Domain Management pane, select the **US - State** domain. On the **Domain Properties** tab, change the domain name to **State**.
14. On the **Domain Values** tab, note the existing values. The leading value for each state is the full state name. Other possible values that should be corrected to the leading value are indented beneath each leading value.
15. Click **Finish**, and then when prompted to publish the knowledge base, click **No**.

Perform Knowledge Discovery

1. In SQL Server Data Quality Services, under **Recent Knowledge Base**, click **Demo KB**, and then click **Knowledge Discovery**.
2. On the **Map** page, in the **Data Source** drop-down list, click **Excel File**.
3. In the **Excel File** box, click **Browse**, browse to **D:\Demofiles\Mod10**, and then double-click **Stores.xls**.
4. In the **Worksheet** drop-down list, ensure **Sheet1\$** is selected, and ensure **Use first row as header** is selected. This worksheet contains a sample of store data that needs to be cleansed.
5. In the first row of the **Mappings** table, in the **Source Column** list, click **State (String)**, and in the **Domain** list, click **State**.
6. In the second row of the **Mappings** table, in the **Source Column** list, click **City (String)**, and then click the **Create a domain** icon.
7. In the **Create Domain** dialog box, in the **Domain Name** box, type **City**, and then click **OK**.
8. Repeat the previous two steps to map the **StoreType (String)** source column to a new domain named **StoreType**.
9. On the **Map** page, click **Next**.
10. On the **Discover** page, click **Start** and wait for the knowledge discovery process to complete.
11. When the process has finished, note that 11 new **City** and **StoreType** records were found and that there were three unique **City** values, five unique **State** values, and four unique **StoreType** values. Click **Next**.
12. On the **Manage Domain Values** page, with the **City** domain selected, note the new values that were discovered.
13. Select the **State** domain and note that no new values were discovered.
14. Clear the **Show Only New** check box and note that all possible values for the **State** domain are shown. Note the **Frequency** column for **California, CA**, and **Washington, WA** indicates that the data was updated.
15. Select the **StoreType** domain and note the values discovered.

16. In the **Value** column, click **Retail**, hold the CTRL key and click **Resale**, and then click **Set selected domain values as synonyms**.
17. Right-click **Retail**, and click **Set as Leading**.
18. In the list of values, note that **Wholesale** has a red spelling checker line. Right-click **Wholesale**, and in the list of suggested spelling corrections, click **Wholesale**. Note that the **Type** for the **Wholesale** value changes to **Error** and the **Correct to** value is set automatically to **Wholesale**.
19. Click **Finish**.
20. When prompted to publish the knowledge base, click **No**.

Perform Domain Management

1. In SQL Server Data Quality Services, under **Recent Knowledge Base**, click **Demo KB**, and then click **Domain Management**.
2. In the **Domain** list, click **StoreType**.
3. On the **Domain Values** tab, note that the values discovered in the previous task are listed with appropriate leading values and correction behavior.
4. Click Add new domain value, and then enter the value Reseller.
Note that the Reseller value now has a star icon, as values were found in the last discovery run (hover over the star).
5. Click the **Retail** leading value, hold the CTRL key and click the new **Reseller** value. Click **Set selected domain values as synonyms** (depending on the screen resolution, this may be in a drop-down list at the end of the toolbar above the table). Note that **Reseller** becomes a valid value that is corrected to the **Retail** leading value.
6. Click **Finish**.
7. When prompted to publish the knowledge base, click **Publish**.
8. When publishing is complete, click **OK**.

Question: What are the three data quality issues that require addressing?

Sequencing Activity

Put the following steps in the correct order by numbering each one.

	Steps
	Create a free account key at the Windows Azure Marketplace.
	Subscribe to a free or paid-for RDS provider's service at the Marketplace.
	Configure the reference data service details in DQS.
	Map your domain to the RDS.
	Use the knowledge base containing the domain that maps to the RDS to cleanse the data.

Lesson 2

Using Data Quality Services to Cleanse Data

One of the major tasks for a data quality management solution is to cleanse data by validating and correcting domain values. This lesson describes how you can use DQS to cleanse data and review the results.

Lesson Objectives

After completing this lesson, you will be able to:

- Create a data cleansing project.
- View cleansed data.
- Use the Data Cleansing transformation in an SSIS data flow.

Creating a Data Cleansing Project

One of the major tasks for a data quality management solution is to cleanse data by validating and correcting domain values. This lesson describes how you can use DQS to cleanse data and review the results.

Data stewards can use the Data Quality Client application to create a data cleansing project that applies the knowledge in a knowledge base to data in a SQL Server database or an Excel® workbook.

When creating a data cleansing project, the data steward must:

- Select the knowledge base to use and specify that cleansing is the action to be performed.
- Select the source containing the data to be cleansed and map the columns in it to the domains in the knowledge base.
- Run the data cleansing process, and then review the suggestions and corrections generated by DQS. The data steward can then approve or reject the suggestions and corrections.
- Export the cleansed data to a database table, comma-delimited file, or Excel workbook.

1. Select a knowledge base
2. Map columns to domains
3. Review suggestions and corrections
4. Export results

Viewing Cleansed Data

The output from a data cleansing project includes the cleansed data, in addition to extra information about the corrections made by DQS. The output columns are named by combining the name of the domain and the type of data. For example, the cleansed output for a domain named **State** is stored in a column named **State_Output**.

Cleansed data output includes the following column types:

- **Output:** the values for all fields after data cleansing
- **Source:** the original value for fields that were mapped to domains and cleansed
- **Reason:** the reason the output value was selected by the cleansing operation
- **Confidence:** an indication of the confidence Data Quality Services estimates for corrected values
- **Status:** the status of the output column (correct or corrected)

- **Output:** the values for all fields after data cleansing. All fields in the original data source generate output columns, even those not mapped to domains in the knowledge base (in which case, they contain the original data values).
- **Source:** the original value for fields that were mapped to domains and cleansed.
- **Reason:** the reason the output value was selected by the cleansing operation. For example, a valid value might be corrected to a leading value defined for the domain, or DQS might have applied a cleansing algorithm and suggested a corrected value.
- **Confidence:** an indication of the confidence DQS estimates for corrected values. For values corrected to leading values defined in the knowledge base, this is usually 1 (or 100%). When DQS uses a cleansing algorithm to suggest a correction, the confidence is a value between 0 and 1.
- **Status:** the status of the output column. A value of "correct" indicates that the original value was already correct, and a value of "corrected" indicates that DQS changed the value.

Demonstration: Cleansing Data

In this demonstration, you will see how to use DQS to cleanse data.

Demonstration Steps

Create a Data Cleansing Project

1. Ensure you have completed the previous demonstration in this module.
2. In SQL Server Data Quality Services, in the **Data Quality Projects** section, click **New Data Quality Project**.
3. In the **Name** box, type **Cleansing Demo**.
4. In the **Use Knowledge Base** list, ensure **Demo KB** is selected.
5. In the **Select Activity** section, ensure **Cleansing** is selected, and then click **Next**.
6. On the **Map** page, in the **Data Source** list, ensure **SQL Server** is selected.
7. In the **Database** list, click **DemoDQS**.
8. In the **Table/View** list, click **Stores**.

9. In the **Mappings** table, create the following mappings, and then click **Next**:

Source Column	Domain
City (varchar)	City
StoreType (varchar)	StoreType
State (varchar)	State

10. On the **Cleanse** page, click **Start**.
11. When the cleansing process has completed, view the data in the **Profiler** tab, noting the number of corrected and suggested values for each domain, and then click **Next**.
12. On the **Manage and View Results** page, ensure that the **City** domain is selected, and on the **Suggested** tab, note that DQS has suggested correcting the value **New Yrk** to **New York**. Click the **Approve** option to accept this suggestion, and then click the **Corrected** tab to verify that the value has been corrected.
13. Click the **State** and **StoreType** domains in turn, and on the **Corrected** tab, note the corrections applied, based on the values defined in the knowledge base. Click **Next**.
14. On the **Export** page, view the output data preview.
15. In the **Export cleansing results** section, in the **Destination Type** list, click **Excel File**.
16. In the **Excel file name** box, type **D:\Demofiles\Mod10\CleansedStores.xls**, ensure that **Standardize Output** is selected, ensure that the **Data and Cleansing Info** option is selected, and then click **Export**.
17. In the **Exporting** dialog box, when the file download has completed, click **Close**.
18. On the **Export** page, click **Finish**.
19. Close SQL Server Data Quality Services.

View Cleansed Data

- Open the file **D:\Demofiles\Mod10\CleansedStores.xls** in Excel.
- Note that the output includes the following types of column:
 - Output:** the values for all fields after data cleansing.
 - Source:** the original value for fields that were mapped to domains and cleansed.
 - Reason:** the reason the output value was selected by the cleansing operation.
 - Confidence:** an indication of the confidence DQS estimates for corrected values.
 - Status:** the status of the output column (correct or corrected).
- Close Excel without saving any changes.

Using the Data Cleansing Data Flow Transformation

In addition to creating data cleansing projects to operate interactively, you can use the Data Cleansing transformation to work in a SQL Server Integration Services (SSIS) data flow. By using the Data Cleansing transformation, you can automate data cleansing as a part of the extract, transform, and load (ETL) processes used to populate your data warehouse.

To add the Data Cleansing transformation to a data flow in an SSIS package, perform the following steps:

1. Add the Data Cleansing transformation to the data flow and drag a data flow connection, from a source or transformation containing the data you want to cleanse, to the input of the Data Cleansing transformation.
2. Edit the settings of the Data Cleansing transformation to connect to the data quality server, specify the knowledge base you want to use, and map the data flow input columns to domains in the knowledge base.
3. Drag the output from the Data Cleansing transformation to the next transformation or destination in the data flow, and map the output columns from the Data Cleansing transformation to the appropriate input columns. The output columns from the Data Cleansing transformation are the same as those generated by an interactive data cleansing project.

- Input data to be cleansed
- Select knowledge base and map columns to domains
- Output cleansed columns



DQS Cleansing Transformation

<http://aka.ms/jkxdbo>

Check Your Knowledge

Question	
Which of these is not a mandatory part of creating a data cleansing project?	
Select the correct answer.	
<input type="checkbox"/>	Select the source containing the data that is to be cleansed and then map its columns to the domains in the knowledge base.
<input type="checkbox"/>	Export the cleansed data to a database table, comma-delimited file, or Excel workbook.
<input type="checkbox"/>	Select the knowledge base to use and specify that cleansing is the action to be performed.
<input type="checkbox"/>	Run the data cleansing process then review the suggestions and corrections generated by DQS. The data steward can then approve or reject the suggestions and corrections.
<input type="checkbox"/>	Create a document to describe the knowledge base.

Lab A: Cleansing Data

Scenario

You have created an ETL solution for the Adventure Works data warehouse, and invited some data stewards to validate the process before putting it into production.

The data stewards have noticed some data quality issues in the staged customer data, and have asked you to provide a way for them to cleanse data, so that the data warehouse is based on consistent and reliable data. The data stewards have given you an Excel workbook containing some examples of the issues found in the data.

Objectives

After completing this lab, you will be able to:

- Create a DQS knowledge base.
- Use DQS to cleanse data.
- Incorporate data cleansing into an SSIS data flow.

Lab Setup

Estimated Time: 30 minutes.

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

The setup script for this lab does not delete any existing DQS knowledge bases or projects. It is recommended that students create snapshots for the 20767C-MIA-DC, 20767C-MIA-CLI, and 20767C-MIA-SQL virtual machines before starting, so that, if necessary, the lab environment can be returned to the starting point.

Exercise 1: Creating a DQS Knowledge Base

Scenario

You have integrated data from many sources into your data warehouse, and this has several benefits. However, users have noticed some quality issues with the data, which you must correct.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. View Existing Data
3. Create a Knowledge Base
4. Perform Knowledge Discovery

► Task 1: Prepare the Lab Environment

1. Ensure that the 20767C-MIA-DC, 20767C-MIA-CLI and 20767C-MIA-SQL virtual machines are all running, and then log on to 20767C-MIA-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the D:\Labfiles\Lab10\Starter folder as **Administrator**.

► Task 2: View Existing Data

1. Log on to 20767C-MIA-CLI as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Open D:\Labfiles\Lab10\Starter\Sample Customer Data.xls in Excel and examine the worksheets in the workbook. Note that:
 - There are multiple names for the same country on the **CountryRegion** and **StateProvince** worksheets.
 - There are multiple names for the same state on the **StateProvince** worksheets.
 - Some customers do not have a gender code of **F** or **M** on the **Gender** worksheet.
3. On the **Sample Customer Data** worksheet, apply column filters to explore the data further and view the source records for the anomalous data.
4. Close Excel without saving any changes to the workbook.

► Task 3: Create a Knowledge Base

1. Start the Data Quality Client application and connect to **MIA-SQL**.
2. Create a new knowledge base with the following properties:
 - **Name:** Customer KB.
 - **Description:** Customer data knowledge base.
 - **Create Knowledge Base From:** Existing Knowledge Base (DQS Data).
 - **Select Activity:** Domain Management.
3. View the domain values for the **Country/Region**, **Country/Region (two-letter leading)**, and **US - State** domains.
4. Change the name of the **US – State** domain to **State**.
5. Create a domain with the following properties:
 - **Domain Name:** Gender.
 - **Description:** Male or female.
 - **Data Type:** String.
 - **Use Leading Values:** Selected.
 - **Normalize String:** Selected.
 - **Format Output to:** Upper Case.
 - **Language:** English.
 - **Enable Speller:** Selected.
 - **Disable Syntax Error Algorithms:** Not selected.
6. View the domain values for the **Gender** domain, and notice that null is allowed.
7. Add new domain values for **F**, **M**, **Female**, and **Male** to the **Gender** domain.
8. Set **F** and **Female** as synonyms, with **F** as the leading value.
9. Set **M** and **Male** as synonyms, with **M** as the leading value.
10. Finish editing the knowledge base, but do not publish it.

► Task 4: Perform Knowledge Discovery

1. Open the **Customer KB** knowledge base for knowledge discovery.
2. Use the **Sample Customer Data\$** worksheet in the **Sample Customer Data.xls** Excel workbook in D:\Labfiles\Lab10\Starter as the source data for mapping. Use the first row as the header.
3. In the **Mappings** table, add the following mappings:

Source Column	Domain
CountryRegionCode (String)	Country/Region (two-letter leading)
CountryRegionName (String)	Country/Region
StateProvinceName (String)	State
Gender (String)	Gender

4. Start the discovery process, and when it is complete, view the new values that have been discovered for the **State** domain. Set **New South Wales** and **NSW** as synonyms with **New South Wales** as the leading value. In the alphabetically ordered list of values, click **New South Wales** first, press the Ctrl key, click **NSW** to select them both, and then click **Set selected domain values as synonyms**.
5. View the new values that have been discovered for the **Country/Region (two-letter leading)** domain, and mark the value **UK** as an error that should be corrected to **GB**.
6. View the new values that have been discovered for the **Gender** domain, mark the value **W** as invalid, and correct it to **F**.
7. View the new values that have been discovered for the **Country/Region** domain, and remove the filter that causes the list to show only new values.
8. Finish and publish the knowledge base.

Results: After this exercise, you should have created a knowledge base and performed knowledge discovery.

Exercise 2: Using a DQS Project to Cleanse Data

Scenario

Now you have a published knowledge base, you can use it to perform data cleansing in a data quality project.

The main tasks for this exercise are as follows:

1. Create a Data Quality Data Project

► Task 1: Create a Data Quality Data Project

1. Create a new data quality project with the following properties:
 - **Name:** Cleanse Customer Data.
 - **Description:** Apply Customer KB to customer data.
 - **Use knowledge base:** Customer KB.
 - **Select Activity:** Cleansing.
2. On the **Map** page, select the **InternetSales** SQL Server database, then select the **Customers** table. Then in the **Mappings** table, add the following mappings:

Source Column	Domain
CountryRegionCode (nvarchar)	Country/Region (two-letter leading)
CountryRegionName (nvarchar)	Country/Region
StateProvinceName (nvarchar)	State
Gender (nvarchar)	Gender

3. Start the cleansing process, review the source statistics in the Profiler pane, and then on the **Manage and View Results** page, note that DQS has found the value **Australia**, which is likely to be a typographical error, and suggested it be corrected to **Australia** on the **Suggested** tab of the **Country** domain.
4. Approve the suggested correction, and note that it is now listed on the **Corrected** tab. View the corrected values for the **Country/Region**, **Country/Region (two-letter leading)**, **Gender**, and **State** domains.
5. On the **Export** page, view the output data, and then export the data and cleansing info to an Excel file named **D:\Labfiles\Lab10\Starter\CleansedCustomers.xls**.
6. When the export is complete, finish the project and view the results of the cleansing process in Excel.

Results: After this exercise, you should have used a DQS project to cleanse data and export it as an Excel workbook.

Exercise 3: Using DQS in an SSIS Package

Scenario

You are happy with the data cleansing capabilities of DQS and the results are accurate enough to be automated. You will edit an SSIS package to include a data cleansing component as part of a dataflow.

The main tasks for this exercise are as follows:

1. Add a DQS Cleansing Transformation to a Data Flow
2. Test the Package

► Task 1: Add a DQS Cleansing Transformation to a Data Flow

1. Return to the 20767C-MIA-SQL virtual machine.
2. Open the D:\Labfiles\Lab10\Starter\AdventureWorksETL.sln solution in Visual Studio®.
3. Open the **Extract Internet Sales Data.dtsx** SSIS package and, if it is not already visible, display the **SSIS Toolbox** (which is available on the **SSIS** menu).
4. Open the **Extract Customers** data flow task, add a **DQS Cleansing** transformation, and rename it to **Cleanse Customer Data**.
5. Remove the data flow between **Customers** and **Staging DB** and add a data flow from **Customers** to **Cleanse Customers**.
6. Configure the following settings for **Cleanse Customer Data**:
 - Create a new Data Quality connection manager for the **MIA-SQL** server and the **Customer KB** knowledge base.
 - Specify the following mappings with the default source, output, and status alias values:

Input Column	Domain
Gender	Gender
StateProvinceName	State
CountryRegionCode	Country/Region (two-letter leading)
CountryRegionName	Country/Region

- Standardize the output.
- Connect the output data flow from **Cleanse Customer Data** to **Staging DB** and change the following column mappings in the **Staging DB** destination (leaving the remaining existing mappings as they are):

Input Column	Destination Column
Gender_Output	Gender
StateProvinceName_Output	StateProvinceName
CountryRegionCode_Output	CountryRegionCode

Input Column	Destination Column
CountryRegionName_Output	CountryRegionName

► **Task 2: Test the Package**

1. Debug the package and observe the **Extract Customers** data flow as it executes, noting the number of rows processed by the **Cleanse Customer Data** transformation.
2. When package execution has completed, stop debugging and close Visual Studio (note that the entire package may take some time to complete after the **Extract Customers** data flow has finished).

Results: After this exercise, you should have created and tested an SSIS package that cleanses data.

Lesson 3

Using Data Quality Services to Match Data

In addition to cleansing data, you can use DQS to identify duplicate data entities. The ability to match data entities is useful when you need to deduplicate data to eliminate errors in reports and analysis, caused by the same entity being counted more than once.

This lesson explains how to create a matching policy, and then use it to find duplicate data entities in a data matching project.

Lesson Objectives

After completing this lesson, you will be able to:

- Create a matching policy.
- Create a data matching project.
- View data matching results.

Creating a Matching Policy

A data warehouse is almost always composed of data from multiple sources and sometimes might be provided by third parties. There are also likely to be many transactions that relate to the same customer or product—unless you have a system that only allows existing customers to buy existing products, duplication is likely.

For example, suppose you sell books on the Internet. To buy a book, a customer must first register with their name, address, username, and password. Customers often forget their details and, although there is a Forgotten

Username/Password link, many choose to register again. You can implement a constraint to stop anyone having the same username, but there might still be many instances when a single customer registers multiple times. This duplication can cause problems for data analysis because you will have more customers in your system than in reality. You might believe there are more people of a particular gender, more between certain age brackets, or more in a particular geographic area, than actually exist. The duplicate entries will also affect sales per customer analysis, which will return less than accurate results.

By providing a constraint to prevent duplicate usernames, you have gone some way to avoiding duplication; however, as you can see from the example, this will only slightly reduce the problem. You could enforce unique names, but that would make it difficult for customers with common names to register. You could enforce unique names at the same address, but that would block someone who has the same name as a partner or child.

Matching Policies

DQS can use a matching policy to assess the likelihood of records being duplicates. In cases with a high likelihood of duplication, data stewards assess the potential duplicates before any changes are made. A data steward can add a matching policy to a knowledge base and create rules that help determine whether multiple data records represent the same business entity. A data matching rule compares one or

- Define matching rules for business entities
- Rules match entities based on domains:
 - **Similarity:** similar or exact match
 - **Weight:** percentage to apply if match succeeds
 - **Prerequisite:** mandatory domain match for rule to succeed
- If the combined weight of all matches meets or exceeds the rule's minimum matching score, the entities are duplicates

more domains across records and applies weighted comparisons to identify matches. For each domain in the matching rule, the data steward defines the following settings:

- **Similarity:** you can specify that the rule should look for similar values, based on fuzzy logic comparison, or an exact match.
- **Weight:** a percentage score to apply if a domain match succeeds.
- **Prerequisite:** indicates that this particular domain must match for the records to be considered duplicates.

For each rule, the data steward specifies a minimum matching score. When the matching process occurs, the individual weightings for each successful domain match comparison are added together. If the total is equal to or greater than the minimum matching score, and all prerequisite domains match, the records are considered to be duplicates.

Creating a Data Matching Project

Data stewards can use the Data Quality Client application to create a data matching project that applies the knowledge in a knowledge base to data in a SQL Server database or an Excel workbook.

To create a data matching project, the data steward must:

1. Select the knowledge base to use and specify that the action to be performed is matching.
2. Select the data source that contains the data that is to be matched and map the columns to the knowledge base domains.
3. Run the data matching process, then review the clusters of matched records that DQS identifies, based on the knowledge base matching policies.
4. Export the matched data to a database table, comma-delimited file, or Excel workbook. Additionally, the data steward can specify a survivorship rule that eliminates duplicate records and exports the surviving records. A data steward can specify the following rules for survivorship:
 - **Pivot record:** a record chosen arbitrarily by DQS in each cluster of matched records.
 - **Most complete and longest record:** the record that has fewest missing data values and the longest values in each field.
 - **Most complete record:** the record that has fewest missing data values.
 - **Longest record:** the record containing the longest values in each field.

1. Select a knowledge base
2. Map columns to domains
3. Review match clusters
4. Export matches and survivors
 - Select survivorship rule:
 - Pivot record
 - Most complete and longest record
 - Most complete record
 - Longest record

Viewing Data Matching Results

After the data matching process is complete, you can view and export the following results:

- **Matches:** the original dataset plus additional columns that indicate clusters of matched records.
- **Survivors:** the resulting dataset, with duplicate records eliminated, based on the selected survivorship rule.

When you export matches, the results include the original data and the following columns:

- **Cluster ID:** a unique identifier for a cluster of matched records.
- **Record ID:** a unique identifier for each matched record.
- **Matching Rule:** the rule that produced the match.
- **Score:** the combined weighting of the matched domains as defined in the matching rule.
- **Pivot Mark:** a matched record chosen arbitrarily by DQS as the pivot record for a cluster.

- **Cluster ID:** identifier for a cluster of matched records
- **Record ID:** identifier for a matched record
- **Matching Rule:** the rule that produced the match
- **Score:** combined weighting of match criteria
- **Pivot Mark:** a matched record arbitrarily chosen by Data Quality Services as the pivot record for a cluster

Demonstration: Matching Data

In this demonstration, you will see how to use DQS to match data.

Demonstration Steps

Create a Matching Policy

1. Ensure you have completed the previous demonstrations in this module.
2. Return to the 20767C-MIA-CLI virtual machine, start the **SQL Server Data Quality Client application**, and connect to **MIA-SQL**.
3. In the **Knowledge Base Management** pane, under **Recent Knowledge Base**, click **Demo KB**, and then click **Matching Policy**.
4. On the **Map** page, in the **Data Source** drop-down list, select **Excel File**, in the **Excel File** box, and browse to **D:\Demofiles\Mod10\Stores.xls**.
5. In the **Worksheet** drop-down list, ensure **Sheet1\$** is selected, and ensure **Use first row as header** is selected. This worksheet contains a sample of store data that needs to be matched.
6. In the **Mappings** table, add the following mappings to existing domains:

Source Column	Domain
City (String)	City
StoreType (String)	StoreType
State (String)	State

- In the **Mappings** table, add the following mappings by selecting a **Source Column** and creating a **New Domain**, then click **Next**.

Note: When the **Mappings** table is full, click **Add a column mapping** to add an additional row.

Source Column	Domain
PhoneNumber (String)	PhoneNumber
StreetAddress (String)	StreetAddress
StoreName (String)	StoreName

- On the **Matching Policy** page, click **Create a matching rule**.
- In the **Rule Details** section, change the rule name to **Is Same Store**.
- In the **Rule Editor** table, click **Add a new domain element**.
- In the **Domain** column, ensure that **StoreName** is selected. In the **Similarity** column, ensure that **Similar** is selected, in the **Weight** column, enter **20**, and leave the **Prerequisite** column clear.
- Repeat the previous step to add the following rules:
 - StreetAddress**: Similar: 20%: Not a prerequisite.
 - City**: Exact: 20%: Not a prerequisite.
 - PhoneNumber**: Exact: 30%: Not a prerequisite.
 - StoreType**: Similar: 10%: Not a prerequisite.
 - State**: Exact: 0%: Prerequisite selected.
- Click **Start**, wait for the matching process to complete, and note the matches that are detected in the sample data (Store 1 is the same as Store One, for example).
- Click **Next**.
- On the **Matching Results** page, view the details in the **Profiler** tab, and then click **Finish**.
- When prompted to publish the knowledge base, click **Publish**.
- When publishing is complete, click **OK**.

Create a Data Matching Project

- In SQL Server Data Quality Services, in the Data Quality Projects pane, click **New Data Quality Project**.
- Create a new project named **Matching Demo** based on the **Demo KB** knowledge base.
- In the **Select Activity** pane, click **Matching**, and then click **Next**.
- On the **Map** page, in the **Data Source** list, ensure **SQL Server** is selected. Then in the **Database** list, click **DemoDQS**, and in the **Table/View** list, click **Stores**.

5. In the **Mappings** table, add the following mappings to existing domains, and then click **Next**:

Source Column	Domain
City (varchar)	City
StoreType (varchar)	StoreType
State (varchar)	State
PhoneNumber (varchar)	PhoneNumber
StreetAddress (varchar)	StreetAddress
StoreName (varchar)	StoreName

6. On the **Matching** page, click **Start**.
7. When matching is complete, note that two matches were detected (Store 1 is the same as Store One and Store 16 is the same as Store Sixteen). Click **Next**.
8. On the **Export** page, in the **Destination Type** drop-down list, select **Excel File**, and then select the following content to export:
- **Matching Results:** D:\Demofiles\Mod10\MatchedStores.xls
 - **Survivorship Results:** D:\Demofiles\Mod10\SurvivingStores.xls
9. Select the **Most complete record** survivorship rule, and then click **Export**. Select **Yes** if prompted to overwrite existing files.
10. When the export has completed successfully, click **Close**.
11. Click **Finish**.
12. Close SQL Server Data Quality Services.

View Data Matching Results

1. Open **D:\Demofiles\Mod10\MatchedStores.xls** in Excel. Note that this file contains all the records in the dataset with additional columns to indicate clusters of matched records. In this case, there are two clusters, each containing two matches.
2. Open **D:\Demofiles\Mod10\SurvivingStores.xls** in Excel. Note this file contains the records that were selected to survive the matching process. The data has been deduplicated by eliminating duplicates and retaining only the most complete record.
3. Close Excel without saving any changes.

Question: What are the rules for survivorship that a data steward can specify?

Verify the correctness of the statement by placing a mark in the column to the right.

Statement	Answer
<p>True or false? For each domain in the matching rule, the data steward defines the following settings:</p> <p>Similarity: you can specify that the rule should look for similar values based on fuzzy logic comparison, or an exact match.</p> <p>Weight: a percentage score to apply if a domain match succeeds.</p> <p>Prerequisite: indicates that this particular domain must match, for the records to be considered duplicates.</p>	

Lab B: Deduplicating Data

Scenario

You have created a DQS knowledge base and used it to cleanse customer data. However, data stewards are concerned that the staged customer data might include duplicate entries. For records to be considered a match, the following criteria must be true:

- The **Country/Region** column must be an exact match.
- A total matching score of **80** or higher must be achieved, based on the following weightings:
 - An exact match of the **Gender** column has a weighting of **10**.
 - An exact match of the **City** column has a weighting of **20**.
 - An exact match of the **EmailAddress** column has a weighting of **30**.
 - A similar **FirstName** column value has a weighting of **10**.
 - A similar **LastName** column value has a weighting of **10**.
 - A similar **AddressLine1** column value has a weighting of **20**.

Objectives

After completing this lab, you will be able to:

- Add a matching policy to a DQS knowledge base.
- Use DQS to match data.

Estimated Time: 30 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Creating a Matching Policy

Scenario

You have implemented a data cleansing solution for the data being staged. However, you have identified that the staged data contains multiple records for the same business entity. You want to use a data matching solution to deduplicate the data.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Create a Matching Policy

► Task 1: Prepare the Lab Environment

- On the 20767C-MIA-SQL virtual machine, run **LabB.cmd** in the D:\Labfiles\Lab10\Starter folder as Administrator.

► Task 2: Create a Matching Policy

1. On the 20767C-MIA-CLI virtual machine, start **SQL Server Data Quality Client** and connect to the **MIA-SQL** server.
2. Open the **Customer KB** knowledge base for the **Matching Policy** activity.

3. Use the **Sheet1\$** worksheet in the **D:\Labfiles\Lab10\Starter\Sample Staged Data.xls** Excel file as the data source for the matching policy.
4. On the **Map** page, map the columns in the Excel worksheet to the following new domains:

Source Column	Domain
FirstName (String)	A new domain named FirstName with a String data type.
LastName (String)	A new domain named LastName with a String data type.
AddressLine1 (String)	A new domain named AddressLine1 with a String data type.
City (String)	A new domain named City with a String data type.
EmailAddress (String)	A new domain named EmailAddress with a String data type.

5. Add more column mappings and map the following fields to existing domains:

Source Column	Domain
Gender (String)	Gender
StateProvinceName (String)	State
CountryRegionCode (String)	Country/Region (two-letter heading)
CountryRegionName (String)	Country/Region

6. On the **Matching Policy** page, create a matching rule with the following properties:
 - **Rule name:** Is Same Customer.
 - **Description:** Checks for duplicate customer records.
 - **Min. matching score:** 80.
 - **Rule Editor.** Add these rules:

Domain	Similarity	Weight	Prerequisite
Country/Region	Exact	0	Selected
Gender	Exact	10	Unselected
City	Exact	20	Unselected
EmailAddress	Exact	30	Unselected
FirstName	Similar	10	Unselected

Domain	Similarity	Weight	Prerequisite
LastName	Similar	10	Unselected
AddressLine1	Similar	20	Unselected

7. Start the matching process and, when it has finished, review the matches found by DQS, noting that there are duplicate records for three customers.
8. When you have finished, publish the knowledge base.

Results: After this exercise, you should have created a matching policy and published the knowledge base.

Exercise 2: Using a DQS Project to Match Data

Scenario

You will now create a data quality project to apply the matching rules from the previous exercise. After this process is complete, you will have exported a deduplicated set of data. Finally, you will apply the deduplication results in the staging database by executing Transact-SQL statements.

The main tasks for this exercise are as follows:

1. Create a Data Quality Project for Matching Data
2. Review and Apply Matching Results

► Task 1: Create a Data Quality Project for Matching Data

1. In SQL Server Data Quality Services, create a new data quality project with the following details:
 - **Name:** Deduplicate Customers.
 - **Description:** Identify customer matches.
 - **Use knowledge base:** Customer KB.
 - **Select Activity:** Matching.
2. Using the **Customers** table in the **Staging** SQL Server database as the data source, map the following columns to domains in the knowledge base:

Source Column	Domain
FirstName (nvarchar)	FirstName
LastName (nvarchar)	LastName
Gender (nvarchar)	Gender

Source Column	Domain
AddressLine1 (nvarchar)	AddressLine1
City (nvarchar)	City
CountryRegionName (nvarchar)	Country/Region
EmailAddress (nvarchar)	EmailAddress

3. Start the matching process and review the results when it is finished.
4. Export the results to the following Excel workbooks, specifying the **Most complete record** survivorship rule:
 - **Matching Results:** D:\Labfiles\Lab10\Starter\Matches.xls
 - **Survivorship Results:** D:\Labfiles\Lab10\Starter\Survivors.xls

► Task 2: Review and Apply Matching Results

1. Open D:\Labfiles\Lab10\Starter\Matches.xls in Excel.
2. Note that the matching process found a match with a score of 90 for the following customer records:
 - **CustomerBusinessKey:** 29261 (Robert Turner).
 - **CustomerBusinessKey:** 29484 (Rob Turner).
3. Open D:\Labfiles\Lab10\Starter\Survivors.xls in Excel.
4. Note that the survivors file contains all the records that should survive deduplication, based on the matches that were found. It contains the record for customer 29261 (Robert Turner), but not for 29484 (Rob Turner).
5. On the 20767C-MIA-SQL virtual machine, open **D:\Labfiles\Lab10\Starter\Fix Duplicates.sql** in SQL Server Management Studio and connect to the **MIA-SQL** instance of the database engine by using Windows authentication.
6. Review the Transact-SQL code and note that it performs the following tasks:
 - It updates the **InternetSales** table so that all sales currently associated with the duplicate customer record become associated with the surviving customer record.
 - It deletes the duplicate customer record.
7. Execute the SQL statement, and then close SQL Server Management Studio.

Results: After this exercise, you should have deduplicated data using a matching project and updated data in your database to reflect these changes.

Module Review and Takeaways

In this module, you have learned how Data Quality Service provides a knowledge-based solution for cleansing and matching data.

Review Question(s)

Question: Who is responsible for ensuring the consistency and accuracy of data in solutions you have implemented or managed?

MCT USE ONLY. STUDENT USE PROHIBITED

Module 11

Master Data Services

Contents:

Module Overview	11-1
Lesson 1: Introduction to Master Data Services	11-2
Lesson 2: Implementing a Master Data Services Model	11-6
Lesson 3: Hierarchies and Collections	11-18
Lesson 4: Creating a Master Data Hub	11-27
Lab: Implementing Master Data Services Model	11-34
Module Review and Takeaways	11-42

Module Overview

Organizations typically use different platforms to store different types of data. For example, sales data might be stored in an online transactional processing (OLTP) database, customer data in a dedicated customer relationship management (CRM) system, and so on. Storing data across multiple, heterogeneous platforms can make it difficult to ensure that the data representing a single instance of a specific business entity is consistent and accurate across the enterprise.

Master Data Services provides a way for organizations to standardize and improve the quality, consistency, and reliability of the data that guides key business decisions. This module introduces Master Data Services and explains the benefits of using it.

Objectives

After completing this module, you will be able to:

- Describe the key concepts of Master Data Services.
- Implement a Master Data Services model.
- Use Master Data Services tools to manage master data.
- Use Master Data Services tools to create a master data hub.

Lesson 1

Introduction to Master Data Services

Master data management can represent a major challenge for organizations. Data representations of a single business entity, such as a specific individual customer, might be recorded in multiple locations in multiple formats. When you consider the number of different types of data a company might own, the potential scale of this problem can be huge. Master Data Services enables organizations to standardize data, which improves the consistency of their key business data and, ultimately, the quality of the decisions they make.

Lesson Objectives

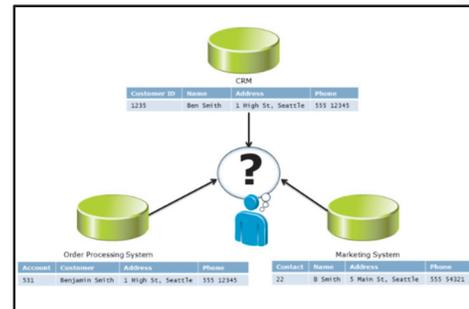
After completing this lesson, you will be able to:

- Explain the need for master data management.
- Explain how Master Data Services helps you meet the challenges of master data management.
- Compare Master Data Services with Data Quality Services.
- Describe the components of Master Data Services.

Master Data Services Concepts

Data owned by an organization is one of its most valuable assets; businesses rely on it for a variety of different reasons. For example, individuals in a company might use data to develop marketing strategies, plan new product lines, or identify areas where they can make efficiency savings. Report writers and data analysts interact directly with this data; decision makers use the data to guide them when making key choices about future business strategy. However, companies often struggle to maintain the quality, consistency, and accuracy of their data across the enterprise for a number of reasons, including:

- **Decentralized data storage.** In modern, complex information ecosystems, data might be stored in multiple systems and formats, perhaps because of departmental differences, or after company mergers or acquisitions. This approach makes it difficult to identify where duplicate data exists and, if it does, how to identify a “master” version.
- **Different methods of handling data changes.** Applications might handle data changes, such as additions, updates, and deletions, by using different rules—this can result in inconsistencies. For example, if one application records addresses without requiring a postal code but others do, the formats of stored addresses will be inconsistent.
- **Human error.** Errors in the insertion and updating of data can lead to inaccuracies. For example, if a user misspells a customer name, applications will accept the input, providing the format is correct and the error is not identified.

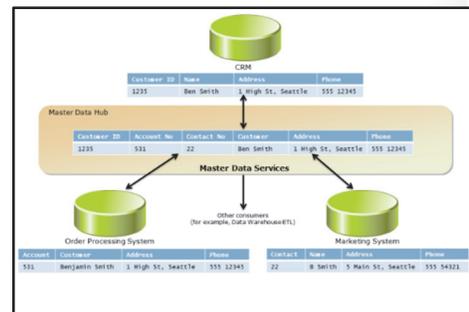


- **Latency and nonpropagation of changes.** Changes to data in one system might not propagate to others where the same data is held—or it might do so with a time delay. Running reports against these diverse systems will yield different results.

Poorly managed data directly affects the quality of reporting and data analysis. Ultimately, it can result in inefficient procedures, missed opportunities, revenue loss, customer dissatisfaction, and increased time spent managing the data to try to solve these problems. This also makes it more difficult for organizations to comply with data regulation requirements. Microsoft SQL Server Master Data Services helps organizations to implement good data stewardship, which can address the issues associated with poorly managed data.

What Is Master Data Services?

Master Data Services is a master data management technology that helps organizations to handle the challenges of data management. Master Data Services can serve as a system of entry for creating and updating master data, in addition to a system of record for making authoritative data available to other applications. With Master Data Services, you can create a master data hub to consolidate and ensure consistency of key business entity data representations across the enterprise.



You can use Master Data Services to enforce data validation rules and track changes to master data through an audit trail that shows the time, date, and author of each change. This promotes accountability and helps organizations to comply with data regulation requirements.

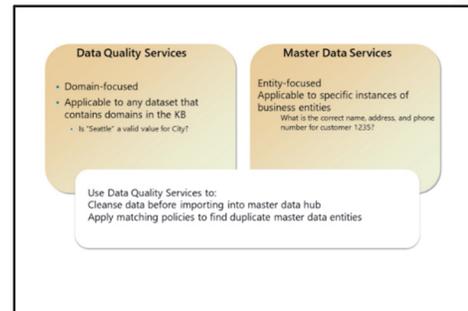
Data Stewards

A data steward is an individual charged with managing master data. Usually a data steward is a business user with a detailed knowledge of the entities used in a particular area—they are responsible for ensuring the integrity of data relating to those entities. Data stewards must ensure that each item of data is unambiguous, has a clear definition, and is used consistently across all systems. Typically, a data steward will use a master data management system, such as Master Data Services, to perform these tasks.

Note: Implementing a successful master data management initiative involves establishing extensive cooperation among the various stakeholders and the owners of the data. This might not always be straightforward. For example, stakeholders might be reluctant to commit to the scheme because they see it as surrendering stewardship of their data; or there could be political issues in an organization that have to be overcome. It is therefore important that all parties fully understand the benefits of master data management and how these are delivered by Master Data Services.

Master Data Services and Data Quality Services

On initial inspection, Master Data Services appears to address many of the same problems as Data Quality Services. Both technologies are designed to help users manage the quality and integrity of business data—but there are some key differences. While Data Quality Services is focused on managing the integrity and consistency of individual domains or columns in a dataset, Master Data Services is designed to manage the integrity and consistency of specific entity instances.



For example, an organization might use Data Quality Services to ensure that all records that include address data use a consistent, approved set of valid values for **City**, **State**, and **Country** fields. This domain-based validation can apply to any record that contains address data, including customer, employee, and supplier records, and so on.

In contrast, the same organization might use Master Data Services to ensure the integrity of a specific customer named Ben Smith. Address data relating to that individual might be entered and maintained in multiple systems across the organization. Master Data Services ensures that there is definitive information about Ben Smith—business users can then be sure which of his multiple addresses, that the organization has spread across various systems, is the correct one. Note that, from a Data Quality Services perspective, all the addresses the organization has for Ben Smith might be valid, but with Master Data Services, users can be sure which one is correct.

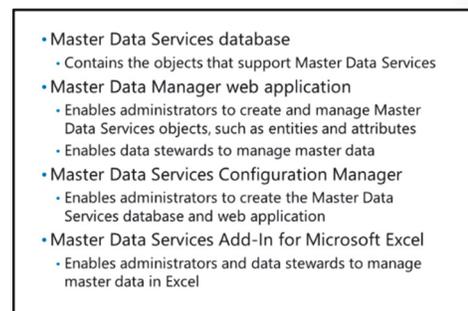
In many scenarios, Master Data Services and Data Quality Services are used together. Examples of this include:

- Using Data Quality Services to cleanse data before loading into Master Data Services.
- Applying Master Data Services matching policies to identify duplicate master data records for business entities.

Components of Master Data Services

Master Data Services includes the following components:

- **Master Data Services database.** This database stores all the database objects that support Master Data Services. The database contains staging tables for processing imported data, views that give client access to master data, and tables that store the master data itself. The database supports additional functionality, including versioning, business rule validation, and email notification.
- **Master Data Manager web application.** You use this application to perform the tasks associated with managing Master Data Services, such as creating models, entities, hierarchies, subscription views, and business rules. You configure the Master Data Manager web application by using the Explorer, Version Management, Integration



Management, System Administration, and User and Group Permissions functional areas. Apart from the Explorer, all functional areas are restricted to Master Data Services administrators. Data stewards can use the Explorer functional area to manage the master data for which they have responsibility.

- **Master Data Services Configuration Manager.** You use this tool to create and manage the Master Data Services database and web application, and to enable the web service. You can also use it to create a Database Mail profile for Master Data Services to use.
- **Master Data Services Add-In for Microsoft Excel®.** Business users and administrators can use this add-in to manage Master Data Services objects, and to work with master data in a Microsoft Excel workbook.

Considerations for Installing Master Data Services

Master Data Services is available only with the following editions of SQL Server:

- SQL Server Business Intelligence.
- SQL Server Enterprise.
- SQL Server Developer.

You can install Master Data Services by using SQL Server Setup, and adding it as a shared feature. You should use Master Data Services Configuration Manager to complete the installation and create the Master Data Services database and web application.

Question: What is the main data issue within an organization that can be resolved using Data Quality Services, and what are the advantages?

Lesson 2

Implementing a Master Data Services Model

At the center of every Master Data Services solution is a Master Data Services model. To create a master data management solution with Master Data Services, you must know how to create and manage a model. This lesson explains the key concepts you need to understand about Master Data Services models, and describes how to create and manage a model, and the master data it contains.

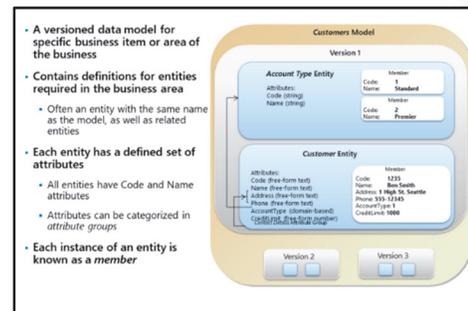
Lesson Objectives

After completing this lesson, you will be able to:

- Describe the key features of a Master Data Services model.
- Create a Master Data Services model.
- Create entities and attributes in a Master Data Services model.
- Add and edit entity members in a Master Data Services model.
- Understand Master Data Services security.
- Deploy a model to another Master Data Services environment.
- Edit a Master Data Services model in Microsoft Excel.

What Is a Master Data Services Model?

A model is the highest level of organization in Master Data Services. It is a container for a related set of business entity definitions. You can create a model for each area of the business for which you want to manage data. For example, to create a definitive set of customer data, you might create a model named **Customers** that contains all customer-related data. The **Customers** model would include data and metadata about the customers themselves, and might include related information, such as customer account types and sales territory data. You might then create a second model named **Products** that contains all product data, and other models as required.



Versions

Master Data Services models are versioned, meaning you can maintain multiple versions of master data at the same time. This can be useful in scenarios where many business applications require a newer definition of a specific business entity; however, some older applications cannot be upgraded to use the new model.

Entities and Attributes

An entity is a data definition for a specific type of item used in the business. For example, a **Customers** model might contain a **Customer** entity and an **Account Type** entity. An entity is analogous to a table in a relational database. In many cases, a model is created primarily to manage a single business entity, so an entity is created with the same name as the model. This scenario is so common that an option to create an identically named entity when creating a model is available in Master Data Services.

Each entity has attributes that describe it. These attributes are analogous to columns in a database table. When you create an entity, Master Data Services automatically adds two attributes to that entity:

- **Code.** The Code attribute can contain only unique values. When you populate the entity, you must provide a Code value for each member. The value of a Code attribute is frequently derived from the primary key column in a relational database table.
- **Name.** The Name attribute does not require a unique value—you can leave this field blank for members if appropriate.

You cannot delete the Code and Name attributes.

In addition to the system generated attributes, you can also use Master Data Services to create one of three types of attribute that describe your data:

- **Free-form.** Use this type of attribute to enter free-form values as text, numbers, dates, or links. You use free-form attributes for the text-based, numerical, and date and time data in your databases.
- **Domain-based.** This type of attribute accepts values only from other entities. You cannot directly enter values into domain-based attributes. You use domain-based attributes to ensure that the values for a particular attribute match code values in an existing entity. For example, suppose you have an entity called Product Category that lists product categories by name, and another entity called Product Subcategory. In the Product Subcategory entity, you can create a domain-based attribute called Category that references the Product Category entity.
- **File.** This type of attribute accepts files, such as documents or images. You can use file attributes to ensure that all files in an attribute have the same file extension.

Attribute Groups

An attribute group is a named grouping of the attributes in an entity. Attribute groups are useful if an entity has a large number of attributes—which makes them difficult to view—or in scenarios where multiple applications will consume entity data from Master Data Services. However, some attributes are only relevant to particular applications. For example, a Customer entity might include attributes that are used only in a CRM application, and other attributes that are only used by an order processing system. By creating application-specific attribute groups, you can simplify the creation of data flows between the master data hub and the applications requiring master data.

Members

Members are individual instances of entities, and are analogous to records in a database table. Each instance of an entity is a member that represents a specific business object, such as an individual customer or product.

Creating a Model

To create a model, open the Master Data Manager web application, and perform the following procedure:

1. Click **System Administration**.
2. On the **Model View** page, on **Manage** menu, click **Models**.
3. On the **Model Maintenance** page, click **Add model**.
4. In the **Model Name** box, type a unique name for the model.
5. Optionally:
 - Select **Create entity with same name as model** to create an entity with the same name as the model.
 - Select **Create explicit hierarchy with same name as model** to create an explicit hierarchy with the same name as the model. This option also enables the entity for collections.
 - Select **Mandatory hierarchy (all leaf members are included)** to make the explicit hierarchy mandatory.
6. Click **Save model**.

- Use the Master Data Services web application
- Specify a unique model name
- Optionally:
 - Create an entity with the same name as the model
 - Create an explicit hierarchy with the same name as the model
 - Make the explicit hierarchy mandatory

Creating Entities and Attributes

After you have created a model, you must add entities to represent the business objects for which you want to manage master data.

Creating an Entity

Use the following procedure to create an entity:

1. In Master Data Manager, click **System Administration**.
2. On the **Model View** page, on the **Manage** menu, click **Entities**.
3. On the **Entity Maintenance** page, in the **Model** list, select the model in which you want to create the entity and click **Add entity**.
4. In the **Entity name** box, type a name that is unique within the model.
5. Optionally, in the **Name for staging tables** box, type a name for the staging table (if you don't enter a name, the entity name is used by default).
6. Optionally, select the **Create Code values automatically** check box so that Code attribute values are generated automatically for members as they are added.

Creating an Entity

- Specify a unique name
- Optionally:
 - Specify a staging table name
 - Enable automatic code values
 - Enable explicit hierarchies and collections

Adding Attributes

- Edit an entity to add leaf member attributes
- Specify attribute type:
 - Free-form
 - Domain-based
 - File

7. In the **Enable explicit hierarchies and collections** list, select **Yes** if you want to enable explicit hierarchies and collections for this entity, or **No** if you do not. If you select **Yes**, you can optionally select **Mandatory hierarchy (all leaf members are included)** to make the explicit hierarchy a mandatory one.
8. Click **Save entity**.

Adding Attributes to an Entity

After you have created an entity, it will contain the mandatory Code and Name attributes. You can then edit the entity to add more attributes by following this procedure:

1. In Master Data Manager, click **System Administration**.
2. On the **Model View** page, on the **Manage** menu, click **Entities**.
3. On the **Entity Maintenance** page, in the **Model** list, select the model that contains the entity to which you want to add attributes.
4. Select the entity to which you want to add attributes, and click **Edit selected entity**.
5. On the **Edit Entity** page, in the leaf member attributes pane, click **Add leaf attribute**.
6. On the **Add Attribute** page, select one of the following options:
 - **Free-form**. This is an option where users can enter attribute values.
 - **Domain-based**. Use this option to create an attribute that is used as a key to look up members in another, related entity.
 - **File**. Use this option to create an attribute that is represented by a file, such as an image.
7. In the **Name** box, type a name for the attribute that is unique within the entity.
8. Depending on the option you chose previously, you can set additional display and data type settings for the attribute. For example, when adding a free-form attribute, you can specify a data type such as Text, Number, or Date. You can then set data type-specific constraints, such as a maximum length for a text attribute or a number of decimal places, for a numeric attribute.
9. Optionally, select **Enable change tracking** to track changes to groups of attributes.
10. Click **Save attribute**.
11. On the **Entity Maintenance** page, click **Save entity**.

Adding and Editing Members

After you have created entities and defined their attributes, you can add members to the model.

To use Master Data Manager to add a member:

1. On the **Master Data Manager** home page, in the **Model** list, select the model to which you want to add a member.
2. In the **Version** list, select the version of the model you want to work with.
3. Click **Explorer**.

- Use Explorer in the Master Data Services web application
- Add, edit, and delete members for each entity in the model
 - Add annotations to document transactions

4. On the **Entities** menu, click the name of the entity to which you want to add a member.
5. Click **Add member**.
6. In the Details pane, enter a value for each attribute.
7. Optionally, in the **Annotations** box, type a comment to document the addition of the member.
8. Click **OK**.

After you have added a member, you can edit its attributes by selecting it and modifying the attribute values in the Details pane.

Adding annotations helps document each change made to the data. You can view a history of all edits, and associated annotations by clicking the **View Transactions** button.

Demonstration: Creating a Master Data Services Model

In this demonstration, you will see how to create a Master Data Services model.

Demonstration Steps

Create a Model

1. Ensure that the **20767C-MIA-DC**, **20767C-MIA-CLI**, and **20767C-MIA-SQL** virtual machines are all running, and log into the **20767C-MIA-SQL** virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod11** folder, right-click **Setup.cmd** and click **Run as administrator**. When prompted, click **Yes**.



Note: The script might report an error and cause the virtual machine to restart. If this happens, log in again once the virtual machine is running and rerun the script. The same error might be reported, but you can ignore it.

3. Start Internet Explorer and browse to **http://mia-sql:81/mds**
4. On the **Master Data Services** home page, click **System Administration**.
5. On the **Manage Models** page, on the **Manage** menu, click **Models**.
6. On the **Model Maintenance** page, click **Add**.
7. In the Add Model pane, in the **Name** box, type **Customers**, clear the **Create entity with same name as model** check box, and then click **Save**.

Create an Entity

1. On the **Manage Models** page, select the **Customers** model, and then click **Entities**.
2. On the **Manage Entities** page, in the **Model** list, click **Customers**, and then click **Add**.
3. In the Add Entity pane, in the **Name** box, type **Customer**, and then click **Save**.

Create Attributes

1. On the **Manage Entities** page, in the **Entity** table, click **Customer**, and then click **Attributes**.
2. On the **Manage Attributes** page, click **Add**.

3. In the Add Attribute pane, in the **Name** box, type **Address**, in the **Attribute Type** list, click **Free-form**, in the **Data Type** list, click **Text**, in the **Length** box, type **400**, and then click **Save**.
4. On the **Manage Attributes** page, click **Add**.
5. In the Add Attribute pane, in the **Name** box, type **Phone**, in the **Attribute Type** list, click **Free-form**, in the **Data Type** list, click **Text**, in the **Length** box, type **20**, and then click **Save**.
6. Click the **SQL Server** text in the page title to return to the home page.

Add and Edit Members

1. On the **Master Data Services** home page, in the **Model** drop-down list, click **Customers**, and then click **Explorer**.
2. Click **Add Member**, and in the Details pane, enter the following data:
 - **Name:** Ben Smith
 - **Code:** 1235
 - **Address:** 1 High St, Seattle
 - **Phone:** 555 12345
 - **Annotations:** Initial data entry
3. Click **OK**, and then click the entity row you have added.
4. In the Details pane, edit the following fields, and then click **OK**:
 - **Phone:** 555 54321
 - **Annotations:** Changed phone number
5. In the right-hand pane, click the **View History** link (not the **View History** menu item), noting the list of transactions for this entity.
6. Click each transaction and view the text in the **Annotations** tab before clicking **Close**.
7. Close Internet Explorer.

Master Data Services Security

Just as you would secure your SQL Server database to prevent unauthorized access, it is essential that you apply the same security principles to Master Data Services. While it is important to ensure users have access to the data they need for their jobs, you must also prevent them from viewing data that should not be available to them. Two types of users exist in MDS: users who can access data in the Explorer functional area, and Administrators who have access to the Explorer functionality and perform administrative tasks in other areas. You can manage security in the **User and Group Permissions** functional area of the Master Data Manager website, or by using the web service.

- Give users access to only what they need in MDS
- Two types of user exist:
 - Users who have access to the Explorer functional area
 - Administrators who have access to other functional areas, and are permitted to perform administrative tasks
- Functional areas:
 - Explorer, Version Management, Integration Management, System Administration, User and Group Permissions, and Super User
- Tree structure of model enables permissions to be inherited at lower levels; individual object permissions can be set underneath
- Permissions: Read, Create, Update, and Deny; explicitly assign a combination to create precise permissions

MDS security is based on local or Active Directory® domain users and groups. This enables you to set permissions at a very detailed level when granting access to a user. However, be aware that permissions granted to users in multiple groups have the potential to overlap.

When granting access to data or functionality in MDS, you need to assign the following:

- **Functional area access:** there are five functional areas within the web user interface that users can have access to:
 - Explorer
 - Version Management
 - Integration Management
 - System Administration
 - User and Group Permissions
 - Super User
- **Model objects permissions:** you can assign a combination of Read, Create, Update, and Deny permissions to model objects. The User and Group Permissions functional area represents the model in a tree structure, so when you assign permission to an object in the tree, all objects below it also inherit the permission. However, you can assign permission to individual objects, which gives you complete control over access.



Best Practice: The recommended best practice is to assign access permission to the model, then assign individual permissions on underlying objects.

You must assign a user or group access to one of the functional areas and one model to enable them to access Master Data Manager.

Deploying a Model

A Master Data Services package is an XML file containing a model structure that you can use to create a new model, copy a model from one MDS environment to another, optionally including the data in the model, or update an existing model. When you deploy a model, all objects are included in the package:

- Entities
- Attributes
- Attribute groups
- Hierarchies
- Collections
- Business rules
- Version flags
- Subscription views

- MDS package is an XML file containing model structure and, optionally, data
- Use to create a new model, clone an existing model, or update an existing previously cloned model
 - All objects deployed apart from file attributes and permissions
- Three tools available for deployment and editing:
 - **MDSModelDeploy:** deploys model, and optionally the data
 - **Model Deployment Wizard:** deploys only the model structure
 - **Model Package Editor Wizard:** edit the model package created by the above tools
- System Administrator permission required on target MDS
- Administrator permissions required on environment from which you are deploying
- Use the command prompt to deploy using MDSModelDeploy.exe

The package excludes file attributes, and user and group permissions, so be aware that you will need to configure security for the new model.

There are three different tools you can use to deploy and edit a model package:

- **MDSModelDeploy**: this executable is installed as part of the MDS installation and can be found in the drive:\Program Files\Microsoft SQL Server\130\Master Data Services\Configuration folder if you installed MDS using the default path. Use this tool to create a new model, or deploy an existing model and, optionally, the data.
- **Model Deployment Wizard**: use the wizard to deploy packages with the model structure only. The wizard can be found in the Master Data Manager web application. It can't be used to deploy data, but you can use it to create new models.
- **Model Package Editor Wizard**: this wizard enables you to edit the model packages created by the above tools. Start the ModelPackageEditor.exe tool to launch the Model Package Editor wizard, which is in the drive:\Program Files\Microsoft SQL Server\130\Master Data Services\Configuration folder if you installed MDS using the default path.

Sample Packages

When you install Master Data Services, sample packages are included. If you use the MDSModelDeploy tool to deploy a sample package, you can also deploy the data. The samples can be found in the drive:\Program Files\Microsoft SQL Server\130\Master Data Services\Samples\Packages folder if you installed MDS to the default location.

Prerequisites

Before updating an existing model with data, ensure the version of the model you are deploying is neither **Locked** nor **Committed**. You will need permission to access the **System Administration** functional area in the target MDS environment, in addition to being an administrator in the environment from which you are deploying the model. Furthermore, the model can only be deployed to the same version of SQL Server—a model created in SQL Server 2012 cannot be deployed to a higher version of SQL Server.

Deploying a Model

1. Decide on the type of deployment: a new model, clone an existing model, or update a previously cloned model.
2. Right-click on **Start**, and select **Command Prompt (Admin)**. In the **User Account Control** dialog, click **Yes** to allow changes.
3. In the command window, type **CD Program Files\Microsoft SQL Server\130\Master Data Services\Configuration** to change directory. You may need to change the drive letter if your installation was not installed using the default settings.
4. For a list of available options, type **MDSModelDeploy** and press Enter. For specific help, type **MDSModelDeploy help deploynew**, replacing **deploynew** with the command you want further information about.
5. The syntax for deploying a new package is: **MDSModelDeploy deploynew -package PackageName -model ModelName -service ServiceName**. If you are unsure of the name of the MDS service, or there are multiple instances installed on the local server, type **MDSModelDeploy listservices**, and press Enter.
6. To deploy the product package in the sample folder, type **MDSModelDeploy deploynew -package "C:\Program Files\Microsoft SQL Server\130\Master Data Services\Samples\Packages\product_en.pkg" -model Products -service MDS1** and press Enter. If there are spaces in the path name of the package, this should be enclosed in double quotes.

7. When the deployment completes, you will see the message **MDSModelDeploy operation completed successfully**.

After deploying the model, you will need to configure object security. For more information on deploying models using the MDS deployment tools, see the following article:



Deploying Models (Master Data Services)

<https://aka.ms/fui3mz>

Editing a Model in Microsoft Excel

SQL Server Master Data Services supports the Master Data Services Add-in for Excel, which you can use to read and manage lists of Master Data Services data. The add-in is a free download for Excel 2007 or later that you can distribute to data stewards and Master Data Services administrators, so they can work with a familiar interface. You can download the Master Data Services Add-In for Excel from the Microsoft download site, or users can install it directly from Master Data Manager.

The Master Data Services Add-in for Excel adds the **Master Data** tab to the Excel ribbon. You can use options on this tab to perform tasks, such as connecting to a Master Data Services server, applying business rules, creating new entities, and publishing changes back to the server. The add-in is security context aware—it only allows users to view and change data for which they have the appropriate permissions.

To view Master Data Services data in Excel, you must first connect to a Master Data Services server. On the **Master Data** tab in Excel, you can use the **Connect** option to create a connection to Master Data Services. After connecting, use the Master Data Explorer to select a model and version to work with from those available on the server. You can then load data into an Excel worksheet from the entities listed in the Master Data Explorer and filter that data, so you only see the actual data you want to work with. After you have loaded the required data, you can browse and edit it, just as you would any other data in Excel. You can create new entities, add columns to existing entities, to define new attributes, and edit member data.

Most data editing operations are performed locally in the Excel worksheet. Changes are propagated to the Master Data Services database only when you explicitly publish the changes made in Excel. When you publish an entity, you can enter annotations to document the changes made.

- Use the Master Data Services Add-In for Excel to connect to a model
- Create entities
- Add columns to create attributes
- Edit entity member data in worksheets
- Publish changes to Master Data Services

Demonstration: Editing a Model in Excel

In this demonstration, you will see how to edit a master data model in Excel.

Demonstration Steps

Connect to a Master Data Services Model in Excel

1. Ensure you have completed the previous demonstration in this module.
2. Log into the **20767C-MIA-CLI** virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
3. Start Microsoft Excel and create a new blank workbook.
4. On the **File** tab, click **Options**.
5. In the **Excel Options** dialog box, on the **Add-ins** tab, in the **Manage** drop-down list, select **COM Add-ins**, and then click **Go**.
6. In the **COM Add-ins** dialog box, if **Master Data Services Add-In for Excel** is not selected, select it. Click **OK**.
7. On the **Master Data** tab of the ribbon, in the **Connect and Load** section, click the drop-down arrow under **Connect**, and then click **Manage Connections**.
8. In the **Manage Connections** dialog box, click **New**.
9. In the **Add New Connection** dialog box, enter the description **Demo MDS Server** and the MDS server address **http://mia-sql:81/mds**, and then click **OK**.
10. Click **Close**.
11. On the **Master Data** tab of the ribbon, in the **Connect and Load** section, click the drop-down arrow under **Connect**, and then click **Demo MDS Server**.
12. In the Master Data Explorer pane, in the **Model** drop-down list, select **Customers**.
13. In the Master Data Explorer pane, click the **Customer** entity.
14. On the ribbon, in the **Connect and Load** section, click **Refresh**.
Note that the Customer entity, including the member you created in the previous demonstration, downloads into a new worksheet.

Add a Member

1. On the **Customer** worksheet, click cell **N4** (which should be an empty cell under Ben Smith).
2. Enter the following details in row 4:
 - **N4**: Andrew Sinclair
 - **O4**: 2600
 - **P4**: 2 Main St, Ontario
 - **Q4**: 555 11111
3. Note that the row for the new member you have created has an orange background to indicate that the data has not been published to Master Data Services.
4. On the ribbon, in the **Publish and Validate** section, click **Publish**. If a **Publish and Annotate** dialog box appears, select the **Do not show this dialog box again** check box, and then click **Publish**. Note that the data is published and the orange background is removed.

Add a Free-Form Attribute to an Entity

1. On the **Customer** worksheet, click cell **R2** (which should be an empty cell to the right of the **Phone** attribute header).
2. Type **CreditLimit** and press Enter. This adds a new attribute named **CreditLimit** (which has a green background because it is not yet saved to the model).
3. Click cell **R2** to reselect it, and on the **Master Data** tab of the ribbon, in the **Build Model** section, click **Attribute Properties**.
4. In the **Attribute Properties** dialog box, in the **Attribute type** drop-down list, click **Number**, note the default **Decimal places** value (2), and then click **OK**. The changes are uploaded to the data model and the cell background changes to blue.
5. In cell **R4**, enter **1000** as the **CreditLimit** value for the **Andrew Sinclair** member, and in cell **R3**, enter **500** as the **CreditLimit** value for the **Ben Smith** member. Note that the cell background is orange to indicate that the data has not been published to Master Data Services.
6. On the ribbon, in the **Publish and Validate** section, click **Publish**.
Note that the data is published and the orange background is removed.

Add a Domain-Based Attribute and Related Entity

1. On the **Customer** worksheet, click cell **S2** (which should be an empty cell to the right of the **CreditLimit** attribute header).
2. Type **AccountType** and press Enter. This adds a new attribute named **AccountType** (which is shown with a green background because it is not yet saved to the model).
3. In cell **S4**, enter **2** as the **AccountType** value for **Andrew Sinclair**. Then in cell **S3**, enter **1** as the **AccountType** value for **Ben Smith**.
4. Click cell **S2** to reselect it, and on the ribbon, in the **Build Model** section, click **Attribute Properties**.
5. In the **Attribute Properties** dialog box, in the **Attribute type** drop-down list, click **Constrained list (Domain-based)**, in the **Populate the attribute with values from** list, ensure **the selected column** is selected, in the **New entity name** box, type **Account Type**, and then click **OK**. Note that the **AccountType** column now contains the values **1 {1}** and **2 {2}**.
6. On the **Master Data** tab of the ribbon, in the **Connect and Load** section, click the drop-down arrow under **Connect**, and then click **Demo MDS Server**.
7. In the Master Data Explorer pane, in the **Model** list, select **Customers**.
8. In the Master Data Explorer pane, click the **Account Type** entity, created when you added a domain-based attribute to the **Customer** entity.
9. On the ribbon, in the **Connect and Load** section, click **Refresh**.
Note that the Account Type entity is downloaded into a new worksheet with two members.
10. Change the **Name** attribute for the existing members as follows:

Name	Code
Standard	1
Premier	2

11. On the ribbon, in the **Publish and Validate** section, click **Publish**.

12. Click the **Customer** worksheet tab, and on the ribbon, in the **Connect and Load** section, click **Refresh**. Note that the **AccountType** attribute values update to show the new Name values you specified for the members in the **Account Type** entity.
13. Click cell **S4**, and note that you can select the **AccountType** value for each member from a list that looks up the related name and code in the **Account Type** entity. Leave the selected **AccountType** for Andrew Sinclair as 2 {Premier}.
14. Close Excel without saving the workbook.

Check Your Knowledge

Question	
What type of attribute would you create if you wanted the values to be dependent on values in another entity?	
Select the correct answer.	
<input type="checkbox"/>	File.
<input type="checkbox"/>	Domain-based.
<input type="checkbox"/>	An entity cannot be dependent on another entity.
<input type="checkbox"/>	Free-form.

Check Your Knowledge

Question	
Which two attributes are automatically created for each new entity?	
Select the correct answer.	
<input type="checkbox"/>	ID and Name.
<input type="checkbox"/>	Code and Name.
<input type="checkbox"/>	Name and Description.
<input type="checkbox"/>	Code and ID.

Lesson 3

Hierarchies and Collections

After you have defined a Master Data Services model and the entities you want to manage, you can use Master Data Services to organize, manage, and maintain entity members. This means you can create master data solutions that meet the requirements of your organization's business applications and ensure the enterprise-wide integrity and consistency of your data.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe how an administrator can use hierarchies and collections to organize master data.
- Create derived hierarchies.
- Create explicit hierarchies.
- Create collections.
- Find duplicate members.
- Use business rules to validate members.

Hierarchies and Collections

Hierarchies and collections are two ways to group related data in a Master Data Services model.

You can use hierarchies to group members and provide users with a structured view of data that is easier to browse. A hierarchy contains all the members from the entity or entities that you add, and each member can appear only once.

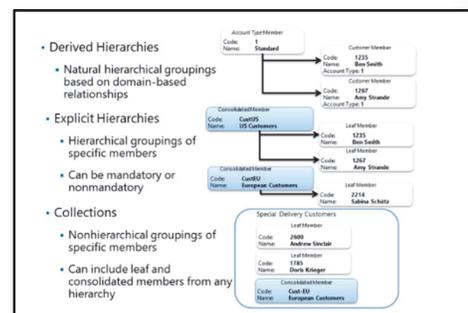
Collections are flat groupings of members within an entity—they have no hierarchical structure, but make it easier to find and manage members that are related in some way.

Derived Hierarchies

Derived hierarchies are based on the relationships that exist between entities. They use domain-based attributes that you create to infer parent-child relationships. For example, suppose that you have two entities, **Customer** and **Account Type**. The **Customer** entity contains a domain-based attribute named **AccountType** that references the **Code** attribute of the **Account Type** entity. You can use these attributes to create a derived hierarchy that groups customers by account type.

Explicit Hierarchies

Explicit hierarchies can contain members from a single entity, and do not use domain-based attributes to determine their structure. Instead, you must create consolidated members to define levels within the hierarchy, and then move the leaf members (members that represent instances of the entity) into the appropriate level. For example, you could create a consolidated member named **US Customers**, and group leaf members for customers with a US address under it. Unlike derived hierarchies, explicit hierarchies can be ragged, so do not need a consistent number of levels. An explicit hierarchy contains all the members in an entity. If you make the hierarchy mandatory, all entity leaf members must be assigned



to a consolidated member. In a nonmandatory hierarchy, any members that you do not group under a consolidated member are grouped together in the hierarchy as “unused”.

Collections

You can use collections to combine leaf members and consolidated members from any existing explicit hierarchies into a group of related members. Collections are flexible and efficient, making it possible to reuse existing hierarchies—there’s no need to create a new one. For example, you could create a collection named **Special Delivery Customers** that contains consolidated members from an explicit hierarchy based on geographic locations—in addition to some individual leaf members that don’t belong to the consolidated members, but which also qualify for special delivery status.

Where possible, you should try to use **Derived Hierarchies** instead of **Explicit Hierarchies** and **Collections**; the Explicit Hierarchy and Collections functionality has been deprecated and might not be included in future versions of Microsoft SQL Server Master Data Services. The addition of features such as an unassigned/unused category for each hierarchy level make using **Derived Hierarchies**, instead of **Explicit Hierarchies** and **Collections**, relatively straightforward.

Deprecated: Explicit Hierarchies and Collections

<http://aka.ms/L5nmut>

 **Note:** You can only create explicit hierarchies and collections for entities that have the Enable Explicit Hierarchies and Collections option enabled.

Creating Derived Hierarchies

Use the following procedure to create a derived hierarchy:

1. In Master Data Manager, click **System Administration**.
2. On the **Model View** page, in the **Manage** menu, click **Derived Hierarchies**.
3. On the **Derived Hierarchy Maintenance** page, from the **Model** list, select the model in which you want to create a derived hierarchy and click **Add**.
4. On the **Add Derived Hierarchy** page, in the **Derived hierarchy name** box, type a name for the hierarchy. Try to make the name meaningful—for example, **Customers by Account Type**.
5. Click **Save derived hierarchy**.
6. On the **Edit Derived Hierarchy** page, in the Available Entities and Hierarchies pane, click the entity that represents the top level of your hierarchy (for example, **Account Type**) and drag it to the Current Levels pane.
7. Drag the entity that represents the next level (for example, Customer) onto the level that you created in the previous step.
Note that there must be a relationship between a domain-based attribute in the entity you are dragging and the Code attribute in the parent entity.
8. Continue to drag further entities until your hierarchy includes all the levels it requires.

- Link hierarchy levels to define a named hierarchy based on domain-based attribute relationships
- Navigate the derived hierarchy in Explorer

Creating Explicit Hierarchies

To create an explicit hierarchy, you must enable explicit hierarchies and collections for the entity, and edit it to create any required consolidated member attributes. You can then define the consolidated members for the hierarchy levels and assign leaf members to them.

Enabling Explicit Hierarchies and Collections

Use the following procedure to enable explicit hierarchies and collections for an entity:

1. In Master Data Manager, click **System Administration**.
2. On the **Model View** page, in the **Manage** menu, click **Entities**.
3. On the **Entity Maintenance** page, in the **Model** list, select the model that contains the entity for which you want to create an explicit hierarchy.
4. Select the entity that you want to update, and then click **Edit selected entity**.
5. In the **Enable explicit hierarchies and collections** list, select **Yes**.
6. In the **Explicit hierarchy** name box, type a name for the explicit hierarchy you want to create. Try to make the name meaningful—for example, **Customers by Geographic Location**.
7. Optionally, clear the **Mandatory hierarchy** check box to create the hierarchy as nonmandatory.
8. Click **Save entity**.

Editing Entities

After you have enabled explicit hierarchies and collections, and created an explicit hierarchy, you can start to define consolidated members. However, by default, consolidated members only have **Code** and **Name** attributes, so you might want to edit the entity to add consolidated member attributes. Additionally, if you need to add another explicit hierarchy to the entity, you can edit it and click **Add explicit hierarchy**.

Defining Hierarchies

After you have created an explicit hierarchy, you must define the consolidated members that form the hierarchy levels and assign leaf members to them. Use the following procedure to define an explicit hierarchy:

1. On the Master Data Manager home page, in the **Model** list, select the model that contains the entity with the explicit hierarchy you want to edit.
2. In the **Version** list, select the version you want to work with.
3. Click **Explorer**.
4. In the **Hierarchies** menu, click the name of the hierarchy you want to define.
5. Above the grid, select either **Consolidated members** or **All consolidated members in hierarchy**.
6. Click **Add**.
7. In the Details pane, enter values for the consolidated member attributes. Optionally, in the **Annotations** box, type a comment to document the creation of the consolidated member.
8. Click **OK**.

- Enable explicit hierarchies and collections, and create an explicit hierarchy
 - Optionally, make the hierarchy mandatory
- Edit entities as required to:
 - Create consolidated member attributes
 - Create additional explicit hierarchies
- In Explorer, define the hierarchies
 - Create consolidated members for levels
 - Move leaf members to consolidated members

9. After creating the consolidated members that define the hierarchy levels, you can move members to create the hierarchy structure. To place members under a consolidated member:
 - a. In the hierarchy tree view, select the check box for each leaf or consolidated member you want to move.
 - b. Click **Cut**.
 - c. Select the consolidated member to which you want to assign the selected members.
 - d. Click **Paste**.

Creating Collections

To create a collection, you must enable explicit hierarchies and collections for the entity as described in the previous topic. You can then use the following procedure to create a collection:

1. On the Master Data Manager home page, in the **Model** list, select the model that contains the entity for which you want to create a collection.
2. In the **Version** list, select the version you want to work with.
3. Click **Explorer**.
4. In the **Collections** area, click the entity for which you want to create a collection.
5. Click **Add collection**.
6. On the **Details** tab, in the **Name** box, type a name for the collection.
7. In the **Code** box, type a unique code for the collection.
8. Optionally, in the **Description** box, type a description for the collection and click **OK**.
9. On the **Collection Members** tab, click **Edit Members**.
10. To filter the list of available members, select from the list on the left.
11. Click each member you want to add and click **Add**.
12. Optionally, rearrange collection members by clicking **Up** or **Down**.

- Enable explicit hierarchies and collections
- Edit entities to create collection attributes if necessary
- In Explorer, create and edit collections
 - Specify a name, code, and other collection attributes for each new collection
 - Add leaf members and consolidated members to the collection

Finding Duplicate Members

When you create a model, the entities might have some initial members to help you structure the hierarchies and collections that the model must support. However, as the model moves into production, you will start to load large volumes of member data; it is important that data stewards can easily manage the integrity of this data.

One common problem is the creation of duplicate members for the same business entity. To help avoid this problem, the Master Data Services Add-in for Excel includes a data matching feature that uses Data Quality Services matching policies to identify possible duplicate members.

- Create a Data Quality Services knowledge base with domains for entity attributes
- Create a matching policy to identify possible duplicates
- Match data in the Master Data Services Add-In for Excel

To use the Master Data Services Add-in for Excel to find duplicate members:

1. Create a Data Quality Services knowledge base that includes domains for the attributes on which you want to compare members.
2. Create a Data Quality Services matching policy that identifies potential duplicates based on domain value matches.
3. Open the entity in the Master Data Services Add-in for Excel, click **Match Data**, and then specify the Data Quality Services server and knowledge base; map the knowledge base domains to entity attributes.

For more information on Data Quality Services knowledge bases, see:



DQS Knowledge Bases and Domains

<http://aka.ms/E4x3kl>

Validating Members with Business Rules

You can use business rules to ensure that the members you add to your models are accurate and meet the criteria that you defined. Business rules help to improve data quality, which increases the accuracy and reliability of the reports and data analyses produced by information workers. A business rule is an IF/THEN statement that you create by using the **Business Rule Maintenance** page in the Master Data Manager web application. To define the terms of a business rule, you use the drag-and-drop interface on the **Edit Business Rule** page.

- Define business rules with the business rule expression editor
- Publish business rules
- Use business rules to validate entity member data in:
 - Explorer
 - Excel



Note: You must be a model administrator to create and publish business rules.

When you create a rule, you define a condition that gives a true or false result. For example, you can create conditions such as "If a customer's credit limit is greater than 1,000". You can use a range of value comparison operators to build conditions, including "is equal to", "starts with", "is between", and "contains". When you validate data, Master Data Services compares the data values against these conditions. Any that match a condition are tested against the actions you define. For example, for the condition "If a customer's credit limit is greater than 1,000", the action could be "The account type must be Premier".

After creating a business rule, you must publish it to make it active. You can then use the rule to validate new members when they are added to Master Data Services. Business rules do not prevent users who have update permission from making changes that violate the rules. Instead, the validation process reveals any data values that violate business rules, so that a data steward can identify problematic data and make the required changes. When you validate data against business rules, you can choose to validate the whole version or just a subset of the data in a version, such as a single entity.



Note: You can use the AND operator and the OR operator to create rules with multiple conditions.

Business Rule Priority

You can set a value for each business rule that defines its priority. You can then control the order in which Master Data Services applies business rules. Rules with a low priority value run before rules that have higher priority values. For example, if a rule that checks the credit limit for customers has a priority of 10, it will run before rules with a priority greater than 10. By default, when you create a new rule, its priority value is higher than the previous rule by a value of 10—so the new rule has a lower priority than the previous rule and runs after it.

Priority is important when business rules depend on the outcome of other business rules. For example, suppose you create a rule that updates the value of the **CreditLimit** attribute of a **Customer** entity member. When you import a new customer, the rule sets the value to -1. A second rule checks the **CreditLimit** attribute for the value -1, changes the status of the member to "is not valid", and notifies the data steward that there is a new customer for whom they need to specify a credit limit. By setting the priority of the second rule to a higher value than the first, you ensure that the rules run in a coordinated and logical order, and the data steward can update credit limits for new customers sooner.

Demonstration: Creating and Applying Business Rules

In this demonstration, you will see how to create and apply business rules.

Demonstration Steps

Create Business Rules

1. Ensure you have completed the previous demonstrations in this module.
2. On the **20767C-MIA-CLI** virtual machine, start Internet Explorer and browse to **http://mia-sql:81/mds**
3. On the **Master Data Services** home page, click **System Administration**.
4. On the **Manage** menu, click **Business Rules**.
5. On the **Manage Business Rules** page, in the **Model** list, click **Customers**, and in the **Entity** list, click **Customer**.

6. In the **Member Type** list, click **Leaf**.
7. Click **Add** to create a new business rule.
8. In the **Add Business Rule** pane, in the **Name** box, type **Check Credit**.
9. In the **Description** box, type **Check that credit limit is valid**. You will use this rule to ensure that all customer credit limits are greater than or equal to zero.
10. Under **Then** click the **Add** link.
11. On the **Create Action** pane, in the **Attribute** list, click **CreditLimit**, in the Operator list, click **must be greater than or equal to**, in the **Must be greater than or equal to** list, click **Attribute value**, in the **Attribute value** box, type **0**, and then click **Save**.
12. In the **Add Business Rule** pane, click **Save**.
13. Click **Add** to create another new business rule.
14. In the **Add Business Rule** pane, in the **Name** box, type **Check Premier Status**.
15. In the **Description** box, type **Check account type for high credit limits**. You will use this rule to check that customers with a credit limit value greater than 1,000 have a premier account type.
16. Under **If**, click the **Add** link.
17. On the **Create Condition** pane, in the **Attribute** box, click **CreditLimit**, in the **Operator** box, click **is greater than**, in the **Is greater than** box, click **Attribute value**, in the **Attribute value** box, type **1000**, and then click **Save**.
18. Under **Then**, click the **Add** link.
19. On the **Create Action** pane, in the **Attribute** list, click **AccountType**, In the **Operator** list, click **must be equal to**, in the **Must be equal to** list, click **Attribute value** and then click the **Attribute value** box.
20. On the **Choose an Attribute Value** pane, in the **Version** list, click **VERSION_1**, in the **Attribute value** list, click **2**, and then click **Save**.
21. On the **Create Action** pane, click **Save** and then on, the **Add Business Rule** pane, click **Save**.

Publish Business Rules

1. On the **Manage Business Rules** page, click **Publish All**.
2. In the **Message from webpage** dialog box, click **OK**.
3. In the **Business Rule State** column, check that the value displayed for both rules is **Active**.
4. On the **Manage Business Rules** page, click the **Microsoft SQL Server** text in the page title to return to the home page.

Apply Business Rules in Explorer

1. On the **Master Data Services** home page, click **Explorer**.
2. In the **Entities** menu, click **Customer**.
3. If the Microsoft Silverlight window appears, complete the following steps:
 - a. Click **Click now to install**.
 - b. In the Internet Explorer message bar, click **Run**.
 - c. In the **User Account Control** dialog box, click **Yes**.

- d. In the **Install Silverlight** dialog box, click **Install now**.
 - e. In the **Enable Microsoft Update** dialog box, click **Next**.
 - f. On the **Installation successful** page, click **Close**.
 - g. Press F5.
4. Click **Apply Rules**, and note that a green tick is displayed next to all valid records.
 5. Click the row for the Andrew Sinclair member, in the Details tab, change the CreditLimit value to -200, and then click OK.
Note that a validation error message is displayed, and that the green tick for this member record changes to a red exclamation mark.
 6. Change the **CreditLimit** value for Andrew Sinclair back to **1000**, and click **OK**, noting that the validation message disappears and the red exclamation mark changes back to a green tick.
 7. Close Internet Explorer.

Apply Business Rules in Excel

1. Start Microsoft Excel and create a new blank workbook.
2. On the **File** tab, click **Options**.
3. In the **Excel Options** dialog box, on the **Add-ins** tab, in the **Manage** drop-down list, select **COM Add-ins**, and then click **Go**.
4. In the **COM Add-ins** dialog box, if **Master Data Services Add-In for Excel** is not selected, select it. Click **OK**.
5. On the **Master Data** tab of the ribbon, in the **Connect and Load** section, click the drop-down arrow under **Connect**, and then click **Demo MDS Server**.
6. In the **Master Data Explorer** pane, in the **Model** list, click **Customers**.
7. In the **Master Data Explorer** pane, click the **Customer** entity.
8. On the ribbon, in the **Connect and Load** section, click **Refresh**. The **Customer** entity members are downloaded into a new worksheet.
9. In the ribbon, in the **Publish and Validate** section, click **Show Status**. This reveals columns that show the validation and input status for all records.
10. In the ribbon, in the **Publish and Validate** section, click **Apply Rules**. This refreshes the validation status for all rows. Currently, validation is successful for all records.
11. In cell **R3**, change the **CreditLimit** value for Ben Smith to **1200**. Note that the **AccountType** value for this record is currently 1 {Standard}.
12. In the ribbon, in the **Publish and Validate** section, click **Publish**.
13. Note that the validation status in cell **B3** is **Validation failed**, and that the row is highlighted in blue.
14. Click cell **S3**, and then in the drop-down list, click **2 {Premier}** to change the **AccountType** value for Ben Smith.
15. In the ribbon, in the **Publish and Validate** section, click **Publish**.
16. Note that validation now succeeds for the Ben Smith member.
17. Close Excel without saving the workbook.

Question: When you have created business rules in Master Data Services, how might you prevent business users from entering data that violates those rules?

Lesson 4

Creating a Master Data Hub

Master data management provides a way to ensure data consistency for key business entities across the enterprise. In many cases, a master data hub is created in which data stewards can centrally manage data definitions for business entities used by multiple applications. This lesson describes the key features of a master data hub and explains how to create a solution in which data can flow into and out of a master data hub implemented in Master Data Services.

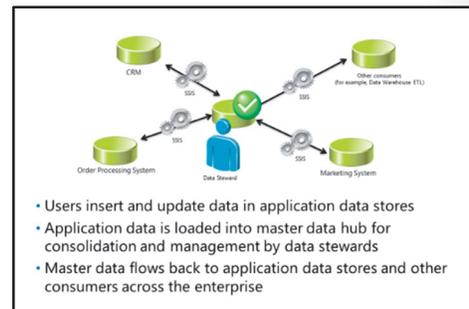
Lesson Objectives

After completing this lesson, you will be able to:

- Describe master data hub architecture.
- Describe the structure of Master Data Services staging tables.
- Import data into Master Data Services.
- Create subscription views so that applications can consume master data.

Master Data Hub Architecture

In some scenarios, data stewards can use Master Data Services as the sole data entry point for records that represent key data entities. This approach ensures that a single, consistent definition of each business entity exists and is used across the organization. However, in many organizations, data relating to the same business entity is entered and updated in multiple applications; a solution that consolidates and harmonizes this data must be developed. In cases like this, you can use a master data hub as a central point where new and updated records from multiple applications are imported, managed, and then flow back to the applications that use them.



For example, an organization might use the following applications to manage customer data:

- A CRM system
- An ERP system
- An order processing system
- A marketing system

As customer records are entered or updated in each of these systems, an SSIS-based solution takes the modified records and imports them into the master data hub, where a data steward can validate and consolidate the data relating to customer entities. An SSIS-based solution can then extract the consolidated customer data from the master data hub and replicate it back to the source systems, to ensure that all systems have consistent definitions for customers. Each SSIS data flow can only transfer the attributes that are required by the individual application being synchronized. The complete set of customer attributes is managed centrally in the master data hub.

Additionally, other applications that require accurate customer data, but do not need to modify it, can consume the data directly from the master data hub. For example, an ETL process for a data warehouse might take customer data from the master data hub to ensure that the data warehouse is populated with consistent, accurate records for customers.

Master Data Services Staging Tables

To facilitate the import of new and updated data from business applications, Master Data Services generates staging tables for each entity defined in a model. You can load data into these staging tables by using an SSIS data flow, the Import and Export Data wizard, Transact-SQL statements, the bulk copy program, or any other technique that can be used to insert data into a SQL Server database table.

Staging Tables for Leaf Members

Master Data Services creates a leaf member staging table for each entity you create. The table is created in the **stg** schema and, by default, is named in the format **EntityName_Leaf**. For example, a staging table for leaf members in a **Customer** entity would be named **stg.Customer_Leaf**. You can override the name of the staging table when you create an entity.

A leaf member staging table contains the following columns:

- **ID**: an automatically assigned identifier.
- **ImportType**: a numeric code that determines what to do when imported data matches an existing member in the model. For example, a value of 1 creates new members, but does not update existing members; a value of 2 replaces existing members with matching staged members.
- **ImportStatus_ID**: the status of the import process. You initially set this to 0 to indicate that the data is staged and ready for import. During the import process, the value is automatically updated to 1 if the import succeeds, or 2 if it fails.
- **Batch_ID**: a unique identifier for a batch of imported records. This value is automatically set when importing data through the Master Data Services web service—otherwise it is not required.
- **BatchTag**: a unique name that identifies a batch of imported records. This is used to group staged records instead of the **Batch_ID** column when not using the Master Data Services web service.
- **ErrorCode**: set by the import process if the import fails.
- **Code**: the unique code attribute for the member.
- **Name**: the name attribute for the member.
- **NewCode**: used to change the code attribute of a member.
- **<AttributeName>**: a column is created for each leaf attribute defined in the entity.

• MDS generates a staging table for each:

- Leaf member
- Consolidated member
- Relationship

stg.EntityName_Leaf	stg.EntityName_Consolidated	stg.EntityName_Relationship
ID	ID	ID
ImportType	ImportType	RelationshipType
ImportStatus_ID	ImportStatus_ID	ImportStatus_ID
Batch_ID	Batch_ID	Batch_ID
BatchTag	BatchTag	BatchTag
ErrorCode	HierarchyName	HierarchyName
Code	ErrorCode	ImportCode
Name	Code	ChildCode
NewCode	Name	SortOrder
<AttributeName>	NewCode	ErrorCode
	<AttributeName>	

Staging Tables for Consolidated Members

By default, staging tables for consolidated members take the name format **stg.EntityName_Consolidated** (for example, **stg.Customer_Consolidated**). Consolidated member staging tables contain the same columns as leaf member staging tables—there's an additional **HierarchyName** column you can use to indicate the explicit hierarchy to which the consolidated member should be imported. Additionally, the **<AttributeName>** columns are created to reflect the consolidated attributes defined in the entity, instead of the leaf attributes.

Staging Tables for Relationships

Relationship staging tables are used to change the location of members in an explicit hierarchy. By default, they take the name format **stg.EntityName_Relationship**, and contain the following columns:

- **ID**: an automatically assigned identifier.
- **RelationshipType**: a numeric code that indicates the type of relationship. Valid values for this column are 1 (parent) and 2 (sibling).
- **ImportStatus_ID**: the status of the import process. You initially set this to 0 to indicate that the data is staged and ready for import. During the import process, the value is automatically updated to 1 if the import succeeds or 2 if it fails.
- **Batch_ID**: a unique identifier for a batch of imported records. This value is automatically set when importing data through the Master Data Services web service. Otherwise it is not required.
- **BatchTag**: a unique name that identifies a batch of imported records. This is used to group staged records instead of the **Batch_ID** column when not using the Master Data Services web service.
- **HierarchyName**: the explicit hierarchy in which the relationship should be defined.
- **ParentCode**: the Code attribute of the member that will be the parent (for **RelationshipType 1** imports) or sibling (for **RelationshipType 2** imports) in the relationship.
- **ChildCode**: the Code attribute of the member that will be the child (for **RelationshipType 1** imports) or sibling (for **RelationshipType 2** imports) in the relationship.
- **SortOrder**: an optional integer that determines the sort order for sibling members under a parent.
- **ErrorCode**: set by the import process if the import fails.

Staging and Importing Data

The process for importing data into Master Data Services consists of the following steps:

1. Load new or updated member records into staging tables. You can use any technique that inserts data into a table in a SQL Server database.
2. Run the appropriate staging stored procedures to process the staged records. Master Data Services generates a stored procedure for each staging table, with names in the following formats:
 - **stg.udp.EntityName_Leaf**

1. Load data into staging tables
 - Using SSIS, Import and Export wizard, Transact-SQL, and so on
2. Run staging stored procedures
 - stg.udp.EntityName_Leaf
 - stg.udp.EntityName_Consolidated
 - stg.udp.EntityName_Relationship
3. View import status in Master Data Manager
4. Match data to identify any duplicates
5. Apply business rules to validate data

- **stg.udp.EntityName_Consolidated**
- **stg.udp.EntityName_Relationship**

When you execute the stored procedures, you must specify values for the **VersionName**, **LogFlag**, and **BatchTag** parameters to indicate the version of the model to be updated; whether or not to log transactions (1) or not (0); and the unique batch tag that identifies the staged records to be processed. Note that, when called by the Master Data Services web service, the **Batch_ID** parameter is used instead of the **BatchTag** parameter.

3. View the import status in Master Data Manager. The import process is asynchronous, so you can use the **Integration Management** area in Master Data Manager to view the start time, end time, and results of each batch that has been processed.
4. To identify and consolidate any duplicate records, data stewards should use the **Match Data** functionality in the Master Data Services Add-In for Excel after importing data.
5. Finally, a data steward should apply the business rules defined in the model for the imported members, to ensure that they are valid, before making them available to applications.

Demonstration: Importing Master Data

In this demonstration, you will see how to import data into a master data model.

Demonstration Steps

Use an SSIS Package to Import Master Data

1. Ensure you have completed the previous demonstrations in this module.
2. On the **20767C-MIA-SQL** virtual machine, start Visual Studio® and open the **MDS Import.sln** solution in the **D:\Demofiles\Mod11** folder.
3. In Solution Explorer, in the **SSIS Packages** folder, double-click the **Import Customers.dtsx** SSIS package.
4. On the control flow surface, double-click the **Load Staging Tables** task.
5. On the data flow surface, double-click **Customers Source**, and in the **Flat File Source Editor** dialog box, on the **Columns** tab, note the columns that will be imported from the source system. Click **Cancel**.
6. On the data flow surface, double-click **Add MDS Columns**, and in the **Derived Column Transformation Editor** dialog box, note that the transformation generates additional columns named **ImportType** (with a value of 0), **ImportStatus_ID** (with a value of 0), and **BatchTag** (with a unique string value derived from the **ExecutionInstanceGUID** system variable). Click **Cancel**.
7. On the data flow surface, double-click **Staging Table**, and in the **OLE DB Destination Editor** dialog box, on the **Connection Manager** tab, note that the data is loaded into the **[stg].[Customer_Leaf]** table in the **MDS** database. On the **Mappings** tab, note the column mappings. Click **Cancel**.
8. Click the **Control Flow** tab, and on the control flow surface, double-click **Load Staged Members**.
9. On the **General** tab, note that the **SqlStatement** property is set to execute the **stg.udp_Customer_Leaf** stored procedure with the values 'VERSION_1', 0, and a parameter. On the **Parameter Mapping** tab, note that the **ExecutionInstanceGUID** system variable is mapped to the **@BatchTag** parameter. Click **Cancel**.
10. On the **Debug** menu, click **Start Debugging**.

- When execution is complete, on the **Debug** menu, click **Stop Debugging** and close Visual Studio without saving any changes.

View Import Status

- Start Internet Explorer and browse to **http://mia-sql:81/mds**
- On the **Master Data Services** home page, click **Integration Management**.
- On the **Integration** page, ensure that the **Customers** model is selected and note the batches listed. There should be a single batch for the data import you performed in the previous task. Note the **Started, Completed, Records, Status** and **Errors** values for this batch.
- On the **Import Data** page, click the **Microsoft SQL Server** text in the page title to return to the home page.

Validate Imported Data

- On the **Master Data Services** home page, click **Explorer**.
- In the **Entities** menu, click **Customer**. Note that nine new member records have been imported, and that their validation status is indicated by a yellow question mark.
- Click **Apply Rules**, and note that the business rules in the model are applied to all records. The validation status for the new records changes to a green tick for records that have passed validation, and an exclamation mark for records that have failed validation (there shouldn't be any failures).
- Close Internet Explorer.

Consuming Master Data with Subscription Views

A subscription view is a standard SQL Server view that enables applications to consume master data. Subscribing applications can use Transact-SQL SELECT statements to access data by using a view, just as they can for tables and views in a typical SQL Server database.

You can create views by using the **Export** page in the Master Data Manager web application. You can use the standard Master Data Services view formats to create views displaying different types of data, including:

- Leaf attributes
- Consolidated attributes
- Collection attributes
- Collections
- Parent/child members in an explicit hierarchy
- Levels in an explicit hierarchy
- Parent/child members in a derived hierarchy
- Levels in a derived hierarchy

- Create subscription views in Master Data Manager
- Formats:
 - Leaf attributes
 - Consolidated attributes
 - Collection attributes
 - Collections
 - Explicit parent child
 - Explicit levels
 - Derived parent child
 - Derived levels

If you modify a data model after creating views that are based on it, you must regenerate the view so that it remains synchronized with the data. You can identify views to be updated by looking at the **Changed** column for each one on the **Subscription Views** page in the Master Data Manager web application. The column shows a value of **True** for any views that are no longer synchronized with the associated model.

Demonstration: Using Subscription Views

In this demonstration, you will see how to use subscription views to consume master data.

Demonstration Steps

Create a Subscription View

1. Ensure you have completed the previous demonstrations in this module.
2. On the **20767C-MIA-SQL** virtual machine, in Internet Explorer, browse to **http://mia-sql:81/mds**
3. On the **Master Data Services** home page, click **Integration Management**.
4. Click **Create Views**.
5. On the **Subscription Views** page, click **Add**.
6. On the **Subscription Views** page, in the **Create Subscription View** pane, in the **Name** box, type **MasterCustomers**. In the **Model** list, click **Customers**, in the **Version** list, click **VERSION_1**, in the **Entity** list, click **Customer**, in the **Format** list, click **Leaf members**, and then click **Save**.
7. Close Internet Explorer.

Query a Subscription View

1. Start SQL Server Management Studio. When prompted, connect to the **MIA-SQL** instance of SQL Server by using Windows authentication.
2. Click **New Query** and enter the following Transact-SQL in the query editor:

```
USE MDS  
  
GO  
  
SELECT * FROM mdm.MasterCustomers
```
3. Click **Execute** to run the query, and review the results.
4. Close SQL Server Management Studio without saving any items.

Check Your Knowledge

Question	
You have a product entity; assuming default table naming was used, in which staging table would you insert leaf member data?	
Select the correct answer.	
<input type="checkbox"/>	dbo.Product_Leaf
<input type="checkbox"/>	stg.Entity_Leaf
<input type="checkbox"/>	stg.product_Leaf
<input type="checkbox"/>	stg.Leaf_Product

Lab: Implementing Master Data Services Model

Scenario

The data warehousing solution you are building for Adventure Works Cycles includes product data from various systems throughout the enterprise. You need to ensure that there is a single, consistent definition for each product.

Objectives

After completing this lab, you will be able to:

- Create a master data model in Master Data Services.
- Use the Master Data Services Add-in for Excel.
- Create and enforce business rules.
- Load data into a master data model.

Estimated Time: 60 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Creating a Master Data Services Model

Scenario

The ETL solution you are building for Adventure Works Cycles consumes product data from an application database. However, product data is created and updated in various systems throughout the enterprise; you need to ensure that there is a single, consistent definition for each product.

To accomplish this, you plan to create a master data hub that includes complete definitions for **Product**, **Product Subcategory**, and **Product Category** entities.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Create a Model
3. Create Entities
4. Create Attributes
5. Add Members

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC**, **20767C-MIA-CLI**, and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab11\Starter** folder as Administrator.

► Task 2: Create a Model

- In the Master Data Manager web application at **http://mia-sql:81/mds**, in the **System Administration** area, create a model named **Product Catalog**. Do not create an entity with the same name.

► Task 3: Create Entities

- In the Product Catalog model, create the following entities. Do not enable explicit hierarchies and collections for these entities:
 - Product
 - Product_Subcategory

► Task 4: Create Attributes

- In the Product entity, create the following leaf attributes:
 - A domain-based attribute named **Product_Subcategory_Key** that references the **Product_Subcategory** entity.
 - A free-form attribute named **Description** that can store a text value with a maximum length of 400.
 - A free-form attribute named **ListPrice** that can store a number with four decimal places and has the input mask (####).

► Task 5: Add Members

1. In the Explorer area of Master Data Manager, view the **Product Subcategory** entity and add the following members:

Name	Code
Mountain Bikes	1
Chains	7
Gloves	20
Helmets	31

2. View the Product entity and add the following members:

Name	Code	ProductSubcategory	Description	ListPrice
Mountain-100 Silver, 42	345	1	Competition Mountain Bike	3399.99
Chain	559	7	Superior Chain	20.24
Full Finger Gloves, S	468	20	Synthetic Gloves	37.99
Sport-100 Helmet, Red	214	31	Universal Fit Helmet	34.99

Results: After this exercise, you should have a Master Data Services model named **Product Catalog** that contains **Product** and **ProductSubcategory** entities. Each of these entities should contain four members.

Exercise 2: Using the Master Data Services Add-in for Excel

Scenario

To make it easier for data stewards to continue to develop the master data model, you plan to use the Master Data Services Add-in for Excel.

The main tasks for this exercise are as follows:

1. Add Free-Form Attributes to an Entity
2. Create an Entity for a Domain-Based Attribute

► Task 1: Add Free-Form Attributes to an Entity

1. Log into the **20767C-MIA-CLI** virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**.
2. Start Microsoft Excel, create a new blank workbook, and enable the Master Data Services add-in.
3. Create a Master Data Services connection named **Local MDS Server** for the MDS server at **http://mia-sql:81/mds**.
4. Connect to the existing **Local MDS Server** and load the **Product** entity from the **Product Catalog** model.
5. Add the following attributes to the **Product** entity by entering the attribute name in the column heading cell (starting with the cell to the right of the existing **ListPrice** attribute heading), and then selecting the new heading cell, clicking **Attribute Properties** in the ribbon, and editing the properties of the attribute:

Attribute Name	Attribute Properties
StandardCost	Number (4 Decimal Places)
Color	Text (Maximum Length: 15)
SafetyStockLevel	Number (0 Decimal Places)
ReorderPoint	Number (0 Decimal Places)
Size	Text (Maximum Length: 50)
Weight	Number (4 Decimal Places)
DaysToManufacture	Number (0 Decimal Places)
ModelName	Text (Maximum Length: 500)

6. After you have added and configured the attribute headings, enter the following attribute values:
- Sport-100 Helmet, Red
 - **StandardCost:** 13.0863
 - **Color:** Red
 - **SafetyStockLevel:** 3
 - **ReorderPoint:** 4
 - **Size:** U
 - **Weight:** 0.2000
 - **DaysToManufacture:** 5
 - **ModelName:** Sport-100
 - Mountain-100 Silver, 42
 - **StandardCost:** 1912.1544
 - **Color:** Silver
 - **SafetyStockLevel:** 100
 - **ReorderPoint:** 75
 - **Size:** L
 - **Weight:** 20.77
 - **DaysToManufacture:** 4
 - **ModelName:** Mountain-100
 - Full-Finger Gloves, S
 - **StandardCost:** 15.6709
 - **Color:** Black
 - **SafetyStockLevel:** 4
 - **ReorderPoint:** 3
 - **Size:** S
 - **Weight:** 0.2000
 - **DaysToManufacture:** 0
 - **ModelName:** Full-Finger Gloves
 - Chain
 - **StandardCost:** 8.9866
 - **Color:** Silver
 - **SafetyStockLevel:** 500
 - **ReorderPoint:** 375
 - **Size:** (leave blank)
 - **Weight:** (leave blank)

- **DaysToManufacture:** 1
 - **ModelName:** Chain
7. Publish the changes you have made to Master Data Services.

► **Task 2: Create an Entity for a Domain-Based Attribute**

1. Connect to the **Local MDS Server** and load data from the **Product_Subcategory** entity in the **Product Catalog** model.
2. Add a new attribute named **Product_Category_Key** to the right of the attribute heading for the **Code** attribute. Do not configure the attribute properties yet.
3. Enter the following **Product_Category_Key** values for the **Product_Subcategory** members:
 - **Mountain Bikes:** 1
 - **Gloves:** 3
 - **Helmets:** 4
 - **Chains:** 2
4. Edit the attribute properties for the **Product_Category_Key** attribute, and configure it as follows:
 - Attribute type: **Constrained list (Domain-based)**.
 - Populate the attribute with values from the selected column.
 - Create a new entity named **Product_Category**.
5. Connect to the **Local MDS Server** connection and load data from the newly created **Product_Category** entity in the **Product_Catalog** model. The **Name** and **Code** values for the entities in this entity have been copied from the distinct **Product_Category_Key** attribute values in the **Product_Subcategory** entity.
6. Change the **Name** value for the **Product_Category** members as follows, and publish the changes:

Name	Code
Bikes	1
Components	2
Clothing	3
Accessories	4

7. On the **Product_Subcategory** worksheet tab, refresh the entity data and verify that the **Product_Category** names are now listed in the **Product_Category_Key** attribute column.

Results: After this exercise, you should have a master data model that contains **Product**, **ProductSubcategory**, and **ProductCategory** entities that contain data entered in Excel.

Exercise 3: Enforcing Business Rules

Scenario

Users noticed inconsistencies in the product data, including missing list prices and invalid safety stock levels. You intend to create business rules that data stewards can use to identify inconsistent data.

The main tasks for this exercise are as follows:

1. Create a Rule for the ListPrice Attribute
2. Create a Rule for the SafetyStockLevel Attribute
3. Publish Business Rules and Validate Data

► Task 1: Create a Rule for the ListPrice Attribute

- In the Master Data Manager web application at <http://mia-sql:81/mds>, in the **System Administration** area, create a business rule called **Validate Price** for the **Product** entity that checks values in the **ListPrice** column are greater than **0**.
 - The business rule does not require an IF clause, only a THEN clause.
 - Use a **must be greater than** validation action that compares the **ListPrice** attribute to the value **0**.

► Task 2: Create a Rule for the SafetyStockLevel Attribute

- Create a rule called **Validate SafetyStockLevel** for the **Product** entity that checks the safety stock level is greater than **ReorderPoint** for products taking a day or more to manufacture.
 - Create an IF condition with an **is greater than** value comparison that checks that the **DaysToManufacture** attribute is greater than **0**.
 - After creating the IF condition, add a **THEN** clause with a **must be greater than** validation action that checks that the **SafetyStockLevel** attribute is greater than the **ReorderPoint** attribute.
 - The completed rule should be based on the following expression:

```
IF DaysToManufacture is greater than 0
THEN SafetyStockLevel must be greater than ReorderPoint
```

► Task 3: Publish Business Rules and Validate Data

1. Publish the business rules.
2. In Master Data Manager, in the **Explorer** area, view the **Product entity** members and click **Apply Rules** to validate the members.
 - Note that the Sport-100 Helmet, Red member is invalid.
 - Fix the invalid product member by changing the SafetyStockLevel attribute value to 5.
3. In Excel, connect to the **Local MDS Server** connection and load the **Product** entity from the **Product Catalog** model.
4. Display the **\$ValidationStatus\$** and **\$InputStatus\$** columns by clicking **Show Status** on the **Master Data** tab of the ribbon, and then click **Apply Rules** to verify that all members are valid.
5. Change the value in the **ListPrice** column for the **Chain** product to **0.00**, and then publish the changes.
6. Note that Excel highlights the invalid member.

7. Fix the member by changing the price back to **20.24** and publishing the changes.

Results: After this exercise, you should have a master data model that includes business rules to validate product data.

Exercise 4: Loading Data into a Model

Scenario

Now that the model is complete, you will populate it by using Transact-SQL scripts.

The main tasks for this exercise are as follows:

1. Load Data into the Model
2. Check the Status of the Data Load
3. Validate Imported Members

► Task 1: Load Data into the Model

1. On the **20767C-MIA-SQL** virtual machine, start SQL Server Management Studio and connect to the **MIA-SQL** database engine instance using Windows authentication.
2. Open the **Import Products into MDS.sql** file in **D:\Labfiles\Lab11\Starter** and examine the Transact-SQL code it contains. The code performs the following tasks:
 - Generates a unique batch ID.
 - Inserts data into the **stg.Product_Category_Leaf** staging table from the **ProductCategory** table in the **Products** database, specifying an **ImportType** value of 2, an **ImportStatus_ID** value of 0, and a **BatchTag** value that is based on the batch ID generated previously.
 - Inserts data into the **stg.Product_Subcategory_Leaf** staging table from the **ProductSubcategory** table in the **Products** database, specifying an **ImportType** value of 2, an **ImportStatus_ID** value of 0, and a **BatchTag** value based on the batch ID generated previously.
 - Inserts data into the **stg.Product_Leaf** staging table from the **Product** table in the **Products** database, specifying an **ImportType** value of 2, an **ImportStatus_ID** value of 0, and a **BatchTag** value based on the batch ID generated previously.
 - Executes the **stg.udp_Product_Category_Leaf** stored procedure to start processing the staged **Product Category** members with the batch tag specified previously.
 - Executes the **stg.udp_Product_Subcategory_Leaf** stored procedure to start processing the staged **Product Subcategory** members with the batch tag specified previously.
 - Executes the **stg.udp_Product_Leaf** stored procedure to start processing the staged **Product** members with the batch tag specified previously.
3. Execute the script to stage the data and start the import process.

► Task 2: Check the Status of the Data Load

- In the Master Data Manager web application, in the **Integration Management** area, verify that the data load process is complete for all entities with no errors.

► Task 3: Validate Imported Members

1. In the Master Data Manager web application, in the **Explorer** area, view the **Product** entity members and apply rules to validate the data.
2. After you have applied the business rules, filter the list of members on the following criteria, so that only invalid records are shown:

Attribute	Operator	Criteria
[Validation Status]	Is equal to	Verification Failed

3. Review the invalid records (which have an invalid **ListPrice** attribute), and resolve them by changing the **ListPrice** attribute value to **1431.50**.
4. When you have resolved all the invalid members, clear the filter.

Results: After this exercise, you should have loaded members into the Master Data Services model.

Module Review and Takeaways

In this module, you have learned how to use Master Data Services to create and manage a master data model for business entity definitions.

Review Question(s)

Question: In your experience, do you think that business users in a data steward capacity will be mostly comfortable using the web-based interface, the Excel add-in, or a combination of both tools to manage master data?

Module 12

Extending SQL Server Integration Services

Contents:

Module Overview	12-1
Lesson 1: Using Scripts in SSIS	12-2
Lesson 2: Using Custom Components in SSIS	12-10
Lab: Using Custom Scripts	12-14
Module Review and Takeaways	12-16

Module Overview

Microsoft® SQL Server® Integration Services (SSIS) provides a comprehensive collection of tasks, connection managers, data flow components, and other items for building an extract, transform, and load (ETL) process for a data warehousing solution. However, there might be some cases where the specific requirements of your organization mean that you must extend SSIS to include some custom functionality.

This module describes the techniques you can use to extend SSIS. The module is not designed as a comprehensive guide to developing custom SSIS solutions, but to show the fundamental steps required to use custom components and scripts in an ETL process, based on SSIS.

Objectives

After completing this module, you will be able to:

- Include custom scripts in an SSIS package.
- Describe how custom components are used to extend SSIS.

Lesson 1

Using Scripts in SSIS

If the built-in SSIS components do not meet your requirements, you might be able to use a script to implement the functionality you need.

This lesson describes the types of functionality you can implement with a script, and how the control flow script task and data flow script component is used in a package.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the types of functionality you can implement with a script.
- Implement a script task in a package control flow.
- Implement a script component in a data flow.

Introduction to Scripting in SSIS

SSIS includes support for adding custom functionality to a package by creating scripts in Microsoft Visual Basic® or Microsoft Visual C#® code. You can add a script to a package, and then configure its properties in the SSIS Designer before opening it in the Microsoft Visual Studio® Tools for Applications (VSTA) development environment to edit it. When you include a script in a package, it is precompiled to optimize execution performance.

SSIS supports the following two scripting scenarios in a package:

- Use of the Script Task to implement a custom control flow task.
- Use of the Script Component to implement a custom data flow source, transformation, or destination.

- Implement custom functionality with Visual Basic or Visual C# code
- Configure properties in the package designer, and develop code using the Visual Studio Tools for Applications (VSTA) environment
- Scripts are precompiled for optimal execution performance
- Two scripting scenarios
 - Use the *Script Task* to implement a custom control flow task
 - Use the *Script Component* to implement a custom data flow source, transformation, or destination

Using the Control Flow Script Task

The script task provides a scriptable wrapper for the **Microsoft.SqlServer.Dts.Runtime.Task** base class—use it to create a custom control flow task without having to create a custom assembly. To use the script task, perform the following steps:

1. Add the script task to the control flow.

To create a scripted custom control flow task, add the script task to a control flow and configure its standard properties, such as **Name**, in the Properties pane. To use the script task to write information to a log file, you must also enable logging for the task and specify the events to be logged. The script task supports the same events as other control flow tasks, in addition to a **ScriptTaskLogEntry** event, which must be enabled to use the script to write log entries programmatically.

2. Configure the task.

Use the Script Task Editor dialog box to configure the following settings for the task:

- **ScriptLanguage:** the programming language in which you write the script. For example, Microsoft Visual Basic or Microsoft Visual C#.
- **EntryPoint:** the method in the script that should be called when the script is executed. By default, this is a method named Main.
- **ReadOnlyVariables** and **ReadWriteVariables:** the package variables and parameters that you want to make accessible to the script. By default, the script task cannot access any package variables, so you must explicitly grant access to any variables you want to use in the script.

3. Implement the code.

After you have configured the script task, click **Edit Script** to open the VSTA environment and implement the code for your script. You might use standard .NET programming structures and syntax in your code, adding references to additional assemblies if required. You can also make use of the **Dts** object, which provides a number of static properties and methods to help you access package and SSIS runtime functionality.

For example, the following code uses the **Dts.Variables** property to access package variables; the **Dts.Log** method to write a custom log entry; and the **Dts.TaskResult** property to specify the control flow result for the script task:

Dts.Variables

```
Public void Main()
{
    Dts.Variables["User::MyVar"].Value = Dts.Variables["System::StartTime"].Value.ToString();
    Dts.Log("Package started:" + (string)Dts.Variables["User::MyVar"].Value, 999, null);
    Dts.TaskResult = (int)ScriptResults.Success;
}
```



Extending the Package with the Script Task

<http://aka.ms/xq34g7>

1. Add the task to the control flow
 - Configure properties and logging settings as required
2. Configure the task
 - Language
 - Entry point method
 - Access to variables
3. Implement the code
 - Use the **Dts** object to access package and runtime functionality

```
Public void Main()
{
    Dts.Variables["User::MyVar"].Value =
        Dts.Variables["System::StartTime"].Value.ToString();
    Dts.Log("Package started:" +
        (string)Dts.Variables["User::MyVar"].Value, 999, null);
    Dts.TaskResult = (int)ScriptResults.Success;
}
```

Demonstration: Implementing a Script Task

In this demonstration, you will see how to:

- Configure a script task.
- Implement a script task.

Demonstration Steps

Configure a Script Task

1. Ensure that **20767C-MIA-DC** and **20767C-MIA-SQL** are started, and log onto **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Start Visual Studio and open the **ScriptDemo.sln** solution in the **D:\Demofiles\Mod12** folder.
3. In Solution Explorer, double-click the **Package.dtsx** SSIS package to open it in the SSIS designer. Notice that the controls flow includes a script task.
4. On the **SSIS** menu, click **Logging**.
5. In the **Configure SSIS** Logs dialog box, on the **Providers and Logs** tab, in the **Provider type** list, note that the **SSIS log provider for Windows Event Log** has already been added.
6. In the **Containers** tree, select the **Script Task** check box.
7. On the **Providers and Logs** tab, ensure the **SSIS log provider for Windows Event Log** check box is selected.
8. On the **Details** tab, ensure the **ScriptTaskLogEntry** check box is selected, and click **OK**.
9. On the control flow surface, double-click **Script Task**.
10. In the **Script Task Editor** dialog box, verify that the following settings are selected:
 - **ScriptLanguage:** Microsoft Visual C#
 - **EntryPoint:** Main
 - **ReadOnlyVariables:** User::Results and System::StartTime
 - **ReadWriteVariables:** User::MyVar

Implement a Script Task

1. In the **Script Task Editor** dialog box, click **Edit Script**, and wait for the Visual Studio VstaProjects editor to open.
2. In the VstaProjects – Microsoft Visual Studio window, scroll through the **ScriptMain.cs** code until you can see the whole of the **public void Main()** method. Note that the code performs the following tasks:
 - Assigns the value of the System::StartTime variable to the User::MyVar variable.
 - Writes the value of the User::MyVar variable to the log.
 - Creates a string that contains the values in the User::Results variable (which is an array of strings), and displays it in a message box.
3. Close the VstaProjects – Microsoft Visual Studio window.
4. In the **Script Task Editor** dialog box, click **OK**.
5. On the **Debug** menu, click **Start Debugging**, and observe the package as it executes. When the **Results** message box is displayed, click **OK**.

6. When package execution has completed, on the **Debug** menu, click **Stop Debugging**, and then minimize Visual Studio.
7. Right-click the **Start** button and click **Event Viewer**.
8. In Event Viewer, expand the **Windows Logs** folder, click **Application**, select the second Information event with the source **SQLISPackage130**, and on the **Details** tab, note that it contains the start time of the package, which was written by the script component.
9. Close Event Viewer.

Using the Data Flow Script Component

To implement a data flow component without the need to create a custom assembly, the script component uses classes in the **Microsoft.SqlServer.Dts.Pipeline**, **Microsoft.SqlServer.Dts.Pipeline Wrapper**, and **Microsoft.SqlServer.Dts.Runtime Wrapper** namespaces.

Using the Script Component to Create a Source

To create a custom source, add the script component to a data flow and, when prompted, select **Source**. Then, as a minimum, you must perform the following actions:

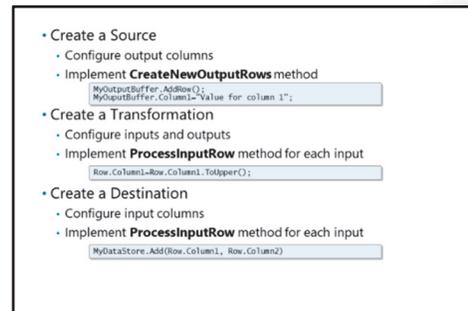
- Configure the output columns. A source must provide at least one output, and each output must contain at least one column. You must specify the name and data type for each column that your source will pass to the next component in the data flow.
- Implement the **CreateNewOutputRows** method. To generate a new row and assign column values, each output provides a buffer object.

Using the Script Component to Create a Transformation

To create a custom transformation, add the script component to a data flow, and when prompted, select **Transformation**. Then, as a minimum, you must perform the following actions:

- Configure inputs and outputs. A transformation can have multiple inputs and outputs, and you must define the columns for each one as required.
- Implement the **ProcessInputRow** method for each input. The script component generates a **ProcessInputRow** method for each input you define. The name of the method is derived from the input name. For example, an input named **Input0** will result in a method named **Input0_ProcessInputRow**. The method is called for each row in the input, and you must add code to process the columns in the row. If your component has multiple outputs, direct them to the appropriate output.

Using the Script Component to Create a Destination



To create a custom destination, add the script component to a data flow, and when prompted, select **Destination**. Then, as a minimum, you must perform the following actions:

- Configure input columns. A destination must have at least one input, and each input must contain at least one column. You must specify the name and data type for each column that your destination will accept from the preceding component in the data flow.
- Implement the **ProcessInputRow** method for each input. Like transformations, the script component generates a **ProcessInputRow** method for each input you define in a destination. The method is called for each row in the input, and you must add code to process a row, usually by adding it to a data store.

Additional Configuration Settings and Methods

In addition to the tasks described above, the following configuration settings and methods are available to enhance the functionality of your scripted data flow component:

- Add variables: for access to variables in a script component, set the **ReadOnlyVariables** and **ReadWriteVariables** properties.
- Add a connection manager: add one or more connection managers to a script component, and use the connection manager in your code to access a database or other data store. For example, you might add a connection manager to a custom source, to retrieve the data for the output from a database.
- Override the **AcquireConnections** method: to open a connection from a connection manager in preparation for execution.
- Override the **PreExecute** method: to perform any preparation tasks for the component. For example, you might use this method in a custom source to retrieve a **SqlDataReader** object from a connection you opened in the **AcquireConnections** method.
- Override the **PostExecute** method: to perform any clean-up tasks when execution is complete. For example, you might use this method to dispose of a **SqlDataReader** object you opened during the **PreExecute** method.
- Override the **ReleaseConnections** method: to close any connections you opened in the **AcquireConnections** method.



Extending the Data Flow with the Script Component

<http://aka.ms/oge7wg>

Demonstration: Using a Script Component in a Data Flow

In this demonstration, you will see how to:

- Implement a Source
- Implement a Transformation
- Implement a Destination

Demonstration Steps

Implement a Source

1. Ensure you have completed the previous demonstration in this module.
2. Return to Visual Studio, and view the control flow of the **Package.dtsx** SSIS package.
3. On the control flow surface, double-click **Data Flow Task**.
4. On the data flow surface, double-click **Script Source**.
5. In the **Script Transformation Editor** dialog box, on the **Inputs and Outputs** page, expand **MyOutput** and the **Output Columns** folder.
6. Select **Column1** and **Column2** in turn and note the **DataType** property of each column.
7. In the **Script Transformation Editor** dialog box, on the **Script** page, click **Edit Script**, and wait for the Visual Studio VstaProjects editor to open.
8. In the VstaProjects – Microsoft Visual Studio window, view the script. Note that the **CreateNewOutputRows** method uses a loop to create six rows of data.
9. Close the VstaProjects – Microsoft Visual Studio window, and in the **Script Transformation Editor** dialog box, click **Cancel**.

Implement a Transformation

1. On the data flow surface, double-click **Script Transformation**.
2. In the **Script Transformation Editor** dialog box, on the **Inputs and Outputs** page, note that the component has a single input named **Input 0**, which consists of two columns named Column1 and Column2. There is no output defined for the component, which means that it will pass the columns in the input buffer to the output buffer without adding any new columns.
3. In the **Script Transformation Editor** dialog box, on the **Input Columns** page, note that both Column1 and Column2 are selected and that the Usage Type for Column1 is ReadOnly, while the Usage Type for Column2 is ReadWrite. This configuration enables the script to make changes to Column2 as it flows through the pipeline.
4. In the **Script Transformation Editor** dialog box, on the **Script** page, click **Edit Script**, and wait for the Visual Studio VstaProjects editor to open.
5. In the VstaProjects – Microsoft Visual Studio window, view the script. Note that the **Input0_ProcessInoutRow** method modifies Column2 by making its value upper case.
6. Close the VstaProjects – Microsoft Visual Studio window, and in the **Script Transformation Editor** dialog box, click **Cancel**.

Implement a Destination

1. On the data flow surface, double-click **Script Destination**.
2. In the **Script Transformation Editor** dialog box, on the **Inputs and Outputs** page, note that the component has a single input named Input 0, which consists of two columns named Column1 and Column2.
3. In the **Script Transformation Editor** dialog box, on the **Input Columns** page, note that both Column1 and Column2 are selected, and that the Usage Type for both columns is ReadOnly. Because the component represents a destination, there is no need to modify the rows in the buffers.
4. In the **Script Transformation Editor** dialog box, on the **Script** page, note that the ReadWriteVariables property gives access to the User::Results variable.
5. Click **Edit Script** and wait for the Visual Studio VstaProjects editor to open.
6. In the VstaProjects – Microsoft Visual Studio window, view the script, and note the following details:
 - o The **PreExecute** method initializes an array variable for six string elements.
 - o The **Input0_ProcessInputRow** method adds the value of Column2 to the next available empty element in the array.
 - o The **PostExecute** method assigns the array variable in the script to the User::Results package variable.
7. Close the VstaProjects – Microsoft Visual Studio window, and in the **Script Transformation Editor** dialog box, click **Cancel**.
8. On the data flow surface, right-click the data flow path between **Script Source** and **Script Transformation**, and click **Enable Data Viewer**.
9. Repeat the previous step for the data flow path between **Script Transformation** and **Script Destination**.
10. On the **Debug** menu, click **Start Debugging**.
11. When the first data viewer window is displayed, view the rows in the pipeline, and then click the green **Continue** button. These rows were generated by the **Script Source** component.
12. When the second data viewer window is displayed, view the rows in the pipeline, and then click the green **Continue** button. Note that Column2 was formatted as upper case by the **Script Transformation** component.
13. When the **Results** message box is displayed, note that it contains the Column2 values that were passed to the Script Destination component. You may need to click the program icon on the taskbar to bring the message box to the front of the open windows.
14. When package execution has completed, on the **Debug** menu, click **Stop Debugging**.
15. Close Visual Studio.

Check Your Knowledge

Question	
You are creating a control flow script task. How do you ensure your script task can write log entries programmatically?	
Select the correct answer.	
<input type="checkbox"/>	Enable logging, and then enable the onError event handler.
<input type="checkbox"/>	No special action is needed as all control flow tasks can write log entries.
<input type="checkbox"/>	Enable logging, and then enable the ScriptTaskLogEntry event.
<input type="checkbox"/>	Write code within the script to write to a log.

Lesson 2

Using Custom Components in SSIS

When the built-in functionality in SSIS does not fully meet your requirements, and you cannot achieve your goal by using a script component, you can obtain or create custom components, add them to the SSIS package designer, and use them in your packages. This lesson describes the types of functionality that you can implement with custom components, and how to create, install, and use them.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe how custom components are used to extend SSIS.
- Describe how to implement a custom component.
- Install and use a custom component.

Introduction to Custom Components

The control flow tasks, connection managers, data flow transformations, and other items in an SSIS package are all managed components. In other words, each item used when creating an SSIS package is implemented as a Microsoft .NET assembly.

SSIS has been designed as an extensible platform, where you can supplement built-in components with custom assemblies that perform tasks that cannot otherwise be achieved. You can extend SSIS with the following kinds of custom components:

- Control flow tasks
- Connection managers
- Log providers
- Enumerators
- Data flow components, including sources, transformations, and destinations

You can develop your own custom SSIS components in a programming tool that targets the .NET Framework, such as Microsoft Visual Studio®. Alternatively, you can commission a professional software developer to create custom components for you, or obtain them from third-party vendors.

Extend SSIS package functionality by using custom:

- Tasks
- Connection managers
- Log providers
- Enumerators
- Data flow components

Implementing a Custom Component

Software developers can use the following process to implement a custom component for SSIS:

1. Create a class library project and add references to the relevant SSIS assemblies. The SSIS assemblies, provided with the SQL Server Client SDK, give access to SSIS task flow and pipeline runtime functionality. This means you can add a reference to them in your custom component project. The specific assemblies you need to reference depend on the type of custom component you are implementing, as described in the following list:

- **Microsoft.SqlServer.ManagedDTS.dll:** the managed runtime engine, used by managed control flow tasks, enumerators, log providers, and connection managers.
- **Microsoft.SqlServer.PipelineHost.dll:** the managed data flow engine, used by managed data flow sources, transformations, and destinations.
- **Microsoft.SqlServer.RuntimeWrapper.dll:** the primary interop assembly (PIA) for the native SSIS runtime engine.
- **Microsoft.SqlServer.PipelineWrapper.dll:** the PIA for the native SSIS data flow engine.

The base class assemblies include namespaces that contain the base classes your custom components must implement. For example, the Microsoft.SqlServer.ManagedDTS.dll includes the Microsoft.SqlServer.Dts.Runtime namespace, which contains the Task base class for control flow tasks.

2. Create a class that inherits from the appropriate SSIS base class, as listed in the following table:

Component Type	Base Class
Task	Microsoft.SqlServer.Dts.Runtime.Task
Connection manager	Microsoft.SqlServer.Dts.Runtime.ConnectionManagerBase
Log provider	Microsoft.SqlServer.Dts.Runtime.LogProviderBase
Enumerator	Microsoft.SqlServer.Dts.Runtime.EachEnumerator
Data flow component	Microsoft.SqlServer.Dts.Pipeline.PipelineComponent

- Add references to required assemblies and declare namespaces
- Use attributes to provide design-time information
- Inherit from the appropriate base class
- Override the base class methods
- Sign the assembly

```
using Microsoft.SqlServer.Dts.Runtime;
namespace MyCustomSSISComponent
{
    [DtsTask(DisplayName="My Custom Task",
            UTypeDisplayName="My Custom Task",
            RequiredProductLevel=
                DTSProductLevel.None)]
    public class MyCustomTask : Task
    {
        public override DTSExecResult Execute
            (Connections ,
             VariableDispenser variableDispenser,
             IDTSComponentEvents componentEvents,
             IDTSLogging log, object transaction)
        {
            //Code to execute the task goes here
        }
    }
}
```

3. Add attributes to the class declaration to provide design-time information about the component that can be displayed in SQL Server Data Tools when using the component in an SSIS package. The class-specific attributes are listed in the following table:

Base Class	Attribute
Microsoft.SqlServer.Dts.Runtime.Task	DtsTaskAttribute
Microsoft.SqlServer.Dts.Runtime.ConnectionManagerBase	DtsConnectionAttribute
Microsoft.SqlServer.Dts.Runtime.LogProviderBase	DtsLogProviderAttribute
Microsoft.SqlServer.Dts.Runtime.ForEachEnumerator	DtsForEachEnumeratorAttribute
Microsoft.SqlServer.Dts.Pipeline.PipelineComponent	DtsPipelineComponentAttribute

4. Override the base class methods to implement the required custom functionality. Each base class includes methods that are called by the runtime engine when the component is used in an SSIS package. Some of the most important methods for each base class are listed in the following table:

Base Class	Methods
Microsoft.SqlServer.Dts.Runtime.Task	Validate, Execute
Microsoft.SqlServer.Dts.Runtime.ConnectionManagerBase	AcquireConnection, ReleaseConnection
Microsoft.SqlServer.Dts.Runtime.LogProviderBase	OpenLog, Log, CloseLog
Microsoft.SqlServer.Dts.Runtime.ForEachEnumerator	GetEnumerator
Microsoft.SqlServer.Dts.Pipeline.PipelineComponent	PipeLineComponentProperties, PrimeOutput, ProcessInput

5. Sign the assembly to generate a strong name. Custom SSIS component assemblies are added to the global assembly cache (GAC), which requires that each assembly has a strong name that uniquely identifies it. A strong name for an assembly is generated by signing the assembly with a cryptographic key when it is compiled. You can easily generate a key and sign an assembly on the Signing tab of the Project Properties dialog box in Visual Studio.



Extending Packages with Custom Objects

<http://aka.ms/raebey>

Installing and Using a Custom Component

Most commercially available custom SSIS components provide a setup program that installs the assemblies— you can then use them in SQL Server Data Tools. However, as a business intelligence (BI) professional, you should be able to perform the necessary manual steps to install a custom component you have developed or have had created for you by a software developer.

You can use the following procedure to install and use a custom SSIS component:

1. Copy the assembly to the DTS folder
2. Install the assembly in the global assembly cache
3. Use the custom component in an SSIS package

1. Copy the assembly file to the appropriate Data Transformation Services (DTS) subfolder.

The SSIS Designer looks for components in a specific subfolder of the DTS folder for the SQL Server installation. By default, the DTS folder is located in C:\Program Files\Microsoft SQL Server\<version>\DTS. For 64-bit installations, this folder is used for 64-bit components, and 32-bit components are stored in subfolders of C:\Program Files (x86)\Microsoft SQL Server\<version>\DTS. The component type-specific subfolders are listed in the following table:

Component Type	Folder
Task	Tasks
Connection manager	Connections
Log provider	LogProviders
Enumerator	ForEachEnumerators
Data flow component	PipelineComponents

2. Install the assembly in the GAC. You can do this by using the **gacutil.exe** command-line tool. The gacutil.exe tool is provided with Visual Studio and Windows® SDK.

You can use the following syntax to install an assembly in the GAC:

Install the assembly in the GAC

```
gacutil /i assembly_name.dll
```

3. Use the custom component in an SSIS package. After you have deployed a custom component to the correct folder and installed it in the GAC, it will be available in SSIS Designer, for use in a package. For example, if you have successfully deployed an assembly for a custom control flow task component, the task will be listed in the SSIS Toolbox pane and can be moved to the control flow surface of a package.

Question: A software developer has developed a custom SSIS component for you that does not have its own installer. How would you install this component so it was available in SQL Server Data Tools?

Lab: Using Custom Scripts

Scenario

The senior database administrator has requested that the data warehouse ETL process logs the number of rows extracted from the staging database and loaded into the data warehouse in the Windows event log. You have decided to use a custom script to accomplish this task.

Objectives

After completing this lab, you will be able to:

- Use a Script task.

Estimated Time: 30 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Using a Script Task

Scenario

You have created an SSIS package that loads data from the staging database into the data warehouse, and writes the number of rows processed in each staging table to a package variable before truncating the staging tables. You want to add a script task that writes the row count variables to the Windows event log when the data warehouse load operation is complete.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Add a Script Task to a Control Flow
3. Enable Logging for the Script Task
4. Configure the Script Task
5. Test the Script

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab12\Starter** folder as Administrator.

► Task 2: Add a Script Task to a Control Flow

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab12\Starter** folder.
2. Open the **Load DW.dtsx** package and review the control flow.
3. View the variables defined in the package. These are used by the **Get record counts and truncate Staging tables** SQL command task to store the row counts for each staging table that was processed during the data warehouse load operation.
4. Add a script task to the control flow and name it **Log Rowcounts** then connect the success precedence arrow from the **Get record counts and truncate Staging tables** task to the **Log Rowcounts** task.

► Task 3: Enable Logging for the Script Task

- Configure logging for the package as follows:
 - Add the **SSIS Log Provider for Windows Event Log** provider.
 - Enable logging for the **Log Rowcounts** script task and select the **SSIS Log Provider for Windows Event Log** provider.
 - Enable logging for the **ScriptTaskLogEntry** event of the **Log Rowcounts** task.

► Task 4: Configure the Script Task

1. Configure the **Log Rowcounts** script task with the following settings:
 - ScriptLanguage: **Microsoft Visual C#**
 - EntryPoint: **Main**
 - ReadOnlyVariables: **User::CustomerCount, User::EmployeeCount, User::InternetSalesCount, User::PaymentCount, User::ResellerCount, and User::ResellerSalesCount**
 - ReadWriteVariables: **None**
2. Edit the script, and replace the comment **//TODO: Add your code here** with the following code (above the existing **Dts.TaskResult = (int)ScriptResults.Success** statement). You can copy and paste this code from **Script.txt** in the **D:\Labfiles\Lab12\Starter** folder:

```
String logEntry = "Data Warehouse records loaded (";
logEntry += Dts.Variables["User::CustomerCount"].Value.ToString() + " customers, ";
logEntry += Dts.Variables["User::ResellerCount"].Value.ToString() + " resellers, ";
logEntry += Dts.Variables["User::EmployeeCount"].Value.ToString() + " employees, ";
logEntry += Dts.Variables["User::PaymentCount"].Value.ToString() + " payments, ";
logEntry += Dts.Variables["User::InternetSalesCount"].Value.ToString() + " Internet
sales, and ";
logEntry += Dts.Variables["User::ResellerSalesCount"].Value.ToString() + " reseller
sales) ";
Dts.Log(logEntry, 999, null);
```

► Task 5: Test the Script

1. Start debugging the **Load DW** package and observe the control flow tab as it executes, noting that each package executed opens in a new window. The entire load process can take a few minutes.
2. When execution is complete, stop debugging and close Visual Studio.
3. View the Windows Application log in the Event Viewer administrative tool and verify that the second most recent information event with a source of **SQLISPackage120** contains the row counts for each staged table.

Results: After this exercise, you should have an SSIS package that uses a script task to log row counts.

Module Review and Takeaways

In this module, you have learned how to extend SSIS with custom script tasks and custom components.

Review Question(s)

Question: What might you consider when deciding whether to implement a custom process as a script or a custom component?

Module 13

Deploying and Configuring SSIS Packages

Contents:

Module Overview	13-1
Lesson 1: Overview of SSIS Development	13-2
Lesson 2: Deploying SSIS Projects	13-5
Lesson 3: Planning SSIS Package Execution	13-14
Lab: Deploying and Configuring SSIS Packages	13-20
Module Review and Takeaways	13-24

Module Overview

Microsoft® SQL Server® Integration Services (SSIS) provides tools for you to deploy packages to another computer. The deployment tools also manage any dependencies, such as configurations or files that the package needs. In this module, you will learn how to use these tools to install packages and their dependencies on a target computer.

Objectives

After completing this module, you will be able to:

- Describe considerations for SSIS deployment.
- Deploy SSIS projects.
- Plan SSIS package execution.

Lesson 1

Overview of SSIS Development

SQL Server supports two deployment models for SSIS packages; choosing the right one can greatly simplify the maintenance and operation of your extract, transform, and load (ETL) solution. This lesson discusses the two deployment models and compares their key features.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the SSIS deployment model infrastructure.
- Describe the project deployment model.
- Explain the package deployment model and the incremental feature.
- Compare the two deployment models.

SSIS Deployment Models

SSIS provides two deployment models, one for packages and another for projects.

The Package Deployment Model

A feature of previous releases of SSIS, this model deploys and manages packages individually. It provides support and continued development of existing SSIS packages from previous versions of SQL Server.

The Project Deployment Model

You use the project deployment model to deploy a project that contains multiple packages as a single unit. You can then manage connection managers and parameters that are shared by all packages in the project. When using the project deployment model, you can create an SSIS catalog on an instance of SQL Server and deploy the project to the SSIS catalog.

Options for Deployment Destination

- Package Deployment Module
 - Deploy to SSIS catalog on an instance of SQL Server.
 - Deploy to file system.
- Project Deployment Module
 - Deploy to SSIS catalog on an instance of SQL Server. You can create an SSIS catalog on an instance of SQL Server and deploy your project to the SSIS catalog.

- Package Deployment Model
 - One or more SSIS Packages are deployed at the same time
- Project Deployment Model
 - Multiple packages are deployed in a single project
- Destination for Deployment
 - Options depend on deployment model

Package Deployment Model

Deployment of Multiple Packages

In previous versions of SQL Server, when using the Package Deployment Model, you could deploy only one package at a time. You can deploy one or more packages at the same time.

Storage

With the package deployment model, packages can be deployed to an instance of SQL Server or to the file system.

When you deploy packages to an instance of Integration Services, you must decide whether to store the packages in SQL Server or in the file system.

In previous versions, SQL Server stored Integration Services packages in the **msdb** system database, and package dependencies in a file folder. You can store packages in **SSISDB**, which is a dedicated database within SQL Server. When you deploy a package to SQL Server, you must specify the SQL Server instance and type of authentication.

When you deploy a package to the file system (as a .dtsx file), Integration Services accesses the package files and their dependencies from the folder where you deployed the packages. Alternatively, you can deploy packages to the SSIS package store. The SSIS package stores are directories related to the SSIS installation.

Package Configurations

A package configuration is a set of property/value pairs added to a completed package to dynamically access property values at run time.

A package can retrieve property settings from a number of different source types, such as an XML file or Microsoft SQL Server database. You can update settings on any executable (including the package object), connection manager, log provider, or variable within a package.

When Integration Services runs the package, it extracts the property value from a source that is based on the configuration type. By using package configurations in projects you intend to deploy with the package deployment model, you can define values at run time. Package configurations make it easier to move packages from development to the production environment.

To create a package configuration, open the package, and on the SSIS menu, click **Package Configurations**. You can then select where you want to save the configuration settings (for example, an XML file) and choose the properties to be included in the configuration.

Package Deployment Utilities

You can simplify package deployment and any configuration files or other resources used by creating a package deployment utility. To do this, first configure the package as required, including any package configurations, and then, in the properties dialog box for the project, on the Deployment page, set the **CreateDeploymentUtility** property to **True**. When you build the project, a folder that contains the packages, all the required files, and a manifest, is generated.

- Storage
 - SSISDB
 - File System
- Package Configurations
 - Property values to be set dynamically at run time
- Package Deployment Utility
 - Generate all required files for easier deployment

Project Deployment Model

Many features, which are only available with the project deployment model, make deployment, security, and execution more straightforward.

The SSIS Catalog

The SSIS catalog is a storage repository for SSIS projects. The catalog contains all folders, projects, packages, environments, and parameters for projects using the project deployment model.

You have to create the catalog before you can use it. First, you must enable CLR integration on the server, then you can right-click **Integration**

Services in SQL Server Management Studio and create the catalog. You can define how data is encrypted by editing the properties in SQL Server Management Studio. You can also define whether logs are cleared periodically, and if so, how often. The maximum time to perform project validation can also be defined.

The catalog is stored as the **SSISDB** database in SQL Server. Although you can perform some administration on this database, and it is useful for listing catalog-specific views and stored procedures, most administration is performed on the SSISDB node inside the Integration Services node.

When you deploy a project with the same name as an existing one, it will be stored as a new version. You can set catalog properties to clean up old versions and define how many versions are kept.

Folders

Folders are used to both organize and secure projects. Putting projects in folders makes it more straightforward to navigate to a specific project. Folders are also securable objects so you can assign someone administrator rights to projects in specific folders without giving that person any rights to projects situated outside them.

Deployment Model Comparison

The project deployment and package deployment models have key differences in terms of:

- Unit of deployment
- Storage
- Dynamic configuration
- Compiled format
- Troubleshooting

The slide for this topic summarizes these differences.

Question: How many different deployment models are available when using SSIS in SQL Server? Give the name and a brief description of each of the models.

- The SSIS catalog
 - Storage and management for SSIS projects on an SQL Server instance
- Folders
 - A hierarchical structure for organizing and securing SSIS projects

Feature	Package Deployment Model	Project Deployment Model
Unit of deployment	One or more packages	Project
Storage	Packages and all associated files can be copied to the file systems of a local or remote computer. They can also be deployed to the SSIS catalog, or a folder within the catalog, of an instance of SQL Server.	A project is deployed to the SSIS catalog, or a folder within the catalog, of an instance of SQL Server.
Dynamic configuration	Both projects and packages can contain parameters and references to environments. You can access the parameters and environment references by using the Configure Dialog box.	Both projects and packages can contain parameters and references to environments. You can access the parameters and environment references by using the Configure Dialog box.
Compiled format	Packages and associated resources are each stored as single files in the file system. The entire project might consist of many files.	The entire project is compiled as a single file (with an .ispac extension).
Troubleshooting	To log events, you have to add a log provider to the package and configure logging for each one individually.	Events are automatically logged and saved to the catalog. These events can then be displayed with views such as catalog execution and catalog event messages.

Lesson 2

Deploying SSIS Projects

In this lesson, you will learn how to deploy and manage an SSIS project.

Lesson Objectives

After completing this lesson, you will be able to:

- Create an SSIS Catalog.
- Use environments and variables to set dynamic values at run time.
- Deploy an SSIS Project.
- View package execution information.

Creating an SSIS Catalog

Prerequisites

You must be running SQL Server 2012 or later.

The first step in deploying an SSIS project is to create an SSIS catalog. However, an SSIS catalog can only be created on a SQL Server instance in which the SQL CLR feature is enabled—if you have not enabled the SQL CLR, this will be done automatically when you create the SSIS catalog.

- Prerequisites
 - SQL Server 2012 or later
 - SQL CLR enabled
- Creating a catalog
 - Use SQL Server Management Studio
 - One SSIS catalog per SQL Server instance
- Catalog Security
 - Folder Security
 - Object Security
 - Catalog Encryption
 - Sensitive Parameters

Creating an SSIS Catalog

To create an SSIS catalog, use SQL Server Management Studio to connect to the SQL Server instance on which you want to deploy SSIS projects, and, in Object Explorer, right-click **Integration Services Catalogs**, and click **Create Catalog**. You must then provide a secure password to create the database master key for the new catalog.

Catalog Security

Only members of the **ssis_admin** or **sysadmin** roles can manage the catalog. All security in the catalog uses Windows® authentication—SQL Server authentication is not supported.

To avoid having to give all SSIS administrators full administration rights, you can grant individuals or roles the **MANAGE_OBJECT_PERMISSIONS** permission on any folders that you want them to administer.

If you right-click a folder, project, or environment, and click **Properties**, you can access the permissions on that object. This will allow you to grant or deny permissions on actions such as **Read** or **Modify**, depending on the object type.

You can also use Transact-SQL to manage object permissions using stored procedures such as **catalog.grant_permission** and views such as **catalog.effective_object_permissions**.

For more information on the SSIS catalog, see *SSIS Catalog* in MSDN (Microsoft Developer Network):

SSIS Catalog

<http://aka.ms/b7xql2>

When you create a catalog, a database master key is automatically created. Using the catalog properties, you can define the encryption type from DES, TRIPLE_DES, TRIPLE_DES_3KEY, DESX, AES_128, AES_192, or AES_256 (the default).

Parameters have a sensitive property. If a parameter stores sensitive data, you should set the sensitive property to TRUE. This causes the parameter value to be encrypted and, if the parameter value is queried directly with Transact-SQL, NULL is returned.

For more information on SSIS security, see *Integration Services Roles (SSIS Service)* in MSDN:



Integration Services Roles (SSIS Service)

<http://aka.ms/bfhrp5>

Environments and Variables

Often, you have to change specific properties or values used in a package dynamically at run time. To do this, SSIS catalog supports multiple environments for each project, in addition to the creation of variables within them that you can map to project parameters and connection managers.

Environments

An environment is an execution context for an SSIS project. You can create multiple environments, and then reference them from a project you have deployed in the SSIS catalog. A project can have multiple environment references, but each time a package is executed, it can use only one. This gives you the ability to switch very quickly from one environment to another and change all associated environment settings from a single property. For example, you could create **Test** and **Production** environments and add references to them both in a deployed project.

Variables

Environments are used to organize and group environment variables that you can map to project-level parameters or properties of project-level connection managers. For example, in the **Test** and **Production** environments, you could create a variable named **DBServer** with the value TEST_SERVER in the **Test** environment and PROD_SERVER in the **Production** environment. You can then map this variable to a project-level parameter or to the **ServerName** property of a connection manager in a deployed project. When you execute a package in the project, the value used for the parameter or **ServerName** property is determined by the selected environment.

- Environments
 - Execution contexts for projects
- Variables
 - Environment-specific values
 - Can be mapped to project parameters and connection manager properties at run time

Deploying an SSIS Project

After you have built your SSIS project, you can use one of two methods to deploy it to your catalog. You can use either SQL Server Data Tools or SQL Server Management Studio.

If you deploy your project to the same location with the same name, it will overwrite the original. Based on the properties of your catalog, it might create an additional version of the project and keep a copy of the original.

- Integration Services Deployment Wizard
 - Visual Studio
 - SQL Server Management Studio
- dtutil.exe to deploy, move, delete, and check the existence of SSIS packages

Deploying an SSIS Project Using SQL Server Data Tools (Visual Studio®)

After building your project in SQL Server Data Tools, you can deploy it by using the Integration Services Deployment Wizard to guide you through the process.

The deployment source is the project that you are currently deploying. The deployment destination is the location where you want to deploy the package, including the server and folder.

Deploying an SSIS Project in SQL Server Management Studio

In SQL Server Management Studio, you can navigate to the folder where you wish to deploy the package. You can then right-click the **Projects** folder within the folder you wish to deploy to and click **Deploy Project**. This will also run the Integration Services Deployment Wizard. The only difference with this method is that you have to find the project deployment file. This will be in the folder created for the Visual Studio project in the **bin** folder, and then in the **Development** folder. It will have an **.ispac** extension.

You can also deploy an existing package within the catalog, or a catalog on another server, to make a copy of a package. Simply choose the existing package as the source in the Integration Services Deployment Wizard.

Deploying an SSIS Package with the dtutil Command Line Executable

You can use the dtutil.exe command line utility to manage your SSIS packages in several ways. In addition to deploying packages, you can check the existence of a package, move, and delete packages. The utility works with packages stored in the SSIS Package Store, msdb, and the file system. When connecting to msdb, you might be asked for your credentials. If the SQL Server is using SQL Server Authentication, you need to provide a username and password, though if these are omitted, dtutil attempts to log in using Windows authentication. When using the utility, you must adhere to the following rules:

- All arguments are string values, which must be enclosed in quotation marks if white space is included.
- Single quotes within strings must be finished off by enclosing them in double quotes.
- Options can be passed in any order, using the pipe character as the OR operator to show possible values.
- All options must start with a forward slash or minus sign. No space must exist between the slash or minus sign and the option text.
- Commands are case insensitive.

For full syntax, examples, and exit code descriptions, see the following article:



dtutil Utility

<https://aka.ms/e96n7q>

Other methods for deploying an SSIS Project are:

- Calling the Integration Services Deployment Wizard from the command line.
- Using the **deploy_packages** stored procedure.
- Using the Management Object Model API.

For more information on deploying an SSIS project, see *Deploy Packages to Integration Services Server* in MSDN:



Deploy Packages to Integration Services Server

<http://aka.ms/vleim0>

Viewing Project Execution Information

SSIS includes a number of tools to troubleshoot packages after deployment. During deployment, you can use various techniques, such as breakpoints and data viewers, which are not available after you have deployed the package. However, you can use other tools for troubleshooting. The Integration Services Dashboard should be the starting point for most debugging as it provides every level of detail for all packages, but other tools can perform specific tasks.

- Integration Services Dashboard provides built-in reports
- Additional sources of information:
 - Event handlers
 - Error outputs
 - Logging
 - Debug dump files

Integration Services Dashboard

The Integration Services Dashboard provides details for each package that has run on the server. If you right-click the **SSISDB** database, in the **Integration Services** node in SQL Server Management Studio, and point to Reports, you can click **Integration Services Dashboard**. From here, you can get an overview of activity in the previous 24 hours and open reports for **All Executions**, **All Validations**, **All Operations**, and **Connections**. Each of these reports gives you the option to drill down further to find specific details for events and view performance data for the package executions.

For more information on Troubleshooting Reports for Package Execution, see *Troubleshooting Tools for Package Execution* in MSDN:



Troubleshooting Tools for Package Execution

<http://aka.ms/jk0r73>

Event Handlers

You can create an event handler for the **OnError** event. This can be used for many tasks, including sending an email message to an administrator, and logging system information to a file. Event handlers have control flow and data flows that are identical to those in a typical package, so you can perform almost unlimited tasks when an error has occurred.

Error Outputs

Most data flow components have outputs for success and failure. You can use the error output to direct rows containing errors to a different destination. Initially, the error output adds columns containing the error code. While you could manually look up the error code, you can add additional information, such as the error description, using the **Script** component.

Error outputs are for a specific component, whereas event handlers are for the whole package.

For more information on enhancing an error output with the Script component, see *Error Handling in Data* in MSDN:



Error Handling in Data

<http://aka.ms/vzcf97>

Logging

You can have logging and view package execution information in a SQL Server table or a file.

Log providers implement logging in packages, containers, and tasks. You can also specify which information you require in the log. To use logging in the package, open it in SQL Server Data Tools, and, on the SSIS menu, click **Logging**. You can now define both how the log information is stored, and which events and information are logged.

For more information on how to enable logging in a package, see *Enable Package Logging in SQL Server Data Tools* in MSDN:



Enable Package Logging in SQL Server Data Tools

<http://aka.ms/yewck9>

Debug Dump Files

If you execute a package with **dtexec**, you can specify that a debug dump file is created. This creates a .tmp file that contains items such as environment information and recent messages. The file is located in <drive>:\Program Files\Microsoft SQL Server\110\Shared>ErrorDumps. You can specify whether to create the debug dump files on an error, a warning, or on information.

For more information on the dtexec utility, see *dtexec Utility* in MSDN:



Dtexec Utility

<http://aka.ms/qeehmk>

Demonstration: Deploying an SSIS Project

In this demonstration, you will see how to:

- Configure the SSIS environment.
- Deploy an SSIS project.
- Create environments and variables.
- Run an SSIS package.
- View execution information.

Demonstration Steps

Configure the SSIS Environment

1. Ensure **20767C-MIA-DC** and **20767C-MIA-SQL** are started, and log onto **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Demofiles\Mod13** folder, run **Setup.cmd** as Administrator.
3. Start SQL Server Management Studio and connect to the **MIS-SQL** database engine using Windows authentication.
4. In Object Explorer, right-click **Integration Services Catalogs**, and then click **Create Catalog**.
5. In the **Create Catalog** dialog box, note that you must enable CLR integration when creating an SSIS catalog. You can also choose to enable automatic execution of the stored procedures used by the catalog when SQL Server starts. Then enter and confirm the password **Pa55w.rd**, and click **OK**.
6. In Object Explorer, expand **Integration Services Catalogs**, and then right-click the **SSISDB** node that has been created, and then click **Create Folder**.
7. In the **Create Folder** dialog box, in the **Folder name** box, type **Demo**, and then click **OK**.
8. Expand the **SSIDB** node to see the folder.
9. Minimize SQL Server Management Studio.

Deploy an SSIS Project

1. Start Visual Studio, and open the **DeploymentDemo.sln** solution in the **D:\Demofiles\Mod13** folder. This project contains the following two packages:
 - **Extract Login List.dtsx**—a package that uses a data flow to extract a list of logons from the master database and save them in a text file.
 - **Extract DB List.dtsx**—a package that extracts a list of databases from the **master** database and saves them in a text file.

Both packages use a project-level connection manager to connect to the **master** database, and a project-level parameter to determine the folder where the text files containing the extracted data should be saved.
2. On the **Build** menu, click **Build Solution**.
3. When the build has succeeded, on the **Project** menu, click **Deploy**.
4. In the **Introduction** page of the **Integration Services Deployment Wizard** dialog box, click **Next**.
5. In the **Select Destination** page, in the **Server name** box, type **MIA-SQL**, and then click **Connect**.

6. In the **Path** box, browse to the **SSISDB\Demo** folder you created earlier, click **OK**, and then click **Next**.
7. On the **Review** page, click **Deploy**.
8. When deployment has completed, click **Close**, and then close Visual Studio.

Create Environments and Variables

1. In SQL Server Management Studio, expand the **Demo** folder you created earlier, and then expand the **Projects** folder. Note that the **DeploymentDemo** project has been deployed.
2. Right-click the **Environments** folder, and then click **Create Environment**.
3. In the **Create Environment** dialog box, enter the environment name **Test**, and then click **OK**.
4. Repeat the previous two steps to create a second environment named **Production**.
5. Expand the **Environments** folder to see the environments you have created.
6. Right-click the **Production** environment, and then click **Properties**.
7. In the **Environment Properties** dialog box, on the **Variables** page, add a variable with the following settings:
 - o **Name:** DBServer
 - o **Type:** String
 - o **Description:** Server
 - o **Value:** MIA-SQL
 - o **Sensitive:** No (Unchecked)
8. Add a second variable with the following settings (making sure to include the trailing “\” in the value), and then click **OK**:
 - o **Name:** FolderPath
 - o **Type:** String
 - o **Description:** Folder
 - o **Value:** D:\Demofiles\Mod13\Production\
 - o **Sensitive:** No
9. Right-click the **Test** environment, and then click **Properties**.
10. In the **Environment Properties** dialog box, on the **Variables** page, add a variable with the following settings:
 - o **Name:** DBServer
 - o **Type:** String
 - o **Description:** Server
 - o **Value:** MIA_SQL
 - o **Sensitive:** No

11. Add a second variable with the following settings (making sure to include the trailing “\” in the value), and then click **OK**:
 - **Name:** FolderPath
 - **Type:** String
 - **Description:** Folder
 - **Value:** D:\Demofiles\Mod13\Test\
 - **Sensitive:** No
12. In the **Projects** folder, right-click **DeploymentDemo**, and then click **Configure**.
13. In the **Configure – DeploymentDemo** dialog box, on the **References** page, click **Add** and add the **Production** environment.
14. Click **Add** again and add the **Test** environment.
15. On the **Parameters** page, in the **Scope** list, select **DeploymentDemo**.
16. On the **Parameters** tab, click the ellipses (...) button for the **OutputFolder** parameter.
17. In the **Set Parameter Value** dialog box, select **Use environment variable**, click **FolderPath** in the list of variables, and then click **OK**.
18. In the **Configure – DeploymentDemo** dialog box, on the **Connection Managers** tab, click the ellipses button (...) for the **ServerName** property.
19. In the **Set Parameter Value** dialog box, select **Use environment variable**, click **DBServer** in the list of variables, and then click **OK**.
20. In the **Configure – DeploymentDemo** dialog box, click **OK**.

Run an SSIS Package

1. In Object Explorer, expand the **DeploymentDemo** package, expand the **Packages** folder.
2. Right-click **Extract DB List.dtsx**, and then click **Execute**.
3. In the **Execute Package** dialog box, select the **Environment** check box, and in the drop-down list, select **.\Test**.
4. View the **Parameters** and **Connection Managers** tabs and note that the **FolderPath** and **DBServer** environment variables are used for the **OutputFolder** parameter and **ServerName** property.
5. Click **OK** to run the package.
6. Click **No** when prompted to open the **Overview Report**.
7. In Object Explorer, right-click **Extract Login List.dtsx**, and then click **Execute**.
8. In the **Execute Package** dialog box, select the **Environment** check box and in the drop-down list, select **.\Production**.
9. Click **OK** to run the package.
10. Click **No** when prompted to open the **Overview Report**.
11. View the contents of the **D:\Demofiles\Mod13\Test** folder and note that it contains a file named **DBs.csv** that was produced by the **Extract DB List.dtsx** package when it was executed in the **Test** environment. You can examine this file by using Notepad; it should contain a list of the databases on the server.

12. View the contents of the **D:\Demofiles\Mod13\Production** folder and note that it contains a file named **Logins.csv** that was produced by the **Extract Login List.dtsx** package when it was executed in the **Production** environment.

View Execution Information

1. In SQL Server Management Studio, in Object Explorer, under **Integration Services Catalogs**, right-click **SSISDB**, point to **Reports**, point to **Standard Reports**, and then click **Integration Services Dashboard**.
2. In the **Packages Detailed Information (Past 24 Hours)** list, notice that the two most recent package executions succeeded, and then click the **Overview** link for the first package in the list.
3. In the **Overview** report, in the **Execution Information** table, note the environment that was used, and in the **Parameters Used** table, note the values you configured with environment variables.
4. At the top of the report, click the **Navigate Backward** button to return to the **Integration Services Dashboard** report.
5. In the **Packages Detailed Information (Past 24 Hours)** list, click the **Overview** link for the second package in the list.
6. In the **Overview** report, in the **Execution Information** table, note the environment that was used, and in the **Parameters Used** table, note the values you configured with environment variables.
7. Close SQL Server Management Studio.

Verify the correctness of the statement by placing a mark in the column to the right.

Statement	Answer
True or false? You can add an environment variable to an SSIS Project, and assign a value to it, which may be, for example, Test or Production .	

Lesson 3

Planning SSIS Package Execution

Before you perform an SSIS deployment, there are a number of factors you should consider. Along with the deployment model, security and variables, all mentioned in previous lessons, you should also consider factors such as scheduling, dependencies, notifications, and which actions take place in SSIS packages and which ones occur in SQL Server Agent jobs.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the methods for running SSIS packages.
- Describe considerations for package scheduling.
- Explain the security context options when running a package.
- Describe the choices when deciding where notifications are handled.
- Explain the options when handling SSIS logging.
- Describe the options for combining SSIS tasks with SQL Server Agent job tasks.
- Explain how to implement a SQL Server Agent job.

Options for Running SSIS Packages

After you deploy a package, you can execute it using one of the following methods:

SQL Server Management Studio

After you deploy a package to SQL Server, you run it from the context menu. You can supply values for the environment and any parameter, modify settings for the connection managers, and specify logging options.

- SQL Server Management Studio
 - Specify values for environment and any parameter
 - Modify settings for the connection managers
 - Specify logging options
- Dtexec and Dtexecui
 - Dtexec provides command line to run SSIS packages
 - Dtexecui is GUI to run Dtexec
- PowerShell
 - Powershell cmdlets make it easier for you to manage, monitor and execute SSIS packages
 - Can integrate SSIS with other Windows PowerShell tasks
- SQL Server Agent
 - Automates tasks in SQL Server
 - Particularly useful for SSIS packages

Dtexec and Dtexecui

Dtexec provides a command line interface to run SSIS packages. You can supply all necessary values for the environment, parameters, connection strings and logging. **Dtexecui** provides a graphical interface to run **dtexec** and has the same functionality.

For more information on the dtexec utility, see *dtexec Utility* in MSDN:



Dtexec Utility

<http://aka.ms/qeehmk>

Windows PowerShell®

SQL Server 2012 introduced Windows PowerShell cmdlets for you to manage, monitor and execute SSIS packages. It's a powerful environment for managing most aspects of your server and integrating SSIS with other Windows PowerShell tasks.

For more information on SSIS and PowerShell, see *SSIS and PowerShell in SQL Server 2012* in the SSIS Team Blog:

 **SSIS and PowerShell in SQL Server 2012**

<http://aka.ms/nvvvar>

SQL Server Agent

SQL Server Agent automates tasks in SQL Server. This is particularly useful for SSIS packages because they are often run on a repeating schedule to load data periodically. You can specify the schedule for the job, the environment, parameters, connection strings and logging.

For more information on SQL Server Agent jobs for packages, see *SQL Server Agent Jobs for Packages* in MSDN:

 **SQL Server Agent Jobs for Packages**

<http://aka.ms/ej7drz>

Scheduling the ETL Process

When scheduling package execution with SQL Server Agent or another automation tool, there are some considerations you should take into account in planning an ETL strategy.

Data Acquisition Time Windows

In a data warehouse scenario, data sources that are used periodically to extract data for your data warehouse are typically live transactional systems. These are often some of the most business-critical systems in your organization. For example, the orders system for an Internet-based bicycle shop holds the most important data for the data warehouse, but is also the most critical one to keep running at optimum performance levels. If the orders system is down, the business cannot operate.

It is crucial, therefore, to ascertain the best time to extract data from the live systems. For a weekly data load, it might be at the weekend, while for a daily load, it may be at night—it's important to check the systems for real-world statistics. Even truly international systems have quiet and busy periods.

If data acquisition is problematic at any time, it might be possible to run queries against mirrored systems that are used to provide high availability. Mirrored systems can provide access to live data at any time with no effect on the live transactional systems.

Package Execution Order

Ensure that packages execute in the correct order and that you complete one step before starting the next. This is typically the case when data is loaded to a staging database. The data load to staging must complete, followed by data transformation and cleansing, and then a data load to the data warehouse. You must either create a system where one step in the ETL process triggers the next one, or create schedules with enough time delay for the previous step to complete.

- Data acquisition time windows
 - Ascertain the best time to extract data from the live systems
- Package execution order
 - Ensure that packages execute in the correct order
 - One step is completed before the next is started
- Execution dependencies
 - Some packages require others to have already run due to data dependencies

A completely linear package execution system, with one package starting the next and so on, sounds straightforward but it might also be very slow, because no packages are run in parallel. Running packages in parallel requires careful scheduling, or more complex logic—to test if all prerequisite packages have run before the next stage starts.

Execution Dependencies

Some packages require others to have already run due to data dependencies. For example, a Salesperson table has a foreign key relationship with a Store table. Imagine you have recently opened a new store with new employees. If the employees are loaded before the store, they will break the foreign key constraint and the data load will fail. However, if the store is loaded before the employees, the data load will be successful.

Configuring Execution Context

Using SQL Server Agent, you can run a job with a different security context from your own. For example, a user has read-only access to the production databases, staging databases, and data warehouse, but occasionally needs to be able to load data from production to data warehouse via staging. In the SQL Server Agent job, this is performed by setting the **Run as** property on the job step to an account that has the required permissions on all three systems. The user does not require the special privileges on the system, only permission to run the SQL Server Agent job.

- Create a credential in the SQL Server instance
- Create a SQL agent proxy that maps to the credential
 - Activate it for the SQL Server Integration Services package job subsystem
- Configure package execution job steps to run as the proxy

There are three fixed database roles that apply to the SQL Server Agent—**SQLAgentUserRole**, **SQLAgentReaderRole**, and **SQLAgentOperatorRole**. **SQLAgentUserRole** and **SQLAgentReaderRole** members can execute only jobs that they have created, whereas **SQLAgentOperatorRole** members can execute, but not edit, a job created by another user. You can add users to these roles in the **MSDB** database.

To use the **Run as** property of a job step:

1. Create a credential in the **Security** node of SQL Server Management Studio that maps to an account with the permissions required for the job step.
2. Create a proxy in **SQL Server Agent** that uses the credential you have just created.
3. Activate the proxy you have just created for the SSIS Package job subsystem.
4. Configure the job step to use the proxy just created.
5. Put the user who needs to execute the job in the appropriate database role.

Where to Handle Notifications

You can send notifications from SSIS packages using the **SendMail** task and other notifications from SQL Server Agent jobs using operators. Be careful not to send too many notifications, as the operator can become overloaded and not notice an important email message amongst many duplicates.

You should consider how the SSIS package is used and which steps the SQL Server Agent job performs. If the SSIS package is executed manually as a SQL Server Agent job from SQL Server Management Studio—and from a Windows PowerShell script—you should consider putting notifications in the SSIS package to ensure that notifications are sent, regardless of the execution method. If the SSIS package is always executed as part of a SQL Server Agent job that has additional steps apart from the SSIS package, you should consider putting notifications in the SQL Server Agent job.

In many situations, the decision is not clear-cut and you should balance the advantages and disadvantages of putting notifications either in SQL Server Agent jobs or in SSIS packages.

- Notifications can be handled from:
 - SSIS package, using SendMail task
 - SQL Server Agent job, using operators

Where to Handle Logging

Logging is less intrusive than email notifications and, therefore, it is less of a problem to implement logging in both SQL Server Agent and SSIS packages. For SSIS, you should consider using the Integration Services Dashboard if that provides sufficient detail. The Integration Services Dashboard is automatic and displays an overview of packages that have executed, giving you the opportunity to drill down into particular items to display the detail.

In SQL Server Agent, logging is enabled by default and this will log any errors or warnings. You can specify the size and maximum age of the job history log in the properties of SQL Server Agent.

For more information on the SQL Server Agent Log, see *SQL Server Agent Error Log* in MSDN:

 **SQL Server Agent Error Log**

<http://aka.ms/t8rq5e>

- Logging can be handled by:
 - SSIS: consider using Integration Services Dashboard
 - SQL Server Agent job: logging enabled by default

In a typical scenario, you would examine the SQL Server Agent log and, if you found there was an error in an SSIS package execution, you would investigate the Integration Services Dashboard.

Combining SSIS Tasks with Other Tasks

SSIS is a powerful tool. You can use it to perform every task a job might require with a single SSIS package and schedule this as a single step in a SQL Server Agent job. Typically, however, this is not the best approach.

If you adopt the policy that each SSIS package performs one defined task, you can make modular packages that run multiple smaller ones. In this way, the work that you have done to create one, perhaps complex, SSIS package, can be easily reused.

Using this approach, you could continue to use SQL Server Agent to execute a single SSIS package, albeit one that runs several child packages. In this scenario, you might want to add another minor step to a package—perhaps executing a Transact-SQL query. To avoid changing any existing packages, you could create a new package to perform this task, and then create a new parent package to run the existing ones, along with your new package, containing the Transact-SQL task. In this scenario, it would be simpler to add a step to the SQL Server Agent job to run the Transact-SQL query.

One factor to consider if you have actions performed in both SQL Server Agent jobs and SSIS packages is documentation. You must ensure that everyone involved in developing the ETL solution is aware of where the actions are performed and, therefore, where they must be monitored and maintained.

- Create modular packages
- Augment SSIS packages with SQL Agent tasks
- Ensure all processes are documented

Implementing SSIS Agent Jobs and Schedules

To use SQL Server Agent to execute jobs, you should perform the following steps:

Select or create a Windows account that the SQL Server Agent will use and ensure that it has sufficient permissions to perform the assigned job steps. The Windows domain account you specify must have the following:

1. Permission to log on as a service (**SeServiceLogonRight**).
2. The following permissions to support SQL Server Agent proxies:
 - a. Permission to bypass traverse checking (**SeChangeNotifyPrivilege**).
 - b. Permission to replace a process-level token (**SeAssignPrimaryTokenPrivilege**).
 - c. Permission to adjust memory quotas for a process (**SeIncreaseQuotaPrivilege**).
 - d. Permission to log on using the batch logon type (**SeBatchLogonRight**).
3. If necessary, create Windows accounts for SQL Server Agent proxies, add credentials for them to the SQL Server instance, and create proxies for the appropriate job subsystems.

1. Select a Windows account that SQL Server Agent will use
2. Create job subsystem proxies if required
3. Set SQL Server Agent service to start automatically
4. Create the required jobs and schedules

4. Set the SQL Server Agent service to start automatically. By default, the SQL Server Agent is configured to start manually and, if there is a server restart, jobs will cease to run.
5. Create jobs for each package execution you want to automate, and specify a schedule on which the job should run.

Check Your Knowledge

Question	
Which of the following methods for running an SSIS package would be most appropriate if you want to schedule the package to run at regular intervals?	
Select the correct answer.	
<input type="checkbox"/>	SQL Server Management Studio
<input type="checkbox"/>	Dtexec
<input type="checkbox"/>	Dtexecui
<input type="checkbox"/>	Windows PowerShell
<input type="checkbox"/>	SQL Server Agent

Lab: Deploying and Configuring SSIS Packages

Scenario

You have completed the SSIS project containing the packages required for the data warehouse ETL process. Now you must deploy the project to an SSIS catalog, configure environments for dynamic configuration, and schedule automatic execution of packages.

Objectives

After completing this lab, you will be able to:

- Create an SSIS catalog.
- Deploy an SSIS project.
- Create environments.
- Run an SSIS package.
- Schedule SSIS package execution.

Estimated Time: 45 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password **Pa55w.rd**

Exercise 1: Creating an SSIS Catalog

Scenario

You will use Transact-SQL to enable the SQL CLR, and then create an SSIS Catalog and a folder to contain the Adventure Works ETL solution.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Create the SSIS Catalog and a Folder

► Task 1: Prepare the Lab Environment

1. Ensure that the **MT17B-WS2016-NAT**, **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab13\Starter** folder as Administrator.

► Task 2: Create the SSIS Catalog and a Folder

1. Use SQL Server Management Studio to create an **SSIS catalog** with a password of **Pa55w.rd** in the **MIA-SQL** instance of SQL Server.
2. In the catalog, create a folder named **DW ETL** with the description **Folder for the Adventure Works ETL SSIS Project**.

Results: A SQL Server instance with an SSIS Catalog.

Exercise 2: Deploying an SSIS Project

Scenario

You now want to deploy an existing SSIS project, created in Visual Studio, to the new folder from the previous exercise.

The main tasks for this exercise are as follows:

1. Deploy an SSIS Project

► Task 1: Deploy an SSIS Project

1. Start Visual Studio, and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab13\Starter** folder.
2. Build the solution.
3. Deploy the **AdventureWorksETL** project to the **DW ETL** folder in the **MIA-SQL** server.
4. Close Visual Studio.

Results: After this exercise, you should have deployed an SSIS project to a folder in your SSIS database.

Exercise 3: Creating Environments for an SSIS Solution

Scenario

You have parameters and connection managers in an SSIS package. You would like to create test and production environments with which to pass the correct values quickly to these parameters when switching from test to production.

The main tasks for this exercise are as follows:

1. Create Environments
2. Create Variables
3. Map Environment Variables

► Task 1: Create Environments

1. In SQL Server Management Studio, create an SSIS environment named **Test** in the DW ETL folder.
2. Create a second environment named **Production**.

► Task 2: Create Variables

1. Add the following variables to the **Production** environment:
 - A string variable named **StgServer** with the value **MIA-SQL**.
 - A string variable named **FolderPath** with the value **D:\Accounts**.
2. Add the following variables to the **Test** environment:
 - A string variable named **StgServer** with the value **MIA-SQL**.
 - A string variable named **FolderPath** with the value **D:\TestAccts**.

► Task 3: Map Environment Variables

1. Configure the **AdventureWorksETL** package in the SSIS Catalog to use the **Production** and **Test** environments.
2. Map the **AccountsFolder** project-level parameter to the **FolderPath** environment variable.
3. Map the **ServerName** property for the **localhost Staging ADO NET** and **localhost.Staging** connection managers to the **StgServer** environment variable.

Results: After this exercise, you should have configured your project to use environments to pass values to package parameters.

Exercise 4: Running an SSIS Package in SQL Server Management Studio

Scenario

Now you have deployed the SSIS project and defined execution environments, you can run the packages in SQL Server Management Studio.

The main tasks for this exercise are as follows:

1. Run a Package

► Task 1: Run a Package

1. Use SQL Server Management Studio to run the **Extract Payment Data.dtsx** package using the **Test** environment.
2. When the package has been executed, review the overview report and verify that it used the environment variable values you specified for the **Test** environment.

Results: After this exercise, you should have performed tests to verify that a value was passed from the environment to the package.

Exercise 5: Scheduling SSIS Packages with SQL Server Agent

Scenario

You've tested the packages by running them interactively. Now you must schedule them to execute automatically.

The main tasks for this exercise are as follows:

1. Create a SQL Server Agent Job
2. View the Integration Services Dashboard

► Task 1: Create a SQL Server Agent Job

1. In SQL Server Management Studio, create a new SQL Server Agent job named **Extract Reseller Data**.
2. Add a step named **Run Extract Reseller Data Package** to the job, and configure it to run the **Extract Reseller Data.dtsx** SQL Server Integration Services package in the **AdventureWorksETL** project you deployed earlier.
3. Configure the job step to use the **Production** environment.
4. Create a schedule for the job to run once, two minutes from the current time.

5. Wait for two minutes, then, in the **SQL Server Agent Job Activity Monitor**, view job activity and verify that the **Extract Reseller Data** job has completed successfully.

► **Task 2: View the Integration Services Dashboard**

1. View the Integration Services Dashboard and verify that the **Extract Reseller Data.dtsx** package was successfully executed.
2. View the overview report for the **Extract Reseller Data.dtsx** package execution and verify the parameter values used.
3. View the messages for the **Extract Reseller Data.dtsx** package execution, and note the diagnostic information that was logged.
4. When you have finished, close SQL Server Management Studio.

Results: After this exercise, you should have created a SQL Server Agent job that automatically runs your SSIS package.

Module Review and Takeaways

In this module, you have learned how to deploy and run SSIS packages.

Review Question(s)

Question: What are the advantages of the project deployment model?

Question: What are the advantages of environments?

Question: Where would you handle notifications?

Question: Describe the new feature available in SQL Server for Deploying SSIS Packages. What are the benefits?

Module 14

Consuming Data in a Data Warehouse

Contents:

Module Overview	14-1
Lesson 1: Introduction to Business Intelligence	14-2
Lesson 2: Introduction to Data Analysis	14-6
Lesson 3: Introduction to Reporting	14-9
Lesson 4: Analyzing Data with Azure SQL Data Warehouse	14-12
Lab: Using a Data Warehouse	14-15
Module Review and Takeaways	14-18

Module Overview

A data warehouse is the foundation of a business intelligence (BI) solution, with which business users, information workers, and data analysts can make faster, better decisions.

This module introduces BI, describing the components of Microsoft® SQL Server® that you can use to create a BI solution, and the client tools with which users can create reports and analyze data.

Objectives

After completing this module, you will be able to:

- Describe BI and common BI scenarios.
- Explain the key features of SQL Server Analysis Services.
- Explain the key features of SQL Server Reporting Services.
- Describe how to analyze data with Azure® SQL Data Warehouse.

Lesson 1

Introduction to Business Intelligence

With BI technologies, you can create decision support systems to help an individual in an organization work more effectively. However, a BI solution is only as good as the data that it processes. Therefore, using a data warehouse as the platform for your BI solution makes a lot of sense.

This lesson explains the benefits of using a data warehouse as a platform for BI, and introduces common reporting and data analysis scenarios.

Lesson Objectives

After completing this lesson, you will be able to:

- Explain the benefits of using a data warehouse as a platform for BI.
- Describe data analysis scenarios and technologies.
- Describe reporting scenarios and technologies.

The Data Warehouse As a Platform for Business Intelligence

The primary purpose of creating a data warehouse is to provide a platform for BI workers to access historical data so that they can perform reporting and data analysis. A well-designed data warehouse will include mechanisms that ensure the quality and accuracy of data, as described in previous modules—they will be optimized to provide the best performance for BI applications. Using a data warehouse as a data source for BI provides many benefits, including:

- Data quality and accuracy
- Data availability
- Complete and up-to-date data
- Query performance

Data quality and accuracy. Effective reporting and data analysis rely on the availability of data that is complete, accurate, consistent, and which does not contain duplicate data, or any other data that might compromise the accuracy of the reports and data analyses created. Using Master Data Services and Data Quality Services ensures that the data in the data warehouse meets the requirements of BI applications, enabling executives, managers, and other information workers to make better strategic and operational business decisions.

Data availability. The data used by information workers must also be easily available, so that individuals do not have to search for information. Centralizing data in a data warehouse makes it much easier for them to find the data they require.

Complete and up-to-date data. To support reporting and data analysis, the data in the data warehouse must be complete and up to date. SQL Server Integration Services provides comprehensive extract, transform, and load (ETL) capabilities to help you populate your data warehouse and maintain it with periodic or incremental updates.

Query performance. BI queries can be highly demanding, requiring a significant amount of processing power. A dedicated data warehouse that is optimized for BI queries—such as a Fast Track Data Warehouse system, Parallel Data Warehouse system, or Microsoft Analytics Platform system—can handle the processing load that BI queries require, even when many users are simultaneously accessing the database.

SQL Server provides a complete platform for data warehousing, and supports reporting tools such as SQL Server Reporting Services (SSRS), Report Builder, and Power View, in addition to data analysis tools like SQL Server Analysis Services (SSAS), Microsoft Excel®, and PerformancePoint Services in Microsoft SharePoint® Server.

Data Analysis

Data analysis involves exploring data to find interesting and useful information that can help companies to identify patterns, opportunities, and inefficiencies. Data in the data warehouse can be used to perform different types of analysis, such as identifying sales patterns over the previous quarter and predicting future performance based on past data.

Key data analysis scenarios include:

Exploring data to identify patterns. Data analysts use various tools to explore data to try and identify useful information, such as the best performing sales regions and correlations between sales figures and other factors. Business decision-makers can use this information to guide strategic planning decisions.

Self-service analysis and data mash-ups. Business users might perform self-service analysis when they want to combine data from the data warehouse and other sources without support from the IT department.

Data mining. Analysts use data mining to discover difficult-to-identify patterns and correlations between data. This information is useful because it can help to predict the success of future activities, such as targeted advertising campaigns.

Data analysis tools combine the data and metadata from data warehouses or other stores to create models that make this data available as useable information. For example, your data warehouse might include a table that contains millions of rows of sales data, a second one for product information, and other tables with details about dates and geographical location. You could create an analysis solution that loads data into a data model, such as a multidimensional cube or a tabular data model. This then calculates and stores aggregations of key data to improve query response times, making all this available to a range of client tools.

- Data analysis improves decision making by revealing patterns in data
- Data analysis scenarios:
 - Exploring data and identifying patterns
 - Self-service analysis
 - Data mining

Reporting

Reporting involves producing documents that summarize specific data. Reports typically include graphical elements, such as charts and tables, that help to make the information easy to understand. There are several common scenarios that report developers and administrators are likely to encounter, such as:

- **Scheduled delivery of standard reports.** BI developers create a set of standard reports that are issued to business users on a regular basis by one or more of the following methods:
 - Email message.
 - Delivery to a file share.
 - Delivery to a SharePoint site.
- **On-demand access to standard reports.** Business users consume reports on demand by browsing to an online report repository.
- **Embedded reports and dashboards.** Reports are integrated into an application or portal to provide contextualized business information at a glance.
- **Request to IT for custom reports.** Business users request specific reports from the IT department, and a BI developer creates these reports to order.
- **Self-service reporting.** Business users can use authoring tools to create their own reports from data models that have been published by the IT department.

- Reports summarize data:
 - Often include graphical elements such as charts and tables
- Reporting scenarios:
 - Scheduled delivery of standard reports
 - On-demand access to standard reports
 - Embedded reports and dashboards
 - Custom reports
 - Self-service reporting

SSRS is a component of SQL Server that provides a platform for building and managing reports in an enterprise. Reporting Services provides comprehensive functionality that supports all these scenarios.

Categorize Activity

Which of the following are possible benefits of a data warehouse? Indicate your answer by writing the category number to the right of each item.

Items	
1	Data quality and accuracy
2	High-speed processing of transactions
3	Data availability
4	User-friendly interface for inputting core data
5	Completeness of data
6	Up-to-date data
7	Query performance
8	Accessibility of data by suitably trained nontechnical users

Category 1	Category 2
Possible Benefits	Not Benefits

Lesson 2

Introduction to Data Analysis

SSAS is the SQL Server component that BI developers can use to create data models; users can then access data for analysis. Excel can help analysts and business users to create PivotTables and other data visualizations from enterprise data models. They can also create personal data models that combine data from the data warehouse with information from other sources.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the role of SSAS in BI.
- Explain the semantic data models that SSAS provides.
- Use Excel as a data analysis tool.

Analysis Services

SSAS is an online analytical processing (OLAP) platform that you can use to provide data for your reports and data analysis. You can use it to design, create, and manage data models, aggregating data from multiple disparate sources for your BI solutions.

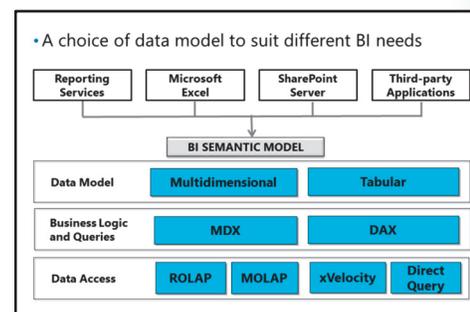
A typical SSAS project might define a data model for analysis, deploy the model as a database on an instance of Analysis Services, and then process the database to load the data from the original data sources. You can then use this new data source for client applications.

You can use SQL Server Data Tools (SSDT) to create the project and the model, utilizing the various wizards that the SSDT provides to quickly generate your new data source.

- A platform for enterprise-scale analytical data models
- Workflow:
 - Define the model
 - Deploy the model as a database
 - Process the database
- Use SQL Server Data Tools to create SSAS projects

Analysis Services Semantic Models

You can use SSAS to create a data model for users to analyze business data without having to understand the complexities of the underlying database, and to achieve optimal performance for BI queries. Data models expose data in a format that users can interact with more easily. With SSAS, you can create enterprise-scale data models to support BI applications including Reporting Services, Excel, PerformancePoint Services in SharePoint Server, and other third-party tools.



SSAS supports a concept known as the BI semantic model (BISM), which abstracts the specific implementation details of the data model so that client tools can use standard protocols and interfaces to access it, regardless of the specific model type. With Analysis Services, you can create two different kinds of data model, each with its own features and performance characteristics:

- **Multidimensional data models.** These have been used in all versions of SSAS, and should be familiar to experienced BI specialists. Multidimensional data models expose data through dimensions and cubes. They use the Multidimensional Expressions (MDX) language to implement business logic, and can provide access to data through relational online analytical processing (ROLAP) storage or multidimensional online analytical processing (MOLAP) storage.
- **Tabular.** Tabular data models were first introduced in SQL Server 2012 Analysis Services. They are constructed from tables and relationships, making them easier to work with for database professionals who have no multidimensional data modeling experience. Tabular data models use the Data Analysis Expressions (DAX) language to implement business logic, and can provide access to data by using the in-memory xVelocity engine. They also use DirectQuery mode to access data in the underlying data source, which is often a data warehouse.

Using Excel As a Data Analysis Tool

Excel is a well-established BI tool that data analysts and information workers use to create workbook-based applications to support business activities. You can use Excel to connect to Analysis Services data models, and explore the data that they contain, by creating PivotTable tables and PivotChart charts. With PivotTable tables, you can display data as a table, a matrix, or as a combination of the two. In a PivotTable table, you can expand and collapse data, color code values to highlight them, and use a range of formatting options to render the data as you require. You can use PivotChart charts to add graphical elements to a data analysis report, such as a column or pie chart.

- Connect to Analysis Services data models
- Create PivotTable tables and PivotChart charts
- Use add-ins for self-service data discovery, modeling, and visualization:
 - PowerPivot
 - Power View
 - Power Map

PowerPivot

The PowerPivot for Excel add-in extends the capabilities of Excel as a data analysis tool—you can create a tabular data model that stores data in the workbook. Using this approach, you can perform complex calculations offline, with excellent performance, even for very large data sets. PowerPivot for Excel adds the PowerPivot ribbon to the Excel interface, from which you can launch the PowerPivot window to create database connections, manage tables, and create DAX measures and calculated columns. In addition to creating PivotTable tables and PivotChart charts, you can add slicers enabling you to group and sort data with just a few clicks. The key features of PowerPivot for Excel include:

- Fast response times when performing complex analyses, even for very large datasets.
- Minimal memory footprint due to the VertiPaq storage engine.
- Advanced features such as hierarchies, perspectives, and relationship management.
- The ability to create custom measures and calculated columns by using DAX.
- The ability to scale up the workbook to a departmental solution by using it as a template for a new SQL Server Analysis Services tabular project.

SharePoint Server provides a central, managed portal through which analysts can securely share and collaborate on their PowerPivot for Excel workbooks. By sharing PowerPivot workbooks in SharePoint sites, users can locate the workbooks they require, rather than having to spend time creating duplicates because they did not know that a similar workbook already existed. PowerPivot for SharePoint helps administrators to establish governance over workbooks, which might be stored on the computers of many different users in an organization. It also makes it easier to provision appropriate resources to optimize application performance.

Power View

Power View is a data visualization application with which you can interactively explore data and create reports with a strong visual impact. You can create Power View reports by choosing from a range of charts and tables, which are called visualizations. These include bar charts, column charts, line charts, scatter charts, cards, tiles, tables, and matrices, which you can use to display data in many ways. You can include images and animations to enhance the appearance of a report and make it easier to understand, adding interactive elements with which you can update a report with a single click. For example, you can create tiled images that represent different categories of products. When you click on one of these images, the report updates to display the data for that category.

Power View reports are very easy to create, and the user interface (UI) is simple and intuitive. When you select a field or measure to add to a visualization, the application automatically uses it in the chart, in the most appropriate way. For example, when you select a field, Power View might add it to the X axis, use it as a measure value, or use it as a chart series. This behavior speeds up the time it takes to build meaningful reports, and makes it easier for inexperienced users to get started with data exploration. Modifying reports is also very easy because, by using just two or three clicks, you can drag visualizations to move and resize them, and replace fields and measures in existing visualizations. Power View will automatically display the recalculated values based on the modifications that you make.

Power Map

Power Map is a Microsoft Excel add-in with which users can render data values with a geographical dimension on a map, and use animation techniques to explore the data visually. For example, you could use Power Map to compare data aggregations such as sales volume by country, average property value by postal code, or population by city. Additionally, you can display the geographical data aggregations over time, to see how values changed or grew during a specific period. Users can view the Power Map visualizations directly in the Excel workbook, or they can be exported as a video.

Check Your Knowledge

Question	
Which of the following is NOT an Analysis Services data model?	
Select the correct answer.	
<input type="checkbox"/>	Multidimensional data model.
<input type="checkbox"/>	Tabular data model.
<input type="checkbox"/>	Entity data model.

Lesson 3

Introduction to Reporting

Many organizations use BI to inform business decision-making and communicate key information to executives, managers, information workers, customers, and other stakeholders. Often, the BI solution used to meet these requirements is managed by the IT department's technical specialists—they create and maintain the data models and reports necessary to support the organization's specific needs.

In this lesson, you will learn about SSRS, or Reporting Services, which information workers can use to create, publish, and subscribe to reports.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the key features of SSRS.
- Create Reporting Services reports.
- Use PowerView to create reports in a web browser.

Reporting Services

Reporting Services is a development environment for professional report developers and information workers who have to occasionally create reports.

You can use Reporting Services to create interactive, visually-interesting reports that can include charts, maps, data bars, and spark lines. Additionally, Reporting Services can support parameters, for users to change or set the context of a report.

SSRS is available in most editions of SQL Server, though not all features are supported in every edition. When you install SSRS, you must choose between two possible deployment modes:

- **SharePoint integrated mode.** In this mode, the report server is installed as a service application in a SharePoint Server farm; users manage and view reports in a SharePoint site.
- **Native mode.** In this mode, Reporting Services provides a management and report viewing UI called Report Manager, which is implemented as a stand-alone web application.

In addition to assisting users to create reports, Reporting Services provides a management environment for you to control the complete report life cycle. Users can publish and subscribe to reports, which can be delivered in a variety of ways. Administrators can define email alerts that inform users of changes to the reports they are interested in, and define permissions to ensure that reports can only be viewed and edited by selected individuals.

- Multiple report authoring environments
- Centralized report management and distribution
- Two installation modes:
 - SharePoint integrated
 - Native mode

Creating Reporting Services Reports

SSRS reports are XML-based documents that include markup defining both the report data and the layout. You can aggregate the report data from a range of data sources and display it in a variety of ways. You can also use Report Builder or Report Designer to create SSRS reports.

Report Designer

Report Designer is part of SSDT, which ships as part of SQL Server, hosted in the Visual Studio® development environment. With Report Designer, you can design reports by hand or by using the Report Service Project Wizard. You can then customize them in the Design pane, and check them in the Preview pane. You can use the query designers to retrieve data from a wide range of data sources and the expression editor to customize the content and appearance.

- Report Designer:
 - Part of SQL Server Data Tools
 - Design and Preview views
 - Query designers
 - Expression editors
 - Report Server Project Wizard
- Report Builder:
 - Microsoft Office-like environment
 - Built-in data providers
 - Query designers
 - Preview and export functionality

Report Builder

Report Builder provides a Microsoft Office-like environment for your power users to design their own reports. You can download the setup from your report server and install it locally on your computer. Report Builder provides built-in data providers for access to data from commonly-used sources, in addition to query designers and expression editors similar to those in Report Designer. It also provides preview functionality and you can export reports to other file formats for further analysis.

After you create a report using either of these tools, you can view the report locally or publish it to a report server or SharePoint site.

Power View

When SQL Server Reporting Services is installed in SharePoint integrated mode, you can enable the Power View SharePoint application so users can create dynamic data visualizations directly in the web browser. Power View reports are easy to create—the UI is simple and intuitive, making it easier for inexperienced users to get started with data exploration. Data sources for Power View in SharePoint Server include PowerPivot workbooks, tabular Analysis Services databases, and multidimensional Analysis Services cubes. With these tools, business users can create self-service data visualizations from their own mash-up data models and managed enterprise data models. Power View visualizations can be saved in SharePoint server as reports, for other users to view and extend, or as PowerPoint presentations, which can be used to bring live data visualizations to a meeting.

- In-browser data visualizations from:
 - PowerPivot workbooks
 - Analysis Services tabular models
 - Analysis Services multidimensional cubes
- Save reports in SharePoint Server or export to PowerPoint

Verify the correctness of the statement by placing a mark in the column to the right.

Statement	Answer
True or false? In general, only IT professionals create Reporting Services reports.	

Lesson 4

Analyzing Data with Azure SQL Data Warehouse

Organizations, individuals, services, and devices generate an ever-increasing volume of data at an ever-increasing rate. The growth of social media, digital devices for photography and video capture, in addition to the use of profile data to personalize user experiences and content, has led to a massive expansion of data processing requirements for business organizations and Internet services.

You can use Azure SQL Data Warehouse with other Azure services to store, process, and analyze these massive volumes of data in the cloud.

Lesson Objectives

After completing this lesson, you will be able to:

- Describe the functionality that Power BI™ for Microsoft Office 365® provides.
- Explain how the Azure Data Factory (ADF) automates and transforms data.
- Understand how you can use Azure Machine Learning to build predictive web services.
- Describe how you can use Azure Stream Analytics (ASA) to perform real-time analytics on large data streams.

Power BI

Power BI for Microsoft Office 365 builds on the self-service features provided by PowerPivot, Power Query, and Power View to create a collaborative environment in which business users can share data visualizations and queries. Users can also explore shared data using natural language query semantics.

Key elements of Power BI for Office 365 include:

- **Power BI Desktop.** A desktop application that you can use to build reports, refine visualizations, and share up-to-date information for each part of your business.
- **Power BI Apps.** Apps for Windows, iOS and Android for users to view reports on mobile devices.
- **Power BI Sites.** SharePoint Online sites where users can share workbooks and Power View visualizations.
- **Power BI Q&A.** A SharePoint site where users can create visualizations from PowerPivot workbooks, using natural language queries.
- **Data Management Gateway.** A service application that organizations can use to make on-premises data sources available to Power BI cloud services.
- **Shared Queries.** Power Query is used to define queries that individuals can then share with other Office 365 Power BI users in their organization.

- Cloud-based collaborative BI environment:
 - Power BI Desktop
 - Power BI Apps
 - Power BI Sites
 - Power BI Q&A
 - Data Management Gateway
 - Shared Queries
- Management Portals
 - Power BI Admin Center
 - Manage Data Portal

Power BI for Office 365 provides the following two web-based portals for managing the service:

- **The Power BI Admin Center.** Administrators can use this portal to define data sources, register data management gateways, manage user roles, and monitor service health statistics.
- **The Manage Data Portal.** Business users can use this portal to manage their own shared queries and data sources, and view usage analytics for the queries they have shared.

Azure Data Factory

ADF is a cloud-based data integration service that automates the movement and transformation of data. It coordinates existing services that collect raw data and transforms the data into ready-to-use information. ADF works across on-premises and cloud data sources to ingest, prepare, transform, analyze, and publish your data. With ADF, you can organize services into managed data flow pipelines to transform your data using services such as Azure HDInsight® (Hadoop) for your big data needs. You can then use Azure Machine Learning to operationalize your analytics solutions.

ADF uses four key components:

- **Activities.** ADF works with activities—the actions to perform on your data. You can use data movement activities to copy data between two locations and data transformation activities to transform and process data into predictions and insights.
- **Datasets.** ADF works with data structures, called datasets, which are pointers to the data you are using as input to or output from a pipeline. You can source this data from many types of data store, including SQL Server, Oracle DB, file shares, and documents.
- **Pipelines.** Pipelines are logical groupings of activities that work together to perform a task. For example, you might create a pipeline that performs three different transformations on a dataset, and then copies it to an alternative location. You can deploy and schedule a pipeline as one unit.
- **Linked services.** Linked services define the information needed to connect to two types of external resource: data stores that hold your data and compute resources that host the execution of an activity.

You can use the Azure Portal to view and monitor your data factories. You can drill down into each factory and view the status of the pipelines, activities, and datasets. Additionally, you can use Azure PowerShell™ cmdlets to pause and resume pipelines if any issues arise.

- Cloud-based data integration service
- Automates movement and transformation of data
- Key components:
 - Activities
 - Datasets
 - Pipelines
 - Linked services
 - Monitoring and management

Azure Machine Learning

Machine learning is a system that runs predictive models, which learn from existing data to forecast future trends. Predictions from machine learning can make apps and devices smarter. When you shop online, apps can use machine learning to recommend other products you might like, based on what you have purchased. When your credit card is swiped, machine learning compares the transaction to a database of transactions and helps the bank to detect fraud.

Azure Machine Learning is a powerful, cloud-based predictive analytics service. It provides tools to model predictive analytics in addition to a fully-managed service that you can use to deploy your predictive models as ready-to-consume web services. Additionally, it provides tools for creating complete predictive analytics solutions in the cloud—you can quickly create, test, operationalize, and manage predictive models.

Azure Machine Learning Studio is the key tool with which you can develop your solutions in the Azure cloud. It provides a development environment for creating, testing, and publishing your learning solutions. You can use the drag-and-drop designer to: create the experiments that generate your predictions; train and evaluate your model; publish the models as a web service for consumption. Alternatively, developers can use Microsoft Visual Studio to create custom ML client applications.

- Machine learning
 - Run predictive models
 - Learn from existing data
 - Forecast future trends
- Azure Machine Learning
 - Cloud-based service
- Azure Machine Learning Studio
 - GUI-based development environment
 - For creating, training, evaluating, and publishing models

Azure Stream Analytics

With Azure Stream Analytics (ASA), you can set up real-time analytic computations on data streams from devices, sensors, websites, social media, applications, and infrastructure systems. It is a fully managed, real-time processing engine that you can use to access insights from your data.

You author Stream Analytics jobs by defining the input source of streaming data, the output location for the results, and any required data transformations. ASA uses a dedicated query language, based on SQL, for the transformations, which nontechnical users can use to create transformations in a browser that supports IntelliSense®.

ASA can process up to 1 GB of data per second and integrates with Azure Event Hubs so that a solution can handle millions of events per second from multiple sources.

Question: How can you use Azure to manage and analyze data in your organization?

- Real-time analytic computations on data streams:

<ul style="list-style-type: none"> Devices Sensors Websites Social media 	<ul style="list-style-type: none"> Applications Infrastructure systems
--------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------
- Dedicated query language based on SQL
- Can process up to 1 GB of data per second

Lab: Using a Data Warehouse

Scenario

In this course, you have created a complete data warehousing solution for the Adventure Works Cycles company. Now you can explore the kinds of BI applications that information workers can use to gain business insights. You will explore an enterprise BI solution and use Excel to perform self-service BI analysis. You will also explore how to drill down into data with Power View.

Objectives

After completing this lab, you will be able to:

- Create and use data models.
- Use PowerPivot slicers to analyze data interactively.
- Use Power View to create charts and drill down into data.

Estimated Time: 45 minutes

Virtual machine: **20767C-MIA-SQL**

User name: **ADVENTUREWORKS\Student**

Password: **Pa55w.rd**

Exercise 1: Exploring a Reporting Services Report

Scenario

You have created a report based on the AdventureWorks sales data in the data warehouse, and now want to deploy it to the report server. After deploying the report, you will export the data to Excel, and filter the data to display values for predetermined criteria.

The main tasks for this exercise are as follows:

1. Prepare the Lab Environment
2. Publish a Reporting Services Report

► Task 1: Prepare the Lab Environment

1. Ensure that the **MT17B-WS2016-NAT**, **20767C-MIA-DC**, **20767C-MIA-CLI**, and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Run **Setup.cmd** in the **D:\Labfiles\Lab14\Starter** folder as Administrator.

► Task 2: Publish a Reporting Services Report

1. Switch to the **20767C-MIA-CLI** virtual machine, and log on as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**.
2. Start Report Builder and open **D:\Labfiles\Lab14\Starter\Internet Sales by State.rdl**.
3. Export the report to Excel as **D:\Labfiles\Lab14\Starter\Internet Sales by State.xlsx**, and then open the spreadsheet that you exported.
4. Add a filter to the data to select a **Product Subcategory** of **Road Bikes**.

Results: After this exercise, you should have deployed a Reporting Services report to the server, rendered the report, and exported the data to Excel.

Exercise 2: Exploring a PowerPivot Workbook

Scenario

You have seen how the data warehouse can be used in an enterprise BI solution, and you now want to explore how the data warehouse can support a self-service BI solution. You will create a visualization to analyze Internet sales by date.

The main tasks for this exercise are as follows:

1. View a PowerPivot Report

► Task 1: View a PowerPivot Report

1. In Excel, open **Analysis.xlsx** from the **D:\Labfiles\Lab14\Starter** folder.
2. Enable the following COM add-in: **Microsoft Office Power Pivot for Excel**.
3. Add a slicer for the **Ship Date** column of the **Reseller Sales** table.
4. Review the functionality of the slicer.
5. Save the workbook in its current location and then save it to the Document Library at **http://mia-sql/sites/adventureworks/**.
6. Close Excel.
7. Open Internet Explorer and navigate to **http://mia-sql/sites/adventureworks/**.
8. Review the report in the Document Library and then close Internet Explorer.

Results: After this exercise, you should have a customized PowerPivot report in an Excel workbook called **Analysis.xlsx**.

Exercise 3: Exploring a Power View Report

Scenario

Finally, you will create a chart for users to drill through to see data specific to their region.

The main tasks for this exercise are as follows:

1. Create a Power View Report

► Task 1: Create a Power View Report

1. In Excel, in the **Analysis.xlsx** workbook, enable **Power View for Excel** COM add-in.
2. Insert a Power View report with the title **Internet Sales**. Note that you might need to install Silverlight first.
3. In the top half of the Power View report, create a stacked column chart showing the sum of the **SalesAmount** field in the **Internet Sales** table by **ProductCategoryName** and **ProductSubcategoryName**.
4. In the bottom half of the Power View report, create a table showing the sum of the **SalesAmount** field in the **Internet Sales** table and the **CountryRegionName** in the **Geography** table.
5. Click bars in the bar chart and note that the relationship between the tables causes the table to be filtered, based on the product that you clicked.
6. Save the workbook and close Excel.

Results: After this exercise, you should have created a Power View report.

Module Review and Takeaways

In this module, you have learned how enterprise and self-service BI solutions can use a data warehouse.

Review Question(s)

Question: What are the issues you should consider when designing a data warehouse that must support self-service BI?

Course Evaluation

Course Evaluation

- Your evaluation of this course will help Microsoft understand the quality of your learning experience.
- Please work with your training provider to access the course evaluation form.
- Microsoft will keep your answers to this survey private and confidential and will use your responses to improve your future learning experience. Your open and honest feedback is valuable and appreciated.

Your evaluation of this course will help Microsoft understand the quality of your learning experience.

Please work with your training provider to access the course evaluation form.

Microsoft will keep your answers to this survey private and confidential and will use your responses to improve your future learning experience. Your open and honest feedback is valuable and appreciated.

Module 1: Introduction to Data Warehousing

Lab: Exploring a Data Warehousing Solution

Exercise 1: Exploring Data Sources

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab01\Starter** folder, right-click **Setup.cmd**, and then click **Run as administrator**.
3. Click **Yes** when prompted to confirm that you want to run the command file, wait for the script to finish, and then press any key to continue. This script can take several minutes to complete.

► Task 2: View the Solution Architecture

1. In the **D:\Labfiles\Lab01\Starter** folder, double-click **Adventure Works DW Solution.png** to open it.
2. Examine the diagram, and note that it shows several data sources on the left, that provide the source data for the data warehouse. You will examine these data sources in the remainder of this exercise.



Note: In addition to the data sources that you will examine in this lab, the diagram includes a Microsoft SQL Server Master Data Services model for product data and a SQL Server Data Quality Services task to cleanse data as it is staged. These elements form part of the complete solution for the lab scenario in this course, but they are not present in this lab.

3. Minimize the image. You will return to it in the next exercise.

► Task 3: View the Internet Sales Data Source

1. Start Microsoft SQL Server Management Studio (SSMS), and when prompted, connect to the **MIA-SQL** database engine instance using Windows® authentication.
2. On the **File** menu, point to **Open** and then click **File**.
3. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab01\Starter** folder, click **View Internet Sales.sql**, and then click **Open**.
4. Click **Execute** to run the query. When the query completes, review the results and note that this data source contains data about customers and the orders they have placed through the Adventure Works e-commerce website.
5. Keep SQL Server Management Studio open for the next task.

► Task 4: View the Reseller Sales Data Source

1. On the **File** menu, point to **Open** and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab01\Starter** folder, click **View Reseller Sales.sql**, and then click **Open**.
3. Click **Execute** to run the query. When the query completes, review the results and note that this data source contains data about resellers and the orders they have placed through Adventure Works reseller account managers.

4. Keep SQL Server Management Studio open for the next task.

► **Task 5: View the Products Data Source**

1. On the **File** menu, point to **Open** and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab01\Starter** folder, click **View Products.sql**, and then click **Open**.
3. Click **Execute** to run the query. When the query completes, review the results and note that this source contains data about products that Adventure Works sells, organized into categories and subcategories.
4. Keep SSMS open for the next task.

► **Task 6: View the Human Resources Data Source**

1. On the **File** menu, point to **Open**, and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab01\Starter** folder, click **View Employees.sql**, and then click **Open**.
3. Click **Execute** to run the query. When the query completes, review the results and note that this source contains data about employees, including the sales representatives associated with reseller sales.
4. Minimize SSMS. You will return to it later in this exercise.

► **Task 7: View the Accounts Data Source**

1. View the contents of the **D:\Accounts** folder that contains several comma-delimited text files.
2. Double-click **Payments – AU.csv** and open the file using Notepad.
3. Review the file contents and note that it contains data about reseller payments processed by the Adventure Works accounting system. Each file in the Accounts folder relates to payments made by resellers in a specific country.
4. Close Notepad without saving any changes.

► **Task 8: View the Staging Database**

1. Return to SQL Server Management Studio.
2. In Object Explorer, expand **Databases**, expand **Staging**, and then expand **Tables**.
3. Right-click the **dbo.Customers** table, and then click **Select Top 1000 Rows**. When the query completes, note that this table is empty.
4. Repeat the previous step to verify that the **dbo.EmployeeDeletes**, **dbo.EmployeeInserts**, **dbo.EmployeeUpdates**, **dbo.InternetSales**, **dbo.Payments**, **dbo.Resellers**, and **dbo.ResellerSales** tables are also empty.

Note: The **dbo.ExtractLog** table contains data that is used to track data extractions from the Internet Sales and Reseller Sales data sources.

5. Minimize SSMS. You will return to it in the next exercise.

Results: After this exercise, you should have viewed data in the InternetSales, ResellerSales, Human Resources, and Products SQL Server databases; viewed payments data in comma-delimited files; and viewed an empty staging database.

Exercise 2: Exploring an ETL Process

► Task 1: View the Solution Architecture

1. Maximize Paint and view the solution architecture diagram.
2. Note that the ETL solution consists of two main phases: a process to **extract** data from the various data sources and load it into a **Staging** database, and another to **load** the data in the **Staging** database into the data warehouse. In this exercise, you will observe these ETL processes as they run.
3. Minimize Paint. You will return to it in the next exercise.

► Task 2: Run the ETL Staging Process

1. Start Visual Studio® from the task bar.
2. On the **File** menu, point to **Open** then click **Project/Solution**.
3. In the **Open Project** dialog box, navigate to the **D:\Labfiles\Lab01\Starter** folder, click **AdventureWorksETL.sln**, and then click **Open**.
4. If the Solution Explorer pane is not visible, on the **View** menu, click **Solution Explorer**. If necessary, click the pin icon to freeze it in position.
5. In Solution Explorer, in the **SSIS Packages** folder, double-click **Stage Data.dtsx** to open it. Note that the staging process consists of five tasks:
 - Stage Internet Sales
 - Stage Reseller Sales
 - Stage Payments
 - Stage Employees
 - Notify Completion
6. On the **Debug** menu, click **Start Debugging**, and then observe that the staging process runs a SQL Server Integration Services package for each task.
7. When the message **Staging process complete** is displayed, click **OK**, and then, on the **Debug** menu, click **Stop Debugging**. Note that the message box might be hidden by the Visual Studio window. Look for a new icon on the taskbar, and then click it to bring the message box to the front.
8. Minimize Visual Studio. You will return to it later in this exercise.

► Task 3: View the Staged Data

1. Return to SQL Server Management Studio.
2. If necessary, in Object Explorer, expand **Databases**, expand **Staging**, and then expand **Tables**.
3. Right-click the **dbo.Customers** table, and then click **Select Top 1000 Rows**. When the query completes, note that this table now contains data that the ETL process has extracted from the data source.
4. Repeat the previous step to verify that the **dbo.EmployeeInserts**, **dbo.InternetSales**, **dbo.Payments**, **dbo.Resellers**, and **dbo.ResellerSales** tables also contain staged data.
5. Minimize SQL Server Management Studio. You will return to it later in this exercise.

► Task 4: Run the ETL Data Warehouse Load Process

1. Return to Visual Studio.
2. In Solution Explorer, in the SSIS Packages folder, double-click **Load DW.dtsx** to open it. Note that the data warehouse loading process consists of a sequence of tasks to load various dimensions and facts. This is followed by a task to determine the number of records loaded from each staging table before truncating the staging tables, and a task to log the row counts.
3. On the **Debug** menu, click **Start Debugging**, and then observe that the data warehouse loading process runs an SSIS package for the dimension or fact table to be loaded. The process might take several minutes to complete.
4. When the message **Data warehouse load complete** is displayed, click **OK**, and then, on the **Debug** menu, click **Stop Debugging**. Note that the message box might be hidden by the Visual Studio window. Look for a new icon on the taskbar, and then click it to bring the message box to the front.
5. Close Visual Studio.

Results: After this exercise, you should have viewed and run the SQL Server Integration Services packages that perform the ETL process for the Adventure Works data warehousing solution.

Exercise 3: Exploring a Data Warehouse

► Task 1: View the Solution Architecture

1. Maximize Paint and view the solution architecture diagram.
2. Note that the data warehouse provides a central data source for business reporting and analysis.
3. Close Paint without saving any changes.

► Task 2: Query the Data Warehouse

1. Return to SQL Server Management Studio, and open the **Query DW.sql** query file in the **D:\Labfiles\Lab01\Starter** folder.
2. Click **Execute** to run the query. When the query completes, review the results and note that the query uses the data warehouse to retrieve:
 - Total sales for each country by fiscal year.
 - Total units sold for each product category by calendar year.
3. Close SQL Server Management Studio without saving any changes.

Results: After this exercise, you should have successfully retrieved business information from the data warehouse.

Module 2: Planning Data Warehouse Infrastructure

Lab: Planning Data Warehouse Infrastructure

Exercise 1: Planning Data Warehouse Hardware

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC**, **20767C-MIA-SQL**, and **20767C-MIA-CLI** virtual machines are both running, and then log on to the **20767C-MIA-SQL** virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab02\Starter** folder, right-click **Setup.cmd**, and then click **Run as administrator**.
3. Click **Yes** when prompted to confirm that you want to run the command file. The script will execute.
4. Wait for the script to complete. You might have to press a key to close the command window.

► Task 2: Measure Maximum Consumption Rate

1. Start SQL Server Management Studio, and then connect to the **MIA-SQL** database engine by using Windows authentication.
2. On the **File** menu, point to **Open**, and then click **File**. Browse to the **D:\Labfiles\Lab02\Starter** folder, click the **Create BenchmarkDB.sql** query file, and then click **Open**.
3. Click **Execute**, and then wait for query execution to complete.
4. On the **File** menu, point to **Open**, and then click **File**. Browse to the **D:\Labfiles\Lab02\Starter** folder, click the **Measure MCR.sql** query file, and then click **Open**.
5. Click **Execute**, and then wait for query execution to complete.
6. In the results pane, click the **Messages** tab.
7. Add the **Logical Reads** values for the two queries together, and then divide the result by two to find the average.
8. Add the **CPU Time** values for the two queries together, and then divide the result by two to find the average. Divide the result by 1000 to convert it to seconds.
9. Calculate MCR by using the following formula:

$$(\text{average logical reads} / \text{average CPU time}) * 8 / 1024$$

10. Calculate the number of cores that are required to support a workload that has an average query size of 500 MB, 10 concurrent users, and a target response time of 20 seconds:

$$((500 / \text{MCR}) * 10) / 20$$

11. Close SQL Server Management Studio without saving any files.

► Task 3: Estimate Server Hardware Requirements

1. Log on to the **20767C-MIA-CLI** virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In Windows Explorer, move to the **D:\Labfiles\Lab02\Starter\Setupfiles** folder, and then double-click **DWHardwareSpec.xlsx** to open it in Microsoft Excel.

- Cells **A3**, **B3**, **C3**, **B6**, **B7**, and **B8** are prepopulated.
- In Cell **B9**, type the following formula to calculate the number of **Cores Required** for the given workload figures:

```
=((B6/C3)*B7)/B8
```

- In Cell **B11**, using the results of the preceding formula, estimate the **Processors Required** value. Use the following estimation:

```
=ROUND(B9/4,0)
```

- In Cell **B14**, enter a value for **Total System RAM** in GB. Note that this figure would depend on your actual deployment:

```
128
```

- Cells **B17**, **B18**, and **B19** are pre-populated. Note that **B19, Fact Data Volume** is derived as follows:

```
=(B17*B18)/100000
```

- In **Cell B20**, enter a formula to calculate the value for **Indexes and Dimensions** storage:

```
=B19*0.1
```

- In Cell **B21**, enter a formula to calculate the **Compressed Data (GB)** value using a typical compression ratio of **3:1**:

```
=SUM(B19:B20)/3
```

- In Cells **B23**, **B24**, and **B25**, enter 50 GB storage each for **Log Space**, **tempdb**, and **Staging**.

- In **Cell B26**, enter a formula to calculate a value for **Total Data Volume (GB)**:

```
=SUM(B21:B25)
```

- In Cell **B28**, enter an annual data growth rate of **150 GB**.

- In Cell **B29**, enter a formula to calculate how much **Data Storage in 3** years may be required:

```
=B26+(B28*3)
```

- Enter a formula in Cell **B30** to calculate a value for **Storage Capacity Required**:

```
=B29*2
```

- In the **D:\Labfiles\Lab02\Starter** folder, double-click **Storage Options.docx**, and then review the available options for storage hardware. Then, based on the storage requirements that you have calculated, select a suitable option for the data warehouse. Switch to Excel, and enter your option choice in Cell **B33**.

- Assess the option you have made and, in Cell **B34**, note the reasons for your choice in Cell **B33**. Discuss SAN, multiple spindles, disk configuration, and any other issues you think are important. For example:

```
* SAN allows for growth and offloads RAID and I/O processing to array controller  
* Multiple spindles helps balance I/O load  
* Disk configuration provides adequate storage capacity at reasonable cost
```

-
17. Save your changes to **DWHardwareSpec.xlsx**, and then close Excel and Word.

Results: After this exercise, you should have a completed worksheet that specifies the required hardware for your data warehouse server.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 3: Designing and Implementing a Data Warehouse

Lab: Implementing a Data Warehouse

Exercise 1: Implementing a Star Schema

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab03\Starter** folder, right-click **Setup.cmd** and then click **Run as administrator**.
3. In the **User Account Control** dialog box, click **Yes**, and then wait for the script to finish.
4. In the Command Prompt window, when prompted, press **y**, and then press Enter.

► Task 2: View a Data Warehouse Schema

1. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the SQL Server database engine using Windows authentication.
2. In Object Explorer, expand **Databases**, expand **AWDataWarehouse**, and then expand **Tables**. Note that the database contains four user tables.
3. Right-click **Database Diagrams**, and then click **New Database Diagram**. If you are prompted to create the required support objects, click **Yes**.
4. In the **Add Table** dialog box, click each table while holding down the CTRL key to select them all, click **Add**, and then click **Close**.
5. In the database diagram, click each table while holding down the CTRL key to select them all.
6. On the **Table Designer** menu, point to **Table View**, and then click **Standard**.
7. Arrange the tables and adjust the zoom level so you can see the entire database schema, and then examine the tables, noting the columns that they contain.

Note that the **FactResellerSales** table contains foreign key columns that relate it to the **DimReseller**, **DimEmployee**, and **DimProduct** tables. It also contains some numerical measures that can be aggregated to provide useful business information, such as the total sales amount (**SalesAmount**) and the total quantity of units ordered (**OrderQuantity**).

8. On the **File** menu, click **Save Diagram_0**.
9. In the **Choose Name** dialog box, in the **Enter a name for the diagram** box, type **AWDataWarehouse Schema**, and then click **OK**.

► Task 3: Create a Dimension Table

1. On the **File** menu, point to **Open**, and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab03\Starter** folder, click **DimCustomer.sql**, and then click **Open**.
3. Review the Transact-SQL code, noting that it creates a table named **DimCustomer** in the **AWDataWarehouse** database.
4. Click **Execute** to create the table, and then close the query window.

► Task 4: Create a Fact Table

1. On the **File** menu, point to **Open**, and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab03\Starter** folder, click **FactInternetSales.sql**, and then click **Open**.
3. Review the Transact-SQL code, noting that it creates a table named **FactInternetSales** in the **AWDataWarehouse** database, and that this table is related to the **DimCustomer** and **DimProduct** tables.
4. Click **Execute** to create the table, and then close the query window.

► Task 5: View the Revised Data Warehouse Schema

1. In the **AWDataWarehouse Schema** window, on the **Database Diagram** menu, click **Add Table**.
2. In the **Add Table** dialog box, click **Refresh**, select the **DimCustomer** and **FactInternetSales** tables, click **Add**, and then click **Close**.
3. Select each of the new tables, and then on the **Table Designer** menu, point to **Table View**, then click **Standard**.
4. Arrange the tables and adjust the zoom level so that you can see the entire database schema.
5. On the **File** menu, click **Save AWDataWarehouse Schema**.
6. Keep SQL Server Management Studio open. You will return to it in the next exercise.

Results: After this exercise, you should have a database diagram in the **AWDataWarehouse** database that shows a star schema consisting of two fact tables (**FactResellerSales** and **FactInternetSales**) and four dimension tables (**DimReseller**, **DimEmployee**, **DimProduct**, and **DimCustomer**).

Exercise 2: Implementing a Snowflake Schema

► Task 1: Create Dimension Tables That Form a Hierarchy

1. On the **File** menu, point to **Open**, and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab03\Starter** folder, click **DimProductCategory.sql**, and then click **Open**.
3. Review the Transact-SQL code, noting that it:
 - Creates a table named **DimProductCategory**.
 - Creates a table named **DimProductSubcategory** that has a foreign key relationship to the **DimProductCategory** table.
 - Drops the **ProductSubcategoryName** and **ProductCategoryName** columns from the **DimProduct** table.
 - Adds a **ProductSubcategoryKey** column to the **DimProduct** table that has a foreign key relationship to the **DimProductSubcategory** table.
4. Click **Execute**, wait for the query to finish, and then close the query window.

► Task 2: Create a Shared Dimension Table

1. On the **File** menu, point to **Open**, and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab03\Starter** folder, click **DimGeography.sql**, and then click **Open**.
3. Review the Transact-SQL code, noting that it:
 - Creates a table named **DimGeography**.
 - Drops the **City**, **StateProvinceName**, **CountryRegionCode**, **CountryRegionName**, and **PostalCode** columns from the **DimReseller** table.
 - Adds a **GeographyKey** column to the **DimReseller** table that has a foreign key relationship to the **DimGeography** table.
 - Drops the **City**, **StateProvinceName**, **CountryRegionCode**, **CountryRegionName**, and **PostalCode** columns from the **DimCustomer** table.
 - Adds a **GeographyKey** column to the **DimCustomer** table that has a foreign key relationship to the **DimGeography** table.
4. Click **Execute**, wait for the query to finish, and then close the query window.

► Task 3: View the Data Warehouse Schema

1. In the Database Diagram window, select the **DimProduct**, **DimReseller**, and **DimCustomer** tables, and then press DELETE to remove them from the diagram (this does not drop the tables from the database).
2. On the **Database Diagram** menu, click **Add Table**.
3. In the **Add Table** dialog box, click **Refresh**, select the **DimCustomer**, **DimGeography**, **DimProduct**, **DimProductCategory**, **DimProductSubcategory**, and **DimReseller** tables, click **Add**, and then click **Close**. If prompted to add a new relationship, click **OK**.
4. Select each of the new tables, and then on the **Table Designer** menu, point to **Table View**, then click **Standard**.
5. Arrange the tables and adjust the zoom level so that you can see the entire database schema.
6. On the **File** menu, click **Save AWDataWarehouse Schema**.
7. Keep SQL Server Management Studio open. You will return to it in the next exercise.

Results: After this exercise, you should have a database diagram in the **AWDataWarehouse** database showing a snowflake schema that contains a dimension consisting of a **DimProduct**, **DimProductSubcategory**, and **DimProductCategory** hierarchy of tables, in addition to a **DimGeography** dimension table that is referenced by the **DimCustomer** and **DimReseller** dimension tables.

Exercise 3: Implementing a Time Dimension Table

► Task 1: Create a Time Dimension Table

1. In SQL Server Management Studio, on the **File** menu, point to **Open**, and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab03\Starter** folder, click **DimDate.sql**, and then click **Open**.
3. Review the Transact-SQL code, noting that it:
 - Creates a table named **DimDate**.
 - Adds **OrderDateKey** and **ShipDateKey** columns to the **FactInternetSales** and **FactResellerSales** tables that have foreign key relationships to the **DimDate** table.
 - Creates indexes on the **OrderDateKey** and **ShipDateKey** foreign key fields in the **FactInternetSales** and **FactResellerSales** tables.
4. Click **Execute**, wait for the query to finish, and then close the query window.

► Task 2: View the Database Schema

1. In the **AWDataWarehouse** Schema window, select the **FactResellerSales** and **FactInternetSales** tables, and then press **DELETE** to remove them from the diagram.
2. On the **Database Diagram** menu, click **Add Table**.
3. In the **Add Table** dialog box, click **Refresh**, select the **DimDate**, **FactInternetSales**, and **FactResellerSales** tables, click **Add**, and then click **Close**.
4. Select each of the new tables, and then on the **Table Designer** menu, point to **Table View**, and then click **Standard**.
5. Arrange the tables and adjust the zoom level so you can see the entire database schema.
6. In Object Explorer, in the **AWDataWarehouse** database, expand **Tables**, right-click **Tables**, and then click **Refresh**.
7. Expand the **dbo.FactInternetSales** table, and then expand **Columns**. Note that each of the columns, **OrderDateKey** and **ShipDateKey**, has a silver key symbol against it, denoting that each has a foreign key relationship. Note that, whereas the primary keys are also indicated by a key symbol on the table (pointing in the opposite direction) in **Object Explorer** and in the **Database Diagram**, the foreign key relationships are *not* indicated by a symbol on the database diagram. However, a line between the **FactResellerSales** table and the **DimDate** table indicates the presence of a foreign key relationship between the two tables, but does not show which columns are the foreign keys, and does not show that there are two foreign key relationships: **OrderDateKey** and **ShipDateKey** in the **FactResellerSales** table linked to **DateKey** on the **DimDate** table.
8. On the **File** menu, click **Save AWDataWarehouse Schema**.

► Task 3: Populate the Time Dimension Table

1. On the **File** menu, point to **Open**, and then click **File**.
2. In the **Open File** dialog box, navigate to the **D:\Labfiles\Lab03\Starter** folder, click **GenerateDates.sql**, and then click **Open**.
3. Review the Transact-SQL code, noting that it:
 - Declares a variable named **@StartDate** with the value 1/1/2000, and a variable named **@EndDate** with the value of the current date.

- Performs a loop to insert appropriate values for each date between **@StartDate** and **@EndDate** into the **DimDate** table.
4. Click **Execute** to populate the table, and then close the query window.
5. In Object Explorer, in the **AWDataWarehouse** database, right-click **Tables**, and then click **Refresh**.
6. Right-click the **dbo.DimDate** table, point to **Script Table as**, point to **SELECT To**, and then click **New Query Editor Window**.
7. Click **Execute** to run the query.
8. View the data in the table, noting that it contains appropriate values for each date.
9. Close SQL Server Management Studio, saving changes if you are prompted.

Results: After this exercise, you should have a database that contains a **DimDate** dimension table that is populated with date values from January 1, 2000, to the current date.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 4: Columnstore Indexes

Lab: Using Columnstore Indexes

Exercise 1: Create a Columnstore Index on the FactProductInventory Table

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab04\Starter** folder, right-click **Setup.cmd**, and then click **Run as administrator**.
3. When you are prompted, click **Yes** to confirm that you want to run the command file, and then wait for the script to finish.

► Task 2: Examine the Existing Size of the FactProductInventory Table and Query Performance

1. On the taskbar, click **SQL Server Management Studio**.
2. In the **Connect to Server** window, ensure the **Server name** box contains **MIA-SQL** and the **Authentication** box contains **Windows Authentication**. Click **Connect**.
3. On the **File** menu, point to **Open**, then click **File**.
4. Navigate to the **D:\Labfiles\Lab04\Starter** folder, click **Query FactProductInventory.sql**, and then click **Open**.
5. On the **Query** menu, click **Include Actual Execution Plan**.
6. Highlight the Transact-SQL between the Step 1 and Step 2 comments. On the toolbar, click **Execute**.
7. Make a note of the data space used.
8. Highlight the code after the Step 2 comment.
9. On the toolbar, click **Execute** and then view the query results.
10. On the **Messages** tab, note the CPU and elapsed times.
11. On the **Execution plan** tab, view the execution plan for the query. Hover the mouse pointer over the icons from right to left to examine their details and note the indexes that were used. Also, note that the query processor has identified that adding a missing index could improve performance.

► Task 3: Create a Columnstore Index on the FactProductInventory Table

1. On the toolbar, click **New Query**, and then in the SQLQuery1.sql pane, type the following Transact-SQL code to create a columnstore index on the **FactProductInventory** table:

```
CREATE NONCLUSTERED COLUMNSTORE INDEX NCI_FactProductInventory_UnitCost_UnitsOut ON
FactProductInventory
(
    ProductKey,
    DateKey,
    UnitCost,
    UnitsOut
)
GO
```

2. On the toolbar, click **Execute** to create the index.

3. Switch back to the **Query FactProductInventory.sql** tab, and then click **Execute** to rerun the query.
4. On the **Execution plan** tab, view the execution plan that is used for the query. Examine the icons from right to left, noting the indexes that were used. Note that the new columnstore index is used.
5. Also note the **Actual Execution Mode** is **Batch** instead of **Row** and that this mode is applied to most of the steps. This is the main reason for the query performance improvement.
6. Close both query windows without saving changes.

Results: After completing this exercise, you will have created a columnstore index and improved the performance of an analytical query. This will have been done in real time without impacting transactional processing.

Exercise 2: Create a Columnstore Index on the FactInternetSales Table

► Task 1: Examine the Existing Size of the FactInternetSales Table and Query Performance

1. On the **File** menu, point to **Open**, then click **File**.
2. Navigate to the **D:\Labfiles\Lab04\Starter** folder, click **Query FactInternetSales.sql**, and then click **Open**.
3. On the **Query** menu, click **Include Actual Execution Plan**.
4. Highlight the Transact-SQL between the Step 1 and Step 2 comments. On the toolbar, click **Execute**.
5. Make a note of the data space used.
6. Highlight the code after the Step 2 comment.
7. On the toolbar, click **Execute** and then view the query results.
8. On the **Messages** tab, note the CPU and elapsed times.
9. On the **Execution plan** tab, view the execution plan for the query. Hold the mouse pointer over the icons from right to left to examine their details and note the indexes that were used. Also note that the query processor has identified that adding a missing index could improve performance.

► Task 2: Create a Columnstore Index on the FactInternetSales Table

1. On the **File** menu, point to **Open**, then click **File**.
2. Navigate to the **D:\Labfiles\Lab04\Starter** folder, click **Drop Indexes on FactInternetSales.sql**, and then click **Open**.
3. Click **Execute** to drop the existing indexes.
4. On the **File** menu, point to **Open**, then click **File**.
5. Click **Create Columnstore Index on FactInternetSales.sql**, and then click **Open**.
6. Click **Execute** to create the index.
7. Switch back to the **Query FactInternetSales.sql** tab, and then click **Execute** to rerun the query.
8. On the **Execution plan** tab, scroll to the end to view the execution plan that is used for the query. Examine the icons from right to left, noting the indexes that were used. Note that the new columnstore index is used.

9. Close the three query windows without saving changes.

Results: After completing this exercise, you will have greatly reduced the disk space taken up by the **FactInternetSales** table, and improved the performance of analytical queries against the table.

Exercise 3: Create a Memory Optimized Columnstore Table

► Task 1: Use the Memory Optimization Advisor

1. In SQL Server Management Studio, in **Object Explorer**, expand **Databases**, expand **AdventureWorksDW**, and then expand **Tables**.
2. Right-click **dbo.FactInternetSales**, and then click **Memory Optimization Advisor**.
3. In the **Table Memory Optimization Advisor** window, on the **Introduction** page, click **Next**.
4. On the **Migration validation** page, scroll down, looking at the descriptions against any row in the table without a green tick.
5. Note the advisor requires that the foreign key relationships are removed before it can continue.
6. Click **Cancel**.

► Task 2: Enable the Memory Optimization Advisor to Create a Memory Optimized FactInternetSales Table

1. On the **File** menu, point to **Open**, then click **File**.
2. Navigate to the **D:\Labfiles\Lab04\Starter** folder, click **Drop Columnstore Indexes on FactInternetSales.sql**, and then click **Open**.
3. Click **Execute** to drop the existing foreign keys.
4. In Object Explorer, right-click **dbo.FactInternetSales**, and then click **Memory Optimization Advisor**.
5. In the **Table Memory Optimization Advisor** window, on the **Introduction** page, click **Next**.
6. On the **Migration validation** page, click **Next**.
7. Review the warnings, and then click **Next**.
8. On the **Migration options** page, select **Also copy table data to the new memory optimized table**, and then click **Next**.
9. On the **Index creation** page, in the column list, select **SalesOrderNumber** and **SalesOrderLineNumber**, and then click **Next**.
10. On the **Summary** page, click **Migrate**.
11. When the wizard has completed, click **Ok**.

► Task 3: Examine the Performance of the Memory Optimized Table

1. On the **File** menu, point to **Open**, and then click **File**.
2. Navigate to the **D:\Labfiles\Lab04\Starter** folder, click **Query FactInternetSales.sql**, and then click **Open**.
3. On the **Query** menu, click **Include Actual Execution Plan**.

4. Highlight the Transact-SQL between the Step 1 and Step 2 comments, and then on the toolbar, click **Execute**.
5. Note that no disk space is used, because this table is now memory optimized.
6. Highlight the code after the Step 2 comment.
7. On the toolbar, click **Execute**, and then view the query results.
8. On the **Execution plan** tab, view the execution plan for the query. Hover the mouse pointer over the icons from right to left to examine their details and note the indexes that were used.
9. Close SQL Server Management Studio without saving changes.

Results: After completing this exercise, you will have created a memory optimized version of the **FactInternetSales** disk-based table, using the Memory Optimization Advisor.

Module 5: Implementing an Azure SQL Data Warehouse

Lab: Implement an Azure SQL Data Warehouse

Exercise 1: Create an Azure SQL Data Warehouse Database

► Task 1: Prepare the Lab Environment

1. Ensure that the **MT17B-WS2016-NAT**, **20767C-MIA-DC**, **20767C-MIA-CLI**, and **20767C-MIA-SQL** virtual machines are running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab05\Starter** folder, right-click **Setup.cmd** and then click **Run as administrator**.
3. Click **Yes** when prompted to confirm that you want to run the command file, and then wait for the script to finish.

► Task 2: Create a Data Warehouse and Server

1. Open Internet Explorer and browse to **https://portal.azure.com/**.
2. Sign in to the Azure portal with your Azure pass or Microsoft account credentials.
3. Click **New**, click **Databases**, and then click **SQL Data Warehouse**.
4. On the **SQL Data Warehouse** blade, in the **Database name** box, type **FoodOrdersDW**.
5. Under **Resource group**, click **Create new**, and then in the name box, type **FoodOrdersRG**.
6. Under **Server**, click **Configure required settings**, and then click **Create a new server**.
7. On the **New server** blade, in the **Server name** box, type **20767c**, followed by your initials, followed by today's date. For example, if your initials are CN and the date is 15 March, 2018, type **20767ccn15mar18**. The server name must be unique; if the name is unique and valid, a green tick appears.
8. In the **Server admin login** box, type **foodordersadmin**.
9. In the **Password** and **Confirm password** boxes, type **Pa55w.rd**
10. Under **Location**, select a region nearest your current geographical location, and then click **Select**.
11. In the **SQL Data Warehouse** blade, click **Performance Tier**.
12. Move the slider to **DW100**, and then click **Apply**.
13. On the **SQL Data Warehouse** blade, verify that **Select source** has the value **Blank database**.
14. Click **Create**. It will take some time for the new server and database to be created. The Azure portal will notify you when this step is finished.
15. Leave Internet Explorer open for the next task.

► Task 3: Configure the Azure Firewall

1. In the Azure portal, click **All resources**.
2. In the **All resources** blade, click **<Server Name>** (where **<Server Name>** is the name of the server you created in the previous task).
3. In the **<Server Name>** blade, click **Firewall / Virtual Networks**.
4. In the firewall settings blade, click **Add client IP**, and then click **Save**.

5. When the firewall changes are complete, click **OK**.
6. Leave Internet Explorer open.

► **Task 4: Connect to Azure SQL Data Warehouse**

1. Start SQL Server Management Studio and connect to the **MIA-SQL** database engine instance using Windows authentication.
2. On the **File** menu, click **Connect Object Explorer**.
3. In the **Connect to Server** dialog box, enter the following values, and then click **Connect**:
 - **Server Name:** <Server Name>.database.windows.net (where <Server Name> is the name of the server you created)
 - **Authentication:** SQL Server Authentication
 - **Login:** foodordersadmin
 - **Password:** Pa55w.rd
4. In **Object Explorer**, under <Server Name>.database.windows.net, expand **Databases**. Note that the **FoodOrdersDW** database is listed.
5. Leave SQL Server Management Studio open.

Results: After this exercise, you should have an Azure SQL Data Warehouse database called **FoodOrdersDW** on a logical server with a name you specified. The firewall settings of the server should be set so that you can connect to the server using SSMS on your computer.

Exercise 2: Migrate to an Azure SQL Data Warehouse Database

► **Task 1: Install the Data Warehouse Migration Utility**

1. In the **D:\Labfiles\Lab05\Starter\Installdwmigrationutility** folder, double-click **Data Warehouse Migration Utility.msi**.
2. In the **Open File - Security Warning** dialog box click **Run**.
3. In the **Data Warehouse Migration Utility** dialog box, on the **Select Installation Folder** page, click **Next**.
4. On the **License Agreement** page, review the License Agreement, click **I Agree**, and then click **Next**.
5. In the **User Account Control** dialog box, click **Yes**.
6. On the **Installation Complete** page, click **Close**.
7. Minimize all open applications and note the **Data Warehouse Migration Utility** icon on the desktop.

► Task 2: Check Compatibility of the Legacy Database

1. In SQL Server Management Studio, in **Object Explorer**, under **MIA-SQL**, expand **Databases**, expand **FoodOrdersDW**, and then expand **Tables**. Right-click each of the following tables, click **Select Top 1000 Rows**, and then review the data they contain:

- **dbo.DimCustomers**
- **dbo.DimProducts**
- **dbo.FactCallCentreOrders**
- **dbo.FactInternetOrders**

Note that the **dbo.FactCallCentreOrders** table is empty.

2. Minimize all windows, and then, on the desktop, double-click the **Data Warehouse Migration Utility** icon.
3. In the **DATA WAREHOUSE MIGRATION UTILITY** dialog box, note the following settings, and then click **Next**:
 - **Source Type**: SQL Server
 - **Destination Type**: Azure SQL Data Warehouse
4. In the **DATA WAREHOUSE MIGRATION UTILITY MIGRATION SOURCE** dialog box, in the **Server Name** box, type **MIA-SQL**, and then click **Connect**.
5. In the **Database List** pane, click **FoodOrdersDW**, and then click **Check Compatibility**.
6. In the **Save As** dialog box, save the **FoodOrdersDW - DatabaseCompatibilityReport.xlsx** file to the **D:\Labfiles\Lab05\Starter\Ex2** folder.
7. When prompted to open the Database Compatibility Report, click **No**.
8. Switch to the **20767C-MIA-CLI** virtual machine, and log in as **Student** with the password **Pa55w.rd**.
9. In File Explorer, in the address bar, type **\\MIA-SQL\D\$\Labfiles\Lab05\Starter\Ex2**, and then press Enter.
10. Double-click **FoodOrdersDW - DatabaseCompatibilityReport.xlsx**, review the data in the Excel spreadsheet (there should be two items assessed to have low impact), and then close Excel without saving changes.
11. Switch to the **20767C-MIA-SQL** virtual machine.

► Task 3: Migrate the Schema

1. In the **DATA WAREHOUSE MIGRATION UTILITY** dialog box, in the **Database List** pane, click **FoodOrdersDW**, and then click **Migrate Selected**.
2. In the **Migrate Database** pane, select the **Table Name** check box. Note that the check boxes for all tables in the database are now selected, and then click **Migrate Schema**.
3. In the **Migrate Schema** pane, note the script that has been generated for each table, and then click **Run Script**.
4. In the **DATA WAREHOUSE MIGRATION UTILITY MIGRATION DESTINATION** dialog box, enter the following values, click **OK**, and then wait for the script to run:
 - **Database**: FoodOrdersDW
 - **Server Name**: <Server Name>.database.windows.net (where <Server Name> is the name of the server you created previously)

- **Authentication:** Standard
 - **User Name:** foodordersadmin
 - **Password:** Pa55w.rd
5. In the **Script Applied Successfully** message box, click **OK**.

► **Task 4: Migrate the Data**

1. In the **Migrate Schema** pane, click **Migrate Data**.
2. In the **Migrate Data** pane, click **Run Migration**.
3. In the **DATA WAREHOUSE MIGRATION UTILITY** dialog box, in the **Package Path** box, type **D:\Labfiles\Lab05\Starter\Ex2**, and then click **Next**.
4. In the **DATA WAREHOUSE MIGRATION UTILITY MIGRATION DESTINATION** dialog box, click **Generate**.
5. In the **Package(s) Generated Successfully** message box, click **OK**.
6. Open a new command prompt window as Administrator and move to the **D:\Labfiles\Lab05\Starter\Ex2** folder.
7. Type **FoodOrdersDW_Export.bat**, to run the export script.
8. Type **FoodOrdersDW_Import.bat**, to run the import script.
9. Close the command prompt window.
10. Close the **DATA WAREHOUSE MIGRATION UTILITY** dialog box, and then close the **DATA WAREHOUSE MIGRATION UTILITY** window.

► **Task 5: Verify That the Data Has Been Migrated Successfully**

1. In SQL Server Management Studio, in **Object Explorer**, under **<Server Name>.database.windows.net** expand **Databases**, right-click **FoodOrdersDW**, and then click **New Query** (where **<Server Name>** is the name of the server you created).
2. In the **Microsoft SQL Server Management Studio** dialog, click **Cancel**.
3. In the query pane, type the following SQL.

```
SELECT * FROM DimCustomers
SELECT * FROM DimProducts
SELECT * FROM FactCallCentreOrders
SELECT * FROM FactInternetOrders
```

4. To check that the data has been migrated to each of the four tables, highlight the first line, click **Execute**, review the data, and then repeat the process for the three other lines. Verify that the **dbo.FactCallCentreOrders** table exists but is empty.

Results: After this exercise, all of the tables in the **FoodOrdersDW** database on the **MIA-SQL** server should be copied to the **FoodOrdersDW** database that you created in Azure SQL Data Warehouse, in the previous exercise. All of the tables should contain data, apart from the **FactCallCentreOrders** table.

Exercise 3: Copy Data with the Azure Data Factory

► Task 1: Create an Azure Data Factory

1. In Internet Explorer, in the Azure portal, click **New**, click **Databases**, and then click **Data Factory**.
2. On the **New data factory** blade, in the **Name** box, type **20767c**, followed by your initials, followed by today's date, followed by **DF**. For example, if your initials are CN and the date is 15 March, 2018, type **20767ccn15mar18DF**. The data factory name must be unique; if the name is unique and valid, a green tick appears.
3. Under **Resource group**, click **Use existing**, then select **FoodOrdersRG**.
4. Under **Version**, select **V1**.
5. Under **Location**, click the region nearest to your current geographical location, and then click **Create**. The Azure portal will notify you when the data factory has been created.

► Task 2: Create a Data Management Gateway

1. In the Azure portal, click **All resources**.
2. On the **All resources** page, click the new data factory.
3. On the data factory blade, click **Author and deploy**.
4. At the top of the **New data store** blade, click **More**, and then click **New integration runtime (gateway)**.
5. In the **Create** blade, in the **Integration runtime name** box, type **MIA-SQLGW**, and then click **OK**.
6. In the **Configure** blade, click **Download and install integration runtime (Self-hosted)**.
7. On the **Download Center** page, click **Download**.
8. On the **Choose the download you want** page, select the 64-bit integration runtime MSI file, and then click **Next**.
9. If Internet Explorer blocks popups from this site, click **Options for this site**, and then click **Always allow**.
10. In the **Do you want to run or save** message box, click **Run**.
11. In the **Microsoft Integration Runtime Setup** wizard, on the welcome page, click **Next**.
12. On the **End-User License Agreement** page, select **I accept the terms in the License Agreement**, and then click **Next**.
13. On the **Destination Folder** page, click **Next**.
14. On the **Ready to install Microsoft Integration Runtime** page, click **Install**.
15. In the **User Account Control** dialog box, click **Yes**.
16. On the **Completed the Microsoft Integration Runtime Setup Wizard** page, click **Finish**.
17. In the **User Account Control** dialog box, click **Yes**. The Microsoft Integration Runtime Configuration Manager will start running.
18. Return to the Azure portal, and click the Copy button adjacent to the **key1** key. Allow the webpage to access the clipboard if prompted.
19. In the Microsoft Integration Runtime Configuration Manager, paste the key into the text box, and then click **Register**.
20. When registration is complete, click **Close**.

- In the Azure portal, on the data factory blade, in the **Essentials** section, click **Linked Services**. Verify that there is a gateway called **MIA-SQLGW**.

► **Task 3: Create a Linked Server for the On-Premises SQL Server Database**

- On the data factory blade, click **Author and deploy**.
- Click **More**, click **New dataset**, and then click **SQL Server table**.
- In the **Drafts/Draft-1** blade, edit the script to read as follows, and then click **Deploy**:

```
{
  "name": "RecentCallCentreOrdersSqlServerLS",
  "properties": {
    "description": "",
    "type": "OnPremisesSqlServer",
    "typeProperties": {
      "connectionString": "Data Source=MIA-SQL;Initial
Catalog=RecentCallCentreOrders;Integrated Security=True",
      "gatewayName": "MIA-SQLGW",
      "userName": "ADVENTUREWORKS\\Student",
      "password": "Pa55w.rd"
    }
  }
}
```

You can copy this script from the **Task3.txt** file in the **D:\Labfiles\Lab05\Starter\Ex3** folder.

- Expand **Linked services**. Note that there is now a linked service called **RecentCallCentreOrdersSqlServerLS**.

► **Task 4: Create a Linked Server for the Azure SQL Data Warehouse Database**

- Click **More**, and then click **New dataset** again, and then click **Azure SQL Data Warehouse**.
- In the **Drafts/Draft-1** blade, edit the script to read as follows, and then click **Deploy**:

```
{
  "name": "FoodOrdersAzureSqlDWLS",
  "properties": {
    "description": "",
    "type": "AzureSqlDW",
    "typeProperties": {
      "connectionString": "Data Source=tcp:<Server
Name>.database.windows.net,1433;Initial Catalog=FoodOrdersDW;Integrated
Security=False;User ID=foodordersadmin@<Server Name>;Password=Pa55w.rd;"
    }
  }
}
```

Where <Server Name> is the name of the server you created previously.

You can copy this script from the **Task4.txt** file in the **D:\Labfiles\Lab05\Starter\Ex3** folder.

- Note that there is now another linked service, called **FoodOrdersAzureSqlDWLS**.

► **Task 5: Create a Dataset for a Table in the On-Premises SQL Server Database**

1. Click **More**, click **New dataset** once more, and then click **SQL Server table**.
2. In the **Drafts/Draft-1** blade, edit the script to read as follows, and then click **Deploy**:

```
{
  "name": "CallCentreOrdersSQLServerDS",
  "properties": {
    "type": "SqlServerTable",
    "linkedServiceName": "RecentCallCentreOrdersSqlServerLS",
    "typeProperties": {
      "tableName": "CallCentreOrders"
    },
    "availability": {
      "frequency": "Minute",
      "interval": 15
    },
    "external": true,
    "policy": {
      "externalData": {
        "retryInterval": "00:01:00",
        "retryTimeout": "00:10:00",
        "maximumRetry": 3
      }
    }
  }
}
```

You can copy this script from the **Task5.txt** file in the **D:\Labfiles\Lab05\Starter\Ex3** folder.

3. Expand **Datasets**. Note that there is now a dataset called **CallCentreOrdersSQLServerDS**.

► **Task 6: Create a Dataset for a Table in the Azure SQL Data Warehouse Database**

1. Click **More**, click **New dataset** for a final time and then click **Azure SQL Data Warehouse**.
2. In the **Drafts/Draft-1** blade, edit the script to read as follows, and then click **Deploy**:

```
{
  "name": "FactCallCentreOrdersAzureSQLDWDS",
  "properties": {
    "type": "AzureSqlDWTable",
    "linkedServiceName": "FoodOrdersAzureSqlDWLS",
    "typeProperties": {
      "tableName": "FactCallCentreOrders"
    },
    "availability": {
      "frequency": "Minute",
      "interval": 15
    }
  }
}
```

You can copy this script from the **Task6.txt** file in the **D:\Labfiles\Lab05\Starter\Ex3** folder.

3. Note that there is now another dataset, called **FactCallCentreOrdersAzureSQLDWDS**.

► Task 7: Create a Pipeline Activity to Copy the Data

1. Click **More**, then click **New pipeline**.
2. In the **Drafts/Draft-1** blade, edit the script to read as follows, and then click **Deploy**:

```
{
  "name": "FoodOrdersPL",
  "properties": {
    "description": "This pipeline has one Copy activity that copies data from an on-prem
SQL to Azure SQL DW",
    "activities": [
      {
        "type": "Copy",
        "typeProperties": {
          "source": {
            "type": "SqlSource",
            "sqlReaderQuery": "select * from CallCentreOrders"
          },
          "sink": {
            "type": "SqlSink",
            "writeBatchSize": 0,
            "writeBatchTimeout": "00:00:00"
          }
        },
        "inputs": [
          {
            "name": "CallCentreOrdersSQLServerDS"
          }
        ],
        "outputs": [
          {
            "name": "FactCallCentreOrdersAzureSQLDWS"
          }
        ],
        "policy": {
          "timeout": "01:00:00",
          "concurrency": 1,
          "executionPriorityOrder": "NewestFirst",
          "style": "StartOfInterval"
        },
        "scheduler": {
          "frequency": "Minute",
          "interval": 15
        },
        "name": "CopyFromSQLtoAzureSQLDW",
        "description": "Copy data from on-prem SQL server to AzureSQLDW"
      }
    ],
    "start": "<Current Date>T<Start Time>Z",
    "end": "<Current Date>T<Finish Time>Z",
    "isPaused": false,
    "pipelineMode": "Scheduled"
  }
}
```

Where *<Current Date>* is today's date in yyyy-mm-dd format, *<Start Time>* is the most recent time that is an exact quarter past the hour in hh:mm:ss format, and *<Finish Time>* is the next time that is a quarter of an hour after *<Start Time>* in hh:mm:ss format. For example, if the date is March 4, 2018 and the time is 11:38, the two lines beginning with the "start" and "end" properties should read:

```
"start": "2018-03-04T11:30:00Z",
```

```
"end": "2018-03-04T11:45:00Z",
```

You can copy this script from the **Task7.txt** file in the **D:\Labfiles\Lab05\Starter\Ex3** folder.

3. Note that there is now a pipeline called **FoodOrdersPL**.

4. Close the data factory blade.

► **Task 8: Check the Progress of the Pipeline**

1. Click **All Resources**, and then click your new data factory.
2. In the **Essentials** section, click **Diagram**.
3. To view the properties of the **FoodOrdersPL** pipeline, double-click **FoodOrdersPL**.
4. Note the properties in the **FoodOrdersPL** blade, and then close the **FoodOrdersPL** blade.
5. To view the properties of the **CallCentreOrdersSQLServerDS** dataset, double-click **CallCentreOrdersSQLServerDS**.
6. Note the properties in the **CallCentreOrdersSQLServerDS** blade. In the **Recently Updated Slices** section, there is a row with the slice start time and slice end time that you specified in the pipeline activity in the previous task. If the current time has not yet passed the slice end time, the status will read **Pending validation**. After the slice end time, the status will read **Ready**. Close the **CallCentreOrdersSQLServerDS** blade.
7. To view the properties of the **FactCallCentreOrdersAzureSQLDWDS** dataset, double-click **FactCallCentreOrdersAzureSQLDWDS**.
8. Note the properties in the **FactCallCentreOrdersAzureSQLDWDS** blade. In the **Recently Updated Slices** section, there is a row with the slice start time and slice end time that you specified in the pipeline in the previous task. If the current time has not yet passed the slice end time, the status will read **Pending execution**. After the slice end time, if the pipeline is still executing, the status will read **In progress**. After the pipeline has executed, the status will read **Ready**. Close the **FactCallCentreOrdersAzureSQLDWDS** blade.
9. After the pipeline has run, in SQL Server Management Studio, in Object Explorer, under **<Server Name>.database.windows.net** (where **<Server Name>** is the name of the server you created), under **Databases**, right-click **FoodOrdersDW**, and then click **New Query**. If the **Microsoft SQL Server Management Studio** message box appears, click **Cancel**.
10. In the query pane, type the following Transact-SQL, highlight the code, and then click **Execute**. Note that the data has been successfully copied from the on-premises table:


```
SELECT * FROM FactCallCentreOrders
```
11. Close SQL Server Management Studio without saving any files.
12. In the Azure portal, click **Resource groups**, click **FoodOrdersRG**, and then click **FoodOrdersDW**.
13. In the **FoodOrdersDW** blade, click **Pause**, and then, in the **Pause 'FoodOrdersDW'** message box, click **Yes**.
14. It will take some time for the database to pause. The Azure portal will notify you when this step is finished.
15. Log out of the Azure portal.
16. Close Internet Explorer.

Results: After this exercise, the data in the **CallCentreOrders** table, in the **RecentCallCentreOrders** database on the **MIA-SQL** server, should be copied to the **FactCallCentreOrders** table in the Azure SQL Data Warehouse database.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 6: Creating an ETL Solution

Lab: Implementing Data Flow in an SSIS Package

Exercise 1: Exploring Source Data

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab06\Starter** folder, right-click **Setup.cmd** and then click **Run as administrator**.
3. Click **Yes** when prompted to confirm that you want to run the command file, and then wait for the script to finish.

► Task 2: Extract and View Sample Source Data

1. Click **Search Windows**, type **Import and Export**, and then click **SQL Server Import and Export Data (64-bit)**.
2. On the **Welcome to SQL Server Import and Export Wizard** page, click **Next**.
3. On the **Choose a Data Source** page, specify the following options, and then click **Next**:
 - **Data source:** SQL Server Native Client 11.0
 - **Server name:** localhost
 - **Authentication:** Windows Authentication
 - **Database:** InternetSales
4. On the **Choose a Destination** page, set the following options, and then click **Next**:
 - **Destination:** Flat File Destination
 - **File name:** D:\Labfiles\Lab06\Starter\Top 1000 Customers.csv
 - **Text qualifier:** " (a quotation mark)
5. On the **Specify Table Copy or Query** page, select **Write a query to specify the data to transfer**, and then click **Next**.
6. On the **Provide a Source Query** page, enter the following Transact-SQL code, and then click **Next**:

```
SELECT TOP 1000 * FROM Customers
```
7. On the **Configure Flat File Destination** page, review the default options, and then click **Next**.
8. On the **Save and Run Package** page, select only **Run immediately**, and then click **Next**.
9. On the **Complete the Wizard** page, click **Finish**.
10. When the data extraction has completed successfully, click **Close**.
11. Open the file **Top 1000 Customers.csv** in the **D:\Labfiles\Lab06\Starter** folder by using Notepad.
12. Examine the data, noting the columns that exist in the **Customers** table and the range of data values they contain, and then close Notepad without saving any changes.

► Task 3: Profile Source Data

1. Start SQL Server Data Tools, and on the **File** menu, point to **New**, and then click **Project**.
2. In the **New Project** dialog box, enter the following, and then click **OK**:
 - **Project Template:** Integration Services Project
 - **Name:** Explore Internet Sales
 - **Location:** D:\Labfiles\Lab06\Starter
 - **Solution name:** Explore Internet Sales
3. In the Solution Explorer pane, right-click **Connection Managers**, and then click **New Connection Manager**.
4. In the **Add SSIS Connection Manager** dialog box, click **ADO.NET**, and then click **Add**.
5. In the **Configure ADO.NET Connection Manager** dialog box, click **New**.
6. In the **Connection Manager** dialog box, enter the following values, and then click **OK**:
 - **Server name:** localhost
 - **Log on to the server:** Windows Authentication
 - **Select or enter a database name:** InternetSales
7. In the **Configure ADO.NET Connection Manager** dialog box, verify that a data connection named **localhost.InternetSales** is listed, and then click **OK**.
8. If the SSIS Toolbox pane is not visible, click the **Package.dtsx [Design]** window, and then on the **SSIS** menu, click **SSIS Toolbox**.
9. In the SSIS Toolbox pane, in the **Common** section, double-click **Data Profiling Task** to add it to the Control Flow surface.
10. Double-click the **Data Profiling Task** icon on the Control Flow surface to open its editor.
11. In the **Data Profiling Task Editor** dialog box, on the **General** tab, in the **Destination** property value drop-down list, click **<New File connection...>**.
12. In the **File Connection Manager Editor** dialog box, in the **Usage type** drop-down list, click **Create file**.
13. In the **File** box, type **D:\Labfiles\Lab06\Starter\Internet Sales Data Profile.xml**, and then click **OK**.
14. In the **Data Profiling Task Editor** dialog box, on the **Profile Requests** page, in the **Profile Type** drop-down list, select **Column Statistics Profile Request**, and then click the **RequestID** column.
15. In the **Request Properties** pane, set the following property values. Do not click **OK** when finished:
 - **ConnectionManager:** localhost.InternetSales
 - **TableOrView:** [dbo].[SalesOrderHeader]
 - **Column:** OrderDate
16. In the **Data Profiling Task Editor** dialog box, on the **Profile Requests** page, in the **Profile Type** drop-down list for the empty row below the profile you just added, select **Column Length Distribution Profile Request**, and then click the **RequestID** column.

17. In the **Request Properties** pane, set the following property values. Do not click **OK** when finished:
 - **ConnectionManager**: localhost.InternetSales
 - **TableOrView**: [dbo].[Customers]
 - **Column**: AddressLine1
 - **IgnoreLeadingSpaces**: False
 - **IgnoreTrailingSpaces**: True
18. In the **Data Profiling Task Editor** dialog box, on the **Profile Requests** page, in the **Profile Type** drop-down list for the empty row below the profile you just added, select **Column Null Ratio Profile Request**, and then click the **RequestID** column.
19. In the **Request Properties** pane, set the following property values. Do not click **OK** when finished:
 - **ConnectionManager**: localhost.InternetSales
 - **TableOrView**: [dbo].[Customers]
 - **Column**: AddressLine2
20. In the **Data Profiling Task Editor** dialog box, on the **Profile Requests** page, in the **Profile Type** drop-down list for the empty row below the profile you just added, select **Value Inclusion Profile Request**, and then click the **RequestID** column.
21. In the **Request Properties** pane, set the following property values:
 - **ConnectionManager**: localhost.InternetSales
 - **SubsetTableOrView**: [dbo].[SalesOrderHeader]
 - **SupersetTableOrView**: [dbo].[PaymentTypes]
 - **InclusionColumns**:
 - **Subset side Columns**: PaymentType
 - **Superset side Columns**: PaymentTypeKey
 - **InclusionThresholdSetting**: None
 - **SupersetColumnsKeyThresholdSetting**: None
 - **MaxNumberOfViolations**: 100
22. In the **Data Profiling Task Editor** dialog box, click **OK**.
23. On the **Debug** menu, click **Start Debugging**.
24. When the Data Profiling task has completed, with the package still running, double-click the **Data Profiling** task, and then click **Open Profile Viewer**.
25. Maximize the **Data Profile Viewer** window, and under the **[dbo].[SalesOrderHeader]** table, click **Column Statistics Profiles**. Review the minimum and maximum values for the **OrderDate** column.
26. Under the **[dbo].[Customers]** table, click **Column Length Distribution Profiles** and click the **AddressLine1** column to view the statistics. Click the bar chart for any of the column lengths, and then click the **Drill Down** button to view the source data that matches the selected column length.
27. Under the **[dbo].[Customers]** table, click **Column Null Ratio Profiles** and view the null statistics for the **AddressLine2** column. Select the **AddressLine2** column, and then click the **Drill Down** button to view the source data.

28. Under the **[dbo].[SalesOrderHeader]** table, click **Inclusion Profiles** and review the inclusion statistics for the **PaymentType** column. Select the inclusion violation for the payment type value of **0**, and then click the **Drill Down** button to view the source data.
29. Close the Data Profile Viewer window, and then in the **Data Profiling Task Editor** dialog box, click **Cancel**.
30. On the **Debug** menu, click **Stop Debugging**.
31. Close SSDT, saving your changes if you are prompted.

Results: After this exercise, you should have a comma-separated text file that contains a sample of customer data, and a data profile report that shows statistics for data in the **InternetSales** database.

Exercise 2: Transferring Data by Using a Data Flow Task

► Task 1: Examine an Existing Data Flow

1. Start SSDT and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab06\Starter\Ex2** folder.
2. In the Solution Explorer pane, expand **SSIS Packages** and double-click **Extract Reseller Data.dtsx** to open it in the designer.
3. On the Control Flow surface, note that the package contains two Data Flow tasks.
4. Double-click the **Extract Resellers** Data Flow task to view it on the **Data Flow** tab and note that it contains a data source component named **Resellers** and a data destination component named **Staging DB**.
5. Double-click the **Resellers** data source component, note the following details, and then click **Cancel**:
 - On the **Connection Manager** page, the data source is configured to use an OLE DB connection manager named **localhost.ResellerSales** and extracts data from the **[dbo].[Resellers]** table.
 - On the **Columns** page, the data source is configured to retrieve every column from the **Resellers** table, and the output columns that the task generates have the same names as the source columns.
6. Double-click the **Staging DB** data destination component, note the following details, and then click **Cancel**:
 - On the **Connection Manager** page, the data destination is configured to use an OLE DB connection manager named **localhost.Staging** and to use a **Table or view - fast load** mode to insert data into the **[dbo].[Resellers]** table.
 - On the **Mappings** page, the data destination is configured to map the input columns (which are the output columns from the **Resellers** data source) to columns in the destination table. The order of the columns in the destination is different from the column order in the source, and the source **ResellerKey** column is mapped to the **ResellerBusinessKey** destination column.
7. Right-click anywhere on the **Data Flow** surface, click **Execute Task**, and then observe the task as it runs. Note the number of rows that have been transferred—this is shown in front of the arrow connecting the two components in the data flow.
8. On the **Debug** menu, click **Stop Debugging**.

► **Task 2: Create a Data Flow Task**

1. In Solution Explorer, right-click **SSIS Packages**, and then click **New SSIS Package**.
2. Right-click **Package1.dtsx**, click **Rename**, and then change the name of the package to **Extract Internet Sales Data.dtsx**.
3. With the **Extract Internet Sales Data.dtsx** package open, and the Control Flow surface visible, in the SSIS Toolbox pane, double-click **Data Flow Task**, and then drag the new Data Flow task to the center of the Control Flow surface.
4. Right-click **Data Flow Task** on the Control Flow surface, click **Rename**, and then change the name of the task to **Extract Customers**.
5. Double-click the **Extract Customers** Data Flow task to view the Data Flow surface.

► **Task 3: Add a Data Source Component to a Data Flow**

1. In Solution Explorer, right-click **Connection Managers**, and then click **New Connection Manager**.
2. In the **Add SSIS Connection Manager** dialog box, click **OLEDB**, and then click **Add**.
3. In the **Configure OLE DB Connection Manager** dialog box, click **New**.
4. In the **Connection Manager** dialog box, enter the following values, and then click **OK**:
 - **Server name:** localhost
 - **Log on to the server:** Windows Authentication
 - **Select or enter a database name:** InternetSales

 **Note:** When you create a connection manager, it is named automatically, based on the server and database name—for example, **localhost.InternetSales**. If you have previously created a connection manager for the same database, a name such as **localhost.InternetSales1** may be generated.

5. In the **Configure OLE DB Connection Manager** dialog box, verify that a new data connection is listed, and then click **OK**.
6. In the SSIS Toolbox pane, in the **Favorites** section, double-click **Source Assistant**.
7. In the **Source Assistant - Add New Source** dialog box, in the **Select source type** list, select **SQL Server**; in the **Select connection manager** list, select the connection manager for the **localhost.InternetSales** database that you created previously, and then click **OK**.
8. Drag the new **OLE DB Source** data source component to the center of the Data Flow surface, right-click **OLE DB Source**, click **Rename**, and then change the name of the data source component to **Customers**.
9. Double-click the **Customers** source component, set the following configuration values, and then click **OK**:
 - On the **Connection Manager** page, ensure that the OLE DB connection manager for the **localhost.InternetSales** database is selected, ensure that the **Table or view** data access mode is selected, and then in the **Name of the table or the view** drop-down list, click **[dbo].[Customers]**.
 - On the **Columns** page, ensure that every column from the **Customers** table is selected, and that the output columns have the same names as the source columns.

► Task 4: Add a Data Destination Component to a Data Flow

1. In the SSIS Toolbox pane, in the **Favorites** section, double-click **Destination Assistant**.
2. In the **Destination Assistant - Add New Destination** dialog box, in the **Select destination type** list, click **SQL Server**. In the **Select connection managers** list, click **localhost.Staging**, and then click **OK**.
3. Drag the new **OLE DB Destination** data destination component below the **Customers** data source component, right-click **OLE DB Destination**, click **Rename**, and then change the name of the data destination component to **Staging DB**.
4. On the Data Flow surface, click the **Customers** source component, and then drag the blue arrow from the **Customers** data source component to the **Staging DB** destination component.
5. Double-click the **Staging DB** destination component, set the following configuration values, and then click **OK**:
 - On the **Connection Manager** page, in the **Name of the table or the view** drop-down list, click **[dbo].[Customers]**, and then select **Keep nulls**.
 - On the **Mappings** page, drag the **CustomerKey** column from the list of available input columns to the **CustomerBusinessKey** column in the list of available destination columns. Then verify that all other input columns are mapped to destination columns of the same name.

► Task 5: Test the Data Flow Task

1. Right-click anywhere on the Data Flow surface, click **Execute Task**, and then observe the task as it runs. Note the number of rows that have been transferred, which is shown in front of the arrow connecting the two components in the data flow.
2. On the **Debug** menu, click **Stop Debugging**.
3. Close SSDT.

Results: After this exercise, you should have an SSIS package that contains a single Data Flow task, which extracts customer records from the **InternetSales** database and inserts them into the **Staging** database.

Exercise 3: Using Transformation Components in a Data Flow

► Task 1: Examine an Existing Data Flow

1. Start SSDT and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab06\Starter\Ex3** folder.
2. In Solution Explorer, expand **SSIS Packages** and double-click **Extract Reseller Data.dtsx**.
3. On the **Data Flow** tab, in the **Data Flow Task** drop-down list, click **Extract Reseller Sales**. Note that this data flow includes a data source component, two transformation components, and two destination components.
4. Double-click the **Reseller Sales** data source component, and in the **OLE DB Source Editor** dialog box, review the following details, and then click **Cancel**:
 - On the **Connection Manager** page, the data source is configured to use an OLE DB connection manager named **localhost.ResellerSales** and to extract data using a Transact-SQL command that queries the **SalesOrderHeader**, **SalesOrderDetail**, and **PaymentTypes** tables. The query includes an ISNULL function to check for a payment type. If none is specified, the value "other" is used.
 - On the **Columns** page, the data source is configured to retrieve several columns including **ProductKey**, **OrderQuantity**, and **UnitPrice**.
5. Double-click the **Calculate Sales Amount** transformation component, in the **Derived Column Transformation Editor** dialog box, review the following details, and then click **Cancel**:
 - The transformation creates a derived column named **SalesAmount**.
 - The derived column is added as a new column to the data flow.
 - The column value is calculated by multiplying the **UnitPrice** column value by the **OrderQuantity** column value.
6. Double-click the **Lookup Product Details** transformation component, and in the **Lookup Transformation Editor** dialog box, review the following details, and then click **Cancel**:
 - On the **General** page, the lookup transformation is configured to use full cache mode and an OLE DB connection manager, and to redirect rows with no matching entries.
 - On the **Connection** page, the lookup transformation is configured to return the results of a Transact-SQL query using the **localhost.Products** OLE DB connection manager. The query returns a table that contains product data from the **Products** database.
 - On the **Columns** page, the lookup transformation is configured to match rows in the data flow with products data based on the **ProductKey** column value, and add all the other columns in the products data as new columns in the data flow.
7. View the arrow between the **Lookup Product Details** transformation component and the **Staging DB** destination component. Note that this arrow represents the lookup match output, so rows where a matching product record was found will follow this data flow path.
8. View the arrow between the **Lookup Product Details** transformation component and the **Orphaned Sales** destination component. Note that this arrow represents the lookup no match output, so rows where no matching product record was found will follow this data flow path.

9. Double-click the **Staging DB** data destination component, review the following details, and then click **Cancel**:
 - On the **Connection Manager** page, the data destination is configured to use an OLE DB connection manager named **localhost.Staging** and to use a **Table or view - fast load** access mode to insert data into the **[dbo].[ResellersSales]** table.
 - On the **Mappings** tab, the order of the columns in the destination is different from the column order in the source, and the **ProductKey** and **ResellerKey** source columns are mapped to the **ProductBusinessKey** and **ResellerBusinessKey** destination columns.
10. Double-click the **Orphaned Sales** destination component, note that it uses a flat file connection manager named **Orphaned Reseller Sales**, and then click **Cancel**.
11. In the Connection Managers pane at the bottom of the design surface, double-click the **Orphaned Reseller Sales** connection manager, note that the rows for sales with no matching product record are redirected to the **Orphaned Reseller Sales.csv** text file in the **D:\ETL** folder, and then click **Cancel**.
12. Right-click anywhere on the Data Flow surface, click **Execute Task**, and then observe the task as it runs. Note the number of rows that have been transferred, which is shown in front of the arrows connecting the components in the data flow. There should be no orphaned sales records.
13. On the **Debug** menu, click **Stop Debugging**.
14. Keep SSDT open for the next task.

► Task 2: Create a Data Flow Task

1. In the Solution Explorer pane, double-click **Extract Internet Sales Data.dtsx**.
2. View the **Control Flow** tab, and then in the SSIS Toolbox pane, in the **Favorites** section, double-click **Data Flow Task**.
3. Drag the new **Data Flow** task below the existing **Extract Customers** task.
4. Right-click the new **Data Flow** task, click **Rename**, and then change the name to **Extract Internet Sales**.
5. Click the **Extract Customers** Data Flow task, and then drag the arrow from the **Extract Customers** task to the **Extract Internet Sales** task.
6. Double-click the **Extract Internet Sales** task to view the Data Flow surface.

► Task 3: Add a Data Source Component to a Data Flow

1. In the SSIS Toolbox pane, in the **Favorites** section, double-click **Source Assistant**.
2. In the **Source Assistant - Add New Source** dialog box, in the **Select source type** list, select **SQL Server**. In the **Select connection manager** list, select **localhost.InternetSales**, and then click **OK**.
3. Drag the new **OLE DB Source** data source component to the center of the Data Flow surface, right-click **OLE DB Source**, click **Rename**, and then change the name of the data source component to **Internet Sales**.
4. Double-click the **Internet Sales** source component, set the following configuration values, and then click **OK**:
 - On the **Connection Manager** page, ensure that the **localhost.InternetSales** OLE DB connection manager is selected. Change the **Data access mode** to **SQL command**. Click **Browse**. In the **Open** dialog box, open the **InternetSales.sql** file in the **D:\Labfiles\Lab06\Starter\Ex3** folder. Click **Preview** to see a data preview, and then click **Close** to close the **Preview Query Results** dialog box.

- On the **Columns** page, ensure that every column from the query is selected, and that the output columns have the same names as the source columns.

► **Task 4: Add a Derived Column Transformation Component to a Data Flow**

1. In the SSIS Toolbox pane, in the **Common** section, double-click **Derived Column**.
2. Drag the new **Derived Column** transformation component below the existing **Internet Sales** data source component, right-click **Derived Column**, click **Rename**, and then change the name of the transformation component to **Calculate Sales Amount**.
3. On the Data Flow surface, click the **Internet Sales** source component, and then drag the blue arrow from the **Internet Sales** data source component to the **Calculate Sales Amount** transformation component.
4. Double-click the **Calculate Sales Amount** transformation component.
5. In the **Derived Column Transformation Editor** dialog box, perform the following steps, and then click **OK**:
 - In the **Derived Column Name** column, type **SalesAmount**.
 - In the **Derived Column** column, ensure that **<add as new column>** is selected.
 - Expand the **Columns** folder, and then drag the **UnitPrice** column to the **Expression** box.
 - In the **Expression** box, after **[UnitPrice]**, type *****, and then drag the **OrderQuantity** column to the Expression box so that the expression looks like the following example:

```
[UnitPrice] * [OrderQuantity]
```

- Ensure that the **Data Type** column contains the value **numeric [DT_NUMERIC]**, the **Precision** column contains the value **25**, and the **Scale** column contains the value **4**.

► **Task 5: Add a Lookup Transformation Component to a Data Flow**

1. In the SSIS Toolbox pane, in the **Common** section, double-click **Lookup**.
2. Drag the new **Lookup** transformation component below the existing **Calculate Sales Amount** transformation component, right-click **Lookup**, click **Rename**, and then change the name of the transformation component to **Lookup Product Details**.
3. On the Data Flow surface, click the **Calculate Sales Amount** transformation component, and then drag the blue arrow from the **Calculate Sales Amount** transformation component to the **Lookup Product Details** transformation component.
4. Double-click the **Lookup Product Details** transformation component.
5. In the **Lookup Transformation Editor** dialog box, perform the following steps, and then click **OK**:
 - On the **General** page, in the **Specify how to handle rows with no matching entries** list, click **Redirect rows to no match output**.
 - On the **Connection** page, in the **OLE DB connection manager** list, select **localhost.Products**, and then click **OK**. Select **Use results of an SQL query**, click **Browse**, and import the **Products.sql** query file from the **D:\Labfiles\Lab06\Starter\Ex3** folder.
 - On the **Columns** page, drag **ProductKey** from the **Available Input Columns** list to **ProductKey** in the **Available Lookup Columns** list.
 - In the **Available Lookup Columns** list, select the check box next to the **Name** column heading to select all columns, and then clear the check box for the **ProductKey** column.

6. In the SSIS Toolbox pane, in the **Other Destinations** section, double-click **Flat File Destination**.
7. Drag the new flat file destination component to the right of the existing **Lookup Product Details** transformation, right-click **Flat File Destination**, click **Rename**, and then change the name of the destination component to **Orphaned Sales**.
8. On the Data Flow surface, click the **Lookup Product Details** transformation component, and then drag the blue arrow from the **Lookup Product Details** transformation component to the **Orphaned Sales** destination component.
9. In the **Input Output Selection** dialog box, in the **Output** list, click **Lookup No Match Output**, and then click **OK**.
10. Double-click the **Orphaned Sales** destination component, and then in the **Flat File Destination Editor** dialog box, next to the **Flat File connection manager** drop-down list, click **New**.
11. In the **Flat File Format** dialog box, click **Delimited**, and then click **OK**.
12. In the **Flat File Connection Manager Editor** dialog box, change the text in the **Connection manager name** box to **Orphaned Internet Sales**.
13. On the **General** page, set the **File name** value to **D:\Labfiles\Lab06\Starter\Ex3\Orphaned Internet Sales.csv**, select **Column names in the first data row**, and then click **OK**.
14. In the **Flat File Destination Editor** dialog box, ensure that **Overwrite data in the file** is selected, and then click **Mappings**. Verify that all input columns are mapped to destination columns with the same name, and then click **OK**.

► **Task 6: Add a Data Destination Component to a Data Flow**

1. In the SSIS Toolbox pane, in the **Favorites** section, double-click **Destination Assistant**.
2. In the **Destination Assistant - Add New Destination** dialog box, in the **Select destination type** list, click **SQL Server**. In the **Select connection manager** list, click **localhost.Staging**, and then click **OK**.
3. Drag the new **OLE DB Destination** data destination component below the **Lookup Product Details** transformation component, right-click **OLE DB Destination**, click **Rename**, and then change the name of the data destination component to **Staging DB**.
4. On the Data Flow surface, click the **Lookup Product Details** transformation component, and then drag the blue arrow from the **Lookup Product Details** transformation component to the **Staging DB** destination component. Note that the **Lookup Match Output** is automatically selected.
5. Double-click the **Staging DB** destination component, set the following configuration values, and then click **OK**:
 - On the **Connection Manager** page, ensure that the **localhost.Staging** OLE DB connection manager is selected, and ensure that the **Table or view – fast load** data access mode is selected.
 - In the **Name of the table or the view** drop-down list, click **[dbo].[InternetSales]**, and select **Keep nulls**.

- On the **Mappings** page, drag the following **Key** columns from the list of available input columns to the **BusinessKey** corresponding columns in the list of available destination columns:

Available Input Columns	Available Destination Columns
ProductKey	ProductBusinessKey
CustomerKey	CustomerBusinessKey
ProductSubcategoryKey	ProductSubcategoryBusinessKey
ProductCategoryKey	ProductCategoryBusinessKey

- Verify that all other input columns are mapped to destination columns of the same name.

► **Task 7: Test the Data Flow Task**

1. Right-click anywhere on the Data Flow surface, click **Execute Task**, and then observe the task as it runs. Note the number of rows that have been transferred, which is shown in front of the arrows connecting the components in the data flow. There should be no orphaned sales records.
2. On the **Debug** menu, click **Stop Debugging**.
3. Close SSDT.

Results: After this exercise, you should have a package that contains a Data Flow task including derived column and lookup transformation components.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 7: Implementing Control Flow in an SSIS Package

Lab A: Implementing Control Flow in an SSIS Package

Exercise 1: Using Tasks and Precedence in a Control Flow

► Task 1: Prepare the Lab Environment

1. Ensure the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab07A\Starter** folder, right-click **Setup.cmd** and click **Run as administrator**.
3. Click **Yes** when prompted to confirm you want to run the command file, and wait for the script to finish.

► Task 2: View the Control Flow

1. Start **SQL Server Data Tools**.
2. On the **File** menu, point to **Open**, and click **Project/Solution**, browse to **AdventureWorksETL.sln** in the **D:\Labfiles\Lab07A\Starter\Ex1** folder, and then click **Open**.
3. In Solution Explorer, in the **SSIS Packages** folder, double-click **Extract Reseller Data.dtsx**.
4. Review the control flow for the **Extract Reseller Data package** and note that it includes two **Send Mail** tasks—one that runs when either the **Extract Resellers** or **Extract Reseller Sales** tasks fail, and one that runs when the **Extract Reseller Sales** task succeeds.
5. Double-click the red dotted arrow connecting the **Extract Resellers** task to the **Send Failure Notification** task. In the **Precedence Constraint Editor**, in the **Multiple constraints** section, note that **Logical OR. One constraint must evaluate to True** is selected so that the **Send Failure Notification** task runs if either of the Data Flow tasks connected should fail. Click **Cancel** to close **Precedence Constraint Editor**.
6. Double-click the **Send Failure Notification** task to view its settings. On the **Mail** tab, note that the task uses an SMTP connection manager named **Local SMTP Server** to send a high-priority email message with the subject **Data Extraction Notification** and the message "The Reseller data extraction process failed!" to Student@adventureworks.msft. Click **Cancel** to close the **Send Mail Task Editor**.
7. Double-click the **Send Success Notification** task to view its settings. On the **Mail** tab, note that the task uses an SMTP connection manager named **Local SMTP Server** to send a normal priority email message with the subject **Data Extraction Notification** and the message "The Reseller data was successfully extracted" to Student@adventureworks.msft. Then click **Cancel** to close the **Send Mail Task Editor**.
8. In the Connection Managers pane, double-click **Local SMTP Server** to view its settings, and note that it connects to the **localhost** SMTP server. Click **Cancel**.
9. On the **Debug** menu, click **Start Debugging**, and observe the control flow as the task executes. When the task has completed, on the **Debug** menu, click **Stop Debugging**.
10. View the contents of the **C:\inetpub\mailroot\Drop** folder and note the email messages that have been received by the local SMTP server.

11. Open the email file by using Notepad and verify that the subject and sender match those specified by the **Send Success Notification** task (the body of the text message will be encoded so you will not be able to read it).

► **Task 3: Add Tasks to a Control Flow**

1. In SQL Server Data Tools, in Solution Explorer, in the **SSIS Packages** folder, double-click **Extract Internet Sales Data.dtsx**. Then view the control flow for the **Extract Internet Sales Data** package.
2. In the **SSIS Toolbox**, in the **Common** section, double-click **Send Mail Task**.
3. Position the new **Send Mail Task** below and to the right of the **Extract Internet Sales** task.
4. Double-click the **Send Mail Task** on the control flow surface, to view its settings.
5. In the **Send Mail Task Editor** dialog box, on the **General** tab, set the **Name** property to **Send Success Notification**.
6. In the **Send Mail Task Editor** dialog box, on the **Mail** tab, in the **SmtpConnection** drop-down list, click **<New connection>**.
7. In the **SMTP Connection Manager Editor** dialog box, enter the following settings, and then click **OK**:
 - **Name**: Local SMTP Server
 - **SMTP Server**: localhost
8. In the **Send Mail Task Editor** dialog box, on the **Mail** tab, enter the following settings, and then click **OK**:
 - **From**: ETL@adventureworks.msft
 - **To**: Student@adventureworks.msft
 - **Subject**: Data Extraction Notification - Success
 - **MessageSourceType**: Direct Input
 - **MessageSource**: The Internet Sales data was successfully extracted
 - **Priority**: Normal
9. On the Control Flow surface, click the **Extract Internet Sales** task, and then drag the green arrow from the **Extract Internet Sales** task to the **Send Success Notification** task.
10. In the SSIS Toolbox, in the **Common** section, double-click **Send Mail Task**.
11. Position the new **Send Mail Task** below and to the left of the **Extract Internet Sales** task.
12. Double-click the **Send Mail Task** on the Control Flow surface, to view its settings.
13. In the **Send Mail Task Editor** dialog box, on the **General** tab, set the **Name** property to **Send Failure Notification**.
14. In the **Send Mail Task Editor** dialog box, on the **Mail** tab, enter the following settings, and then click **OK**:
 - **SmtpConnection**: Local SMTP Server
 - **From**: ETL@adventureworks.msft
 - **To**: Student@adventureworks.msft
 - **Subject**: Data Extraction Notification - Failure
 - **MessageSourceType**: Direct Input

- **MessageSource:** The Internet Sales data extraction process failed
 - **Priority:** High
15. On the Control Flow surface, click the **Extract Customers** task, and then drag the green arrow from the **Extract Customers** task to the **Send Failure Notification** task. Then right-click the arrow and click **Failure**.
 16. On the Control Flow surface, click the **Extract Internet Sales** task, and then drag the green arrow from the **Extract Internet Sales** task to the **Send Failure Notification** task. Then right-click the arrow and click **Failure**.
 17. Double-click the red arrow connecting the **Extract Customers** task to the **Send Failure Notification** task.
 18. In the **Precedence Constraint Editor**, in the **Multiple Constraints** section, select **Logical OR. One constraint must evaluate to True**, and then click **OK**.

► Task 4: Test the Control Flow

1. In **SQL Server Data Tools**, on the Control Flow surface for the **Extract Internet Sales Data** package, click the **Extract Customers** task, and press F4.
2. In the Properties pane, set the **ForceExecutionResult** property to **Failure**.
3. On the **Debug** menu, click **Start Debugging**, and observe the control flow as the task executes, noting that the **Extract Customer task** fails.
4. When the task has completed, on the **Debug** menu, click **Stop Debugging**.
5. View the contents of the **C:\inetpub\mailroot\Drop** folder and note the email messages that have been received by the local SMTP server.
6. Open the most recent email file by using Notepad and note that the Subject line indicates data extraction failure. Close Notepad.
7. In SQL Server Data Tools, on the Control Flow surface for the **Extract Internet Sales Data** package, click the **Extract Customers** task. In the Properties pane, set the **ForceExecutionResult** property to **None**.
8. On the **Debug** menu, click **Start Debugging**, and observe the control flow as the task executes, noting that the **Extract Customer** task succeeds. Then, when the task has completed, on the **Debug** menu, click **Stop Debugging**.
9. View the contents of the **C:\inetpub\mailroot\Drop** folder and note the email messages that have been received by the local SMTP server.
10. Open the most recent message using Notepad, and note that the subject line indicates data transfer success. Close Notepad.
11. Close SQL Server Data Tools, saving your changes if prompted.

Results: After this exercise, you should have a control flow that sends an email message if the **Extract Internet Sales** task succeeds, or sends an email message if either the **Extract Customers** or **Extract Internet Sales** tasks fail.

Exercise 2: Using Variables and Parameters

► Task 1: View the Control Flow

1. View the contents of the **D:\Accounts** folder and note the files it contains. In this exercise, you will modify an existing package to create a dynamic reference to one of these files.
2. Start SQL Server Data Tools.
3. On the **File** menu, point to **Open**, and click **Project/Solution**, browse to **AdventureWorksETL.sln** in the **D:\Labfiles\Lab07A\Starter\Ex2** folder, and then click **Open**.
4. In Solution Explorer, in the **SSIS Packages** folder, double-click **Extract Payment Data.dtsx**.
5. Review the control flow for the **Extract Payment Data** package, and note that it contains a single Data Flow task named **Extract Payments**.
6. Double-click the **Extract Payments** task to view it in the **Data Flow** tab, and note that it contains a flat file source named **Payments File**, and an OLE DB destination named **localhost.Staging**. This destination connects to the Staging database.
7. Double-click the **Payments File** source and note that it uses a connection manager named **Payments File**. Click **Cancel**.
8. In the Connection Managers pane, double-click **Payments File**, and note that it references the **Payments.csv** file in the **D:\Labfiles\Lab07A\Starter\Ex2** folder. Click **Cancel**.
9. On the **Debug** menu, click **Start Debugging** and observe the data flow while the package runs. When the package has completed, on the **Debug** menu, click **Stop Debugging**.
10. On the **Output** tab, find the following line in the package execution log:

```
[Payments File [2]: The processing of the file "D:\Labfiles\Lab07A\Starter\Ex2\Payments.csv" has started
```

11. Click the **Data Flow** tab to return to the data flow design surface.

► Task 2: Create a Variable

1. In SQL Server Data Tools, with the **Extract Payment Data.dtsx** package open, on the **View** menu, point to **Other Windows**, and then click **Variables**.
2. In the Variables pane, click **Add Variable**, and create a variable with the following properties:
 - **Name:** fName
 - **Scope:** Extract Payment Data
 - **Data type:** String
 - **Value:** Payments - US.csv

Note that the Value property includes a space on either side of the "-" character.

► Task 3: Create a Parameter

1. In SQL Server Data Tools, in **Solution Explorer**, double-click **Project.params**.
2. In the **Project.params [Design]** window, click **Add Parameter**, and add a parameter with the following properties:
 - **Name:** AccountsFolderPath
 - **Data type:** String
 - **Value:** D:\Accounts\

- **Sensitive:** False
- **Required:** True
- **Description:** Path to accounts files



Note: Be sure to include the trailing “\” in the Value property.

3. On the **File** menu, click **Save All**, and then close the Project.params [Design] window.

► **Task 4: Use a Variable and a Parameter in an Expression**

1. On the **Data Flow** design surface for the **Extract Payment Data.dtsx** package, in the Connection Managers pane, click **Payments File**. Press F4 to view the Properties pane.
2. In the Properties pane, in the **Expressions** property box, click the ellipsis (...) button.
3. In the **Property Expressions Editor** dialog box, in the **Property** box, click **ConnectionString** and in the **Expression** box, click the ellipsis (...) button.
4. In the **Expression Builder** dialog box, expand the **Variables and Parameters** folder, and drag the **\$Project::AccountsFolderPath** parameter to the **Expression** box.
5. In the **Expression** box, type a plus (+) symbol after the **\$Project::AccountsFolderPath** parameter.
6. Drag the **User::fName** variable to the **Expression** box to create the following expression:


```
@[ $Project::AccountsFolderPath ]+ @[ User::fName ]
```
7. In the **Expression Builder** dialog box, click **Evaluate Expression** and verify that the expression produces the result **D:\Accounts\Payments - US.csv**.
8. Click **OK** to close the **Expression Builder** dialog box, and in the **Property Expressions Editor** dialog box, click **OK**.
9. On the **Debug** menu, click **Start Debugging** and observe the data flow while the package runs. When the package has completed, on the **Debug** menu, click **Stop Debugging**.
10. On the **Output** tab, find the following line in the log, noting that the default values for the **fName** variable and **AccountsFolderPath** parameter were used:

```
[Payments File [2]: The processing of the file "D:\Accounts\Payments - US.csv" has started
```

11. Close the SQL Server Data Tools, saving the changes if you are prompted.

Results: After this exercise, you should have a package that loads data from a text file based on a parameter that specifies the folder path where the file is stored, and a variable that specifies the file name.

Exercise 3: Using Containers

► Task 1: Add a Sequence Container to a Control Flow

1. Start **SQL Server Data Tools** and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab07A\Starter\Ex3** folder.
2. In Solution Explorer, in the **SSIS Packages** folder, double-click **Extract Internet Sales Data.dtsx**.
3. Review the control flow for the **Extract Internet Sales Data** package. You created the control flow for this package in Exercise 1 of this lab.
4. Right-click the red dotted arrow connecting the **Extract Customers** task to the **Send Failure Notification** task, and click **Delete**. Repeat this step to delete the red dotted line connecting the **Extract Internet Sales** task to the **Send Failure Notification** task, and the green arrow connecting the **Extract Internet Sales** task to the **Send Success notification** task.
5. Drag a Sequence container from the **Containers** section of the SSIS Toolbox to the Control Flow surface.
6. Right-click the new Sequence container, click **Rename**, and change the container name to **Extract Customer Sales Data**.
7. Click the **Extract Customers** task, hold the CTRL key, and click the **Extract Internet Sales** task, then drag both tasks into the **Extract Customer Sales Data** sequence container.
8. Click the **Extract Customer Sales Data** sequence container, and then drag the green arrow from the **Extract Customer Sales Data** sequence container to the **Send Success Notification** task.
9. Click the **Extract Customer Sales Data** sequence container, and then drag the green arrow from the **Extract Customer Sales Data** sequence container to the **Send Failure Notification** task.
10. Right-click the green arrow connecting the **Extract Customer Sales Data** sequence container to the **Send Failure Notification** task, and click **Failure**.
11. On the **Debug** menu, click **Start Debugging**, and observe the control flow as the package executes.
12. When package execution is complete, on the **Debug** menu, click **Stop Debugging**.

► Task 2: Add a Foreach Loop Container to a Control Flow

1. In **SQL Server Data Tools**, in **Solution Explorer**, in the **SSIS Packages** folder, double-click **Extract Payment Data.dtsx**.
2. Review the control flow for the **Extract Payment Data** package, and note that it contains a single Data Flow task named **Extract Payments**. This is the same Data Flow task you updated in the previous exercise.
3. In the SSIS Toolbox, in the **Containers** section, double-click **Foreach Loop Container**.
4. On the Control Flow surface, click the **Extract Payments** task and drag it into the **Foreach Loop Container**.
5. Double-click the title area at the top of the **Foreach Loop Container** to view the **Foreach Loop Editor** dialog box.
6. In the **Foreach Loop Editor** dialog box, on the **Collection** tab, in the **Enumerator** list, click **Foreach File Enumerator**.
7. In the **Expressions** box, click the ellipsis (...) button.
8. In the **Property Expressions Editor** dialog box, in the **Property** list, select **Directory** and in the **Expression** box click the ellipsis (...) button.

9. In the **Expression Builder** dialog box, expand the **Variables and Parameters** folder and drag the **\$Project::AccountsFolderPath** parameter to the **Expression** box.
10. Click **OK** to close the **Expression Builder**, and click **OK** again to close the **Property Expression Editor**.
11. In the **Foreach Loop Editor** dialog box, on the **Collection** tab, in the **Retrieve file name** section, select **Name and extension**.
12. In the **Foreach Loop Editor** dialog box, on the **Variable Mappings** tab, in the **Variable** list, select **User::fName** and in the **Index column** ensure that **0** is specified. Click **OK**.
13. On the **Debug** menu, click **Start Debugging** and observe the data flow while the package runs.
14. When the package has completed, on the **Debug** menu, click **Stop Debugging**.
15. On the **Execution Results** tab, scroll through the log, noting that the data flow was executed once for each file in the **D:\Lab07A\Solution\Ex2\Accounts** folder. The following files should have been processed:
 - Payments – AU.csv
 - Payments – CA.csv
 - Payments – DE.csv
 - Payments – FR.csv
 - Payments – GB.csv
 - Payments – US.csv
16. Close SQL Server Data Tools, saving your changes if prompted.

Results: After this exercise, you should have one package that encapsulates two Data Flow tasks in a Sequence container, and another that uses a Foreach Loop to iterate through the files in a folder specified in a parameter—and uses a Data Flow task to load their contents into a database.

Lab B: Using Transactions and Checkpoints

Exercise 1: Using Transactions

► Task 1: Prepare the Lab Environment

1. Ensure the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab07B\Starter** folder, right-click **Setup.cmd** and click **Run as administrator**.
3. Click **Yes** when prompted to confirm you want to run the command file. Press **y**, when prompted, and then press **Enter** and wait for the script to finish.

► Task 2: View the Data in the Database

1. Start **SQL Server Management Studio**, and connect to the **MIA-SQL** database engine using Windows authentication.
2. In Object Explorer, expand **Databases**, expand **Staging**, and then expand **Tables**.
3. Right-click **dbo.Customers** and click **Select Top 1000 Rows**. Note that the table is empty.
4. Right-click **dbo.InternetSales** and click **Select Top 1000 Rows**. Note that the table is also empty.
5. Minimize **SQL Server Management Studio**.

► Task 3: Run a Package to Extract Data

1. Start **SQL Server Data Tools** and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab07B\Starter\Ex1** folder.
2. In Solution Explorer, in the **SSIS Packages** folder, double-click **Extract Internet Sales Data.dtsx**.
3. Review the control flow for the **Extract Internet Sales Data** package.
4. On the **Debug** menu, click **Start Debugging**, and observe the control flow as the package is executed, noting that the **Extract Customers** task succeeds, but the **Extract Internet Sales** task fails.
5. When package execution has completed, on the **Debug** menu, click **Stop Debugging**.
6. Return to **SQL Server Management Studio** and re-execute the queries you created earlier to view the top 1,000 rows in the **dbo.Customers** and **dbo.InternetSales** tables. Verify that the **dbo.InternetSales** table is still empty but the **dbo.Customers** table now contains customer records.
7. In **SQL Server Management Studio**, click **New Query**.
8. In the new query window, enter the following Transact-SQL and click **Execute**:

```
TRUNCATE TABLE Staging.dbo.Customers;
```

9. Close the query tab containing the **TRUNCATE TABLE** statement without saving it, and minimize **SQL Server Management Studio**.

► Task 4: Implement a Transaction

1. In **SQL Server Data Tools**, on the Control Flow surface for the **Extract Internet Sales Data.dtsx** package, click the **Extract Customer Sales Data** sequence container and press **F4** to view the Properties pane.
2. In the Properties pane, set the **TransactionOption** property of the **Extract Customer Sales Data** sequence container to **Required**.

3. Click the **Extract Customers** task, and in the Properties pane, ensure that the **TransactionOption** property value is set to **Supported**, and set the **FailParentOnFailure** property to **True**.
4. Repeat the previous step for the **Extract Internet Sales** task.

► **Task 5: Observe Transaction Behavior**

1. In SQL Server Data Tools, on the **Debug** menu, click **Start Debugging**, and observe the control flow as the package is executed, noting that the **Extract Customers** task succeeds again, but the **Extract Internet Sales** task fails.
2. When the package execution has completed, on the **Debug** menu, click **Stop Debugging**.
3. Return to **SQL Server Management Studio** and re-execute the queries you created earlier to view the top 1,000 rows in the **dbo.Customers** and **dbo.InternetSales** tables, and verify that both tables are empty.
4. In **SQL Server Data Tools**, on the Control Flow surface for the **Extract Internet Sales Data.dtsx** package, double-click the **Extract Internet Sales** task to view it in the Data Flow tab.
5. Double-click the **Calculate Sales Amount** transformation and modify the **Expression** value to resemble that shown below, and then, click **OK**:

```
UnitPrice * OrderQuantity
```

6. Click the **Control Flow** tab, and on the **Debug** menu, click **Start Debugging**. Observe the control flow as the package is executed, noting that both the **Extract Customers** and **Extract Internet Sales** tasks succeed.
7. When package execution has completed, on the **Debug** menu, click **Stop Debugging**.
8. Return to **SQL Server Management Studio** and re-execute the queries you created earlier to view the top 1,000 rows in the **dbo.Customers** and **dbo.InternetSales** tables, and verify that both tables now contain data.
9. Leave **SQL Server Management Studio** open as you will use it again in the next exercise.
10. Close SQL Server Data Tools, saving changes if prompted.

Results: After this exercise, you should have a package that uses a transaction to ensure that all Data Flow tasks succeed or fail as an atomic unit of work.

Exercise 2: Using Checkpoints

► **Task 1: View the Data in the Database**

1. Return to **SQL Server Management Studio**, and in Object Explorer, ensure that **Databases**, **Staging**, and **Tables** are expanded for the **MIA-SQL** database engine instance.
2. Right-click **dbo.Resellers** and click **Select Top 1000 Rows**. Note that the table is empty.
3. Right-click **dbo.ResellerSales** and click **Select Top 1000 Rows**. Note that this table is also empty.

► Task 2: Run a Package to Extract Data

1. Start **SQL Server Data Tools** and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab07B\Starter\Ex2** folder.
2. In Solution Explorer, in the **SSIS Packages** folder, double-click **Extract Reseller Data.dtsx**, and then examine the control flow.
3. On the **Debug** menu, click **Start Debugging**, and observe the control flow as the package is executed, noting that the **Extract Resellers** task succeeds, but the **Extract Reseller Sales** task fails.
4. When package execution has completed, on the **Debug** menu, click **Stop Debugging**.
5. Return to **SQL Server Management Studio** and re-execute the queries you created earlier to view the top 1,000 rows in the **dbo.Resellers** and **dbo.ResellerSales** tables. Verify that the **dbo.ResellerSales** table is still empty but the **dbo.Resellers** table now contains records.
6. In **SQL Server Management Studio**, click **New Query**. In the new query window, enter the following Transact-SQL code and click **Execute**:

```
TRUNCATE TABLE Staging.dbo.Resellers;
```

7. Close the query tab containing the TRUNCATE TABLE statement without saving it.

► Task 3: Implement Checkpoints

1. In **SQL Server Data Tools**, click any empty area on the Control Flow surface for the **Extract Reseller Data.dtsx** package, and press F4 to view the Properties pane.
2. In the Properties pane, set the following properties of the **Extract Reseller Data** package:
 - **CheckpointFileName**: D:\ETL\CheckPoint.chk
 - **CheckpointUsage**: IfExists
 - **SaveCheckpoints**: True
3. Click the **Extract Resellers** task, and in the Properties pane, set the **FailPackageOnFailure** property to **True**.
4. Repeat the previous step for the **Extract Reseller Sales** task.

► Task 4: Observe Checkpoint Behavior

1. View the contents of the **D:\ETL** folder and verify that no file named **CheckPoint.chk** exists.
2. In **SQL Server Data Tools**, on the **Debug** menu, click **Start Debugging**. Observe the control flow as the **Extract Reseller Data.dtsx** package is executed, noting the **Extract Resellers** task succeeds again, but the **Extract Reseller Sales** task fails.
3. When package execution is complete, on the **Debug** menu, click **Stop Debugging**.
4. View the contents of the **D:\ETL** folder and verify that a file named **Checkpoint.chk** has been created.
5. Return to **SQL Server Management Studio** and re-execute the queries you created earlier to view the top 1,000 rows in the **dbo.Resellers** and **dbo.ResellerSales** tables; verify that the **dbo.ResellerSales** table is still empty, but the **dbo.Resellers** table now contains reseller records.
6. In SQL Server Data Tools, on the Control Flow surface for the **Extract Reseller Data.dtsx** package, double-click the **Extract Reseller Sales** task to view it in the Data Flow tab.

7. Double-click the **Calculate Sales Amount** transformation and modify the **Expression** value so that it matches the following code, then click **OK**:

```
UnitPrice * OrderQuantity
```

8. Click the **Control Flow** tab, and on the **Debug** menu, click **Start Debugging**. Observe the control flow as the package is executed, noting that the **Extract Resellers** task is not re-executed, and package execution starts with the **Extract Reseller Sales** task, which failed on the last attempt.
9. When package execution has completed, on the **Debug** menu, click **Stop Debugging**.
10. View the contents of the **D:\ETL** folder and verify that the **Checkpoint.chk** file has been deleted.
11. Return to **SQL Server Management Studio** and re-execute the queries you created earlier to view the top 1,000 rows in the **dbo.Resellers** and **dbo.ResellerSales** tables. Verify that both now contain data.
12. Close **SQL Server Management Studio** without saving changes.
13. Close SQL Server Data Tools, saving changes if prompted.

Results: After this exercise, you should have a package that uses checkpoints, so that execution can be restarted at the point of failure of the previous execution.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 8: Debugging and Troubleshooting SSIS Packages

Lab: Debugging and Troubleshooting an SSIS Package

Exercise 1: Debugging an SSIS Package

► Task 1: Prepare the Lab Environment

1. Ensure the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab08\Starter** folder, right-click **Setup.cmd** and click **Run as administrator**.
3. Click **Yes** when prompted to confirm you want to run the command file, and wait for the script to finish.

► Task 2: Run an SSIS Package

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex1** folder.
2. In Solution Explorer, double-click the **Extract Payment Data.dtsx** SSIS package to open it in the designer. Note that this package includes a control flow that iterates through files in a folder, and loads payments data from each file into a staging database.
3. On the **Debug** menu, click **Start Debugging** and observe the package as it executes, noting that it fails. When execution has completed, on the **Debug** menu, click **Stop Debugging**.

► Task 3: Add a Breakpoint

1. On the control flow surface for the **Extract Payment Data.dtsx** package, right-click the **Foreach Loop Container**, and then click **Edit Breakpoints**.
2. In the **Set Breakpoints – Foreach Loop Container** dialog box, select **Enabled** for the **Break at the beginning of every iteration of the loop** condition, and then click **OK**.

► Task 4: Add a Data Viewer

1. On the control flow surface for the **Extract Payment Data.dtsx** package, double-click the **Extract Payments** data flow task to view the data flow design surface.
2. Double-click the data flow path between the **Payments File** source and the **Staging DB** destination to open its editor.
3. In the **Data Flow Path Editor** dialog box, on the **Data Viewer** tab, select **Enable data viewer** and note that all available columns are included in the **Displayed columns** list. Click **OK**.

► Task 5: View Breakpoints

1. Click the **Control Flow** tab to view the control flow design surface for the **Extract Payment Data.dtsx** package.
2. On the **Debug** menu, point to **Windows**, and then click **Breakpoints**.
3. In the Breakpoints pane, view the breakpoints. Note that they include the one you added to the Foreach Loop container and the data viewer you defined in the data flow.

► Task 6: Observe Variables while Debugging

1. On the **Debug** menu, click **Start Debugging** and note that execution stops at the first breakpoint.
2. On the **Debug** menu, point to **Windows**, and then click **Locals**.
3. In the Locals pane, expand **Variables**, find the **User::fName** variable, and note its value (which should be **Payments – AU.csv**). Then right-click the **User::fName** variable and click **Add Watch**. Note that the **User::fName** variable has been added to the Watch 1 pane, which is now shown.
4. Click the **Locals** tab to view the Locals pane again and add a watch for the **\$Project::AccountsFolderPath** parameter.
5. Ensure that the Watch 1 pane is visible and that you can see the values for the **User::fName** variable and the **\$Project::AccountsFolderPath** parameter. Then, on the **Debug** menu, click **Continue**. Note that execution stops at the next breakpoint, which is the data viewer you added to the data flow.
6. Review the data in the data viewer pane, noting that this represents the contents of the Payments – AU.csv file. Drag the data viewer pane and position it so you can see all the following user interface elements:
 - The Foreach Loop Container on the control flow surface.
 - The variables in the Watch 1 pane.
 - The data viewer pane.
7. In the data viewer pane, click **Continue** (the green arrow). Note that execution continues until the next breakpoint and that the **Extract Payments** task completed successfully for the first loop iteration. In the Watch 1 window, note that the value of the **User::fName** variable has changed to **Payments – CA.csv**, but that the contents of the data viewer pane still reflect the **Payments – AU.csv** file, because the data flow for the current loop iteration has not yet started.
8. On the **Debug** menu, click **Continue**, and note that execution continues to the next breakpoint. The data viewer pane now shows the contents of the **Payments – CA.csv** file.
9. In the data viewer pane, click the **Continue** button. Note that execution continues until the next breakpoint and that the **Extract Payments** task completed successfully for the second loop iteration. In the Watch 1 window, note that the value of the **User::fName** variable has changed to **Payments – DE.csv**, indicating that it is the next file to be processed by the **Extract Payments** task.
10. On the **Debug** menu, click **Continue** and note that execution continues to the next breakpoint. The data viewer pane now shows the contents of the **Payments – DE.csv** file.
11. In the data viewer pane, click the **Continue** button. Note that execution continues until the next breakpoint and that the **Extract Payments** task completed successfully for the **Payments – DE.csv** file. In the Watch 1 window, note that the value of the **User::fName** variable has changed to **Payments – EU.csv**, indicating that it is the next file to be processed by the **Extract Payments** task.
12. On the **Debug** menu, click **Continue** and note that execution continues to the next breakpoint. The data viewer pane now shows the contents of the **Payments – EU.csv** file (which contains a mixture of country codes including DE, FR, and GB).
13. In the data viewer pane, click the **Continue** button. Note that execution continues until the next breakpoint and that the **Extract Payments** task fails for the **Payments – FR.csv** file.
14. In the data viewer pane, note that the contents of the **Payments – EU.csv** file are still shown. Click **Copy Data** to copy the data to the clipboard.
15. On the **Debug** menu, click **Stop Debugging**.

16. Close Visual Studio.

► **Task 7: View Data Copied from a Data Viewer**

1. Start Notepad.
2. On the **Edit** menu, click **Paste**.
3. Examine the data to see why the SSIS package failed to load from the **Payments – EU.csv** file. Note that the data contains the following issues:
 - Payment 4074 has an invalid payment date ("Last Tuesday").
 - Payment 4102 has an invalid payment amount (300 USD).
 - Payment 4124 has an invalid payment date (13/13/2004).
4. Close Notepad without saving any changes.

Results: After this exercise, you should have observed the variable values and data flows for each iteration of the loop in the Extract Payment Data.dtsx package. You should also have identified the file that caused the data flow to fail and examined its contents to find the data errors that triggered the failure.

Exercise 2: Logging SSIS Package Execution

► **Task 1: Configure SSIS Logs**

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex2** folder.
2. In Solution Explorer, double-click the **Extract Payment Data.dtsx** SSIS package to open it in the designer.
3. On the **SSIS** menu, click **Logging**.
4. In the **Configure SSIS Logs: Extract Payment Data** dialog box, in the **Containers** tree, select **Extract Payment Data**.
5. On the **Providers and Logs** tab, in the **Provider type** drop-down list, select **SSIS log provider for Windows Event Log**, and then click **Add**.
6. In the **Select the logs to use for the container** list, select **SSIS log provider for Windows Event Log**.
7. On the **Details** tab, select the **OnError** and **OnTaskFailed** events.
8. Click **OK**.

► **Task 2: View Logged Events**

1. On the **Debug** menu, click **Start Debugging**, and observe the execution of the package, noting that it fails. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
2. On the **SSIS** menu, click **Log Events** and review the events that were logged during the package execution (if there are none, rerun the package). Close the **Log Events** tab and close Visual Studio.
3. On the Windows taskbar, click **Search Windows**, type **Event Viewer**, and then press Enter.
4. In Event Viewer, expand **Windows Logs** and click **Application**.

5. In Event Viewer, click each event with the source **SQLISPackage140** and review the event information in the pane at the bottom of the Event Viewer window.
6. Close Event Viewer.

Results: After this exercise, you should have a package that logs event details to the Windows Event Log.

Exercise 3: Implementing an Event Handler

► Task 1: Create an Event Handler

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex3** folder.
2. In Solution Explorer, double-click the **Extract Payment Data.dtsx** SSIS package to open it in the designer.
3. Click the **Event Handlers** tab.
4. In the **Executable** list, ensure **Extract Payment Data** is selected and, in the **Event handler** list, ensure **OnError** is selected.
5. Click the **Click here to create an 'OnError' event handler for executable 'Extract Payment Data'** link.

► Task 2: Add a File System Task

1. In the SSIS Toolbox, double-click **File System Task** to add a file system task to the event handler design surface for the OnError event handler.
2. On the design surface, right-click the **File System Task**, click **Rename**, and change the name to **Copy Failed File**.
3. Double-click the **Copy Failed File** data flow task to view its task editor.
4. On the **General** tab, ensure the **Copy file** operation is selected.
5. In the **SourceConnection** drop-down list, select **Payments File**.
6. In the **DestinationConnection** drop-down list, select **<New connection>**.
7. In the **File Connection Manager Editor** dialog box, in the **Usage type** drop-down list, select **Create file**.
8. In the **File** box, type **D:\ETL\FailedPayments.csv**, and then click **OK**.
9. In the **File System Task Editor** dialog box, in the **OverwriteDestination** drop-down list, select **True**, and then click **OK**.
10. In the **Connection Managers** area at the bottom of the design surface, click **FailedPayments.csv** and press F4.
11. In the Properties pane, in the **Expressions** property value, click the ellipsis (...).
12. In the **Property Expressions Editor** dialog box, in the **Property** drop-down list, select **ConnectionString**, and, in the **Expression** box, click the ellipsis (...) to open the **Expression Builder**.

- In the **Expression Builder** dialog box, in the **Expression** text area, enter the following expression, and then click **OK**:

```
"D:\\ETL\\" + @[System::ExecutionInstanceGUID] + @[User::fName]
```

- In the **Property Expressions Editor** dialog box, click **OK**.

► Task 3: Add a Send Mail Task

- In the SSIS Toolbox, double-click the **Send Mail** task to add a send mail task to the event handler design surface.
- On the design surface, right-click the **Send Mail** task, click **Rename**, and change the name to **Send Notification**.
- Move the **Send Notification** task below the **Copy Failed File** task, click the **Copy Failed File** task and drag the green precedence connection to the **Send Notification** task.
- Right-click the **precedence** connection and click **Completion**.
- Double-click the **Send Notification** task to open its task editor.
- In the **Send Mail Task Editor** dialog box, on the **Mail** tab, set the following properties:
 - **SmtplibConnection:** Local SMTP Server
 - **From:** etl@adventureworks.msft
 - **To:** student@adventureworks.msft
 - **Subject:** An error occurred
 - **Priority:** High
- In the **Send Mail Task Editor** dialog box, on the **Expressions** tab, click the **Expressions** property, and then click the ellipsis (...).
- In the **Property Expressions Editor** dialog box, in the **Property** drop-down list, select **MessageSource** and, in the **Expression** box, click the ellipsis (...) to open the **Expression Builder**.
- In the **Expression Builder** dialog box, in the **Expression** text area, enter the following expression, and then click **OK**:

```
@[User::fName] + " failed to load. " + @[System::ErrorDescription]
```

- In the **Property Expressions Editor** dialog box, click **OK**.
- In the **Send Mail Task Editor**, click **OK**.

► Task 4: Test the Event Handler

- Click the **Control Flow** tab.
- On the **Debug** menu, click **Start Debugging**, and observe the execution of the package, noting that the package fails.
- When execution is complete, click the **Event Handlers** tab to verify that the event handler has been executed then, on the **Debug** menu, click **Stop Debugging**.
- Close Visual Studio. Save your changes if you are prompted.

5. View the contents of the **D:\ETL** folder and note that a file with a name similar to {1234ABCD-1234-ABCD-1234-ABCD1234}Payments - EU.csv has been created. Open this file in Notepad, review the contents looking for any anomalous data (it should contain the same issues noted earlier), and then close Notepad without saving any changes.
6. View the contents of the C:\inetpub\mailroot\Drop folder, and note that several email messages were sent at the same time.
7. Open the most recent email file by using Notepad, examine it; it should contain the subject line "An error occurred".
8. Close Notepad.

Results: After this exercise, you should have a package that includes an event handler for the OnError event. The event handler should create a copy of files that contain invalid data and send an email message.

Exercise 4: Handling Errors in a Data Flow

► Task 1: Redirect Data Flow Errors

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab08\Starter\Ex4** folder.
2. In Solution Explorer, double-click the **Extract Payment Data.dtsx** SSIS package to open it in the designer.
3. On the control flow surface, double-click the **Extract Payments** task to view its data flow.
4. Double-click the **Staging DB** destination to view its editor.
5. On the **Error Output** tab, in the **Error** column for **OLE DB Destination Input**, select **Redirect row**, and then click **OK**.
6. In the SSIS Toolbox, in the **Other Destinations** section, double-click **Flat File Destination**.
7. On the design surface, right-click the **Flat File Destination** task, click **Rename**, and change the name to **Invalid Rows**.
8. Drag **Invalid Rows** to the right of **Staging DB**, and then connect the red data flow path from **Staging DB** to **Invalid Rows**.
9. In the **Configure Error Output** dialog box, ensure that **Redirect row** is selected in the **Error** column, and then click **OK**.
10. Double-click the **Invalid Rows** task.
11. In the **Flat File Destination Editor** dialog box, click **New**.
12. In the **Flat File Format** dialog box, ensure **Delimited** is selected, and then click **OK**.
13. In the **Flat File Connection Manager Editor** dialog box, change the following properties, and then click **OK**:
 - **Connection manager name:** Invalid Payment Records
 - **File name:** D:\ETL\InvalidPaymentsLog.csv

14. In the **Flat File Destination Editor** dialog box, on the **Connection Manager** tab, clear the **Overwrite data in the file** check box.
15. On the **Mappings** tab, note that the input columns include the columns from the payments file and two additional columns for the error code and error column, and then click **OK**.

► **Task 2: View Invalid Data Flow Rows**

1. Click the **Control Flow** tab.
2. On the **Debug** menu, click **Start Debugging**, and observe the execution of the package, noting that it succeeds.
3. When execution is complete, on the **Debug** menu, click **Stop Debugging** and close Visual Studio, saving any changes if prompted.
4. View the contents of the **D:\ETL** folder and note that a file named **InvalidPaymentsLog.csv** has been created.
5. Start Notepad and open **InvalidPaymentsLog.csv** to view its contents. The file should contain the three rows that contain invalid data.
6. Close Notepad without saving any changes.

Results: After this exercise, you should have a package that includes a data flow where rows containing errors are redirected to a text file.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 9: Implementing a Data Extraction Solution

Lab A: Extracting Modified Data

Exercise 1: Using a Datetime Column to Incrementally Extract Data

► Task 1: Prepare the Lab Environment

1. Ensure that the 20767C-MIA-DC and 20767C-MIA-SQL virtual machines are both running, and then log on to 20767C-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab09\Starter** folder, right-click **Setup.cmd** and click **Run as administrator**.
3. Click **Yes** when prompted to confirm you want to run the command file, and wait for the script to finish.

► Task 2: View Extraction Data

1. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the database engine by using Windows authentication.
2. In Object Explorer, expand **Databases**, expand **Staging**, and then expand **Tables**.
3. Right-click **dbo.ExtractLog**, and then click **Select Top 1000 Rows**.
4. Review the data in the **ExtractLog** table, noting the values in the **LastExtract** column, which indicate the date and time of the last extract operations for the **InternetSales** and **ResellerSales** databases. This is initially set to the date and time "1900-01-01 00:00:00.000".
5. In Object Explorer, under **Databases**, expand **InternetSales** and then expand **Tables**.
6. Right-click **dbo.SalesOrderHeader** and then click **Select Top 1000 Rows**. Note that the **LastModified** column indicates the date and time that the sales record was last modified.
7. Minimize SQL Server Management Studio.

► Task 3: Examine an Incremental Data Extraction

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex1** folder. If a Security warning dialog box appears click **OK**.
2. In Solution Explorer, under **SSIS Packages**, double-click **Extract Reseller Data.dtsx**.
3. On the **SSIS** menu, click **Variables**.
4. In the **Variables** window, note that the following variables have been defined with a data type of **DateTime**, and then close the window:
 - CurrentTime
 - ResellerSalesLastExtract
5. On the **Control Flow** tab of the design surface, double-click the **Get Current Time** task, and note that it uses the GETDATE() function to assign the current data and time to the **CurrentTime** variable.
6. Click **Cancel** to close the **Expression Builder** dialog box.
7. Double-click the **Get Last Extract Time** task to open the **Execute SQL Task Editor** dialog box.

8. Note the following configuration settings, and then click **Cancel**:
 - o On the **General** tab, the task is configured to return a single row from the **Staging** database by executing the following Transact-SQL statement:

```
SELECT MAX>LastExtract) LastExtract
FROM ExtractLog
WHERE DataSource = 'ResellerSales'
```

- o On the **Result Set** tab, the **LastExtract** column in the result returned by the query is assigned to the **User::ResellerSalesLastExtract** variable.
9. Double-click the **Extract Reseller Sales** task to view its data flow.
 10. Double-click the **Reseller Sales** source to view its settings.
 11. In the **OLE DB Source Editor** dialog box, note that the SQL command text used to retrieve the reseller sales data includes the following parameterized WHERE clause:

```
WHERE LastModified > ?
AND LastModified <= ?
```

12. Click **Parameters**, and note that the parameters in the WHERE clause are mapped to the **User::ResellerSalesLastExtract** and **User::CurrentTime** variables.
13. Click **Cancel** to close the **Set Query Parameters** dialog box, and click **Cancel** again to close the **OLE DB Source Editor** dialog box.
14. Click the **Control Flow** tab, and then double-click the **Update Last Extract Time** task.
15. Note the following configuration settings, and then click **Cancel**:
 - o On the **General** tab, the task is configured to update the **Staging** database by executing the following parameterized Transact-SQL statement:

```
UPDATE ExtractLog
SET LastExtract = ?
WHERE DataSource = 'ResellerSales'
```

- o On the **Parameter Mapping** tab, the parameter in the query is mapped to the **User::CurrentTime** variable.
16. On the **Debug** menu, click **Start Debugging**, and observe the control flow as the package runs.
 17. When package execution is complete, on the **Debug** menu, click **Stop Debugging**.
 18. Return to SQL Server Management Studio.
 19. In Object Explorer, in the **Tables** folder for the **Staging** database, right-click **dbo.ExtractLog**, and then click **Select Top 1000 Rows**. Note that the last extract time for the **ResellerSales** data source has been updated.
 20. Leave SQL Server Management Studio open.

► Task 4: Define Variables for Execution Times

1. In Visual Studio, in Solution Explorer, double-click **Extract Internet Sales Data.dtsx** within the **AdventureWorksETL.sln** project that you loaded from the **D:\Labfiles\Lab09\Starter\Ex1** folder.
2. On the **SSIS** menu, click **Variables**.

3. In the Variables window, click **Add Variable** and add a variable with the following settings:
 - **Name:** CurrentTime
 - **Data type:** DateTime
4. Click **Add Variable** again, add a variable with the following settings, and then close the Variables window:
 - **Name:** InternetSalesLastExtract
 - **Data type:** DateTime
5. On the **SSIS** menu, click **SSIS Toolbox**.
6. In the **SSIS** Toolbox, drag an **Expression Task** inside the **Extract Customer Sales Data** sequence container on the control surface. Arrange the tasks in this container so that the new expression task is above the **Extract Customers** task.
7. Right-click the expression task, click **Rename**, and change the name to **Get Current Time**.
8. Double-click the **Get Current Time** task.
9. In the **Expression Builder** dialog box, in the **Expression** box, type the following expression, and then click **OK**:

```
@[User::CurrentTime] = GETDATE()
```

10. In the **SSIS** Toolbox, drag an **Execute SQL** task to the control flow surface and drop it inside the **Extract Customer Sales Data** sequence container, immediately below the **Get Current Time** task.
11. Right-click the **Execute SQL** task, click **Rename**, and change the name to **Get Last Extract Time**.
12. Double-click the **Get Last Extract Time** task.
13. In the **Execute SQL Task Editor** dialog box, configure the following settings, and then click **OK**:
 - On the **General** tab, in the **Connection** drop-down list, select **localhost.Staging**.
 - On the **General** tab, in the **SQLStatement** box, click the ellipsis (...) button. In the **Enter SQL Query** dialog box, type the following Transact-SQL query, and then click **OK**:

```
SELECT MAX>LastExtract) LastExtract
FROM ExtractLog
WHERE DataSource = 'InternetSales'
```

- On the **General** tab, in the **ResultSet** drop-down list, select **Single** row.
 - On the **Result Set** tab, click **Add**, and then in the **Result Name** column, change **NewResultName** to **LastExtract**, and in the **Variable Name** drop-down list, select **User::InternetSalesLastExtract**.
14. On the control flow surface, click the **Get Current Time** task and drag its green precedence constraint to the **Get Last Extract Time** task.
 15. Click the **Get Last Extract Time** task and drag its green precedence constraint to the **Extract Customers** task.
- **Task 5: Modify a Data Source to Filter Data**
1. On the control flow surface, double-click the **Extract Internet Sales** task to display its data flow.
 2. On the data flow surface, double-click the **Internet Sales** source.

3. In the **OLE DB Source Editor** dialog box, review the existing SQL command text used to extract sales data.
4. Add the following parameterized WHERE clause to the SQL command text:

```
WHERE LastModified > ?
AND LastModified <= ?
```

5. Click **Parameters**.
6. In the **Set Query Parameters** dialog box, specify the following parameter mappings with a **Param direction** of **Input**, and then click **OK**:
 - **Parameter0**: User::InternetSalesLastExtract
 - **Parameter1**: User::CurrentTime
7. In the **OLE DB Source Editor** dialog box, click **OK**.

► Task 6: Add a Task to Update the Extraction Log

1. Click the **Control Flow** tab, and then in the SSIS Toolbox, drag an **Execute SQL Task** to the **Extract Customer Sales Data** sequence under the **Extract Internet Sales** task, on the control flow surface.
2. Right-click **Execute SQL Task** and click **Rename**. Change the name to **Update Last Extract Time**.
3. Double-click **Update Last Extract Time** and configure the following settings, and then click **OK**:
 - On the **General** tab, in the **Connection** drop-down list, select **localhost.Staging**.
 - On the **General** tab, in the **SQLStatement** box, click the ellipsis (...) button.
 - In the **Enter SQL Query** dialog box, enter the following Transact-SQL query, and then click **OK**:

```
UPDATE ExtractLog
SET LastExtract = ?
WHERE DataSource = 'InternetSales'
```

- On the **Parameter Mapping** tab, click **Add** and create the following parameter mapping:
 - **Variable Name**: User::CurrentTime
 - **Direction**: Input
 - **Data Type**: DATE
 - **Parameter Name**: 0
 - **Parameter Size**: -1
- 4. On the control flow surface, click the **Extract Internet Sales** task, and then drag its green precedence constraint to the **Update Last Extract Time** task.

► Task 7: Test the Package

1. Click the **Data Flow** tab and view the **Extract Internet Sales** data flow.
2. On the **Debug** menu, click **Start Debugging** and observe the package as it executes, noting the number of rows transferred.
3. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
4. Return to SQL Server Management Studio.
5. In Object Explorer, in the **Staging** database, right-click the **dbo.ExtractLog** table, and then click **Select Top 1000 Rows**.

6. View the data in the **ExtractLog** table, noting the value in the **LastExtract** column for the **InternetSales** database has been updated to the date and time when you ran the package.
7. Right-click the **dbo.InternetSales** table, and then click **Select Top 1000 Rows**. The sales records in this table were extracted from the **InternetSales** database, where the **SalesOrderHeader** table had a **LastModified** column value between the previous **LastExtract** value, and the date and time when the package was executed.
8. Minimize SQL Server Management Studio.
9. In Visual Studio, with the **Extract Internet Sales** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging** to execute the package again, noting that no rows are transferred during this execution.
10. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
11. Close Visual Studio.

Results: After this exercise, you should have an SSIS package that uses the high water mark technique to extract only records that have been modified since the previous extraction.

Exercise 2: Using Change Data Capture

► Task 1: Enable Change Data Capture

1. Return to SQL Server Management Studio, and open the **Enable CDC.sql** file in the **D:\Labfiles\Lab09\Starter\Ex2** folder.
2. Examine the Transact-SQL code in this script, noting that it enables CDC in the **InternetSales** database, and for the **Customers** table.
3. Click **Execute** to run the script. Two jobs should be started.
4. Open the **Test CDC.sql** file in the **D:\Labfiles\Lab09\Starter\Ex2** folder, and examine the query, noting that it contains statements to perform the following tasks:
 - Retrieve data changes between 1/1/1900 and the current date by using a CDC function.
 - Modify the data in the **Customers** table.
 - Retrieve data changes between 1/1/1900 and the current date again.
5. Select the code under the comment **Select all changed customer records between 1/1/1900 and today** and click **Execute**. Note that no records are returned because there have been no changes in the database since CDC was enabled.
6. Select the two UPDATE statements under the comment **Make a change to all customers (to create CDC records)** and click **Execute**. This statement modifies the data in the **Customers** table by reversing the **FirstName** value, and then reversing it back to its original value.
7. Select the code under the comment **Now see the net changes** and click **Execute**. Note that the query returns all records in the **Customers** table, because they have all been changed within the specified time period.

► Task 2: Create a Stored Procedure to Retrieve Modified Rows

1. In SQL Server Management Studio, open the **Create SP.sql** file in the **D:\Labfiles\Lab09\Starter\Ex2** folder.
2. Examine the Transact-SQL code in the query file, and note that it creates a stored procedure with **StartDate** and **EndDate** parameters. The stored procedure performs the following tasks:
 - Retrieves the log sequence numbers for the dates specified in the **StartDate** and **EndDate** parameters.
 - If neither of the log sequence numbers is null, then at least one transaction has occurred in the database within the specified time period. The stored procedure uses a CDC function to return all records that have changed in the **Customers** table.
 - If no transactions have taken place in the specified time period, the stored procedure returns an empty rowset.
3. Select the code under the comment **Drop existing stored procedure if it exists** and then click **Execute**.
4. Select the code under the comment **Create a stored procedure to return customer changes between two given dates**, and then click **Execute** to run the Transact-SQL code and create the stored procedure.
5. Click **New Query**, and type the following Transact-SQL in the new query window. Then click **Execute** to test the stored procedure:

```
USE InternetSales
GO
EXEC GetChangedCustomers '1/1/1900', '1/1/2099';
```

6. GO

► Task 3: Use the Stored Procedure in a Data Flow

1. In SQL Server Management Studio, open the **Reset Staging.sql** file in the **D:\Labfiles\Lab09\Starter\Ex2** folder.
2. Click **Execute** to reset the **Staging** database.
3. Minimize SQL Server Management Studio.
4. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex2** folder. If a Security warning dialogue box appears click **OK**.
5. In Solution Explorer, under **SSIS Packages**, double-click **Extract Internet Sales Data.dtsx**.
6. On the control flow surface, double-click the **Extract Customers** task.
7. On the data flow surface, double-click the **Customers** source.
8. In the **OLE DB Source Editor** dialog box, make the following changes to the configuration of the **Customers** source, and then click **OK**:
 - In the **Data access mode** drop-down list, select **SQL Command**.
 - In the **SQL command** text box, type the following Transact-SQL statement:

```
EXEC GetChangedCustomers ?, ?
```

- Click **Parameters**, and in the **Set Query Parameters** dialog box, create the following parameter mappings with a **Param direction** of **Input**, and then click **OK**.
 - **@StartDate:**
User::InternetSalesLastExtract
 - **@EndDate:**
User::CurrentTime

► Task 4: Test the Package

1. With the **Extract Customers** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging** and observe the package as it executes, noting the number of rows transferred.
2. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
3. Return to SQL Server Management Studio.
4. In Object Explorer, in the **Staging** database, right-click the **dbo.ExtractLog** table, and then click **Select Top 1000 Rows**.
5. View the data in the **ExtractLog** table, noting the value in the **LastExtract** column for the **InternetSales** database has been updated to the date and time when you ran the package.
6. Right-click the **dbo.Customers** table, and then click **Select Top 1000 Rows**. The customer records in this table were extracted from the **InternetSales** database, where no row has been changed between the previous **LastExtract** value and the date and time when the package was executed.
7. Minimize SQL Server Management Studio.
8. In Visual Studio, with the **Extract Customers** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging** to execute the package again, noting that no rows are transferred during this execution.
9. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
10. Close Visual Studio.

Results: After this exercise, you should have a database in which CDC is enabled, and an SSIS package that uses a stored procedure to extract modified rows, based on changes monitored by CDC.

Exercise 3: Using the CDC Control Task

► Task 1: Enable Change Data Capture

1. Return to SQL Server Management Studio, and open the **Enable CDC.sql** file in the **D:\Labfiles\Lab09\Starter\Ex3** folder.
2. Examine the query, noting that it enables CDC in the **HumanResources** database, and for the **Employee** table.
3. Click **Execute** to run the query. Two jobs should be started.
4. Open the **Fix CDC.sql** file in the **D:\Labfiles\Lab09\Starter\Ex3** folder, and then click **Execute** to run this file. This script corrects an incompatibility between the CDC functions created when enabling CDC in the database and the CDC SSIS controls that you will use in Visual Studio.

► Task 2: View Staging Tables

1. In SQL Server Management Studio, in Object Explorer, under the **Tables** folder for the **Staging** database, right-click the **dbo.EmployeeDeletes** table and click **Select Top 1000 Rows**. Note that the table is empty.
2. Repeat the previous step for the **dbo.EmployeeInserts** and **dbo.EmployeeUpdates** tables.
3. Minimize SQL Server Management Studio.

► Task 3: Create Connection Managers for CDC Components

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex3** folder.
2. In Solution Explorer, right-click the **Connection Managers** folder, and then click **New Connection Manager**.
3. In the **Add SSIS Connection Manager** dialog box, click the **ADO.NET** connection manager type, and then click **Add**.
4. In the **Configure ADO.NET Connection Manager** dialog box, click **New**.
5. In the **Connection Manager** dialog box, in the **Server name** box, type **localhost**, ensure **Windows Authentication** is selected, in the **Select or enter a database name** drop-down list, select **HumanResources**, and then click **OK**.
6. In the **Configure ADO.NET Connection Manager** dialog box, click **OK**.
7. In Solution Explorer, right-click the **localhost.HumanResources.commgr** connection manager and click **Rename**. Then rename the connection manager to **localhost.HumanResources.ADO.NET.commgr**.
8. In Solution Explorer, right-click the **Connection Managers** folder and click **New Connection Manager**.
9. In the **Add SSIS Connection Manager** dialog box, click the **ADO.NET** connection manager type, and then click **Add**.
10. In the **Configure ADO.NET Connection Manager** dialog box, click **New**.
11. In the **Connection Manager** dialog box, in the **Server name** box, type **localhost**, ensure **Windows Authentication** is selected, in the **Select or enter a database name** drop-down list, select **Staging**, and then click **OK**.
12. In the **Configure ADO.NET Connection Manager** dialog box, click **OK**.
13. In Solution Explorer, right-click the **localhost.Staging 1.commgr** connection manager and click **Rename**. Then rename the connection manager to **localhost.Staging.ADO.NET.commgr**

► Task 4: Create a Package for Initial Data Extraction

1. In Solution Explorer, right-click the **SSIS Packages** folder, and then click **New SSIS Package**.
2. In Solution Explorer, right-click **Package1.dtsx** and click **Rename**. Rename the package to **Extract Initial Employee Data.dtsx**.
3. In the **SSIS Toolbox**, in the **Other Tasks** section, drag a **CDC Control Task** to the control flow surface of the **Extract Initial Employee Data.dtsx** package.
4. On the control flow surface, right-click **CDC Control Task** and click **Rename**. Rename it to **Mark Initial Load Start**.
5. Double-click the **Mark Initial Load Start** task.

6. In the **CDC Control Task Editor** dialog box, set the following properties, and then click **OK**:
 - **SQL Server CDC database ADO.NET connection manager:** **localhost HumanResources ADO NET.**
 - **CDC control operation:** mark initial load start.
 - **Variable containing the CDC state:** click **New**. In the **Add New Variable** dialog box, click **OK** to create a variable named **CDC_State** in the **Extract Initial Employee Data** container.
 - **Automatically store state in a database table:** selected.
 - **Connection manager for the database where the state is stored:** localhost Staging ADO NET.
 - **Table to use for storing state:** click **New**. In the **Create New State Table** dialog box, click **Run to create a table named [dbo].[cdc_states] in the Staging database.**
 - **State name:** CDC_State.
7. In the SSIS Toolbox, in the **Favorites** section, drag a **Data Flow Task** to the control flow surface of the **Extract Initial Employee Data.dtsx** package, under the **Mark Initial Load Start** task.
8. On the control flow surface, right-click **Data Flow Task** and click **Rename**. Rename it to **Extract Initial Employee Data**.
9. Drag the green precedence constraint from the **Mark Initial Load Start** task to the **Extract Initial Employee Data** task.
10. Double-click the **Extract Initial Employee Data** task to view its data flow surface.
11. In the SSIS Toolbox, in the **Other Sources** section, drag an **ADO NET Source** to the data flow surface.
12. On the data flow surface, right-click **ADO NET Source**, click **Rename**, and rename it to **Employees**.
13. Double-click **Employees**.
14. In the **ADO.NET Source Editor** dialog box, set the following properties, and then click **OK**:
 - **ADO.NET connection manager:** localhost HumanResources ADO NET.
 - **Data access mode:** table or view.
 - **Name of the table or view:** "dbo"."Employee".
15. In the SSIS Toolbox, in the **Other Destinations** section, drag an **ADO NET Destination** to the data flow surface, below the **Employees** data source.
16. On the data flow surface, right-click **ADO NET Destination**, click **Rename**, and rename it to **Employee Inserts**.
17. Click the **Employees** source and drag the blue data flow path connection to the **Employee Inserts** destination.
18. Double-click **Employee Inserts**.
19. In the **ADO.NET Destination Editor** dialog box, set the following properties, and then click **OK**:
 - **Connection manager:** **localhost Staging ADO NET.**
 - **Use a table or view:** **"dbo"."EmployeeInserts".**
 - **Mappings:** on the **Mappings** tab, ensure all available input columns are mapped to destination columns of the same name.

20. Click the **Control Flow** tab.
21. In the **SSIS Toolbox**, in the **Other Tasks** section, drag a **CDC Control Task** to the control flow surface, below the **Extract Initial Employee Data** task.
22. On the control flow surface, right-click **CDC Control Task** and click **Rename**. Rename it to **Mark Initial Load End**.
23. Drag a “success” precedence constraint from the **Extract Initial Employee Data** task to the **Mark Initial Load End** task.
24. Double-click the **Mark Initial Load End** task.
25. In the **CDC Control Task Editor** dialog box, set the following properties, and then click **OK**:
 - **SQL Server CDC database ADO.NET connection manager**: localhost HumanResources ADO NET
 - **CDC control operation**: mark initial load end
 - **Variable containing the CDC state**: User:: CDC_State
 - **Automatically store state in a database table**: selected
 - **Connection manager for the database where the state is stored**: localhost Staging ADO NET
 - **Table to use for storing state**:
[dbo].[cdc_states]
 - **State name**: CDC_State
26. On the **File** menu, click **Save All**.

► Task 5: Test Initial Extraction

1. In Visual Studio, ensure that the control flow for the **Extract Initial Employee Data.dtsx** package is open, and on the **Debug** menu, click **Start Debugging**.
2. When package execution is complete, on the **Debug** menu, click **Stop Debugging**.
3. Return to SQL Server Management Studio.
4. In Object Explorer, right-click the **dbo.EmployeeInserts** table in the **Staging** database, and then click **Select Top 1000 Rows**. Note that the table now contains employee records.
5. In Object Explorer, under the **Staging** database, right-click the **Tables** folder, and then click **Refresh**. Note that a new table named **dbo.cdc_states** has been created in the **Staging** database.
6. Right-click the **dbo.cdc_states** table, and then click **Select Top 1000 Rows**. Note that the table contains an encoded string that indicates the CDC state.
7. Minimize SQL Server Management Studio.

► Task 6: Create a Package for Incremental Data Extraction

1. Return to Visual Studio.
2. In Solution Explorer, right-click the **SSIS Packages** folder, and then click **New SSIS Package**.
3. In Solution Explorer, right-click **Package1.dtsx** and click **Rename**. Rename the package to **Extract Changed Employee Data.dtsx**.
4. In the SSIS Toolbox, in the **Other Tasks** section, drag a **CDC Control Task** to the control flow surface of the package.

5. On the control flow surface, right-click **CDC Control Task** and click **Rename**. Rename it to **Get Processing Range**.
6. Double-click the **Get Processing Range** task.
7. In the **CDC Control Task Editor** dialog box, set the following properties, and then click **OK**:
 - **SQL Server CDC database ADO.NET connection manager**: localhost HumanResources ADO NET.
 - **CDC control operation**: get processing range.
 - **Variable containing the CDC state**: click **New**, and then click **OK** to create a user variable named **CDC_State**.
 - **Automatically store state in a database table**: selected.
 - **Connection manager for the database where the state is stored**: localhost Staging ADO NET.
 - **Table to use for storing state**: **[dbo].[cdc_states]**.
 - **State name**: **CDC_State**.
8. In the SSIS Toolbox, drag a **Data Flow Task** to the control flow surface, and drop it under the **Get Processing Range** task.
9. On the control flow surface, right-click **Data Flow Task** and click **Rename**. Rename it to **Extract Changed Employee Data**.
10. Click **Get Processing Range** and drag the green precedence constraint to the **Extract Changed Employee Data** task.
11. Double-click the **Extract Changed Employee Data** task to view its data flow surface.
12. In the **SSIS Toolbox**, in the **Other Sources** section, drag a **CDC Source** to the data flow surface of the **Extract Changed Employee Data** task.
13. Right-click **CDC Source**, click **Rename**, and rename it to **Employee Changes**.
14. Double-click **Employee Changes**.
15. In the **CDC Source** dialog box, set the following properties, and then click **OK**:
 - **ADO.NET connection manager**: localhost HumanResources ADO NET
 - **CDC enabled table**: **[dbo].[Employee]**
 - **Capture instance**: **dbo_Employee**
 - **CDC processing mode**: Net
 - **Variable containing the CDC state**: User::CDC_State.
16. In the **SSIS Toolbox**, in the **Other Transforms** section, drag a **CDC Splitter** to the data flow surface, below the **Employee Changes** data source.
17. Click **Employee Changes** and drag the blue data flow path connection to the **CDC Splitter** transformation.
18. In the **SSIS Toolbox**, in the **Other Destinations** section, drag an **ADO NET Destination** to the data flow surface, below and to the left of the **CDC Splitter** transformation.

19. On the data flow surface, right-click **ADO NET Destination**, click **Rename**, and rename it to **Employee Inserts**.
20. Click the **CDC Splitter** transformation and drag the blue data flow path connection to the **Employee Inserts** destination.
21. In the **Input Output Selection** dialog box, select the **InsertOutput** output, and then click **OK**.
22. Double-click **Employee Inserts**.
23. In the **ADO.NET Destination Editor** dialog box, set the following properties, and then click **OK**:
 - **Connection manager:** localhost Staging ADO NET.
 - **Use a table or view:** "dbo"."EmployeeInserts".
 - **Mappings:** on the **Mappings** tab, verify that all available input columns other than **__\$start_Isn**, **__\$operation**, **__\$update_mask** and **__\$command_id** are mapped to destination columns of the same name.
24. In the **SSIS Toolbox**, in the **Other Destinations** section, drag an **ADO NET Destination** to the data flow surface, below the **CDC Splitter** transformation.
25. On the data flow surface, right-click **ADO NET Destination**, click **Rename**, and rename it to **Employee Updates**.
26. Click the **CDC Splitter** transformation and drag the blue data flow path connection to the **Employee Updates** destination.
27. In the **Input Output Selection** dialog box, select the **UpdateOutput** output, and then click **OK**.
28. Double-click **Employee Updates**.
29. In the **ADO.NET Destination Editor** dialog box, set the following properties, and then click **OK**:
 - **Connection manager:** localhost Staging ADO NET.
 - **Use a table or view:** "dbo"."EmployeeUpdates"
 - **Mappings:** on the **Mappings** tab, verify that all available input columns other than **__\$start_Isn**, **__\$operation**, **__\$update_mask** and **__\$command_id** are mapped to destination columns of the same name.
30. In the **SSIS Toolbox**, in the **Other Destinations** section, drag an **ADO NET Destination** to the data flow surface, below and to the right of the **CDC Splitter** transformation.
31. On the data flow surface, right-click **ADO NET Destination**, click **Rename**, and rename it to **Employee Deletes**.
32. Click the **CDC Splitter** transformation and drag the blue data flow path connection to the **Employee Deletes** destination. The **DeleteOutput** output should be selected automatically.
33. Double-click **Employee Deletes**.
34. In the **ADO.NET Destination Editor** dialog box, set the following properties, and then click **OK**:
 - **Connection manager:** localhost Staging ADO NET.
 - **Use a table or view:** "dbo"."EmployeeDeletes".

- **Mappings:** on the **Mappings** tab, verify that all available input columns other than **__\$start_Isn**, **__\$operation**, **__\$update_mask** and **__\$command_id** are mapped to destination columns of the same name.
35. Click the **Control Flow** tab.
 36. In the **SSIS Toolbox**, in the **Other Tasks** section, drag a **CDC Control Task** to the control flow surface, below the **Extract Changed Employee Data** task.
 37. On the control flow surface, right-click **CDC Control Task** and click **Rename**. Rename it to **Mark Processed Range**.
 38. Click **Extract Changed Employee Data** and drag its green precedence constraint to the **Mark Processed Range** task.
 39. Double-click the **Mark Processed Range** task.
 40. In the **CDC Control Task Editor** dialog box, set the following properties, and then click **OK**:
 - **SQL Server CDC database ADO.NET connection manager:** localhost HumanResources ADO NET.
 - **CDC control operation:** mark processed range.
 - **Variable containing the CDC state:** User:: CDC_State.
 - **Automatically store state in a database table:** selected.
 - **Connection manager for the database where the state is stored:** localhost Staging ADO NET.
 - **Table to use for storing state:** [dbo].[cdc_states]
 - **State name:** CDC_State.
 41. On the **File** menu, click **Save All**.

► Task 7: Test Incremental Extraction

1. In Visual Studio, ensure that the control flow for the **Extract Changed Employee Data.dtsx** package is open, and on the **Debug** menu, click **Start Debugging**.
2. When package execution is complete, double-click the **Extract Changed Employee Data** task to verify that no rows were extracted (because no changes have been made to the source data since the initial extraction).
3. On the **Debug** menu, click **Stop Debugging**.
4. Return to SQL Server Management Studio, and open the **Change Employees.sql** file in the **D:\Labfiles\Lab09\Starter\Ex3** folder.
5. Review the Transact-SQL code and note that it truncates the **dbo.EmployeeInserts**, **dbo.EmployeeUpdates**, and **dbo.EmployeeDeletes** tables in the **Staging** database, and then makes the following changes to the **dbo.Employee** table in the **HumanResources** database:
 - Inserts a new employee record.
 - Updates employee 281 to change the Title column value.
 - Deletes employee 273.
6. Click **Execute** to run the Transact-SQL code.
7. In Visual Studio, ensure that the data flow for the **Extract Changed Employee Data** task is open, and on the **Debug** menu, click **Start Debugging**.

8. When package execution is complete, on the **Debug** menu, click **Stop Debugging**.
9. Close Visual Studio.
10. Return to SQL Server Management Studio.
11. In Object Explorer, under the **Tables** folder for the **Staging** database, right-click the **dbo.EmployeeInserts** table, and then click **Select Top 1000 Rows**. Note that the table contains the row that was inserted.
12. Repeat the previous step for the **dbo.EmployeeUpdates** and **dbo.EmployeeDeletes** tables, and verify that they contain the updated and deleted records respectively.

Results: After this exercise, you should have a **HumanResources** database in which CDC has been enabled, and an SSIS package that uses the CDC Control to extract the initial set of employee records. You should also have an SSIS package that uses the CDC Control and CDC data flow components to extract modified employee records, based on changes recorded by CDC.

Exercise 4: Using Change Tracking

► Task 1: Enable Change Tracking

1. In SQL Server Management Studio, open the **Enable CT.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder.
2. Examine the query, noting that it enables change tracking in the **ResellerSales** database, and for the **Resellers** table.
3. Click **Execute** to run the query.
4. Open the **Test CT.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder, and note that it contains statements to perform the following tasks:
 - Get the current Change Tracking version number.
 - Retrieve all data from the **Resellers** table.
 - Store the current version number as the previously retrieved version.
 - Update the **Resellers** table.
 - Get the new current version number.
 - Get all changes between the previous and current versions.
 - Store the current version number as the previously retrieved version.
 - Update the **Resellers** table again.
 - Get the new current version number.
 - Get all changes between the previous and current versions.
5. Click **Execute** and view the results. Note that:
 - The first resultset shows all reseller records.
 - The second resultset indicates that the previously retrieved version was numbered 0, and the current version is numbered 1.

- The third resultset indicates that the previously retrieved version was numbered 1, and the current version is numbered 2.

► **Task 2: Create a Stored Procedure to Retrieve Modified Rows**

1. In SQL Server Management Studio, open the **Create SP.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder.
2. Examine the Transact-SQL code in the query file, and note that it enables snapshot isolation and creates a stored procedure with a single parameter named **LastVersion**. The stored procedure performs the following tasks:
 - Sets the isolation level to snapshot.
 - Retrieves the current Change Tracking version number.
 - If the **LastVersion** parameter is -1, the stored procedure assumes that no previous versions have been retrieved, and returns all records from the **Resellers** table.
 - If the **LastVersion** parameter is not -1, the stored procedure retrieves all changes between **LastVersion** and the current version.
 - The stored procedure updates the **LastVersion** parameter to the current version, so the calling application can store the last version retrieved for next time.
 - Sets the isolation level back to read "committed".
3. Click **Execute** to run the Transact-SQL code and create the stored procedure.
4. Click **New Query**, and type the following Transact-SQL in the new query window:

```
USE ResellerSales
GO
DECLARE @p BigInt = -1;
EXEC GetChangedResellers @p OUTPUT;
SELECT @p LastVersionRetrieved;
EXEC GetChangedResellers @p OUTPUT;
```

5. Click **Execute** to test the stored procedure.

► **Task 3: Modify a Data Flow to Use the Stored Procedure**

1. In SQL Server Management Studio, open the **Reset Staging.sql** file in the **D:\Labfiles\Lab09\Starter\Ex4** folder.
2. Click **Execute** to reset the **Staging** database.
3. Minimize SQL Server Management Studio.
4. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex4** folder. If a Security warning dialogue box appears click **OK**.
5. In Solution Explorer, under **SSIS Packages**, double-click **Extract Reseller Data.dtsx**.
6. On the **SSIS** menu, click **Variables**.
7. In the Variables pane, click **Add Variable**, and add a variable with the following settings:
 - **Name:**
PreviousVersion
 - **Data Type:**
Decimal

- **Value:**
0
8. In the **SSIS Toolbox**, drag an **Execute SQL** task to the control flow surface, above the **Extract Resellers** task.
 9. Right-click the **Execute SQL** task, click **Rename**, and change the name to **Get Previously Extracted Version**.
 10. Double-click the **Get Previously Extracted Version** task.
 11. In the **Execute SQL Task Editor** dialog box, configure the following settings, and then click **OK**:
 - On the **General** tab, in the **ResultSet** drop-down list, select **Single row**.
 - On the **General** tab, in the **Connection** drop-down list, select **localhost.Staging**.
 - On the **General** tab, in the **SQLStatement** box, click the ellipsis (...) button. In the **Enter SQL Query** dialog box, enter the following Transact-SQL query, and then click **OK**:


```
SELECT MAX>LastVersion) LastVersion
FROM ExtractLog
WHERE DataSource = 'ResellerSales'
```
 - On the **Result Set** tab, click **Add**. In the **Result Name** column, change **NewResultName** to **LastVersion**, in the **Variable Name** drop-down list, select **User::PreviousVersion**, and then click **OK**.
 12. On the control flow surface, right-click the green precedence constraint between **Get Last Extract Time** and **Extract Resellers**, and click **Delete**.
 13. Click the **Get Last Extract Time** task and drag its green precedence constraint to the **Get Previously Extracted Version** task.
 14. Click the **Get Previously Extracted Version** task and drag its green precedence constraint to the **Extract Resellers** task.
 15. In the **SSIS Toolbox**, drag an **Execute SQL Task** under the **Update Last Extract Time** task on the control flow surface.
 16. Right-click **Execute SQL Task** and click **Rename**. Change the name to **Update Previous Version**.
 17. Double-click **Update Previous Version**, configure the following settings, and then click **OK**:
 - On the **General** tab, in the **Connection** drop-down list, select **localhost.Staging**.
 - On the **General** tab, in the **SQLStatement** box, click the ellipsis (...) button. In the **Enter SQL Query** dialog box, enter the following Transact-SQL query, and then click **OK**:


```
UPDATE ExtractLog
SET LastVersion = ?
WHERE DataSource = 'ResellerSales'
```
 - On the **Parameter Mapping** tab, click **Add** and create the following parameter mapping:
 - **Variable Name:** User::PreviousVersion
 - **Direction:** Input
 - **Data Type:** LARGE_INTEGER
 - **Parameter Name:** 0
 - **Parameter Size:** -1

18. On the control flow surface, right-click the green precedence constraint between **Update Last Extract Time** and **Send Success Notification**, and click **Delete**.
19. Click the **Update Last Extract Time** task and drag its green precedence constraint to the **Update Previous Version** task.
20. Click the **Update Previous Version** task and drag its green precedence constraint to the **Send Success Notification** task.
21. On the control flow surface, double-click the **Extract Resellers** task.
22. On the data flow surface, double-click the **Resellers** source.
23. In the **OLE DB Source Editor** dialog box, make the following changes to the configuration of the **Resellers** source, and then click **OK**:
 - In the **Data access mode** drop-down list, select **SQL Command**.
 - In the **SQL command** text box, type the following Transact-SQL statement:


```
EXEC GetChangedResellers ? OUTPUT
```
 - Click **Parameters**. In the **Set Query Parameters** dialog box, create the following parameter mappings, and then click **OK**:
 - **Parameters:**
@LastVersion
 - **Variables:**
User::PreviousVersion
 - **Param direction:**
InputOutput

► Task 4: Test the Package

1. In Visual Studio, with the **Extract Resellers** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging** and observe the package as it executes, noting the number of rows transferred.
2. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
3. Return to SQL Server Management Studio.
4. In Object Explorer, in the **Staging** database, right-click the **dbo.ExtractLog** table, and then click **Select Top 1000 Rows**.
5. View the data in the **ExtractLog** table, noting the value in the **LastVersion** column for the **ResellerSales** database has been updated to the latest version retrieved from the source database.
6. Right-click the **dbo.Resellers** table, and then click **Select Top 1000 Rows**. The customer records in this table were extracted from the **ResellerSales** database, where no row has been changed between the previous **LastVersion** value, and the current version.
7. Close SQL Server Management Studio without saving any changes.
8. In Visual Studio, with the **Extract Resellers** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging** to execute the package again, noting that no rows are transferred during this execution.
9. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
10. Close Visual Studio.

Results: After this exercise, you should have a database in which CT has been enabled, and an SSIS package that uses a stored procedure to extract modified rows, based on changes recorded by CT.

Lab B: Loading a Data Warehouse

Exercise 1: Loading Data from CDC Output Tables

► Task 1: Prepare the Lab Environment

1. Ensure that the 20767C-MIA-DC and 20767C-MIA-SQL virtual machines are both running, and then log on to MIA-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab09\Starter** folder, right-click **SetupB.cmd** and click **Run as administrator**.
3. When prompted, click **Yes** to confirm you want to run the command file, and wait for the script to finish.

► Task 2: Create a Data Flow for Inserts

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex5** folder. If a Security warning dialogue box appears click **OK**.
2. In Solution Explorer, in the **Connection Managers** folder, note that a connection manager for the **AWDataWarehouse** database has been created.
3. In Solution Explorer, right-click the **SSIS Packages** folder and click **New SSIS Package**.
4. In Solution Explorer, right-click **Package1.dtsx** and click **Rename**, then rename the package to **Load Employee Data.dtsx**.
5. In the SSIS Toolbox, in the **Favorites** section, drag a **Data Flow Task** to the control flow surface of the package.
6. On the control flow surface, right-click **Data Flow Task** and click **Rename**. Then rename it to **Insert Employees**.
7. Double-click **Insert Employees** to view its data flow surface.
8. In the SSIS Toolbox, drag a **Source Assistant** to the data flow surface.
9. In the **Add New Source** dialog box, in the **Select source type** list, select **SQL Server**. In the **Select connection managers** list, select **localhost.Staging**, and then click **OK**.
10. Right-click the **OLE DB Source**, click **Rename**, and change the name to **Staged Employee Inserts**.
11. Double-click the **Staged Employee Inserts** source.
12. In the **OLE DB Source Editor** dialog box, configure the following settings, and then click **OK**:
 - **OLE DB connection manager**: localhost.Staging.
 - **Data access mode**: table or view.
 - **Name of the table or view**: [dbo].[EmployeeInserts].
13. In the SSIS Toolbox, drag a **Destination Assistant** to the data flow surface below the **Staged Employee Inserts** source.
14. In the **Add New Destination** dialog box, in the **Select destination type** list, select **SQL Server**. In the **Select connection managers** list, select **localhost.AWDataWarehouse**, and click **OK**.
15. Right-click **OLE DB Destination**, click **Rename**, and change the name to **New Employees**.
16. Click the **Staged Employee inserts** source and drag the blue data flow arrow to the **New Employees** destination.
17. Double-click the **New Employees** destination.

18. In the **OLE DB Source Editor** dialog box, configure the following settings, and then click **OK**:
 - **OLE DB connection manager**: localhost.AWDDataWarehouse.
 - **Data access mode**: table or view; fast load.
 - **Name of the table or view**: [dbo].[DimEmployee].
 - **Mappings**: on the **Mappings** tab, drag the **EmployeeID** input column to the **EmployeeAlternateKey** destination column, and verify that all other input columns are mapped to destination columns of the same name, and that the **EmployeeKey** and **Deleted** destination columns are not mapped.
19. On the **File** menu, click **Save All**.

► **Task 3: Create a Data Flow for Updates**

1. In Visual Studio, click the **Control Flow** tab of the **Load Employee Data.dtsx** package.
2. In the SSIS Toolbox, in the **Favorites** section, drag a **Data Flow Task** to the control flow surface, directly below the **Insert Employees** data flow task.
3. On the control flow surface, right-click the **Data Flow Task** and click **Rename**. Then rename it to **Update Employees**.
4. Click the **Insert Employees** data flow task, and then drag its green precedence constraint to the **Update Employees** data flow task.
5. Double-click the **Update Employees** task to view its data flow surface.
6. In the SSIS Toolbox, drag a **Source Assistant** to the data flow surface.
7. In the **Add New Source** dialog box, in the **Select source types** list, select **SQL Server**. In the **Select connection managers** list, select **localhost.Staging**, and then click **OK**.
8. Right-click the **OLE DB Source**, click **Rename**, and change the name to **Staged Employee Updates**.
9. Double-click the **Staged Employee Updates** source.
10. In the **OLE DB Source Editor** dialog box, configure the following settings then click **OK**:
 - **OLE DB connection manager**: localhost.Staging.
 - **Data access mode**: table or view.
 - **Name of the table or view**: [dbo].[EmployeeUpdates].
11. In the SSIS Toolbox, drag an **OLE DB Command** to the data flow surface below the **Staged Employee Updates** source.
12. Right-click **OLE DB Command**, click **Rename**, and change the name to **Update Existing Employees**.
13. Click the **Staged Employee Updates** data source, and then drag its blue data flow path to the **Update Existing Employees** transformation.
14. Double-click the **Update Existing Employees** transformation.
15. In the **Advanced Editor** for **Update Existing Employees** dialog box, configure the following settings, and then click **OK**:
 - On the **Connection Managers** tab, in the **Connection Manager** drop-down list, select **localhost.AWDDataWarehouse**.

- On the **Component Properties** tab, set the **SqlCommand** property to the following Transact-SQL statement:

```
UPDATE dbo.DimEmployee
SET FirstName = ?, LastName = ?, EmailAddress = ?, Title = ?, HireDate = ?
WHERE EmployeeAlternateKey = ?
```

- On the **Column Mappings** tab, create the following mappings:
 - **FirstName:** Param_0
 - **LastName:** Param_1
 - **EmailAddress:** Param_2
 - **Title:** Param_3
 - **HireDate:** Param_4
 - **EmployeeID:** Param_5

16. On the **File** menu, click **Save All**.

► Task 4: Create a Data Flow for Deletes

1. In Visual Studio, click the **Control Flow** tab of the **Load Employee Data.dtsx** package.
2. In the SSIS Toolbox, in the **Favorites** section, drag a **Data Flow Task** to the control flow surface, directly below the **Update Employees** data flow task.
3. On the control flow surface, right-click the **Data Flow Task** and click **Rename**, then rename it to **Delete Employees**.
4. Click the **Update Employees** data flow task, and then drag its green precedence constraint to the **Delete Employees** data flow task.
5. Double-click the **Delete Employees** task to view its data flow surface.
6. In the SSIS Toolbox, drag a **Source Assistant** to the data flow surface.
7. In the **Add New Source** dialog box, in the **Select source types** list, select **SQL Server**. In the **Select connection managers** list, select **localhost.Staging**, and then click **OK**.
8. Right-click **OLE DB Source**, click **Rename**, and change the name to **Staged Employee Deletes**.
9. Double-click the **Staged Employee Deletes** source.
10. In the **OLE DB Source Editor** dialog box, configure the following settings, and then click **OK**:
 - **OLE DB connection manager:** localhost.Staging.
 - **Data access mode:** table or view.
 - **Name of the table or view:** [dbo].[EmployeeDeletes].
11. In the SSIS Toolbox, drag an **OLE DB Command** to the data flow surface below the **Staged Employee Deletes** source.
12. Right-click the **OLE DB Command**, click **Rename**, and change the name to **Delete Existing Employees**.
13. Click the **Staged Employee Deletes** data source, and then drag its blue data flow path to the **Delete Existing Employees** transformation.
14. Double-click the **Delete Existing Employees** transformation.

15. In the **Advanced Editor for Delete Existing Employees** dialog box, configure the following settings, and then click **OK**:

- On the **Connection Managers** tab, in the **Connection Managers** drop-down list, select localhost.AWDDataWarehouse.
- On the **Component Properties** tab, set the **SqlCommand** property to the following Transact-SQL statement:

```
UPDATE dbo.DimEmployee
SET Deleted = 1
WHERE EmployeeAlternateKey = ?
```

- On the **Column Mappings** tab, map the **EmployeeID** column to **Param_0**.
16. On the **File** menu, click **Save All**.

► Task 5: Test the Package

1. In Visual Studio, click the **Control Flow** tab of the **Load Employee Data.dtsx** package.
2. On the **Debug** menu, click **Start Debugging** and observe the package as it executes.
3. When execution is complete, double-click the **Insert Employees** data flow task and review the number of rows inserted.
4. In the **Data Flow Task** drop-down list at the top of the data flow surface designer, select **Update Employees** and **Delete Employees** in turn, noting the number of rows processed by these tasks.
5. On the **Debug** menu, click **Stop Debugging**.
6. Close Visual Studio, saving the changes if prompted.

Results: After this exercise, you should have an SSIS package that uses data flows to apply inserts, updates, and logical deletes in the data warehouse, based on staging tables extracted by the CDC Control Task and data flow components.

Exercise 2: Using a Lookup Transformation to Insert or Update Dimension Data

► Task 1: View Data Flows

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex6** folder. If a Security warning dialogue box appears click **OK**.
2. In Solution Explorer, in the **SSIS Packages** folder, double-click **Load Products Data.dtsx**.
3. On the control flow surface, in the **Load Product Dimension Hierarchy** sequence container, double-click the **Load Product Category Dimension** task to view its data flow.
4. Examine the data flow for the **Load Product Category Dimension** task, and note the following features:
 - The **Staged Product Category Data** source extracts product category data from the **InternetSales** and **ResellerSales** tables in the **Staging** database.
 - The **Lookup Existing Product Categories** task retrieves the **ProductCategoryKey** value for product categories that exist in the **DimProductCategory** table in the **AWDataWarehouse**

database by matching the **Product Category** business key in the **Staging** database to the **Product Category** alternative key in the data warehouse.

- The **Lookup No Match Output** data flow path from the **Lookup Existing Product Categories** task connects to the **New Product Categories** destination, and the **Lookup Match Output** data flow path connects to the **Update Existing Product Categories** task.
 - The **New Product Categories** destination loads new product category records into the **DimProductCategory** table.
 - The **Update Existing Product Categories** task executes a Transact-SQL statement to update the **ProductCategoryName** column in the **DimProductCategory** table for an existing row based on the **ProductCategoryKey**.
5. In the **Data Flow Task** drop-down list, select **Load Product Subcategory Dimension**, and note that this data flow inserts or updates product subcategory dimension data using a similar approach to the **Load Product Category Dimension** data flow. Additionally, there is a lookup task to retrieve the **ProductCategoryKey** in **AWDataWarehouse** for the parent category, which should have already been loaded.

► Task 2: Create a Data Flow

1. In Visual Studio, in the **Load Products Data.dtsx** package designer, click the **Control Flow** tab.
2. In the SSIS Toolbox, drag a **Data Flow Task** to the control flow surface, under the **Load Product Subcategory Dimension** task in the **Load Product Dimension Hierarchy** sequence container.
3. Right-click the data flow task you added, click **Rename**, and change the name to **Load Product Dimension**.
4. Click the **Load Product Subcategory Dimension** task, and drag the green precedence constraint to the **Load Product Dimension** task.
5. Double-click the **Load Product Dimension** task to view the data flow surface.
6. In the SSIS Toolbox, drag a **Source Assistant** to the data flow surface.
7. In the **Add New Source** dialog box, in the **Select source type** list, select **SQL Server**. In the **Connection Managers** list, select **localhost.Staging**, and then click **OK**.
8. Right-click **OLE DB Source**, click **Rename**, and change the name to **Staged Product Data**.
9. Double-click the **Staged Product Data** source, and in the **OLE DB Source Editor** dialog box, in the **Data access mode** drop-down list, select **SQL Command**.
10. In the **OLE DB Source Editor** dialog box, in the **SQL command text:** box, type the following Transact-SQL statement, and then click **OK**:

```
SELECT DISTINCT ProductSubcategoryBusinessKey, ProductBusinessKey, ProductName,
StandardCost, Color, ListPrice, Size, Weight, Description
FROM dbo.InternetSales
UNION
SELECT DISTINCT ProductSubcategoryBusinessKey, ProductBusinessKey, ProductName,
StandardCost, Color, ListPrice, Size, Weight, Description FROM dbo.ResellerSales
```

► Task 3: Add a Lookup Transformation for Parent Keys

1. In the SSIS Toolbox, drag a **Lookup** transformation to the data flow surface below the **Staged Product Data** source.
2. Right-click the **Lookup** transformation, click **Rename**, and change the name to **Lookup Parent Subcategory**.

3. Click the **Staged Product Data** source and drag the blue data flow arrow to the **Lookup Parent Subcategory** transformation.
4. Double-click the **Lookup Parent Subcategory** transformation.
5. In the **Lookup Transformation Editor** dialog box, configure the following settings, and then click **OK**:
 - On the **General** tab, in the **Specify how to handle rows with no matching entries** drop-down list, ensure **Fail component** is selected.
 - On the **Connection** tab, in the **OLE DB connection manager** drop-down list, select **localhost.AWDataWarehouse**.
 - On the **Connection** tab, in the **Use a table or a view** drop-down list, select **[dbo].[DimProductSubcategory]**.
 - On the **Columns** tab, drag the **ProductSubcategoryBusinessKey** column in the **Available Input Columns** list to the **ProductSubcategoryAlternateKey** column in the **Available lookup columns** list. Then, in the **Available lookup columns** list, select the check box for the **ProductSubcategoryKey** column.

► **Task 4: Configure a Lookup Transformation**

1. In the SSIS Toolbox, drag another **Lookup** transformation to the data flow surface below the **Lookup Parent Subcategory** transformation.
2. Right-click **Lookup**, click **Rename**, and change the name to **Lookup Existing Products**.
3. Click the **Lookup Parent Subcategory** transformation and drag the blue data flow arrow to the **Lookup Existing Products** transformation.
4. In the **Input Output Selection** dialog box, in the **Output** drop-down list, select **Lookup Match Output**, and then click **OK**.
5. Double-click the **Lookup Existing Products** transformation.
6. In the **Lookup Transformation Editor** dialog box, configure the following settings, and then click **OK**:
 - On the **General** tab, in the **Specify how to handle rows with no matching entries** drop-down list, select **Redirect rows to no match output**.
 - On the **Connection** tab, in the **OLE DB connection manager** drop-down list, select **localhost.AWDataWarehouse**.
 - On the **Connection** tab, in the **Use a table or a view** drop-down list, select **[dbo].[DimProduct]**.
 - On the **Columns** tab, drag the **ProductBusinessKey** column in the **Available Input Columns** list to the **ProductAlternateKey** column in the **Available lookup columns** list. Then in the **Available lookups columns** list, select the check box for the **ProductKey**.

► **Task 5: Add a Destination for New Products**

1. In the SSIS Toolbox, drag a **Destination Assistant** to the data flow surface below and to the right of the **Lookup Existing Products** transformation.
2. In the **Add New Destination** dialog box, in the **Types** list, select **SQL Server**, in the **Connection Managers** list, select **localhost.AWDataWarehouse**, and then click **OK**.
3. Right-click **OLE DB Destination**, click **Rename**, and change the name to **New Products**.

4. Click the **Lookup Existing Products** transformation and drag the blue data flow arrow to the **New Products** destination.
5. In the **Input Output Selection** dialog box, in the **Output** drop-down list, select **Lookup No Match Output**, and then click **OK**.
6. Double-click the **New Products** destination.
7. In the **OLE DB Destination Editor** dialog box, in the **Name of the table or the view** drop-down list, select **[dbo].[DimProduct]**.
8. On the **Mappings** tab, create the following mappings, and then click **OK**.
 - **<ignore>**: ProductKey
 - **ProductBusinessKey**: ProductAlternateKey
 - **ProductName**: ProductName
 - **StandardCost**: StandardCost
 - **Color**: Color
 - **ListPrice**: ListPrice
 - **Size**: Size
 - **Weight**: Weight
 - **Description**: Description
9. **<ignore>**: ProductSubcategoryKey

► **Task 6: Add an OLE DB Command for Updated Product Records**

1. In the SSIS Toolbox, drag an **OLE DB Command** to the data flow surface below and to the left of the **Lookup Existing Products** transformation.
2. Right-click **OLE DB Command**, click **Rename**, and change the name to **Update Existing Products**.
3. Click the **Lookup Existing Products** transformation and drag the blue data flow arrow to the **Update Existing Products** transformation. The **Lookup Match Output** is automatically selected.
4. Double-click the **Update Existing Products** transformation.
5. In the **Advanced Editor for Update Existing Products** dialog box, configure the following settings, and then click **OK**:
 - On the **Connection Managers** tab, in the **Connection Manager** drop-down list, select **localhost.AWDDataWarehouse**.
 - On the **Component Properties** tab, set the **SqlCommand** property to the following Transact-SQL statement:


```
UPDATE dbo.DimProduct
SET ProductName = ?, StandardCost = ?, Color = ?, ListPrice = ?, Size = ?,
Weight = ?, Description = ?, ProductSubcategoryKey = ?
WHERE ProductKey = ?
```
 - On the **Column Mappings** tab, create the following mappings:
 - **ProductName**: Param_0
 - **StandardCost**: Param_1
 - **Color**: Param_2

- **ListPrice:** Param_3
- **Size:** Param_4
- **Weight:** Param_5
- **Description:** Param_6
- **ProductSubcategoryKey:** Param_7
- **ProductKey:** Param_8

► Task 7: Test the Package

Note: When debugging, the rows do not always appear on the data flow path. An alternate check is to open the database and ensure the changes have been successful.

1. With the **Load Product Dimension** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging** and observe the package as it executes, noting that all rows flow to the **New Products** destination (because the data warehouse contained no existing product records). When execution is complete, on the **Debug** menu, click **Stop Debugging**.
2. On the **Debug** menu, click **Start Debugging** and observe the package as it executes again, noting that this time, all rows flow to the **Update Existing Products** transformation. This is because all staged product records were loaded to the data warehouse during the previous execution, so they all match existing records. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
3. Close Visual Studio, saving your changes if prompted.

Results: After this exercise, you should have an SSIS package that uses a Lookup transformation to determine whether product records already exist, and updates them or inserts them as required.

Exercise 3: Implementing a Slowly Changing Dimension

► Task 1: Execute a Package to Load a Non-Changing Dimension

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex7** folder. If a Security warning dialogue box appears click **OK**.
2. In Solution Explorer, under the **SSIS Packages** folder, double-click **Load Geography Data.dtsx**.
3. Review the control flow and data flow defined in the package. This package includes a simple data flow to load staged geography data into the data warehouse. Only new rows are loaded, and rows that match existing data are discarded.
4. On the **Debug** menu, click **Start Debugging**, and observe the package execution as it loads geography data into the data warehouse.
5. When package execution has completed, on the **Debug** menu, click **Stop Debugging**.

► Task 2: Observe a Data Flow for a Slowly Changing Dimension

1. In Solution Explorer, under the **SSIS Packages** folder, double-click **Load Reseller Data.dtsx**.
2. On the control flow surface, double-click the **Load Reseller Dimension** task to view its data flow.
3. Examine the data flow for the **Load Reseller Dimension** task, and note the following features:
 - The **Staged Reseller Data** source extracts data from the **Resellers** table in the **Staging** database.

- The **Lookup Geography Key** transformation looks up the geography key for the reseller in the **DimGeography** table in the **AWDataWarehouse** database.
 - The **Reseller SCD** is an SCD transformation that has generated the remaining transformations and destinations. You can double-click the **Reseller SCD** transformation to view the wizard used to configure the SCD, and then click **Cancel** to avoid making any unintentional changes.
 - The **Reseller SCD** transformation maps the **ResellerBusinessKey** input column to the **ResellerAlternateKey** dimension column and uses it as a business key to identify existing records.
 - The **Reseller SCD** transformation treats **AddressLine1**, **AddressLine2**, **BusinessType**, **GeographyKey**, and **NumberEmployees** as historical attributes, **Phone** and **ResellerName** as changing attributes, and **YearOpened** as a fixed attribute.
4. On the **Debug** menu, click **Start Debugging** and observe the data flow as it executes.
 5. When execution is complete, on the **Debug** menu, click **Stop Debugging**.

► Task 3: Implement a Slowly Changing Dimension Transformation

1. In Solution Explorer, under the **SSIS Packages** folder, double-click **Load Customer Data.dtsx**.
2. On the control flow surface, double-click the **Load Customer Dimension** task to view its data flow. Note that a source to extract customer data from the **Staging** database and a Lookup transformation that retrieves a **GeographyKey** value from the **AWDataWarehouse** database have already been added to the data flow.
3. In the SSIS Toolbox pane, drag a **Slowly Changing Dimension** transformation to the data flow surface, below the **Lookup Geography Key** transformation.
4. Right-click **Slowly Changing Dimension**, click **Rename**, and change the name to **Customer SCD**.
5. Click the **Lookup Geography Key** transformation and drag the blue data flow arrow to the **Customer SCD** transformation.
6. In the **Input Output Selection** dialog box, in the **Output** drop-down list, select **Lookup Match Output**, and then click **OK**.
7. Double-click the **Customer SCD** transformation to start the SCD wizard.
8. On the **Welcome to the Slowly Changing Dimension Wizard** page, click **Next**.
9. On the **Select a Dimension Table and Keys** page, in the **Connection manager** drop-down list, select **localhost.AWDataWarehouse**, in the **Table or view** drop-down list, select **[dbo].[DimCustomer]**, specify the following column mappings, and then click **Next**:

Input Columns	Dimension Columns	Key Type
AddressLine1	AddressLine1	Not a key column
AddressLine2	AddressLine2	Not a key column
BirthDate	BirthDate	Not a key column
CommuteDistance	CommuteDistance	Not a key column
	CurrentRecord	

Input Columns	Dimension Columns	Key Type
CustomerBusinessKey	CustomerAlternateKey	Business key
EmailAddress	EmailAddress	Not a key column
FirstName	FirstName	Not a key column
Gender	Gender	Not a key column
GeographyKey	GeographyKey	Not a key column
HouseOwnerFlag	HouseOwnerFlag	Not a key column
LastName	LastName	Not a key column
MaritalStatus	MaritalStatus	Not a key column
MiddleName	MiddleName	Not a key column
NumberCarsOwned	NumberCarsOwned	Not a key column
Occupation	Occupation	Not a key column
Phone	Phone	Not a key column
Suffix	Suffix	Not a key column
Title	Title	Not a key column

10. On the **Slowly Changing Dimension Columns** page, specify the following change types, and then click **Next**:

Dimension Columns	Change Type
AddressLine1	Historical attribute
AddressLine2	Historical attribute
BirthDate	Changing attribute
CommuteDistance	Historical attribute
EmailAddress	Changing attribute
FirstName	Changing attribute

Dimension Columns	Change Type
Gender	Historical attribute
GeographyKey	Historical attribute
HouseOwnerFlag	Historical attribute
LastName	Changing attribute
MaritalStatus	Historical attribute
MiddleName	Changing attribute
NumberCarsOwned	Historical attribute
Occupation	Historical attribute
Phone	Changing attribute
Suffix	Changing attribute
Title	Changing attribute

11. On the **Fixed and Changing Attribute Options** page, leave both options clear, and then click **Next**.
12. On the **Historical Attribute Options** page, select **Use a single column** to show current and expired records. In the **Column to indicate current record** drop-down list, select **CurrentRecord**. In the **Value when current** drop-down list, select **True**, and in the **Expiration** value drop-down list, select **False**. Click **Next**.
13. On the **Inferred Dimension Members** page, clear the **Enable inferred member support** option, and then click **Next**.
14. On the **Finish the Slowly Changing Dimension Wizard** page, click **Finish**. Note that a number of transformations and a destination are created.

► Task 4: Test the Package

Note: When debugging, the rows do not always appear on the data flow path. An alternative check is to open the database and ensure the changes have been successful.

1. With the **Load Customer Dimension** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging**.
2. Observe the package as it executes, noting that all rows pass through the **New Output** data flow path.
3. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
4. On the **Debug** menu, click **Start Debugging** and observe the package as it executes again, noting that no rows pass through the **New Output** data flow path, because they already exist and no changes have been made.
5. When execution is complete, on the **Debug** menu, click **Stop Debugging**.

6. Start SQL Server Management Studio and connect to the **MIA-SQL** database engine instance by using Windows authentication.
7. Open the **Update Customers.sql** file in the **D:\Labfiles\Lab09\Starter\Ex7** folder.
8. Examine the script and note that it updates two records in the **Staging** database, changing one customer's phone number and another's marital status.
9. Click **Execute**.
10. When the query has completed, close SQL Server Management Studio without saving changes.
11. In Visual Studio, on the **Debug** menu, click **Start Debugging** and observe the package as it executes, noting that one row passes through the **Historical Attribute Inserts Output** data flow path, and another passes through **Changing Attributes Updates Output**.
12. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
13. Close Visual Studio, saving your changes if prompted.

Results: After this exercise, you should have an SSIS package that uses an SCD transformation to load data into a dimension table.

Exercise 4: Using the MERGE Statement

► Task 1: Examine a Control Flow That Uses the MERGE Statement

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab09\Starter\Ex8** folder. If a Security warning dialogue box appears click **OK**.
2. In Solution Explorer, under the **SSIS Packages** folder, double-click **Load Reseller Sales Data.dtsx**.
3. On the control flow surface, double-click the **Merge Reseller Sales** task.
4. In the **Execute SQL Task Editor** dialog box, note the following configuration settings, and then click **Cancel**:
 - The task uses the **localhost.Staging** connection manager to connect to the **Staging** database.
 - The task executes a Transact-SQL MERGE statement that retrieves reseller sales and related dimension keys from the **Staging** and **AWDataWarehouse** databases. It then matches these records with the **FactResellerSales** table, based on the **SalesOrderNumber** and **SalesOrderLineNumber** columns, updates rows that match, and inserts new records for rows that do not.
5. On the **Debug** menu, click **Start Debugging** and observe the package as it executes.
6. When execution is complete, on the **Debug** menu, click **Stop Debugging**.

► Task 2: Create a Package That Uses the MERGE Statement

1. In Solution Explorer, right-click the **SSIS Packages** folder and click **New SSIS Package**.
2. Right-click **Package1.dtsx**, click **Rename**, and change the name to **Load Internet Sales Data.dtsx**.
3. In the SSIS Toolbox pane, drag an **Execute SQL Task** to the control flow surface.
4. Right-click **Execute SQL Task**, click **Rename**, and change the name to **Merge Internet Sales Data**.
5. Double-click **Merge Internet Sales Data**.

6. In the **Execute SQL Task Editor** dialog box, configure the following settings, and then click **OK**:
 - In the **General** tab in the **Connection** drop-down list, select **localhost.Staging**.
 - Click **Browse**, and then open D:\Labfiles\Lab09\Starter\Ex8\Merge Internet Sales.sql.

► **Task 3: Test the Package**

1. On the **Debug** menu, click **Start Debugging** and observe the package as it executes.
2. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
3. Close Visual Studio.

Results: After this exercise, you should have an SSIS package that uses an Execute SQL task to execute a MERGE statement that inserts or updates data in a fact table.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 10: Enforcing Data Quality

Lab A: Cleansing Data

Exercise 1: Creating a DQS Knowledge Base

► Task 1: Prepare the Lab Environment

1. Ensure that the 20767C-MIA-DC, 20767C-MIA-CLI, and 20767C-MIA-SQL virtual machines are all running, and then log on to 20767C-MIA-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the D:\Labfiles\Lab10\Starter folder, right-click **Setup.cmd** and click **Run as administrator**.
3. In the **User Account** Control dialog box, click **Yes**.

► Task 2: View Existing Data

1. Log on to 20767C-MIA-CLI as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the D:\Labfiles\Lab10\Starter folder, double-click **Sample Customer Data.xls** to open it in Excel.
3. In Excel, on the **CountryRegion** worksheet, note that the data shows the number of customers by country name and country code. Also note that the data contains the following quality issues:
 - The list of country names includes both **America** and **United States**. For the purposes of Adventure Works sales, these represent the same sales territory.
 - The list of country names includes both **Great Britain** and **United Kingdom**. Again, for the purposes of Adventure Works sales, these represent the same sales territory.
 - The list of country codes includes both **GB** and **UK**. For the purposes of Adventure Works sales, these represent the same sales territory.
4. In Excel, on the **StateProvince** worksheet, note that the data shows the number of customers by country and state. In addition to the concerns previously identified on the **CountryRegion** worksheet, the data exposes the following quality issues:
 - The states **Oregon** and **California** exist in both **America** and **United States**.
 - The states **Wales** and **England** exist in both **Great Britain** and **United Kingdom**.
 - Australia includes **New South Wales** and a state named **NSW**; both represent the same state.
 - **United States** includes **Washington** and **WA**, which represent the same state, as do **California** and **CA**.
5. In Excel, on the **Gender** worksheet, note that two customers have a gender code of **W**. Valid values for the gender code are **F** and **M**.
6. On the **Sample Customer Data** worksheet, apply column filters to explore the data further and view the source records for the anomalous data.
7. Close Excel without saving any changes to the workbook.

► Task 3: Create a Knowledge Base

1. Click **Search Windows**, type **Data Quality Client**, and then press Enter
2. When prompted, enter the server name **MIA-SQL**, and click **Connect**.

3. In SQL Server Data Quality Services, in the **Knowledge Base Management** section, click **New Knowledge Base**.
4. On the **New Knowledge Base** page enter the following details, and then click **Next**:
 - **Name**: Customer KB.
 - **Description**: Customer data knowledge base.
 - **Create Knowledge Base From**: Existing Knowledge Base.
 - **Select Knowledge Base**: DQS Data.
 - **Select Activity**: Domain Management.
5. In the **Domain** list, click **Country/Region**.
6. In the Country/Region pane, click the **Domain Values** tab and note the country values that are defined in this knowledge base. Leading domain values are shown in bold, with synonym values indented below.
7. In the **Domain** list, click **Country/Region (two-letter leading)**, and on the **Domain Values** tab, note the country code values that are defined in this knowledge base. This defines the same values as the **Country/Region** domain, but with the two-character country code designated as the leading value.
8. In the **Domain** list, click **US - State**, and on the **Domain Values** tab, note the state values that are defined in this knowledge base. Note that state codes (such as **CA** and **OR**) are valid, but are shown as synonyms for leading state name values (such as **California** and **Oregon**). State values commonly entered in error are shown with a red cross and a valid value to which they are automatically corrected.
9. With the **US - State** domain selected, on the **Domain Properties** tab, change **Domain Name** to **State**. Later in this lab, you will use this domain for states and provinces in countries other than the United States.
10. Click **Create a domain**.
11. In the **Create Domain** dialog box, enter the following details, and then click **OK**:
 - **Domain Name**: Gender.
 - **Description**: Male or Female.
 - **Data Type**: String.
 - **Use Leading Values**: Selected.
 - **Normalize String**: Selected.
 - **Format Output to**: Upper Case.
 - **Language**: English.
 - **Enable Speller**: Selected.
 - **Disable Syntax Error Algorithms**: Not selected.
12. In the **Domain** list, click **Gender**. Click the **Domain Values** tab and note that the value **DQS_NULL** (which represents a null value) has already been defined as a valid value for this domain.
13. On the **Domain Values** tab, click **Add new domain value**, and then enter the value **F** and press Enter.
14. Repeat the previous step to add the values **M**, **Female**, and **Male**.

15. Click the value **F**, and then hold the Ctrl key and click the value **Female**. Then, with both values selected, click **Set selected domain values as synonyms**. Depending on the screen resolution, this may be in a drop-down list at the end of the toolbar above the table. Note that **F** becomes the leading value, with **Female** indented below it.
16. Repeat the previous step to set **M** and **Male** as synonyms with **M** as the leading value.
17. Click **Finish**.
18. When prompted to publish the knowledge base, click **No**.

► **Task 4: Perform Knowledge Discovery**

1. In SQL Server Data Quality Services, under **Recent knowledge base**, click **Customer KB**, and then click **Knowledge Discovery**.
2. On the **Map** page, in the **Data Source** drop-down list, select **Excel File**, click **Browse** and open the **Sample Customer Data.xls** file in the **D:\Labfiles\Lab10\Starter** folder.
3. After the file has loaded, in the **Worksheet** drop-down list, select **'Sample Customer Data\$'** and ensure that **Use first row as header** is selected.
4. In the **Mappings** table, add the following mappings, and then click **Next**:

Source Column	Domain
CountryRegionCode (String)	Country/Region (two-letter leading)
CountryRegionName (String)	Country/Region
StateProvinceName (String)	State
Gender (String)	Gender

5. On the **Discover** page, click **Start**.
6. When discovery analysis is complete, click **Next**.
7. On the **Manage Domain Values** page, in the **Domain** list, click **State** and view the new values that have been discovered.
8. In the list of values, click **New South Wales**, press Ctrl, click **NSW**, and then click **Set selected domain values as synonyms**. Note that **NSW** now has a **Correct to** value of **New South Wales**.
9. On the **Manage Domain Values** page, in the **Domain** list, click **Country/Region (two-letter leading)** and view the new values that were discovered.
10. In the **Type** column for the value **UK**, select the symbol for **Error**, and in the **Correct To** column, type **GB**.
11. Clear the **Show Only New** check box and note that **UK** now appears as an error under the **GB** leading value.
12. On the **Manage Domain Values** page, in the **Domain** list, click **Gender** and view the new values that have been discovered.
13. In the **Type** column for the value **W**, click the symbol for **Invalid**, and in the **Correct To** column, type **F**.

14. Clear the **Show Only New** check box and note that **W** now appears as an invalid value under the **F** leading value.
15. On the **Manage Domain Values** page, in the **Domain** list, click **Country/Region** and view the new values that have been discovered. Clear the **Show Only New** check box to show all values for this domain.
16. Click **Finish**.
17. When prompted to publish the knowledge base, click **Publish**, and when notified that publication was successful, click **OK**.
18. Keep SQL Server Data Quality Services open for the next exercise.

Results: After this exercise, you should have created a knowledge base and performed knowledge discovery.

Exercise 2: Using a DQS Project to Cleanse Data

► Task 1: Create a Data Quality Data Project

1. In SQL Server Data Quality Services, in the **Data Quality Projects** section, click **New Data Quality Project**.
2. On the **New Data Quality Project** page, enter the following details, and then click **Next**:
 - **Name:** Cleanse Customer Data.
 - **Description:** Apply Customer KB to customer data.
 - **Use knowledge base:** Customer KB.
 - **Select Activity:** Cleansing.
3. On the **Map** page, in the **Data Source** list, select **SQL Server**, in the **Database** list select **InternetSales**, and in the **Table/View** list select **Customers**.
4. In the **Mappings** table, add the following mappings, and click **Next**:

Source Column	Domain
CountryRegionCode (nvarchar)	Country/Region (two-letter leading)
CountryRegionName (nvarchar)	Country/Region
StateProvinceName (nvarchar)	State
Gender (nvarchar)	Gender

5. On the **Cleanse** page, click **Start**.
6. When cleansing is complete, review the source statistics in the Profiler pane, and click **Next**.
7. On the **Manage and View Results** page, select the **Country/Region** domain, and on the **Suggested** tab, note that **DQS** has found the value **Australia**, which is likely to be a typographical error, and suggested it be corrected to **Australia**.

8. Click the **Approve** option to accept the suggested correction.
9. Click the **Corrected** tab to see all the corrections made for the **Country/Region** domain.
10. In the list of domains, click **Country/Region (two-letter leading)** and on the **Corrected** tab, note the corrections made for this domain.
11. In the list of domains, click **Gender** and on the **Corrected tab**, note the corrections that made for this domain.
12. In the list of domains, click **State** and on the **New** tab, note the new values that have been found for this domain.
13. Click the **Approve** option for each new value, and then click the **Corrected** tab to see all the corrections made for the **State** domain.
14. Click **Next**.
15. On the **Export** page, view the output data, which contains the following columns:
 - **Output**. For each column in the source that was not mapped to a domain, containing the original source value.
 - **Source**. For each domain in the knowledge base containing the original source value.
 - **Output**. For each domain containing the output value.
 - **Reason**. For each domain containing the reason for the output value.
 - **Confidence**. For each domain indicating the confidence level (0 to 1) for any corrections to the original value.
 - **Status**. For each domain indicating the status of the output value.
16. In the **Export cleansing results** section, in the **Destination Type** list, select **Excel File**, in the **Excel file name** box, enter **D:\Labfiles\Lab10\Starter\CleansedCustomers.xls**, ensure that **Data and Cleansing Info** is selected, and then click **Export**.
17. When you are notified that the file download is complete, click **Close**, and then click **Finish**,
18. Close SQL Server Data Quality Services.
19. In the D:\Labfiles\Lab10\Starter folder, double-click **CleansedCustomers.xls** to open it with Excel.
20. Review the output from the cleansing process, and then close Excel.

Results: After this exercise, you should have used a DQS project to cleanse data and export it as an Excel workbook.

Exercise 3: Using DQS in an SSIS Package

► Task 1: Add a DQS Cleansing Transformation to a Data Flow

1. Return to the 20767C-MIA-SQL virtual machine and start Visual Studio®.
2. On the **File** menu, point to **Open**, click **Project/Solution**, browse to the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab10\Starter** folder, and then click **Open**.
3. In Solution Explorer, in the **SSIS Packages** folder, double-click **Extract Internet Sales Data.dtsx**.
4. On the **SSIS** menu, click **SSIS Toolbox** to display the toolbox.
5. On the control flow surface, double-click the **Extract Customers** task to view its data flow.
6. In the **SSIS Toolbox**, in the **Other Transforms** section, drag a **DQS Cleansing** transformation onto the data flow surface.
7. On the data flow surface, right-click **DQS Cleansing**, click **Rename**, and change the name to **Cleanse Customer Data**.
8. Right-click the data flow path between **Customers** and **Staging DB**, and then click **Delete**.
9. Reposition **Cleanse Customer Data** below **Customers**, click **Customers**, and drag the data flow arrow to **Cleanse Customer Data**.
10. Double-click **Cleanse Customer Data**.
11. In the **DQS Cleansing Transformation Editor** dialog box, configure the following settings, and then click **OK**:
 - On the **Connection Manager** tab, next to the **Data quality connection manager** drop-down list, click **New**. In the **DQS Cleansing Connection Manager** dialog box, in the **Server Name** box, type **MIA-SQL**, and then click **OK**.
 - On the **Connection Manager** tab, in the **Data Quality Knowledge Base** drop-down list, select **Customer KB**.
 - On the **Mapping** tab, select the check boxes for the **Gender**, **StateProvinceName**, **CountryRegionCode**, and **CountryRegionName** input columns. Create the following mappings with the default source, output, and status alias values:

Input Column	Domain
Gender	Gender
StateProvinceName	State
CountryRegionCode	Country/Region (two-letter leading)
CountryRegionName	Country/Region

- On the **Advanced** tab, ensure that **Standardize output** is selected.
12. Click **Cleanse Customer Data**, and drag the data flow arrow to **Staging DB**.
 13. Double-click **Staging DB**.

14. In the **OLE DB Destination Editor** dialog box, on the **Mappings** tab, change the current column mappings for the following destination columns, and then click **OK**:

Input Column	Destination Column
Gender_Output	Gender
StateProvinceName_Output	StateProvinceName
CountryRegionCode_Output	CountryRegionCode
CountryRegionName_Output	CountryRegionName

► **Task 2: Test the Package**

1. With the **Extract Customers** data flow displayed in the designer, on the **Debug** menu, click **Start Debugging**.
2. Observe the package as it executes, noting the number of rows processed by the **Cleanse Customer Data** transformation. Execution may take several minutes.
3. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
4. Close Visual Studio, and save the solution files if prompted.

Results: After this exercise, you should have created and tested an SSIS package that cleanses data.

Lab B: Deduplicating Data

Exercise 1: Creating a Matching Policy

► Task 1: Prepare the Lab Environment

1. Ensure that the 20767C-MIA-DC, 20767C-MIA-CLI, and 20767C-MIA-SQL virtual machines are all running, and then log on to 20767C-MIA-SQL as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab10\Starter** folder, right-click **LabB.cmd** and click **Run as administrator**.
3. In the **User Account Control** dialog box, click **Yes**.

► Task 2: Create a Matching Policy

1. Log on to 20767C-MIA-CLI as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. Click **Search Windows**, type **Data Quality Client**, and then press Enter.
3. When prompted, enter the server name **MIA-SQL**, and click **Connect**.
4. In the **Knowledge Base Management** pane, under **Recent knowledge base**, click **Customer KB**, and then click **Matching Policy**.
5. On the **Map** page, in the **Data Source** drop-down list, select **Excel File**, in the **Excel File** box, click **Browse**, and then select **Sample Staged Data.xls** in the **D:\Labfiles\Lab10\Starter** folder.
6. In the **Worksheet** drop-down list, ensure that **Sheet1\$** and **Use first row as header** are selected.
7. In the **Mappings** table, in the first **Source Column** row, select **FirstName (String)**. Click **Create a domain**, and create a domain using the following properties, and then click **OK**:
 - **Domain Name:** FirstName
 - **Description:** First Name.
 - **Data Type:** String
 - **Use Leading Values:** Selected
 - **Normalize String:** Selected
 - **Format Output to:** None
 - **Language:** English
 - **Enable Speller:** Selected
 - **Disable Syntax Error Algorithms:** Not selected
8. In the **Mappings** table, in the second **Source Column** row, select **LastName (String)**. Click **Create a domain**, create a domain using the following properties, and then click **OK**:
 - **Domain Name:** LastName
 - **Description:** Last Name.
 - **Data Type:** String
 - **Use Leading Values:** Selected
 - **Normalize String:** Selected
 - **Format Output to:** None

- **Language:** English
 - **Enable Speller:** Selected
 - **Disable Syntax Error Algorithms:** Not selected
9. In the **Mappings** table, in the third **Source Column** row, select **AddressLine1 (String)**. Click **Create a domain**, create a domain using the following properties, and then click OK:
- **Domain Name:** AddressLine1
 - **Description:** First line of address.
 - **Data Type:** String
 - **Use Leading Values:** Selected
 - **Normalize String:** Selected
 - **Format Output to:** None
 - **Language:** English
 - **Enable Speller:** Selected
 - **Disable Syntax Error Algorithms:** Not selected
10. In the **Mappings** table, in the fourth **Source Column** row, select **City (String)**. Click **Create a domain**, create a domain using the following properties, and then click **OK**:
- **Domain Name:** City
 - **Description:** City.
 - **Data Type:** String
 - **Use Leading Values:** Selected
 - **Normalize String:** Selected
 - **Format Output to:** None
 - **Language:** English
 - **Enable Speller:** Selected
 - **Disable Syntax Error Algorithms:** Not selected
11. In the **Mappings** table, in the fifth **Source Column** row, select **EmailAddress (String)**. Click **Create a domain**, create a domain using the following properties, and then click **OK**:
- **Domain Name:** EmailAddress **Description:** Email address.
 - **Data Type:** String
 - **Use Leading Values:** Selected
 - **Normalize String:** Selected
 - **Format Output to:** None
 - **Language:** English
 - **Enable Speller:** Selected
 - **Disable Syntax Error Algorithms:** Not selected

12. Click **Add a column mapping**, and create the following mappings for the following existing domains:

Source Column	Domain
Gender (String)	Gender
StateProvinceName (String)	State
CountryRegionCode (String)	Country/Region (two-letter leading)
CountryRegionName (String)	Country/Region

13. Click **Next**.
14. On the **Matching Policy** page, click **Create a matching rule**.
15. Create a matching rule with the following properties:
- **Rule name:** Is Same Customer
 - **Description:** Checks for duplicate customer records.
 - **Min. matching score:** 80
 - **Rule Editor.** Create rules by clicking **Add a new domain** element, set the following weightings, and then click **Next**:

Domain	Similarity	Weight	Prerequisite
Country/Region	Exact		Selected
Gender	Exact	10	Clear
City	Exact	20	Clear
EmailAddress	Exact	30	Clear
FirstName	Similar	10	Clear
LastName	Similar	10	Clear
AddressLine1	Similar	20	Clear

16. On the **Matching Results** page, click **Start**.
17. When the matching process has finished, review the matches found by DQS, noting that there are duplicate records for three customers.
18. Click **Finish**.
19. When prompted to publish the knowledge base, click **Publish**.
20. When notified that the knowledge base has been published successfully, click **OK**.

Results: After this exercise, you should have created a matching policy and published the knowledge base.

Exercise 2: Using a DQS Project to Match Data

► Task 1: Create a Data Quality Project for Matching Data

1. In SQL Server Data Quality Services, in the **Data Quality Projects** pane, click **New Data Quality Project**.
2. Enter the following details, and then click **Next**.
 - **Name:** Deduplicate Customers
 - **Description:** Identify customer matches.
 - **Use knowledge base:** Customer KB
 - **Select Activity:** Matching
3. On the **Map** page in the **Data Source** list select **SQL Server**, in the **Database** list select **Staging**, and in the **Table/View** list, select **Customers**.
4. In the **Mappings** table, add the following mappings, and then click **Next**:

Source Column	Domain
FirstName (nvarchar)	FirstName
LastName (nvarchar)	LastName
Gender (nvarchar)	Gender
AddressLine1 (nvarchar)	AddressLine1
City (nvarchar)	City
CountryRegionName (nvarchar)	Country/Region
EmailAddress (nvarchar)	EmailAddress

5. On the **Matching** page, click **Start**.

6. When the matching process has completed, review the matches that have been found, and then click **Next**.
7. On the **Export** page, in the **Destination Type** drop-down list, select **Excel File**.
8. Select the **Matching Results** and **Survivorship Results** check boxes, and then specify the following file names:
 - **Matching Results:** D:\Labfiles\Lab10\Starter\Matches.xls
 - **Survivorship Results:** D:\Labfiles\Lab10\Starter\Survivors.xls
9. In the **Survivorship Rule** options, select **Most complete record**, and then click **Export**. Select **Yes** if prompted to overwrite existing files.
10. When the file download is complete, in the **Matching Export** dialog box, click **Close**.
11. Click **Finish**, and then close SQL Server Data Quality Services.

► Task 2: Review and Apply Matching Results

1. In the **D:\Labfiles\Lab10\Starter** folder, double-click **Matches.xls** to open the file in Excel.
2. Note that the matching process found a match with a score of 90 for the following customer records:
 - **CustomerBusinessKey:** 29484 (Rob Turner).
 - **CustomerBusinessKey:** 29261 (Robert Turner).
3. In the **D:\Labfiles\Lab10\Starter** folder, double-click **Survivors.xls** to open the file in Excel.
4. Note that the survivors file contains all the records that should survive deduplication, based on the matches found. It contains the record for customer 29261 (Robert Turner), but not for 29484 (Rob Turner).
5. Close all instances of **Excel** without saving any changes.
6. Switch to the 20767C-MIA-SQL virtual machine.
7. In the **D:\Labfiles\Lab10\Starter** folder, double-click **Fix.Duplicates.sql** to open the file in **SQL Server Management Studio**.
8. When prompted, connect to the **MIA-SQL** instance of the database engine by using Windows authentication.
9. Review the Transact-SQL code and note that it performs the following tasks:
 - It updates the **InternetSales** table so that all sales currently associated with the duplicate customer record become associated with the surviving customer record.
 - It deletes the duplicate customer record.
10. Click **Execute**.
11. When the query has completed, close SQL Server Management Studio without saving any changes.

Results: After this exercise, you should have deduplicated data using a matching project and updated data in your database to reflect these changes.

Module 11: Master Data Services

Lab: Implementing Master Data Services Model

Exercise 1: Creating a Master Data Services Model

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC**, **20767C-MIA-CLI**, and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab11\Starter** folder, right-click **Setup.cmd**, and then click **Run as administrator**.
3. In the **User Account Control** dialog box, click **Yes**.

► Task 2: Create a Model

1. Start Internet Explorer, and navigate to **http://mia-sql:81/MDS**
2. On the **Master Data Services** home page, click **System Administration**.
3. On the **Manage Models** page, click **Add**.
4. In the **Add Model** pane, in the **Name** box, type **Product Catalog**, clear the **Create entity with same name as model** check box, and then click **Save**.

► Task 3: Create Entities

1. On the **Manage Models** page, click the **Product Catalog** model, and then click **Entities**.
2. On the **Manage Entities** page, click **Add**.
3. In the **Add Entity** pane, in the **Name** box, type **Product**, and then click **Save**.
4. On the **Manage Entities** page, click **Add**.
5. In the **Add Entity** pane, in the **Name** box, type **Product_Subcategory**, and then click **Save**.

► Task 4: Create Attributes

1. On the **Manage Entities** page, in the **Entity** table, click **Product**, and then click **Attributes**.
2. On the **Manage Attributes** page, click **Add**.
3. In the **Add Attribute** pane, create an attribute with the following settings:
 - **Name:** Product_Subcategory_Key
 - **Attribute Type:** Domain-based
 - **Domain Entity:** Product_Subcategory
4. Leave all other settings at their defaults, and then click **Save**.
5. On the **Manage Attributes** page, click **Add**.

6. In the **Add Attribute** pane, create an attribute with the following settings:
 - **Name:** Description
 - **Attribute Type:** Free-form
 - **Data Type:** Text
 - **Length:** 400
7. Leave all other settings at their defaults, and then click **Save**.
8. On the **Manage Attributes** page, click **Add**.
9. In the **Add Attribute** pane, create an attribute with the following settings:
 - **Name:** ListPrice
 - **Attribute Type:** Free-form
 - **Data type:** Number
 - **Decimals:** 4
 - **Input Mask:** (####)
10. Leave all other settings at their defaults, and then click **Save**.

► Task 5: Add Members

1. In Internet Explorer, click the **SQL Server Master Data Services** text to return to the Master Data Manager home page.
2. In the **Model** list, click **Product Catalog**, and then click **Explorer**.
3. On the **Entities** menu, click **Product_Subcategory**.
4. Click **Add Member**, and in the Details pane, enter the following attribute values:
 - **Name:** Mountain Bikes
 - **Code:** 1
5. In the **Annotations** box, type **New Subcategory**, and then click **OK**.
6. Repeat the previous two steps to add the following **Product Subcategory** members:

Name	Code
Chains	7
Gloves	20
Helmets	31

7. On the **Entities** menu, click **Product**.
8. Click **Add Member**, and in the Details pane, enter the following values:
 - **Name:** Mountain-100 Silver, 42
 - **Code:** 345
 - **Product_Subcategory_key:** 1

- **Description:** Competition Mountain Bike
 - **ListPrice:** 3399.99
9. In the **Annotations** box, type **New product**, and then click **OK**.
 10. Repeat the previous two steps to add the following **Product** members:

Name	Code	Product_Subcategory_key	Description	ListPrice
Chain	559	7	Superior Chain	20.24
Full Finger Gloves, S	468	20	Synthetic Gloves	37.99
Sport-100 Helmet, Red	214	31	Universal Fit Helmet	34.99

11. Close Internet Explorer.

Results: After this exercise, you should have a Master Data Services model named **Product Catalog** that contains **Product** and **ProductSubcategory** entities. Each of these entities should contain four members.

Exercise 2: Using the Master Data Services Add-in for Excel

► Task 1: Add Free-Form Attributes to an Entity

1. Log into the **20767C-MIA-CLI** virtual machine as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**.
2. Start Microsoft Excel and create a new blank workbook.
3. On the **File** tab, click **Options**.
4. In the **Excel Options** dialog box, on the **Add-ins** tab, in the **Manage** drop-down list, select **COM Add-ins**, and then click **Go**.
5. In the **COM Add-ins** dialog box, if **Master Data Services Add-in for Excel** is not selected, select it, and then click **OK**.
6. In Excel, on the ribbon, click the **Master Data** tab.
7. In the **Connect and Load** section, click the drop-down arrow under **Connect**, and then click **Manage Connections**.
8. On the ribbon, in the **Manage Connections** dialog box, click **New**.
9. In the **Add New Connection** dialog box, in the **Description** box, type **Local MDS Server**, in the **MDS server address** box, type **http://mia-sql:81/mds**, and then click **OK**.
10. In the **Manage Connections** dialog box, click **Close**.
11. On the ribbon, in the **Connect and Load** area, click the drop-down arrow under **Connect**, and then click **Local MDS Server**.

12. In the **Master Data Explorer** pane, in the **Model** drop-down list, click **Product Catalog**.
13. In the list of entities, click **Product**.
14. On the ribbon, in the **Connect and Load** area, click **Refresh**. The Product members you created in the previous exercise are displayed in a worksheet named **Product**.
15. In cell **S2** (to the right of the **ListPrice** attribute header), type **StandardCost** and press Enter.
16. Click cell **S2** (in which you just entered **StandardCost**), and on the ribbon, in the Master Data tab in the **Build Model** area, click **Attribute Properties**.
17. In the **Attribute Properties** dialog box, in the **Attribute type** list, select **Number**, in the **Decimal places** box, type **4**, and then click **OK**. This creates a new free-form attribute named **StandardCost** and publishes it to the model.
18. Repeat the previous three steps to create the following free-form attributes:

Cell	Attribute Name	Attribute Properties
T2	Color	<ul style="list-style-type: none"> • Text • Maximum length: 15
U2	SafetyStockLevel	<ul style="list-style-type: none"> • Number • Decimal places: 0
V2	ReorderPoint	<ul style="list-style-type: none"> • Number • Decimal places: 0
W2	Size	<ul style="list-style-type: none"> • Text • Maximum length: 50
X2	Weight	<ul style="list-style-type: none"> • Number • Decimal places: 4
Y2	DaysToManufacture	<ul style="list-style-type: none"> • Number • Decimal places: 0
Z2	ModelName	<ul style="list-style-type: none"> • Text • Maximum length: 500

19. In cell **S3**, (the **StandardCost** value for the Sport-100 Helmet, Red product), enter **13.0863**.
20. In cell **T3** (the **Color** value for the Sport-100 Helmet, Red product), enter **Red**.
21. In cell **U3**, (the **SafetyStockLevel** value for the Sport-100 Helmet, Red product), enter **3**.
22. In cell **V3**, (the **ReorderPoint** value for the Sport-100 Helmet, Red product), enter **4**.
23. In cell **W3** (the **Size** value for the Sport-100 Helmet, Red product), enter **U**.
24. In cell **X3**, (the **Weight** value for the Sport-100 Helmet, Red product), enter **0.2000**.
25. In cell **Y3**, (the **DaysToManufacture** value for the Sport-100 Helmet, Red product), enter **5**.
26. In cell **Z3**, (the **ModelName** value for the Sport-100 Helmet, Red product), enter **Sport-100**.

27. Enter the following attribute values for the remaining products:

- Full-Finger Gloves, S
 - **StandardCost:** 15.6709
 - **Color:** Black
 - **SafetyStockLevel:** 4
 - **ReorderPoint:** 3
 - **Size:** S
 - **Weight:** 0.2000
 - **DaysToManufacture:** 0
 - **ModelName:** Full-Finger Gloves
- Chain
 - **StandardCost:** 8.9866
 - **Color:** Silver
 - **SafetyStockLevel:** 500
 - **ReorderPoint:** 375
 - **Size:** (leave blank)
 - **Weight:** (leave blank)
 - **DaysToManufacture:** 1
 - **ModelName:** Chain
- Mountain-100 Silver, 42
 - **StandardCost:** 1912.1544
 - **Color:** Silver
 - **SafetyStockLevel:** 100
 - **ReorderPoint:** 75
 - **Size:** L
 - **Weight:** 20.77
 - **DaysToManufacture:** 4
 - **ModelName:** Mountain-100

28. On the ribbon, in the **Publish and Validate** area, click **Publish**.

► **Task 2: Create an Entity for a Domain-Based Attribute**

1. On the ribbon, in the **Connect and Load** area, click the drop-down arrow under **Connect**, and then click **Local MDS Server**.
2. In the **Master Data Explorer** pane, in the **Model** list, select **Product Catalog**, in the list of entities, click **Product_Subcategory**.
3. On the ribbon, in the **Connect and Load** area, click **Refresh**. The **Product_Subcategory** members you created in the previous exercise are displayed in a worksheet named **Product_Subcategory**.

4. In cell **P2** (to the right of the **Code** attribute header), enter **Product_Category_Key**
5. In cell **P6**, enter the value **2** as the **Product_Category_Key** for the **Chains** subcategory.
6. Enter the following **Product_Category_Key** values for the remaining subcategories:
 - o **Gloves:** 3
 - o **Helmets:** 4
 - o **Mountain Bikes:** 1
7. Click cell **P2** (the header for the **Product_Category_Key** attribute), and on the ribbon, in the **Build Model** area, click **Attribute Properties**.
8. In the **Attribute Properties** dialog box, in the **Attribute type** list, click **Constrained list (Domain-based)**. In the **Populate the attribute with values from** list, ensure that **the selected column** is selected, in the **New entity name** box, type **Product_Category**, and then click **OK**.
9. On the ribbon, in the **Connect and Load** area, click the drop-down arrow under **Connect**, and then click **Local MDS Server**.
10. In the Master Data Explorer pane, in the **Model** drop-down list, click **Product Catalog**, in the list of entities, click **Product_Category**, and then on the ribbon, in the **Connect and Load** area, click **Refresh**. The newly created **Product_Category** entity members are displayed in a worksheet named **Product_Category**.
11. In cell **N3** (the **Name** value for the first member, which currently contains the value **1**), enter **Bikes**.
12. Enter the following **Name** values for the remaining categories:

Name	Code
Bikes	1
Components	2
Clothing	3
Accessories	4

13. On the ribbon, in the **Publish and Validate** area, click **Publish**. Note: Ensure you have not clicked inside the columns when clicking publish, otherwise the operation will fail.
14. Click the tab for the **Product_Subcategory** worksheet, and on the ribbon, in the **Connect and Load** area, click **Refresh**.
Note that the **Product_Category** member names are displayed in the **Product_Category_Key** column.
15. Minimize Excel. You will return to it in a later exercise.

Results: After this exercise, you should have a master data model that contains **Product**, **ProductSubcategory**, and **ProductCategory** entities that contain data entered in Excel.

Exercise 3: Enforcing Business Rules

► Task 1: Create a Rule for the ListPrice Attribute

1. Start Internet Explorer, and navigate to **http://mia-sql:81/MDS**.
2. On the **Master Data Services** home page, click **System Administration**.
3. On the **Manage** menu, click **Business Rules**.
4. On the **Business Rules** page, in the **Model** list, ensure that **Product Catalog** is selected, in the **Entity** list, ensure that **Product** is selected and in the **Member Types** list, ensure that **Leaf** is selected, and then click **Add**.
5. In the **Add Business Rule** window, in the **Name** box, type **Validate Price**, in the **Description** box, type **Check that prices are non-zero**.
6. Under **Then**, click the **Add** link.
7. In the **Create Action** window, in the **Attribute** drop-down list, select **ListPrice**, in the **Operator** list select **must be greater than**, in the **Must be greater than** list select **Attribute value**, in the **Attribute value** box, type **0**, and then click **Save**.
8. In the **Add Business Rule** window, click **Save**.

► Task 2: Create a Rule for the SafetyStockLevel Attribute

1. On the **Business Rules** page, click **Add**.
2. In the **Add Business Rule** window, in the **Name** box, type **Validate SafetyStockLevel**, in the **Description** box, type **Check that safety stock level is greater than reorder point for products that take a day or more to manufacture**, and then under **If** click **Add**.
3. In the **Create Condition** window, in the **Attribute** drop-down list, click **DaysToManufacture**, in the **Operator** drop-down list, click **is greater than**, in the **Is greater than** drop-down list, click **Attribute value**, in the **Attribute value** box, type **0**, and then click **Save**.
4. In the **Add Business Rule** window, under **Then**, click **Add**.
5. In the **Create Action** window, in the **Attribute** drop-down list, click **SafetyStockLevel**, in the **Operator** drop-down list, click **must be greater than**, in the **Must be greater than** list, click **Attribute**, in the **Attribute** drop-down list select **ReorderPoint** and then click **Save**.
6. In the **Add Business Rule** window, click **Save**.

► Task 3: Publish Business Rules and Validate Data

1. On the **Business Rules** page, click **Publish All**.
2. In the **Message from webpage** dialog box, click **OK**.
3. In the **Business Rule State** column, check that the value displayed for both rules is **Active**.
4. Click the **SQL Server** text in the page title to return to the home page.
5. In the **Model** list, click **Product Catalog**.
6. Click **Explorer**.
7. If the Microsoft Silverlight window appears, complete the following steps:
 - a. Click **Click now to install**.
 - b. In the Internet Explorer message bar, click **Run**.

- c. In the **User Account Control** dialog box, click **Yes**.
 - d. In the **Install Silverlight** dialog box, click **Install now**.
 - e. In the **Enable Microsoft Update** dialog box, click **Next**.
 - f. On the **Installation successful** page, click **Close**.
 - g. Press F5.
8. On the **Entities** menu, click **Product** to ensure that you are viewing the **Product** entity members, and then click **Apply Rules**. Note that the **Sport-100 Helmet, Red** member has a red exclamation mark, indicating that a business rule has been violated.
 9. Click the **Sport-100 Helmet, Red** member, and at the bottom of the **Details** section, note the validation error.
 10. In the **Details** tab, change the **SafetyStockLevel** attribute for the **Sport-100 Helmet, Red** member to **5**, and then click **OK**. Note that the validation error and red exclamation mark disappear.
 11. Close Internet Explorer.
 12. Return to Excel, and click the **Product** worksheet tab.
 13. On the **Master Data** tab of the ribbon, in the **Connect and Load** area, click **Refresh**.
 14. On the ribbon, on the **Master Data** tab, in the **Publish and Validate** area, click **Show Status** to display the **\$ValidationStatus\$** and **\$InputStatus\$** columns.
 15. In the **Publish and Validate** area, click **Apply Rules**. Note that validation succeeds for all members.
 16. Change the value in the **ListPrice** column for the **Chain** product to **0.00**.
 17. On the ribbon, in the **Publish and Validate** section, click **Publish**. Note: Ensure you have not clicked inside the columns when clicking publish, otherwise the operation will fail.
 18. In the **Publish and Annotate** dialog box, click **Publish**. Note that validation failed for this member.
 19. Change the value in the **ListPrice** column for the **Chain** product back to **20.24**.
 20. On the ribbon, in the **Publish and Validate** section, click **Publish**. Note: Ensure you have not clicked inside the columns when clicking publish, otherwise the operation will fail.
 21. Close Excel without saving the workbook.

Results: After this exercise, you should have a master data model that includes business rules to validate product data.

Exercise 4: Loading Data into a Model

► Task 1: Load Data into the Model

1. On the **20767C-MIA-SQL** virtual machine, start SQL Server Management Studio and connect to the **MIA-SQL** database engine instance using Windows authentication.
2. Open the **Import Products into MDS.sql** file in the **D:\Labfiles\Lab11\Starter** folder.
3. Review the Transact-SQL code in this script. Note that it performs the following tasks:
 - Generates a unique batch ID.
 - Inserts data into the **stg.Product_Category_Leaf** staging table from the **ProductCategory** table in the **Products** database, specifying an **ImportType** value of 2, an **ImportStatus_ID** value of 0, and a **BatchTag** value based on the batch ID generated previously.
 - Inserts data into the **stg.Product_Subcategory_Leaf** staging table from the **ProductSubcategory** table in the **Products** database, specifying an **ImportType** value of 2, an **ImportStatus_ID** value of 0, and a **BatchTag** value based on the batch ID generated previously.
 - Inserts data into the **stg.Product_Leaf** staging table from the **Product** table in the **Products** database, specifying an **ImportType** value of 2, an **ImportStatus_ID** value of 0, and a **BatchTag** value based on the batch ID generated previously.
 - Executes the **stg.udp_Product_Category_Leaf** stored procedure to start processing the staged **Product Category** members with the batch tag specified previously.
 - Executes the **stg.udp_Product_Subcategory_Leaf** stored procedure to start processing the staged **Product Subcategory** members with the batch tag specified previously.
 - Executes the **stg.udp_Product_Leaf** stored procedure to start processing the staged **Product** members with the batch tag specified previously.
4. Click **Execute** and note the numbers of rows affected by the statements.

► Task 2: Check the Status of the Data Load

1. Start Internet Explorer, and navigate to **http://mia-sql:81/mds**
2. On the **Master Data Services** home page, click **Integration Management**, and then click **Import Data**.
3. Review the information on the **Integration** page. In the **Status** column, check that the value displayed is **Completed**, and in the **Errors** column, check that the value displayed is **0**.

► Task 3: Validate Imported Members

1. Click the **SQL Server** text in the page title to return to the home page.
2. Click **Explorer**.
3. On the **Entities** menu, click **Product**.
4. Click **Apply Rules** to apply business rules to the imported members.
5. Click **Filter**, in the **Attribute** column, select **[Validation Status]**, in the **Operator** column, ensure that **Is equal to** is selected, and in the **Criteria** column, select **Validation failed**, and then click **Apply**. The list is filtered to show only invalid members.
6. Click each member and review the validation errors at the bottom of the Details pane. Some members have failed because they have no **ListPrice** value. Resolve this problem by setting the **ListPrice** of each invalid member to **1431.50**.

7. Click **Filter**, and then click **Clear Filter**.
8. Click the **SQL Server** text in the page title to return to the home page.
9. Close Internet Explorer

Results: After this exercise, you should have loaded members into the Master Data Services model.

Module 12: Extending SQL Server Integration Services

Lab: Using Custom Scripts

Exercise 1: Using a Script Task

► Task 1: Prepare the Lab Environment

1. Ensure that the **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are both running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab12\Starter** folder, right-click **Setup.cmd**, and then click **Run as administrator**.
3. Click **Yes** when prompted to confirm that you want to run the command file.
4. When prompted, press **y**, then press **Enter**, and then wait for the script to finish.

► Task 2: Add a Script Task to a Control Flow

1. Start Visual Studio and open the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab12\Starter** folder.
2. In Solution Explorer, double-click the **Load DW.dtsx** SSIS package to open it in the designer.
3. Review the control flow, noting that it runs multiple packages to load data from the staging database into the data warehouse. It then uses an Execute SQL task named **Get record counts and truncate Staging tables** to determine the number of rows that were processed before truncating the staging tables, ready for the next refresh cycle.
4. On the **SSIS** menu, click **Variables**, and note the variables in this package. **The Get record counts and truncate Staging tables** task assigns the row counts in the staging tables to these variables.
5. Close the Variables window.
6. On the **SSIS** menu, click **SSIS Toolbox**.
7. Add a **Script Task** to the control flow, under the **Get record counts and truncate Staging tables** task.
8. Right-click **Script Task**, click **Rename**, and change the name to **Log Rowcounts**.
9. Click the **Get record counts and truncate Staging tables** task, and connect the green precedence constraint to the **Log Rowcounts** task.

► Task 3: Enable Logging for the Script Task

1. On the **SSIS** menu, click **Logging**.
2. On the **Providers and Logs** tab, in the **Provider type** list, select **SSIS log provider for Windows Event Log**, and then click **Add**.
3. In the **Containers** tree, select the **Log Rowcounts** check box, and in the **Providers and Logs** tab, select the **SSIS log provider for Windows Event Log** check box.
4. On the **Details** tab, select the **ScriptTaskLogEntry** check box, and then click **OK**.

► Task 4: Configure the Script Task

1. On the control flow surface, double-click the **Log Rowcounts** script task.
2. In the **Script Task Editor** dialog box, on the **Script** page, in the **ReadOnlyVariables** box, click the ellipsis (...) button.
3. In the **Select Variables** dialog box, select the following variables, and then click **OK**:
 - **User::CustomerCount**
 - **User::EmployeeCount**
 - **User::InternetSalesCount**
 - **User::PaymentCount**
 - **User::ResellerCount**
 - **User::ResellerSalesCount**
4. Click **Edit Script**, and wait for the Visual Studio VstaProjects editor to open.
5. In the **VstaProjects editor**, in the **Main()** function, replace the comment **//TODO: Add your code here** with the following code (above the existing **Dts.TaskResult = (int)ScriptResults.Success** statement). You can copy and paste this code from **Script.txt** in the **D:\Labfiles\Lab12\Starter** folder.

```
String logEntry = "Data Warehouse records loaded (" ;
logEntry += Dts.Variables["User::CustomerCount"].Value.ToString() + " customers, ";
logEntry += Dts.Variables["User::ResellerCount"].Value.ToString() + " resellers, ";
logEntry += Dts.Variables["User::EmployeeCount"].Value.ToString() + " employees, ";
logEntry += Dts.Variables["User::PaymentCount"].Value.ToString() + " payments, ";
logEntry += Dts.Variables["User::InternetSalesCount"].Value.ToString() + " Internet
sales, and ";
logEntry += Dts.Variables["User::ResellerSalesCount"].Value.ToString() + " reseller
sales) ";
Dts.Log(logEntry, 999, null);
```

6. Save the script and close the VstaProjects – Microsoft Visual Studio window.
7. In the **Script Task Editor** dialog box, click **OK**.

► Task 5: Test the Script

1. On the **Debug** menu, click **Start Debugging** and observe the control flow tab as the package executes, noting that each package executed opens in a new window. The entire load process can take a few minutes.
2. When execution is complete, on the **Debug** menu, click **Stop Debugging**.
3. Close Visual Studio, saving your work if prompted.
4. Right-click the **Start** button, and then click **Event Viewer**.
5. In **Event Viewer**, expand the **Windows Logs** folder and select the **Application** log.
6. In the list of application events, select the second information event with a source of **SQLISPackage120** and read the information in the **Details** tab.
7. Close Event Viewer.

Results: After this exercise, you should have an SSIS package that uses a script task to log row counts.

Module 13: Deploying and Configuring SSIS Packages

Lab: Deploying and Configuring SSIS Packages

Exercise 1: Creating an SSIS Catalog

► Task 1: Prepare the Lab Environment

1. Ensure that the **MT17B-WS2016-NAT**, **20767C-MIA-DC** and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab13\Starter** folder, right-click **Setup.cmd** and click **Run** as administrator.
3. Click **Yes** when prompted to confirm you want to run the command file, and wait for the script to finish.
4. If a window opens with a message, "Do you want to continue this operation?" type **Y** and then press Enter.

► Task 2: Create the SSIS Catalog and a Folder

1. Start SQL Server Management Studio and connect to the **MIA-SQL** instance of the database engine using **Windows authentication**.
2. In Object Explorer, right-click **Integration Services Catalogs**, and then click **Create Catalog**.
3. In the **Password** box, type **Pa55w.rd**, in the **Retype Password** box type **Pa55w.rd**, and then click **OK**.
4. In Object Explorer, expand **Integration Services Catalogs**.
5. Right-click **SSISDB**, and then click **Create Folder**.
6. In the **Create Folder** dialog box, in the **Folder name** box, type **DW ETL**, in the **Folder description** box, type **Folder for the Adventure Works ETL SSIS Project**, and then click **OK**.

Results: A SQL Server instance with an SSIS Catalog.

Exercise 2: Deploying an SSIS Project

► Task 1: Deploy an SSIS Project

1. Start Visual Studio.
2. On the **File** menu, point to **Open**, click **Project/Solution**, select to the **AdventureWorksETL.sln** solution in the **D:\Labfiles\Lab13\Starter** folder, and then click **Open**.
3. On the **Build** menu, click **Build Solution**.
4. When the build has succeeded, on the **Project** menu, click **Deploy**.
5. In the **Integration Services Deployment Wizard** dialog box, on the **Introduction** page, click **Next**.
6. On the **Select Destination** page, in the **Server name** box, type **MIA-SQL**, and then click **Connect**.
7. In the **Path** box, click **Browse**, browse to the **SSISDB\DW ETL** folder you created earlier, and then click **OK**.
8. On the **Select Destination** page, click **Next**.
9. On the **Review** page, click **Deploy**.
10. When deployment has completed, click **Close**, and then close Visual Studio.

Results: After this exercise, you should have deployed an SSIS project to a folder in your SSIS database.

Exercise 3: Creating Environments for an SSIS Solution

► Task 1: Create Environments

1. In SQL Server Management Studio, expand the DW ETL folder you created earlier, and expand the **Projects** folder. Note that the AdventureWorksETL project has been deployed.
2. Right-click the **Environments** folder, and then click **Create Environment**.
3. In the **Create Environment** dialog box, enter the environment name **Test**, and then click **OK**.
4. Repeat the previous two steps to create a second environment named **Production**.

► Task 2: Create Variables

1. Expand the **Environments** folder, right-click the **Production** environment, and then click **Properties**.
2. In the **Environment Properties** dialog box, on the **Variables** page, add a variable with the following settings:
 - **Name:** StgServer
 - **Type:** String
 - **Description:** Staging server
 - **Value:** MIA-SQL
 - **Sensitive:** No

3. Add a second variable with the following settings (making sure to include the trailing “\” in the value), and then click **OK**:
 - **Name:** FolderPath
 - **Type:** String
 - **Description:** Accounts folder
 - **Value:** D:\Accounts\
 - **Sensitive:** No
4. Right-click the **Test** environment, and then click **Properties**.
5. In the **Environment Properties** dialog box, on the **Variables** page, add a variable with the following settings:
 - **Name:** StgServer
 - **Type:** String
 - **Description:** Staging server
 - **Value:** MIA-SQL
 - **Sensitive:** No
6. Add a second variable with the following settings (making sure to include the trailing “\” in the value), and then click **OK**:
 - **Name:** FolderPath
 - **Type:** String
 - **Description:** Accounts folder
 - **Value:** D:\TestAccts\
 - **Sensitive:** No

► Task 3: Map Environment Variables

1. In the **Projects** folder, right-click **AdventureWorksETL**, and then click **Configure**.
2. In the **Configure – AdventureWorksETL** dialog box, on the **References** page, click **Add** and add the **Production** environment.
3. Click **Add** again and add the **Test** environment.
4. On the **Parameters** page, in the **Scope** list, select **AdventureWorksETL**.
5. On the **Parameters** tab, click the ellipses (...) button for the **AccountsFolderPath** parameter.
6. In the **Set Parameter Value** dialog box, select **Use environment variable**, select **FolderPath** in the list of variables, and then click **OK**.
7. In the **Configure – AdventureWorksETL** dialog box, on the **Parameters** page, on the **Connection Managers** tab, select the **localhost Staging ADO NET** connection manager, and then click the ellipses (...) button for the **ServerName** property.
8. In the **Set Parameter Value** dialog box, select **Use environment variable**, select **StgServer** in the list of variables, and then click **OK**.
9. In the **Configure – AdventureWorksETL** dialog box, click **OK**.

Results: After this exercise, you should have configured your project to use environments to pass values to package parameters.

Exercise 4: Running an SSIS Package in SQL Server Management Studio

► Task 1: Run a Package

1. In Object Explorer, expand the **AdventureWorksETL** project, and then expand **Packages**.
2. Right-click **Extract Payment Data.dtsx**, and then click **Execute**.
3. Select **Environment**, and then select the **.\Test** environment.
4. Click **OK**.
5. When prompted to open **Overview** report, wait 10 seconds for the package to run, and then click **Yes**.
6. In the **Overview** report, in the **Execution Information** table, note the environment used.
7. In the **Parameters Used** table, note the values you configured with environment variables.

Results: After this exercise, you should have performed tests to verify that a value was passed from the environment to the package.

Exercise 5: Scheduling SSIS Packages with SQL Server Agent

► Task 1: Create a SQL Server Agent Job

1. In Object Explorer, right-click **SQL Server Agent**, point to **New**, and click **Job**.
2. In the **New Job** dialog box, in the **Name** box, type **Extract Reseller Data**.
3. In the **New Job** dialog box, on the **Steps** page, click **New**.
4. In the **New Job Step** dialog box, in the **Step name** box, type **Run Extract Reseller Data Package**.
5. In the **Type** list, select **SQL Server Integration Services Package**.
6. On the **Package** tab, in the **Package source** list, select **SSIS Catalog**.
7. On the **Package** tab, in the **Server** box, type **MIA-SQL**.
8. Click the ellipses (...) button for the **Package** box, and in the **Select an SSIS Package** dialog box, expand **SSISDB**, expand **DW ETL**, expand **AdventureWorksETL**, select **Extract Reseller Data.dtsx**, and then click **OK**.
9. In the **New Job Step** dialog box, on the **Configuration** tab, select **Environment**, and select the **.\Production** environment.
10. In the **New Job Step** dialog box, on the **Advanced** page, in the **On success action** list, select **Quit the job reporting success**. Then click **OK**.
11. In the **New Job** dialog box, on the **Schedules** page, click **New**.
12. In the **Name** box, type **ETL Schedule**.
13. In the **Schedule type** list, select **One time**.

14. Ensure **Enabled** is selected.
15. In the **One-time occurrence** section, select the current day and enter a time that is two minutes later than the current time.
16. Click **OK** to close the **New Job Schedule** dialog box, and then click **OK** to close the **New Job** dialog box.
17. Wait for two minutes.
18. In Object Explorer, expand **SQL Server Agent**, right-click **Job Activity Monitor**, and click **View Job Activity**.
19. In the Job Activity Monitor – localhost window, note the **Status** and **Last Run Outcome** values for the **Extract Reseller Data** job. If the job has not yet completed, wait a minute and click **Refresh** until the job has completed successfully, and then click **Close**.

► **Task 2: View the Integration Services Dashboard**

1. In SQL Server, in Object Explorer, under **Integration Services Catalogs**, right-click **SSISDB**, point to **Reports**, point to **Standard Reports**, and click **Integration Services Dashboard**.
2. In the **Packages Detailed Information (Past 24 Hours)** list, notice that the most recent two package executions succeeded. Click the **Overview** link for the first package in the list—it should be the **Extract Reseller Data.dtsx** package.
3. In the **Overview** report, in the **Execution Information** table, note the environment that was used. In the **Parameters Used** table, note the values you configured with environment variables.
4. Click the **Navigate Backward** button at the top left of the Overview pane.
5. In the **Packages Detailed Information (Past 24 Hours)** list, click the **All Messages** link for the first package in the list. Look at the messages logged during package execution.
6. In the **Execution Information** table, note the execution **Duration**, **Start Time** and **End Time**. In the **Error Messages** table, note that the list is empty. Finally, note that in the **All Messages** table a large number of messages is listed; this is because running an SSIS package involves a large number of steps, many of which the user may be unaware of.
7. Close SQL Server Management Studio.

Results: After this exercise, you should have created a SQL Server Agent job that automatically runs your SSIS package.

MCT USE ONLY. STUDENT USE PROHIBITED

Module 14: Consuming Data in a Data Warehouse

Lab: Using a Data Warehouse

Exercise 1: Exploring a Reporting Services Report

► Task 1: Prepare the Lab Environment

1. Ensure that the **MT17B-WS2016-NAT**, **20767C-MIA-DC**, **20767C-MIA-CLI**, and **20767C-MIA-SQL** virtual machines are all running, and then log on to **20767C-MIA-SQL** as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**
2. In the **D:\Labfiles\Lab14\Starter** folder, right-click **Setup.cmd**, and then click **Run as administrator**.
3. In the **User Account Control** dialog box, click **Yes**.

► Task 2: Publish a Reporting Services Report

1. Switch to the **20767C-MIA-CLI** virtual machine, and log on as **ADVENTUREWORKS\Student** with the password **Pa55w.rd**.
2. On the Windows taskbar, click **Search Windows**, type **Report Builder**, and then click **Report Builder**.
3. Close the **Getting Started** wizard.
4. On the **File** menu, click **Open**.
5. Browse to **D:\Labfiles\Lab14\Starter**, click **Internet Sales by State.rdl**, and then click **Open**.
6. Click **Run**, and wait for the report to be generated.
7. Click **Export**, and then click **Excel**.
8. Browse to **D:\Labfiles\Lab14\Starter**, and then click **Save**. Click **Yes** to overwrite the existing report file.
9. Close Report Builder.
10. On the taskbar, click **File Explorer**.
11. Browse to **D:\Labfiles\Lab14\Starter**, and then double-click **Internet Sales by State.xlsx**.
12. In the Excel report, highlight the header row of the table.
13. On the ribbon, in the **Editing** group, click **Sort & Filter**, and then click then **Filter**.
14. Click the filter icon in the **Product Subcategory** column, then click **(Select All)** to clear all values. Click **Road Bikes**, and then click **OK**.
15. Review the worksheet, save the workbook, and then close Excel.

Results: After this exercise, you should have deployed a Reporting Services report to the server, rendered the report, and exported the data to Excel.

Exercise 2: Exploring a PowerPivot Workbook

► Task 1: View a PowerPivot Report

1. Start Excel and open the **Analysis.xlsx** workbook in the **D:\Labfiles\Lab14\Starter** folder. In the **PROTECTED VIEW** section, click **Enable Editing**. In the **SECURITY WARNING** section, click **Enable Content**.
2. On the ribbon, click **File**, and then click **Options**.
3. In the **Excel Options** dialog box, on the **Add-ins** page, in the **Manage** drop-down list, select **COM Add-Ins**, and then click **Go**.
4. In the **COM Add-ins** dialog box, select **Microsoft Power Pivot for Excel**, and then click **OK**.
5. On the ribbon, on the **Analyze** tab, click **Insert Slicer**.
6. In the **Insert Slicers** dialog box, under **FactResellerSales**, select **ShipDate**, and then click **OK**. **ShipDate** selector appears, listing dates by month.



Note: This slicer only applies to the **FactResellerSales** table. If you also want to filter by Internet sales dates, then you must add another slicer for the **FactInternetSales** table.

7. In the **ShipDate** selector, click **11/8/2005** and note that data changes to show sales from November 2005.
8. On the **File** menu, click **Save**.
9. On the **File** menu, click **Save As**, and then click **Browse**.
10. In the **Save As** dialog box, in the **Filename** box, type **http://mia-sql/sites/adventureworks/**, and then press Enter.
11. Under **Document Libraries**, click **Documents**, click **Open**, and then click **Save**.
12. Close Excel.
13. Open Internet Explorer and navigate to **http://mia-sql/sites/adventureworks/**.
14. In the quick launch area on the left side of the page, click **Documents**
15. In the document list, click **Analysis**.
16. When the worksheet appears, in the menu bar, click **Open In Excel**.
17. Review the report and verify that it functions as before. Note that you may need to enable editing and active content before the slicer works correctly.
18. Close Excel, and then close Internet Explorer.

Results: After this exercise, you should have a customized PowerPivot report in an Excel workbook called **Analysis.xlsx**.

Exercise 3: Exploring a Power View Report

► Task 1: Create a Power View Report

1. Start Excel and open the **Analysis.xlsx** workbook in the **D:\Labfiles\Lab14\Starter** folder.
2. On the ribbon, click **File**, and then click **Options**.
3. In the **Excel Options** dialog box, on the **Add-ins** page, in the **Manage** drop-down list, select **COM Add-ins**, and then click **Go**.
4. In the **COM Add-Ins** dialog box, select **Microsoft Power View for Excel**, and then click **OK**.
5. On the ribbon, click **File**, and then click **Options**.
6. In the **Excel Options** dialog box, click **Customize Ribbon**.
7. Under **Main Tabs** (the list on the right), in the list click **Insert**, and then click **New Group**.
8. Click **Rename**, type **Reports**, and then click **OK**.
9. In the **Choose commands from** list, select **Commands Not in the Ribbon**.
10. Scroll down the list of commands, click **Insert a Power View Report**, click **Add**, and then click **OK**.
11. On the **Insert** tab, in the **Reports** group, click **Power View**. Click **Enable** to enable the feature.
12. If prompted, click **Install Silverlight** and install Silverlight using the default options, and then click **Reload**.
13. Wait for the Power View sheet to be added to the workbook.
14. In the Power View report, close the **Filters** pane.
15. Click **Click here to add a title**, and then type **Internet Sales**.
16. In the **Power View Fields** pane, expand **FactInternetSales** and select **SalesAmount**. Expand **DimProductCategory** and select **EnglishProductCategoryName**. Expand **DimProductSubcategory** and select **EnglishProductSubcategoryName**.
17. In the **Switch Visualization** area of the ribbon, in the **Column Chart** list, select **Stacked Column** then resize the column chart to fill the top half of the report.
18. Click below the chart and then, in the Power View Fields pane, expand **DimGeography** and select **EnglishCountryRegionName**. Expand **FactInternetSales** and select **SalesAmount**.
19. In the chart, click the large blue block in the **Bikes** column. Note that the table updates to show only the results for the sales of Mountain Bikes.
20. Click the brown block above this and note again that the table updates, showing only the results for Road Bikes.
21. Click the brown block again to redisplay the sales figures for all products.
22. Close Excel, saving the changes to the workbook if prompted.

Results: After this exercise, you should have created a Power View report.

MCT USE ONLY. STUDENT USE PROHIBITED