

31.2 Discuss how Online Transaction Processing (OLTP) systems differ from data warehousing systems.

Aspect	OLTP Systems	Data Warehousing Systems
Purpose	Day-to-day transactions	Analytical and reporting tasks
Data Type	Current, detailed data	Historical, summarized data
Schema Design	Normalized	Denormalized (star/snowflake)
Query Type	Simple, transactional	Complex, analytical
Users	Operational staff, end-users	Analysts, data scientists
Performance Focus	Speed, low latency	Query efficiency, read speed
Use Cases	Banking, e-commerce	BI platforms, financial analysis

31.3 Discuss the main benefits and problems associated with data warehousing.

Benefits of Data Warehousing

1. Enhanced Data Analysis and Business Intelligence (BI)

- Benefit: Data warehouses integrate data from multiple sources (e.g., CRM, ERP, transactional databases) into a single, unified system. This enables comprehensive analysis, supports business intelligence (BI) tools, and helps in generating valuable insights.
- Example: Companies can identify sales trends, customer preferences, and market opportunities more effectively.

2. Improved Data Quality and Consistency

- Benefit: Data warehousing involves processes like Extract, Transform, Load (ETL) that clean and standardize data before storage. This results in higher data quality, consistency, and reliability across the organization.
- Example: A unified view of customer data reduces inconsistencies (e.g., duplicate records), providing more accurate reporting.

3. Historical Data Storage

- Benefit: Data warehouses store historical data, making it possible to analyze trends over time. Businesses can track performance metrics and make data-driven predictions based on historical patterns.
- Example: Analyzing yearly sales data to forecast demand for the next fiscal year.

4. Enhanced Decision-Making

- Benefit: By providing a centralized repository of well-organized data, data warehouses help stakeholders make informed decisions quickly. They support complex queries and generate reports that guide strategic planning.
- Example: A marketing team can access reports to analyze the success of a campaign and decide on the next steps based on data insights.

5. Increased Query Performance

- Benefit: Data warehouses are optimized for complex, read-heavy queries, significantly improving query performance compared to operational databases (OLTP systems). They use indexing, partitioning, and parallel processing to handle large datasets efficiently.
- Example: Running a sales report across multiple years can be completed in seconds rather than hours.

6. Enhanced Security and Data Governance

- Benefit: Data warehouses provide robust security features, including access control, auditing, and encryption, ensuring that sensitive data is protected. They also facilitate compliance with data governance standards and regulations.
 - Example: Ensuring compliance with GDPR by providing controlled access to customer data.
-

Problems of Data Warehousing

1. High Implementation and Maintenance Costs

- Problem: Building and maintaining a data warehouse can be costly, involving hardware, software, and specialized personnel (e.g., data engineers, ETL developers). Ongoing maintenance costs can also be significant, especially with growing data volumes.
- Example: A small company may find it challenging to allocate the budget for the infrastructure and expertise required to build a robust data warehouse.

2. Complex Data Integration

- Problem: Integrating data from various sources (e.g., legacy systems, external APIs, different databases) can be complex and time-consuming. Inconsistent data formats, missing data, and discrepancies between sources add to the challenge.
- Example: Merging sales data from an old ERP system with customer data from a new CRM may require extensive data transformation.

3. Data Latency

- Problem: Data warehouses typically use batch processing for data updates, leading to some level of data latency. This means the data is not always up to date, which can be an issue for real-time analytics.
- Example: In a retail business, real-time stock updates may not be reflected immediately in the data warehouse, affecting inventory reports.

4. Scalability Challenges

- Problem: As data volumes grow, scaling a traditional on-premises data warehouse can be challenging. It requires additional storage, computing power, and changes to the ETL processes, which can be costly and complex.
- Example: An e-commerce platform experiencing rapid growth may struggle to scale its data warehouse infrastructure to handle the increasing volume of transaction data.

5. Complex Query Design and Maintenance

- Problem: Complex queries and reports can be difficult to design and may require advanced SQL skills and optimization techniques. Poorly designed queries can lead to slow performance and resource contention.
- Example: An inefficient query joining multiple large tables may take hours to run, slowing down the entire system and impacting other users.

6. Data Governance and Privacy Issues

- Problem: Centralizing data from various sources raises concerns about data privacy, governance, and compliance. Managing sensitive information and ensuring compliance with regulations like GDPR or HIPAA can be challenging.
- Example: A healthcare data warehouse must carefully manage patient data to comply with HIPAA regulations, requiring robust access controls and data masking.

31.11 *Discuss how data marts differ from data warehouses and discuss the main reasons for implementing a data mart.*

Aspect	Data Mart	Data Warehouse
Purpose	Provides targeted, subject-specific data	Provides a comprehensive view of all data
Scope	Department-specific (e.g., sales, HR)	Enterprise-wide (integrates all data)
Data Volume	Smaller, limited to a specific function	Large, covering the entire organization
Design	Simple, often denormalized	Complex, can be normalized or denormalized
Implementation Time	Faster (weeks to months)	Longer (months to years)
Complexity	Lower complexity	Higher complexity
User Base	Specific department or user group	Broad, organization-wide access
Maintenance	Easier to maintain	Requires more resources and effort

32.4 *Discuss the concepts associated with dimensionality modeling.*

In **dimensional modeling (DM)**, the structure consists of a main table called the **fact table**, which contains a composite primary key, and several smaller tables known as **dimension tables**. Each dimension table has a simple primary key, which corresponds directly to one of the elements in the composite key of the fact table. Essentially, the primary key of the fact table is made up of two or more foreign keys linked to these dimension tables. This design forms a distinctive "star-like" layout, known as a **star schema** or **star join**.

A key feature of dimensional modeling is the use of **surrogate keys** instead of natural keys. Surrogate keys are artificial keys, typically represented as simple integers, that replace the natural keys. As a result, all the joins between the fact and dimension tables are performed using these surrogate keys. This approach helps the data warehouse maintain a degree of independence from the data structures used in the underlying OLTP systems.

33.1 Discuss what Online Analytical Processing (OLAP) represents.

Online Analytical Processing (OLAP) refers to a category of technologies designed for querying and analyzing complex, multi-dimensional data efficiently. OLAP systems are used primarily in business intelligence and data warehousing to support decision-making processes by allowing users to perform advanced analytics, such as trend analysis, forecasting, and data exploration.

33.2 Discuss the relationship between data warehousing and OLAP.

In the past few years, relational DBMS vendors have targeted the data warehousing market and have promoted their systems as tools for building data warehouses. A data warehouse stores operational data and is expected to support a wide range of queries from the relatively simple to the highly complex. However, the ability to answer queries is dependent on the types of end-user access tools available for use on the data warehouse. General-purpose tools such as reporting and query tools can easily support ‘who?’ and ‘what?’ questions about past events. A typical query submitted directly to a data warehouse is: ‘What was the total revenue for Scotland in the third quarter of 2001?’. In this section we focus on a tool that can support more advanced queries, namely online analytical processing (OLAP).

While OLAP systems can easily answer ‘who?’ and ‘what?’ questions, it is their ability to answer, ‘what if?’ and ‘why?’ type questions that distinguish them from general-purpose query tools. OLAP enables decision-making about future actions.

34.1 Discuss what data mining represents.

Data mining is the process of discovering patterns, trends, and valuable information from large datasets using techniques from statistics, machine learning, and database systems. It is a key component of **data analysis** and is used in fields like business intelligence, market research, and predictive analytics to extract actionable insights.

Data mining helps organizations extract valuable knowledge from large datasets, enabling data-driven decision-making, improved predictions, and strategic insights. It combines statistical methods, algorithms, and machine learning techniques to uncover meaningful information that can lead to competitive advantages.

34.2 *Provide examples of data mining applications.*

Examples of Data Mining Applications (List)

1. Retail and E-commerce:

- Market Basket Analysis
- Customer Segmentation

2. Banking and Finance:

- Fraud Detection
- Credit Scoring and Risk Analysis

3. Healthcare:

- Disease Prediction and Diagnosis
- Healthcare Management

4. Telecommunications:

- Churn Analysis
- Network Optimization

5. Manufacturing and Supply Chain:

- Quality Control
- Demand Forecasting

6. Marketing and Advertising:

- Personalized Recommendations
- Campaign Effectiveness Analysis

7. Education:

- Student Performance Prediction
- Personalized Learning

8. Social Media and Web Analytics:

- Sentiment Analysis
- Trend Analysis

34.6 *Discuss the relationship between data warehousing and data mining.*

One of the major challenges for organizations seeking to exploit data mining is identifying suitable data to mine. Data mining requires a single, separate, clean, integrated, and self-consistent source of data. A data warehouse is well equipped for providing data for mining for the following reasons:

- Data quality and consistency is a prerequisite for mining to ensure the accuracy of the predictive models. Data warehouses are populated with clean, consistent data.
- It is advantageous to mine data from multiple sources to discover as many interrelationships as possible. Data warehouses contain data from several sources.
- Selecting the relevant subsets of records and fields for data mining requires the query capabilities of the data warehouse.
- The results of a data mining study are useful if there is some way to further investigate the uncovered patterns. Data warehouses provide the capability to go back to the data source.