

## 2. CONCEPTS

### 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.