

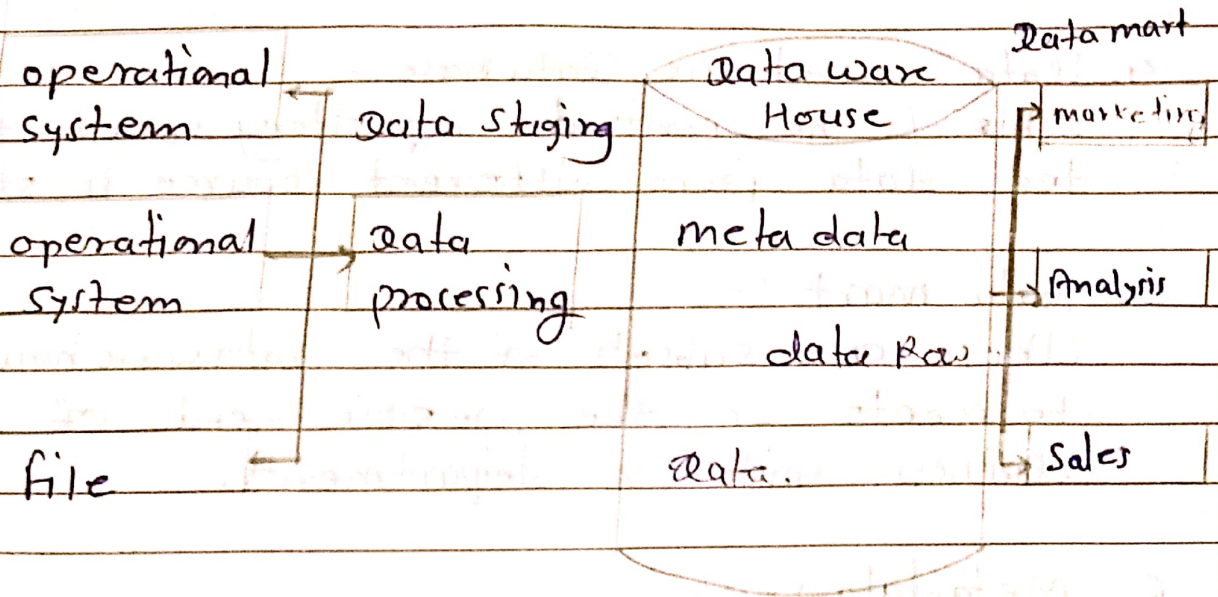
Assignment 2.

ENRCE No.	
DATE	/ /

Q1 Explain the main component of a typical Data warehouse architecture and their roles.

→

Data Source



Data warehouse

- A typical data warehouse architecture consists of several key components, each playing a crucial role in the data management and analytic process:

1. Source system:-

These are the operational database and external data source where raw data originates. They include CRM, ERP and Flat file.

2. ETL (Extract Transform Load) Tools:-

These tools extract data from source systems, transform it into a suitable format and load it into the data warehouse.

ETL Process handles data cleaning, integration, etc.

3. Data Staging Area :-

This is a temporary storage area where data is held before it is transformed and loaded into the data warehouse.

4. Data warehouse Database :-

This is the central repository where the integrated data from different source is stored.

5. Data mart :-

These are subsets of the data warehouse, tailored to meet the specific needs of different business units or departments.

6. Metadata :-

This stores metadata, which is data about data. It includes information about data source, transformations, structures and business rules.

7. OLAP (Online Analytical Process) Tools :-

These tools enable to perform multiple-dimensional analysis of the data stored in the data warehouse.

Q2. What is the difference between the data staging area and the data presentation area in a data warehouse architecture? How do they contribute to the overall functionality of the system.

Data staging area

- The data staging area is a temporary storage area where raw data from various sources is collected for querying.

- It serves as a workspace for ETL (Extract Transform efficient data and to access load) processes.

- Raw data is extracted from different source and stored in structured system, which can include format, often in a star database, flat files, APIs or snowflake schema.

- Data Transformed into a suitable format which can include aggregation, normalization, and organization. Data is a way that is understandable by business users often through dashboard reported and visualization.

- Contribution :-

- Data Staging area :-

Ensures data quality and consistency before it is loaded into the warehouse.

- Provides a controlled environment for complex data processing tasks without impacting operational system.

- Data presentation Area :- Facilitates easy and efficient access to data for decision-making process.

Data presentation area

- Enhances the performance of queries by providing optimized structures and indexing.

Q3 Describe the concept of data normalization and denormalization in the context of data warehouse architecture. What are the benefits and challenges of each approach?

→ * Data Normalization :-

- Normalization involves organizing data into tables and columns to reduce redundancy and improve data integrity. It follows a set of rules called normal forms.
- It reduces data redundancy, improved data integrity, easier data maintenance, and minimized update anomalies.
- Can lead to complex queries, may impact query performance due to the need for multiple table joins.

* Data Denormalization :-

- Denormalization involves combining tables to reduce the number of joins required for queries. It increases redundancy to optimize query performance.
- Improved query performance, simple queries, faster data retrieval for read-heavy operations.
- Increased data redundancy, potential for data anomalies, more complex data updates and maintenance.

Q4. In a Data warehouse architecture, what is the purpose of an OLAP? How does it differ from the OLTP (Online Transactional Processing) layer, and why is it important for analytical tasks?

→ OLAP (Online Analytical Processing) Layer :-

* Purpose :-

The OLAP layer is designed to support complex analytical queries and multi-dimensional analysis. It allows users to analyze data across multiple dimensions.

* Importance :-

Essential for decision-making and strategic planning, enables users to gain insights from large volumes of data, supports advanced analytical functions such as forecasting and trend analysis.

* Difference from OLTP :-

- Optimized for transactions processing, focuses on data integrity and consistency, handles a large number of short online transactions, supports day-to-day operations.
- Optimized for data analysis, focuses on query performance and data aggregation, handles complex queries and large volumes of data, supports business intelligence and decision making.

Q.5. Discuss the various data distribution strategies used in Data warehouse architecture. Compare and contrast hash-based distribution and range-based distribution, highlighting their respective advantage and potential use case.

→

Hash-based Distribution:-

Concept - Data is distributed across nodes based on a hash function applied to a distribution key. Each value of the key is assigned to a specific node.

Advantages:-

- Even distribution of data
- Prevents data skew.
- Efficient for equality-based queries.

* Range-based Distribution:-

Concepts - Data is distributed across nodes based on predefined ranges of the distribution key. Each range is assigned to a specific node.

Advantages:-

- Efficient for range-based queries
- Allow for localized queries processing
- Can optimize data storage for specific access patterns.

8/8/21