# DWH
## data warehousing

# tutorialspoint
## SIMPLYEASYLEARNING

## About the Tutorial

A data warehouse is constructed by integrating data from multiple heterogeneous sources. It supports analytical reporting, structured and/or ad hoc queries and decision making. This tutorial adopts a step-by-step approach to explain all the necessary concepts of data warehousing.

## Audience

This tutorial will help computer science graduates to understand the basic-to-advanced concepts related to data warehousing.

## Prerequisites

Before proceeding with this tutorial, you should have an understanding of basic database concepts such as schema, ER model, structured query language, etc.

## Copyright & Disclaimer

# Table of Contents

About the Tutorial.................................................................................................................................i

Audience ..............................................................................................................................................i

Prerequisites .......................................................................................................................................i

Copyright & Disclaimer........................................................................................................................i

Table of Contents ...............................................................................................................................ii

1.   DWH – OVERVIEW...........................................................................................................1

Understanding a Data Warehouse .....................................................................................................1

Why a Data Warehouse is Separated from Operational Databases ....................................................1

Data Warehouse Features...................................................................................................................2

Data Warehouse Applications ............................................................................................................2

Types of Data Warehouse ..................................................................................................................2

2.   DWH – CONCEPTS ..........................................................................................................4

What is Data Warehousing? ...............................................................................................................4

Using Data Warehouse Information ...................................................................................................4

Integrating Heterogeneous Databases ..............................................................................................4

Functions of Data Warehouse Tools and Utilities ..............................................................................5

3.   DWH – TERMINOLOGIES ................................................................................................6

Metadata ...........................................................................................................................................6

Metadata Repository .........................................................................................................................6

Data Cube ..........................................................................................................................................6

Data Mart...........................................................................................................................................8

Virtual Warehouse .............................................................................................................................9

4.   DWH – DELIVERY PROCESS...........................................................................................10

Delivery Method ..............................................................................................................................10

The term "Data Warehouse" was first coined by Bill Inmon in 1990. According to Inmon, a data warehouse is a subject-oriented, integrated, time-variant, and non-volatile collection of data. This data helps analysts to take informed decisions in an organization.

An operational database undergoes frequent changes on a daily basis on account of the transactions that take place. Suppose a business executive wants to analyze previous feedback on any data such as a product, a supplier, or any consumer data, then the executive will have no data available to analyze because the previous data has been updated due to transactions.

A data warehouses provides us generalized and consolidated data in multidimensional view. Along with generalized and consolidated view of data, a data warehouses also provides us Online Analytical Processing (OLAP) tools. These tools help us in interactive and effective analysis of data in a multidimensional space. This analysis results in data generalization and data mining.

Data mining functions such as association, clustering, classification, prediction can be integrated with OLAP operations to enhance the interactive mining of knowledge at multiple level of abstraction. That's why data warehouse has now become an important platform for data analysis and online analytical processing.

## Understanding a Data Warehouse

- A data warehouse is a database, which is kept separate from the organization's operational database.

- There is no frequent updating done in a data warehouse.

- It possesses consolidated historical data, which helps the organization to analyze its business.

- A data warehouse helps executives to organize, understand, and use their data to take strategic decisions.

- Data warehouse systems help in the integration of diversity of application systems.

- A data warehouse system helps in consolidated historical data analysis.

## Why a Data Warehouse is Separated from Operational Databases

A data warehouses is kept separate from operational databases due to the following reasons:

- An operational database is constructed for well-known tasks and workloads such as searching particular records, indexing, etc. In contrast, data warehouse queries are often complex and they present a general form of data.

- Operational databases support concurrent processing of multiple transactions. Concurrency control and recovery mechanisms are required for operational databases to ensure robustness and consistency of the database.

- An operational database query allows to read and modify operations, while an OLAP query needs only **read only** access of stored data.

- An operational database maintains current data. On the other hand, a data warehouse maintains historical data.

## Data Warehouse Features

The key features of a data warehouse are discussed below:

- **Subject Oriented** - A data warehouse is subject oriented because it provides information around a subject rather than the organization's ongoing operations. These subjects can be product, customers, suppliers, sales, revenue, etc. A data warehouse does not focus on the ongoing operations, rather it focuses on modelling and analysis of data for decision making.

- **Integrated** – A data warehouse is constructed by integrating data from heterogeneous sources such as relational databases, flat files, etc. This integration enhances the effective analysis of data.

- **Time Variant** - The data collected in a data warehouse is identified with a particular time period. The data in a data warehouse provides information from the historical point of view.

- **Non-volatile** - Non-volatile means the previous data is not erased when new data is added to it. A data warehouse is kept separate from the operational database and therefore frequent changes in operational database is not reflected in the data warehouse.

**Note:** A data warehouse does not require transaction processing, recovery, and concurrency controls, because it is physically stored and separate from the operational database.

## Data Warehouse Applications

As discussed before, a data warehouse helps business executives to organize, analyze, and use their data for decision making. A data warehouse serves as a sole part of a plan-execute-assess "closed-loop" feedback system for the enterprise management. Data warehouses are widely used in the following fields:

- Financial services

- Banking services

- Consumer goods

- Retail sectors

- Controlled manufacturing

# Types of Data Warehouse

Information processing, analytical processing, and data mining are the three types of data warehouse applications that are discussed below:

- **Information Processing** – A data warehouse allows to process the data stored in it. The data can be processed by means of querying, basic statistical analysis, reporting using crosstabs, tables, charts, or graphs.

- **Analytical Processing** – A data warehouse supports analytical processing of the information stored in it. The data can be analyzed by means of basic OLAP operations, including slice-and-dice, drill down, drill up, and pivoting.

- **Data Mining** - Data mining supports knowledge discovery by finding hidden patterns and associations, constructing analytical models, performing classification and prediction. These mining results can be presented using visualization tools.

| Data Warehouse (OLAP) | Operational Database(OLTP) |
|---|---|
| It involves historical processing of information. | It involves day-to-day processing. |
| OLAP systems are used by knowledge workers such as executives, managers, and analysts. | OLTP systems are used by clerks, DBAs, or database professionals. |
| It is used to analyze the business. | It is used to run the business. |
| It focuses on Information out. | It focuses on Data in. |
| It is based on Star Schema, Snowflake Schema, and Fact Constellation Schema. | It is based on Entity Relationship Model. |
| It focuses on Information out. | It is application oriented. |
| It contains historical data. | It contains current data. |
| It provides summarized and consolidated data. | It provides primitive and highly detailed data. |
| It provides summarized and multidimensional view of data. | It provides detailed and flat relational view of data. |

| | |
|---|---|
| The number of users is in hundreds. | The number of users is in thousands. |
| The number of records accessed is in millions. | The number of records accessed is in tens. |
| The database size is from 100GB to 100 TB. | The database size is from 100 MB to 100 GB. |
| These are highly flexible. | It provides high performance. |

## What is Data Warehousing?

Data warehousing is the process of constructing and using a data warehouse. A data warehouse is constructed by integrating data from multiple heterogeneous sources that support analytical reporting, structured and/or ad hoc queries, and decision making. Data warehousing involves data cleaning, data integration, and data consolidations.

## Using Data Warehouse Information

There are decision support technologies that help utilize the data available in a data warehouse. These technologies help executives to use the warehouse quickly and effectively. They can gather data, analyze it, and take decisions based on the information present in the warehouse. The information gathered in a warehouse can be used in any of the following domains:

- **Tuning Production Strategies** - The product strategies can be well tuned by repositioning the products and managing the product portfolios by comparing the sales quarterly or yearly.

- **Customer Analysis** - Customer analysis is done by analyzing the customer's buying preferences, buying time, budget cycles, etc.

- **Operations Analysis** - Data warehousing also helps in customer relationship management, and making environmental corrections. The information also allows us to analyze business operations.

## Integrating Heterogeneous Databases

To integrate heterogeneous databases, we have two approaches:

- Query-driven Approach

- Update-driven Approach

### Query-Driven Approach

This is the traditional approach to integrate heterogeneous databases. This approach was used to build wrappers and integrators on top of multiple heterogeneous databases. These integrators are also known as mediators.

### Process of Query-Driven Approach

1. When a query is issued to a client side, a metadata dictionary translates the query into an appropriate form for individual heterogeneous sites involved.

2. Now these queries are mapped and sent to the local query processor.

3. The results from heterogeneous sites are integrated into a global answer set.

## Disadvantages

- Query-driven approach needs complex integration and filtering processes.

- This approach is very inefficient.

- It is very expensive for frequent queries.

- This approach is also very expensive for queries that require aggregations.

## Update-Driven Approach

This is an alternative to the traditional approach. Today's data warehouse systems follow update-driven approach rather than the traditional approach discussed earlier. In update-driven approach, the information from multiple heterogeneous sources are integrated in advance and are stored in a warehouse. This information is available for direct querying and analysis.

## Advantages

This approach has the following advantages:

- This approach provides high performance.

- The data is copied, processed, integrated, annotated, summarized and restructured in semantic data store in advance.

- Query processing does not require an interface to process data at local sources.

# Functions of Data Warehouse Tools and Utilities

The following are the functions of data warehouse tools and utilities:

- **Data Extraction** - Involves gathering data from multiple heterogeneous sources.

- **Data Cleaning** - Involves finding and correcting the errors in data.

- **Data Transformation** - Involves converting the data from legacy format to warehouse format.

- **Data Loading** - Involves sorting, summarizing, consolidating, checking integrity, and building indices and partitions.

- **Refreshing** - Involves updating from data sources to warehouse.

**Note:** Data cleaning and data transformation are important steps in improving the quality of data and data mining results.

# 3. DWH — TERMINOLOGIES

In this chapter, we will discuss some of the most commonly used terms in data warehousing.

## Metadata

Metadata is simply defined as data about data. The data that are used to represent other data is known as metadata. For example, the index of a book serves as a metadata for the contents in the book. In other words, we can say that metadata is the summarized data that leads us to the detailed data.

In terms of data warehouse, we can define metadata as following:

- Metadata is a roadmap to data warehouse.

- Metadata in data warehouse defines the warehouse objects.

- Metadata acts as a directory. This directory helps the decision support system to locate the contents of a data warehouse.

## Metadata Repository

Metadata repository is an integral part of a data warehouse system. It contains the following metadata:

- **Business metadata** - It contains the data ownership information, business definition, and changing policies.

- **Operational metadata** - It includes currency of data and data lineage. Currency of data refers to the data being active, archived, or purged. Lineage of data means history of data migrated and transformation applied on it.

- **Data for mapping from operational environment to data warehouse** - It metadata includes source databases and their contents, data extraction, data partition, cleaning, transformation rules, data refresh and purging rules.

- **The algorithms for summarization** - It includes dimension algorithms, data on granularity, aggregation, summarizing, etc.

## Data Cube

A data cube helps us represent data in multiple dimensions. It is defined by dimensions and facts. The dimensions are the entities with respect to which an enterprise preserves the records.

### Illustration of Data Cube

Suppose a company wants to keep track of sales records with the help of sales data warehouse with respect to time, item, branch, and location. These dimensions allow to keep track of monthly sales and at which branch the items were sold. There is a table associated with each dimension. This table is known as dimension table. For example, "item" dimension table may have attributes such as item_name, item_type, and item_brand.

The following table represents the 2-D view of Sales Data for a company with respect to time, item, and location dimensions.

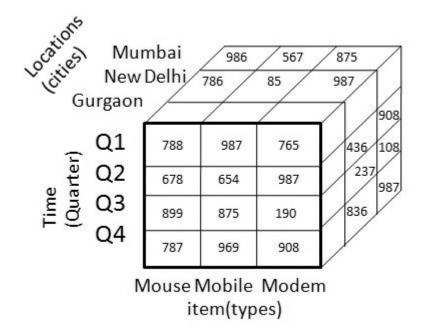| Location="New Delhi" | | | | |
|---|---|---|---|---|
| | Item(type) | | | |
| Time(quarter) | Entertainment | Keyboard | Mobile | Locks |
| Q1 | 500 | 700 | 10 | 300 |
| Q2 | 769 | 765 | 30 | 476 |
| Q3 | 987 | 489 | 18 | 659 |
| Q4 | 666 | 976 | 40 | 539 |

But here in this 2-D table, we have records with respect to time and item only. The sales for New Delhi are shown with respect to time, and item dimensions according to type of items sold. If we want to view the sales data with one more dimension, say, the location dimension, then the 3-D view would be useful. The 3-D view of the sales data with respect to time, item, and location is shown in the table below:

| Time | Location="Gurgaon" | | | Location="New Delhi" | | | Location="Mumbai" | | |
|------|-------|--------|-------|-------|--------|-------|-------|--------|-------|
| | Item | | | Item | | | Item | | |
| | Mouse | Mobile | Modem | Mouse | Mobile | Modem | Mouse | Mobile | Modem |
| Q1 | 788 | 987 | 765 | 786 | 85 | 987 | 986 | 567 | 875 |
| Q2 | 678 | 654 | 987 | 659 | 786 | 436 | 980 | 876 | 908 |
| Q3 | 899 | 875 | 190 | 983 | 909 | 237 | 987 | 100 | 1089 |
| Q4 | 787 | 969 | 908 | 537 | 567 | 836 | 837 | 926 | 987 |

The above 3-D table can be represented as 3-D data cube as shown in the following figure:

# Data Mart

Data marts contain a subset of organization-wide data that is valuable to specific groups of people in an organization. In other words, a data mart contains only those data that is specific to a particular group. For example, the marketing data mart may contain only data related to items, customers, and sales. Data marts are confined to subjects.

## Points to Remember About Data Marts

- Windows-based or Unix/Linux-based servers are used to implement data marts. They are implemented on low-cost servers.

- The implementation cycle of a data mart is measured in short periods of time, i.e., in weeks rather than months or years.

- The life cycle of data marts may be complex in the long run, if their planning and design are not organization-wide.

- Data marts are small in size.

- Data marts are customized by department.

- The source of a data mart is departmentally structured data warehouse.

- Data marts are flexible.

16

tutorialspoint
SIMPLYEASYLEARNING

End of ebook preview
If you liked what you saw…
Buy it from our store @ **https://store.tutorialspoint.com**