

# **TDWI Data Modeling**

Data Analysis and Design for BI and Analytics Solutions



Previews of TDWI course books offer an opportunity to see the quality of our material and help you to select the courses that best fit your needs. The previews cannot be printed.

TDWI strives to provide course books that are contentrich and that serve as useful reference documents after a class has ended.

This preview shows selected pages that are representative of the entire course book; pages are not consecutive. The page numbers shown at the bottom of each page indicate their actual position in the course book. All table-of-contents pages are included to illustrate all of the topics covered by the course.

# This page intentionally left blank.

TDWI takes pride in the educational soundness and technical accuracy of all of our courses. Please send us your comments—we'd like to hear from you. Address your feedback to:

info@tdwi.org

Publication Date:

April 2021

© Copyright 2005-2021 by TDWI. All rights reserved. No part of this document may be reproduced in any form, or by any means, without written permission from TDWI.

| S   | Module 1   | Data Modeling Concepts         | 1-1 |
|-----|------------|--------------------------------|-----|
|     | Module 2   | Business Data Models           | 2-1 |
|     | Module 3   | Logical Data Models            | 3-1 |
| E   | Module 4   | Physical Data Models           | 4-1 |
| Z   | Module 5   | Summary and Conclusion         | 5-1 |
| ŭ   | Appendix A | The Data University Case Study | A-1 |
| L   | Appendix B | Exercises                      | B-1 |
| 0   | Appendix C | Solutions                      | C-1 |
| Щ   | Appendix D | Bibliography and References    | D-1 |
| ABL |            |                                |     |

# CTIVE COUR

To learn:

- Differences in modeling techniques for business transactions, business events, and business metrics
- ✓ Different types of data and their implications
- How modeling processes differ based on analytics objectives and data management technology
- ✓ Application of business context to modeling activities
- ✓ The role of business requirements in BI data modeling
- ✓ The role of source data analysis in data modeling
- ✓ Use of normalized modeling techniques for data warehouse analysis and design
- ✓ Use of dimensional modeling techniques for BI and data mart analysis and design
- ✓ The roles of generalization and abstraction in data warehouse design
- ✓ The roles of identity and hierarchy management in data model design
- ✓ How time-variant data is represented in data models
- Implementation and optimization considerations for data stores



# Module 1

# Data Modeling Concepts

| Торіс  | Page |
|--|------|
| The Data Modeling Life Cycle                 | 1-2  |
| Kinds of Data Systems                        | 1-6  |
| Data Taxonomy                                | 1-8  |
| Data Modeling Framework for BI and Analytics | 1-12 |

© TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY. 1-1

# The Data Modeling Life Cycle Where Data Modeling Begins and Ends



# The Data Modeling Life Cycle Where Data Modeling Begins and Ends

| THE BEGINNING | A <i>data model</i> is an abstract representation of the data in an enterprise or of the information that is derived from that data. Data and information can (and should) be represented at multiple levels of abstraction, each providing a different perspective and understanding of the data. The highest level of abstraction is a business context view with both external (outside looking in) and internal (looking from within) perspectives. Thus data modeling begins with business activities and the information needs of those activities. This view describes the scope of and the context for business information requirements—a sensible start to <i>modeling the right data</i> . |
|---------------|---|
| THE END       | The natural conclusion of data modeling is implemented data—data files<br>and database tables. Far from the top level of abstraction, implemented<br>data is beyond the bottom tier; it is no longer abstract but real and<br>physical. At this level attention turns to <i>the right implementation for the</i><br><i>data</i> .   |
| THE MIDDLE    | The complexities and challenges of data modeling lie between the top<br>layer of context and the bottom tier of implementation. Getting from the<br>right data to the right implementation involves an understanding of the<br>business and many information systems, ranging from those that capture<br>data in the course of business activities to those that turn data into<br>information and supply that information to the business.   |
| IN REVERSE    | In business analytics, it is often necessary to explore data sets that are not<br>fully understood. This search for value reverses the process described<br>above. Analysis starts with a physical data structure that is the subject of<br>exploration; the process ends with an understanding of its content and<br>applicability in the enterprise.  |

© TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY. 1-3

# The Data Modeling Life Cycle Between Business Needs and Implemented Data



# The Data Modeling Life Cycle Between Business Needs and Implemented Data

| BUSINESS AND<br>SYSTEMS VIEWS | Business context alone is insufficient to describe business requirements<br>for data and information. Similarly, implementation is not adequate to<br>design and deploy information systems and databases. Multiple levels of<br>abstraction are needed to manage complexity, understand and document<br>from multiple perspectives, communicate effectively, and provide a<br>natural progression from need to solution. The six layers of modeling<br>abstraction are based on the Zachman Framework ( <i>A Framework for</i><br><i>Information Systems Architecture, IBM Los Angeles Scientific Center</i><br><i>Report</i> , John A. Zachman, 1986). Zachman's framework addresses many<br>aspects of information systems architecture. This course is concerned<br>solely with the data. To learn more about the Zachman Framework, visit<br>www.zachman.com. The six levels of data modeling are: |
|-------------------------------|---|
| CONTEXTUAL                    | Context modeling provides a view of the scope of the planned data<br>warehousing program. Context models communicate understanding of<br>the business requirements and establish a context for analysis. This level<br>corresponds to the <i>Executive Perspective</i> of the Zachman Framework.  |
| CONCEPTUAL                    | Models at this level are about understanding the required set of data stores. Conceptual models describe data requirements from a business point of view without the burden of technical details. This corresponds to the <i>Business Management Perspective</i> of the framework.  |
| LOGICAL                       | Models at this level refine conceptual models by documenting entities, their attributes, and their relationships. These models are technology-<br>oriented designs, although they are platform-independent. This level corresponds to the <i>Architect Perspective</i> of Zachman's framework.  |
| STRUCTURAL                    | Structural models move a step closer to implementation. They specify<br>the design necessary for the warehouse/marts to maintain history, to<br>distribute data, and to provide for ease of use. Structural modeling<br>corresponds to the <i>Engineer Perspective</i> of the Zachman framework.  |
| PHYSICAL                      | Physical models represent the detailed specification of what is physically implemented using specific technology. This level corresponds to the <i>Technician Perspective</i> of Zachman's framework.   |
| FUNCTIONAL                    | This level does not describe models, but implemented and instantiated data stores. This is the tangible result of modeling activity and the product of the development process. This level corresponds to the <i>Enterprise Perspective</i> of the Zachman framework.   |

© TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY. 1-5



# Module 2

#### Business Data Models

| Торіс                             | Page |
|-----------------------------------|------|
| Business Context                  | 2-2  |
| Business Data Model Development   | 2-14 |
| Gathering Business Questions      | 2-18 |
| Analyzing Business Questions      | 2-26 |
| Qualifier Analysis and Refinement | 2-38 |
| Fact-Qualifier Analysis Results   | 2-40 |
| Business Dimensional Modeling     | 2-42 |

#### Business Context Business Drivers, Goals, and Strategies



# **Business Context**

Business Drivers, Goals, and Strategies



- WHY MODEL Business context determines the nature of data and information services: the business processes to be affected, the kinds of applications to be implemented, and the information services to provide. Business context provides the means to align data with business goals. The contextual representation is generally determined prior to initiating a project. Project team members should be aware of and use the context to establish a more complete and relevant set of requirements.
- **BUSINESS DRIVERS Business drivers** are those things that are strategically important in positioning the business to achieve its short- and long-term goals. They are the external forces that have significant influence on the operation and performance of a business. Drivers create the need to take action, but they don't dictate the actions to be taken. Common business drivers include changing economic, political, social, and technological factors.
- **BUSINESS GOALS** *Business goals* are the things that the business wants to accomplish to respond to business drivers. Drivers create the need to act. Goals describe the desired outcomes of taking action. Goals are commonly related to financial or operational performance (i.e., cost reduction, generation of revenue, increased market share, etc.). Goals are most effective in setting data management priorities and directions when they are: (1) described by clear, concise, understandable statements, (2) specific enough that the level of achievement can be measured, and (3) of high business priority.
- **BUSINESS STRATEGIES Business strategies** are plans to turn the defined business goals into reality. They include the things that the business will do to shape its future and achieve the business goals. The range of strategies is broad introducing new products, exploiting new sales channels, pricing competitively, optimizing business processes, etc. Strategies help to determine which business processes and organizations most need to be information-enabled.

# Business Context

#### Business Drivers, Goals, and Strategies



- Provide incentives to professors pursuing funded research projects.
- Choose new degree program to introduce.

# Business Context

#### Business Drivers, Goals, and Strategies

BUSINESS<br/>DRIVERSThroughout this course, examples will be presented based on a<br/>hypothetical university that provides post-secondary education to students<br/>and has a large faculty which engages in research and publication.

Sample business drivers for the Data University could include:

- Tuition costs are having an increasing impact on university selection
- Donors are demanding innovative programs
- Many universities are offering online programs

**BUSINESS GOALS** Sample business goals for the Data University could include:

- Become recognized as a leading university in our fields
- Create and publicize innovative programs
- Increase donor involvement

#### BUSINESS STRATEGIES

Sample business strategies for the Data University could include:

- Engage field luminaries as advisors or professors
- Create donor focus groups
- Implement alumni relationship organization
- Expand unique online offerings
- Increase the proportion of students enrolled in online programs by 20 percent



- What are some of the business drivers for your company?
- What goals has your company created as a result?
- What are some of your business strategies to achieve those goals?



# Module 3

# Logical Data Models

| Торіс                                     | Page |
|---|------|
| What to Model                             | 3-2  |
| Understanding Data Sources                | 3-4  |
| Logical E-R Modeling                      | 3-10 |
| Logical Dimensional Modeling              | 3-20 |
| Logical Models and Business Metrics       | 3-34 |
| Logical Models and Business Analytics     | 3-40 |
| Logical Models and Master Data Management | 3-44 |
| Logical Models and Nonrelational Data     | 3-48 |

# What to Model The Data and Information Pipeline



3-2 © TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY.

# What to Model

#### The Data and Information Pipeline

#### THE MODELING FRAMEWORK



# FLOW OF DATA In a typical data warehousing environment data flows from sources, through a sequence of intermediate processes and data stores, to be delivered ultimately as information that is useful to business. The data modeler's challenge—with data stores including sources, staging, data warehouses, and data marts—is to determine which data to model and at what level to model the data.

Every project must determine the right amount of data modeling to be done. Compromises are often necessary to achieve appropriate balance between data complexity and the real project constraints of time and resources. A practical minimum standard includes:

- User view modeling that is sufficient to understand end-user requirements for data.
- Modeling of managed data that is detailed enough to meet the data delivery requirements for user products, including data flows, cleansing, and integration.
- Source data modeling as needed to understand the content and structure of data sources.

## Understanding Data Sources Why Sources Matter



# Understanding Data Sources

#### Why Sources Matter

#### CHOOSING THE RIGHT DATA

Choosing the right data sources is a critical step to getting the right data in the warehouse and delivering the right information to the business. Making the right choices depends on real understanding of the contents of each data source, followed by careful qualification of sources based on the criteria that matter to your project and your data warehouse. Some of the common qualification criteria are listed below.

| Qualifying<br>Criteria | Assessment Questions  |
|------------------------|---|
| Availability           | How available and accessible is the data? Are there technical obstacles to access? Or ownership and access authority issues?                            |
| Understandability      | How easily understood is the data? Is it well documented? Does someone in the organization have depth of knowledge? Who works regularly with this data? |
| Stability              | How frequently do data structures change? What is the history of change for the data? What is the expected life span of the potential data source?      |
| Accuracy               | How reliable is the data? Do the business people who work with the data trust it?   |
| Timeliness             | When and how often is the data updated? How current is the data? How much history is available? How available is it for extraction?                     |
| Completeness           | Does the scope of data correspond to the scope of the data warehouse? Is any data missing?  |
| Granularity            | Is the source the lowest available grain (most detailed level) for this data?   |

#### BEYOND ENTERPRISE DATA

The advent of big data analytics has pushed the boundary for what is considered an acceptable data source. Although past solutions were limited to enterprise data extracted from relational databases, some analytics applications may leverage data that originates outside the enterprise and information that is managed in a nonrelational format.

In this environment, the concept of data quality takes on increased nuance. Rather than asking, "is the data high quality," it becomes necessary to ask, "is the quality of this data suited to this business need." An enterprise dashboard likely requires high-quality data, while less reliable data may be acceptable for an analytic model that forecasts market demand.

© TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY. 3-5



# Module 4

# Physical Data Models

| Торіс                                      | Page |
|--|------|
| Defining Physical Modeling                 | 4-2  |
| Data Structure in Transaction Systems      | 4-4  |
| Structural Modeling and Data Integration   | 4-6  |
| Structural Modeling and Business Analytics | 4-28 |
| Physical Design Overview                   | 4-46 |
| Some Optimization Techniques               | 4-48 |
| Physical Design and Implementation         | 4-60 |

© TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY. 4-1

# Defining Physical Modeling Both Structural and Physical



# Defining Physical Modeling

Both Structural and Physical

| PHYSICAL AND<br>STRUCTURAL<br>MODELS | The enterprise architecture framework introduced in Module 1 divided<br>modeling activities across six levels of abstraction. These layers ranged<br>from a contextual representation of the business situation to the<br>implementation of actual data.   |  |
|--------------------------------------|--|--|
|                                      | Once data modelers have completed logical modeling activities, they<br>often move to a stage called "physical modeling" that actually combines<br>the physical and structural layers of this framework.  |  |
|                                      | This simplification is both acceptable and commonplace. In fact, most<br>modeling tools combine these categories into a single modeling activity.<br>In this module, physical and structural modeling will sometimes be<br>collapsed under the "physical" moniker, in acknowledgement of industry<br>practices.  |  |
|                                      | It is useful to note that many teams divide responsibilities at this stage,<br>even though the result is a single model. For example, a data modeler<br>may define the structural features of the solution, such as the tables,<br>columns, and keys, while a database administrator takes care of the<br>detailed specifications, such as data types, indexing, and so forth. |  |
| IMPLEMENTATION                       | The bottom-most layer of the enterprise architecture framework is the implementation layer. This component represents the data itself. An additional model is not produced for this layer.   |  |

© TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY. 4-3

# Data Structure in Transaction Systems Extracting the Structure of Existing Data

| Table   | Column                 | Description                 | ID  | Entity  |
|---------|------------------------|-----------------------------|-----|---------|
| Grades  | Student ID             | Unique identifier           | Yes |         |
|         | Student Name           | Name of student             |     | Student |
|         | Student Number         | Assigned student number     |     | Student |
|         | Course Number          | Catalog number of course    |     | Course  |
|         | Grade                  | Letter grade for the course |     |         |
|         | Credit Hours           | Credit hours for the course |     | Course  |
|         | Program Name           | Program student is enrolled |     | Program |
| Student | Student ID             | Unique identifier           | Yes |         |
|         | Student Name           | Name of student             |     |         |
|         | Student Birth Date     | Birth date of student       |     |         |
|         | Student Marital Status | Marital status of student   |     |         |

![](_page_24_Figure_4.jpeg)

4-4 © TDWI. All rights reserved. Reproductions in whole or in part are prohibited except by written permission. DO NOT COPY.

#### Data Structure in Transaction Systems Extracting the Structure of Existing Data

#### LEVELS OF TRANSACTION SYSTEM MODELS

Four levels of transaction systems data and models are of interest for the value that they contribute to data warehouse modeling:

- *Implemented Data* already exists as application data files and tables.
- *Physical models* usually exist in part as descriptions of data structures (DDL, COBOL copy code, etc.). Multiple potentially confusing descriptions may exist for a single file or table. This is especially true of older legacy systems.
- *Structural Models* may need to be created to identify and resolve conflict where multiple physical descriptions exist for a file or table.
- *Logical Models* exist infrequently. Where logical models are present, they are often incomplete or out-of-date, and are almost certainly limited to a single application's view of the data.

#### BOTTOM-UP DATA SOURCE MODELING

It is valuable to build a source data model in reverse—working from implemented data backward to an E-R model. The process is one of examining every data element (column or field) in every data store (file or table) to answer the questions:

- What business fact (attribute) does this data element contain?
- What thing (entity) is it a fact about?
- Does the data element identify the entity that it describes?
- Does the data element indicate a relationship to another entity?

The answers can be represented in a table such as that on the opposite page, which is readily translated into an E-R diagram if necessary.

This bottom-up understanding of source data is referred to as *data profiling*, and it can be supported by automated tools that analyze data and identify patterns, domains of values, apparent key relationships, and more. A data profiling process reveals the truth about the data, which often contradicts existing artifacts such as schema designs or E-R documentation which has not been kept up-to-date.

# This page intentionally left blank.