
Introduction to Data Warehousing Concepts

This chapter provides an overview of the Oracle data warehousing implementation. It contains:

- [What Is a Data Warehouse?](#)
- [Contrasting OLTP and Data Warehousing Environments](#)
- [Common Data Warehouse Tasks](#)
- [Data Warehouse Architectures](#)

What Is a Data Warehouse?

A **data warehouse** is a database designed to enable business intelligence activities: it exists to help users understand and enhance their organization's performance. It is designed for query and analysis rather than for transaction processing, and usually contains historical data derived from transaction data, but can include data from other sources. Data warehouses separate analysis workload from transaction workload and enable an organization to consolidate data from several sources. This helps in:

- Maintaining historical records
- Analyzing the data to gain a better understanding of the business and to improve the business

In addition to a relational database, a data warehouse environment can include an extraction, transportation, transformation, and loading (ETL) solution, statistical analysis, reporting, data mining capabilities, client analysis tools, and other applications that manage the process of gathering data, transforming it into useful, actionable information, and delivering it to business users.

To achieve the goal of enhanced business intelligence, the data warehouse works with data collected from multiple sources. The source data may come from internally developed systems, purchased applications, third-party data syndicators and other sources. It may involve transactions, production, marketing, human resources and more. In today's world of big data, the data may be many billions of individual clicks on web sites or the massive data streams from sensors built into complex machinery.

Data warehouses are distinct from online transaction processing (OLTP) systems. With a data warehouse you separate analysis workload from transaction workload. Thus data warehouses are very much read-oriented systems. They have a far higher amount of data reading versus writing and updating. This enables far better analytical performance and avoids impacting your transaction systems. A data warehouse system can be optimized to consolidate data from many sources to achieve a key goal: it becomes your organization's "single source of truth". There is great value in having a

consistent source of data that all users can look to; it prevents many disputes and enhances decision-making efficiency.

A data warehouse usually stores many months or years of data to support historical analysis. The data in a data warehouse is typically loaded through an extraction, transformation, and loading (ETL) process from multiple data sources. Modern data warehouses are moving toward an extract, load, transformation (ELT) architecture in which all or most data transformation is performed on the database that hosts the data warehouse. It is important to note that defining the ETL process is a very large part of the design effort of a data warehouse. Similarly, the speed and reliability of ETL operations are the foundation of the data warehouse once it is up and running.

Users of the data warehouse perform data analyses that are often time-related. Examples include consolidation of last year's sales figures, inventory analysis, and profit by product and by customer. But time-focused or not, users want to "slice and dice" their data however they see fit and a well-designed data warehouse will be flexible enough to meet those demands. Users will sometimes need highly aggregated data, and other times they will need to drill down to details. More sophisticated analyses include trend analyses and data mining, which use existing data to forecast trends or predict futures. The data warehouse acts as the underlying engine used by middleware business intelligence environments that serve reports, dashboards and other interfaces to end users.

Although the discussion above has focused on the term "data warehouse", there are two other important terms that need to be mentioned. These are the data mart and the operation data store (ODS).

A data mart serves the same role as a data warehouse, but it is intentionally limited in scope. It may serve one particular department or line of business. The advantage of a data mart versus a data warehouse is that it can be created much faster due to its limited coverage. However, data marts also create problems with inconsistency. It takes tight discipline to keep data and calculation definitions consistent across data marts. This problem has been widely recognized, so data marts exist in two styles. Independent data marts are those which are fed directly from source data. They can turn into islands of inconsistent information. Dependent data marts are fed from an existing data warehouse. Dependent data marts can avoid the problems of inconsistency, but they require that an enterprise-level data warehouse already exist.

Operational data stores exist to support daily operations. The ODS data is cleaned and validated, but it is not historically deep: it may be just the data for the current day. Rather than support the historically rich queries that a data warehouse can handle, the ODS gives data warehouses a place to get access to the most current data, which has not yet been loaded into the data warehouse. The ODS may also be used as a source to load the data warehouse. As data warehousing loading techniques have become more advanced, data warehouses may have less need for ODS as a source for loading data. Instead, constant trickle-feed systems can load the data warehouse in near real time.

A common way of introducing data warehousing is to refer to the characteristics of a data warehouse as set forth by William Inmon:

- [Subject Oriented](#)
- [Integrated](#)
- [Nonvolatile](#)
- [Time Variant](#)

Subject Oriented

Data warehouses are designed to help you analyze data. For example, to learn more about your company's sales data, you can build a data warehouse that concentrates on sales. Using this data warehouse, you can answer questions such as "Who was our best customer for this item last year?" or "Who is likely to be our best customer next year?" This ability to define a data warehouse by subject matter, sales in this case, makes the data warehouse subject oriented.

Integrated

Integration is closely related to subject orientation. Data warehouses must put data from disparate sources into a consistent format. They must resolve such problems as naming conflicts and inconsistencies among units of measure. When they achieve this, they are said to be integrated.

Nonvolatile

Nonvolatile means that, once entered into the data warehouse, data should not change. This is logical because the purpose of a data warehouse is to enable you to analyze what has occurred.

Time Variant

A data warehouse's focus on change over time is what is meant by the term time variant. In order to discover trends and identify hidden patterns and relationships in business, analysts need large amounts of data. This is very much in contrast to **online transaction processing (OLTP)** systems, where performance requirements demand that historical data be moved to an archive.

Key Characteristics of a Data Warehouse

The key characteristics of a data warehouse are as follows:

- Data is structured for simplicity of access and high-speed query performance.
- End users are time-sensitive and desire speed-of-thought response times.
- Large amounts of historical data are used.
- Queries often retrieve large amounts of data, perhaps many thousands of rows.
- Both predefined and ad hoc queries are common.
- The data load involves multiple sources and transformations.

In general, fast query performance with high data throughput is the key to a successful data warehouse.

Contrasting OLTP and Data Warehousing Environments

There are important differences between an OLTP system and a data warehouse. One major difference between the types of system is that data warehouses are not exclusively in **third normal form (3NF)**, a type of data normalization common in OLTP environments.

Data warehouses and OLTP systems have very different requirements. Here are some examples of differences between typical data warehouses and OLTP systems:

- Workload

Data warehouses are designed to accommodate *ad hoc* queries and data analysis. You might not know the workload of your data warehouse in advance, so a data warehouse should be optimized to perform well for a wide variety of possible query and analytical operations.

OLTP systems support only predefined operations. Your applications might be specifically tuned or designed to support only these operations.

- Data modifications

A data warehouse is updated on a regular basis by the ETL process (run nightly or weekly) using bulk data modification techniques. The end users of a data warehouse do not directly update the data warehouse except when using analytical tools, such as data mining, to make predictions with associated probabilities, assign customers to market segments, and develop customer profiles.

In OLTP systems, end users routinely issue individual data modification statements to the database. The OLTP database is always up to date, and reflects the current state of each business transaction.

- Schema design

Data warehouses often use partially denormalized schemas to optimize query and analytical performance.

OLTP systems often use fully normalized schemas to optimize update/insert/delete performance, and to guarantee data consistency.

- Typical operations

A typical data warehouse query scans thousands or millions of rows. For example, "Find the total sales for all customers last month."

A typical OLTP operation accesses only a handful of records. For example, "Retrieve the current order for this customer."

- Historical data

Data warehouses usually store many months or years of data. This is to support historical analysis and reporting.

OLTP systems usually store data from only a few weeks or months. The OLTP system stores only historical data as needed to successfully meet the requirements of the current transaction.

Common Data Warehouse Tasks

As an Oracle data warehousing administrator or designer, you can expect to be involved in the following tasks:

- Configuring an Oracle database for use as a data warehouse
- Designing data warehouses
- Performing upgrades of the database and data warehousing software to new releases
- Managing schema objects, such as tables, indexes, and materialized views
- Managing users and security
- Developing routines used for the extraction, transformation, and loading (ETL) processes

- Creating reports based on the data in the data warehouse
- Backing up the data warehouse and performing recovery when necessary
- Monitoring the data warehouse's performance and taking preventive or corrective action as required

In a small-to-midsized data warehouse environment, you might be the sole person performing these tasks. In large, enterprise environments, the job is often divided among several DBAs and designers, each with their own specialty, such as database security or database tuning.

These tasks are illustrated in the following:

- For more information regarding partitioning, see *Oracle Database VLDB and Partitioning Guide*.
- For more information regarding database security, see *Oracle Database Security Guide*.
- For more information regarding database performance, see *Oracle Database Performance Tuning Guide* and *Oracle Database SQL Tuning Guide*.
- For more information regarding backup and recovery, see *Oracle Database Backup and Recovery User's Guide*.
- For more information regarding ODI, see *Oracle Fusion Middleware Developer's Guide for Oracle Data Integrator*.

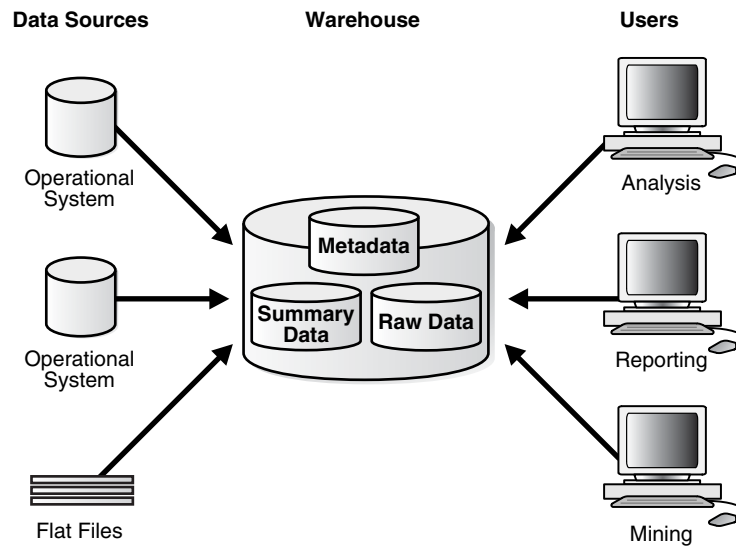
Data Warehouse Architectures

Data warehouses and their architectures vary depending upon the specifics of an organization's situation. Three common architectures are:

- [Data Warehouse Architecture: Basic](#)
- [Data Warehouse Architecture: with a Staging Area](#)
- [Data Warehouse Architecture: with a Staging Area and Data Marts](#)

Data Warehouse Architecture: Basic

[Figure 1-1](#) shows a simple architecture for a data warehouse. End users directly access data derived from several source systems through the data warehouse.

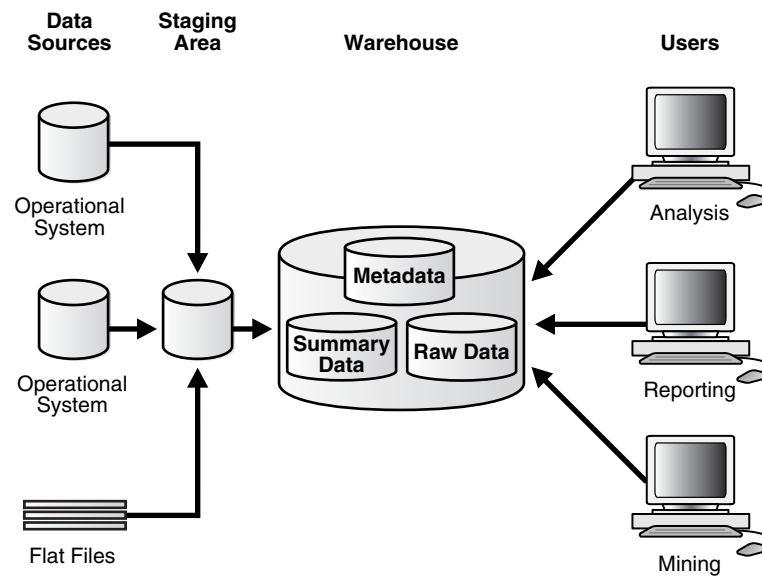
Figure 1–1 Architecture of a Data Warehouse

In [Figure 1–1](#), the metadata and raw data of a traditional OLTP system is present, as is an additional type of data, summary data. Summaries are a mechanism to pre-compute common expensive, long-running operations for sub-second data retrieval. For example, a typical data warehouse query is to retrieve something such as August sales. A summary in an Oracle database is called a **materialized view**.

The consolidated storage of the raw data as the center of your data warehousing architecture is often referred to as an Enterprise Data Warehouse (EDW). An EDW provides a 360-degree view into the business of an organization by holding all relevant business information in the most detailed format.

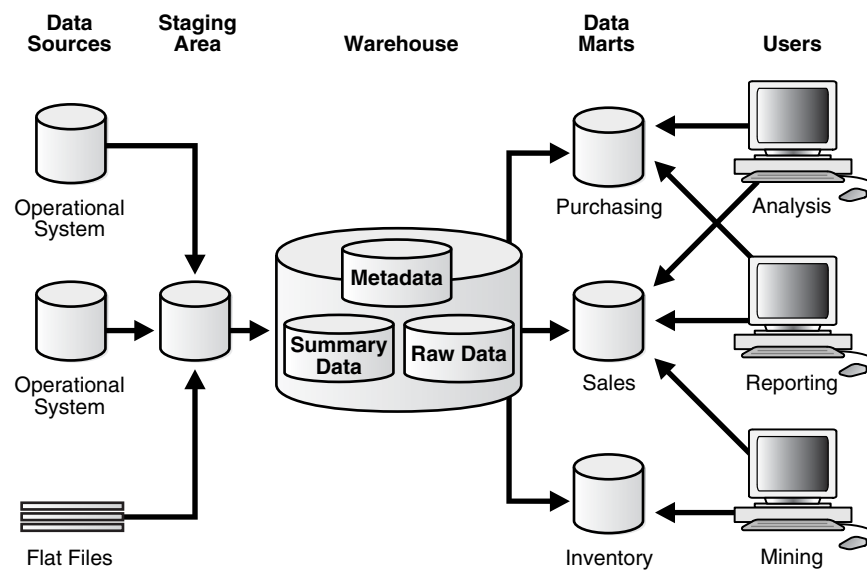
Data Warehouse Architecture: with a Staging Area

You must clean and process your operational data before putting it into the warehouse, as shown in [Figure 1–2](#). You can do this programmatically, although most data warehouses use a **staging area** instead. A staging area simplifies data cleansing and consolidation for operational data coming from multiple source systems, especially for enterprise data warehouses where all relevant information of an enterprise is consolidated. [Figure 1–2](#) illustrates this typical architecture.

Figure 1-2 Architecture of a Data Warehouse with a Staging Area

Data Warehouse Architecture: with a Staging Area and Data Marts

Although the architecture in [Figure 1-2](#) is quite common, you may want to customize your warehouse's architecture for different groups within your organization. You can do this by adding **data marts**, which are systems designed for a particular line of business. [Figure 1-3](#) illustrates an example where purchasing, sales, and inventories are separated. In this example, a financial analyst might want to analyze historical data for purchases and sales or mine historical data to make predictions about customer behavior.

Figure 1-3 Architecture of a Data Warehouse with a Staging Area and Data Marts

Note: Data marts can be physically instantiated or implemented purely logically through views. Furthermore, data marts can be co-located with the enterprise data warehouse or built as separate systems. Building an end-to-end data warehousing architecture with an enterprise data warehouse and surrounding data marts is not the focus of this book.

Data Warehousing Logical Design

This chapter explains how to create a logical design for a data warehousing environment and includes the following topics:

- [Logical Versus Physical Design in Data Warehouses](#)
- [Creating a Logical Design](#)
- [About Third Normal Form Schemas](#)
- [About Star Schemas](#)
- [About the Oracle In-Memory Column Store](#)
- [Automatic Big Table Caching to Improve the Performance of In-Memory Parallel Queries](#)
- [About In-Memory Aggregation](#)

Logical Versus Physical Design in Data Warehouses

Your organization has decided to build an enterprise data warehouse. You have defined the business requirements and agreed upon the scope of your business goals, and created a conceptual design. Now you need to translate your requirements into a system deliverable. To do so, you create the logical and physical design for the data warehouse. You then define:

- The specific data content
- Relationships within and between groups of data
- The system environment supporting your data warehouse
- The data transformations required
- The frequency with which data is refreshed

The logical design is more conceptual and abstract than the physical design. In the logical design, you look at the logical relationships among the objects. In the physical design, you look at the most effective way of storing and retrieving the objects as well as handling them from a transportation and backup/recovery perspective.

Orient your design toward the needs of the end users. End users typically want to perform analysis and look at aggregated data, rather than at individual transactions. However, end users might not know what they need until they see it. In addition, a well-planned design allows for growth and changes as the needs of users change and evolve.

By beginning with the logical design, you focus on the information requirements and save the implementation details for later.

Creating a Logical Design

A logical design is conceptual and abstract. You do not deal with the physical implementation details yet. You deal only with defining the types of information that you need.

One technique you can use to model your organization's logical information requirements is entity-relationship modeling. Entity-relationship modeling involves identifying the things of importance (entities), the properties of these things (attributes), and how they are related to one another (relationships).

The process of logical design involves arranging data into a series of logical relationships called entities and attributes. An entity represents a chunk of information. In relational databases, an entity often maps to a table. An attribute is a component of an entity that helps define the uniqueness of the entity. In relational databases, an attribute maps to a column.

To ensure that your data is consistent, you must use unique identifiers. A unique identifier is something you add to tables so that you can differentiate between the same item when it appears in different places. In a physical design, this is usually a primary key.

Entity-relationship modeling is purely logical and applies to both OLTP and data warehousing systems. It is also applicable to the various common physical schema modeling techniques found in data warehousing environments, namely normalized (3NF) schemas in Enterprise Data Warehousing environments, star or snowflake schemas in data marts, or hybrid schemas with components of both of these classical modeling techniques.

See Also:

- *Oracle Fusion Middleware Developer's Guide for Oracle Data Integrator* for more details regarding ODI

What is a Schema?

A schema is a collection of database objects, including tables, views, indexes, and synonyms. You can arrange schema objects in the schema models designed for data warehousing in a variety of ways. Most data warehouses use a dimensional model.

The model of your source data and the requirements of your users help you design the data warehouse schema. You can sometimes get the source model from your company's enterprise data model and reverse-engineer the logical data model for the data warehouse from this. The physical implementation of the logical data warehouse model may require some changes to adapt it to your system parameters—size of computer, number of users, storage capacity, type of network, and software. A key part of designing the schema is whether to use a third normal form, star, or snowflake schema, and these are discussed later.

About Third Normal Form Schemas

Third Normal Form design seeks to minimize data redundancy and avoid anomalies in data insertion, updates and deletion. 3NF design has a long heritage in online transaction processing (OLTP) systems. OLTP systems must maximize performance and accuracy when inserting, updating and deleting data. Transactions must be handled as quickly as possible or the business may be unable to handle the flow of events, perhaps losing sales or incurring other costs. Therefore, 3NF designs avoid redundant data manipulation and minimize table locks, both of which can slow

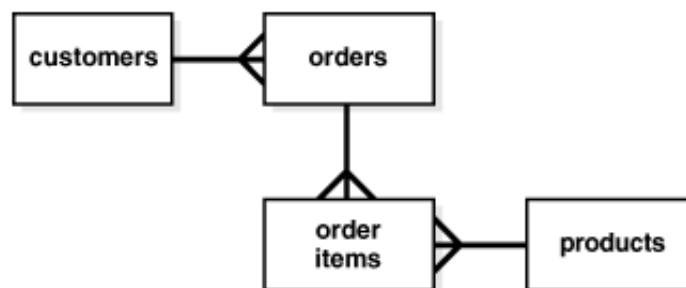
inserts, updates and deletes. 3NF designs also work well to abstract the data from specific application needs. If new types of data are added to the environment, you can extend the data model with relative ease and minimal impact to existing applications. Likewise, if you have completely new types of analyses to perform in your data warehouse, a well-designed 3NF schema will be able to handle them without requiring redesigned data structures.

3NF designs have great flexibility, but it comes at a cost. 3NF databases use very many tables and this requires complex queries with many joins. For full scale enterprise models built in 3NF form, over one thousand tables are commonly encountered in the schema. With the kinds of queries involved in data warehousing, which will often need access to many rows from many tables, this design imposes understanding and performance penalties. It can be complex for query builders, whether they are humans or business intelligence tools and applications, to choose and join the tables needed for a given piece of data when there are very large numbers of tables available. Even when the tables are readily chosen by the query generator, the 3NF schema often requires that a large number of tables be used in a single query. More tables in a query mean more potential data access paths, which makes the database query optimizer's job harder. The end result can be slow query performance.

The issue of slow query performance in a 3NF system is not necessarily limited to the core queries used to create reports and analyses. It can also show up in the simpler task of users browsing subsets of data to understand the contents. Similarly, the complexity of a 3NF schema may impact generating the pick-lists of data used to constrain queries and reports. Although these may seem relatively minor issues, speedy response time for such processes makes a big impact on user satisfaction.

Figure 2-1 presents a tiny fragment of a 3NF Schema. Note how order information is broken into order and order items to avoid redundant data storage. The "crow's feet" markings on the relationship between tables indicate one-to-many relationships among the entities. Thus, one order may have multiple order items, a single customer may have many orders, and a single product may be found in many order items. Although this diagram shows a very small case, you can see that minimizing data redundancy can lead to many tables in the schema.

Figure 2-1 *Fragment of a Third Normal Form Schema*



About Normalization

Normalization is a data design process that has a high level goal of keeping each fact in just one place to avoid data redundancy and insert, update, and delete anomalies. There are multiple levels of normalization, and this section describes the first three of them. Considering how fundamental the term third normal form (3NF) term is, it only makes sense to see how 3NF is reached.

Consider a situation where you are tracking sales. The core entity you track is sales orders, where each sales order contains details about each item purchased (referred to

as a line item): its name, price, quantity, and so on. The order also holds the name and address of the customer and more. Some orders have many different line items, and some orders have just one.

In first normal form (1NF), there are no repeating groups of data and no duplicate rows. Every intersection of a row and column (a field) contains just one value, and there are no groups of columns that contain the same facts. To avoid duplicate rows, there is a primary key. For sales orders, in first normal form, multiple line items of each sales order in a single field of the table are not displayed. Also, there will not be multiple columns showing line items.

Then comes second normal form (2NF), where the design is in first normal form and every non-key column is dependent on the complete primary key. Thus, the line items are broken out into a table of sales order line items where each row represents one line item of one order. You can look at the line item table and see that the names of the items sold are not dependent on the primary key of the line items table: the sales item is its own entity. Therefore, you move the sales item to its own table showing the item name. Prices charged for each item can vary by order (for instance, due to discounts) so these remain in the line items table. In the case of sales order, the name and address of the customer is not dependent on the primary key of the sales order: customer is its own entity. Thus, you move the customer name and address columns out into their own table of customer information.

Next is third normal form, where the goal is to ensure that there are no dependencies on non-key attributes. So the goal is to take columns that do not directly relate to the subject of the row (the primary key), and put them in their own table. So details about customers, such as customer name or customer city, should be put in a separate table, and then a customer foreign key added into the orders table.

Another example of how a 2NF table differs from a 3NF table would be a table of the winners of tennis tournaments that contained columns of tournament, year, winner, and winner's date of birth. In this case, the winner's date of birth is vulnerable to inconsistencies, as the same person could be shown with different dates of birth in different records. The way to avoid this potential problem is to break the table into one for tournament winners, and another for the player dates of birth.

Design Concepts for 3NF Schemas

The following section discusses some basic concepts when modeling for a data warehousing environment using a 3NF schema approach. The intent is not to discuss the theoretical foundation for 3NF modeling (or even higher levels of normalization), but to highlight some key components relevant for data warehousing.

Identifying Candidate Primary Keys

A primary key is an attribute that uniquely identifies a specific record in a table. Primary keys can be identified through single or multiple columns. It is normally preferred to achieve unique identification through as little columns as possible - ideally one or two - and to either use a column that is most likely not going to be updated or even changed in bulk. If your data model does not lead to a simple unique identification through its attributes, you would require too many attributes to uniquely identify a single records, or the data is prone to changes, the usage of a surrogate key is highly recommended.

Specifically, 3NF schemas rely on proper and simple unique identification since queries tend to have many table joins and all columns necessary to uniquely identify a record are needed as join condition to avoid row duplication through the join.

Foreign Key Relationships and Referential Integrity Constraints

3NF schemas in data warehousing environments often resemble the data model of its OLTP source systems, in which the logical consistency between data entities is expressed and enforced through primary key - foreign key relationships, also known as parent-child relationship. A foreign key resolves a 1-to-many relationship in relational system and ensures logical consistency: for example, you cannot have an order line item without an order header, or an employee working for a non-existent department.

While such referential are always enforced in OLTP system, data warehousing systems often implement them as declarative, non-enforced conditions, relying on the ETL process to ensure data consistency. Whenever possible, foreign keys and referential integrity constraints should be defined as non-enforced conditions, since it enables better query optimization and cardinality estimates.

Denormalization

Proper normalized modelling tends to decompose logical entities - such as a customer, a product, or an order - into many physical tables, making even the retrieval of perceived simple information requiring to join many tables. While this is not a problem from a query processing perspective, it can put some unnecessary burden on both the application developer (for writing code) as well as the database (for joining information that is always used together). It is not uncommon to see some sensible level of denormalization in 3NF data warehousing models, in a logical form as views or in a physical form through slightly denormalized tables.

Care has to be taken with the physical denormalization to preserve the subject-neutral shape and therefore the flexibility of the physical implementation of the 3NF schema.

About Star Schemas

Star schemas are often found in data warehousing systems with embedded logical or physical data marts. The term star schema is another way of referring to a "dimensional modeling" approach to defining your data model. Most descriptions of dimensional modeling use terminology drawn from the work of Ralph Kimball, the pioneering consultant and writer in this field. Dimensional modeling creates multiple star schemas, each based on a business process such as sales tracking or shipments. Each star schema can be considered a data mart, and perhaps as few as 20 data marts can cover the business intelligence needs of an enterprise. Compared to 3NF designs, the number of tables involved in dimensional modeling is a tiny fraction. Many star schemas will have under a dozen tables. The star schemas are knit together through conformed dimensions and conformed facts. Thus, users are able to get data from multiple star schemas with minimal effort.

The goal for star schemas is structural simplicity and high performance data retrieval. Because most queries in the modern era are generated by reporting tools and applications, it's vital to make the query generation convenient and reliable for the tools and application. In fact, many business intelligence tools and applications are designed with the expectation that a star schema representation will be available to them.

Discussions of star schemas are less abstracted from the physical database than 3NF descriptions. This is due to the pragmatic emphasis of dimensional modeling on the needs of business intelligence users.

Note how different the dimensional modeling style is from the 3NF approach that minimizes data redundancy and the risks of update/inset/delete anomalies. The star schema accepts data redundancy (denormalization) in its dimension tables for the sake

of easy user understanding and better data retrieval performance. A common criticism of star schemas is that they limit analysis flexibility compared to 3NF designs. However, a well designed dimensional model can be extended to enable new types of analysis, and star schemas have been successful for many years at the largest enterprises.

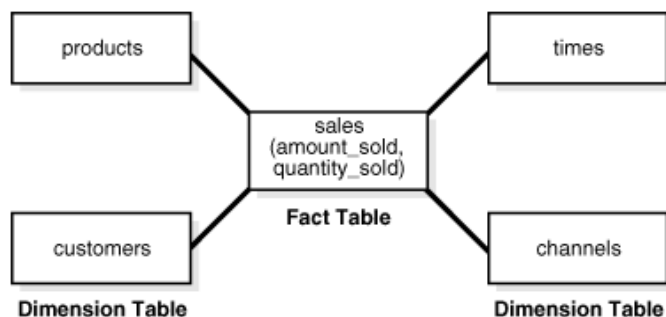
As noted earlier, the modern approach to data warehousing does not pit star schemas and 3NF against each other. Rather, both techniques are used, with a foundation layer of 3NF - the Enterprise Data Warehouse of 3NF, acting as the bedrock data, and star schemas as a central part of an access and performance optimization layer.

Facts and Dimensions

Star schemas divide data into facts and dimensions. Facts are the measurements of some event such as a sale and are typically numbers. Dimensions are the categories you use to identify facts, such as date, location, and product.

The name "star schema" comes from the fact that the diagrams of the schemas typically show a central fact table with lines joining it to the dimension tables, so the graphic impression is similar to a star. [Figure 2–2](#) is a simple example with sales as the fact table and products, times, customers, and channels as the dimension table.

Figure 2–2 Star Schema



Fact Tables

Fact tables have measurement data. They have many rows but typically not many columns. Fact tables for a large enterprise can easily hold billions of rows. For many star schemas, the fact table will represent well over 90 percent of the total storage space. A fact table has a composite key made up of the primary keys of the dimension tables of the schema.

A fact table contains either detail-level facts or facts that have been aggregated. Fact tables that contain aggregated facts are often called summary tables. A fact table usually contains facts with the same level of aggregation. Though most facts are additive, they can also be semi-additive or non-additive. Additive facts can be aggregated by simple arithmetical addition. A common example of this is sales. Non-additive facts cannot be added at all. An example of this is averages. Semi-additive facts can be aggregated along some of the dimensions and not along others. An example of this is inventory levels stored in physical warehouses, where you may be able to add across a dimension of warehouse sites, but you cannot aggregate across time.

In terms of adding rows to data in a fact table, there are three main approaches:

- Transaction-based

Shows a row for the finest level detail in a transaction. A row is entered only if a transaction has occurred for a given combination of dimension values. This is the most common type of fact table.

- Periodic Snapshot

Shows data as of the end of a regular time interval, such as daily or weekly. If a row for the snapshot exists in a prior period, a row is entered for it in the new period even if no activity related to it has occurred in the latest interval. This type of fact table is useful in complex business processes where it is difficult to compute snapshot values from individual transaction rows.

- Accumulating Snapshot

Shows one row for each occurrence of a short-lived process. The rows contain multiple dates tracking major milestones of a short-lived process. Unlike the other two types of fact tables, rows in an accumulating snapshot are updated multiple times as the tracked process moves forward.

Dimension Tables

Dimension tables provide category data to give context to the fact data. For instance, a star schema for sales data will have dimension tables for product, date, sales location, promotion and more. Dimension tables act as lookup or reference tables because their information lets you choose the values used to constrain your queries. The values in many dimension tables may change infrequently. As an example, a dimension of geographies showing cities may be fairly static. But when dimension values do change, it is vital to update them fast and reliably. Of course, there are situations where data warehouse dimension values change frequently. The customer dimension for an enterprise will certainly be subject to a frequent stream of updates and deletions.

A key aspect of dimension tables is the hierarchy information they provide. Dimension data typically has rows for the lowest level of detail plus rows for aggregated dimension values. These natural rollups or aggregations within a dimension table are called hierarchies and add great value for analyses. For instance, if you want to calculate the share of sales that a specific product represents within its specific product category, it is far easier and more reliable to have a predefined hierarchy for product aggregation than to specify all the elements of the product category in each query. Because hierarchy information is so valuable, it is common to find multiple hierarchies reflected in a dimension table.

Dimension tables are usually textual and descriptive, and you will use their values as the row headers, column headers and page headers of the reports generated by your queries. While dimension tables have far fewer rows than fact tables, they can be quite wide, with dozens of columns. A location dimension table might have columns indicating every level of its rollup hierarchy, and may show multiple hierarchies reflected in the table. The location dimension table could have columns for its geographic rollup, such as street address, postal code, city, state/province, and country. The same table could include a rollup hierarchy set up for the sales organization, with columns for sales district, sales territory, sales region, and characteristics.

See Also: [Chapter 9, "Dimensions"](#) for further information regarding dimensions

Design Concepts in Star Schemas

Here we touch on some of the key terms used in star schemas. This is by no means a full set, but is intended to highlight some of the areas worth your consideration.

Data Grain

One of the most important tasks when designing your model is to consider the level of detail it will provide, referred to as the grain of the data. Consider a sales schema: will the grain be very fine, storing every single item purchased by each customer? Or will it be a coarse grain, storing only the daily totals of sales for each product at each store? In modern data warehousing there is a strong emphasis on providing the finest grain data possible, because this allows for maximum analytic power. Dimensional modeling experts generally recommend that each fact table store just one grain level. Presenting fact data in single-grain tables supports more reliable querying and table maintenance, because there is no ambiguity about the scope of any row in a fact table.

Working with Multiple Star Schemas

Because the star schema design approach is intended to chunk data into distinct processes, you need reliable and performant ways to traverse the schemas when queries span multiple schemas. One term for this ability is a data warehouse bus architecture. A data warehouse bus architecture can be achieved with conformed dimensions and conformed facts.

Conformed Dimensions

Conformed dimensions means that dimensions are designed identically across the various star schemas. Conformed dimensions use the same values, column names and data types consistently across multiple stars. The conformed dimensions do not have to contain the same number of rows in each schema's copy of the dimension table, as long as the rows in the shorter tables are a true subset of the larger tables.

Conformed Facts

If the fact columns in multiple fact tables have exactly the same meaning, then they are considered conformed facts. Such facts can be used together reliably in calculations even though they are from different tables. Conformed facts should have the same column names to indicate their conformed status. Facts that are not conformed should always have different names to highlight their different meanings.

Surrogate Keys

Surrogate or artificial keys, usually sequential integers, are recommended for dimension tables. By using surrogate keys, the data is insulated from operational changes. Also, compact integer keys may allow for better performance than large and complex alphanumeric keys.

Degenerate Dimensions

Degenerate dimensions are dimension columns in fact tables that do not join to a dimension table. They are typically items such as order numbers and invoice numbers. You will see them when the grain of a fact table is at the level of an order line-item or a single transaction.

Junk Dimensions

Junk dimensions are abstract dimension tables used to hold text lookup values for flags and codes in fact tables. These dimensions are referred to as junk, not because they have low value, but because they hold an assortment of columns for convenience, analogous to the idea of a "junk drawer" in your home. The number of distinct values (cardinality) of each column in a junk dimension table is typically small.

Embedded Hierarchy

Classic dimensional modeling with star schemas advocates that each table contain data at a single grain. However, there are situations where designers choose to have multiple grains in a table, and these commonly represent a rollup hierarchy. A single sales fact table, for instance, might contain both transaction-level data, then a day-level rollup by product, then a month-level rollup by product. In such cases, the fact table will need to contain a level column indicating the hierarchy level applying to each row, and queries against the table will need to include a level predicate.

Factless Fact Tables

Factless fact tables do not contain measures such as sales price or quantity sold. Instead, the rows of a factless fact table are used to show events not represented by other fact tables. Another use for factless tables is as a "coverage table" which holds all the possible events that could have occurred in a given situation, such as all the products that were part of a sales promotion and might have been sold at the promotional price.

Slowly Changing Dimensions

One of the certainties of data warehousing is that the way data is categorized will change. Product names and category names will change. Characteristics of a store will change. The areas included in sales territories will change. The timing and extent of these changes will not always be predictable. How can these slowly changing dimensions be handled? Star schemas treat these in three main ways:

- Type 1

The dimension values that change are simply overwritten, with no history kept. This creates a problem for time-based analyses. Also, it invalidates any existing aggregates that depended on the old value of the dimension.

- Type 2

When a dimension value changes, a new dimension row showing the new value and having a new surrogate key is created. You may choose to include date columns in our dimension showing when the new row is valid and when it is expired. No changes need be made to the fact table.

- Type 3

When a dimension value is changed, the prior value is stored in a different column of the same row. This enables easy query generation if you want to compare results using the current and prior value of the column.

In practice, Type 2 is the most common treatment for slowly changing dimensions.

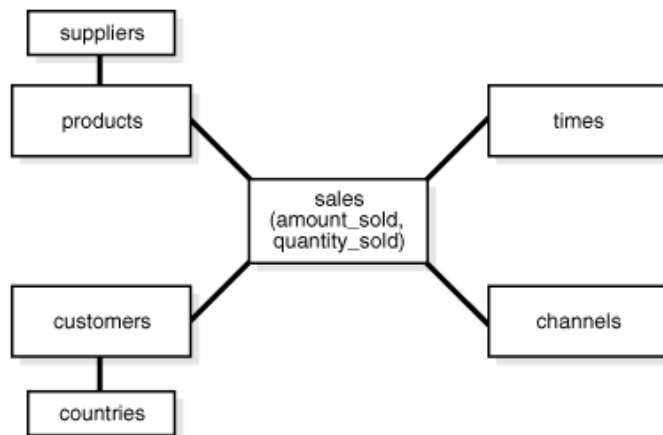
About Snowflake Schemas

The snowflake schema is a more complex data warehouse model than a star schema, and is a type of star schema. It is called a snowflake schema because the diagram of the schema resembles a snowflake.

Snowflake schemas normalize dimensions to eliminate redundancy. That is, the dimension data has been grouped into multiple tables instead of one large table. For example, a product dimension table in a star schema might be normalized into a `products` table, a `product_category` table, and a `product_manufacturer` table in a snowflake schema. While this saves space, it increases the number of dimension tables and requires more foreign key joins. The result is more complex queries and reduced query performance. [Figure 2–3](#) presents a graphical representation of a snowflake

schema.

Figure 2–3 Snowflake Schema



About the Oracle In-Memory Column Store

Note: This feature is available starting with Oracle Database 12c Release 1 (12.1.0.2).

Traditional analytics has certain limitations or requirements that need to be managed to obtain good performance for analytic queries. You need to know user access patterns and then customize your data structures to provide optimal performance for these access patterns. Existing indexes, materialized views, and OLAP cubes need to be tuned. Certain data marts and reporting databases have complex ETL and thus need specialized tuning. Additionally, you need to strike a balance between performing analytics on stale data and slowing down OLTP operations on the production databases.

The Oracle In-Memory Column Store (IM column store) within the Oracle Database provides improved performance for both ad-hoc queries and analytics on live data. The live transactional database is used to provide instant answers to queries, thus enabling you to seamlessly use the same database for OLTP transactions and data warehouse analytics.

The IM column store is an optional area in the SGA that stores copies of tables, table partitions, and individual columns in a compressed columnar format that is optimized for rapid scans. Columnar format lends itself to easily to vector processing thus making aggregations, joins, and certain types of data retrieval faster than the traditional on-disk formats. The columnar format exists only in memory and does not replace the on-disk or buffer cache format. Instead, it supplements the buffer cache and provides an additional, transaction-consistent, copy of the table that is independent of the disk format.

See Also: *Oracle Database Concepts* for conceptual information the about IM column store

Configuring the Oracle In-Memory Column Store

Configuring the IM column store is simple. You set the `INMEMORY_SIZE` initialization parameter, and then use DDL to specify the tablespaces, tables, partitions, or columns to be populated into the IM column store.

See Also: *Oracle Database Administrator's Guide* for information about configuring the IM column store

Populating the Oracle In-Memory Column Store

You can specify that the database populates data into the IM column store from row storage either at database instance startup or when the data is accessed for the first time.

See Also: *Oracle Database Concepts* for detailed information about how the IM column store is populated

In-Memory Columnar Compression

The Oracle Databases uses special compression formats that are optimized for access speeds rather than storage reductions to store data in the IM column store. You can select different compression options for each table, partition, or column.

See Also:

- *Oracle Database Concepts*
- *Oracle Database Administrator's Guide*

Benefits of Using the Oracle In-Memory Column Store

The IM column store enables the Oracle Database to perform scans, joins, and aggregates much faster than when it uses the on-disk format exclusively. Business applications, ad-hoc analytic queries, and data warehouse workloads benefit most. Pure OLTP databases that perform short transactions using index lookups benefit less.

The IM column store seamlessly integrates with the Oracle Database. All existing database features, including High Availability features, are supported with no application changes required. Therefore, by configuring the IM column store, you can instantly improve the performance of existing analytic workloads and ad-hoc queries.

The Oracle Optimizer is aware of the IM column store making it possible for the Oracle Database to seamlessly send analytic queries to the IM column store while OLTP queries and DML are sent to the row store.

The advantages offered by the IM column store for data warehousing environments are:

- Faster scanning of large number of rows and applying filters that use operators such as `=`, `<`, `>`, and `IN`.
- Faster querying of a subset of columns in a table, for example, selecting 5 of 100 columns. See ["Faster Performance for Analytic Queries"](#) on page 2-12.
- Enhanced performance for joins by converting predicates on small dimension tables to filters on a large fact table. See ["Enhanced Join Performance Using Vector Joins"](#) on page 2-12.
- Efficient aggregation by using `VECTOR GROUP BY` transformation and vector array processing. See ["Enhanced Aggregation Using VECTOR GROUP BY Transformations"](#) on page 2-12.