

Research Assignment

Section A

① Relational Databases (RDBMS)

- Store data in tables (row and columns)
- Use SQL to manage and query data.

NOSQL Databases

- Store unstructured or semi-structured data

Cloud Databases

- Hosted and managed on cloud platforms
- Provide scalability, security and remote access

② A Relational Database Management is a type of database software that stores, manage and organize data in tables made up of rows and columns each table represents a specific type of data

③ A primary key is a unique identifier for each record (row) in a database table. It ensures that no two rows have the same value in the column. (No duplicates)

④ Database normalization is the process of organizing data in a database to reduce data redundancy (repetition) and improve data integrity.

⑤ A database schema is the structure of a database. It defines how data is organized including the tables. It basically the plan that shows what data exists, how it's stored and how different part of the database connect.

⑥ Structured Data - organized in fixed formats like tables with rows and columns easy to store and query using SQL
Example: Databases, excel sheet

Semi-structured Data - Has some structured but not in a strict table format uses tag or marker to separate data.
Example: CSV, email.

Unstructured Data - No predefined structure difficult to store in tables
Example: video, image and analyze directly

→ It is connected to dimension tables using foreign keys.

⑦ Fact Table = It stores measurable quantitative data

Example:

Record
each
Sale

Order ID

Product ID

Customer ID

Sale Amount

→ It contains numbers like sales amount, quantity or profit

→ Each record represents a specific event (e.g. one sale)

Dimension Table = It contains descriptive attributes (like customer

Example:

Name, product categories or dates]

contain Primary key [contain the details who, what, where, when]

⑧ A data model is a visual or logical way of showing how data is structured, stored and related in a database. It defines WHAT DATA will be stored (entities or tables)

How data relates, what rules apply (keys)

Importance

→ Helps organize data clearly before building the database

→ Makes it easier to understand relationships between data

→ Guides analysts when designing or querying the database

It helps plan and structure data properly so that database works efficiently and accurately

⑨ Database = store current, day-to-day operational data. It is mostly

structured data and is organized in tables [row and column]

Example: Banking system, online shopping apps

Data warehouse = store historical data for analysis and reporting. It is structured and cleaned data. It is organized and optimized for queries

Example: Business intelligence dashboards, trend analysis

Data lake = store all types of data (raw, unprocessed), it is structured semi-structured and unstructured. NO fixed structure, very flexible [Example: Big data analytics, machine learning]

⑩ A data mart is a smaller, focused part of a data warehouse that stores data for a specific business area or department.

Difference between Data Mart and Data Warehouse

Aspect	Data Warehouse	Data Mart
① Scope	cover the entire organization data	Focuses on a specific department or function [e.g. Sales, Finance]
② Data size	very large - store data from many sources	smaller - contain only relevant data for one team
③ Purpose	used for company-wide analysis and reporting	used for department-level analysis.
④ Data source	integrate data from multiple systems	usually gets data from the data warehouse

Section B SQL and Data Processing

- ⑪ A query language is a language used to communicate with a database to retrieve add, update, or delete data
 → Structured Query Language is commonly used because it's the standard language for most relational databases.
- ⑫ Indexes in databases are special data structures that make it faster to find and retrieve data from a table [Just like a book index helps you find information quickly trading a bit of space for much better speed.]
- ⑬ A transaction in a database is a single unit of work that groups one or more operations (like inserts, update or delete) so they either all happen successfully or none happen at all

ACID Properties

Atomicity - All parts of a transaction happen or none do if one step fails the whole transaction is rolled back

Consistency - The database remains valid before and after the transaction. Balance must still add up correctly after transfer

Isolation - Transaction don't interfere with each other

Durability - Once committed, data is permanently saved even after a crash. After a crash data won't be lost.

⑭ A database engine is like the brain of the database it manages how data is stored, processed and retrieved in a database. It's which execute SQL queries and handles reading and writing data.

→ The impact on performance is manage locking and memory use, reducing disk access. A good engine optimizes queries and uses index efficiently making searches faster.

⑮ Views = is a virtual table created from the results of an SQL query
- used to simplify complex queries, secure data and present customized results

→ Stored Procedure = A stored procedure is a predefined set of SQL statements save in the database
- It can take inputs, perform operations and return results

Example → CREATE PROCEDURE GetCustomerOrders (Q. CustomerID INT)
AS

SELECT * FROM orders WHERE customer_id = Q.CustomerID;

→ Triggers = a piece of SQL code that automatically runs when a specific event happens on a table like INSERT, UPDATE, or DELETE (automatically executes when data changes)

⑯ ETL = Transform before loading (clean first, then store)

ELT = load first then transform inside the warehouse

⑰ Batch processing = handles data in bulk, later it works with historical data
Stream processing = handles data continuously in real time it works with real-time

⑱ A Join in SQL is used to combine data from two or more tables based on a related column between them (1) Inner Join = Returns only the rows that match in both tables (2) Left Join = Returns all rows from the left table and matching rows from the right table (3) Right Join = returns all rows from the right table

(19) Referential integrity means that the relationship between tables in a relational database remain consistent and accurate.

→ If maybe there is a foreign key in a table it prevents and lists it as something that doesn't exist

(20) Data Redundancy means storing the same piece of data in multiple places within a database. e.g.: if a customer's address appears in both the customer table and the order table that's redundancy. The same data is repeated unnecessarily.

→ Storage = It repeated data consume extra disk space especially in large database with million of records

→ Performance = More data mean slower queries, inserts, updates, and backups. The database engine has to process more rows

Section C:

(21) Cloud Database is hosted on remote servers managed by a cloud provider (e.g. AWS, Google Cloud). The cloud provider handles hardware, backups, scaling and updates.

→ On-Premise Database is installed and maintained on organization own physical servers. The organization IT team is responsible for setup.

(22) Data governance is the process of managing the availability, usability, integrity and security of data within an organization. It is important because it maintains data quality, supports compliance, improve decision-making and reduce risk related to data misuse or errors.

(23) Data integrity means the accuracy, consistency and reliability of data throughout its entire lifecycle from when it's created to when it's used or deleted. It can be maintained by using validation rule, primary and foreign keys, access control, backup and audit trail to ensure that data remain reliable and unchanged.

- (24) Data quality refers to how accurate, complete, consistent and reliable data is for its intended purpose. It is critical for analytics because reliable data leads to accurate insights better decision and greater trust in analytical results.
- (25) A Data Analyst manages and analyzes data from databases by collecting, cleaning and interpreting information to find trends and supports decision-making.
- (26) A Database Administrator is responsible for managing and maintaining databases. Their key duties include installing and configuring database systems, ensuring data security, performing backups and recovery.
- (27) The main steps in designing a data pipeline include identifying data source, defining data requirements, extracting data, transforming it for consistency and accuracy, loading it into the target system, automating the workflow and continuously monitoring and maintaining the pipeline.
- (28) Common challenges in managing large-scale databases include maintaining performance and scalability, ensuring data security and integrity, managing backups and storage, handling concurrency and integrating data from multiple sources.
- (29) MySQL = open-source relational database widely used for web applications (Wordpress)
PostgreSQL = Advanced open-source relational database that supports complex queries
Snowflake = cloud-based data warehouse built for big data and analytics.
Handles structured and semi-structured data with scalability and high performance
- (30) CSV (comma-separated value) = simple text format for tabular data easy to use and read.
Parquet = columnar format optimized for analytics and large datasets
JSON (JavaScript Object Notation) flexible text format for semi-structured data widely in API and web applications.

Unstructured data

Example: videos, images

and analyze directly