

# Information Technology Policy

## Data Modeling Basics

<b>ITP Number</b> STD-INF003B	<b>Effective Date</b> August 2, 2005
<b>Category Information</b>	<b>Supersedes</b> None
<b>Contact</b> <a href="mailto:RA-ITCentral@pa.gov">RA-ITCentral@pa.gov</a>	<b>Scheduled Review</b> May 2022

## Data Modeling Basics

Office of Data & Digital Technology

### 1. Definitions

**Aggregation:** a technique used to consolidate/collect and present summarized KPIs data; one that optimizes data retrieval by summarizing rows of a fact table according to a specific dimensions/attributes.

**Business Rule:** stipulates specific business-related definition and information that is linked to database objects. The information is always associated with the business facts or descriptions; or it might be formulas or algorithms, either client-based or destined for the server. Once defined, Business Rules can be applied common across the enterprise across various enterprise and analytical applications.

**Cardinality:** indicates the relationship for linked or associated entities (one or many) of an Entity in relation to another Entity. You can select the following values for Cardinality:

- One-to-one - One instance of the first Entity can correspond to only one instance of the second Entity.
- One-to-many - One instance of the first Entity can correspond to more than one instance of the second Entity.
- Many-to-one - More than one instance of the first Entity can correspond to the same one instance of the second Entity.
- Many-to-many - More than one instance of the first Entity can correspond to more than one instance of the second Entity.

**Data Attribute:** A term used in Logical Data Models to describe a kind of fact common to all or most instances of an Entity. Student ID is an attribute of the Entity Student. The corresponding Physical Data Model generally implements the attribute as a database column or field.

**Data Element:** An Entity, attribute, database table, or database column used to represent business information in Logical or Physical Data Models. Data element primarily defines the metadata and represents data atomicity. Users should be aware that the literature also defines Data Element to explicitly mean an attribute of an Entity. However, as defined in this document, the term encompasses both Entities and attributes in Logical Data Models as well as tables and columns in Physical Data Models.

**Data Entity:** A term used in Logical Data Models to describe a business entity eg. class of persons, places, things, concepts, or events of interest to the business, about which the business intends to keep facts. The corresponding Physical Data Model generally implements the Entity in a database table or view.

**Dimension:** Dimensions refers to the business objects signifying non values fields/attributes used to attribute the axis of investigation of a fact. The dimensions define the significance of the KPIs/Facts.

**Domain:** A way of identifying and grouping the types of data items in the model. This makes it easier to standardize data characteristics for attributes/columns in different Entities/tables. Some database management systems (DBMSs) will implement Domains as "User Defined Datatypes". Another feature of Domains is in the maintenance of similar columns. If all "name" columns (LastName, CityName, ProductName, etc.) are defined as a common Domain, then changing the datatype from char(40) to char(50) is a one-step procedure, rather than having to visit each table and search for the correct columns.

**Enterprise Class Database Management System:** integrates multiple business processes or applications into a single DBMS and hardware platform. This contrasts with creating application specific database management systems.

**Entity:** A business object like Person, place, thing, or concept that has characteristics of interest to the enterprise and about which you want to store information.

**Inheritance:** Inheritance allows you to define an Entity as a special case of a more general entity. The Entities involved in an Inheritance have many similar characteristics but are nonetheless different. The general entity is known as a supertype (or parent) entity and contains all the common characteristics. The special case entity is known as a subtype (or child) entity and contains all the particular characteristics.

**Logical Data Model:** Data model indicating the relationships of entities, representing a structured representation of the data of importance to the business, in terms of Entities, attributes, and their Relationships including the Business Rules that govern them. The representation includes both graphical depictions and textual definitions. Logical Data Models are used to translate business requirements into data representations that are understandable to information systems professionals.

**Logical Data Name:** A unique identifier of an Entity or attribute as stored within a Logical Data Model or data dictionary. Logical Data Names should consist of English words and must be understandable by the end user. Also known as the Business Name or Functional Name.

**Physical Data Model:** A structured representation of the data of importance to the business, in terms of database tables and columns along with their Relationships, formats, and Business Rules that govern the data. The representation includes both graphical depictions and textual definitions. Physical Data Models are used exclusively by information systems professionals to deploy database systems using appropriate database software.

**Physical Data Name:** A unique identifier of an Entity or attribute as implemented within one or more database systems. Physical Data Names are generally constrained by the limitations of the database software.

**Referential Integrity:** Referential Integrity refers to rules governing data consistency, specifically the interaction between primary keys and foreign keys in different tables. Referential Integrity dictates what happens when you update or delete a value in a referenced column in the parent table and when you delete a row containing a referenced column from the parent table.

**Relationship:** A Relationship is a named connection or association between Entities. Each Relationship is drawn as a line connecting the two Entity types; each Relationship is given a name that indicates what information it imparts (Relationships are named in both directions); the *type* of Relationship (*Cardinality* and *optionality*) is specified as follows: the line style (dash or solid) indicates optionality, and the Relationship ends indicate Cardinality.

**Source:** Materials referenced and used from datamodel.org and Ambysoft's data modeling website.

## Why is Data Modeling Important?

Data modeling allows to define the common business entities, their relationships and use them to define

a standardized common enterprise model. The foundations of correct modeling practices will benefit long term in establishing a

Data modeling is a very vital part in the development process. One can compare this to creating a blueprint to build a house before the actual building takes place. As much as the blueprint takes time to prepare, it can also go through multiple iterations of validation to ensure the foundation, structure, and aesthetics of the building plan conform to intended objectives and quality standards. Therefore, data modeling is an intensive process which consumes a major part of the development time.

The model is built in a phased manner and goes through several iterations of validation to ensure that the structure and content of the model addresses the business objectives of the enterprise or application, meets quality standards, is modular, provides a solid foundation for future extensions and data reuse for other enterprise applications, etc. That being the case, less time spent in this effort will only produce a weak unstructured model that will be very expensive to maintain in the future, may produce inconsistent and incorrect results and will be unfit for reporting or future extensions.

Major events in data modeling include:

- Identifying business requirements, domains and business entities
- Identifying the entity relationships and cardinality
- Data definitions Segregation – dimensions / attributes, individual entities, KPIs/ Measures and derived entities such as formulas etc
- Identifying Entities, data requirements and processes.
- Defining attributes of the data such as data types, sizes, defaults.
- Applying validation and Business Rules to ensure data integrity.
- Defining data management and security processes.
- Specifying data archival and storage.

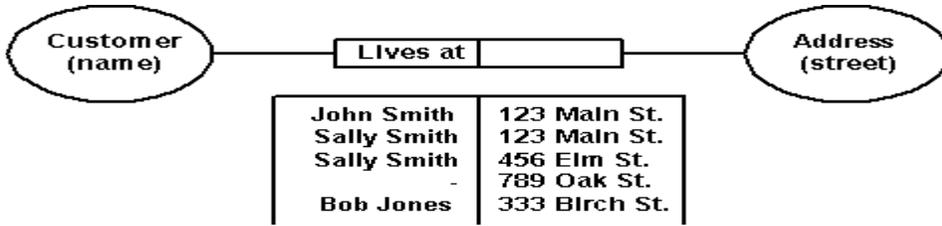
## How are Data Models Used in Practice?

Generally, there are three basic types of data models:

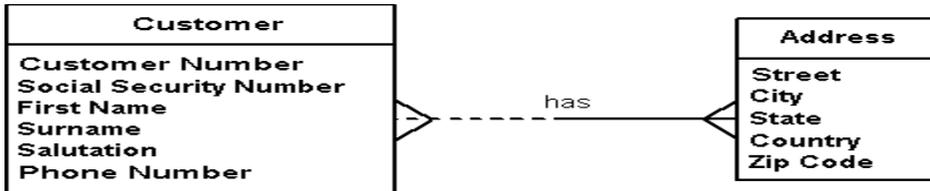
- **Conceptual data models.** These models, sometimes called domain models, are typically used to identify and document business (Domain) concepts with project stakeholders. Conceptual data models are often created as the precursor to Logical Data Models (LDMs) or as alternatives to LDMs.
- **Logical Data Models (LDMs).** Logical Data Models are used to further explore the Domain concepts, and their Relationships and Relationship Cardinalities. This could be done for the scope of a single project or for your entire enterprise. Logical Data Models depict the logical entity types, typically referred to simply as entity types, the Data Attributes describing those Entities, and the Relationships between the Entities. DDL can be generated at this level.
- **Physical Data Models (PDMs).** Physical Data Models are used to design the internal schema of a database, depicting the data tables (derived from the logical data entities), the data columns of those tables (derived from the Entity attributes), and the Relationships between the tables derived from the Entity Relationships).

The level of detail that is modeled is significantly different for each model type. This is because the goals and audience for each diagram are different. You can use a Logical Data Model to explore Domain concepts with your stakeholders and the Physical Data Model to define your database design. Each of the various models should also reflect your organization's naming standards. A Physical Data Model should also indicate the data types for the columns, such as integer or character.

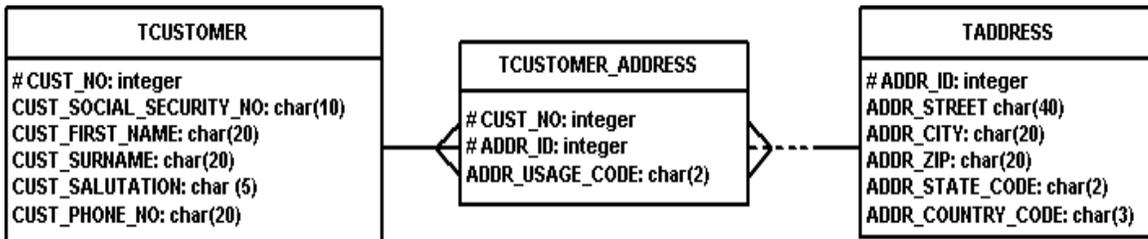
## A simple conceptual data model:



## A simple Logical Data Model:



## A simple Physical Data Model:



Data models can be used effectively at both the enterprise level and on individual projects. Enterprise architects will often create one or more high-level Logical Data Models that depict the data structures that support the enterprise. These models are typically referred to as enterprise data models or enterprise information models. An enterprise data model is one of several critical views that the organization's enterprise architects will maintain and support; other views may explore network/hardware infrastructure, organization structure, software infrastructure, and business processes (to name a few). Enterprise data models provide information that a project team can use -- both as a set of constraints and as important insights into the structure of their system.

Project teams will typically create Logical Data Models as a primary analysis artifact when their implementation environment is predominantly procedural in nature. Logical Data Models are also a good choice when a project is data-oriented in nature -- perhaps a data warehouse or reporting system is being developed. However, Logical Data Models are often a poor choice when a project team is using object-oriented or component-based technologies where the developers typically prefer UML diagrams or when the project is not data-oriented in nature.

When a relational database is used for data storage, project teams are best advised to create a Physical Data Model to model its internal schema. A Physical Data Model is often one of the critical design artifacts for business application development projects.

## Data Modeling vs. Class Modeling:

Data modeling is different from class modeling because it focuses solely on data.

It is important to perform data modeling and develop the ERD (Entity Relationship Diagram) to ensure the relational database is properly designed.

From the point of view of an object-oriented developer, class modeling is a useful approach. Class

models allow you to explore both the behavior and data aspects of your Domain. Similar to data entities, there are associations between classes; Relationships, Inheritance, composition, and Aggregation are all applicable concepts in modeling.

It is very important to get project stakeholders actively involved in the creation of the model. Instead of a traditional analyst-led drawing session, you can instead engage stakeholders by facilitating the creation of models.

### Data Perspectives:

Data Models are a valuable source of information, providing a graphical depiction of data at different levels of abstraction. For example, the owner of a business process is interested in the conceptual view of data – the conceptual model. The designer of the data is interested in the logical view – the logical model. This view is sometimes referred to as the transformation layer. The data administrator is typically concerned with this model. The database administrator is typically more concerned with the physical model and the physical implementation of a relational database.

This table summarizes the different perspectives.

Model	Perspective	Model Description	Type of Model	Entity	Type of Relation
Business Model	Owner	Semantic Model	Conceptual	Business Entity	Business
System Model	Designer	Logical Data Model	Logical	Data Entity	Data
Technology Model	Builder	Data Design	Physical	Table/Segments	Key/Pointer
Detailed Representation	Developer	Data Definitions		Field	Address

### Data Modeling Standards Supported:

There are three common data modeling notations: Information Engineering (IE), IDEF1X, and the Unified Modeling Language (UML).

From a notation standard, current product standards should support both IDEF1X and IE modeling standards as well as naming standards that allow you to create glossaries of approved words and enforce the way words are used in naming tables and columns, etc.

XML (which has a very loose industry standard) should also be supported in the tool. It is likely that UML based tools will co-exist with other data modeling notations for the near future. UML provides all of the syntax needed to perform data modeling, as well as behavioral modeling. Some developers see an advantage to using one modeling language for all modeling purposes.

### Common data modeling notations:

Notation	Comments
IE	The IE notation (Finkelstein 1989) is simple, easy to read, and well suited for high-level logical and enterprise data modeling.
IDEF1X	This notation is more complex. It was originally intended for physical modeling but has been applied for logical modeling as well.

UML	The Object Management Group (OMG) adopted UML as a standard in 1997.
-----	--

The current standards either support the use of IE and IDEF1X notations or UML. (Reference [Data Modeling Products and Standards](#) for the latest version.)

**Model Reuse:**

A Logical Data Model facilitates data re-use and sharing. Data is stable over time; therefore, the model remains stable over time. As additional project teams scope out their areas, they can re-use model components that are shared by the business. This leads to physical data sharing and less storage of redundant data. It also helps the organization recognize that information is an organization-wide resource, not the property of one department or another. Data sharing makes the organization more cohesive and increases the quality of service.

**Model Review Process:**

A review process with appropriate team members and stakeholders should follow each stage of model development. This process will ensure that the respective data model correctly and accurately represents the business requirements and maximizes data integrity. NOTE: A separate effort is currently underway to develop a formal data modeling methodology, including templates and checklists to validate Logical and Physical Data Models.

**Data Modeling Best Practice Standards Supported:**

A list of Data Modeling Best Practices has been compiled by the Office of Data & Digital Technology. These standards have applicability across all current standard products and are required to be used for all application development efforts of sufficient size and scope. If a specific standard applies only to mission- critical applications, it will be identified as such. Reference BPD-INF003C [Data Modeling Best Practices](#) for the latest version.

<b>Data Modeling Basic Steps</b>
<p>Pre-requisite:</p> <ul style="list-style-type: none"> <li>- Business Vision and Business process modeling</li> <li>- Define business entities and conceptual model</li> </ul> <p><b>1. Identify Entity Types</b> - an entity type represents a collection of similar objects. An Entity could represent a collection of people, places, things, events, or concepts. Examples of Entities in an order entry system would include <i>Customer, Address, Order, Item, and Tax</i>. If you were class modeling you would expect to discover classes with the exact same names. However, the difference between a class and an entity type is that classes have both data and behavior whereas entity types just have data. Ideally, an Entity should be "normal" -- the data modeling world's version of cohesive. A normal Entity depicts one concept, just like a cohesive class models one concept. For example, customer and order are clearly two different concepts. Therefore, it makes sense to model them as separate Entities.</p> <p><b>2. Identify Attributes</b> - each entity type will have one or more Data Attributes. For example, the <i>Customer</i> Entity has attributes such as <i>First Name</i> and <i>Surname</i> and the <i>TCUSTOMER</i> table had corresponding data columns <i>CUST_FIRST_NAME</i> and <i>CUST_SURNAME</i> (a column is the implementation of a Data Attribute within a relational database). Attributes should also be cohesive from the point of view of your Domain, something that is often a judgment call. If you wanted to model the fact that people had both first and last names instead of just a name (e.g. "John" and</p>

“Doe” vs. “John Doe”) whereas we did not distinguish between the sections of a zip code (e.g. 90210-1234-5678). Getting the level of detail right can have a significant impact on your development and maintenance efforts.

**3. Establish Data Naming Conventions** - Standards and guidelines applicable to data modeling should be set and enforced. Commonly, this would be the responsibility of a data administrator. These guidelines should include naming conventions for both logical and physical modeling. The logical naming conventions should be focused on human readability whereas the physical naming conventions will reflect technical considerations. The basic idea is that developers should agree to and follow a common set of modeling standards on a software project. Just like there is value in following common coding conventions, clean code that follows your chosen coding guidelines is easier to understand and evolve than code that doesn't, there is similar value in following common modeling conventions.

**4. Identify Relationships** - Entities have Relationships with other Entities. For example, customers PLACE orders, customers LIVE AT addresses, and line items ARE PART OF orders. Place, live at, and are part of are all terms that define Relationships between Entities. The Relationships between Entities are conceptually identical to the Relationships (associations) between objects.

**5. Assign Keys** - A key is one or more Data Attributes that uniquely identify an Entity. A key that consists of two or more attributes is called a *composite key*. A key that is formed of attributes that already exist in the real world is called a *natural key*. An Entity type in a Logical Data Model will have zero or more *candidate keys*, also referred to simply as *unique identifiers*. Both of these keys are called candidate keys because they are candidates to be chosen as the *primary key*, an *alternate key* (also known as a *secondary key*), or perhaps not even a key at all within a Physical Data Model. A primary key is the preferred key for an entity type, whereas an alternate key (also known as a secondary key) is an alternative way to access rows within a table. In a physical database, a key would be formed of one or more table columns whose value(s) uniquely identify a row within a relational table.

**6. Normalize Data** - Normalization is a process in which Data Attributes within a data model are organized to increase the cohesion of entity types. In other words, the goal of data normalization is to reduce, and even eliminate, data redundancy.

**7. Optimize Performance** - Normalized data schemas, when put into production, may suffer from performance problems. This makes sense – the rules of data normalization focus on reducing data redundancy, not on improving performance of data access. It may be necessary to denormalize portions of your data schema to improve database access efficiency. It should be documented why changes were made to the model.

This chart contains a history of this publication's revisions:

Version	Date	Purpose of Revision
Original	08/2/2005	Base Document
Revision	11/18/2010	ITP Refresh
Revision	05/13/2021	ITP Refresh