# Databases

By : Lakshmikant Deshpande

# Databases

- Databases are organized collections of data that are crucial for storing, managing, and retrieving information in various applications and systems.
- Databases are structured repositories for storing, managing, and organizing data.
- A Database Management System (DBMS) acts as the software layer that facilitates efficient data storage, retrieval, management, and security. It plays a crucial role in ensuring data reliability and accessibility for users and applications.

**Main Types**

- **Relational Databases**: Store data in tables with rows and columns, using SQL for querying (e.g., MySQL, PostgreSQL, Oracle).
- **NoSQL Databases**: Flexible, non-tabular data storage, suitable for unstructured or semi-structured data (e.g., MongoDB, Cassandra, Redis).
- **Graph Databases**: Designed for managing data with complex relationships (e.g., Neo4j).
- **In-memory Databases**: Store data in main memory for fast access (e.g., Redis, Memcached).

# Databases contd…

- **Components**:
    - **Tables**: Organized data structures in relational databases.
    - **Rows**: Individual records within tables.
    - **Columns**: Define the type of data stored in a table.
    - **Indexes**: Speed up data retrieval by creating efficient access paths.
    - **Queries**: Commands used to retrieve, update, or manipulate data.
- **ACID Properties**:
    - Databases ensure data integrity using ACID (Atomicity, Consistency, Isolation, Durability) properties.
- **Primary Keys**:
    - Unique identifiers for each row in a table to ensure data integrity and establish relationships between tables.

# Databases contd…

- **Normalization**:
    - Process of organizing data in a database to minimize data redundancy and improve data integrity.
- **Transactions**:
    - Sets of database operations treated as a single unit, ensuring data consistency.
- **Backup and Recovery**:
    - Regularly backing up data to prevent loss and implementing recovery procedures.
- **Scalability**:
    - Databases can scale vertically (adding more resources to a single server) or horizontally (distributing data across multiple servers) to handle increased loads.
- **Security**:
    - Access control, encryption, and auditing are essential for protecting data in databases.

# Databases contd…

- **Data Warehousing**:
  - Storing and managing large volumes of historical data for analysis and reporting purposes.
- **Cloud Databases**:
  - Databases hosted in the cloud, offering scalability, flexibility, and ease of management (e.g., Amazon RDS, Azure SQL Database).
- **Big Data Databases**:
  - Designed to handle massive volumes of data, often with distributed architectures (e.g., Hadoop, Apache Cassandra).
- **NoSQL Categories**:
  - Document-oriented, key-value, column-family, and graph databases are common types in the NoSQL ecosystem.
- **Query Languages**:
  - SQL (Structured Query Language) is prevalent for querying relational databases, while NoSQL databases may use their query languages.

# Databases contd…

- **Data Models**:
    - Relational databases use a tabular model, while NoSQL databases can adopt various data models like document-based, key-value, or graph.
- **Data Consistency Models**:
    - Different databases may implement various consistency models, such as strong consistency, eventual consistency, or causal consistency.
- **Database Management Systems (DBMS)**:
    - Software systems that provide tools and interfaces for interacting with databases.
- **Data Migration and ETL**:
    - Extract, Transform, Load (ETL) processes move data between databases and systems, often used in data integration and migration projects.
- **Data Warehousing**:
    - Dedicated databases for storing and analyzing historical data for business intelligence and reporting purposes.

# Components

- **Tables**:
    - Relational databases use tables to organize data.
    - Each table consists of rows (records) and columns (fields), and each column has a specific data type.
- **Rows**:
    - Rows represent individual records or entries within a table.
    - Each row contains data values corresponding to the columns.
- **Columns**:
    - Columns define the type of data that can be stored in a table.
    - Each column has a name and a data type (e.g., integer, text, date).
- **Indexes**:
    - Indexes are data structures that improve data retrieval efficiency.
    - They provide a quick way to locate specific rows based on indexed columns.
- **Queries**:
    - Queries are commands or statements used to retrieve, update, or manipulate data within a database.
    - SQL is the primary language for querying relational databases

# ACID Properties

- ACID stands for Atomicity, Consistency, Isolation, and Durability. These properties ensure the reliability and integrity of data in a database.
- **Atomicity**:
  a. Atomicity ensures that a transaction is treated as a single, indivisible unit.
  b. Either all changes within a transaction are committed, or none are.
- **Consistency**:
  a. Consistency ensures that a database transitions from one valid state to another.
  b. It enforces data integrity constraints defined in the schema.
- **Isolation**:
  a. Isolation ensures that transactions do not interfere with each other.
  b. Concurrent transactions are isolated from each other to prevent data inconsistencies.
- **Durability**:
  a. Durability guarantees that once a transaction is committed, its effects are permanent and survive system failures.

# Primary Keys and Foreign Key

- **Primary Keys**:
  - Primary keys are unique identifiers for each row in a table.
  - They ensure that each record can be uniquely identified.
  - Primary keys are often used to establish relationships between tables.
- **Foreign Keys**:
  - Foreign keys are references to primary keys in other tables.
  - They establish relationships between tables and maintain referential integrity.
  - Foreign keys ensure that data consistency is maintained when related records are updated or deleted.

# Normalization

- Normalization is the process of organizing data in a database to reduce data redundancy and improve data integrity.
- It involves dividing a database into smaller, related tables and defining relationships between them.
- Normalization techniques, such as First Normal Form (1NF) and Third Normal Form (3NF), help eliminate data anomalies.
- Normalization is used to minimize the redundancy from a relation or set of relations.
- It is also used to eliminate undesirable characteristics like Insertion, Update, and Deletion Anomalies.
- Normalization divides the larger table into smaller ones and establishes relationships between them.

# Normal Forms & Types

- **1NF** : A relation is in 1NF if it contains an atomic value.

- **2NF** : A relation will be in 2NF if it is in 1NF and all non-key attributes are fully functional dependent on the primary key.

- **3NF** : A relation will be in 3NF if it is in 2NF and no transition dependency exists.

- **BCNF**: A stronger definition of 3NF is known as Boyce Codd's normal form.

- **4NF** : A relation will be in 4NF if it is in Boyce Codd's normal form and has no multi-valued dependency.

- **5NF** : A relation is in 5NF. If it is in 4NF and does not contain any join dependency, joining should be lossless.

# Advantages of Normalization

- Normalization helps to minimize data redundancy.

- Greater overall database organization.

- Data consistency within the database.

- Much more flexible database design.

- Enforces the concept of relational integrity.

# Disadvantages of Normalization

- You cannot start building the database before knowing what the user needs.

- The performance degrades when normalizing the relations to higher normal forms, i.e., 4NF, 5NF.

- It is very time-consuming and difficult to normalize relations of a higher degree.

- Careless decomposition may lead to a bad database design, leading to serious problems.

# 1NF

- A relation will be 1NF if it contains an atomic value.

- It states that an attribute of a table cannot hold multiple values. It must hold only single-valued attribute.

- First normal form disallows the multi-valued attribute, composite attribute, and their combinations.

# 2NF

- In the 2NF, relational must be in 1NF.
- In the second normal form, all non-key attributes are fully functional dependent on the primary key

# 3NF

- A relation will be in 3NF if it is in 2NF and not contain any transitive partial dependency.
- 3NF is used to reduce the data duplication. It is also used to achieve the data integrity.
- If there is no transitive dependency for non-prime attributes, then the relation must be in third normal form.

# Boyce Codd normal form (BCNF)

- BCNF is the advance version of 3NF. It is stricter than 3NF.

- A table is in BCNF if every functional dependency X → Y, X is the super key of the table.

- For BCNF, the table should be in 3NF, and for every FD, LHS is super key.

# 4NF

- A relation will be 1NF if it contains an atomic value.

- It states that an attribute of a table cannot hold multiple values. It must hold only single-valued attribute.

- First normal form disallows the multi-valued attribute, composite attribute, and their combinations.

# 5NF

- A relation is in 5NF if it is in 4NF and not contains any join dependency and joining should be lossless.
- 5NF is satisfied when all the tables are broken into as many tables as possible in order to avoid redundancy.
- 5NF is also known as Project-join normal form (PJ/NF).

# Transactions

- Transactions are sets of database operations treated as a single unit of work.

- They follow the ACID properties to ensure data consistency and integrity.

- Examples of transactional operations include transferring money between bank accounts or updating inventory levels in an e-commerce system.

# Backup and Recovery

- Regularly backing up data is crucial to prevent data loss due to hardware failures, human errors, or disasters.
- Database systems provide mechanisms for creating backups and restoring data in case of failure.
- Backup strategies may include full backups, incremental backups, and point-in-time recovery.

# Scalability

- Databases need to scale to handle increased loads and growing data volumes.
- Scalability can be achieved by either vertically scaling (adding more resources to a single server) or horizontally scaling (distributing data across multiple servers or nodes).
- Horizontal scaling is often associated with NoSQL databases.

# Security

- Database security involves access control, encryption, and auditing.

- Access control mechanisms restrict who can access and manipulate data.

- Encryption protects data at rest and in transit.

- Auditing tracks database activities for security and compliance purposes.

# Data Warehousing

- Data warehousing involves the storage and management of large volumes of historical data.

- It is used for business intelligence and reporting, allowing organizations to analyze trends and make data-driven decisions.

# Cloud Databases

- Cloud databases are hosted in the cloud, offering scalability, flexibility, and ease of management.
- Cloud providers like AWS, Azure, and Google Cloud offer database services that can be easily provisioned and scaled as needed.

# Big Data Databases

- Big Data databases are designed to handle massive volumes of data, often with distributed architectures.
- They are commonly used in data analytics and processing, as well as machine learning.
- Examples include Hadoop (for distributed storage and processing) and Apache Cassandra (for distributed NoSQL storage).

# Query Languages

- SQL (Structured Query Language) is the standard language for querying relational databases.
- NoSQL databases often have their query languages tailored to their data models (e.g., MongoDB's query language).

# Data Models

- Relational databases use a tabular data model where data is organized into tables with rows and columns.
- NoSQL databases can adopt various data models, including document-based, key-value, column-family, and graph.

# Data Consistency Models

- Different databases may implement various consistency models, such as strong consistency (immediate and strict consistency), eventual consistency (eventual agreement among distributed nodes), or causal consistency (based on causal relationships between operations).

# Database Management Systems (DBMS)

- DBMS software provides tools and interfaces for creating, managing, and interacting with databases.

- Popular DBMS include MySQL, Oracle Database, Microsoft SQL Server, and MongoDB.

# Data Migration and ETL

- Data migration involves moving data between databases or systems.

- ETL (Extract, Transform, Load) processes are used to extract data from source systems, transform it into the desired format, and load it into a target database or data warehouse.

# DDL, DML, DCL, and TCL

DDL, DML, DCL, and TCL are four categories of SQL (Structured Query Language) commands that are used for various database management tasks

**DDL (Data Definition Language):** DDL defines the structure of a database, enabling the creation, alteration, and deletion of database objects like tables and indexes.

**DML (Data Manipulation Language):** DML is used for manipulating data within a database, allowing for the retrieval, insertion, updating, and deletion of records in tables.

**DCL (Data Control Language):** DCL commands control access and permissions in a database by granting or revoking privileges to users, ensuring data security.

**TCL (Transaction Control Language):** TCL manages transactions, enabling the control of their execution, rollback, and commit, ensuring data integrity and consistency.

## DDL - Data Definition Language:

- DDL commands are used to define the structure of a database, including tables, indexes, constraints, and other database objects.
- Key DDL commands include:
    - CREATE: Used to create database objects such as tables, indexes, and views.
    - ALTER: Used to modify the structure of existing database objects, e.g., adding or dropping columns.
    - DROP: Used to delete database objects, including tables and indexes.
    - TRUNCATE: Removes all data from a table but retains the table structure.
    - COMMENT: Adds comments or descriptions to database objects.
- DDL commands do not manipulate data but rather define the database schema.

# DML - Data Manipulation Language

- DML commands are used to manipulate data stored in the database. These commands include:
  - SELECT: Retrieves data from one or more tables.
  - INSERT: Adds new records (rows) to a table.
  - UPDATE: Modifies existing records in a table.
  - DELETE: Removes records from a table.
- DML commands are used for querying, inserting, updating, and deleting data.

# DCL - Data Control Language

- DCL commands are used to control access and permissions in a database. These commands include:
  - GRANT: Assigns specific privileges to database users or roles, allowing them to perform actions on database objects.
  - REVOKE: Removes previously granted privileges, restricting user access.
- DCL commands are essential for managing database security and access control.

# TCL - Transaction Control Language

- TCL commands are used to manage transactions, which are sequences of one or more SQL statements treated as a single unit of work.
- Key TCL commands include:
    - COMMIT: Saves all changes made during a transaction and makes them permanent.
    - ROLLBACK: Undoes all changes made during a transaction, restoring the database to its previous state.
    - SAVEPOINT: Creates a point within a transaction to which you can later roll back.
    - SET TRANSACTION: Sets transaction-specific properties, such as isolation level and transaction name.
- TCL commands help maintain data consistency and integrity, especially in multi-user database environments.
- In summary, DDL is used for defining the database structure, DML is for data manipulation, DCL is for controlling access and permissions, and TCL is for managing transactions. Understanding and using these categories of SQL commands is fundamental to effective database management and application development.

# Subqueries

- **Definition**: A subquery, also known as a nested query or inner query, is a SQL query nested within another query. It is used to retrieve data that will be used as a condition for the main query.
- **Types**: Subqueries can be categorized into two main types: scalar subqueries (returning a single value) and table subqueries (returning a result set).
- **Use Cases**:
    - Filtering: Subqueries are often used to filter results based on a condition from another table or subquery.
    - Aggregation: Subqueries can provide aggregated values, such as counts, sums, or averages.
    - Comparison: You can compare values from the main query with those in a subquery.
    - IN/NOT IN: Subqueries can be used with the IN and NOT IN operators to filter data based on a list of values.
- **Syntax**: The general syntax of a subquery depends on its purpose, but it typically appears within parentheses and can be used in various clauses like WHERE, FROM, HAVING, or SELECT.
- **Performance**: Subqueries can impact performance, especially if they are correlated or involve large datasets. Careful optimization may be necessary.

# Subqueries Example

```sql
-- Create the employees table with Indian names for illustration
CREATE TABLE employees (
    employee_id INT AUTO_INCREMENT PRIMARY KEY,
    first_name VARCHAR(50),
    last_name VARCHAR(50),
    department VARCHAR(50),
    salary DECIMAL(10, 2)
);
-- Insert sample data with Indian names
INSERT INTO employees (first_name, last_name, department, salary) VALUES
('Amit', 'Sharma', 'HR', 55000.00),
('Priya', 'Verma', 'Finance', 60000.00),
('Rajesh', 'Yadav', 'IT', 70000.00),
('Sneha', 'Patel', 'IT', 75000.00),
('Kiran', 'Reddy', 'HR', 58000.00);
-- Find employees with the highest salary in each department
SELECT department, first_name, last_name, salary
FROM employees e
WHERE salary = (
    SELECT MAX(salary)
    FROM employees
    WHERE department = e.department
);
```

# Correlated Subqueries

- **Definition**: A correlated subquery is a subquery that references one or more columns from the outer query. It's executed once for each row processed by the outer query.
- **Relationship with Outer Query**: Correlated subqueries establish a connection with the outer query by using values from the outer query's current row in their conditions.
- **Use Cases**:
  - Data Validation: Correlated subqueries are useful for validating data based on related information in another table.
  - Aggregation by Group: They can be used to calculate aggregates per group in the outer query.
- **Syntax**: The syntax of a correlated subquery is similar to a regular subquery, but it includes references to columns from the outer query in its WHERE clause.
- **Performance**: Correlated subqueries can be less efficient than non-correlated subqueries because they are executed multiple times (once for each row in the outer query). Proper indexing and query optimization are crucial for improving performance.
- **Example**: An example of a correlated subquery might involve finding all employees whose salary is greater than the average salary in their department. The subquery would calculate the average salary for each department and compare it to the salary of each employee in that department.

# Correlated Subqueries Example

```sql
-- Create the employees table with Indian names for illustration
CREATE TABLE employees (
    employee_id INT AUTO_INCREMENT PRIMARY KEY,
    first_name VARCHAR(50),
    last_name VARCHAR(50),
    department VARCHAR(50),
    salary DECIMAL(10, 2)
);
-- Insert sample data with Indian names
INSERT INTO employees (first_name, last_name, department, salary) VALUES
('Amit', 'Sharma', 'HR', 55000.00),
('Priya', 'Verma', 'Finance', 60000.00),
('Rajesh', 'Yadav', 'IT', 70000.00),
('Sneha', 'Patel', 'IT', 65000.00),
('Kiran', 'Reddy', 'HR', 58000.00);
-- Find employees with salaries greater than the average salary in their department
SELECT first_name, last_name, department, salary
FROM employees e
WHERE salary > (
    SELECT AVG(salary)
    FROM employees
    WHERE department = e.department
);
```

# Views

- **Definition**: Views in a database are virtual tables created by defining a SQL query. They don't store data themselves but provide a dynamic way to access and present data from one or more underlying tables.
- **Abstraction Layer**: Views serve as an abstraction layer, allowing users to interact with the data in a simplified and controlled manner without needing to know the underlying database structure.
- **Data Security**: Views can restrict access to sensitive or confidential data by exposing only the necessary columns and rows to specific users or roles. This enhances data security.
- **Data Consistency**: Views help maintain data consistency by presenting a consistent and standardized view of the data, even if the underlying tables undergo changes.
- **Complex Queries**: Views are useful for simplifying complex SQL queries. They allow you to encapsulate complex joins, aggregations, or calculations into a single, easy-to-use object.
- **Performance**: While views simplify queries, they can impact performance if not used carefully. Complex views may execute multiple joins and calculations, which can slow down query execution. Optimizing views and underlying queries is crucial for maintaining good performance.

# Views Example

```sql
-- Create the employees table for illustration
CREATE TABLE employees (
    employee_id INT AUTO_INCREMENT PRIMARY KEY,
    first_name VARCHAR(50),
    last_name VARCHAR(50),
    salary DECIMAL(10, 2)
);
-- Insert sample data with Indian names
INSERT INTO employees (first_name, last_name, salary) VALUES
('Amit', 'Sharma', 50000.00),
('Priya', 'Verma', 55000.00),
('Rajesh', 'Yadav', 60000.00),
('Sneha', 'Patel', 48000.00);
-- Create a view to display employee names and salaries
CREATE VIEW employee_view AS
SELECT CONCAT(first_name, ' ', last_name) AS employee_name, salary
FROM employees;
-- Query the view
SELECT * FROM employee_view;
```

# Stored Procedures

- **Definition**: A stored procedure is a precompiled and reusable SQL code block stored in a database. It can contain a series of SQL statements, control-of-flow statements, and variables. Stored procedures are typically created to perform specific tasks or operations.
- **Reusability**: Stored procedures promote code reusability by allowing developers to define a set of actions that can be executed multiple times without rewriting the same SQL code. This simplifies maintenance and ensures consistency.
- **Security**: Stored procedures can enhance security by controlling access to data. Users can be granted permission to execute a stored procedure while being restricted from directly accessing underlying tables. This helps prevent SQL injection and unauthorized data access.
- **Performance**: Stored procedures can improve database performance because they are precompiled and stored in a compiled form. This reduces the overhead of parsing and optimizing SQL statements each time they are executed.
- **Transaction Management**: Stored procedures can be used to encapsulate a series of SQL statements within a single transaction. This ensures that a set of operations either all succeed or all fail, maintaining data integrity.

# Stored Procedure Example

```
-- Create the employees table for illustration
CREATE TABLE employees (
    employee_id INT AUTO_INCREMENT PRIMARY KEY,
    first_name VARCHAR(50),
    last_name VARCHAR(50),
    department VARCHAR(50)
);
-- Insert sample data with Indian names
INSERT INTO employees (first_name, last_name, department) VALUES
('Amit', 'Sharma', 'HR'),
('Priya', 'Verma', 'Finance'),
('Rajesh', 'Yadav', 'IT'),
('Sneha', 'Patel', 'IT');
-- Create a stored procedure
DELIMITER //
CREATE PROCEDURE GetEmployeesByDepartment(IN dept VARCHAR(50))
BEGIN
    SELECT * FROM employees WHERE department = dept;
END;
//
DELIMITER ;
-- Call the stored procedure
CALL GetEmployeesByDepartment('IT');
```

# Triggers

- **Definition**: A trigger is a database object that automatically executes a predefined action (a set of SQL statements) in response to a specific event that occurs within a database table. Triggers are typically used to enforce data integrity, log changes, or perform certain actions when data modifications occur.
- **Event-Based Execution**: Triggers are activated based on specific events, such as INSERT, UPDATE, DELETE, or even a combination of these actions, occurring in a table. These events are defined by the trigger's type (BEFORE or AFTER).
- **Data Integrity**: Triggers are often used to enforce data integrity constraints, ensuring that data modifications conform to predefined rules. For example, a trigger can prevent the deletion of a record related to other records or calculate and update derived values automatically.
- **Logging and Auditing**: Triggers are valuable for logging changes to data. They can be used to maintain an audit trail by recording who made changes, what changes were made, and when these changes occurred. This is crucial for compliance and troubleshooting.
- **Complex Business Logic**: Triggers can encapsulate complex business logic within the database itself. This is particularly useful when the same logic needs to be applied consistently across various parts of an application.
- **Performance Impact**: While triggers offer powerful functionality, they can also impact performance, especially if they involve complex operations or are triggered frequently. Careful design and testing are necessary to ensure that triggers do not cause performance bottlenecks.

# Trigger Example

```sql
-- Creating a table for illustration
CREATE TABLE orders (
    order_id INT AUTO_INCREMENT PRIMARY KEY,
    order_date TIMESTAMP,
    order_total DECIMAL(10, 2)
);
-- Creating a trigger
DELIMITER //
CREATE TRIGGER update_order_date_trigger BEFORE INSERT ON orders
FOR EACH ROW
BEGIN
    SET NEW.order_date = NOW(); -- Update the order_date with the current timestamp
END;
//
DELIMITER ;
-- Inserting a new order
INSERT INTO orders (order_total) VALUES (100.00);
-- Check the updated order
SELECT * FROM orders;
```

# Functions

- **Purpose**: Functions in databases are primarily used for encapsulating logic related to data manipulation, transformation, and validation.
- **Usage**: They are employed when there is a need for code reusability, especially for complex queries, and they facilitate maintaining data integrity.
- Database functions provide flexibility via parameter acceptance, while abstracting underlying details, and enhancing security through controlled data access.
- They optimize query performance by encapsulating complex operations, promote code maintenance, and offer portability for migrating data across systems.
- **Examples**: Functions can be utilized in triggers for automating actions on data events, for data validation before insertion or update operations, and for generating reports and summaries through aggregation and calculations.

# Function Example

```sql
-- Create a function in PostgreSQL
CREATE OR REPLACE FUNCTION calculate_square(input_number numeric)
RETURNS numeric AS
$$
BEGIN
  RETURN input_number * input_number;
END;
$$
LANGUAGE plpgsql;
-- Usage example:
SELECT calculate_square(5); -- This will return 25
```

# ER (Entity-Relationship) diagrams

- **Entities**: Entities represent real-world objects or concepts within a database. They are typically nouns, such as "Customer," "Product," or "Employee."
- **Attributes**: Attributes are properties or characteristics of entities. They describe the data we want to store about each entity. For example, a "Customer" entity might have attributes like "CustomerID," "Name," "Email," and "Phone."
- **Relationships**: Relationships define how entities are connected to each other. They indicate how data from one entity relates to data in another entity. Common relationship types include one-to-one, one-to-many, and many-to-many.
- **Cardinality**: Cardinality specifies the number of instances of one entity that are related to one instance of another entity in a given relationship. Common cardinalities are "1" (one), "0..1" (zero or one), "0..n" (zero to many), and "1..n" (one to many).
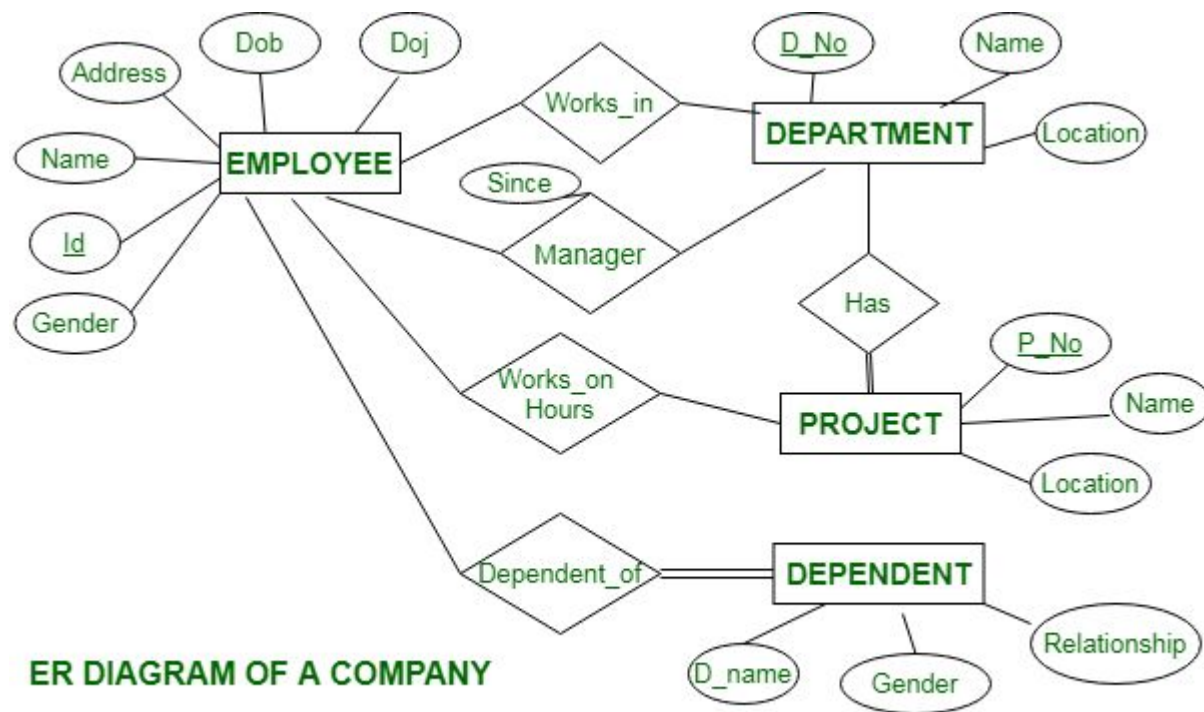
# ER (Entity-Relationship) diagrams contd…

- **Primary Key**: A primary key uniquely identifies each instance (row) of an entity. It is usually one or more attributes that ensure the uniqueness of the data within that entity.
- **Foreign Key**: A foreign key is an attribute in one entity that is used to establish a link to the primary key in another entity. It creates relationships between entities.
- **Weak Entities**: Some entities do not have a primary key on their own and rely on another entity for identification. These are called weak entities, and they are identified by a partial key.
- **Supertype and Subtype**: In some cases, entities may have a common set of attributes but also unique attributes. The common attributes are grouped in a supertype entity, while the unique ones are in subtype entities.
- **Inheritance**: In the context of ER diagrams, inheritance refers to the concept of creating a new entity (subtype) based on an existing entity (supertype). Subtypes inherit attributes and relationships from their supertype.

# ER (Entity-Relationship) diagrams contd…

- **Associative Entities**: Sometimes, you may need to introduce an intermediary entity to represent a many-to-many relationship between two entities. This intermediary entity is often called an associative entity.
- **Relationship Attributes**: In some cases, you may want to add attributes to a relationship itself, rather than to the participating entities. These are known as relationship attributes and are used to store information specific to the relationship.
- **Diagram Notation**: ER diagrams use specific symbols and notations to represent entities, attributes, relationships, cardinalities, and keys. Common symbols include rectangles for entities, ovals for attributes, and lines with various markings for relationships.
- **Normalization**: Understanding the concept of normalization is essential for designing efficient databases. It involves organizing data to minimize redundancy and dependency issues.
- **ERD Tools**: Familiarize yourself with software tools like Microsoft Visio, Lucidchart, or online ERD generators that help create and edit ER diagrams.

# ER Diagram Example



ER DIAGRAM OF A COMPANY

# Databases and when to use them

1) Relational Databases (SQL):

- Use when your data is structured and consistent.

- Supports ACID transactions and complex queries.

- Examples: MySQL, PostgreSQL, Oracle.

2) Key-Value Store:

- Use when data model is based on key-value pairs and requires high scalability and availability.

- Provides lightning-fast data retrieval and high throughput.

- Examples: Aerospike, DynamoDB

# Types of Databases contd…

3) Document Databases:

- Handles semi-structured data with varying fields.

- Provides schema flexibility and horizontal scaling.

- Examples: MongoDB, Couchbase

4) Graph Databases:

- Perfect for data with complex relationships.

- Used in applications like social networks and recommendation engines.

- Examples: Neo4j, Amazon Neptune.

# Types of Databases contd…

5) Columnar Databases:

- Data is stored by columns instead of rows to optimize reading from a column.

- Great for applications that involve storing massive data sets and running analytical queries.

- Examples: HBase, Redshift.

6) Time-Series Databases:

- Opt for this when dealing with time-series data like IoT sensor readings or server logs.

- Great for efficient storage and retrieval of time-stamped data.

- Examples: InfluxDB, Prometheus.

# Types of Databases contd…

7) In-Memory Databases:

- When speed is of the essence, and you can afford to sacrifice persistence.

- Ideal for caching, real-time analytics, and high-frequency trading.

- Redis and Memcached are popular choices.

8) Wide-Column Stores:

- Use in applications with large volumes of data and high write throughput.

- Best suited for read-heavy, analytical workloads.

- Apache Cassandra is a prominent example.

# Types of Databases contd…

9) Search Engines:

- When your primary use case revolves around full-text search.

- Essential for applications that require searching of data content.

- Elasticsearch and Solr are popular choices.


10) Spatial Databases:

- Used to store geographical and location-based data.

- Choose for applications that require Spatial indexing, geospatial analytics.

- Examples include PostGIS, CartoDB

# Types of Databases contd…

11) Blob Datastore:

- Use in applications that requires storing large documents, images, audio and video files.

- Provides high availability, durability and cost effective storage.

- Examples include HDFS, Amazon S3

12) Ledger Databases:

- Used for maintaining a transparent, immutable, and cryptographically verifiable transaction log.

- Useful for applications dealing with financial transactions and supply chain systems

- Examples: Amazon QLDB, Azure SQL Ledger

# Types of Databases contd…

13) Vector Database:

- Designed for efficient semantic search, excelling in similarity-based data retrieval tasks

- Crucial for applications like recommendation systems and content discovery.

- Examples: Milvus and Faiss are leading vector databases, optimized for high-dimensional vector

indexing and similarity search in fields like NLP and recommendation.

Tech world offers a rich ecosystem of databases and there's no one-size-fits-all solution.

The choice of a database depends on your specific use case, data model, scalability needs, and budget.

# Error Handling

Error handling in database management systems (DBMS) like MySQL is crucial to ensure the reliability and integrity of data and to provide meaningful feedback to users and developers when something goes wrong. Error handling typically involves the identification, reporting, and resolution of errors that can occur during database operations.

Here are some key aspects of error handling in DBMS, with a focus on MySQL:

- Error Codes and Messages:
  - DBMS, including MySQL, use error codes and error messages to communicate the nature of an error to the user or application. Each error has a unique error code and a corresponding error message.
- Exception Handling:
  - In MySQL, you can use the DECLARE, BEGIN, END, and EXCEPTION keywords to create stored procedures or triggers for handling exceptions. This allows you to define custom error handling logic.

# Error Handling contd…

- SQLSTATE Codes:
  - MySQL uses SQLSTATE codes to categorize errors. You can use these codes to catch specific types of errors and handle them accordingly.
- TRY...CATCH Blocks (Starting from MySQL 8.0):
  - MySQL 8.0 introduced support for structured error handling using TRY...CATCH blocks in stored procedures. This allows you to write more structured and readable error handling code.
- Error Logging:
  - MySQL provides error logs where you can find detailed information about errors, warnings, and other events. You can configure error logging settings in the MySQL configuration file (my.cnf) or through the MySQL command-line client.

# Error Handling contd…

- Transactions:
    - Transactions are essential for maintaining data consistency. You can use START TRANSACTION, COMMIT, and ROLLBACK statements to control transactions. Proper error handling is crucial within transactions to ensure that data is not left in an inconsistent state.
- Application-Level Error Handling:
    - In addition to DBMS-specific error handling, you should implement error handling in your application code that interacts with the database. This includes checking for error codes and messages returned by MySQL and taking appropriate actions, such as logging errors and notifying users.
- Security Considerations:
    - When handling errors, be cautious not to expose sensitive information, such as database structure or data, to users. Use generic error messages and log detailed error information securely.

# Error Handling Example

```sql
-- Create the employees table with Indian names for illustration
CREATE TABLE employees (
    employee_id INT AUTO_INCREMENT PRIMARY KEY,
    first_name VARCHAR(50),
    last_name VARCHAR(50),
    department VARCHAR(50),
    salary DECIMAL(10, 2)
);
-- Insert sample data with Indian names
INSERT INTO employees (first_name, last_name, department, salary) VALUES
('Amit', 'Sharma', 'HR', 55000.00),
('Priya', 'Verma', 'Finance', 60000.00);
-- Create a stored procedure with error handling
DELIMITER //
CREATE PROCEDURE InsertEmployee(
    IN in_first_name VARCHAR(50),
    IN in_last_name VARCHAR(50),
    IN in_department VARCHAR(50),
    IN in_salary DECIMAL(10, 2)
)
BEGIN
    DECLARE duplicate_error CONDITION FOR SQLSTATE '23000';

    DECLARE EXIT HANDLER FOR duplicate_error
    BEGIN
        SIGNAL SQLSTATE '45000'
        SET MESSAGE_TEXT = 'Employee with the same name already exists in the department.';
    END;
    INSERT INTO employees (first_name, last_name, department, salary)
    VALUES (in_first_name, in_last_name, in_department, in_salary);
END;
//
DELIMITER ;
-- Attempt to insert a duplicate employee record
CALL InsertEmployee('Amit', 'Sharma', 'HR', 58000.00);
```

# Data Import

- **File Formats:**
  - Data can be stored in various file formats, such as CSV (Comma-Separated Values), Excel spreadsheets, JSON, XML, and more. Choose the format that best suits your data and the target system.
- **Database Management Systems (DBMS):**
  - Most DBMS provide utilities or commands to import data from external files into database tables. For example, in MySQL, you can use the LOAD DATA INFILE statement for CSV or TSV files.
- **ETL (Extract, Transform, Load) Tools:**
  - ETL tools like Apache NiFi, Talend, or Informatica are designed for complex data integration tasks. They can extract data from various sources, transform it as needed, and load it into a target database.
- **Data Mapping:**
  - Ensure that the data in your source file aligns with the structure of the target database table. Map columns in the source to corresponding columns in the destination table.

# Data Import contd…

- **Data Validation:**
    - Perform data validation during the import process to identify and handle errors. This includes checking for data type mismatches, missing values, and duplicate records.
- **Batch Processing:**
    - For large datasets, consider breaking the import process into smaller batches to prevent memory and performance issues.
- **Error Handling:**
    - Implement error handling mechanisms to gracefully handle errors during data import. Log errors for debugging and auditing purposes.
- **Data Transformation:**
    - If necessary, apply data transformations during import to format or manipulate data to match the target schema or business requirements.
- **Concurrency and Locking:**
    - Be aware of database concurrency and locking issues when importing data into a live production database. Schedule imports during low-traffic periods if possible.

# Data Export

- **Export Formats:**
    - Choose the appropriate format for data export based on your needs. Common export formats include CSV, Excel, JSON, XML, and database-specific dump files.
- **Export Commands:**
    - Most DBMS provide commands or utilities for exporting data. For example, in MySQL, you can use the SELECT ... INTO OUTFILE statement to export data to a file.
- **Query Filters:**
    - Use SQL queries to filter the data you want to export. You can select specific columns, apply WHERE conditions, and use JOINs to retrieve data from multiple tables.
- **Scheduled Exports:**
    - For regular reporting or data backup purposes, consider scheduling automated exports using tools like cron jobs or task schedulers.
- **Compression:**
    - Compress exported files to reduce storage space and improve data transfer efficiency, especially when dealing with large datasets.

# Data Export contd…

- **Data Privacy and Security:**
  - Ensure that sensitive data is appropriately masked or encrypted during export to protect sensitive information.
- **Data Serialization:**
  - When exporting structured data (e.g., objects), consider serializing it into a suitable format (e.g., JSON or XML) for storage or transfer.
- **Metadata Export:**
  - Include metadata such as column descriptions, data dictionaries, and schema information in your export if needed for documentation or analysis.
- **Data Export Logs:**
  - Keep logs of data exports, including timestamps, export parameters, and destination details, for auditing and traceability.
- **Data Validation:**
  - After exporting data, perform validation checks to ensure that the exported data matches the expected results.

# Performance Monitoring

- **Set Performance Goals:**
    - Define clear performance goals and metrics that are relevant to your database and application. These goals could include response time, throughput, resource utilization, and more.
- **Choose the Right Monitoring Tools:**
    - Select appropriate monitoring tools and software that are compatible with your database system. Many RDBMSs offer built-in monitoring features, and there are also third-party monitoring solutions available.
- **Monitor Key Metrics:**
    - Continuously monitor critical performance metrics, including:
        - **Query performance:** Analyze query execution times and identify slow queries.
        - **Resource utilization:** Monitor CPU, memory, disk I/O, and network usage.
        - **Locks and deadlocks:** Keep an eye on locking issues that can impact concurrency.
        - **Buffer pool and cache usage:** Understand how efficiently your database utilizes memory.
        - **Storage space:** Track database size and available storage.

# Performance Monitoring contd…

- **Alerting and Thresholds:**
    - Implement alerting mechanisms that notify you when predefined performance thresholds are breached. This allows you to take proactive action before performance problems affect users.
- **Historical Data and Trends:**
    - Store historical performance data and analyze trends over time to identify long-term issues or performance degradation.
- **Performance Baselines:**
    - Establish performance baselines under normal operating conditions. This helps you quickly detect deviations from the norm when issues arise.
- **Anomaly Detection:**
    - Utilize anomaly detection techniques and machine learning algorithms to automatically identify unusual patterns or deviations in performance metrics.
- **Query Analysis:**
    - Use query analysis tools to identify and optimize slow or inefficient database queries. Tools like MySQL's EXPLAIN statement can help you understand query execution plans.

# Performance Monitoring contd…

- **Indexing and Optimization:**
  - Regularly review and optimize database indexes to improve query performance. Indexes can significantly impact the speed of data retrieval.
- **Database Configuration:**
  - Review and fine-tune database configuration settings based on workload and performance requirements. Adjust parameters related to memory allocation, caching, and thread handling as needed.
- **Concurrency Management:**
  - Monitor and manage database concurrency, particularly in multi-user environments. Identify and address locking and contention issues.
- **Backup and Maintenance Monitoring:**
  - Keep an eye on database backups, maintenance tasks, and data integrity checks. Ensure that backups are running as scheduled and that they can be restored when needed.
- **Scaling and Load Balancing:**
  - Plan for scalability and implement load balancing strategies to distribute database workloads effectively, especially in high-traffic scenarios.

# Performance Monitoring contd…

- **Logging and Auditing:**
  - Enable database logging and auditing to track changes and access to sensitive data. Analyze logs for security and performance insights.
- **Documentation and Reporting:**
  - Document your monitoring procedures, findings, and actions taken. Generate regular performance reports for stakeholders.
- **Collaboration and Communication:**
  - Foster collaboration among database administrators, developers, and system administrators to address performance issues collectively.
- **Capacity Planning:**
  - Forecast future growth and resource requirements based on historical data and trends. Ensure that your database can scale to meet future demands.
- **Security and Compliance:**
  - Consider security and compliance implications when monitoring database performance. Ensure that monitoring tools and practices adhere to security best practices.

# Database Security Best Practices

- **Authentication & Authorization**
  - Strong user authentication
  - Least privilege principle
- **Encryption**
  - Data at rest: TDE
  - Data in transit: SSL/TLS
- **Patch Management**
  - Regular updates
  - Vulnerability mitigation
- **Access Control**
  - Firewalls & network segmentation
  - Intrusion detection systems
- **Audit & Monitoring**
  - Activity logging
  - Real-time monitoring

# Database Security Best Practices contd…

- **Data Protection**
    - Masking & redaction
    - Encryption within applications
- **Backup & Recovery**
    - Regular backups
    - Disaster recovery plans
- **Secure Configuration**
    - Minimize services
    - Strong system-level account control
- **Parameterized Queries**
    - SQL injection prevention
- **User Training**
    - Security awareness
    - Secure coding practices
- **Compliance & Regulations**
    - GDPR, HIPAA, PCI DSS
- **Security Testing**
    - Assessments & penetration testing
- **Incident Response**
    - Plan & procedures
- **Continuous Improvement**
    - Adapt to evolving threats