

Research

Assignment

SECTION A: DATABASE Fundamentals

1. There are various kinds of databases:
 - * Relational database - which stores data in tables with rows and columns (e.g. MS SQL, PostgreSQL)
 - * Non-relational database are designed for flexibility and scale
 - Data warehouse - for analytical queries and BI (Snowflake, Redshift)
 - * In memory database: Stores data in RAM for fast access
 - * Time-series databases: optimize storage for time-stamped data.
 - * Object-oriented database: Stores data as objects

2. A relational database management system is → a software that stores & manages data using tables with relationships defined by keys. It supports SQL queries, and structured tabular data e.g. MS SQL Server

3. A primary key: is a unique identifier for each record in the table. It cannot be null.

A foreign key is a field that links one table to another by referencing the primary key of the related table.

⑧ A data model is a blueprint that shows how data is structured and related

- ↳ It guides a database design
- ↳ Supports efficient queries &
- ↳ Ensures consistency.

⑨ Database: Supports daily operations / or transactions (e.g. OLTP) 11
Maris scanning at checkpoint, all records will record on an OLT & that is a database. A relational database

Data warehouse: Is designed for analytic processes (An OLAP).

- ↳ Stores current & historical data & data is stored & refreshed from source system
- ↳ It stores cleared and structured data (Snowflake)

Data lake: Stores raw, unstructured, semi-structured and structured for large scale analytics

⑩ A data mart is a smaller, (maybe) department-based / department-specific subset of a data warehouse (e.g. HR Mart in a company database).
• It differs from a database in that a database is an enterprise wide analytical storage serving multiple functions

Section B : SQL and Data Processing

① A query language is the way in which users retrieve and manipulate data in SQL - This will be syntax (how we communicate with the database) to make sense of it.

SQL is the most common query language for structured data.
It is standardized
& works on most RDBMS

② Indexes in a database are data structures (like pointers) that speed up search queries by avoiding full table scans.
They improve read performance

③ A transaction in a database is a unit of work executed as a whole (e.g. insert, update, delete)

ACID Properties

↳ A set of guarantees that ensures transactions are reliable. e.g.

ACID has the following elements

- 1- Atomicity (All or nothing)
- 2- Consistency (Valid before & after)
- 3- Isolation (No interference between transaction)
- 4- Durability (Permanent results after commit)

14. A database engine is the underlying software responsible for storing, indexing and retrieving data. It affects

- ↳ Speed
- ↳ Storage method
- ↳ Concurrency control
- ↳ Query optimization

15. A view : A virtual table based on a query
Stored procedure : A saved SQL script executed on demand

Trigger : A procedure executed when certain events occur (e.g. insert, delete)

16. The difference between ETL & ELT

This is a pipeline through which data gets taken from OLTP to OLAP

ETL : Extract - Transform - Load.

ETL : Transformation is done before data is loaded. Most used in traditional warehouses

ELT : Extract - Load - Transform.
: Transformation of data is done inside the data warehouse, after data is loaded. Most used with Cloud systems like BigQuery & Snowflake.

17) Batch processing: Processes large chunks of data at intervals - (nearly, 1 daily)

Stream processing: (processes data in real time) (This could be in an OLTP)

18) A Join in SQL is a way to combine rows from two or more tables based on related columns. There are various types of Joins

↳ Inner Join : matches all common records

↳ Left Join : All from left table & only matched from right

↳ Right Join : All from right & only matched from left.

↳ Full Join : All records when match exists

↳ Cross Join : Cartesian product - (meaning all rows possible combination of rows)

19) Referential integrity : ensures relationships between tables remain valid by enforcing rules via foreign keys

Importance : L prevents orphan values

L Ensures consistency

L Maintains data accuracy

20 Redundancy causes increased storage costs, slower query performance and causes more higher chance of inconsistencies and complex updates.

Section C: Data Management and Analytics Concepts.

21 Cloud database management: is scalable on demand, accessible anywhere and its managed services with lower upfront costs

On-premises database management: has/ offers full control, high setup/ maintenance costs and hard ware dependency.

22 Data governance: A framework of policies and processes ensuring the correct management of data.

- ↳ It ensures compliance.
- ↳ Protects data quality.
- ↳ Manages risk &
- ↳ Defines roles & responsibility in terms of accountability.

23 Data integrity refers to accuracy, consistency, and reliability of data. It is maintained through validation rules, keys, normalizer and

24 Data quality

- ↳ means data is accurate, complete, consistent, timely and valid. This is critical because poor data affects decision making, providing misleading insights and impacts reporting accuracy.

25 Role of Data Analysts

- ↳ Extract & clean data
- ↳ Perform SQL queries
- ↳ Build dashboards that provide visualisation & Storytelling of data for insights for decision making
- ↳ Analyze trends
- ↳ Supports decision making with insights

26 A database Administrator (DBA) -

- ↳ installs & maintains databases
- ↳ performs tuning
- ↳ Back up + recovery
- ↳ Access control of database
- ↳ Security management
- ↳ Ensures data integrity / assurance

27 Steps in designing a data pipeline

- ① requirements gathering
- ② source identification
- ③ data extraction
- ④ data transformation / cleaning

- ⑤ Load into Storage (warehouse/lake)
- ⑥ Scheduling
- G Monitoring & Maintenance

28

There are various challenges that could present in managing large scale databases including:

- l Performance Bottlenecks
- l Storage Scalability
- l Data quality issues
- l Concurrency conflicts
- l Backup complexity
- l Security threats
- l Cost management (Cloud compute)

29

Popular databases are amongst others

- ↳ MSSQL - web applications, OLTP systems
- l Snowflake - Advanced analytics of big data, cloud warehousing
- l Oracle - Large enterprise transaction systems

30

Main types of data storage formats used in analytics include-

- CSV (Simple text format) & it is the mostly used
- JSON - semi-structured, API's & NoSQL