**Microsoft**

**EXAM 70-463**

# Implementing a Data Warehouse with Microsoft SQL Server 2012

Dejan Sarka
Matija Lah
Grega Jerkič

# Training Kit

# How to access
# your CD files

The print edition of this book includes a CD. To access the CD files, go to http://aka.ms/666092/files, and look for the Downloads tab.

Note: Use a desktop web browser, as files may not be accessible from all ereader devices.

Questions? Please contact: mspinput@microsoft.com

## Microsoft Press

# Exam 70-463: Implementing a Data Warehouse with Microsoft SQL Server 2012

| OBJECTIVE | CHAPTER | LESSON |
|---|---|---|
| **1. DESIGN AND IMPLEMENT A DATA WAREHOUSE** | | |
| 1.1  Design and implement dimensions. | Chapter 1 | Lessons 1 and, 2 |
| | Chapter 2 | Lessons 1, 2, and 3 |
| 1.2  Design and implement fact tables. | Chapter 1 | Lesson 3 |
| | Chapter 2 | Lessons 1, 2, and 3 |
| **2.  EXTRACT AND TRANSFORM DATA** | | |
| 2.1  Define connection managers. | Chapter 3 | Lessons 1 and 3 |
| | Chapter 4 | Lesson 1 |
| | Chapter 9 | Lesson 2 |
| 2.2  Design data flow. | Chapter 3 | Lesson 1 |
| | Chapter 5 | Lessons 1, 2, and 3 |
| | Chapter 7 | Lesson 1 |
| | Chapter 10 | Lesson 2 |
| | Chapter 13 | Lesson 2 |
| | Chapter 18 | Lessons 1, 2, and 3 |
| | Chapter 19 | Lesson 2 |
| | Chapter 20 | Lesson 1 |
| 2.3  Implement data flow. | Chapter 3 | Lesson 1 |
| | Chapter 5 | Lessons 1, 2, and 3 |
| | Chapter 7 | Lessons 1 and 3 |
| | Chapter 13 | Lesson 1 and 2 |
| | Chapter 18 | Lesson 1 |
| | Chapter 20 | Lessons 2 and 3 |
| 2.4  Manage SSIS package execution. | Chapter 8 | Lessons 1 and 2 |
| | Chapter 12 | Lesson 1 |
| 2.5  Implement script tasks in SSIS. | Chapter 19 | Lesson 1 |
| **3.  LOAD DATA** | | |
| 3.1  Design control flow. | Chapter 3 | Lessons 2 and 3 |
| | Chapter 4 | Lessons 2 and 3 |
| | Chapter 6 | Lessons 1 and 3 |
| | Chapter 8 | Lessons 1, 2, and 3 |
| | Chapter 10 | Lesson 1 |
| | Chapter 12 | Lesson 1 and 2 |
| | Chapter 19 | Lesson 1 |
| 3.2  Implement package logic by using SSIS variables and parameters. | Chapter 6 | Lessons 1 and 2 |
| | Chapter 9 | Lessons 1 and 2 |
| 3.3  Implement control flow. | Chapter 4 | Lessons 2 and 3 |
| | Chapter 6 | Lesson 3 |
| | Chapter 8 | Lessons 1 and 2 |
| | Chapter 10 | Lesson 3 |
| | Chapter 13 | Lessons 1, 2, and 3 |
| 3.4  Implement data load options. | Chapter 7 | Lesson 2 |
| 3.5  Implement script components in SSIS. | Chapter 19 | Lesson 2 |

| OBJECTIVE | CHAPTER | LESSON |
|---|---|---|
| **4. CONFIGURE AND DEPLOY SSIS SOLUTIONS** | | |
| 4.1  Troubleshoot data integration issues. | Chapter 10 | Lesson 1 |
| | Chapter 13 | Lessons 1, 2, and 3 |
| 4.2  Install and maintain SSIS components. | Chapter 11 | Lesson 1 |
| 4.3  Implement auditing, logging, and event handling. | Chapter 8 | Lesson 3 |
| | Chapter 10 | Lessons 1 and 2 |
| 4.4  Deploy SSIS solutions. | Chapter 11 | Lessons 1 and 2 |
| | Chapter 19 | Lesson 3 |
| 4.5  Configure SSIS security settings. | Chapter 12 | Lessons 1 and 2 |
| **5. BUILD DATA QUALITY SOLUTIONS** | | |
| 5.1  Install and maintain Data Quality Services. | Chapter 14 | Lessons 1, 2, and 3 |
| 5.2  Implement master data management solutions. | Chapter 15 | Lessons 1, 2, and 3 |
| | Chapter 16 | Lessons 1, 2, and 3 |
| 5.3  Create a data quality project to clean data. | Chapter 14 | Lesson 1 |
| | Chapter 17 | Lessons 1, 2, and 3 |
| | Chapter 20 | Lessons 1 and 2 |

# Exam 70-463: Implementing a Data Warehouse with Microsoft® SQL Server® 2012

Training Kit

**Dejan Sarka**
**Matija Lah**
**Grega Jerkič**

# Contents at a Glance

# Contents

**What do you think of this book? We want to hear from you!**

Microsoft is interested in hearing your feedback so we can continually improve our
books and learning resources for you. To participate in a brief online survey, please visit:

**www.microsoft.com/learning/booksurvey/**

## PART II    DEVELOPING SSIS PACKAGES

**PART III     ENHANCING SSIS PACKAGES**

**Chapter 7    Enhancing Data Flow                                       283**

**PART IV     MANAGING AND MAINTAINING SSIS PACKAGES**

**Chapter 11  Installing SSIS and Deploying Packages     421**

## Chapter 13   Troubleshooting and Performance Tuning     497

## Chapter 20  Identity Mapping and De-Duplicating                  735

# Introduction

This Training Kit is designed for information technology (IT) professionals who support or plan to support data warehouses, extract-transform-load (ETL) processes, data quality improvements, and master data management. It is designed for IT professionals who also plan to take the Microsoft Certified Technology Specialist (MCTS) exam 70-463. The authors assume that you have a solid, foundation-level understanding of Microsoft SQL Server 2012 and the Transact-SQL language, and that you understand basic relational modeling concepts.

The material covered in this Training Kit and on Exam 70-463 relates to the technologies provided by SQL Server 2012 for implementing and maintaining a data warehouse. The topics in this Training Kit cover what you need to know for the exam as described on the Skills Measured tab for the exam, available at:

*http://www.microsoft.com/learning/en/us/exam.aspx?id=70-463*

By studying this Training Kit, you will see how to perform the following tasks:

- Design an appropriate data model for a data warehouse
- Optimize the physical design of a data warehouse
- Extract data from different data sources, transform and cleanse the data, and load it in your data warehouse by using SQL Server Integration Services (SSIS)
- Use advanced SSIS components
- Use SQL Server 2012 Master Data Services (MDS) to take control of your master data
- Use SQL Server Data Quality Services (DQS) for data cleansing

Refer to the objective mapping page in the front of this book to see where in the book each exam objective is covered.

## System Requirements

The following are the minimum system requirements for the computer you will be using to complete the practice exercises in this book and to run the companion CD.

## SQL Server and Other Software Requirements

This section contains the minimum SQL Server and other software requirements you will need:

- **SQL Server 2012**   You need access to a SQL Server 2012 instance with a logon that has permissions to create new databases—preferably one that is a member of the sysadmin role. For the purposes of this Training Kit, you can use almost any edition of

on-premises SQL Server (Standard, Enterprise, Business Intelligence, and Developer), both 32-bit and 64-bit editions. If you don't have access to an existing SQL Server instance, you can install a trial copy of SQL Server 2012 that you can use for 180 days. You can download a trial copy here:

> *http://www.microsoft.com/sqlserver/en/us/get-sql-server/try-it.aspx*

- **SQL Server 2012 Setup Feature Selection**   When you are in the Feature Selection dialog box of the SQL Server 2012 setup program, choose at minimum the following components:

  - Database Engine Services

  - Documentation Components

  - Management Tools - Basic

  - Management Tools – Complete

  - SQL Server Data Tools

- **Windows Software Development Kit (SDK) or Microsoft Visual Studio 2010**   The Windows SDK provides tools, compilers, headers, libraries, code samples, and a new help system that you can use to create applications that run on Windows. You need the Windows SDK for Chapter 19, "Implementing Custom Code in SSIS Packages" only. If you already have Visual Studio 2010, you do not need the Windows SDK. If you need the Windows SDK, you need to download the appropriate version for your operating system. For Windows 7, Windows Server 2003 R2 Standard Edition (32-bit x86), Windows Server 2003 R2 Standard x64 Edition, Windows Server 2008, Windows Server 2008 R2, Windows Vista, or Windows XP Service Pack 3, use the Microsoft Windows SDK for Windows 7 and the Microsoft .NET Framework 4 from:

  > *http://www.microsoft.com/en-us/download/details.aspx?id=8279*

## Hardware and Operating System Requirements

You can find the minimum hardware and operating system requirements for SQL Server 2012 here:

> *http://msdn.microsoft.com/en-us/library/ms143506(v=sql.110).aspx*

## Data Requirements

The minimum data requirements for the exercises in this Training Kit are the following:

- **The AdventureWorks OLTP and DW databases for SQL Server 2012**   Exercises in this book use the AdventureWorks online transactional processing (OLTP) database, which supports standard online transaction processing scenarios for a fictitious bicycle

manufacturer (Adventure Works Cycles), and the AdventureWorks data warehouse (DW) database, which demonstrates how to build a data warehouse. You need to download both databases for SQL Server 2012. You can download both databases from:

*http://msftdbprodsamples.codeplex.com/releases/view/55330*

You can also download the compressed file containing the data (.mdf) files for both databases from MS Press website here:

*http://www.microsoftpressstore.com/title/9780735666092.*

# Using the Companion CD

A companion CD is included with this Training Kit. The companion CD contains the following:

- **Practice tests**   You can reinforce your understanding of the topics covered in this Training Kit by using electronic practice tests that you customize to meet your needs. You can practice for the 70-463 certification exam by using tests created from a pool of over 200 realistic exam questions, which give you many practice exams to ensure that you are prepared.

- **An eBook**   An electronic version (eBook) of this book is included for when you do not want to carry the printed book with you.

- **Source code**   A compressed file called TK70463_CodeLabSolutions.zip includes the Training Kit's demo source code and exercise solutions. You can also download the compressed file from  website here:

   *http://www.microsoftpressstore.com/title/9780735666092.*

   For convenient access to the source code, create a local folder called **C:\TK463\** and extract the compressed archive by using this folder as the destination for the extracted files.

- **Sample data**   A compressed file called AdventureWorksDataFiles.zip includes the Training Kit's demo source code and exercise solutions. You can also download the compressed file from  website here:

   *http://www.microsoftpressstore.com/title/9780735666092.*

   For convenient access to the source code, create a local folder called **C:\TK463\** and extract the compressed archive by using this folder as the destination for the extracted files. Then use SQL Server Management Studio (SSMS) to attach both databases and create the log files for them.

# How to Install the Practice Tests

To install the practice test software from the companion CD to your hard disk, perform the following steps:

1. Insert the companion CD into your CD drive and accept the license agreement. A CD menu appears.

> **NOTE**  **IF THE CD MENU DOES NOT APPEAR**
>
> If the CD menu or the license agreement does not appear, AutoRun might be disabled on your computer. Refer to the Readme.txt file on the CD for alternate installation instructions.

2. Click Practice Tests and follow the instructions on the screen.

# How to Use the Practice Tests

To start the practice test software, follow these steps:

1. Click Start | All Programs, and then select Microsoft Press Training Kit Exam Prep.

    A window appears that shows all the Microsoft Press Training Kit exam prep suites installed on your computer.

2. Double-click the practice test you want to use.

When you start a practice test, you choose whether to take the test in Certification Mode, Study Mode, or Custom Mode:

- **Certification Mode**   Closely resembles the experience of taking a certification exam. The test has a set number of questions. It is timed, and you cannot pause and restart the timer.
- **Study Mode**   Creates an untimed test during which you can review the correct answers and the explanations after you answer each question.
- **Custom Mode**   Gives you full control over the test options so that you can customize them as you like.

In all modes, when you are taking the test, the user interface is basically the same but with different options enabled or disabled depending on the mode.

When you review your answer to an individual practice test question, a "References" section is provided that lists where in the Training Kit you can find the information that relates to that question and provides links to other sources of information. After you click Test Results

to score your entire practice test, you can click the Learning Plan tab to see a list of references for every objective.

## How to Uninstall the Practice Tests

To uninstall the practice test software for a Training Kit, use the Program And Features option in Windows Control Panel.

# Acknowledgments

A book is put together by many more people than the authors whose names are listed on the title page. We'd like to express our gratitude to the following people for all the work they have done in getting this book into your hands: Miloš Radivojević (technical editor) and Fritz Lechnitz (project manager) from SolidQ, Russell Jones (acquisitions and developmental editor) and Holly Bauer (production editor) from  Kathy Krause (copyeditor) and Jaime Odell (proofreader) from OTSI. In addition, we would like to give thanks to Matt Masson (member of the SSIS team), Wee Hyong Tok (SSIS team program manager), and Elad Ziklik (DQS group program manager) from Microsoft for the technical support and for unveiling the secrets of the new SQL Server 2012 products. There are many more people involved in writing and editing practice test questions, editing graphics, and performing other activities; we are grateful to all of them as well.

# Support & Feedback

The following sections provide information on errata, book support, feedback, and contact information.

## Errata

We've made every effort to ensure the accuracy of this book and its companion content. Any errors that have been reported since this book was published are listed on our Microsoft Press site:

> *http://www.microsoftpressstore.com/title/9780735666092.*

If you find an error that is not already listed, you can report it to us through the same page.

If you need additional support, email Microsoft Press Book Support at:

> *mspinput@microsoft.com*

Please note that product support for Microsoft software is not offered through the addresses above.

## We Want to Hear from You

At Microsoft Press, your satisfaction is our top priority, and your feedback our most valuable asset. Please tell us what you think of this book at:

*http://www.microsoft.com/learning/booksurvey*

The survey is short, and we read every one of your comments and ideas. Thanks in advance for your input!

## Stay in Touch

Let's keep the conversation going! We are on Twitter: *http://twitter.com/MicrosoftPress*.

# Preparing for the Exam

Microsoft certification exams are a great way to build your resume and let the world know about your level of expertise. Certification exams validate your on-the-job experience and product knowledge. While there is no substitution for on-the-job experience, preparation through study and hands-on practice can help you prepare for the exam. We recommend that you round out your exam preparation plan by using a combination of available study materials and courses. For example, you might use the training kit and another study guide for your "at home" preparation, and take a Microsoft Official Curriculum course for the classroom experience. Choose the combination that you think works best for you.

Note that this training kit is based on publicly available information about the exam and the authors' experience. To safeguard the integrity of the exam, authors do not have access to the live exam.

# Data Warehouse Logical Design

## Exam objectives in this chapter:

- ■ Design and Implement a Data Warehouse
  - ■ Design and implement dimensions.
  - ■ Design and implement fact tables.

Analyzing data from databases that support line-of-business (LOB) applications is usually not an easy task. The normalized relational schema used for an LOB application can consist of thousands of tables. Naming conventions are frequently not enforced. Therefore, it is hard to discover where the data you need for a report is stored. Enterprises frequently have multiple LOB applications, often working against more than one database. For the purposes of analysis, these enterprises need to be able to merge the data from multiple databases. Data quality is a common problem as well. In addition, many LOB applications do not track data over time, though many analyses depend on historical data.

A common solution to these problems is to create a *data warehouse (DW)*. A DW is a centralized data silo for an enterprise that contains merged, cleansed, and historical data. DW schemas are simplified and thus more suitable for generating reports than normalized relational schemas. For a DW, you typically use a special type of logical design called a Star schema, or a variant of the Star schema called a Snowflake schema. Tables in a Star or Snowflake schema are divided into dimension tables (commonly known as *dimensions*) and fact tables.

Data in a DW usually comes from LOB databases, but it's a transformed and cleansed copy of source data. Of course, there is some latency between the moment when data appears in an LOB database and the moment when it appears in a DW. One common method of addressing this latency involves refreshing the data in a DW as a nightly job. You use the refreshed data primarily for reports; therefore, the data is mostly read and rarely updated.

Queries often involve reading huge amounts of data and require large scans. To support such queries, it is imperative to use an appropriate physical design for a DW.

DW logical design seems to be simple at first glance. It is definitely much simpler than a normalized relational design. However, despite the simplicity, you can still encounter some advanced problems. In this chapter, you will learn how to design a DW and how to solve some of the common advanced design problems. You will explore Star and Snowflake schemas, dimensions, and fact tables. You will also learn how to track the source and time for data coming into a DW through auditing—or, in DW terminology, *lineage information*.

### Lessons in this chapter:

- Lesson 1: Introducing Star and Snowflake Schemas
- Lesson 2: Designing Dimensions
- Lesson 3: Designing Fact Tables

## Before You Begin

To complete this chapter, you must have:

- An understanding of normalized relational schemas.
- Experience working with Microsoft SQL Server 2012 Management Studio.
- A working knowledge of the Transact-SQL language.
- The AdventureWorks2012 and AdventureWorksDW2012 sample databases installed.

## Lesson 1: Introducing Star and Snowflake Schemas

Before you design a data warehouse, you need to understand some common design patterns used for a DW, namely the Star and Snowflake schemas. These schemas evolved in the 1980s. In particular, the Star schema is currently so widely used that it has become a kind of informal standard for all types of business intelligence (BI) applications.

**After this lesson, you will be able to:**

- Understand why a normalized schema causes reporting problems.
- Understand the Star schema.
- Understand the Snowflake schema.
- Determine granularity and auditing needs.

**Estimated lesson time: 40 minutes**

# Reporting Problems with a Normalized Schema

This lesson starts with normalized relational schema. Let's assume that you have to create a business report from a relational schema in the AdventureWorks2012 sample database. The report should include the sales amount for Internet sales in different countries over multiple years. The task (or even challenge) is to find out which tables and columns you would need to create the report. You start by investigating which tables store the data you need, as shown in Figure 1-1, which was created with the diagramming utility in SQL Server Management Studio (SSMS).



**FIGURE 1-1** A diagram of tables you would need for a simple sales report.

Even for this relatively simple report, you would end up with 10 tables. You need the sales tables and the tables containing information about customers. The AdventureWorks2012 database schema is highly normalized; it's intended as an example schema to support LOB applications. Although such a schema works extremely well for LOB applications, it can cause problems when used as the source for reports, as you'll see in the rest of this section.

Normalization is a process in which you define entities in such a way that a single table represents exactly one entity. The goal is to have a complete and non-redundant schema. Every piece of information must be stored exactly once. This way, you can enforce data integrity. You have a place for every piece of data, and because each data item is stored only once, you do not have consistency problems. However, after a proper normalization, you typically wind up with many tables. In a database that supports an LOB application for an enterprise, you might finish with thousands of tables!

Finding the appropriate tables and columns you need for a report can be painful in a normalized database simply because of the number of tables involved. Add to this the fact that nothing forces database developers to maintain good naming conventions in an LOB database. It's relatively easy to find the pertinent tables in AdventureWorks2012, because the tables and columns have meaningful names. But imagine if the database contained tables named *Table1*, *Table2*, and so on, and columns named *Column1*, *Column2*, and so on. Finding the objects you need for your report would be a nightmare. Tools such as SQL Profiler might help. For example, you could create a test environment, try to insert some data through an LOB application, and have SQL Profiler identify where the data was inserted. A normalized schema is not very narrative. You cannot easily spot the storage location for data that measures something, such as the sales amount in this example, or the data that gives context to these measures, such as countries and years.

In addition, a query that joins 10 tables, as would be required in reporting sales by countries and years, would not be very fast. The query would also read huge amounts of data—sales over multiple years—and thus would interfere with the regular transactional work of inserting and updating the data.

Another problem in this example is the fact that there is no explicit lookup table for dates. You have to extract years from date or date/time columns in sales tables, such as *OrderDate* from the *SalesOrderHeader* table in this example. Extracting years from a date column is not such a big deal; however, the first question is, does the LOB database store data for multiple years? In many cases, LOB databases are purged after each new fiscal year starts. Even if you have all of the historical data for the sales transactions, you might have a problem showing the historical data correctly. For example, you might have only the latest customer address (from which you extract customer's current country), which might prevent you from calculating historical sales by country correctly.

The AdventureWorks2012 sample database stores all data in a single database. However, in an enterprise, you might have multiple LOB applications, each of which might store data in its own database. You might also have part of the sales data in one database and part in another. And you could have customer data in both databases, without a common identification. In such cases, you face the problems of how to merge all this data and how to identify which customer from one database is actually the same as a customer from another database.

Finally, data quality could be low. The old rule, "garbage in garbage out," applies to analyses as well. Parts of the data could be missing; other parts could be wrong. Even with good data, you could still have different representations of the same data in different databases. For example, gender in one database could be represented with the letters F and M, and in another database with the numbers 1 and 2.

The problems listed in this section are indicative of the problems that led designers to create different schemas for BI applications. The Star and Snowflake schemas are both simplified and narrative. A data warehouse should use Star and/or Snowflake designs. You'll also sometimes find the term *dimensional model* used for a DW schema. A dimensional model actually consists of both Star and Snowflake schemas. This is a good time to introduce the Star and Snowflake schemas.

# Star Schema

Often, a picture is worth more than a thousand words. Figure 1-2 shows a *Star schema*, a diagram created in SSMS from a subset of the tables in the AdventureWorksDW2012 sample database.

In Figure 1-2, you can easily spot how the Star schema got its name—it resembles a star. There is a single central table, called a *fact table*, surrounded by multiple tables called *dimensions*. One Star schema covers one business area. In this case, the schema covers Internet sales. An enterprise data warehouse covers multiple business areas and consists of multiple Star (and/or Snowflake) schemas.



**FIGURE 1-2**  A Star schema example.

The fact table is connected to all the dimensions with foreign keys. Usually, all foreign keys taken together uniquely identify each row in the fact table, and thus collectively form a unique key, so you can use all the foreign keys as a composite primary key of the fact table. You can also add a simpler key. The fact table is on the "many" side of its relationships with the dimensions. If you were to form a proposition from a row in a fact table, you might express it with a sentence such as, "Customer A purchased product B on date C in quantity D for amount E." This proposition is a fact; this is how the fact table got its name.

The Star schema evolved from a conceptual model of a cube. You can imagine all sales as a big box. When you search for a problem in sales data, you use a divide-and-conquer technique: slicing the cube over different categories of customers, products, or time. In other words, you slice the cube over its dimensions. Therefore, customers, products, and time represent the three dimensions in the conceptual model of the sales cube. Dimension tables (dimensions) got their name from this conceptual model. In a logical model of a Star schema, you can represent more than three dimensions. Therefore, a Star schema represents a multi-dimensional hypercube.

As you already know, a data warehouse consists of multiple Star schemas. From a business perspective, these Star schemas are connected. For example, you have the same customers in sales as in accounting. You deal with many of the same products in sales, inventory, and production. Of course, your business is performed at the same time over all the different business areas. To represent the business correctly, you must be able to connect the multiple Star schemas in your data warehouse. The connection is simple – you use the same dimensions for each Star schema. In fact, the dimensions should be shared among multiple Star schemas. Dimensions have foreign key relationships with multiple fact tables. Dimensions with connections to multiple fact tables are called *shared* or *conformed dimensions*. Figure 1-3 shows a conformed dimension from the AdventureWorksDW2012 sample database with two different fact tables sharing the same dimension.



**FIGURE 1-3** *DimProduct* is a shared dimension.

In the past, there was a big debate over whether to use shared or private dimensions. Private dimensions are dimensions that pertain to only a single Star schema. However, it is quite simple to design shared dimensions; you do not gain much from the design-time perspective by using private dimensions. In fact, with private dimensions, you lose the connections between the different fact tables, so you cannot compare the data in different fact tables over the same dimensions. For example, you could not compare sales and accounting data for the same customer if the sales and accounting fact tables didn't share the same customer dimension. Therefore, unless you are creating a small proof-of-concept (POC) project that covers only a single business area where you do not care about connections with different business areas, you should always opt for shared dimensions.

A data warehouse is often the source for specialized analytical database management systems, such as SQL Server Analysis Services (SSAS). SSAS is a system that performs specialized analyses by drilling down and is used for analyses that are based on the conceptual model of a cube. Systems such as SSAS focus on a single task and fast analyses, and they're considerably more optimized for this task than general systems such as SQL Server. SSAS enables analysis in real time, a process called *online analytical processing* (OLAP). However, to get such performance, you have to pay a price. SSAS is out of the scope of this book, but you have to know the limitations of SSAS to prepare a data warehouse in a way that is useful for SSAS. One thing to remember is that in an SSAS database, you can use shared dimensions only. This is just one more reason why you should prefer shared to private dimensions.

## Snowflake Schema

Figure 1-4 shows a more detailed view of the *DimDate* dimension from the AdventureWorks-DW2012 sample database.

The highlighted attributes show that the dimension is denormalized. It is not in third normal form. In third normal form, all non-key columns should nontransitively depend on the key. A different way to say this is that there should be no functional dependency between non-key columns. You should be able to retrieve the value of a non-key column only if you know the key. However, in the *DimDate* dimension, if you know the month, you obviously know the calendar quarter, and if you know the calendar quarter, you know the calendar semester.

In a Star schema, dimensions are denormalized. In contrast, in an LOB normalized schema, you would split the table into multiple tables if you found a dependency between non-key columns. Figure 1-5 shows such a normalized example for the *DimProduct*, *DimProduct-Subcategory* and *DimProductCategory* tables from the AdventureWorksDW2012 database.

**FIGURE 1-4** The *DimDate* denormalized dimension.



**FIGURE 1-5** The *DimProduct* normalized dimension.

The *DimProduct* dimension is not denormalized. The *DimProduct* table does not contain the subcategory name, only the ProductSubcategoryKey value for the foreign key to the *DimProductSubcategory* lookup table. Similarly, the *DimProductSubcategory* table does not contain a category name; it just holds the foreign key ProductCategoryKey from the *Dim-ProductCategory* table. This design is typical of an LOB database schema.

You can imagine multiple dimensions designed in a similar normalized way, with a central fact table connected by foreign keys to dimension tables, which are connected with foreign keys to lookup tables, which are connected with foreign keys to second-level lookup tables.

In this configuration, a star starts to resemble a snowflake. Therefore, a Star schema with normalized dimensions is called a *Snowflake schema*.

In most long-term projects, you should design Star schemas. Because the Star schema is simpler than a Snowflake schema, it is also easier to maintain. Queries on a Star schema are simpler and faster than queries on a Snowflake schema, because they involve fewer joins. The Snowflake schema is more appropriate for short POC projects, because it is closer to an LOB normalized relational schema and thus requires less work to build.

---

*EXAM TIP*

**If you do not use OLAP cubes and your reports query your data warehouse directly, then using a Star instead of a Snowflake schema might speed up the reports, because your reporting queries involve fewer joins.**

---

In some cases, you can also employ a hybrid approach, using a Snowflake schema only for the first level of a dimension lookup table. In this type of approach, there are no additional levels of lookup tables; the first-level lookup table is denormalized. Figure 1-6 shows such a partially denormalized schema.



**FIGURE 1-6** Partially denormalized dimensions.

In Figure 1-6, the *DimCustomer* and *DimReseller* dimensions are partially normalized. The dimensions now contain only the GeographyKey foreign key. However, the *DimGeography* table is denormalized. There is no additional lookup table even though a city is in a region and a region is in a country. A hybrid design such as this means that geography data is written only once and needs to be maintained in only a single place. Such a design is appropriate

when multiple dimensions share the same attributes. In other cases, you should use the simpler Star schema. To repeat: you should use a Snowflake schema only for quick POC projects.

> **✔ Quick Check**
> - How do you connect multiple Star schemas in a DW?
>
> **Quick Check Answer**
> - You connect multiple Star schemas through shared dimensions.

## Granularity Level

The number of dimensions connected with a fact table defines the level of granularity of analysis you can get. For example, if no products dimension is connected to a sales fact table, you cannot get a report at the product level—you could get a report for sales for all products only. This kind of granularity is also called the *dimensionality* of a Star schema.

But there is another kind of granularity, which lets you know what level of information a dimension foreign key represents in a fact table. Different fact tables can have different granularity in a connection to the same dimension. This is very typical in budgeting and planning scenarios. For example, you do not plan that customer A will come on date B to store C and buy product D for amount E. Instead, you plan on a higher level—you might plan to sell amount E of products C in quarter B in all stores in that region to all customers in that region. Figure 1-7 shows an example of a fact table that uses a higher level of granularity than the fact tables introduced so far.

In the AdventureWorksDW2012 database, the *FactSalesQuota* table is the fact table with planning data. However, plans are made for employees at the per-quarter level only. The plan is for all customers, all products, and so on, because this Star schema uses only the *DimDate* and *DimEmployee* dimensions. In addition, planning occurs at the quarterly level. By investigating the content, you could see that all plans for a quarter are bound to the first day of a quarter. You would not need to use the DateKey; you could have only *CalendarYear* and *CalendarQuarter* columns in the *FactSalesQuota* fact table. You could still perform joins to *DimDate* by using these two columns—they are both present in the *DimDate* table as well. However, if you want to have a foreign key to the *DimDate* dimension, you do need the DateKey. A foreign key must refer to unique values on the "one" side of the relationship. The combination of *CalendarYear* and *CalendarQuarter* is, of course, not unique in the *DimDate* dimension; it repeats approximately 90 times in each quarter.

**FIGURE 1-7** A fact table with a higher level of granularity.

## Auditing and Lineage

In addition to tables for reports, a data warehouse may also include auditing tables. For every update, you should audit who made the update, when it was made, and how many rows were transferred to each dimension and fact table in your DW. If you also audit how much time was needed for each load, you can calculate the performance and take action if it deteriorates. You store this information in an auditing table or tables. However, you should realize that auditing does not help you unless you analyze the information regularly.

Auditing tables hold batch-level information about regular DW loads, but you might also want or need to have more detailed information. For example, you might want to know where each row in a dimension and/or fact table came from and when it was added. In such cases, you must add appropriate columns to the dimension and fact tables. Such detailed auditing information is also called *lineage* in DW terminology. To collect either auditing or lineage information, you need to modify the extract-transform-load (ETL) process you use for DW loads appropriately.

If your ETL tool is SQL Server Integration Services (SSIS), then you should use SSIS logging. SSIS has extensive logging support. In addition, SSIS also has support for lineage information.

**Reviewing the AdventureWorksDW2012 Internet Sales Schema**

The AdventureWorksDW2012 sample database is a good example of a data warehouse. It has all the elements needed to allow you to see examples of various types of dimensional modeling.

### EXERCISE 1    Review the AdventureWorksDW2012 Database Schema

In this exercise, you review the database schema.

1. Start SSMS and connect to your instance of SQL Server. Expand the Databases folder and then the AdventureWorksDW2012 database.

2. Right-click the Database Diagrams folder and select the New Database Diagram option. If no diagrams were ever created in this database, you will see a message box informing you that the database has no support objects for diagramming. If that message appears, click Yes to create the support objects.

3. From the Add Table list, select the following tables (click each table and then click the Add button):
   - *DimCustomer*
   - *DimDate*
   - *DimGeography*
   - *DimProduct*
   - *DimProductCategory*
   - *DimProductSubcategory*
   - *FactInternetSales*

   Your diagram should look similar to Figure 1-8.

**FIGURE 1-8** The AdventureWorksDW2012 Internet Sales Schema.

4. Thoroughly analyze the tables, columns, and relationships.

5. Save the diagram with the name **Practice_01_01_InternetSales**.

### EXERCISE 2 Analyze the Diagram

Review the AdventureWorksDW2012 schema to note the following facts:

- The *DimDate* dimension has no additional lookup tables associated with it and therefore uses the Star schema.

- The *DimProduct* table is snowflaked; it uses the *DimProductSubcategory* lookup table, which further uses the *DimProductCategory* lookup table.

- The *DimCustomer* dimension uses a hybrid schema—the first level of the Snowflake schema only through the *DimGeography* lookup table. The *DimGeography* table is denormalized; it does not have a relationship with any other lookup table.

- There are no specific columns for lineage information in any of the tables.

Close the diagram.

---

*NOTE* **CONTINUING WITH PRACTICES**

**Do not exit SSMS if you intend to continue immediately with the next practice.**

---

## Lesson Summary

- The Star schema is the most common design for a DW.
- The Snowflake schema is more appropriate for POC projects.
- You should also determine the granularity of fact tables, as well as auditing and lineage needs.

## Lesson Review

Answer the following questions to test your knowledge of the information in this lesson. You can find the answers to these questions and explanations of why each answer choice is correct or incorrect in the "Answers" section at the end of this chapter.

1. Reporting from a Star schema is simpler than reporting from a normalized online transactional processing (OLTP) schema. What are the reasons for wanting simpler reporting? (Choose all that apply.)

   A. A Star schema typically has fewer tables than a normalized schema. Therefore, queries are simpler because they require fewer joins.

   B. A Star schema has better support for numeric data types than a normalized relational schema; therefore, it is easier to create aggregates.

   C. There are specific Transact-SQL expressions that deal with Star schemas.

   D. A Star schema is standardized and narrative; you can find the information you need for a report quickly.

2. You are creating a quick POC project. Which schema is the most suitable for this kind of a project?

   A. Star schema

   B. Normalized schema

   C. Snowflake schema

   D. XML schema

3. A Star schema has two types of tables. What are those two types? (Choose all that apply.)

   A. Lookup tables

   B. Dimensions

   C. Measures

   D. Fact tables

# Lesson 2: Designing Dimensions

Star and Snowflake schemas are the de facto standard. However, the standard does not end with schema shapes. Dimension and fact table columns are part of this informal standard as well and are introduced in this lesson, along with natural hierarchies, which are especially useful as natural drill-down paths for analyses. Finally, the lesson discusses a common problem with handling dimension changes over time.

> **After this lesson, you will be able to:**
>
> - Define dimension column types.
> - Use natural hierarchies.
> - Understand and resolve the slowly changing dimensions problem.
>
> **Estimated lesson time: 40 minutes**

## Dimension Column Types

Dimensions give context to measures. Typical analysis includes pivot tables and pivot graphs. These pivot on one or more dimension columns used for analysis—these columns are called *attributes* in DW and OLAP terminology. The naming convention in DW/OLAP terminology is a little odd; in a relational model, every column represents an attribute of an entity. Don't worry too much about the correctness of naming in DW/OLAP terminology. The important point here is for you to understand what the word "attribute" means in a DW/OLAP context.

Pivoting makes no sense if an attribute's values are continuous, or if an attribute has too many distinct values. Imagine how a pivot table would look if it had 1,000 columns, or how a pivot graph would look with 1,000 bars. For pivoting, discrete attributes with a small number of distinct values is most appropriate. A bar chart with more than 10 bars becomes difficult to comprehend. Continuous columns or columns with unique values, such as keys, are not appropriate for analyses.

If you have a continuous column and you would like to use it in analyses as a pivoting attribute, you should discretize it. *Discretizing* means grouping or binning values to a few discrete groups. If you are using OLAP cubes, SSAS can help you. SSAS can discretize continuous attributes. However, automatic discretization is usually worse than discretization from a business perspective. Age and income are typical attributes that should be discretized from a business perspective. One year makes a big difference when you are 15 years old, and much less when you are 55 years old. When you discretize age, you should use narrower ranges for younger people and wider ranges for older people.

Columns with unique values identify rows. These columns are *keys*. In a data warehouse, you need keys just like you need them in an LOB database. Keys uniquely identify entities. Therefore, keys are the second type of columns in a dimension.

After you identify a customer, you do not refer to that customer with the key value. Having only keys in a report does not make the report very readable. People refer to entities by using their names. In a DW dimension, you also need one or more columns that you use for naming an entity.

A customer typically has an address, a phone number, and an email address. You do not analyze data on these columns. You do not need them for pivoting. However, you often need information such as the customer's address on a report. If that data is not present in a DW, you will need to get it from an LOB database, probably with a distributed query. It is much simpler to store this data in your data warehouse. In addition, queries that use this data perform better, because the queries do not have to include data from LOB databases. Columns used in reports as labels only, not for pivoting, are called *member properties*.

You can have naming and member property columns in multiple languages in your dimension tables, providing the translation for each language you need to support. SSAS can use your translations automatically. For reports from a data warehouse, you need to manually select columns with appropriate language translation.

In addition to the types of dimension columns already defined for identifying, naming, pivoting, and labeling on a report, you can have columns for lineage information, as you saw in the previous lesson. There is an important difference between lineage and other columns: lineage columns are never exposed to end users and are never shown on end users' reports.

To summarize, a dimension may contain the following types of columns:

- **Keys**   Used to identify entities
- **Name columns**   Used for human names of entities
- **Attributes**   Used for pivoting in analyses
- **Member properties**   Used for labels in a report
- **Lineage columns**   Used for auditing, and never exposed to end users

# Hierarchies

Figure 1-9 shows the *DimCustomer* dimension of the AdventureWorksDW2012 sample database.



**FIGURE 1-9**  The *DimCustomer* dimension.

In the figure, the following columns are attributes (columns used for pivoting):

- *BirthDate* (after calculating age and discretizing the age)
- *MaritalStatus*
- *Gender*
- *YearlyIncome* (after discretizing)
- *TotalChildren*
- *NumberChildrenAtHome*
- *EnglishEducation* (other education columns are for translations)
- *EnglishOccupation* (other occupation columns are for translations)
- *HouseOwnerFlag*
- *NumberCarsOwned*
- *CommuteDistance*

All these attributes are unrelated. Pivoting on *MaritalStatus*, for example, is unrelated to pivoting on *YearlyIncome*. None of these columns have any functional dependency between them, and there is no natural drill-down path through these attributes. Now look at the *Dim-Date* columns, as shown in Figure 1-10.

**DimDate**

| | |
|---|---|
| 🔑 | DateKey |
| | FullDateAlternateKey |
| | DayNumberOfWeek |
| | EnglishDayNameOfWeek |
| | SpanishDayNameOfWeek |
| | FrenchDayNameOfWeek |
| | DayNumberOfMonth |
| | DayNumberOfYear |
| | WeekNumberOfYear |
| | EnglishMonthName |
| | SpanishMonthName |
| | FrenchMonthName |
| | MonthNumberOfYear |
| | CalendarQuarter |
| | CalendarYear |
| | CalendarSemester |
| | FiscalQuarter |
| | FiscalYear |
| | FiscalSemester |

**FIGURE 1-10** The *DimDate* dimension.

Some attributes of the *DimDate* dimension include the following (not in the order shown in the figure):

- *FullDateAlternateKey* (denotes a date in date format)
- *EnglishMonthName*
- *CalendarQuarter*
- *CalendarSemester*
- *CalendarYear*

You will immediately notice that these attributes are connected. There is a functional dependency among them, so they break third normal form. They form a hierarchy. Hierarchies are particularly useful for pivoting and OLAP analyses—they provide a natural drill-down path. You perform divide-and-conquer analyses through hierarchies.

Hierarchies have levels. When drilling down, you move from a parent level to a child level. For example, a calendar drill-down path in the *DimDate* dimension goes through the following levels: *CalendarYear* → *CalendarSemester* → *CalendarQuarter* → *EnglishMonthName* → *FullDateAlternateKey*.

At each level, you have members. For example, the members of the month level are, of course, January, February, March, April, May, June, July, August, September, October, November, and December. In DW and OLAP jargon, rows on the leaf level—the actual dimension

rows—are called *members*. This is why dimension columns used in reports for labels are called *member properties*.

In a Snowflake schema, lookup tables show you levels of hierarchies. In a Star schema, you need to extract natural hierarchies from the names and content of columns. Nevertheless, because drilling down through natural hierarchies is so useful and welcomed by end users, you should use them as much as possible.

Note also that attribute names are used for labels of row and column groups in a pivot table. Therefore, a good naming convention is crucial for a data warehouse. You should always use meaningful and descriptive names for dimensions and attributes.

## Slowly Changing Dimensions

There is one common problem with dimensions in a data warehouse: the data in the dimension changes over time. This is usually not a problem in an OLTP application; when a piece of data changes, you just update it. However, in a DW, you have to maintain history. The question that arises is *how* to maintain it. Do you want to update only the changed data, as in an OLTP application, and pretend that the value was always the last value, or do you want to maintain both the first and intermediate values? This problem is known in DW jargon as the *Slowly Changing Dimension (SCD)* problem.

The problem is best explained in an example. Table 1-1 shows original source OLTP data for a customer.

**TABLE 1-1** Original OLTP Data for a Customer

| CustomerId | FullName | City | Occupation |
|---|---|---|---|
| 17 | Bostjan Strazar | Vienna | Professional |

The customer lives in Vienna, Austria, and is a professional. Now imagine that the customer moves to Ljubljana, Slovenia. In an OLTP database, you would just update the *City* column, resulting in the values shown in Table 1-2.

**TABLE 1-2** OLTP Data for a Customer After the City Change

| CustomerId | FullName | City | Occupation |
|---|---|---|---|
| 17 | Bostjan Strazar | Ljubljana | Professional |

If you create a report, all the historical sales for this customer are now attributed to the city of Ljubljana, and (on a higher level) to Slovenia. The fact that this customer contributed to sales in Vienna and in Austria in the past would have disappeared.

In a DW, you can have the same data as in an OLTP database. You could use the same key, such as the business key, for your *Customer* dimension. You could update the *City* column when you get a change notification from the OLTP system, and thus overwrite the history.

This kind of change management is called *Type 1 SCD*. To recapitulate, Type 1 means over-writing the history for an attribute and for all higher levels of hierarchies to which that attribute belongs.

But you might prefer to maintain the history, to capture the fact that the customer contributed to sales in another city and country or region. In that case, you cannot just overwrite the data; you have to insert a new row containing new data instead. Of course, the values of other columns that do not change remain the same. However, that creates a new problem. If you simply add a new row for the customer with the same key value, the key would no longer be unique. In fact, if you tried to use a primary key or unique constraint as the key, the constraint would reject such an insert. Therefore, you have to do something with the key. You should not modify the business key, because you need a connection with the source system. The solution is to introduce a new key, a *data warehouse key*. In DW terminology, this kind of key is called a *surrogate key*.

Preserving the history while adding new rows is known as *Type 2 SCD*. When you implement Type 2 SCD, for the sake of simpler querying, you typically also add a flag to denote which row is current for a dimension member. Alternatively, you could add two columns showing the interval of validity of a value. The data type of the two columns should be Date, and the columns should show the values Valid From and Valid To. For the current value, the *Valid To* column should be NULL. Table 1-3 shows an example of the flag version of Type 2 SCD handling.

**TABLE 1-3** An SCD Type 2 Change

| DWCId | CustomerId | FullName | City | Occupation | Current |
|-------|-----------|----------|------|-----------|---------|
| 17 | 17 | Bostjan Strazar | Vienna | Professional | 0 |
| 289 | 17 | Bostjan Strazar | Ljubljana | Professional | 1 |

You could have a mixture of Type 1 and Type 2 changes in a single dimension. For example, in Table 1-3, you might want to maintain the history for the *City* column but overwrite the history for the *Occupation* column. That raises yet another issue. When you want to update the *Occupation* column, you may find that there are two (and maybe more) rows for the same customer. The question is, do you want to update the last row only, or all the rows? Table 1-4 shows a version that updates the last (current) row only, whereas Table 1-5 shows all of the rows being updated.

**TABLE 1-4** An SCD Type 1 and Type 2 Mixture, Updating the Current Row Only

| DWCId | CustomerId | FullName | City | Occupation | Current |
|-------|-----------|----------|------|-----------|---------|
| 17 | 17 | Bostjan Strazar | Vienna | Professional | 0 |
| 289 | 17 | Bostjan Strazar | Ljubljana | Management | 1 |

**TABLE 1-5** An SCD Type 1 and Type 2 Mixture, Updating All Rows

| DWCId | CustomerId | FullName | City | Occupation | Current |
|-------|-----------|----------|------|-----------|---------|
| 17 | 17 | Bostjan Strazar | Vienna | Management | 0 |
| 289 | 17 | Bostjan Strazar | Ljubljana | Management | 1 |

Although Type 1 and Type 2 handling are most common, other solutions exist. Especially well-known is *Type 3 SCD*, in which you manage a limited amount of history through additional historical columns. Table 1-6 shows Type 3 handling for the *City* column.

**TABLE 1-6** SCD Type 3

| CustomerId | FullName | CurrentCity | PreviousCity | Occupation |
|-----------|----------|-------------|--------------|-----------|
| 17 | Bostjan Strazar | Ljubljana | Vienna | Professional |

You can see that by using only a single historical column, you can maintain only one historical value per column. So Type 3 SCD has limited usability and is far less popular than Types 1 and 2.

Which solution should you implement? You should discuss this with end users and subject matter experts (SMEs). They should decide for which attributes to maintain the history, and for which ones to overwrite the history. You should then choose a solution that uses Type 2, Type 1, or a mixture of Types 1 and 2, as appropriate.

However, there is an important caveat. To maintain customer history correctly, you must have some attribute that uniquely identifies that customer throughout that customer's history, and that attribute must not change. Such an attribute should be the original—the business key. In an OLTP database, business keys should not change.

Business keys should also not change if you are merging data from multiple sources. For merged data, you usually have to implement a new, surrogate key, because business keys from different sources can have the same value for different entities. However, business keys should not change; otherwise you lose the connection with the OLTP system. Using surrogate keys in a data warehouse for at least the most common dimensions (those representing customers, products, and similar important data), is considered a best practice. Not changing OLTP keys is a best practice as well.

---

*EXAM TIP*

**Make sure you understand why you need surrogate keys in a data warehouse.**

---

The AdventureWorksDW2012 sample database has many dimensions. In this practice, you will explore some of them.

### EXERCISE 1    Explore the AdventureWorksDW2012 Dimensions

In this exercise, you create a diagram for the dimensions.

1. If you closed SSMS, start it and connect to your SQL Server instance. Expand the Databases folder and then the AdventureWorksDW2012 database.

2. Right-click the Database Diagrams folder, and then select the New Database Diagram option.

3. From the Add Table list, select the following tables (click each table and then click the Add button):

   - *DimProduct*
   - *DimProductCategory*
   - *DimProductSubcategory*

   Your diagram should look like Figure 1-11.

4. Try to figure out which columns are used for the following purposes:

   - Keys
   - Names
   - Translations
   - Attributes
   - Member properties
   - Lineage
   - Natural hierarchies

5. Try to figure out whether the tables in the diagram are prepared for a Type 2 SCD change.

6. Add the *DimSalesReason* table to the diagram.

7. Try to figure out whether there is some natural hierarchy between attributes of the *DimSalesReason* dimension. Your diagram should look like Figure 1-12.

8. Save the diagram with the name **Practice_01_02_Dimensions**.

**FIGURE 1-11**  *DimProduct* and related tables.



**FIGURE 1-12**  Adding *DimSalesReason*.

In this exercise, review the database schema from the previous exercise to learn more:

■ The *DimProduct* dimension has a natural hierarchy: *ProductCategory* → *ProductSubcategory* → *Product.*

■ The *DimProduct* dimension has many additional attributes that are useful for pivoting but that are not a part of any natural hierarchy. For example, Color and Size are such attributes.

■ Some columns in the *DimProduct* dimension, such as the *LargePhoto* and *Description* columns, are member properties.

■ *DimSalesReason* uses a Star schema. In a Star schema, it is more difficult to spot natural hierarchies. Though you can simply follow the lookup tables in a Snowflake schema and find levels of hierarchies, you have to recognize hierarchies from attribute names in a Star schema. If you cannot extract hierarchies from column names, you could also check the data. In the *DimSalesReason* dimension, it seems that there is a natural hierarchy: *SalesReasonReasonType* → *SalesReasonName.*

Close the diagram.

> *NOTE*    **CONTINUING WITH PRACTICES**
>
> **Do not exit SSMS if you intend to continue immediately with the next practice.**

## Lesson Summary

■ In a dimension, you have the following column types: keys, names, attributes, member properties, translations, and lineage.

■ Some attributes form natural hierarchies.

■ There are standard solutions for the Slowly Changing Dimensions (SCD) problem.

## Lesson Review

Answer the following questions to test your knowledge of the information in this lesson. You can find the answers to these questions and explanations of why each answer choice is correct or incorrect in the "Answers" section at the end of this chapter.

1.   You implement a Type 2 solution for an SCD problem for a specific column. What do you actually do when you get a changed value for the column from the source system?

   **A.**  Add a column for the previous value to the table. Move the current value of the updated column to the new column. Update the current value with the new value from the source system.

**B.** Insert a new row for the same dimension member with the new value for the updated column. Use a surrogate key, because the business key is now duplicated. Add a flag that denotes which row is current for a member.

**C.** Do nothing, because in a DW, you maintain history, you do not update dimension data.

**D.** Update the value of the column just as it was updated in the source system.

2. Which kind of a column is not a part of a dimension?

    **A.** Attribute

    **B.** Measure

    **C.** Key

    **D.** Member property

    **E.** Name

3. How can you spot natural hierarchies in a Snowflake schema?

    **A.** You need to analyze the content of the attributes of each dimension.

    **B.** Lookup tables for each dimension provide natural hierarchies.

    **C.** A Snowflake schema does not support hierarchies.

    **D.** You should convert the Snowflake schema to the Star schema, and then you would spot the natural hierarchies immediately.

# Lesson 3: Designing Fact Tables

Fact tables, like dimensions, have specific types of columns that limit the actions that can be taken with them. Queries from a DW aggregate data; depending on the particular type of column, there are some limitations on which aggregate functions you can use. Many-to-many relationships in a DW can be implemented differently than in a normalized relational schema.

> **After this lesson, you will be able to:**
> - Define fact table column types.
> - Understand the additivity of a measure.
> - Handle many-to-many relationships in a Star schema.
>
> **Estimated lesson time: 30 minutes**

# Fact Table Column Types

Fact tables are collections of measurements associated with a specific business process. You store measurements in columns. Logically, this type of column is called a *measure*. Measures are the essence of a fact table. They are usually numeric and can be aggregated. They store values that are of interest to the business, such as sales amount, order quantity, and discount amount.

From Lesson 1 in this chapter, you already saw that a fact table includes foreign keys from all dimensions. These foreign keys are the second type of column in a fact table. A fact table is on the "many" side of the relationships with dimensions. All foreign keys together usually uniquely identify each row and can be used as a composite primary key.

You often include an additional surrogate key. This key is shorter and consists of one or two columns only. The surrogate key is usually the business key from the table that was used as the primary source for the fact table. For example, suppose you start building a sales fact table from an order details table in a source system, and then add foreign keys that pertain to the order as a whole from the *Order Header* table in the source system. Tables 1-7, 1-8, and 1-9 illustrate an example of such a design process.

Table 1-7 shows a simplified example of an *Orders Header* source table. The *OrderId* column is the primary key for this table. The *CustomerId* column is a foreign key from the *Customers* table. The *OrderDate* column is not a foreign key in the source table; however, it becomes a foreign key in the DW fact table, for the relationship with the explicit date dimension. Note, however, that foreign keys in a fact table can—and usually are—replaced with DW surrogate keys of DW dimensions.

**TABLE 1-7** The Source Orders Header Table

| OrderId | CustomerId | Orderdate |
|---------|------------|-----------|
| 12541 | 17 | 2012/02/21 |

Table 1-8 shows the source *Order Details* table. The primary key of this table is a composite one and consists of the *OrderId* and *LineItemId* columns. In addition, the source *Order Details* table has the *ProductId* foreign key column. The *Quantity* column is the measure.

**TABLE 1-8** The Source Order Details Table

| OrderId | LineItemId | ProductId | Quantity |
|---------|------------|-----------|----------|
| 12541 | 2 | 5 | 47 |

Table 1-9 shows the *Sales Fact* table created from the *Orders Header* and *Order Details* source tables. The *Order Details* table was the primary source for this fact table. The *OrderId*,

*LineItemId*, and *Quantity* columns are simply transferred from the source *Order Details* table. The *ProductId* column from the source *Order Details* table is replaced with a surrogate DW *ProductKey* column. The *CustomerId* and *OrderDate* columns take the source *Orders Header* table; these columns pertain to orders, not order details. However, in the fact table, they are replaced with the surrogate DW keys CustomerKey and OrderDateKey.

**TABLE 1-9** The Sales Fact Table

| OrderId | LineItemId | CustomerKey | OrderDateKey | ProductKey | Quantity |
|---------|-----------|-------------|--------------|------------|----------|
| 12541 | 2 | 289 | 444 | 25 | 47 |

You do not need the *OrderId* and *LineItemId* columns in this sales fact table. For analyses, you could create a composite primary key from the *CustomerKey*, *OrderDateKey*, and *ProductKey* columns. However, you should keep the *OrderId* and *LineItemId* columns to make quick controls and comparisons with source data possible. In addition, if you were to use them as the primary key, then the primary key would be shorter than one composed from all foreign keys.

The last column type used in a fact table is the lineage type, if you implement the lineage. Just as with dimensions, you never expose the lineage information to end users. To recapitulate, fact tables have the following column types:

- Foreign keys
- Measures
- Lineage columns (optional)
- Business key columns from the primary source table (optional)

## Additivity of Measures

Additivity of measures is not exactly a data warehouse design problem. However, you should consider which aggregate functions you will use in reports for which measures, and which aggregate functions you will use when aggregating over which dimension.

The simplest types of measures are those that can be aggregated with the SUM aggregate function across all dimensions, such as amounts or quantities. For example, if sales for product A were $200.00 and sales for product B were $150.00, then the total of the sales was $350.00. If yesterday's sales were $100.00 and sales for the day before yesterday were $130.00, then the total sales amounted to $230.00. Measures that can be summarized across all dimensions are called *additive measures*.

Some measures are not additive over any dimension. Examples include prices and percentages, such as a discount percentage. Typically, you use the AVERAGE aggregate function for such measures, or you do not aggregate them at all. Such measures are called *non-additive measures*. Often, you can sum additive measures and then calculate non-additive measures from the additive aggregations. For example, you can calculate the sum of sales amount and then divide that value by the sum of the order quantity to get the average price. On higher

levels of aggregation, the calculated price is the average price; on the lowest level, it's the data itself—the calculated price is the actual price. This way, you can simplify queries.

For some measures, you can use SUM aggregate functions over all dimensions but time. Some examples include levels and balances. Such measures are called *semi-additive measures*. For example, if customer A has $2,000.00 in a bank account, and customer B has $3,000.00, together they have $5,000.00. However, if customer A had $5,000.00 in an account yesterday but has only $2,000.00 today, then customer A obviously does not have $7,000.00 altogether. You should take care how you aggregate such measures in a report. For time measures, you can calculate average value or use the last value as the aggregate.

> ✔ **Quick Check**
> - You are designing an accounting system. Your measures are debit, credit, and balance. What is the additivity of each measure?
>
> **Quick Check Answer**
> - Debit and credit are additive measures, and balance is a semi-additive measure.

## Additivity of Measures in SSAS

SSAS is out of the scope of this book; however, you should know some facts about SSAS if your data warehouse is the source for SSAS databases. SSAS has support for semi-additive and non-additive measures. The SSAS database model is called the *Business Intelligence Semantic Model (BISM)*. Compared to the SQL Server database model, BISM includes much additional metadata.

SSAS has two types of storage: *dimensional* and *tabular*. Tabular storage is quicker to develop, because it works through tables like a data warehouse does. The dimensional model more properly represents a cube. However, the dimensional model includes even more metadata than the tabular model. In BISM dimensional processing, SSAS offers semi-additive aggregate functions out of the box. For example, SSAS offers the LastNonEmpty aggregate function, which properly uses the SUM aggregate function across all dimensions but time, and defines the last known value as the aggregate over time. In the BISM tabular model, you use the Data Analysis Expression (DAX) language. The DAX language includes functions that let you build semi-additive expressions quite quickly as well.

## Many-to-Many Relationships

In a relational database, the many-to-many relationship between two tables is resolved through a third intermediate table. For example, in the AdventureWorksDW2012 database, every Internet sale can be associated with multiple reasons for the sale—and every reason can be associated with multiple sales. Figure 1-13 shows an example of a many-to-many rela-

tionship between *FactInternetSales* and *DimSalesReason* through the *FactInternetSalesReason* intermediate table in the AdventureWorksDW2012 sample database.



**FIGURE 1-13** A classic many-to-many relationship.

For a data warehouse in a relational database management system (RDBMS), this is the correct model. However, SSAS has problems with this model. For reports from a DW, it is you, the developer, who writes queries. In contrast, reporting from SSAS databases is done by using client tools that read the schema and only afterwards build a user interface (UI) for selecting measures and attributes. Client tools create multi-dimensional expression (MDX) queries for the SSAS dimensional model, and DAX or MDX queries for the SSAS tabular model. To create the queries and build the UI properly, the tools rely on standard Star or Snowflake schemas. The tools expect that the central table, the fact table, is always on the "many" side of the relationship.

A quick look at Figure 1-13 reveals that the *FactInternetSales* fact table is on the "one" side of its relationship with the *FactInternetSalesReason* fact table. SSAS with a BISM tabular model does not support many-to-many relationships at all in its current version. In SSAS with a BISM dimensional model, you can solve the problem by creating an intermediate dimension between both fact tables. You create it from the primary key of the *FactInternetSales* table. Let's call this dimension *DimFactInternetSales*. Then you put it on the "one" side of the relationships with both fact tables. This way, both fact tables are always on the "many" side of any relationship. However, you have to realize that the relationship between the *FactInternetSales* and the new *DimFactInternetSales* dimension is de facto one to one.

You can generate such intermediate dimensions in your data warehouse and then just inherit them in your SSAS BISM dimensional database. (Note that SSAS with BISM in a tabular model does not recognize many-to-many relationships, even with an additional intermediate dimension table.) This way, you can have the same model in your DW as in your BISM dimensional database. In addition, when you recreate such a dimension, you can expose it to end users for reporting. However, a dimension containing key columns only is not very useful for reporting. To make it more useful, you can add additional attributes that form a hierarchy. Date variations, such as year, quarter, month, and day are very handy for drilling down. You can get these values from the *DimDate* dimension and enable a drill-down path of year → quarter → month → day → sales order in this dimension. Figure 1-14 shows a many-to-many relationship with an additional intermediate dimension.



**FIGURE 1-14**  A many-to-many relationship with two intermediate tables.

Note that SSMS created the relationship between *DimFactInternetSales* and *FactInternetSales* as one to one.

**PRACTICE**  **Reviewing the AdventureWorksDW2012 Fact Tables**

The AdventureWorksDW2012 sample database has many types of fact tables as well, in order to show all possible measures. In this practice, you are going to review one of them.

**EXERCISE 1**  **Create a Diagram for an AdventureWorksDW2012 Fact Table**

In this exercise, you create a database diagram for a fact table and two associated dimensions.

1.  If you closed SSMS, start it and connect to your SQL Server instance. Expand the Databases folder and then the AdventureWorksDW2012 database.

2.  Right-click the Database Diagrams folder and select the New Database Diagram option.

3.  From the Add Table list, select the following tables (click each table and then click the Add button):

    ■  *DimProduct*

    ■  *DimDate*

    ■  *FactProductInventory*

Your diagram should look like Figure 1-15.



**FIGURE 1-15** FactProductInventory and related tables.

**Analyze Fact Table Columns**

In this exercise, you learn more details about the fact table in the schema you created in the previous exercise. Note that you have to conclude these details from the names of the measure columns; in a real-life project, you should check the content of the columns as well.

- Knowing how an inventory works, you can conclude that the UnitsIn and UnitsOut are additive measures. Using the SUM aggregate function for these two columns is reasonable for aggregations over any dimension.

- The UnitCost measure is a non-additive measure. Summing it over any dimension does not make sense.

- The UnitsBalance measure is a semi-additive measure. You can use the SUM aggregate function over any dimension but time.

Save the diagram using the name **Practice_01_03_ProductInventory**. Close the diagram and exit SSMS.

## Lesson Summary

- Fact tables include measures, foreign keys, and possibly an additional primary key and lineage columns.

- Measures can be additive, non-additive, or semi-additive.

- For many-to-many relationships, you can introduce an additional intermediate dimension.

## Lesson Review

Answer the following questions to test your knowledge of the information in this lesson. You can find the answers to these questions and explanations of why each answer choice is correct or incorrect in the "Answers" section at the end of this chapter.

1. Over which dimension can you not use the SUM aggregate function for semi-additive measures?

   A. Customer

   B. Product

   C. Date

   D. Employee

2. Which measures would you expect to be non-additive? (Choose all that apply.)

   A. Price

   B. Debit

   C. SalesAmount

   D. DiscountPct

   E. UnitBalance

3. Which kind of a column is not part of a fact table?

   A. Lineage

   B. Measure

   C. Key

   D. Member property

# Case Scenarios

In the following case scenarios, you apply what you've learned about Star and Snowflake schemas, dimensions, and the additivity of measures. You can find the answers to these questions in the "Answers" section at the end of this chapter.

## Case Scenario 1: A Quick POC Project

You are hired to implement a quick POC data warehousing project. You have to prepare the schema for sales data. Your customer's SME would like to analyze sales data over customers, products, and time. Before creating a DW and tables, you need to make a couple of decisions and answer a couple of questions:

1.  What kind of schema would you use?

2.  What would the dimensions of your schema be?

3.  Do you expect additive measures only?

## Case Scenario 2: Extending the POC Project

After you implemented the POC sales data warehouse in Case Scenario 1, your customer was very satisfied. In fact, the business would like to extend the project to a real, long-term data warehouse. However, when interviewing analysts, you also discovered some points of dissatisfaction.

### Interviews

Here's a list of company personnel who expressed some dissatisfaction during their interviews, along with their statements:

- **Sales SME**  "I don't see correct aggregates over regions for historical data."
- **DBA Who Creates Reports**  "My queries are still complicated, with many joins."

You need to solve these issues.

### Questions

1.  How would you address the Sales SME issue?

2.  What kind of schema would you implement for a long-term DW?

3.  How would you address the DBA's issue?

# Suggested Practices

To help you successfully master the exam objectives presented in this chapter, complete the following tasks.

## Analyze the AdventureWorksDW2012 Database Thoroughly

To understand all kind of dimensions and fact tables, you should analyze the AdventureWorksDW2012 sample database thoroughly. There are cases for many data warehousing problems you might encounter.

- **Practice 1**  Check all fact tables. Find all semi-additive measures.
- **Practice 2**  Find all hierarchies possible for the *DimCustomer* dimension. Include attributes in the dimension and attributes in the lookup *DimGeography* table.

## Check the SCD and Lineage in the AdventureWorksDW2012 Database

Although the AdventureWorksDW2012 database exemplifies many cases for data warehousing, not all possible problems are covered. You should check for what is missing.

- **Practice 1**  Is there room for lineage information in all dimensions and fact tables? How would you accommodate this information?

- **Practice 2**  Are there some important dimensions, such as those representing customers and products, that are not prepared for a Type 2 SCD solution? How would you prepare those dimensions for a Type 2 SCD solution?

# Answers

This section contains answers to the lesson review questions and solutions to the case scenarios in this chapter.

## Lesson 1

1. **Correct Answers: A and D**

   A. **Correct:** A Star schema typically has fewer tables than a normalized schema.

   B. **Incorrect:** The support for data types depends on the database management system, not on the schema.

   C. **Incorrect:** There are no specific Transact-SQL expressions or commands for Star schemas. However, there are some specific optimizations for Star schema queries.

   D. **Correct:** The Star schema is a de facto standard for data warehouses. It is narrative; the central table—the fact table—holds the measures, and the surrounding tables, the dimensions, give context to those measures.

2. **Correct Answer: C**

   A. **Incorrect:** The Star schema is more suitable for long-term DW projects.

   B. **Incorrect:** A normalized schema is appropriate for OLTP LOB applications.

   C. **Correct:** A Snowflake schema is appropriate for POC projects, because dimensions are normalized and thus closer to source normalized schema.

   D. **Incorrect:** An XML schema is used for validating XML documents, not for a DW.

3. **Correct Answers: B and D**

   A. **Incorrect:** Lookup tables are involved in both Snowflake and normalized schemas.

   B. **Correct:** Dimensions are part of a Star schema.

   C. **Incorrect:** Measures are columns in a fact table, not tables by themselves.

   D. **Correct:** A fact table is the central table of a Star schema.

## Lesson 2

1. **Correct Answer: B**

   A. **Incorrect:** This is Type 3 SCD management.

   B. **Correct:** This is how you handle changes when you implement a Type 2 SCD solution.

   C. **Incorrect:** Maintaining history does not mean that the content of a DW is static.

   D. **Incorrect:** This is Type 1 SCD management.

2. **Correct Answer: B**

   A. **Incorrect:** Attributes are part of dimensions.

   B. **Correct:** Measures are part of fact tables.

   C. **Incorrect:** Keys are part of dimensions.

   D. **Incorrect:** Member properties are part of dimensions.

   E. **Incorrect:** Name columns are part of dimensions.

3. **Correct Answer: B**

   A. **Incorrect:** You need to analyze the attribute names and content in order to spot the hierarchies in a Star schema.

   B. **Correct:** Lookup tables for dimensions denote natural hierarchies in a Snowflake schema.

   C. **Incorrect:** A Snowflake schema supports hierarchies.

   D. **Incorrect:** You do not need to convert a Snowflake to a Star schema to spot the hierarchies.

## Lesson 3

1. **Correct Answer: C**

   A. **Incorrect:** You can use SUM aggregate functions for semi-additive measures over the *Customer* dimension.

   B. **Incorrect:** You can use SUM aggregate functions for semi-additive measures over the *Product* dimension.

   C. **Correct:** You cannot use SUM aggregate functions for semi-additive measures over the *Date* dimension.

   D. **Incorrect:** You can use SUM aggregate functions for semi-additive measures over the *Employee* dimension.

2. **Correct Answers: A and D**

   A. **Correct:** Prices are not additive measures.

   B. **Incorrect:** Debit is an additive measure.

   C. **Incorrect:** Amounts are additive measures.

   D. **Correct:** Discount percentages are not additive measures.

   E. **Incorrect:** Balances are semi-additive measures.

3. **Correct Answer: D**

    **A.** **Incorrect:** Lineage columns can be part of a fact table.

    **B.** **Incorrect:** Measures are included in a fact table.

    **C.** **Incorrect:** A fact table includes key columns.

    **D.** **Correct:** Member property is a type of column in a dimension.

## Case Scenario 1

1. For a quick POC project, you should use the Snowflake schema.

2. You would have customer, product, and date dimensions.

3. No, you should expect some non-additive measures as well. For example, prices and various percentages, such as discount percentage, are non-additive.

## Case Scenario 2

1. You should implement a Type 2 solution for the slowly changing customer dimension.

2. For a long-term DW, you should choose a Star schema.

3. With Star schema design, you would address the DBA's issue automatically.

# Implementing Dynamic Packages

## Exam objectives in this chapter:

- Extract and Transform Data
  - Define connection managers.
- Load Data
  - Implement package logic by using SSIS variables and parameters.

When you are developing Microsoft SQL Server Integration Services (SSIS) packages, it is a good practice to make each task or transformation as dynamic as possible. This enables you to move your packages from one environment to another (for example, from development to a test environment, and then to a production environment) without opening and changing the package. You can also configure your packages to set different properties at run time.

To do this, you can use the new features in SQL Server 2012, such as parameters and project-level connection managers, or you can use the package configurations that were first made available in the package deployment model in earlier versions of SQL Server. This chapter discusses both possibilities for designing and configuring your package to dynamically set values at run time. Using dynamic packages eliminates the need to make changes as you move from one environment to the other or to open the project when you want to run your packages using different property or variable values.

## Lessons in this chapter:

- Lesson 1: Package-Level and Project-Level Connection Managers and Parameters
- Lesson 2: Package Configurations

# Before You Begin

To complete this chapter, you must have:

- Basic knowledge of SQL Server 2012 SISS control flow and data flow features and components.
- Experience working with SQL Server 2012 Management Studio (SSMS).
- Experience working with SQL Server Data Tools (SSDT) or SQL Server Business Intelligence Development Studio (BIDS).
- An understanding of the basics of file copying and security.
- Knowledge of system environment variables.

# Lesson 1: Package-Level and Project-Level Connection Managers and Parameters

In SQL Server 2012, SSIS introduces parameters and project-level connection managers. Parameters enable you to assign values to properties within packages at the time of package execution. Project-level connection managers allow you to define the data source connection once and use it in all packages that are part of the project. Both features are available only to projects developed for the project deployment model. This means that SSIS packages created with prior versions of SQL Server have to be upgraded to the new project deployment model if you want to use these new features.

Both functionalities are a replacement for the package configurations used in previous versions of SQL Server. In SQL Server 2012, SSIS also takes advantage of *build configurations* from Microsoft Visual Studio, which enable you to define multiple configurations and set the values of your parameters per configuration. This makes it easier to debug and deploy packages in SQL Server Data Tools against different SSIS servers.

---

**After this lesson, you will be able to:**

- Understand the difference between package-level and project-level connection managers.
- Implement parameters.
- Use property expressions to make your packages dynamic.

**Estimated lesson time: 30 minutes**

---

# Using Project-Level Connection Managers

Project-level connection managers allow you to set up a connection manager on the project level. This means that you can define a connection for the whole project and use it in all packages that are part of the project. In prior versions of SQL Server, the connections were always contained within each package.

To create a project-level connection, define a new connection manager by right-clicking Connection Managers in Solution Explorer under your project and selecting the New Connection Manager option. An alternative method is to convert an existing package-level connection to a project-level connection. You can do this by opening the package, right-clicking the connection you want to convert in the Connection Managers window, and selecting Convert To Project Connection, as shown in Figure 9-1.



**FIGURE 9-1** Converting a package-level connection to a project-level connection.

Project-level connections are a very useful new feature when it comes to developing and maintaining your packages. Note that in cases when you need to have a connection available only for the specific package, you can still create a connection at the package scope.

# Parameters

In SQL Server 2012, SSIS introduces parameters. Parameters allow you to assign values to properties within packages at the time of package execution. They can also be used in SSIS expressions—for example, to use an Expression task to set a variable value based on the specific value of the parameter. There are two types of parameters: project parameters, which are created on the project level, and package parameters, which are created at the package level. You can use and set parameters only in projects developed for the project deployment model. When you are using the project deployment model, projects are deployed to the *Integration Services catalog*. Details regarding deployment possibilities are explained later in Chapter 11, "Installing SSIS and Deploying Packages."

## Using Parameters

You can use a single parameter to assign a value to multiple package properties, and you can use a parameter in multiple SSIS expressions. Depending on where the project is in the project deployment life cycle, you can assign up to three different types of values to each parameter. The three types are listed in the order in which they can be applied to the parameter:

- **Design default value**  The default value assigned when the project is created or edited in SQL Server Data Tools (SSDT). This value persists with the project.
- **Server default value**  The default value assigned during project deployment or later, while the project resides in the SSIS catalog. This value overrides the design default.
- **Execution value**  The value that is assigned in reference to a specific instance of package execution. This assignment overrides all other values but applies to only a single instance of package execution.

Note that the last two types of values become relevant after the SSIS project has been deployed to the Integration Services catalog.

## Defining Parameters

You add project or package parameters by using SSDT. Usually you create a project parameter by selecting Project.params in Solution Explorer under your project name, as shown in Figure 9-2. To add a package parameter, you must first open the package and then select the parameters tab in the package design area.



**FIGURE 9-2** Solution Explorer with Project.params selected.

When you open the package or project parameters window, you can define a new parameter by clicking the Add Parameter icon (the first icon on the left; it looks like a blue cube). Figure 9-3 shows a parameter window with one parameter called pTK463DWConnectionString.



**FIGURE 9-3**  The parameters window.

When you create a parameter in SSDT, there are several properties to specify:

- **Name**   The name of the parameter. The first character of the name must be a letter or an underscore.

- **Data type**   The data type of the parameter.

- **Value**   The default value for the parameter assigned at design time. This is also known as the design default.

- **Sensitive**   If the value of this property is True, parameter values are encrypted in the catalog and appear as NULL when viewed with Transact-SQL or SQL Server Management Studio.

- **Required**   Requires that a value other than the design default be specified before the package can be executed.

- **Description**   For maintainability, the description of the parameter. In SSDT, set the parameter description in the Visual Studio Properties window when the parameter is selected in the applicable parameters window.

---

*EXAM TIP*

**Parameter design values are stored in the project file.**

---

You can edit parameters in the list in the parameter window, or you can use the Properties window to modify the values of parameter properties. You can delete a parameter by using the Delete Parameter toolbar button. By using the Add Parameters To Configurations toolbar button (the last button on the right), you can specify a value for a parameter for a specific build configuration that is used only when you execute the package in SQL Server Data Tools. Build configurations are explained in the next topic.

Another approach to creating a parameter is to do so implicitly by right-clicking a control flow task or a connection manager and selecting the Parameterize option. A Parameterize dialog box opens, as shown in Figure 9-4. Here you can specify which property should be dynamically evaluated at run time and whether the parameter should be created (the latter is specified by selecting the Create New Parameter option).



**FIGURE 9-4**  The Parameterize dialog box.

*EXAM TIP*

**You cannot change the value of a parameter while a package is running.**

# Build Configurations in SQL Server 2012 Integration Services

Build configurations in Visual Studio provide a way to store multiple versions of solution and project properties. The active configuration can be quickly accessed and changed, allowing you to easily build multiple configurations of the same project. Two levels of build configurations can be defined in Visual Studio: solution configurations and project configurations.

The new project deployment model in SQL Server 2012 takes advantage of more possibilities for build configurations from Visual Studio than earlier versions did. You can now set project and package parameters to get values from build configurations. Also, some of the project-level properties can be set via build configurations (for example, deployment server name and server project path).

## Creating Build Configurations

A project can have a set of defined project properties for every unique combination of a configuration and platform. You can create an additional project or solution configuration by using the Configuration Manager:

1. Select the Configuration Manager option from the Configurations drop-down list on the Standard toolbar to open the Configuration Manager dialog box. Alternatively, you can first open the project's Property Pages dialog box (right-click the project name in Solution Explorer and select Properties), as shown in Figure 9-5, and click the Configuration Manager button.



**FIGURE 9-5** The project Property Pages dialog box.

2. In the Active Solution Configuration drop-down list, select New.

3. In the New Solution Configuration dialog box, enter the name of the solution configuration you would like to create. Select the Create New Project Configurations check box to also create the project configuration.

You can now change the active configuration directly from the Solution Configurations drop-down list on the Standard toolbar or from the Configuration Manager dialog box.

## Using Build Configurations

You can bind parameter values to build configurations by using the parameter window:

1. Open the project or package parameters window.

2. In the parameters window, select the last toolbar button on the left (Add Parameters To Configurations).

3. In the Manage Parameter Values dialog box that appears, you can add values for each parameter for each configuration. Click the Add button to add a parameter to the configuration.

4. In the Add Parameters window, select parameters and click OK.

5. Set the appropriate values of parameters for your configurations, as shown in Figure 9-6.

Figure 9-6 shows that the pNoOfRows parameter will have a value of 10,000 when the package is executed using the Production configuration and 100 when using the Development configuration. This means that if you have multiple parameters and need to change the value of the parameters at design time when you are developing or debugging the SSIS project, you just need to switch between build configurations to have all parameter values changed at once.



**FIGURE 9-6** Managing parameter values.

You can also use build configurations to set build, deployment, and debugging configuration properties for a project. To assign different project properties based on configuration, right-click the project in Solution Explorer and select Properties. Figure 9-7 shows the deployment properties inside the project's Property Pages dialog box. If you have to deploy your project to multiple servers, you can set different values for the Server Name and Server Project Path properties for each configuration, so that you can quickly switch between deployment environments at design time by using configurations.

**FIGURE 9-7** Deployment properties in a project's Property Pages dialog box.

## Property Expressions

In previous chapters, you saw some examples of setting connection managers or control flow task properties at run time. The SSIS expressions used to update properties of the control flow during package execution are called *property expressions*. You can apply a property expression in two ways. First, you can set the property as an expression through the properties window by clicking the ellipsis (…) button of the Expressions property. This will open the Property Expression Editor, as shown in Figure 9-8. In this dialog box, you can select a property from the drop-down list and then type an expression.



**FIGURE 9-8** The Property Expressions Editor.

The second way to set property expressions is to use the task or container editors, which offer one or more ways to set properties through expressions. Figure 9-9 shows the Execute SQL Task Editor. On the Expressions tab, you can specify property expressions.



**FIGURE 9-9** The Property Expressions Editor opened from the Expressions tab.

---

✔ **Quick Check**

- **When are property expressions evaluated as a package is running?**

**Quick Check Answer**

- **Unlike parameters that are read at the start of package execution, property expressions are updated when the property is accessed by the package during package execution. A property expression can change the value of a property in the middle of package execution, so that the new value is read when the property is needed by the package.**

PRACTICE  **Implementing Parameters**

In this practice, you will use parameters to make your packages dynamic. In the first exercise
you will parameterize the connection string, and in the second exercise you will use a param-
eter value to filter the source query in the data flow task. The third exercise focuses on setting
up an additional build configuration to test project execution against another database.

If you encounter a problem completing an exercise, you can install the completed projects
from the Solution folder that is provided with the companion content for this chapter and
lesson.

### EXERCISE 1   Set a Parameter for a Connection String

In this exercise, you parameterize a connection string by using a project parameter.

1. If you are missing the database objects from Chapter 5, "Designing and Implementing
   Data Flow," execute the needed SQL code from that chapter to have all the stage and
   dimension tables available in the TK463DW database.

2. Start SQL Server Data Tools, open the TK 463 Chapter 9 project in the Starter folder,
   and then open the FillStageTablesParameters.dtsx package.

   Notice that this package is using two connections, AdventureWorks2012 and TK463DW.

3. In the Connection Managers window, right-click the TK463DW connection and select
   Convert To Project Connection. You have changed the connection from the package
   level to the project level and can now see this connection in the Solution Explorer win-
   dow in the Connection Managers folder. The TK463DW connection can now be used
   by any other package within the same SSIS project.

4. Right-click Project.params in Solution Explorer, and then click Open or double-click
   Project.params to open it.

5. Click the Add Parameter button on the toolbar.

6. Name the parameter by setting the Name property to **pTK463DWConnString**
   and set the Data Type property to String. In the Value property field, type **Data
   Source=localhost;Initial Catalog=TK463DW;Provider=SQLNCLI11.1;Integrated
   Security=SSPI;**.

7. Close the Project.params window.

8. Inside the FillStageTablesParameters package, right-click the (project) TK463DW connection in the Connection Managers window and select Parameterize.

9. In the Parameterize dialog box, select the Use Existing Parameter check box and select the pTK463DWConnString parameter in the drop-down list. Notice that the project parameter name starts with $Project::. Click OK to close the Parameterize dialog box.

10. Look at the (project) TK463DW connection in the Connection Managers window and notice a small icon next to it reflecting that this connection is parameterized.

**EXERCISE 2**   **Use a Parameter in the Data Flow Task**

In this exercise, you create a package-level parameter that will be used to filter source data. You will use this parameter in the data flow task to filter only rows from the source for which the year of the modified date is equal to or greater than the parameter value.

1. If necessary, start SQL Server Data Tools, open the TK 463 Chapter 9 project, and then open the FillStageTablesParameters.dtsx package from the previous exercise for editing.

2. Select the Parameters tab and click the Add Parameter button on the toolbar.

3. Name the parameter by setting the Name property to **pYear** and set the Data Type property to Int16. For the Value property, enter **2002**.

4. Click the Data Flow tab and open the Person OLE DB Source adapter.

5. In the OLE DB Source Editor, change the data access mode to SQL Command and enter the following SELECT statement to retrieve the necessary rows and use a parameter placeholder inside the query.

```
SELECT
  BusinessEntityID,
  PersonType,
  NameStyle,
  Title,
  FirstName,
  MiddleName,
  LastName,
  Suffix,
  EmailPromotion,
  AdditionalContactInfo,
  Demographics,
  rowguid,
    ModifiedDate
FROM
  Person.Person
WHERE
  YEAR(ModifiedDate) >= ?
```

6. Click the Parameters button to open the Set Query Parameters dialog box.

7. For the Variables property, select the $Package::pYear parameter as the source for the query parameter. Click OK twice to close the window and the OLE DB Source Editor.

8. Execute the package and observe the number of rows displayed in the data flow area.

9. Change the parameter value to **2008** and execute the package. Notice that fewer rows are read from the OLE DB Source adapter in the data flow area.

**EXERCISE 3** Use Build Configurations

In this exercise, you create an additional database, TK463DWProd, and then create a new build configuration in Visual Studio to use this database when running the SSIS package build from the previous exercise.

1. Start SSMS and connect to your SQL Server instance. Open a new query window by clicking the New Query button.

2. You will create a new database and the *stg.Person* table so that you can execute the SSIS package created in the previous exercise. Execute the provided T-SQL code to create the database and the table.

```
USE master;
IF DB_ID('TK463DWProd') IS NOT NULL
  DROP DATABASE TK463DWProd;
GO
CREATE DATABASE TK463DWProd
 ON PRIMARY
 (NAME = N'TK463DWProd', FILENAME = N'C:\TK463\TK463DWProd.mdf',
  SIZE = 307200KB , FILEGROWTH = 10240KB )
 LOG ON
 (NAME = N'TK463DWProd_log', FILENAME = N'C:\TK463\TK463DWProd_log.ldf',
  SIZE = 51200KB , FILEGROWTH = 10%);
GO
ALTER DATABASE TK463DWProd SET RECOVERY SIMPLE WITH NO_WAIT;
GO
USE TK463DWProd;
GO
CREATE SCHEMA stg AUTHORIZATION dbo;
GO
CREATE TABLE stg.Person
(
 BusinessEntityID INT         NULL,
 PersonType       NCHAR(2)    NULL,
 Title            NVARCHAR(8) NULL,
 FirstName        NVARCHAR(50) NULL,
 MiddleName       NVARCHAR(50) NULL,
 LastName         NVARCHAR(50) NULL,
 Suffix           NVARCHAR(10) NULL,
 ModifiedDate     DATETIME    NULL
);
```

3. If necessary, start SQL Server Data Tools, open the TK 463 Chapter 9 project, and then open the FillStageTablesParameters.dtsx package from the previous exercise for editing.

4. Select the Configuration Manager option in the Solution Configurations drop-down list.

5. In the Configuration Manager dialog box, select the New option in the Active Solution Configuration drop-down list.

6. In the New Solution Configuration dialog box, enter **Production** as the name of the configuration and click OK. Close the Configuration Manager dialog box.

7. Right-click Project.params in Solution Explorer, and then click Open (or double-click Project.params).

8. Click the Add Parameter To Configurations button on the toolbar.

9. In the Manage Parameter Values dialog box, click Add. Select the pTK463DWConnString parameter and click OK. Notice that the value of the parameter was copied to both configurations.

10. Change the Production configuration to use the newly created database, TK463DWProd. The value should look like this—**Data Source=localhost;Initial Catalog=TK463DWProd;Provider=SQLNCLI11.1;Integrated Security=SSPI;**.

11. Click OK and save the SSIS project to store the values for the configurations.

12. Execute the package first under the Development configuration and then under the Production configuration by selecting the different values in the Solution Configurations drop-down list. Look at the *stg.Person* table in the TK463DWProd database to see if it contains data.

---

*NOTE* **CONTINUING WITH PRACTICES**

**Do not exit SSMS or SSDT if you intend to continue immediately with the next practice.**

---

## Lesson Summary

■ Use parameters to set up connection properties at run time.

■ Parameters and project-level connection mangers can only be used with the new project deployment model introduced with SSIS in SQL Server 2012.

■ Use property expressions to change the control flow properties at run time.

## Lesson Review

Answer the following questions to test your knowledge of the information in this lesson. You can find the answers to these questions and explanations of why each answer choice is correct or incorrect in the "Answers" section at the end of this chapter.

1. Which parameter types are available in SSIS in SQL Server 2012? (Choose all that apply.)

    A. Project-level parameters

    B. Solution-level parameters

    C. Control flow–level parameters

    D. Package parameters

2. Which properties can be set by using build configurations? (Choose all that apply.)

    A. Parameter values

    B. Variable values

    C. Control flow task properties

    D. The Deployment Server Name property

3. Which properties can be set by using property expressions? (Choose all that apply.)

    A. SQL statement for the Execute SQL task

    B. Variable values

    C. Data flow task properties

    D. The Lookup transformation SqlCommand property

# Lesson 2: Package Configurations

In versions of SSIS before SQL Server 2012, you had to use package configurations to update properties, variable values, and connections at run time. You could have the package look to an external source for configuration information that changed the settings within the package when it executed. Package configurations are optional, but in real-life scenarios they are almost mandatory, because they provide a way for you to update package settings without having to open each package in Business Intelligence Development Studio (BIDS) or SSDT. For example, by using package configurations, you can maintain connection strings and variable settings for all of your packages in a single location.

> **After this lesson, you will be able to:**
>
> ■ Implement package configurations.
>
> **Estimated lesson time: 40 minutes**

# Implementing Package Configurations

When you are executing a package in a package deployment model, the first action the package takes is to look at its configurations and overwrite the package's current settings with the new settings from the configurations. Common elements that are configured by using package configurations are:

- **Connection properties**   These include properties that set the connection string, the server name, the user name, and the password.

- **Package variable properties**   You can set variable values, variable descriptions, and the Raise Change Event property.

- **Package properties**   These include any property on the package level, such as checkpoint and security settings.

Before you can enable package configuration, you must convert your SSIS project to the package deployment model. By default, a new package in the SQL Server 2012 version of SSIS is set up for the project deployment model. You can change this by selecting Convert To Package Deployment Model under the project's name on the main toolbar. Note that you can convert only projects that do not have any parameters or project-level connection managers defined.

By default, each package has its package configuration turned off. To enable and set up configurations, you use the Package Configuration Organizer, with which you can perform the following tasks:

- Enable or disable a package's package configurations

- Add or remove configurations assigned to the package

- Define the order in which the configurations are applied

To open the Package Configurations Organizer, open the package for which you want to turn on configurations, and then choose SSIS Configurations from the SSIS menu. To enable configurations, select the Enable Package Configurations check box at the top of the dialog box. Figure 9-10 shows the Package Configurations Organizer dialog box with package configurations enabled.

**FIGURE 9-10**  The Package Configurations Organizer.

## Creating a Configuration

To create a new configuration, click the Add button in the Package Configurations Organizer dialog box to start the Package Configuration Wizard. First you must specify the configuration type by selecting the appropriate value from the drop-down list, as shown in Figure 9-11. SSIS supports different package configuration types; Table 9-1 describes the configuration types you can use.

**FIGURE 9-11** Selecting a configuration type by using the Package Configuration Wizard.

**TABLE 9-1** Package Configuration Types

| Type | Description |
|---|---|
| XML Configuration File | Stores configuration settings in an XML file in the file system. Use this option if you are comfortable working with configuration files and your project re-quirements let you store configuration information in a file system file. Note that you can store multiple configurations in a single XML file. |
| Environment Variable | Saves the configuration information inside the system's global variables col-lection, which is called an *environment variable.* Only one property can be stored in each Environment Variable configuration. |
| Registry Entry | Lets you save package properties and settings in your computer's registry. |
| Parent Package Variable | Provides a way for your package to inherit the value of a variable from a parent package. When a package is executed from another SSIS package by using the Execute Package task, the values of its variables are available to the child package through the Parent Package Variable configuration. With this configuration type, you can choose only one package property setting at a time. |
| SQL Server | Stores configuration settings in a SQL Server table. You can store multiple configurations in a single table. |

Choose the most appropriate configuration for your environment and your project requirements. Ensure that you consider how the package will be supported in a production environment and how other technologies are supported and configured. Take care to evaluate any security and compliance requirements when you are storing connection information such as server name, user name, or password information.

> **IMPORTANT** **PARAMETERS OR PACKAGE CONFIGURATIONS?**
>
> If you are using SQL Server 2012, use the new project deployment model with parameters and project-level connection managers to support moving your solution from one environment to another. These new features provide better package management and flexibility in package development compared to package configurations. These new features are positioned as an evolution of Integration Services deployment and administration in SQL Server.

The most commonly used configuration types are the XML Configuration File and SQL Server configurations. The next section looks more closely at each of these types.

## Creating an XML File Configuration

When you choose the XML Configuration File type, you can specify the location for your configuration file. There are two ways to specify the location of the XML file:

- Enter the location of the file directly in the Configuration File Name text box. Use this when you intend to always use the same location and name for your configuration file.

- Use an environment variable that contains the location and name of the configuration file. To use this approach, you must create a system variable in your computer's system properties. The value of the variable must contain the full path, name, and extension of the file.

  Using an environment variable for the file location pointer is called the *indirect file location approach* and is very valuable if your XML file location or file name might change in the future or already changes between environments. If you choose to use the environment variable, be sure to add it to the servers on which the package will run.

As with all of the configuration types, more than one package can use the same XML Configuration File. If you have several packages that have common properties, such as connection strings, you might want to have all of them use one XML file for configuration.

After you have defined the location and name of the file, you define the server settings and properties that the XML Configuration File should contain. Because these are common among all configuration types, this chapter reviews the SQL Configuration setup before describing the server settings and property definitions.

## Creating a SQL Server Configuration

To store your package configurations in a SQL Server table, select SQL Server from the Configuration Type drop-down list in the Package Configuration Wizard. Using SQL Server as the storage mechanism for your configurations requires a different group of settings than those used by the other configuration types, such as XML Configuration File. Figure 9-12 shows the SQL Server configuration options available for setting up configurations.



**FIGURE 9-12** The Package Configuration Wizard for a SQL Server table configuration.

Just as with the XML Configuration File type, you can specify an environment variable as the location of your configuration (for example, the data source name for the SQL Server configuration), or you can specify configuration settings directly. There are three settings that define the table location details:

- **Connection**   This must be a SQL Server–based connection that sets the server and database in which your configurations will be stored and read. If you did not define the connection you need, you can click New next to Connection to open the Configure OLE DB Connection Manager dialog box.

- **Configuration Table**   This is the name of the table in which the configurations will reside. This table has predefined column name and data type definitions that cannot be changed. To create the table, you click New next to the Configuration Table text box to open the Create Table dialog box, in which you can change the name of the table and execute the table creation statement for the connection that you specified in the previous setting.

- **Configuration Filter**   Multiple SQL Server configurations can share the same table, and you can specify the configuration you want by using the Configuration Filter drop-down list. You can enter a new filter or use an existing one. The name you select or enter for this property is used as a value in the Configuration Filter column in the underlying table.

## Adding Properties to Your Configuration

No matter which SSIS configuration type you are using, you can select Properties To Export on the next page of the wizard to select the SSIS package and object properties that are to be used in the configuration. After you define the configuration type properties in the Package Configuration Wizard, click Next.

At this point, SSIS prompts you to verify whether configuration entries already exist for the configuration type you selected. If they do, SSIS prompts you to either reuse the configuration entries or overwrite them. If you see this dialog box, you will probably want to share the existing configurations between packages. If you do, click the Reuse Existing button. If you want to clear the existing entries and create new ones, click Overwrite.

If configuration entries do not already exist in this configuration, or if you clicked Overwrite, you will see the Select Properties To Export page, as shown in Figure 9-13.

**FIGURE 9-13** The Select Properties To Export page of the Package Configuration Wizard.

The Select Properties To Export page uses a tree view structure of your package properties, allowing you to select the properties for the SSIS configuration you have selected. Properties are grouped within the following folders:

- **Variables**   Lists all of the package variables you can select for configuration entries, along with their properties.
- **Connection Managers**   Lists all of the package connections, from which you can choose the specific properties for your connections.
- **Log Providers**   Lets you dynamically set the log configuration.
- **Properties**   Displays a list of all package-level properties that you can use to configure your package.
- **Executables**   Contains the tree structure of your tasks and containers. By navigating through this tree, you can configure the specific properties of your tasks and containers.

If you are using an XML Configuration File, SQL Server, or Registry Entry configuration, you can set multiple configuration properties at one time by selecting multiple property check boxes.

## Sharing, Ordering, and Editing Your Configurations

If you have several configurations in your list, you can define the order in which configurations are applied in a package. The configurations are called in the order in which they are listed in the Package Configuration Organizer. This is an important consideration if you have multiple configurations that will update the same property or if you have configurations that have a dependency on prior configurations. For example, you might have a configuration that updates a connection string, which is then used as the location of the configuration entries in a second configuration. Note that the last-applied property update will be the value that is used in the package.

A common approach is to share configurations between packages. If you do this, you might have configuration entries that apply to one package and not another. This does not affect package execution, but you will receive a warning to indicate that a configuration property does not exist in the package.

As a final note, you can modify all SSIS configuration entries you have made by simply editing the file, SQL Server, registry, or environment variable value. Look for the Configured Value property and change it as necessary.

---

**REAL WORLD**   **MULTIPLE CONFIGURATIONS**

In previous versions of SQL Server, the use of configurations was almost mandatory for moving from a development to a test environment and then to a production environment. Because user names and passwords for connection strings should not be stored as clear text in an XML file, most of the clients store configuration properties in SQL Server. To have the flexibility to move from one environment to another, you also need to put the information about the SQL Server instance used for storing configurations in a configuration setup. A common solution is to first use an XML configuration that contains the location of the file stored as an environment variable and that includes only the connection string property for the SQL Server configuration that will hold other configuration values. Then you create a second configuration that is a SQL Server configuration, and you use that configuration for all configuration values.

---

**PRACTICE**   **Using Package Configurations**

In this practice, you will create an XML Configuration File and share it between two packages.

If you encounter a problem completing an exercise, you can install the completed projects from the Solution folder that is provided with the companion content for this chapter and lesson.

**EXERCISE** Create an XML Configuration

In this exercise, you use SSIS configurations to create an SSIS XML Configuration File that contains the connection string property of the AdventureWorks2012 and TK463 databases. You then share this configuration with another package.

1. Start SQL Server Data Tools, open the TK 463 Chapter 9 project, and open the FillStageTables_1.dtsx package.

2. Now you need to convert the project to a package deployment model. Click Project on the toolbar and select Convert To Project Deployment Model. Click OK in the dialog box.

3. In the Convert To Package Deployment Model dialog box, every step should show the Result value as Passed. Click OK to convert the project to a package deployment model.

4. Choose Package Configurations from the SSIS menu.

5. Select the Enable Package Configurations check box in the Package Configurations Organizer dialog box.

6. Click Add to create a new configuration.

7. Click Next on the Welcome To The Package Configuration Wizard page.

8. In the Configuration Type drop-down list, select XML Configuration File.

9. Click the Browse button next to the Configuration File Name box, browse to the \Chapter09\Lesson2\Starter\ installed files folder, and then type **SSIS_Conn.dtsConfig**. Click Save to save the file name and path.

10. Click Next in the Package Configuration Wizard to go to the Select Properties To Export page.

11. Under Objects, expand the Connection Managers folder, expand the Properties folder for the AdventureWorks2012 connection, and select the check box next to the ConnectionString property. Repeat the process for the TK463DW connection. Click Next.

12. Name the configuration **MainXMLConfiguration** and close the Configuration Organizer dialog box.

13. Save and close the FillStageTables_1.dtsx package.

14. Open the FillStageTables_2.dtsx package and repeat steps 4 through 9. Select the file you created in step 9, and click Next. You will be prompted to overwrite the existing file or reuse the configuration that it contains. Click the Reuse Existing button.

15. Name the configuration **MainXMLConfiguration** and close the Configuration Organizer dialog box.

16. Save and close the FillStageTables_2.dtsx package.

17. Execute the packages and try to change the XML file so that the connection string points to the TK463DWProd database created in the previous lesson, and execute the packages again.

## Lesson Summary

- Package configurations are available in the package deployment model.
- Use package configurations if you are using previous versions of SSIS to set connection properties at run time.
- Use a combination of XML and SQL Server configurations to provide additional portability for your packages.

## Lesson Review

Answer the following questions to test your knowledge of the information in this lesson. You can find the answers to these questions and explanations of why each answer choice is correct or incorrect in the "Answers" section at the end of this chapter.

1. Which configuration types can you use to store configuration values? (Choose all that apply.)

   A. XML Configuration File

   B. SQL Server

   C. Any relational database system

   D. Registry entry

2. On which objects can you set dynamic properties by using package configurations? (Choose all that apply.)

   A. Parameters

   B. Variables

   C. Data flow transformations

   D. The Sequence Container task

3. Which SSIS elements can be configured by using a package configuration? (Choose all that apply.)

   A. Connection properties

   B. Package variable properties

   C. Parameter properties

   D. Package properties

# Case Scenario

In the following case scenarios, you apply what you've learned about implementing dynamic packages. You can find the answers to these questions in the "Answers" section at the end of this chapter.

## Case Scenario: Making SSIS Packages Dynamic

You are creating a set of SSIS packages to move data from flat files and different databases to a data warehouse. Because of strict development, test, and production procedures, your SSIS project must support these possibilities:

1. You need to test different deployment options for your SSIS packages in the development phase by using SSDT.

2. You have a development, test, and production environment and would like to minimize any changes to packages when they are ready for the test phase.

3. You need to parameterize the location of the flat files but still be able to dynamically set the correct file name when using the Foreach Loop container to read all the files inside the specified folder.

How would you address these issues? What would your solution look like?

# Suggested Practices

To help you successfully master the exam objectives presented in this chapter, complete the following tasks.

## Use a Parameter to Incrementally Load a Fact Table

In order to practice what you have learned in this chapter, you will create a package to load a fact table that will store Internet sales data based on the AdventureWorks2012 database.

- **Practice 1**  Create the necessary package to load the fact table. At each execution of the package, do the full load and truncate the table at the beginning.

- **Practice 2**  Add a project parameter that will accept the incremental date and use it to load only data that is newer than the supplied value. Replace the TRUNCATE T-SQL statement at the beginning with the DELETE T-SQL statement.

- **Practice 3**  Parameterize all connection strings.

# Answers

This section contains answers to the lesson review questions and solutions to the case scenario in this chapter.

## Lesson 1

1. **Correct Answers: A and D**

   A. **Correct:** Project-level parameters are available.

   B. **Incorrect:** Solution-level parameters are not available.

   C. **Incorrect:** Parameters can be defined on the project or package level.

   D. **Correct:** SSIS supports package-level parameters.

2. **Correct Answers: A and D**

   A. **Correct:** Parameter values can be bound to build configurations.

   B. **Incorrect:** Variables cannot be bound to build configurations.

   C. **Incorrect:** Control flow task properties cannot be bound to build configurations.

   D. **Correct:** Project-level properties can be set for each build configuration.

3. **Correct Answers: A and C**

   A. **Correct:** You can set the SQL statement of the Execute SQL task at run time by using property expressions.

   B. **Incorrect:** Variable properties cannot be set by using property expressions.

   C. **Correct:** General data flow task properties can be set by using property expressions.

   D. **Incorrect:** You can only change the SQL command of the Lookup task when you are using property expressions for the data flow task that has a Lookup task inside it.

## Lesson 2

1. **Correct Answers: A, B, and D**

   A. **Correct:** You can use XML as a configuration file.

   B. **Correct:** You can store the configuration values in a SQL Server database.

   C. **Incorrect:** You can only use a SQL Server database.

   D. **Correct:** You can use registry entries to store configuration values.

2. **Correct Answers: B and D**

    A. **Incorrect:** Parameters can be used in the project deployment model.

    B. **Correct:** You can set variable properties.

    C. **Incorrect:** You can only set data flow task properties, and not for a specific transformation.

    D. **Correct:** You can dynamically set properties for a Sequence Container task by using configurations.

3. **Correct Answers: A, B, and D**

    A. **Correct:** Connection properties can be set by using package configurations.

    B. **Correct:** You can set variable properties.

    C. **Incorrect:** You cannot use parameters when using package configurations.

    D. **Correct:** You can dynamically set properties for the package by using package configurations.

## Case Scenario

1. Add parameters for the connection strings. Create development, test, and production build configurations and bind parameters to each of them with different data source values. This will allow you to execute the package from SSDT against different configurations without manually changing the value of parameters.

2. Add parameters and parameterize all needed connection managers.

3. Create a parameter that will hold the value of the file location folder. Create a new variable to store the current file name. Using SSIS expression language, set the new variable value to combine the value from the parameter with the value you get from the Foreach Loop container. Use the new variable as the property expression for the Foreach Loop to dynamically change the fully qualified file name while the package is running.

# Index

## Symbols

## A

# B

# D

# E

# F

# N

# P

# S

# U