

Research Assignment 1

* Section A

① The main types of databases are structured and unstructured data

② A Relational Database Management System is a type of database management system that stores and manages data in a structured and organized way, using tables, rows and columns.

NB It is based on the relational model, which means that data is stored in tables with well-defined relationships between them

③ A Primary key is a column or set of columns in a table that uniquely identifies each row in the table. It is a unique identifier for each record, and it cannot be NULL or duplicated.

A foreign key is a column or set of columns in a table that references the primary key of another table. It is used to establish relationships between tables and ensure data consistency

④ Data normalization is the process of organising data in a database to minimize data redundancy and dependency

It is important because it improves data quality, reduces data anomalies, supports data

security, improves database performance

⑤ A Schema is the overall structure or organization of a database, including the relationships between different tables, columns, and data types.

⑥ The difference between these 3 types of databases is that:

* Structured data

- is highly organized
- Fits into a schema
- is easily searchable

e.g. Relational database (customer information), financial transactions, Spreadsheets and CSV files

* Semi structured data

- Partially organized
- Flexible schema
- Self-describing
- e.g. JSON (JavaScript Object Notation) XML (Extensible Markup language), CSV files with variable columns and key-value stores.

* Unstructured data

- lacks organization
- No schema
- Human readable
- e.g. Text documents (e.g emails, articles, social media posts, images, videos, audio files and sensor data)

④ The difference between a fact table and a dimension table is that a fact table measures business performance, it is also central to the datawarehouse, it is as well typically large. for example sales, orders, website interactions, customer transactions.
While a dimension ~~table~~ is a table that provides context, describes attributes and is typically smaller. for example date, time, customer, product, location, product category

⑤ A data model is a conceptual representation of the structure and relationships of data in a database. It's a blueprint or a map that defines ~~the~~ how data is organised, stored, and related

It is important because it:

- improves data integrity (ensures data consistency and accuracy)
- supports scalability (enable easy modification and extension of the database)
- facilitates communication
- Reduces errors.

⑥ The difference between the 3 is that a database:

- stores structured data
- Optimized for transactions
- Schema driven
- Eg MySQL, PostgreSQL, Microsoft SQL server.

Data warehouse:

- stores historical data
- optimized for analytics
- schema-on-read (datawarehouse often uses a schema-on-read approach, where the schema is applied when the data is queried.)

E.g. Amazon Redshift, Google BigQuery
Snowflake.

Data lake

- stores raw, unprocessed data like videos, texts and images.
 - schema-on-write
 - flexible and scalable
- e.g. Hadoop Distributed file System (HDFS) Amazon S3, Azure Data Lake Storage

⑩ A data mart is a storage solution which:

- subset of data (is a subset of a datawarehouse that focuses on a specific business area or department)
- departmental
- contains simplified data, making it easier for users to access and analyze

The key differences are:

- scope (data warehouse (enterprise wide))

vs Data Mart (departmental or functional))

- Data complexity
(data warehouse (complex data) vs Data Mart (simplified data))
- purpose (Data warehouse (strategic enterprise-wide analytics) vs Data Mart (tactical departmental analytics))
- Design (Data warehouses (designed for flexibility and scalability) ~~versus~~ Data Mart (designed for fast query performance and ease of use))

Section B (SQL and Data Processing)

⑪ A query language is a programming language designed to manage and manipulate data stored in a database management system (DBMS). It's used to retrieve, update, and manipulate data in a database.

SQL is a commonly used because it offers

- Standardization
- Portability
- Flexibility
- Easy to use
- Industry adoption
- Data security
- Improved data management
- Data analysis

(12) Indexes in databases are data structures that improve the speed of data retrieval operations, such as SELECT, INSERT, UPDATE and DELETE. An index is a data structure that facilitates quick lookup and access to data in a table.

They improve performance by:

- faster query performance (reduces number of rows to be scanned).
- Improved data retrieval (quickly locate specific data, reducing the time it takes to retrieve data)
- Reduced I/O operations (minimize the number of disk I/O operations, improving overall system performance)

(13) Transactions in databases are sequences of operations that are executed as a single, all-or-nothing unit of work. They ensure data consistency and integrity by guaranteeing that either all changes are committed or none are, maintaining the database in a consistent state.

ACID (Atomicity, Consistency, Isolation, Durability) properties are a set of standards that ensure database transactions are processed reliably.

(14) A database engine, also known as a database management system (DBMS), is a software component that manages and interacts with a database. It is responsible for storing, retrieving, and manipulating data in a database.

It impacts performance by:

- Query performance
- Data retrieval speed
- Concurrency
- Scalability

(15) Views are virtual tables, dynamic data, simplified queries.

- * Stored procedures are precompiled SQL, modularity, Improve performance.
- * Triggers are automated actions that occurs in response to specific specific database events data integrity and auditing.

(16) key differences

- * Transformation location: ETL transforms data before loading, while ELT transforms data after loading
- * Data processing: ETL processes data in a separate environment, while ELT leverages the target system's processing capabilities
- * Data storage: ETL typically requires a staging area for transformation, while ELT stores raw data in the target system

- (7) *Batch Processing is a type of an approach to processing data, it includes:
- Batching data
 - Scheduled processing
 - High - throughput while
- * Stream processing includes
- real time processing
 - continuous processing
 - low latency

- (8) A join is a clause used to combine rows from two or more tables based on a related column between them.

Types of joins

Inner join

Left join

full outer join

Cross join

Self join

- (9) Referential integrity is a database concept that ensures the relationships between tables in a relational database are consistent and valid.

It is important to databases because

(1) Data consistency

- ② Data accuracy
- ③ Prevents orphaned records
- ④ Improves data quality

⑩ Data redundancy affects database performance by

- increased storage requirements
- Data inconsistency (redundant data can lead to inconsistencies if not properly synchronized, causing errors and inaccuracies.)
- Slower query performance
- Update anomalies.

Data redundancy affects database storage by

- Increased storage capacity
- Data duplication
- Data fragmentation.

* Section C

⑪ Cloud and on-site (^(premise)) database management differs in the following ways:

- Location: cloud databases are hosted on cloud servers and accessed via the internet, while on-premise databases are hosted on local servers within an organization's physical premise.
- Scalability: cloud databases offer easy scalability, allowing businesses to quickly adjust resources to meet changing demands. On-site (premise) databases require manual upgrades and hardware purchases to scale.

- Cost : cloud databases operate on a pay-as-you-go model, reducing upfront costs. On-premise databases require significant upfront investment in hardware, software, and IT staff.
- Accessibility : cloud databases enable remote access and global collaboration, whereas on-premise databases are typically limited to local networks or VPNs.
- Maintenance : cloud providers manage maintenance and updates freeing up IT resources. On-premise databases require in-house IT teams for maintenance, backups, and updates.

② Data Governance is the process of establishing and enforcing policies, procedures and standards for the effective and efficient use of data across an organization. It ensures that data is accurate, complete, consistent, and secure.

IT is important because it :

- Improve data quality
- Reduced Risk
- Increased Efficiency
- Better Decision making
- Regulatory Compliance

(23) Data integrity refers to the accuracy, completeness, and consistency of data throughout its entire life cycle. It ensures that data is reliable, trustworthy and consistent across different systems and applications.

It can be maintained by :

- Data validation
- Data normalization
- Data Quality Checks
- Access Controls
- Backup and Recovery
- Data Standardization
- Data governance

(24) Data quality is the accuracy, completeness, consistency and reliability of data. High-quality data is fit for its intended use and meets the requirements of the organization.

It is critical for analytics because it helps :

- with informed decision making
- with accurate insights
- with business performance
- with regulatory compliance
- with operational efficiency

25) A Data Analyst's role in database management and analysis is helping organizations make informed decisions via:

- Data collection
- Data cleaning
- Data organization
- Data modeling
- Data visualization
- Insight generation
- Reporting
- Storytelling
- Recommendations

26) Key responsibilities of a DBA are:

- Database design
- Database performance
- Data security
- Data back up and recovery
- Database maintenance
- trouble shooting
- Capacity planning
- Data Migration
- Documentation
- Collaboration

27) The main steps involved in designing a pipeline are:

- ① Define requirements
- ② Data ingestion
- ③ Data processing
- ④ Data storage
- ⑤ Data output

* General pipeline steps

- ① Define pipeline purpose
- ② Determine pipeline requirements
- ③ Design pipeline architecture
- ④ Choose tools and technologies
- ⑤ Develop pipeline code
- ⑥ Test and validate
- ⑦ Deploy and Monitor

(28) Common challenges in Managing large scale data :

- ① Scalability - handling increasing amounts of data and user traffic.
- ② Performance - ensuring optimal database performance, query optimization and indexing.
- ③ Data security - Protecting sensitive data from unauthorized access, breaches, cyber threats.
- ④ Data integrity - ensuring data accuracy, completeness, and consistency across the database.
- ⑤ Data backup and recovery - ensuring data availability and recoverability in case of failures or disasters.

- ⑥ Data Management - managing large amounts of data, including data storage, archiving, and purging
- ⑦ Query Optimization - optimizing complex queries to improve performance and reduce latency.
- ⑧ Concurrency Control - Managing multiple users and transactions to prevent data inconsistencies
- ⑨ Data Integration - integrating data from multiple sources, formats and systems
- ⑩ Cost Management - Managing costs associated with database infrastructure, maintenance and personnel

29 Popular Database Platforms

- ① MySQL → web applications, e-commerce platforms, and content-management systems.
- ② PostgreSQL → Enterprise applications, data warehousing and Analytics
- ③ Microsoft SQL → Enterprise applications, business intelligence, data analytics
- ④ Amazon Aurora → Cloud native relational databases for MySQL, SQL Server, PostgreSQL
- ⑤ Snowflake
- ⑥ Oracle

(use cases)

① Web Applications → Relational databases

like MySQL and PostgreSQL

For web applications

② Big data Analytics → NoSQL databases

like MongoDB and Cassandra

for big data storage and analytics

③ Real-time data processing → In-memory

databases like Redis for real-

time data processing and

cacheing.

④ IoT IOT Data storage → time series

databases like InfluxDB for

IoT data storage and analytics

⑤ Artificial Intelligence and Machine learning →

Graph databases like Neo4j for

knowledge graphs and recommendation

engines

⑥ E-commerce platforms → Relational databases

like MySQL for e-commerce platforms

⑦ Content Management Systems → relational

databases like MySQL and

PostgreSQL for content management systems

30 Main data storage formats used in Analytics

① ~~structured~~ Structured data formats

- CSV (comma separated values) :

plain text format for tabular data

- JSON (JavaScript Object Notation) - a light weight, human-readable format for data interchange
- Avro - a binary format for big data storage and processing

② Columnar Storage formats

- Parquet - a columnar storage format for big data analytics
- ORC - ~~(Opt)~~ a columnar storage format for hadoop

③ Specialised data formats

- time-series data formats - formats like influx DB's line protocol and TSDB's data format

④ Semi structured data format

- XML - a markup language for storing and transporting data

⑤ Unstructured data formats.

- txt files.