

Table of Contents

CHAPTER 1.....	6
AN OVERVIEW OF DATABASE SYSTEMS AND MANAGEMENT	6
<i>Executive Summary</i>	6
1. <i>Foundational Concepts: Database vs. DBMS</i>	6
2. <i>The Rationale for DBMS: Overcoming File System Limitations</i>	7
3. <i>Core Characteristics and Objectives of the Database Approach</i>	7
4. <i>The System Catalog: The Database's Self-Description</i>	8
5. <i>Transaction Management and the ACID Properties</i>	9
6. <i>Database Security Imperatives</i>	10
7. <i>The Database Design and Implementation Lifecycle</i>	11
8. <i>Database Users and Roles</i>	11
9. <i>A Pragmatic View: When a DBMS Is Not Necessary</i>	12
INTRODUCTION TO DATABASE SYSTEMS: A STUDY GUIDE.....	12
<i>Short-Answer Quiz</i>	12
<i>Essay Questions</i>	14
<i>Answer Key (Short)</i>	15
GLOSSARY OF KEY TERMS	16
CHAPTER 2.....	20
CONCEPTUAL DATABASE DESIGN: THE ENTITY-RELATIONSHIP MODEL	20
<i>Executive Summary</i>	20
1. <i>Foundational Concepts in Database Design</i>	20
2. <i>The Entity-Relationship (ER) Model Components</i>	22
3. <i>Relationships in the ER Model</i>	24
4. <i>Constraints on Relationships</i>	24

CS 340 Summary

<i>5. Advanced ER Concepts.....</i>	26
<i>6. Summary of ER Diagram Notations.....</i>	27
STUDY GUIDE: CONCEPTUAL DESIGN AND THE ENTITY-RELATIONSHIP MODEL	27
QUIZ	27
<i>Essay Questions.....</i>	29
<i>Answer Key (Short).....</i>	30
GLOSSARY OF KEY TERMS AND CONCEPTS.....	32
ER DIAGRAM NOTATION SUMMARY.....	35
CHAPTER 3.....	36
ENHANCED ENTITY-RELATIONSHIP (EER) MODELING CONCEPTS	36
<i>Executive Summary</i>	36
<i>1. Core Concepts: Superclasses and Subclasses</i>	36
<i>2. Methodologies for Hierarchy Design</i>	37
<i>3. Constraints and Classifications in EER.....</i>	38
<i>4. Visualizing EER Models: Diagrammatic Representations</i>	39
STUDY GUIDE FOR THE ENHANCED ENTITY-RELATIONSHIP (EER) MODEL.....	40
<i>Short-Answer Quiz.....</i>	40
<i>Essay Questions.....</i>	41
<i>Answer Key (Short).....</i>	42
GLOSSARY OF KEY TERMS	44
CHAPTER 4.....	45
BRIEFING: THE RELATIONAL DATABASE MODEL	45
<i>Executive Summary</i>	45
<i>I. Core Concepts of the Relational Model.....</i>	45
<i>II. Fundamental Properties of Relations.....</i>	46

CS 340 Summary

<i>III. Keys and Integrity Constraints</i>	47
<i>IV. Introduction to Structured Query Language (SQL)</i>	48
RELATIONAL MODEL STUDY GUIDE	49
<i>Quiz: Test Your Knowledge</i>	49
<i>Essay Questions</i>	50
<i>Answer Key (Short).....</i>	51
GLOSSARY OF KEY TERMS	52
CHAPTER 6.....	54
BRIEFING DOCUMENT: FUNDAMENTALS OF BASIC SQL	54
<i>Executive Summary</i>	54
<i>Introduction to the SQL Query Language.....</i>	54
<i>Data Definition Language (DDL) Overview.....</i>	55
<i>Core Data Retrieval: The SELECT Statement.....</i>	55
<i>Filtering and Selection with the WHERE Clause</i>	56
<i>Advanced Data Retrieval Operations.....</i>	57
<i>Set Operators for Combining Queries.....</i>	58
<i>Data Manipulation Language (DML): Modification Commands</i>	58
<i>INSERT Statement</i>	58
<i>Conclusion and Next Steps.....</i>	59
STUDY GUIDE FOR BASIC SQL (CS 340).....	60
<i>Short-Answer Quiz.....</i>	60
<i>Essay Questions</i>	61
<i>Answer Key (Short).....</i>	62
GLOSSARY OF KEY TERMS	64
CHAPTER 7.....	66

CS 340 Summary

ADVANCED SQL CONCEPTS AND IMPLEMENTATION	66
<i>Executive Summary</i>	66
1. <i>Handling NULL Values and Three-Valued Logic</i>	66
2. <i>Complex Query Construction</i>	67
3. <i>Data Aggregation and Grouping</i>	68
4. <i>Views: Virtual Tables for Simplification and Security</i>	69
5. <i>Database Schema Modification</i>	70
6. <i>Integrating SQL with Application Code</i>	70
ADVANCED SQL STUDY GUIDE.....	71
<i>Short-Answer Quiz</i>	71
<i>Essay Questions</i>	72
<i>Answer Key (Short)</i>	73
GLOSSARY OF KEY TERMS	75
CHAPTER 9.....	77
A BEGINNER'S GUIDE TO DATABASE NORMALIZATION: TAMING YOUR DATA.....	77
1. <i>Introduction: Why Do We Need to Organize Data?</i>	77
2. <i>The Building Blocks: Understanding Keys and Dependencies</i>	78
3. <i>The Step-by-Step Journey: The Normal Forms</i>	79
4. <i>The Final Picture: A Well-Organized Database</i>	81
5. <i>Conclusion: The Power of Clean Data</i>	81
STUDY GUIDE: DATABASE NORMALIZATION	82
<i>Short-Answer Quiz</i>	82
<i>Essay Questions</i>	83
<i>Answer Key (Short)</i>	84
GLOSSARY OF KEY TERMS	86

CS 340 Summary

CHAPTER 10.....	88
NoSQL AND OBJECT-ORIENTED DATABASES: A BRIEFING	88
<i>Executive Summary</i>	88
<i>The Rationale for Non-Relational Databases</i>	89
<i>Introduction to the NoSQL Movement</i>	90
<i>Primary NoSQL Data Models</i>	90
<i>The Document Model in Detail: JSON and BSON</i>	91
STUDY GUIDE: NoSQL AND OBJECT-ORIENTED DATABASES	94
<i>Short-Answer Quiz</i>	95
<i>Essay Questions</i>	96
<i>Answer Key (Short)</i>	97
GLOSSARY OF KEY TERMS.....	99

Chapter 1

An Overview of Database Systems and Management

Executive Summary

A database is an organized collection of related data, while a Database Management System (DBMS) is the software that manages it, providing reliability, security, concurrency, and efficiency. A DBMS represents a significant advancement over older file-based systems by centralizing data management and solving inherent issues such as data redundancy, security vulnerabilities, and concurrent access conflicts. The core advantages of the database approach are rooted in its main characteristics: a self-describing nature via a system catalog containing metadata, program-data independence, data abstraction, and robust support for multiple user views.

To ensure data integrity during operations, transactions are governed by the ACID properties—Atomicity, Consistency, Isolation, and Durability. These properties guarantee that transactions are processed as a single, complete unit, leave the database in a correct state, do not interfere with one another, and have permanent effects. Comprehensive database security is a critical parallel concern, encompassing access control, encryption, auditing, and other measures to protect data, the DBMS, and associated applications from both internal and external threats. The development of a database follows a structured lifecycle, including conceptual modeling, logical design, optimization, and implementation of querying capabilities, often extending to advanced applications like data warehousing and data mining.

1. Foundational Concepts: Database vs. DBMS

The terms “database” and “Database Management System” (DBMS) are foundational to understanding data management, representing the data itself and the software that controls it, respectively.

- **Database:** Defined as “an organized collection of related data, stored on disk, sometimes accessed by multiple people.” A database is composed of entities (e.g., Student, Course) and the relationships between them (e.g., a Student *takes* a Course). Examples include a university database containing data on students and faculty, or a movie database with information on actors and genres.
- **Database Management System (DBMS):** The software system that manages databases. It facilitates crucial operations such as accessing, modifying, backing up, and recovering the database.

A DBMS serves as the intermediary between the users/applications and the physical database. Prominent examples include Oracle, PostgreSQL, MySQL, and MongoDB.

The fundamental distinction is that the **database is the data**, while the **DBMS is the software that manages that data**.

2. The Rationale for DBMS: Overcoming File System Limitations

Before the advent of the DBMS, information was stored using file systems provided by the operating system, where each application program defined and managed its own data. This approach, known as a file-based system, presented significant challenges that a DBMS is designed to solve. The central question is, “Why ‘design’ a database and use a management system? Why not just throw my data into a file and be done with it?”

The answer lies in the inherent limitations of file-based systems compared to the comprehensive functionality of a DBMS.

Feature	File-Based System Challenges	DBMS Solution
Security	Difficult to enforce granular security and protect against unauthorized access.	Provides robust mechanisms for access control, user authentication, and authorization.
Concurrent Access	Prone to conflicts and data inconsistency when multiple users access data simultaneously.	Manages concurrent access to ensure transactions are executed correctly without interference (e.g., OLTP).
Recovery	Offers limited to no automated recovery from failures like hard drive crashes or power loss.	Includes recovery subsystems to ensure completed transactions are permanently recorded and the database is durable.
Performance	Lacks optimized lookup mechanisms; cannot guarantee access speeds.	Employs sophisticated query optimization and indexing to ensure efficient data retrieval.
Redundancy	Leads to data duplication and inconsistency, as data is repeated across different files.	Reduces redundancy through careful design (e.g., normalization), improving data integrity and saving space.

3. Core Characteristics and Objectives of the Database Approach

The database approach is defined by a set of core objectives and characteristics that enable powerful, flexible, and reliable data management.

Objectives of a DBMS

- **Data Availability:** To make data available to various users in a meaningful format at a reasonable cost.
- **Data Integrity:** To ensure the data in the database is reliable and correct.
- **Data Security:** To ensure only authorized users can access the data, often enforced with passwords and access controls.
- **Data Independence:** To hide the details of how data is stored and maintained, allowing users to store, update, and retrieve data efficiently without needing to know the physical storage details.

Main Characteristics

- **Self-Describing Nature:** A database system stores not only the data but also a complete description of the database structure, types, and constraints. This “data about data,” known as **meta-data**, is stored in a **System Catalog**. This allows the DBMS software to be generic and work with many different database applications.
- **Insulation Between Programs and Data (Program-Data Independence):** The structure of data files is stored in the DBMS catalog separately from the programs that access the data. This allows the database’s structure or storage organization to be changed without requiring changes to the application programs.
- **Data Abstraction:** This characteristic allows for program-data independence. A DBMS provides a conceptual representation of data through a **data model**, which hides complex storage details from users and applications.
- **Support of Multiple Views:** Each user can be presented with a different view of the database that shows only the data relevant to them, enhancing security and simplifying user interaction.
- **Sharing of Data and Multi-User Transaction Processing:** A DBMS is designed to allow multiple concurrent users to retrieve and update the database. Concurrency control mechanisms guarantee that each transaction is executed correctly or aborted, while a recovery subsystem ensures the permanent recording of completed transactions. This is a major part of **Online Transaction Processing (OLTP)** applications.

4. The System Catalog: The Database’s Self-Description

The System Catalog, often used interchangeably with the term “data dictionary,” is a critical component of a relational DBMS that embodies its self-describing nature. It is a collection of tables and views that contains metadata, defining the structure of the database itself.

Key Information Stored in the System Catalog

- **Schema Definitions:** Relation (table) names, attribute (column) names, and attribute data types (domains).
- **Constraints:** Definitions of primary keys, foreign keys, secondary keys, and other data integrity rules (e.g., NOT NULL).
- **Views and Indexes:** Definitions of user views and the structure of indexes used for performance optimization.
- **Security and Authorization:** Information on authorized users, their passwords, and their specific privileges to access or modify data.
- **Statistical and Descriptive Data:** Information used for query processing, such as the number of tuples in a relation and storage methods.

Because the catalog is stored as a set of relations within the database, both the DBMS software and authorized users can query it using standard languages like SQL to understand the database's structure, their permissions, or other metadata.

5. Transaction Management and the ACID Properties

A **transaction** is a single logical unit of work, which may consist of one or more operations that access or modify the contents of a database. To maintain data consistency and correctness, especially in a multi-user environment, every DBMS must support the **ACID properties**.

Property	Description
A - Atomicity	Known as the “all or nothing rule.” The entire transaction is treated as a single unit that either fully completes (commits) or is fully undone (aborts). Partial completion is not possible.
C - Consistency	A transaction must bring the database from one valid, consistent state to another. All data integrity constraints must be maintained.
I - Isolation	Ensures that concurrently executing transactions do not interfere with each other. The result of concurrent execution is equivalent to a state achieved by executing the transactions serially. Changes in one transaction are not visible to others until it is committed.
D - Durability	Once a transaction has been successfully committed, its changes are permanent and will survive any subsequent system failure, such as a power loss or crash. Updates are stored in non-volatile memory.

These properties collectively provide a mechanism to ensure the correctness of the database, guaranteeing that data is not lost or corrupted during operations like inserts, updates, and deletes.

6. Database Security Imperatives

Database security is a comprehensive set of measures designed to protect a database, its DBMS, and all accessing applications from misuse, damage, and unauthorized access. Its primary goal is to maintain the **confidentiality, integrity, and availability** of sensitive information.

Sources of Security Risks

Security threats can originate from multiple sources, each with potential access to the database:

- A **malicious insider** with intent to cause harm.
- A **negligent insider** who exposes the database to attack through careless actions.
- An **outsider** who gains access through methods like social engineering.

Insider threats are noted as one of the most typical causes of security breaches, often because many employees are granted privileged user access.

Key Components of Database Security

1. **Access Control:** Restricts database access to authorized users through authentication, authorization, and role-based controls.
2. **Encryption:** Converts data into a secure, unreadable format that requires a decryption key, protecting it even if unauthorized access occurs.
3. **Audit Trails:** Records and monitors all database activities, allowing administrators to track who did what and when, which helps in identifying suspicious behavior.
4. **Data Masking and Redaction:** Hides sensitive information from users who lack the necessary permissions to view it.
5. **Backup and Recovery:** Involves regular backups of database contents and a robust recovery plan to restore data in case of loss, corruption, or a security incident.
6. **Patch Management:** Entails regularly updating the DBMS and related software to address known vulnerabilities and close security loopholes.
7. **Firewalls and Network Security:** Implements network-level defenses such as firewalls and intrusion detection systems to protect the database server from external threats.

7. The Database Design and Implementation Lifecycle

Designing and implementing a database is a structured, multi-step process that aligns with the core topics of database systems education.

1. **Conceptual Modeling:** This initial step involves defining the main data objects and their relationships. **Entity-Relationship Diagrams** are a common tool, where an **entity** is a real-world object (a noun, e.g., Student) described by **attributes** (e.g., ID, Name), and a **relationship** is an association between entities (a verb, e.g., Takes).
2. **Logical Modeling:** The conceptual model is translated into a logical structure using a specific data model. The **relational model**, which organizes data into tables, is the most widely used model today. This provides a logical view separate from the physical storage details.
3. **Optimization:** The database design is refined to ensure it captures all necessary information, minimizes storage space, and is easy to maintain. This primarily involves **normalization**.
4. **Querying:** A language is chosen to retrieve and manipulate data. **Structured Query Language (SQL)** is the standard declarative language for this purpose. It allows users to describe the desired result rather than the computational steps.
5. **Advanced Database Applications:** The database can be used for more than simple record-keeping. Key applications include:
 - **Data Warehousing & OLAP (Online Analytical Processing):** Collecting and analyzing large volumes of current and historical data to identify patterns.
 - **OLTP (Online Transaction Processing):** Managing large numbers of concurrent transactions, as seen in e-commerce or banking.
 - **Data Mining & Knowledge Discovery:** Exploring large datasets to find novel patterns and knowledge, drawing on fields like statistics and machine learning.
6. **Application Creation:** Finally, an application with a user interface is built to interact with the database. This often involves using connectivity tools like **Java Database Connectivity (JDBC)** to link a programming language (e.g., Java, Python) with the database.

8. Database Users and Roles

Different individuals interact with a database system in distinct roles, each with specific responsibilities.

- **Database Administrators (DBAs):** Responsible for the overall management of the database. Their tasks include authorizing access, coordinating and monitoring its use, acquiring necessary hardware and software, and monitoring operational efficiency.
- **Database Designers:** Responsible for defining the database's structure. They define the content, constraints, and functions or transactions. This role requires close communication with end-users to understand their needs.
- **End-Users:** The people who use the data for queries, reports, and updates. They can be categorized further:
 - **Casual End-Users:** Access the database occasionally and may require different information each time. They are often middle- or high-level managers.
 - **Naive or Parametric End-Users:** This group makes up a large portion of database users. Their jobs involve constantly querying and updating the database using predefined, standard queries known as "**canned transactions**." Examples include bank tellers, airline reservation agents, and social media users posting content.

9. A Pragmatic View: When a DBMS Is Not Necessary

Despite their immense benefits, a full-scale DBMS is not always the appropriate solution. The process of designing and implementing a DBMS involves significant overhead.

A database may not be necessary under the following conditions:

- The data is static and never changes.
- There are severe hardware limitations, such as in some embedded systems.
- There is no requirement for multiple users to access the data concurrently.

In these simpler scenarios, the complexity and cost of a DBMS may outweigh its advantages.

Introduction to Database Systems: A Study Guide

Short-Answer Quiz

Instructions: Answer the following ten questions in two to three sentences each, based on the provided source materials.

1. What is the fundamental difference between a database and a database management system (DBMS)?

CS 340 Summary

2. Identify three key problems inherent in a file-based system that a DBMS is designed to solve.
3. What are the four ACID properties that a DBMS guarantees for transactions?
4. Explain the “self-describing nature” of a database system and the role of the system catalog.
5. Who are the three main types of database users, and what is the primary role of each?
6. Define data abstraction in the context of a DBMS and explain how it relates to program-data independence.
7. What is the purpose of database security, and what three areas does a security program aim to protect?
8. Briefly describe what a transaction is and how it differs from a query.
9. According to the course materials, what is SQL, and what does it mean for it to be a “declarative language”?
10. List three scenarios in which using a full-fledged database system might not be necessary.

Essay Questions

Instructions: The following questions are designed to test a deeper, more integrated understanding of the source materials. Formulate a comprehensive response for each prompt.

1. Compare and contrast the file-based system approach with the database management system (DBMS) approach. Elaborate on how a DBMS addresses the specific challenges of security, concurrent access, recovery, performance, and data redundancy that are problematic in file-based systems.
2. Explain the concepts of “program-data independence” and the “self-describing nature” of a database system. Discuss how the system catalog, meta-data, and data abstraction work together to achieve these key characteristics of the database approach.
3. Describe the four ACID properties (Atomicity, Consistency, Isolation, Durability) in detail. Using the examples provided in the text, explain why each property is critical for ensuring the correctness and consistency of a database during online transaction processing (OLTP).
4. Outline the six-step process for designing and implementing a database application as presented in the course materials. For each step, describe its primary purpose and mention the key concepts or tools involved (e.g., entity-relationship diagrams, relational model, normalization, SQL).
5. Discuss the critical importance of database security. Identify the three primary sources of security risks mentioned in the text and describe at least four key components of a database security program (e.g., Access Control, Encryption, Audit Trails), explaining how each component mitigates potential threats.

Answer Key (Short)

1. A database is an organized collection of related data stored on disk. In contrast, a database management system (DBMS) is the software that manages this data, allowing users to access, modify, back up, and recover the database.
2. A file-based system struggles with security (unauthorized access), concurrent access by multiple users, and recovery from failures like a hard drive crash. A DBMS is designed to handle all of these functions reliably and efficiently.
3. The four ACID properties are Atomicity (transactions are all-or-nothing), Consistency (the database remains valid before and after a transaction), Isolation (concurrent transactions do not interfere with each other), and Durability (committed changes persist through system failures).
4. The self-describing nature of a database means the system stores not only the data but also a complete definition of the database structure, types, and constraints. This descriptive information, called metadata, is stored in the system catalog.
5. The three main users are: Database Administrators (DBAs), who authorize access and monitor use; Database Designers, who define the database's structure and constraints; and End-Users, who query, report on, and update the data.
6. Data abstraction is a characteristic that allows for program-data independence by providing users with a conceptual representation of data, hiding the details of how it is stored. Programs refer to this conceptual data model rather than physical storage details, allowing storage to change without affecting the access programs.
7. Database security is a set of measures designed to maintain the confidentiality, integrity, and availability of sensitive information. A security program protects the data itself, the DBMS, and all applications that access the database.
8. A transaction is a single logical unit of work, which can be a collection of requests that are executed as a single unit. A query is a specific type of request that generates a result based on the stored data.
9. SQL (Structured Query Language) is the language used to query data from a database. It is a declarative language, which means statements describe the desired result of a computation rather than describing the computation steps directly.
10. A database is not always needed. Scenarios where it may be unnecessary include when data is static and never changes, when working with significant hardware limitations like in embedded systems, or when there is no need for multiple users to access the data.

Glossary of Key Terms

Term	Definition
ACID Properties	A set of properties (Atomicity, Consistency, Isolation, Durability) that guarantee database transactions are processed reliably.
Access Control	A component of database security that restricts access to the database, ensuring only authorized users have permission to view or modify data.
Atomicity	The property that a transaction is treated as a single unit, which either executes completely or not at all. Also known as the ‘All or nothing rule’.
Attribute	A property that describes an entity. For example, a “Student” entity can be described by attributes like ID, Name, and Address.
Audit Trails	A security component that records and monitors database activities, allowing administrators to track who accessed the database, what operations were performed, and when.
Backup and Recovery	A security component involving regularly backing up database contents to prevent data loss and having a robust plan to restore data in case of a breach.
Casual End-Users	Users who occasionally access the database, may need different information each time, and are typically middle- or high-level managers.
Conceptual Modeling	The first step in database design, which involves creating a high-level description of the data. Entity-relationship diagrams are used in this stage.
Concurrency Control	A mechanism within a DBMS that guarantees that transactions executed concurrently are correctly executed or aborted, ensuring they don’t lead to inconsistency.
Consistency	The property ensuring that integrity constraints are maintained so that the database is correct and valid both before and after a transaction.
Data Abstraction	The characteristic that allows program-data independence by hiding storage details and presenting users with a conceptual view of the database.

Data Availability	A DBMS objective to make data available to a wide variety of users in a meaningful format at a reasonable cost.
Data Dictionary	A term often used interchangeably with System Catalog; it is an implementation of the system catalog.
Data Independence	A DBMS objective that allows the user to store, update, and retrieve data efficiently while hiding details of how the data is stored and maintained.
Data Integrity	A DBMS objective ensuring that the data available in the database is reliable and correct.
Data Masking and Redaction	A security component that conceals sensitive information when displayed to users without the necessary permissions.
Data Mining	The exploration and analysis of large datasets to find knowledge in the form of patterns.
Data Model	The tools used to describe data, relationships, and constraints. It's a type of data abstraction used to hide storage details.
Data Security	A DBMS objective ensuring that only authorized users can access the data, often enforced by methods like passwords.
Data Warehousing	The process where organizations collect and organize current and historical data, often from multiple sources, to identify patterns and perform analysis of complex queries.
Database	An organized collection of related data, stored on disk, and sometimes accessed by multiple people.
Database Administrator (DBA)	A user responsible for authorizing access, monitoring database use, and controlling the efficiency of operations.
Database Designer	A user responsible for defining the content, structure, constraints, and functions of the database.
Database Management System (DBMS)	Software that manages databases, allowing for accessing, modifying, backing up, and recovering data. Examples include Oracle, MySQL, and MongoDB.
Database Security	A set of measures implemented to safeguard databases, the DBMS, and associated applications from unauthorized access, misuse, damage, and intrusion.

Declarative Language	A type of language where statements describe the desired result of a computation, rather than describing the steps to perform the computation directly. SQL is an example.
Durability	The property ensuring that once a transaction has completed, its updates are permanently stored on disk and persist even if a system failure occurs.
End-Users	People who use the data for queries, reports, and sometimes update the database content.
Entity	A real-world object or noun that is described by attributes, such as a “Student” or “Course”.
File-Based System	A pre-DBMS method where the operating system’s file system was used to store information, with each application program defining and managing its own data.
Insider Threat	A security risk from a malicious or negligent person within an organization, often one who has been granted privileged user access.
Insulation between programs and data	Also called program-data independence, this is a characteristic where the structure of data files is stored separately from the access programs.
Isolation	The property that ensures multiple transactions can occur concurrently without interfering with each other or leading to an inconsistent database state.
Java Database Connectivity (JDBC)	An API used to connect Java programs to various databases like Oracle, MySQL, and IBM DB2.
Logical Modeling	The second step in database design, where a data model (like the relational model) is used to describe data, relationships, and constraints.
Meta-data	The description of a database (data structures, types, constraints) that is stored in the System Catalog. It is “data about data.”
Naive or Parametric End-Users	Users whose main job function involves constantly querying and updating the database using standard, pre-programmed queries called canned transactions.
Normalization	A process in database design optimization to ensure data is structured to take up less space and is easy to maintain.

Online Analytical Processing (OLAP)	A database application focused on the analysis of complex queries on large datasets, often used with data warehousing.
Online Transaction Processing (OLTP)	A major part of database applications that allows hundreds of concurrent transactions to execute per second.
Program-data independence	The ability to change data structures and storage organization without having to change the DBMS access programs.
Queries	Requests sent to the database that generate a result based on the stored data.
Recovery Subsystem	A component of a DBMS that ensures each completed transaction has its effect permanently recorded in the database.
Relational Model	The most used data model in the world today for logical modeling of databases.
Relationship	An association between two or more entities, often represented as a verb, such as a “Student” takes a “Course”.
SQL (Structured Query Language)	A declarative language used to query data from a database.
System Catalog	A collection of tables and views containing schema metadata, such as information about tables, columns, constraints, and users. It is the self-describing component of a database system.
Transactions	A single logical unit of work, often a collection of requests, that is executed entirely or not at all to ensure data integrity.
Views	A subset of the database that describes only the data of interest to a particular user. Each user may see a different view.

Chapter 2

Conceptual Database Design: The Entity-Relationship Model

Executive Summary

This document synthesizes the core principles of conceptual database design as articulated through the Entity-Relationship (ER) Model. The central thesis is that careful, abstract design is a mandatory prerequisite for building efficient, reliable, and useful database systems. The ER model provides the framework for this design process by identifying key entities, their attributes, and the relationships between them.

A foundational concept is **data abstraction**, which separates the user's view from the physical storage details through three distinct levels: physical, logical, and external (views). This separation enables **data independence**, allowing changes to the underlying data structure (physical independence) or logical schema (logical independence) without affecting user applications.

The ER model's primary components are **Entities** (real-world objects), **Attributes** (properties describing entities), and **Relationships** (associations between entities). A rigorous design process involves defining these components and applying **constraints** to ensure data integrity. Key constraints include cardinality ratios (one-to-one, one-to-many, many-to-many), which define the number of associations an entity can have, and participation constraints (total or partial), which dictate whether an entity's existence depends on a relationship. Advanced constructs such as weak entities, recursive relationships, and higher-degree relationships provide tools for modeling more complex scenarios. The entire conceptual schema is visualized through an Entity-Relationship (ER) diagram, which uses standardized notation to represent these components and constraints.

1. Foundational Concepts in Database Design

Before implementing a database, a conceptual design phase is critical. This phase is guided by several key principles that ensure the database is robust, flexible, and serves its intended purpose.

1.1. Data Abstraction

Data abstraction is a primary characteristic of the database approach, defined as the “suppression of details of data organization and storage” to highlight essential features for users. This is achieved through a three-level architecture:

- **Physical Level:** The lowest level, describing *how* data is physically stored on storage devices.
- **Logical Level:** The next level up, describing *what* data is stored in the database and the relationships that exist among that data. This is the level where the database schema is defined.
- **External Level (Views):** The highest level, which describes only a portion of the entire database tailored to specific user groups. For example, a student view of a university database is different from that of a registrar or a database administrator.

This tiered structure is a key advantage of relational databases, allowing different users to perceive data at their preferred level of detail.

1.2. Schemas and Instances

The design of a database is captured in a schema, while the data itself is an instance.

- **Schema:** The logical structure or blueprint of the database. It is defined once and changes very infrequently. There are three types of schemas corresponding to the levels of abstraction:
 - **Physical Schema:** Describes the database design at the physical storage level.
 - **Logical Schema:** Describes the database design at the logical level (e.g., tables and columns).
 - **External Schema (Subschema):** Describes the various external views for different users.
- **Instance:** The actual data contained in the database at a particular moment in time. The database state, or instance, changes every time the data is updated.

1.3. Data Independence

Data independence is the capacity to modify a schema definition at one level without affecting the schema at the next higher level. It protects applications from changes in data organization and definition.

- **Physical Data Independence:** The ability to change the internal (physical) schema without having to change the conceptual (logical) schema. This is crucial for performance tuning, such as changing

file organization, storage structures, or indexing strategies, without impacting the logical data view or applications.

- **Logical Data Independence:** The ability to change the conceptual (logical) schema without having to change external schemas or application programs. For example, adding or removing a column from a table should not break an existing application that does not use that column. This is more difficult to achieve than physical independence because applications are often tightly coupled to the logical structure of the data they access.

2. The Entity-Relationship (ER) Model Components

The ER model is the primary tool for conceptual design. It involves identifying entities, the attributes that describe them, and the relationships that connect them.

2.1. Entities

An **entity** is a real-world object that is distinguishable from other objects (e.g., a specific student, a particular course). Entities are typically represented by nouns.

- **Entity Set:** A collection of similar entities that share the same set of attributes (e.g., the set of all EMPLOYEES).
- **Entity Type:** The schema or description for an entity set, displayed as a rectangular box in an ER diagram.

2.2. Attributes

Attributes are the properties that describe an entity. Each attribute has a domain of possible values (e.g., integer, string).

Attribut Type	Description	ER Diagram Notation	Example
Simple	An attribute that has a single, atomic value for an entity.	Single oval	gender, Salary

CS 340 Summary

Composite	An attribute that can be divided into smaller sub-parts, which represent more basic attributes with independent meanings. This can form a hierarchy.	Oval connected to other ovals	Name composed of Fname, Minit, Lname; Address composed of Street, City, State
Multi-valued	An attribute that can hold a set of values for the same entity. It may have lower and upper bounds on the number of values.	Double oval	{Color} for a car; {Locations} for a department.
Stored	A standard attribute whose value is stored directly in the database.	Single oval	Birth_date
Derived	An attribute whose value can be calculated or derived from another attribute or related entity. It is not stored directly.	Dashed oval	Age derived from Birth_date; Number_of_employees derived by counting employees.
Complex	Attributes that are nested composites of multi-valued and composite attributes. This is rare and not considered good design.	N/A	{PreviousDegrees(College, Year, Degree, Field)}

2.3. Keys

A **key attribute** is an attribute of an entity type for which each entity must have a unique value. Keys are used to uniquely identify an entity within an entity set.

CS 340 Summary

- A key can be **composite**, consisting of multiple attributes (e.g., `VehicleTagNumber` composed of `Number` and `State`).
- An entity type may have **more than one key**. For example, a `CAR` entity might have `VehicleIdentificationNumber` and `VehicleTagNumber` as two distinct keys.
- In ER diagrams, all key attributes are **underlined**.

3. Relationships in the ER Model

A **relationship** associates two or more distinct entities with a specific meaning. Relationships are typically represented by verbs (e.g., an `EMPLOYEE` *works on* a `PROJECT`).

- **Relationship Type:** The schema description of a relationship, identifying its name and participating entity types. It is represented by a diamond shape in ER diagrams.
- **Relationship Set:** The current set of relationship instances in the database.
- **Degree (Arity):** The number of participating entity types in a relationship.
 - **Unary (Recursive):** A relationship between entities of the same type, but in different roles (e.g., an `EMPLOYEE` *supervises* another `EMPLOYEE`). Role names are used to distinguish the participations.
 - **Binary:** A relationship between two entity types (most common).
 - **Ternary:** A relationship among three entity types.

A relationship type can also have its own attributes. This is most common for many-to-many relationships, where the attribute's value depends on the combination of participating entities (e.g., `Hours` in a `WORKS_ON` relationship between `EMPLOYEE` and `PROJECT`).

4. Constraints on Relationships

Constraints are rules that ensure data integrity by limiting the possible combinations of entities in a relationship set.

4.1. Cardinality Ratios

The cardinality ratio specifies the maximum number of relationship instances an entity can participate in. For binary relationships, the common types are:

- **One-to-One (1:1):** An entity in one set is associated with at most one entity in the other set, and vice versa. (e.g., an EMPLOYEE *manages* at most one DEPARTMENT).
- **Many-to-One (N:1):** Many entities in one set can be associated with at most one entity in the other set. (e.g., many EMPLOYEES *work for* one DEPARTMENT).
- **One-to-Many (1:N):** One entity in the first set can be associated with many entities in the second, but an entity in the second set is associated with at most one in the first.
- **Many-to-Many (M:N):** An entity in one set can be associated with any number of entities in the other set, and vice versa. (e.g., an EMPLOYEE *works on* many PROJECTS, and a PROJECT has many EMPLOYEES).

4.2. Participation Constraints (Existence Dependency)

This constraint specifies the minimum number of relationship instances each entity must participate in.

- **Total Participation (Mandatory):** Every entity in the entity set must participate in at least one relationship instance. This indicates an existence dependency and is represented by a **double line** in an ER diagram. (e.g., Every STUDENT must be enrolled in at least one COURSE).
- **Partial Participation (Optional):** An entity in the entity set may or may not participate in a relationship instance. This is the default and is represented by a **single line**. (e.g., A COURSE may have no students enrolled).

4.3. The (min, max) Notation

A more precise method for specifying structural constraints uses (min, max) notation on each participation of an entity type in a relationship.

- **min:** The minimum number of relationship instances an entity *must* participate in. $\text{min}=0$ denotes optional participation, while $\text{min}>=1$ denotes mandatory (total) participation.
- **max:** The maximum number of relationship instances an entity *can* participate in. $\text{max}=1$ specifies a single participation, while $\text{max}=N$ (or n) specifies no limit.

Examples:

- An EMPLOYEE can manage at most one DEPARTMENT ((0, 1)) and a DEPARTMENT must have exactly one manager ((1, 1)).

- An EMPLOYEE must work for exactly one DEPARTMENT ((1, 1)) and a DEPARTMENT can have any number of employees ((1, N) or (4, N) if there's a minimum).

5. Advanced ER Concepts

5.1. Weak Entity Types

A **weak entity** is an entity that does not have its own key attribute and is identification-dependent on another “owner” entity type.

- It is identified by the combination of its own **partial key** (an attribute that is unique in relation to the owner) and the key of its owner entity.
- A weak entity must have **total participation** in an **identifying relationship** with its owner.
- **Notation:** Weak entities are shown in a double-lined rectangle, and their identifying relationship is shown in a double-lined diamond. The partial key is underlined with a dashed line.
- **Example:** A DEPENDENT entity is weak. Its partial key is Name. It is identified by its name *and* the SSN of the EMPLOYEE it is related to through the DEPENDENTS_OF identifying relationship.

5.2. Higher-Degree (N-ary) Relationships

While most relationships are binary, some scenarios require **ternary** (degree 3) or higher **n-ary** relationships.

- An n-ary relationship is generally **not equivalent** to n separate binary relationships, as the n-ary relationship often represents a single, indivisible fact involving all n entities.
- **Example:** A SUPPLY relationship can be ternary, linking a SUPPLIER, a PART, and a PROJECT. A single instance (s1, p1, j1) means supplier s1 supplies part p1 to project j1. This cannot be accurately broken down into three binary relationships without losing information.
- Specifying constraints on n-ary relationships is significantly more complex and potentially ambiguous than for binary relationships.

6. Summary of ER Diagram Notations

Symbol	Meaning
Rectangle	Entity Type
Double Rectangle	Weak Entity Type
Diamond	Relationship Type
Double Diamond	Identifying Relationship for a Weak Entity
Oval	Attribute
Underlined Oval	Key Attribute
Double Oval	Multi-valued Attribute
Dashed Oval	Derived Attribute
Oval connected to other ovals	Composite Attribute
Single Line to Relationship	Partial Participation
Double Line to Relationship	Total Participation
Numbers (1, N, M) on relationship lines	Cardinality Ratio
(min, max) on relationship line	Structural Constraint (min/max participation)

Study Guide: Conceptual Design and the Entity-Relationship Model

Quiz

1. Describe the three levels of data abstraction in a database system and the purpose of each level.
2. What is the difference between a database schema and a database instance? Explain which one changes more frequently and why.
3. Define the two types of data independence. Which type is considered more difficult to achieve and why?

CS 340 Summary

4. What are the core components of conceptual design using the Entity-Relationship (ER) model? Briefly describe what entities and relationships typically represent.

5. Explain the difference between a simple attribute, a composite attribute, and a multi-valued attribute, providing an example for each.

6. What is a key attribute in an entity type? Can an entity type have more than one key, and can a key be composite?

7. Define a “weak entity type” and explain how it is identified within an ER model.

8. What is a recursive relationship? Provide an example and explain the importance of using “roles” in this context.

9. Explain the difference between total and partial participation constraints in a relationship. How are these constraints represented in an ER diagram?

10. What is the degree (or arity) of a relationship? Name the three degrees mentioned in the source context.

Essay Questions

1. Using the “COMPANY Database” example, describe the process of refining an initial conceptual design. Discuss how attributes from initial entity types (e.g., “Manager” in DEPARTMENT, “Works_on” in EMPLOYEE) can be transformed into distinct relationship types in a more refined ER diagram.
 2. Compare and contrast the four main types of relationship cardinality constraints (one-to-one, one-to-many, many-to-one, and many-to-many). For each type, provide a clear definition and an example from the provided source material.
 3. Explain the concept of an attribute in detail. Your explanation should cover simple, composite, multi-valued, stored, derived, and complex attributes, providing examples for each and discussing how they are represented in ER diagrams.
 4. Discuss the complexities of modeling higher-degree (n-ary) relationships, specifically ternary relationships. Explain why a ternary relationship is generally not equivalent to multiple binary relationships, using the “SUPPLY” example to illustrate your points.

5. Describe the (min, max) notation for relationship structural constraints. Explain how this notation provides more detailed information than traditional cardinality ratios and participation constraints, using the `MANAGES` and `WORKS FOR` relationships as examples.

Answer Key (Short)

1. The three levels of abstraction are the physical, logical, and external (view) levels. The physical level describes how data is actually stored, the logical level describes what data is stored and the relationships between them, and the external level shows different parts of the database to different users.
2. A schema is the logical structure of the database, which is created during the design phase and changes infrequently. An instance is the actual data in the database at a particular moment in time, which changes every time the database is updated.
3. The two types are physical and logical data independence. Physical data independence is the capacity to change the internal (physical) schema without changing the conceptual schema. Logical data independence is the ability to change the conceptual schema without altering external schemas or application programs. Logical data independence is more difficult to achieve because applications are heavily dependent on the logical structure of the data they access.
4. The core components are entities, relationships, attributes, and integrity constraints. Entities are typically nouns representing real-world objects (e.g., Student, Course), while relationships are typically verbs representing statements about two or more objects (e.g., a Student *takes* a Course).
5. A simple attribute has a single atomic value, like `ssn`. A composite attribute can be divided into smaller parts, such as an `Address` attribute composed of Number, Street, and City. A multi-valued attribute can hold a set of values for the same entity, for instance, a `Color` attribute for a car which could be {red, black}.
6. A key attribute is an attribute of an entity type for which each entity must have a unique value, such as `ssn` for an employee. An entity type can have more than one key, and a key can be composite, meaning it is composed of multiple attributes, like `VehicleTagNumber` (Number, State).
7. A weak entity type is an entity that does not have its own key attribute and is identification-dependent on another “owner” entity type. It is identified by the combination of its partial key and the specific owner entity it is related to through an identifying relationship.
8. A recursive (or unary) relationship is a relationship between the same participating entity type in distinct roles. An example is a `SUPERVISION` relationship where an `EMPLOYEE` entity participates twice: once in the `supervisor` role and once in the `supervisee` role. Roles are necessary to distinguish the meaning of each participation.
9. Total participation means every entity in an entity set must participate in at least one instance of the relationship; it is represented by a double line. Partial participation means an entity may or may not participate in the relationship; it is represented by a single line.

- 10.** The degree or arity of a relationship is the number of entity types that participate in it. The source context names three types: unary (degree 1), binary (degree 2), and ternary (degree 3).

Glossary of Key Terms and Concepts

Term	Definition
Attribute	A property that describes an entity. Displayed in ovals in ER diagrams.
Binary Relationship	A relationship type of degree two, involving two participating entity types.
Cardinality Ratio	Specifies the number of relationships in a set in which an entity can participate (e.g., 1:1, 1:N, M:N).
Complex Attribute	An attribute where composite and multi-valued attributes are nested.
Composite Attribute	An attribute that may be divided into smaller sub-parts. For example, Name (FirstName, LastName).
Conceptual Design	The process of identifying entities, relationships, and integrity constraints for an enterprise to be modeled in a database.
Data Abstraction	The suppression of details of data organization and storage, highlighting essential features for an improved understanding of data.
Data Independence	The protection of user applications from changes made in the definition and organization of data.
Degree (Arity)	The number of entity types that participate in a relationship.
Derived Attribute	An attribute whose value can be determined from other attributes or related entities. For example, Age can be derived from Birth_date.
Entity	A real-world object distinguishable from other objects, described by a set of attributes.
Entity Set	A collection of similar entities that all have the same set of attributes.
Entity Type	A grouping of entities with the same basic attributes.

ER Diagram	Entity-Relationship Diagram; the encoding of entities, attributes, and relationships in the relational database realm.
Existence Dependency	Also called a participation constraint. Specifies the minimum participation of an entity in a relationship (optional or mandatory).
External Level (View)	The highest level of data abstraction, which describes only a part of the entire database for specific users.
External Schema (Subschema)	Describes the different views of the database for various external users.
Identifying Relationship	The relationship type that relates a weak entity type to its owner entity type. Represented by a double diamond.
Instances	The actual data in the database at a particular moment in time.
Key Attribute	An attribute of an entity type for which each entity must have a unique value. It is underlined in ER diagrams.
Logical Data Independence	The capacity to change the conceptual schema without having to change external schemas and associated application programs.
Logical Level	The level of data abstraction that describes what data is stored in the database and what relationships exist among that data.
Logical Schema	Describes the database design at the logical level.
Multi-valued Attribute	An attribute that can have a set of values for the same entity. Represented by a double oval in ER diagrams.
Partial Key	An attribute of a weak entity type that helps identify it in combination with its owner entity.

Partial Participation	An existence dependency where entities in an entity set may or may not participate in a relationship instance. Represented by a single line.
Physical Data Independence	The capacity to change the internal (physical) schema without having to change the conceptual schema.
Physical Level	The lowest level of data abstraction, which describes how data is actually stored.
Physical Schema	Describes the database design at the physical (storage) level.
Recursive Relationship	A relationship type where the same entity type participates more than once in distinct roles. Also called a unary relationship or self-referencing relationship.
Relationship	Relates two or more distinct entities with a specific meaning. Usually represented by verbs.
Relationship Set	The current set of relationship instances represented in the database at a specific moment in time.
Relationship Type	The schema description of a relationship, identifying the relationship name and participating entity types.
Schema	The logical structure of the database.
Simple Attribute	An attribute where each entity has a single atomic value.
Stored Attribute	An attribute that is directly stored in the database, from which a derived attribute's value may be obtained.
Ternary Relationship	A relationship type of degree three, involving three participating entity types.
Total Participation	An existence dependency where every entity in the entity set must participate in at least one instance of the relationship. Represented by a double line.
Unary Relationship	A relationship type of degree one. See Recursive Relationship.

Weak Entity Type	An entity that does not have a key attribute of its own and is identification-dependent on an owner entity. Represented by a double rectangle.
-------------------------	--

ER Diagram Notation Summary

Symbol	Meaning
	Entity
	Weak Entity
	Relationship
	Identifying Relationship
	Attribute
	Key Attribute
	Multivalued Attribute
	Composite Attribute
	Derived Attribute
	Total Participation
	Cardinality Ratio (e.g., 1:N)
	Structural Constraint (min, max)

Chapter 3

Enhanced Entity-Relationship (EER) Modeling Concepts

Executive Summary

The Enhanced Entity-Relationship (EER) model extends the foundational concepts of the ER model to provide more sophisticated tools for data modeling, particularly for complex applications. The core of the EER model is the introduction of superclass/subclass relationships, which allow for the creation of entity hierarchies through the processes of specialization and generalization. Specialization involves a top-down approach of defining distinct subgroups within a general entity type, while generalization is a bottom-up process of identifying common attributes among multiple entity types to create a more general superclass.

These hierarchical relationships are governed by specific constraints that ensure the logical integrity of the model. The **disjoint/overlapping** constraint determines whether an entity can belong to one or multiple subclasses simultaneously. The **total/partial participation** constraint specifies whether every entity in the superclass must belong to a subclass. Together, these concepts enable the creation of semantically rich and precise database schemas that accurately represent inheritance and complex relationships within the data.

1. Core Concepts: Superclasses and Subclasses

The fundamental innovation of the EER model is the ability to represent superclass/subclass relationships, also known as IS-A relationships. This structure allows for attribute inheritance, where a subclass entity inherits all attributes and relationships of its superclass.

- **Superclass:** A general entity type that has one or more distinct subgroups with unique attributes or relationships.
- **Subclass:** A more specific entity type that represents a subgroup of the superclass. A subclass inherits all features of the superclass and also includes its own specific attributes or relationships.

Example: The **EMPLOYEE** Hierarchy

The EER model can represent an **EMPLOYEE** entity as a superclass with multiple, distinct specializations. As shown in Figure 4.1, an **EMPLOYEE** has general attributes like **Ssn**, **Birth_date**, **Address**, and **Name**. This **EMPLOYEE** superclass can then be specialized into several subclass groupings:

- A job-role specialization consisting of SECRETARY, TECHNICIAN, and ENGINEER.
- A management specialization for the MANAGER subclass.
- A pay-type specialization consisting of HOURLY_EMPLOYEE and SALARIED_EMPLOYEE.

Each of these subclasses inherits the attributes of EMPLOYEE while also possessing its own specific attributes. For instance, SECRETARY has a Typing_speed attribute, MANAGER participates in a MANAGES relationship with PROJECT, and HOURLY_EMPLOYEE has a Pay_scale attribute.

2. Methodologies for Hierarchy Design

The EER model supports two primary methodologies for creating superclass/subclass hierarchies: specialization and generalization.

2.1 Specialization (*Top-Down Approach*)

Specialization is the process of starting with a single, general entity type and identifying subgroups that have unique characteristics. These subgroups become the subclasses in the hierarchy. This is a top-down design process.

For example, an EMPLOYEE entity type can be specialized based on job roles. If some employees are engineers and have a unique Eng_type attribute, while others are secretaries with a unique Typing_speed attribute, these groups can be modeled as ENGINEER and SECRETARY subclasses of the EMPLOYEE superclass.

2.2 Generalization (*Bottom-Up Approach*)

Generalization is the reverse, bottom-up process. It begins with two or more entity types that possess common attributes. These shared features are “generalized” into a new, higher-level entity type, which becomes the superclass. The original entity types become its subclasses.

Example: Generalizing CAR and TRUCK into VEHICLE

As illustrated in Figure 4.3, consider two separate entity types: CAR and TRUCK.

- **Initial State:** CAR has attributes Vehicle_id, License_plate_no, Price, Max_speed, and No_of_passengers. TRUCK has attributes Vehicle_id, License_plate_no, Price, No_of_axles, and Tonnage.

- **Generalization Process:** The common attributes (`Vehicle_id`, `License_plate_no`, `Price`) are used to create a new `VEHICLE` superclass.
- **Final State:** `CAR` and `TRUCK` become subclasses of `VEHICLE`, inheriting the common attributes and retaining only their specific attributes (`Max_speed` for `CAR`, `Tonnage` for `TRUCK`, etc.).

3. Constraints and Classifications in EER

To ensure logical consistency, EER specializations are defined by two key constraints: the disjoint/overlapping constraint and the participation constraint.

3.1 Disjoint vs. Overlapping Constraint

This constraint defines whether an entity from the superclass can be a member of one or multiple subclasses within the same specialization.

- **Disjoint (d):** An entity can be a member of **at most one** of the subclasses in the specialization. In Figure 4.4, the specialization of `EMPLOYEE` into {`SECRETARY`, `TECHNICIAN`, `ENGINEER`} is disjoint, meaning an employee cannot be both a secretary and an engineer.
- **Overlapping (o):** An entity can be a member of **more than one** subclass. For example, a `PERSON` superclass can be specialized into {`EMPLOYEE`, `ALUMNUS`, `STUDENT`}. This can be an overlapping specialization, allowing a single person to be a student and an employee simultaneously.

3.2 Total vs. Partial Participation Constraint

This constraint specifies whether every entity in the superclass must belong to a subclass.

- **Total Participation (double line):** Every entity in the superclass **must** be a member of at least one subclass. In the `PERSON` example, if the specialization is total, it means every individual in the database must be classified as an `EMPLOYEE`, `ALUMNUS`, or `STUDENT`.
- **Partial Participation (single line):** An entity in the superclass is **not required** to be a member of any of the subclasses. The `EMPLOYEE` specialization in Figure 4.1 is partial; an employee might have a role other than `SECRETARY`, `TECHNICIAN`, `ENGINEER`, or `MANAGER`.

3.3 Defining Specializations

The criteria for membership in a subclass can be explicitly defined or determined by the application user.

- **Attribute-Defined Specialization:** Subclass membership is determined by the value of a specific attribute in the superclass, known as the specialization attribute. In Figure 4.4, the `Job_type` attribute of the `EMPLOYEE` entity is used to define membership in the `SECRETARY`, `TECHNICIAN`, and `ENGINEER` subclasses. The condition (e.g., `Job_type = 'Secretary'`) is written on the diagram line.
- **User-Defined Specialization:** When no attribute exists to define the specialization, membership is determined externally by the application or database user. This is the default mechanism if a specialization attribute is not specified.

4. Visualizing EER Models: Diagrammatic Representations

The EER model uses specific notation to represent these advanced concepts, building upon the standard ER diagram. The diagrams analyzed provide clear examples of this notation.

Concept	EER Notation	Description	Example
Superclass/Subclass	A line connects the superclass to a circle, which in turn connects to the subclasses.	The subclasses are shown below the superclass, inheriting from it.	<code>EMPLOYEE</code> connected to subclasses like <code>SECRETARY</code> and <code>ENGINEER</code> .
Disjoint Constraint	A 'd' inside the circle.	Indicates an entity can belong to at most one subclass in the specialization.	The <code>{SECRETARY, TECHNICIAN, ENGINEER}</code> specialization of <code>EMPLOYEE</code> .
Overlapping Constraint	An 'o' inside the circle.	Indicates an entity can belong to multiple subclasses.	The <code>{EMPLOYEE, ALUMNUS, STUDENT}</code> specialization of <code>PERSON</code> .
Total Participation	A double line from the superclass to the circle.	Every entity in the superclass must be in a subclass.	The <code>PERSON</code> specialization, indicating every person must be in a subclass.

Partial Participation	A single line from the superclass to the circle.	An entity in the superclass is not required to be in a subclass.	The <code>EMPLOYEE</code> specializations, where an employee might not fit any defined subclass.
Attribute-Defined	The defining attribute is written next to the line from the superclass.	The value of this attribute determines subclass membership.	The <code>Job_type</code> attribute defining the specialization of <code>EMPLOYEE</code> .

Study Guide for the Enhanced Entity-Relationship (EER) Model

This guide is designed to review and test your understanding of the concepts presented in the Enhanced Entity-Relationship (EER) model, focusing on specialization, generalization, and related constraints as illustrated in the provided diagrams.

Short-Answer Quiz

Instructions: Answer the following ten questions in 2-3 complete sentences each, based on the concepts demonstrated in the source diagrams.

- 1.** What is the primary difference between the process of specialization and generalization in EER modeling?

- 2.** Explain the relationship between a superclass and a subclass, using the `EMPLOYEE` and `SECRETARY` entities from the diagrams as an example.

- 3.** What does the concept of attribute inheritance mean in an EER model?

- 4.** Describe what a disjointness constraint on specialization signifies. How is this represented in the EER diagrams provided?

CS 340 Summary

5. What is an overlapping constraint on specialization, and what does it allow?
6. Using Figure 4.4, explain how an attribute-defined specialization works.
7. In the diagram featuring the PERSON superclass, is it possible for one person to be classified as both a STUDENT and an ALUMNUS? Explain your reasoning.
8. According to Figure 4.3, what is the advantage of generalizing the CAR and TRUCK entity types into a single VEHICLE superclass?
9. In Figure 4.1, the EMPLOYEE entity has three different specializations. List these three specializations.
10. If a MANAGER is a subclass of EMPLOYEE, what attributes would an entity in the MANAGER subclass possess?

Essay Questions

Instructions: Prepare a detailed response for each of the following essay questions. Answers are not provided.

1. Compare and contrast the top-down design process of specialization with the bottom-up process of generalization. Provide detailed scenarios where a database designer would choose one approach over the other, referencing the diagrams for examples.

CS 340 Summary

- 2.** Explain the concepts of disjointness/overlapping constraints and completeness (total/partial) constraints in specialization hierarchies. Using the provided diagrams as a foundation, discuss how the application of these different constraints fundamentally changes the business rules and the data that can be represented in the database.
- 3.** Discuss the concept of attribute inheritance and its implications for database design, particularly in complex specialization lattices where a subclass can have multiple superclasses. How does inheritance simplify or complicate data querying and maintenance?
- 4.** Analyze the three distinct specializations of the EMPLOYEE superclass shown in Figure 4.1. Explain why a designer might choose to model these as separate specializations rather than one large, complex specialization. Discuss the potential for conflicts or ambiguities if an employee can be, for example, both an ENGINEER and a MANAGER.
- 5.** Describe the process of mapping an EER model that includes specialization and generalization to a relational database schema. Discuss at least two distinct mapping strategies, detailing their respective advantages and disadvantages in terms of data integrity, query performance, and storage efficiency.

Answer Key (Short)

- 1. What is the primary difference between the process of specialization and generalization in EER modeling?** Specialization is a top-down process where a general entity type (superclass) is broken down into more specific subgroupings (subclasses). Generalization is the reverse, a bottom-up process where the common features of multiple entity types are combined to form a more general superclass.
- 2. Explain the relationship between a superclass and a subclass, using the EMPLOYEE and SECRETARY entities from the diagrams as an example.** A superclass is a general entity type, while a subclass is a specialized grouping of entities from that superclass. In the example, EMPLOYEE is the superclass, and SECRETARY is one of its subclasses, meaning every secretary is also an employee and shares the general attributes of an employee.
- 3. What does the concept of attribute inheritance mean in an EER model?** Attribute inheritance is the principle that an entity in a subclass automatically includes all the attributes

and relationships of its superclass. For instance, a TECHNICIAN entity inherits attributes like Name, Ssn, and Birth_date from the EMPLOYEE superclass, in addition to its own specific attribute, Tgrade.

4. Describe what a disjointness constraint on specialization signifies. How is this represented in the EER diagrams provided? A disjointness constraint specifies that the subclasses of a specialization are mutually exclusive, meaning an entity can belong to at most one of the subclasses. This is represented in the diagrams by a circle containing the letter 'd', as seen in the specialization of EMPLOYEE into job types.

5. What is an overlapping constraint on specialization, and what does it allow? An overlapping constraint allows an entity from the superclass to be a member of more than one subclass simultaneously. This is useful when the subclass categories are not mutually exclusive. It is represented by a circle containing the letter 'o'.

6. Using Figure 4.4, explain how an attribute-defined specialization works. In an attribute-defined specialization, membership in a subclass is determined by the value of a specific attribute in the superclass. In Figure 4.4, the Job_type attribute of the EMPLOYEE entity explicitly determines whether an employee belongs to the SECRETARY, TECHNICIAN, or ENGINEER subclass.

7. In the diagram featuring the PERSON superclass, is it possible for one person to be classified as both a STUDENT and an ALUMNUS? Explain your reasoning. Yes, it is possible for a person to be both a STUDENT and an ALUMNUS. The diagram shows an overlapping constraint (an 'o' in the circle) for this specialization, which means the subclasses EMPLOYEE, ALUMNUS, and STUDENT are not mutually exclusive.

8. According to Figure 4.3, what is the advantage of generalizing the CAR and TRUCK entity types into a single VEHICLE superclass? The primary advantage is the reduction of redundancy and the creation of a cleaner data model. By moving common attributes such as Vehicle_id, Price, and License_plate_no into the VEHICLE superclass, these attributes only need to be defined once.

9. In Figure 4.1, the EMPLOYEE entity has three different specializations. List these three specializations. The three specializations of the EMPLOYEE entity shown are: {SECRETARY, TECHNICIAN, ENGINEER}, {MANAGER}, and {HOURLY_EMPLOYEE, SALARIED_EMPLOYEE}. Each specialization groups employees based on a different criterion.

10. If a MANAGER is a subclass of EMPLOYEE, what attributes would an entity in the MANAGER subclass possess? A MANAGER entity would inherit all the attributes of the EMPLOYEE superclass, such as Name, Ssn, Birth_date, and Address. In addition, it would participate in any relationships specific to the MANAGER subclass, such as the MANAGES relationship with PROJECT.

Glossary of Key Terms

Term	Definition
Attribute Inheritance	The property by which entities in a subclass inherit all attributes and relationship participations of their superclass.
Attribute-Defined Specialization	A specialization where membership in a subclass is determined by the value of a specific attribute of the superclass (the defining attribute).
Disjointness Constraint	A constraint on a specialization specifying that its subclasses are mutually exclusive. An entity of the superclass can be a member of at most one of the subclasses.
Enhanced Entity-Relationship (EER) Model	A conceptual data model that extends the basic ER model to include the concepts of superclasses, subclasses, specialization, and generalization for more accurate modeling.
Generalization	A bottom-up modeling process where the common attributes and relationships of multiple entity types are grouped to form a more general superclass.
Overlapping Constraint	A constraint on a specialization specifying that its subclasses are not mutually exclusive. An entity of the superclass may be a member of more than one subclass.
Specialization	A top-down modeling process where an entity type is divided into more specific subgroupings, or subclasses, based on distinguishing characteristics.
Subclass	A meaningful subgrouping of entities from a superclass that may have its own specific attributes and relationships. Every member of a subclass is also a member of the superclass.
Superclass	A general entity type that includes one or more distinct subgroupings of its entities, known as subclasses.

Chapter 4

Briefing: The Relational Database Model

Executive Summary

The relational model is a foundational concept in database management, structuring data into tables known as relations. This model, grounded in set theory, organizes information into rows (tuples) and columns (attributes), creating a simple yet powerful paradigm for data storage and retrieval. Each relation is defined by a schema, which specifies its name and attributes. The integrity and consistency of the data are maintained through a system of constraints.

Key to this model is the concept of keys. A primary key is a designated attribute or set of attributes that uniquely identifies each tuple within a relation, ensuring no duplicates and that the key value is never null (Entity Integrity). Data is defined, manipulated, and queried using Structured Query Language (SQL), which is divided into subsets like the Data Definition Language (DDL) for creating and modifying the database structure, and the Data Manipulation Language (DML) for retrieving and updating the data itself.

I. Core Concepts of the Relational Model

The relational model represents a database as a collection of relations. In practical terms, a relation can be visualized as a table. This model is built upon a formal mathematical foundation, which provides a clear and unambiguous framework for data management.

Key Terminology

The model uses specific terminology, which often corresponds to more common, informal terms.

Formal Relational Term	Informal/Common Term	Description
Relation	Table	A two-dimensional structure containing data about a specific entity.
Tuple	Row or Record	A single entry in a relation, representing a set of related data values.
Attribute	Column or Field	A named property of a relation that describes a piece of data for each tuple.

Domain	-	A set of permissible, atomic values for one or more attributes.
Relation Schema	Table Definition	The name of the relation and its set of attributes (e.g., STUDENT (Name, Ssn, Address)).
Relation State	Table Instance	The set of tuples present in a relation at a specific moment in time.

Example: The STUDENT Relation

The STUDENT relation serves as a clear example of these core concepts.

Relation Name: STUDENT

Attributes (Columns): Name, Ssn, Home_phone, Address, Office_phone, Age, Gpa

Tuples (Rows): Each row represents a specific student, such as ‘Benjamin Bayer’ or ‘Chung-cha Kim’.

Name	Ssn	Home_phone	Address	Office_phone	Age	Gpa
Benjamin Bayer	305-61- 2435	373-1616	2918 Bluebonnet Lane	NULL	19	3.21
Chung-cha Kim	381-62- 1245	375-4409	125 Kirby Road	NULL	18	2.89
Dick Davidson	422-11- 2320	NULL	3452 Elgin Road	749-1253	25	3.53
Rohan Panchal	489-22- 1100	376-9821	265 Lark Lane	749-6492	28	3.93
Barbara Benson	533-69- 1238	839-8461	7384 Fontana Lane	NULL	19	3.25

II. Fundamental Properties of Relations

Relations in this model have distinct characteristics derived from set theory.

- **Ordering of Tuples is Irrelevant:** The tuples (rows) in a relation are not considered to be ordered. A relation is a set of tuples, and in set theory, the order of elements is not significant. Therefore, rearranging the rows does not change the relation itself.

- **Ordering of Attributes is Irrelevant:** Similar to tuples, the attributes (columns) in a relation schema are also considered a set. The order in which columns are defined or displayed has no bearing on the meaning of the data.
- **Atomic Values:** Each cell in a relation, at the intersection of a row and column, must contain a single, atomic (indivisible) value.
- **Distinct Tuples:** Every tuple in a relation must be distinct. Because a relation is a set of tuples, duplicate entries are not permitted. This is enforced by the key constraints.

III. Keys and Integrity Constraints

Constraints are rules that ensure the validity and consistency of data within the database. Keys are a primary mechanism for enforcing these constraints.

Types of Keys

- **Superkey:** A set of one or more attributes that, taken collectively, allows for the unique identification of a tuple within a relation.
- **Candidate Key:** A minimal superkey. This means it is a superkey from which no attribute can be removed without it losing its uniqueness property. A relation may have multiple candidate keys.
- **Primary Key:** One candidate key that is chosen by the database designer to be the principal means of identifying tuples. Its value is used to uniquely reference each record. Primary keys are conventionally underlined in a schema definition.

Example: The CAR Relation The `CAR` relation has two candidate keys: `License_number` and `Engine_serial_number`. Either could be chosen as the primary key.

License_number	Engine_serial_number	Make	Model	Year
Texas ABC-739	A69352	Ford	Mustang	02
Florida TVP-347	B43696	Oldsmobile	Cutlass	05
New York MPO-22	X83554	Oldsmobile	Delta	01
California 432-TFY	C43742	Mercedes	190-D	99
California RSK-629	Y82935	Toyota	Camry	04
Texas RSK-629	U028365	Jaguar	XJS	04

Core Integrity Constraints

- **Key Constraint:** This fundamental constraint states that no two tuples in a relation can have the same value for their primary key attributes. This ensures that every tuple is unique.
- **Entity Integrity Constraint:** This constraint mandates that no primary key value can be `NULL`. This is essential because a `NULL` value would imply that there is a tuple that cannot be uniquely identified.

IV. Introduction to Structured Query Language (SQL)

SQL is the standard language for interacting with relational databases. It includes sub-languages for defining the database structure and for manipulating the data stored within it.

Data Definition Language (DDL)

DDL commands are used to create, modify, and delete database structures like tables and indexes.

- **CREATE TABLE:** Defines a new relation (table), specifying its attributes and their data types, as well as constraints like `PRIMARY KEY` and `NOT NULL`.
- **ALTER TABLE:** Modifies an existing table's structure, such as adding a new column.
- **DROP TABLE:** Permanently removes a table and all its data from the database.
- **CREATE INDEX:** Creates an index on one or more columns to improve query performance.

Data Manipulation Language (DML) - The `SELECT` Statement

The `SELECT` statement is the primary tool for retrieving data from one or more tables.

- **Selecting Specific Columns:** A query can retrieve a subset of columns from a table.
 - **Query:**
 - **Result:** | CUSTOMER_ID | NAME | CREDIT_LIMIT | | :--- | :--- | :--- | | 35 | Kimberly-Clark | 400 | | 36 | Hartford Financial Services Group | 400 | | 38 | Kraft Heinz | 500 | | 40 | Fluor | 500 | | 41 | AECOM | 500 | | ... | ... | ... |
- **Joining Multiple Tables:** The `INNER JOIN` clause combines rows from two or more tables based on a related column between them.
 - **Query:**

- **Result:** This query combines data from the `orders` and `order_items` tables, matching rows where the `order_id` is the same, and presents a comprehensive view of each order item alongside its parent order's details.

Relational Model Study Guide

Quiz: Test Your Knowledge

Answer each of the following questions in two to three complete sentences, based on the provided source materials.

- 1.** In the context of the relational model, what is a “relation” and what is a “tuple”?

- 2.** Based on the `STUDENT` relation example, explain what an “attribute” is and provide two examples.

- 3.** What is the purpose of a primary key in a database table?

- 4.** Differentiate between a candidate key and a superkey.

- 5.** What does the `NULL` value signify within a relation, as seen in the `STUDENT` table?
- 6.** What is the function of Data Definition Language (DDL) in SQL? Provide an example of a DDL command.

- 7.** Explain the purpose of the `ALTER TABLE` command and how it might be used.

- 8.** Describe what the `DROP TABLE` command does.

9. How does an INNER JOIN clause function in a SELECT statement?

10. According to the principles of the relational model, why is the order of tuples within a relation considered irrelevant?

Essay Questions

The following questions are designed for longer, more detailed responses. Formulate your answers by synthesizing information from across the source materials.

1. Discuss the fundamental properties that define a relation in the relational database model. Using the `STUDENT` table as a reference, explain why characteristics like unordered tuples, atomic attribute values, and the use of `NULL` are crucial for maintaining data consistency and integrity.

2. Explain the hierarchy and relationship between a superkey, a candidate key, and a primary key. Using the `CAR` relation as an example, describe why a relation might have multiple candidate keys and outline the considerations for choosing one as the primary key.

3. Describe the lifecycle of a database table using the three main Data Definition Language (DDL) commands shown in the source (`CREATE`, `ALTER`, `DROP`). Explain the specific syntax and purpose of each command, referencing the `persons` and `members` table examples.

4. Analyze the structure of the provided `SELECT` statement that joins the `orders` and `order_items` tables. Break down the function of each clause (`SELECT * , FROM, INNER JOIN ... ON, ORDER BY ... DESC`) and explain how they work in sequence to produce the final, sorted result set.

5. Define what a “relation schema” is and contrast it with a “relation instance.” Use the `CREATE TABLE persons (...)` command to illustrate the components of a schema and use the `STUDENT` table data to illustrate the concept of an instance.

Answer Key (Short)

1. A relation is a two-dimensional table used to store data, such as the `STUDENT` table. A tuple is a single row within that table, representing a single record or data entry, like the complete record for “Benjamin Bayer”.
2. An attribute is a named column in a relation that describes a specific property of each tuple. In the `STUDENT` relation, examples of attributes include `Name`, which stores the student’s name, and `Ssn`, which stores their social security number.
3. A primary key is a designated candidate key whose value is used to uniquely identify each tuple within a relation. For instance, the `person_id` is specified as the primary key for the `persons` table, ensuring no two people have the same ID.
4. A superkey is a set of one or more attributes that, when taken together, uniquely identifies a tuple in a relation. A candidate key is a minimal superkey, meaning it is a superkey from which no attributes can be removed without it losing its uniqueness property.
5. The `NULL` value indicates that the value for a specific attribute in a particular tuple is either unknown or not applicable. In the `STUDENT` table, some tuples have a `NULL` value for `Office_phone`, meaning that student’s office phone number is not available in the database.
6. Data Definition Language (DDL) is a set of SQL commands used to define and manage the structure of database objects like tables and indexes. An example is the `CREATE TABLE` command, which is used to build a new table and define its attributes and their data types.
7. The `ALTER TABLE` command is used to modify the structure of an existing table. For example, the command `ALTER TABLE members ADD birth_date DATE NOT NULL;` would add a new column named `birth_date` with a `DATE` data type to the `members` table.
8. The `DROP TABLE` command is a DDL command used to permanently delete an entire table, including its structure, data, and associated indexes from the database. The command `DROP TABLE persons;` would remove the `persons` table completely.
9. An `INNER JOIN` clause is used in a `SELECT` statement to combine rows from two or more tables based on a related column between them. For example, it can be used to combine the `orders` and `order_items` tables where the `order_id` in both tables match, creating a result set with data from both.
10. The order of tuples is irrelevant because a relation is defined as a set of tuples. Mathematical set theory does not impose any order on the elements within a set, so queries should not depend on the sequence of rows in a table.

Glossary of Key Terms

Term	Definition
ALTER TABLE	A DDL command used to modify the structure of an existing database table, such as adding a new column.
Attribute	A named column of a relation that corresponds to a specific property of the data being stored (e.g., Name, Ssn, Age).
Candidate Key	A minimal superkey. It is an attribute or set of attributes that uniquely identifies a tuple, and from which no attribute can be removed without losing the uniqueness property. A relation can have multiple candidate keys (e.g., License_number and Engine_serial_number in the CAR table).
CREATE INDEX	A DDL command that creates an index on a table to improve the speed of data retrieval operations on a database.
CREATE TABLE	A DDL command used to create a new table in a database, defining its name, attributes, and data types.
Data Definition Language (DDL)	A subset of SQL commands used for creating, modifying, and deleting database structures, but not the data itself. Examples include CREATE TABLE, ALTER TABLE, and DROP TABLE.
DROP TABLE	A DDL command that permanently deletes a table and all its data from the database.
INNER JOIN	A SQL clause used to combine rows from two or more tables based on a matching value in a common column.
NULL	A special marker used to indicate that a data value does not exist in the database for a specific attribute in a tuple. It represents a missing or inapplicable value.
Primary Key	The candidate key that is chosen to be the main identifier for tuples within a relation. Its value must be unique for each tuple.
Relation	The fundamental structure in the relational model, represented as a two-dimensional table containing a set of tuples (rows).

CS 340 Summary

Relation Name	The name given to a relation (e.g., STUDENT, CAR).
SELECT	The primary SQL command for retrieving data from one or more database tables.
Superkey	An attribute or a set of attributes that, taken together, allows for the unique identification of a tuple within a relation.
Tuple	A single row in a relation, representing a set of related data values. Tuples in a relation are not ordered.

Chapter 6

Briefing Document: Fundamentals of Basic SQL

Executive Summary

This document provides a comprehensive synthesis of the fundamentals of Structured Query Language (SQL), the standard declarative language for interacting with relational database systems. The core focus is on the Data Manipulation Language (DML) for retrieving and modifying data, while also touching upon the Data Definition Language (DDL) for schema creation.

Key takeaways include the declarative nature of SQL, where users specify *what* data is required rather than *how* to retrieve it. A standard SQL query for data retrieval consists of mandatory `SELECT` and `FROM` clauses, with optional `WHERE` and `ORDER BY` clauses for filtering and sorting. SQL treats tables as multisets (or bags), meaning duplicate rows are permitted by default; the `DISTINCT` keyword is required to enforce uniqueness in results. The language also provides a suite of set operators (`UNION`, `UNION ALL`, `INTERSECT`, `EXCEPT`) for combining the results of multiple queries. Finally, data modification is handled through three primary commands: `INSERT` to add new data, `UPDATE` to modify existing data, and `DELETE` to remove data, all of which respect the integrity constraints defined in the database schema.

Introduction to the SQL Query Language

SQL serves as the standardized language for relational database management systems, bridging the gap between conceptual data models (like ER diagrams) and practical data interaction. It is fundamentally a declarative language, contrasting with procedural languages by focusing on defining the desired result rather than the step-by-step process to achieve it.

The language is formally divided into several components, with two being central to basic operations:

- **Data Definition Language (DDL):** Used to define and manage the database schema. Commands like `CREATE TABLE` fall into this category.
- **Data Manipulation Language (DML):** Used to interact with the data within the schema. This includes:
 - **Data Retrieval:** Primarily handled by the `SELECT` statement.

- **Data Modification:** Handled by `INSERT`, `UPDATE`, and `DELETE` statements.

Data Definition Language (DDL) Overview

The DDL is used to construct the database schema. The `CREATE TABLE` command defines a new table, its attributes with their corresponding data types, and critical integrity constraints.

Key Integrity Constraints in DDL

- **Primary Key:** `PRIMARY KEY` (attribute) uniquely identifies each record in a table.
- **Foreign Key:** `FOREIGN KEY` (attribute) `REFERENCES` table(attribute) establishes a link between tables and enforces referential integrity.
- **Referential Actions:** DDL allows for specifying actions to take when referenced data is modified:
 - `ON DELETE SET NULL`: Sets the foreign key value to NULL if the referenced primary key is deleted.
 - `ON UPDATE CASCADE`: Updates the foreign key value if the referenced primary key is changed.

Example: DDL for an EMPLOYEE Table

```
CREATE TABLE EMPLOYEE (
    name      VARCHAR(15),
    Ssn       CHAR(2),
    Super_ssn CHAR(2),
    PRIMARY KEY (Ssn),
    FOREIGN KEY (Super_ssn) REFERENCES EMPLOYEE(Ssn)
        ON DELETE SET NULL
        ON UPDATE CASCADE
);
```

Core Data Retrieval: The `SELECT` Statement

The fundamental query for data retrieval follows a standard structure. While a query can have up to four primary clauses, only `SELECT` and `FROM` are mandatory.

Basic Query Structure

```
SELECT <attribute list>
FROM <table list>
[ WHERE <condition> ]
[ ORDER BY <attribute list> ];
```

- **SELECT Clause:** Corresponds to the **projection** operation in relational algebra. It specifies the columns (attributes) to be included in the final result.
- **FROM Clause:** Specifies the source tables for the query. If multiple tables are listed, a cross-product of these tables is conceptually formed as the initial data set.
- **WHERE Clause:** Corresponds to the **selection** operation. It applies a conditional filter to the rows, retaining only those that satisfy the specified criteria.
- **ORDER BY Clause:** Sorts the final result set based on the values of one or more attributes.

Filtering and Selection with the WHERE Clause

The `WHERE` clause is a powerful tool for filtering data using a variety of operators and logical combinations.

Operator Type	Description	Example
Comparison	Standard operators like <code>></code> , <code><</code> , <code>=</code> , <code>!=</code> .	<code>WHERE DoB > 2012</code>
Pattern Matching	<code>LIKE</code> is used for string matching. <code>%</code> matches zero or more characters; <code>_</code> matches a single character.	<code>WHERE Name LIKE 'A%</code>
Range Queries	<code>BETWEEN</code> checks if a value falls within a specified range (inclusive).	<code>WHERE (Salary BETWEEN 30000 AND 40000)</code>
Date/Time	Special formats and functions are used to compare date and time values.	<code>WHERE date < to_date('1-1-1943', 'DD-MON-YYYY')</code>
Logical Operators	<code>AND</code> and <code>OR</code> are used to combine multiple conditions.	<code>WHERE (Salary > 30000) AND Dno = 5</code>

Selection Examples

Based on the STUDENT table below:

Sid	Name	DoB	Major
123	Ali	2010	SE
124	Saleh	2010	CS
125	Ahmd	2010	SE
126	Fahad	2013	CS
127	Omar	2012	CS
128	Mhmd	2015	IS
129	Khalid	2015	IS

- **Query:** `SELECT * FROM STUDENT WHERE DoB > 2012;`
 - **Result:** Returns students Fahad, Mhmd, and Khalid.
- **Query:** `SELECT * FROM STUDENT WHERE Name LIKE 'A%';`
 - **Result:** Returns students Ali and Ahmd.

Advanced Data Retrieval Operations

Handling Duplicates: Sets vs. Multisets

While the formal relational model defines relations as sets (with no duplicate tuples), SQL treats tables as **multisets** (or bags), where duplicates are allowed.

- `SELECT (default)`: Retains all rows, including duplicates. It is implicitly `SELECT ALL`.
- `SELECT DISTINCT`: Scans the results and eliminates any duplicate rows, returning a true set.

Example: A query joining `sailors` and `reserves` to find sailor names might return a name multiple times if that sailor made multiple reservations.

- `SELECT sname ...`: Returns “dustin” four times.
- `SELECT DISTINCT sname ...`: Returns “dustin” only once.

Renaming Attributes and Performing Calculations

- **Renaming:** The `AS` keyword can be used to assign a new name (alias) to a column in the result set.

CS 340 Summary

- **Example:** `SELECT DISTINCT sname AS sailor_name FROM sailors s, reserves r WHERE s.sid = r.sid;`
- **Arithmetic Operations:** Calculations can be performed directly within the `SELECT` clause.
 - **Example:** `SELECT Fname, Lname, 1.1 * E.Salary FROM EMPLOYEE;`

Sorting Results with ORDER BY

The `ORDER BY` clause sorts the final query output. It can sort by multiple columns and in ascending (`ASC`, the default) or descending (`DESC`) order.

- **Example:** `... ORDER BY D.Dname, E.Lname, E.Fname;` This query sorts results first by department name, then by employee last name, and finally by first name.

Set Operators for Combining Queries

SQL provides set operators to combine the result sets of two or more `SELECT` statements into a single result set.

- **UNION:** Combines the results and removes all duplicate rows.
- **UNION ALL:** Combines the results and retains all duplicate rows. This is generally more performant as it does not need to check for duplicates.
- **INTERSECT:** Returns only the rows that are common to both query result sets.
- **EXCEPT / MINUS:** Returns the unique rows from the first query that are not present in the second query's result set.

Data Manipulation Language (DML): Modification Commands

SQL provides three commands to modify the data stored in tables: `INSERT`, `DELETE`, and `UPDATE`.

INSERT Statement

The `INSERT` command adds new tuples (rows) to a table.

- **Basic Form:** Inserts a single tuple with values provided in the order of the table's columns.
 - `INSERT INTO sailors VALUES(22, 'dustin', 7, 45.0);`

- **Specifying Columns:** Allows inserting values for specific columns; unlisted columns will be assigned their default value or NULL.
 - `INSERT INTO sailors(sid, sname, rating) VALUES(22, 'dustin', 7);`
- **Inserting from a Query:** Inserts the entire result set of a `SELECT` query into a table.
 - `INSERT INTO WORKS_ON_INFO (...) SELECT E.Lname, P.Pname, W.Hours FROM ...;`

All `INSERT` operations must satisfy the table's integrity constraints (e.g., primary key uniqueness, valid foreign keys).

DELETE Statement

The `DELETE` command removes tuples from a table based on a condition.

- **Syntax:** `DELETE FROM <table> WHERE <condition>;`
- **Example:** `DELETE FROM EMPLOYEE WHERE Lname = 'Brown';`
- **Propagation:** If referential integrity constraints like `ON DELETE CASCADE` are defined, deleting a row in one table can trigger the automatic deletion of related rows in other tables.

UPDATE Statement

The `UPDATE` command modifies the attribute values of existing tuples.

- **Syntax:** `UPDATE <table> SET <attribute = new_value> WHERE <condition>;`
- **Example:** `UPDATE EMPLOYEE SET Salary = Salary * 1.1 WHERE Dno = 5;`
- **Propagation:** Similar to `DELETE`, an `UPDATE` to a primary key can propagate to foreign keys in other tables if `ON UPDATE CASCADE` is specified.

Conclusion and Next Steps

The commands and concepts outlined—including DDL for schema creation and the core DML commands `SELECT`, `INSERT`, `UPDATE`, and `DELETE`—form the foundation of interacting with relational databases. Mastery of basic queries, filtering, sorting, and data modification is essential for any database work.

CS 340 Summary

Future topics will build upon this foundation to explore more advanced SQL features, including complex retrieval queries, the use of views, the implementation of triggers and assertions, and commands for schema modification.

Study Guide for Basic SQL (CS 340)

Short-Answer Quiz

Instructions: Answer the following questions in 2-3 sentences each, based on the provided course material.

1. What are the two primary components of the SQL query language, and what is the main function of each?

2. What is the basic form of an SQL data retrieval query, and which clauses are mandatory?

3. Explain the difference between SQL being a “declarative” language versus a “procedural” one.

4. How does SQL’s treatment of tables differ from the strict relational model regarding duplicate data?

5. What is the purpose of the `DISTINCT` keyword in a `SELECT` clause, and how does its behavior compare to the default?

6. Describe the two wildcard characters used for pattern matching with the `LIKE` operator.

CS 340 Summary

7. What are the three SQL commands used to modify the data within a database?

8. Explain the function of the `AS` keyword within a `SELECT` statement.

9. What is the purpose of the `ORDER BY` clause?

10. Differentiate between the `UNION` and `UNION ALL` set operators.

Essay Questions

Instructions: The following questions are designed for longer-form answers. Prepare a detailed response for each, synthesizing concepts from the course material.

1. Discuss the relationship between the relational model's formal prohibition of duplicates and SQL's default behavior of treating tables as multisets. Explain how and why SQL allows for this difference and what tools it provides to manage duplicates in query results.

2. Describe the complete conceptual evaluation process of a multi-table `SELECT` query that includes a `WHERE` clause join condition. Use the "Sailors" and "Reserves" tables from the lecture example to illustrate the three main steps: the cross-product, the application of qualifications, and the final projection.

CS 340 Summary

3. Compare and contrast the four set operators discussed: UNION, UNION ALL, INTERSECT, and EXCEPT (MINUS). Explain the function of each operator and provide a practical example scenario where each would be the most appropriate choice.

4. Explain the distinct roles of Data Definition Language (DDL) and Data Manipulation Language (DML) within SQL. Provide examples of commands for each category and discuss how they work together in the lifecycle of creating, populating, and querying a relational database.

5. Explain how referential integrity constraints, such as ON DELETE SET NULL and ON UPDATE CASCADE, connect DDL commands like CREATE TABLE with DML commands like DELETE and UPDATE. Use the EMPLOYEE table example, with its self-referencing Super_ssn foreign key, to illustrate the potential cascading impact of these constraints.

Answer Key (Short)

1. The SQL query language consists of the Data Definition Language (DDL) and the Data Manipulation Language (DML). The DDL is used to define the database schema, with commands like CREATE TABLE. The DML is used for data retrieval (with SELECT) and data modification (with INSERT, UPDATE, DELETE).
2. The basic form of a data retrieval query is `SELECT <attribute list> FROM <table list> WHERE <condition>;`. Of these clauses, only the `SELECT` and `FROM` clauses are mandatory for a simple retrieval query. The `WHERE` clause is optional and used for filtering.
3. SQL is a declarative language, which means the user specifies *what* data they want to retrieve without specifying *how* to retrieve it. This contrasts with a procedural language, where the user would need to outline the specific steps or procedure to obtain the data.

CS 340 Summary

4. While the formal relational model does not allow for duplicate tuples, SQL usually treats a table as a multiset (or bag), where duplicate tuples can appear more than once. The user may want to see these duplicates in the result of a query.
5. The `DISTINCT` keyword is used in the `SELECT` clause to eliminate duplicate tuples from the result of an SQL query. The default behavior, `SELECT ALL` (which is implied if neither is specified), does not remove duplicates and returns all matching rows.
6. The `LIKE` operator uses the percent sign (%) and the underscore (_) for pattern matching. The % character replaces an arbitrary number of zero or more characters, while the _ character replaces a single character.
7. The three SQL commands used to modify a database are `INSERT`, `DELETE`, and `UPDATE`. `INSERT` adds new tuples (rows), `DELETE` removes existing tuples, and `UPDATE` modifies the attribute values of existing tuples.
8. The `AS` keyword is used to rename an attribute in the result set of a query. This allows the output column to have a different, often more descriptive, name than the original attribute name in the table.
9. The `ORDER BY` clause allows the user to sort the tuples in the result of a query. The sorting can be based on the values of one or more attributes that appear in the query result.
10. The `UNION` operator combines the result sets of two `SELECT` statements and returns only the unique rows. The `UNION ALL` operator also combines the result sets but retains all duplicate rows from both queries.

Glossary of Key Terms

Term	Definition
AS	A keyword used in a SELECT statement to rename an attribute (column) in the result set.
BETWEEN	A comparison operator used in a WHERE clause to select values within a given range.
Data Definition Language (DDL)	The part of SQL used to define the database schema. An example command is CREATE TABLE.
Data Manipulation Language (DML)	The part of SQL used for data retrieval and modification. Commands include SELECT, INSERT, UPDATE, and DELETE.
Declarative Language	A language where the user specifies what data they want to retrieve without specifying the procedure for how to retrieve it. SQL is declarative.
DELETE	A DML command that removes specified tuples from a relation.
DISTINCT	A keyword used in the SELECT clause to eliminate duplicate tuples from the result of a query.
EXCEPT (MINUS)	A set operator used to retrieve all unique records from the left query that are not present in the result set of the right query.
Foreign Key	A key used to link two tables together, which acts as a cross-reference between tables by referencing the primary key of another table.
FROM	A mandatory clause in a SELECT statement that specifies the table or tables from which to retrieve data.
INSERT	A DML command used to add a single tuple or the results of a query into a relation.
Integrity Constraints	Rules enforced on data columns, such as defining primary keys, candidate keys, and foreign keys.
INTERSECT	A set operator used to retrieve the records that are common or identical between the result sets of two or more queries.
JDBC	Java Database Connectivity, mentioned as a tool that will be used later to execute SQL queries from Java applications.
LIKE	An operator used in a WHERE clause for pattern matching in string values. It uses wildcard characters % and _.
Multiset (Bag)	A collection of objects where members are allowed to appear more than once. SQL tables are treated as multisets by default, allowing duplicate tuples.
ON DELETE CASCADE	A referential integrity action that specifies if a referenced row is deleted, all referencing rows should also be deleted.

CS 340 Summary

ON DELETE SET NULL	A referential integrity action specifying that when a referenced row is deleted, the foreign key values in referencing rows are set to NULL.
ON UPDATE CASCADE	A referential integrity action specifying that if a referenced primary key is updated, all corresponding foreign key values should also be updated.
ORDER BY	An optional clause used to sort the tuples in the result of a query by the values of one or more attributes.
Primary Key	An integrity constraint that uniquely identifies each record in a table.
Referential Integrity	A constraint ensuring that a value that appears in one relation for a given set of attributes also appears for a certain set of attributes in another relation.
SELECT	The DML command that begins a data retrieval query. It specifies the attributes (columns) to be returned.
Set Operators	Operators like UNION, UNION ALL, INTERSECT, and EXCEPT that combine the result sets of two or more SELECT statements.
SQL	Structured Query Language, the standard declarative language for managing and querying relational databases.
UNION	A set operator that combines the result sets of two or more SELECT statements into a single result set, removing any duplicate rows.
UNION ALL	A set operator that combines the result sets of two or more SELECT statements, retaining all duplicate rows.
UPDATE	A DML command used to modify the attribute values of one or more selected tuples in a relation.
WHERE	An optional clause in a query used to specify a condition to filter tuples. Only tuples that satisfy the condition are included in the result.

Chapter 7

Advanced SQL Concepts and Implementation

Executive Summary

This document provides a comprehensive overview of advanced SQL features, synthesizing key concepts from the CS 340 course on Introduction to Database Systems. The analysis covers complex query construction, the nuanced handling of `NULL` values through three-valued logic, and the powerful capabilities of various `JOIN` operations, particularly `OUTER JOINS`, which preserve non-matching data. It details the use of aggregation functions with the `GROUP BY` clause to perform calculations on data subsets. Furthermore, the document explores the role of Views as virtual tables for simplifying queries and enforcing security, distinguishing between non-materialized and materialized views for performance considerations. Key commands for schema modification, `DROP` and `ALTER`, are defined, and the practical integration of SQL within application code is demonstrated through APIs like JDBC for Java and DB-API for Python.

1. Handling NULL Values and Three-Valued Logic

The concept of `NULL` in SQL extends beyond simple absence of data, introducing a more complex logical framework for query evaluation.

Interpretations of NULL

A `NULL` value can represent several scenarios:

- **Unknown Value:** The value exists but is not known (e.g., an unknown date of birth).
- **Unavailable or Withheld Value:** The value is known but intentionally not provided (e.g., an unlisted phone number).
- **Not Applicable Attribute:** The attribute does not apply to the specific entity (e.g., the `LastCollegeDegree` for a person with no degrees).

Three-Valued Logic in SQL

Instead of standard two-valued Boolean logic (`TRUE`, `FALSE`), SQL employs a three-valued logic system:

- **TRUE**
- **FALSE**
- **UNKNOWN**

This third value, `UNKNOWN`, arises from comparisons involving `NULL` values. The results of logical operations `AND`, `OR`, and `NOT` are defined in the following truth table:

Operator		TRUE	FALSE	UNKNOWN
AND	TRUE	TRUE	FALSE	UNKNOWN

	FALSE	FALSE	FALSE	FALSE
	UNKNOWN	UNKNOWN	FALSE	UNKNOWN
OR	TRUE	TRUE	TRUE	TRUE
	FALSE	TRUE	FALSE	UNKNOWN
	UNKNOWN	TRUE	UNKNOWN	UNKNOWN
NOT	TRUE	FALSE		
	FALSE	TRUE		
	UNKNOWN	UNKNOWN		

Querying for NULL Values

To properly handle `NULLS` in queries, SQL provides specific operators, as standard comparison operators like `=` or `<>` are insufficient.

- **`IS NULL / IS NOT NULL`:** These are the correct operators to check if an attribute's value is `NULL`.
- **Distinctness:** SQL considers each `NULL` value to be distinct from every other `NULL` value.

2. Complex Query Construction

Advanced SQL allows for the creation of sophisticated queries through nesting and a variety of join operations that can combine data from multiple tables in flexible ways.

Nested Queries (Subqueries)

A nested query, or subquery, is a `SELECT` statement embedded inside another SQL statement. This allows for complex, multi-step data retrieval to be performed in a single query.

Example: To find all workshops offered by the same presenter as “Robot Operating System,” a nested query can be used to first find the presenter and then use that result to find all their workshops.

```
SELECT *
FROM Workshop
WHERE Presenter = (SELECT Presenter
                    FROM Workshop
                    WHERE Title = 'Robot Operating System');
```

JOIN Operations

Joins are fundamental to relational databases, enabling the combination of rows from two or more tables based on a related column.

INNER JOIN

The standard `JOIN` (or `INNER JOIN`) returns only the rows where the join condition is met in both tables. Tuples from one table that do not have a matching tuple in the other are excluded from the result.

OUTER JOINs

Outer joins are used to retrieve matching rows *and* preserve rows from one or both tables that do not have a match in the other table. Missing values from the non-matching side are filled with `NULL`.

- **LEFT OUTER JOIN:** All tuples from the left table are included in the result. If a tuple has no match in the right table, the attributes from the right table are padded with `NULL`.
- **RIGHT OUTER JOIN:** All tuples from the right table are included in the result. If a tuple has no match in the left table, the attributes from the left table are padded with `NULL`.
- **FULL OUTER JOIN:** All tuples from both the left and right tables are included. Where matches do not exist, the corresponding attributes are padded with `NULL`.

Note on Syntax: The keyword `OUTER` is optional. `LEFT JOIN` is identical to `LEFT OUTER JOIN`, `RIGHT JOIN` to `RIGHT OUTER JOIN`, and `FULL JOIN` to `FULL OUTER JOIN`.

NATURAL JOIN

A `NATURAL JOIN` is a specialized join where the join condition is implicitly created on all pairs of attributes that share the same name across the tables.

- No join condition is explicitly specified.
- The result includes only one column for each pair of identically-named columns.
- This can be combined with outer joins (e.g., `NATURAL LEFT OUTER JOIN`).

CROSS JOIN

The `CROSS JOIN` operation produces the Cartesian product of the tables involved, combining every row from the first table with every row from the second table. This operation should be used with caution as it can generate very large result sets.

Multiway JOIN

SQL allows for the joining of three or more tables in a single statement by nesting `JOIN` specifications. This is known as a multiway join.

Example:

```
SELECT Pnumber, Dnum, Lname, Address, Bdate
FROM ((PROJECT JOIN DEPARTMENT ON Dnum = Dnumber)
      JOIN EMPLOYEE ON Mgr_ssn = Ssn)
WHERE Plocation = 'Stafford';
```

3. Data Aggregation and Grouping

SQL provides aggregate functions to perform calculations on a set of values and return a single summary value.

Aggregate Functions

The primary SQL aggregate functions are:

- **COUNT([DISTINCT] A)**: Returns the number of (unique) values.
- **SUM([DISTINCT] A)**: Returns the sum of (unique) values.
- **AVG([DISTINCT] A)**: Returns the average of (unique) values.
- **MAX(A)**: Returns the maximum value.
- **MIN(A)**: Returns the minimum value.

The GROUP BY Clause

The `GROUP BY` clause is used with aggregate functions to partition the result set into groups based on the values in one or more columns. The aggregate function is then applied to each group.

Example: Retrieve the number of employees and the average salary for each department.

```
SELECT Dno, COUNT(*), AVG(Salary)
FROM EMPLOYEE
GROUP BY Dno;
```

Usage and Restrictions

- When using `GROUP BY`, the `SELECT` list can only include the grouping attributes and aggregate functions.
- Grouping can be performed on multiple attributes to create nested sub-groups.
- If `NULL` values exist in the grouping attribute, all tuples with a `NULL` value are collected into a separate group.

4. Views: Virtual Tables for Simplification and Security

A view is a virtual table based on the result-set of an SQL statement. It provides a powerful mechanism for abstracting complexity and controlling data access.

Core Purpose of Views

- **Simplify Complex Queries:** A complex query involving multiple joins and conditions can be saved as a view, which can then be queried like a simple table.
- **Enhance Security:** Views can be used to restrict data access. Users can be granted permission to query a view that exposes only specific columns or rows, effectively hiding sensitive data.
- **Present Data in a Specific Format:** Views can reformat or rename columns to present data in a more user-friendly way.

Materialized vs. Non-Materialized Views

- **Non-Materialized Views (Standard Views):** These do not store data physically. The underlying SQL query is executed every time the view is accessed, ensuring the data is always current. They can be inefficient for complex queries.
- **Materialized Views:** These store the query result as a physical table. This can significantly improve performance for complex queries on large datasets, but the data must be periodically refreshed to reflect changes in the underlying tables.

View Updatability and Access Control

- **View Updates:** Not all views are updatable. An update to a view must unambiguously map to an update on its underlying base table. Views based on complex transformations like joins or aggregations are generally not updatable.
- **Access Control:** Views are a key tool for database security. They allow administrators to create a public interface to the data while restricting access to the underlying schema, thus enhancing data integrity and security.

5. Database Schema Modification

SQL provides commands to alter the database schema after its initial creation.

- **DROP:** This command is used to permanently remove database objects. It can delete tables (`DROP TABLE tableName`), entire databases (`DROP DATABASE databaseName`), views, or indexes. This action is irreversible.
- **ALTER:** This command modifies the structure of an existing database object, most commonly a table. Its functions include:
 - Adding a column: `ALTER TABLE tableName ADD columnName dataType;`
 - Modifying a column: `ALTER TABLE tableName MODIFY COLUMN columnName newDataType;`
 - Dropping a column: `ALTER TABLE tableName DROP COLUMN columnName;`

6. Integrating SQL with Application Code

To build applications, SQL queries must be executed from within a programming language. This is accomplished using database connectivity APIs.

Key APIs: JDBC and DB-API

- **JDBC (Java Database Connectivity):** An API for the Java programming language that defines how a client application can access a database.
- **DB-API (Database Application Programming Interface):** A standard interface for Python applications to interact with relational databases. Each database requires its own DB-API compliant driver.

Connection and Execution Flow

The general process for executing queries from an application, as demonstrated by the provided Java JDBC example, is:

1. **Register the JDBC Driver:** Load the specific driver for the target database.
2. **Open a Connection:** Establish a connection to the database using credentials (URL, user, password).
3. **Create a Statement:** Create a `Statement` object to execute SQL commands.
4. **Execute a Query:** Run the SQL query using a method like `executeQuery()`. This returns a `ResultSet` object.
5. **Process the Result Set:** Iterate through the `ResultSet` to extract the data row by row.
6. **Close Resources:** Close the `ResultSet`, `Statement`, and `Connection` objects in a `finally` block to release database resources.

Advanced SQL Study Guide

Short-Answer Quiz

Instructions: Please answer the following questions in 2-3 sentences each, based on the provided course materials.

1. Describe the three primary interpretations of a `NULL` value in a database.
2. What is the purpose of a nested query (or subquery) in SQL? Provide a conceptual example of its use.
3. Explain the fundamental difference between an `INNER JOIN` and a `LEFT OUTER JOIN`.
4. What is a `NATURAL JOIN`, and how does its join condition differ from a standard `JOIN`?
5. What is the role of the `GROUP BY` clause when used with SQL aggregation functions?
6. List the five main SQL aggregation functions and the purpose of the optional `DISTINCT` keyword.
7. What is a view in SQL, and what are two primary purposes for using one?

8. Differentiate between a materialized view and a non-materialized view.
9. Explain the difference between the `DROP TABLE` and `ALTER TABLE` commands in SQL.
10. What is SQL's three-valued logic system, and what are the three values it uses?

Essay Questions

1. Discuss the role of views as an access control mechanism in a Database Management System (DBMS). Explain how a view can be created to hide sensitive data and how this enhances both security and data integrity.
2. Compare and contrast the three types of outer joins: `LEFT OUTER JOIN`, `RIGHT OUTER JOIN`, and `FULL OUTER JOIN`. For each join type, provide a conceptual scenario where it would be the most effective choice and explain why.
3. Describe the behavior of `NULL` values in SQL, particularly in the context of three-valued logic. Using the provided truth tables, explain how the logical operators `AND`, `OR`, and `NOT` evaluate expressions involving the `UNKNOWN` value.
4. Explain the concept of aggregation in SQL. Detail the purpose of the `COUNT`, `AVG`, and `MAX` functions, and then describe how the `GROUP BY` clause allows for more granular analysis. What restrictions are placed on the `SELECT` list when a `GROUP BY` clause is present?

5. Imagine you are a database administrator tasked with evolving a database schema over time without losing critical data. Describe how you would use the `ALTER` command to add a new column, modify an existing column's data type, and remove an unnecessary column. Contrast this with the `DROP` command and explain the critical importance of using `DROP` with caution.

Answer Key (Short)

1. A `NULL` value can have three interpretations. It can represent an “unknown value,” such as a date of birth that is not known. It can also mean an “unavailable or withheld value,” like a phone number a person does not wish to list. Finally, it can signify a “not applicable attribute,” for example, a `LastCollegeDegree` for a person who has never attended college.
2. A nested query, or subquery, is a query embedded within another SQL query. Its purpose is to perform a multi-step data retrieval process in a single statement. For example, one can find all workshops presented by the same person who teaches ‘Robot Operating System’ by nesting the query that finds the presenter’s name inside the query that selects the workshops.
3. An `INNER JOIN` returns only the rows where the join condition is met in both tables, meaning there is a match. A `LEFT OUTER JOIN` returns all rows from the left table, regardless of whether there is a match in the right table; if no match is found, the columns from the right table are filled with `NULL` values.
4. A `NATURAL JOIN` is a type of join where the join condition is not explicitly specified. Instead, it implicitly joins tables based on all pairs of attributes that share the same name. In the result, each pair of matching attributes is included only once.
5. The `GROUP BY` clause is used to partition data into groups based on the values in one or more columns. When used with aggregation functions like `COUNT()` or `AVG()`, the function is applied to each group individually rather than to the entire table. This allows for calculations on subsets of data, such as finding the average salary for each department.
6. The five main aggregation functions are `COUNT` (number of values), `SUM` (sum of values), `AVG` (average of values), `MAX` (maximum value), and `MIN` (minimum value). The optional `DISTINCT` keyword, when used with `COUNT`, `SUM`, or `AVG`, ensures that the function operates only on unique values within the specified column.
7. A view is a virtual table based on the result-set of an SQL statement. Two primary purposes are simplifying complex queries by encapsulating joins and calculations, and ensuring data security by restricting user access to only specific rows or non-sensitive columns of a table.
8. A non-materialized view does not store its data physically; it runs the underlying query each time it is accessed, always showing the most current data. A materialized view stores the query result as a physical table, which must be refreshed periodically but can significantly improve performance for complex queries on large datasets.

CS 340 Summary

9. The `DROP TABLE` command is used to permanently delete an entire table, including its structure and all of its data. The `ALTER TABLE` command is used to modify the structure of an existing table, such as by adding, deleting, or modifying columns and their data types.
10. SQL uses a three-valued logic system instead of standard two-valued Boolean logic. The three logical values it uses are `TRUE`, `FALSE`, and `UNKNOWN`. The `UNKNOWN` value typically arises from comparisons involving `NULL` values.

Glossary of Key Terms

Term	Definition
Aggregation Function	A function that operates on a set of values to return a single summary value. SQL's primary aggregation functions are <code>COUNT</code> , <code>SUM</code> , <code>AVG</code> , <code>MAX</code> , and <code>MIN</code> .
ALTER	An SQL command used to modify the structure of an existing database object, such as a table, by adding, deleting, or modifying columns or constraints.
CROSS JOIN	An SQL join operation that produces the Cartesian product of the rows from the tables in the join.
DB-API	The Database Application Programming Interface for Python. It provides a standard interface for Python applications to interact with relational databases.
DROP	An SQL command used to permanently remove database objects such as tables, views, or entire databases. This action is irreversible.
FULL OUTER JOIN	A join that returns all rows when there is a match in either the left or the right table. If there is no match, the result is <code>NULL</code> on the side that does not have a match. Also written as <code>FULL JOIN</code> .
GROUP BY	An SQL clause used with aggregate functions to group rows that have the same values in specified columns into summary rows.
INNER JOIN	A join that returns records that have matching values in both tables. This is the default join type in many database systems. Also written as <code>JOIN</code> .
JDBC	Java Database Connectivity. An application programming interface (API) for Java that defines how a client may access a database.
LEFT OUTER JOIN	A join that returns all records from the left table, and the matched records from the right table. The result is <code>NULL</code> from the right side if there is no match. Also written as <code>LEFT JOIN</code> .
Materialized View	A view that stores its query result as a physical table. It offers better performance for complex queries but must be periodically refreshed to reflect changes in the underlying data.
Multiway Join	A join operation involving three or more tables, often specified by nesting join operations.

CS 340 Summary

NATURAL JOIN	A type of join that creates an implicit join clause based on all common columns in the two tables being joined. Common columns appear only once in the result set.
Nested Query	A query that is embedded inside another SQL query. Also known as a subquery.
Non-Materialized View	A standard view that does not store data physically. The underlying query is re-executed every time the view is accessed, ensuring the data is always current.
NULL	A special marker in SQL used to indicate that a data value does not exist in the database. It can mean the value is unknown, unavailable, or not applicable.
RIGHT OUTER JOIN	A join that returns all records from the right table, and the matched records from the left table. The result is <code>NULL</code> from the left side if there is no match. Also written as <code>RIGHT JOIN</code> .
Subquery	See Nested Query .
Three-Valued Logic	The logic system used by SQL that includes <code>TRUE</code> , <code>FALSE</code> , and <code>UNKNOWN</code> as possible values for logical expressions.
UNKNOWN	The logical value in SQL's three-valued logic that results from comparisons involving <code>NULL</code> values.
View	A virtual table based on the result-set of an SQL statement. Views can be used to simplify complex queries, join data from multiple tables, and implement access control.

Chapter 9

A Beginner's Guide to Database Normalization: Taming Your Data

1. Introduction: Why Do We Need to Organize Data?

If you've heard the term "database normalization," you might think it's a highly technical, complex process. In reality, it's simply a methodical way of organizing data in a database to avoid common pitfalls. The goal is to make your data structures logical, efficient, and reliable.

The single biggest problem that normalization solves is **data redundancy**, which means unnecessarily repeating information. Consider a single table, `EMP_DEPT`, that stores information about both employees and the departments they work in. For every employee in the "Research" department, the department's name ('Research') and manager's SSN ('333445555') are repeated. This repetition might seem harmless at first, but it creates a host of problems as you try to use and maintain your data.

Common Data Problems (Anomalies)

Redundant data leads to a set of issues known as data anomalies. These are errors that can occur when you try to add, update, or delete information.

- **Insertion Anomaly:** What happens if you want to add a new department that has no employees yet? In a combined table like `EMP_DEPT`, you can't. This is because the employee's Social Security Number (`Ssn`) is the Primary Key, and Primary Keys cannot be empty (NULL). Without an employee, you can't create the row, so you can't add the department.
- **Deletion Anomaly:** What happens if you delete the record for the last employee in a particular department? You would accidentally lose all information about the department itself, because it only existed in that employee's row.
- **Modification Anomaly:** What happens if a department's manager changes? You would have to find every single employee record for that department and update the manager's information, risking inconsistency if you miss one.

Normalization is a formal, step-by-step process designed to break down large, problematic tables into smaller, well-structured ones, fixing these very problems along the way.

2. The Building Blocks: Understanding Keys and Dependencies

Before we can start organizing our data, we need to understand two fundamental concepts: how we uniquely identify a piece of data (Keys) and how different pieces of data relate to one another (Dependencies).

2.1. Identifying Your Data: Primary Keys (PK)

A **Primary Key (PK)** is an attribute (or a set of attributes) that uniquely identifies a single record in a table. Think of it as a record's unique ID number. In the diagrams, primary keys are underlined.

- In the `EMPLOYEE` table, `Ssn` is the primary key because no two employees will have the same Social Security Number.
- Sometimes, a primary key requires more than one column. For example, in the `WORKS_ON` table, the combination of `{Ssn, Pnumber}` is the primary key. This is necessary because a single employee can work on many projects, and a single project can have many employees; only the combination of both uniquely identifies a specific work assignment.

2.2. The Core Concept: Functional Dependencies (FD)

A **functional dependency** exists when the value of one attribute uniquely determines the value of another. We use the notation $X \rightarrow Y$ to say, “X determines Y.” If you know the value of X, you can find the one and only corresponding value for Y.

The `EMP_DEPT` table provides clear examples of this relationship:

1. $Ssn \rightarrow Ename$: An employee's Social Security Number (`Ssn`) uniquely determines their name (`Ename`).
2. $Dnumber \rightarrow Dname$: A department's number (`Dnumber`) uniquely determines its name (`Dname`).

Understanding these dependencies is the key to correctly organizing your data. They act as the clues that tell us which pieces of information belong together in their own tables.

3. The Step-by-Step Journey: The Normal Forms

The Normal Forms are a series of rules or checks we apply to our database tables to reduce redundancy and improve data integrity. We'll look at the first three, which are the most important for good database design.

3.1. First Normal Form (1NF): Making Data Atomic

The first rule is the simplest. For a table to be in First Normal Form (1NF), it must meet two conditions: it must not have any “repeating groups,” and every attribute in every row must hold a single, “atomic” value. In simple terms, this means you can't have lists of items in a single cell.

For example, a poorly designed table might try to store all of a department's locations in one field, which violates 1NF:

Bad Design (Not 1NF) | Dnumber | Dname | Dlocations || :--- | :--- | :--- || 5 | Research | Bellaire, Sugarland, Houston |

To fix this, we create a separate row for each department-location pair. The `DEPT_LOCATIONS` table is a perfect example of a 1NF structure, where the composite primary key `{Dnumber, Dlocation}` ensures each row is unique.

Good Design (1NF) | Dnumber | Dlocation | | :--- | :--- | | 1 | Houston | | 4 | Stafford | | 5 | Bellaire | | 5 | Sugarland | | 5 | Houston |

3.2. Second Normal Form (2NF): Removing Partial Dependencies

The rule for Second Normal Form (2NF) is: the table must already be in 1NF, and every non-key attribute must be **fully dependent** on the *entire* primary key.

This rule is most relevant for tables with composite primary keys (keys made of more than one column). It solves the problem of **partial dependency**, which occurs when a non-key attribute depends on only a part of the primary key, not the whole thing.

The `EMP_PROJ` table is the classic example. Its primary key is `{ssn, Pnumber}`. However, some attributes in the table don't depend on both parts of this key.

Dependency	Label	Explanation for a Beginner
$\text{Ssn} \rightarrow \text{Ename}$	(FD2)	The employee's name depends <i>only</i> on the Ssn , not the project number.
$\text{Pnumber} \rightarrow \text{Pname}, \text{Plocation}$	(FD3)	The project's name and location depend <i>only</i> on the Pnumber , not the employee's Ssn .

The solution to partial dependencies is **decomposition**. We pull out the attributes that are partially dependent and put them into their own tables. This process turns the single `EMP_PROJ` table into three smaller, more focused tables: `EMPLOYEE`, `PROJECT`, and `WORKS_ON`.

3.3. Third Normal Form (3NF): Eliminating Transitive Dependencies

The rule for Third Normal Form (3NF) is: the table must already be in 2NF, and there must be no **transitive dependencies**.

A transitive dependency is an indirect relationship where a non-key attribute depends on another non-key attribute, instead of depending directly on the primary key. You can think of it like a chain reaction: $\text{A} \rightarrow \text{B} \rightarrow \text{C}$, where A is the primary key. C depends on B, and B depends on A.

The original `EMP_DEPT` table demonstrates this problem. The primary key is Ssn .

- We know that $\text{Ssn} \rightarrow \text{Dnumber}$ (the employee's SSN determines their department number).
- We also know that $\text{Dnumber} \rightarrow \text{Dname}$ (the department number determines the department name).

This creates the transitive dependency: $\text{Ssn} \rightarrow \text{Dnumber} \rightarrow \text{Dname}$. The department's name (Dname) doesn't depend directly on the employee's Ssn ; it depends on the department number (Dnumber), which is another non-key attribute.

The solution, again, is **decomposition**. We move the transitively dependent attributes (Dname , Dmgr_ssn) into their own new `DEPARTMENT` table. The original `EMPLOYEE` table no longer contains the department name or manager, but it keeps the Dnumber column, which now serves as a **Foreign Key** to link to the new `DEPARTMENT` table.

This act of decomposition is the heart of normalization. By methodically identifying and resolving these dependencies, we transform our table from a state of logical inconsistency into two tables, each with a single, clear purpose.

4. The Final Picture: A Well-Organized Database

After applying the normalization rules, we transform our initial, clunky tables into a set of lean, logical, and interconnected tables. This final structure for the COMPANY database is clean, efficient, and free from the insertion, deletion, and modification anomalies we started with.

Here is the resulting normalized database schema:

- **EMPLOYEE** (Ename, Ssn, Bdate, Address, Dnumber)
 - PK: Ssn
 - FK: Dnumber (references DEPARTMENT)
- **DEPARTMENT** (Dname, Dnumber, Dmgr_ssn)
 - PK: Dnumber
 - FK: Dmgr_ssn (references EMPLOYEE)
- **DEPT_LOCATIONS** (Dnumber, Dlocation)
 - PK: {Dnumber, Dlocation}
 - FK: Dnumber (references DEPARTMENT)
- **PROJECT** (Pname, Pnumber, Plocation, Dnum)
 - PK: Pnumber
 - FK: Dnum (references DEPARTMENT)
- **WORKS_ON** (Ssn, Pnumber, Hours)
 - PK: {Ssn, Pnumber}
 - FK: Ssn (references EMPLOYEE)
 - FK: Pnumber (references PROJECT)

The key insight is that by breaking down large, multi-purpose tables, we have created a set of smaller tables where each one is responsible for a single entity—employees, departments, projects. The relationships between them are now clear, intentional, and enforced by keys, eliminating the ambiguity that caused our initial problems.

5. Conclusion: The Power of Clean Data

Database normalization is more than just a technical exercise; it's a powerful technique for creating stable, reliable, and efficient databases. By systematically applying the normal forms, we directly eliminate

CS 340 Summary

the anomalies that make data difficult to manage. The process provides a clear path from a problematic design to a robust one.

Let's see how our journey solves the problems from the introduction:

1. **Reduces Redundancy:** Storing information like a department name or project location in a single place (its own table) is the core outcome. The **Modification Anomaly** is solved because a change only needs to be made in one row.
2. **Prevents Errors:** The **Insertion Anomaly** is solved because we can now add a new department to the `DEPARTMENT` table without needing an employee. The **Deletion Anomaly** is solved because deleting the last employee from a department in the `EMPLOYEE` table doesn't remove the department's record from the `DEPARTMENT` table.
3. **Improves Data Integrity:** By creating focused tables for each entity and linking them with keys, the logical structure ensures that relationships between data are clear and reliable, leading to higher-quality data overall.

Study Guide: Database Normalization

This guide provides a comprehensive review of database normalization principles, including key definitions, the process of advancing through normal forms, and the problems that normalization aims to solve. It includes a short-answer quiz with an answer key, suggested essay questions for deeper study, and a glossary of essential terms.

Short-Answer Quiz

Instructions: Answer the following questions in two to three sentences based on the provided source material.

1. What is database normalization and what is its primary goal?
2. Describe the three main types of update anomalies that can occur in a poorly designed database.

3. What is a functional dependency (FD) in the context of a database relation?
4. Explain the difference between a partial dependency and a full functional dependency.
5. What are the requirements for a relation to be in First Normal Form (1NF)?
6. Define Second Normal Form (2NF) and identify the type of dependency it eliminates.
7. What is a transitive dependency?
8. Define Third Normal Form (3NF) and state which type of dependency it is designed to remove.
9. Using the provided database examples, what is the purpose of a primary key?
10. What are the two desirable properties of a database decomposition?

Essay Questions

1. Using the `EMP_PROJ` relation as a case study, walk through the process of normalizing it to Third Normal Form (3NF). Identify all functional dependencies, explain which normal forms are violated and why, and show the final, decomposed 3NF relations.

2. Discuss the problems caused by data redundancy in detail. Provide specific examples of insertion, deletion, and modification anomalies using the `EMP_DEPT` relation from the source material.
3. Explain the relationship between a superkey, a candidate key, and a primary key. Why is the careful selection of a primary key a critical step in relational database design?
4. Describe the process of database design through normalization, starting from an unnormalized state. What specific issues are addressed at the 1NF, 2NF, and 3NF stages?
5. Evaluate the importance of the lossless join and dependency preservation properties in database decomposition. What are the potential consequences of a decomposition that fails to satisfy one or both of these properties?

Answer Key (Short)

1. **What is database normalization and what is its primary goal?** Normalization is the process of analyzing and decomposing complex relations into smaller, more manageable, and well-structured relations. Its primary goal is to minimize data redundancy, which in turn helps to avoid update anomalies and ensure data integrity.
2. **Describe the three main types of update anomalies that can occur in a poorly designed database.** The three types are insertion, deletion, and modification anomalies. An insertion anomaly occurs when certain data cannot be added unless another unrelated piece of data is also added. A

deletion anomaly happens when deleting a record inadvertently removes other essential information. A modification anomaly requires changing the same piece of information in multiple records, leading to potential inconsistencies.

3. **What is a functional dependency (FD) in the context of a database relation?** A functional dependency is a constraint between two sets of attributes in a relation, denoted as $X \rightarrow Y$. It specifies that the value of the attribute set X uniquely determines the value of the attribute set Y.
4. **Explain the difference between a partial dependency and a full functional dependency.** A full functional dependency means an attribute depends on an entire composite primary key and not on any of its subsets. In contrast, a partial dependency occurs when a non-prime attribute is functionally dependent on only a part of the composite primary key.
5. **What are the requirements for a relation to be in First Normal Form (1NF)?** A relation is in 1NF if the domain of each attribute contains only atomic (indivisible) values. This means the relation must not have any multi-valued attributes, composite attributes, or repeating groups.
6. **Define Second Normal Form (2NF) and identify the type of dependency it eliminates.** A relation is in Second Normal Form (2NF) if it is already in 1NF and every non-prime attribute is fully functionally dependent on the primary key. This process eliminates partial dependencies, where a non-prime attribute depends on only a part of a composite primary key.
7. **What is a transitive dependency?** A transitive dependency is an indirect functional dependency where $X \rightarrow Z$ exists because $X \rightarrow Y$ and $Y \rightarrow Z$. In this case, a non-prime attribute (Z) is dependent on another non-prime attribute (Y), which is itself dependent on the primary key (X).
8. **Define Third Normal Form (3NF) and state which type of dependency it is designed to remove.** A relation is in Third Normal Form (3NF) if it is already in 2NF and no non-prime attribute is transitively dependent on the primary key. 3NF is achieved by removing transitive dependencies, ensuring all attributes depend only on the primary key.
9. **Using the provided database examples, what is the purpose of a primary key?** A primary key is a candidate key chosen to uniquely identify each tuple (row) within a relation. For example, in the EMPLOYEE relation, the `ssn` attribute is the primary key because it ensures every employee record is unique and can be reliably referenced.
10. **What are the two desirable properties of a database decomposition?** The two desirable properties are the lossless join property and the dependency preservation property. Lossless join ensures that no information is lost and no spurious tuples are created when the decomposed tables are rejoined. Dependency preservation ensures that all original functional dependencies can be enforced on the new, smaller relations.

Glossary of Key Terms

Term	Definition
Normalization	The process of decomposing complex or “bad” relations by breaking them into smaller, well-structured relations to minimize redundancy and avoid update anomalies.
Redundancy	The unnecessary storage of the same data in multiple places within a database, leading to wasted space and potential data inconsistencies.
Update Anomaly	A data inconsistency that results from redundant data when performing insertion, deletion, or modification operations.
Insertion Anomaly	The inability to add a new tuple to a relation because a required attribute value is not yet known or available.
Deletion Anomaly	The unintended loss of data about one entity when a tuple representing another entity is deleted.
Modification Anomaly	An issue where changing a data value in one tuple requires making the same change in multiple other tuples to maintain consistency.
Functional Dependency (FD)	A constraint between two sets of attributes (X and Y) in a relation, written as $X \rightarrow Y$, where the value of X uniquely determines the value of Y.
Full Functional Dependency	A dependency where a set of attributes Y is functionally dependent on a composite key X, but not on any proper subset of X.
Partial Dependency	A dependency where a non-prime attribute is functionally dependent on only a part of a composite primary key. This violates Second Normal Form.
Transitive Dependency	An indirect dependency where $A \rightarrow C$ because $A \rightarrow B$ and $B \rightarrow C$, and B is not part of the primary key. This violates Third Normal Form.
First Normal Form (1NF)	A property of a relation where all attribute domains contain only atomic (indivisible) values, and each tuple is unique.
Second Normal Form (2NF)	A property of a relation that is in 1NF and where every non-prime attribute is fully functionally dependent on the primary key, meaning there are no partial dependencies.
Third Normal Form (3NF)	A property of a relation that is in 2NF and where no non-prime attribute is transitively dependent on the primary key.
Superkey	A set of one or more attributes that, taken collectively, can uniquely identify a tuple in a relation.
Candidate Key	A minimal superkey; a superkey from which no attribute can be removed without losing its unique identification property. A relation may have multiple candidate keys.

Primary Key (PK)	A candidate key selected by the database designer to be the main identifier for tuples in a relation. It is typically underlined in schema diagrams.
Foreign Key (FK)	An attribute or a set of attributes in one relation that refers to the primary key of another (or the same) relation, used to link the two relations.
Decomposition	The process of breaking down a single relation into two or more smaller relations as part of normalization.
Lossless Join Property	A property of decomposition that ensures the original relation can be perfectly reconstructed by joining the decomposed relations, without creating any spurious (extra) tuples.
Dependency Preservation Property	A property of decomposition that ensures all functional dependencies from the original relation are still enforceable by checking constraints on the smaller, decomposed relations.

Chapter 10

NoSQL and Object-Oriented Databases: A Briefing

Executive Summary

This document provides a comprehensive overview of non-relational database technologies, primarily focusing on the NoSQL movement and the document-oriented database MongoDB. The rise of these technologies is a direct response to the limitations of the traditional relational database model, particularly the “impedance mismatch” with modern object-oriented programming paradigms.

Key takeaways include:

- **Motivation for Alternatives:** The rigid, tabular structure of relational databases creates a fundamental conflict with the complex, nested object structures common in application development, necessitating alternative data storage solutions like Object-Oriented DBMS (OODBMS) and, more prominently, NoSQL.
- **The NoSQL Paradigm:** NoSQL, or “Not Only SQL,” represents a broad category of databases that are non-relational, distributed, open-source, and designed for horizontal scalability. They often operate on the BASE principle (Basically Available, Soft state, Eventually consistent), prioritizing availability and performance over the strict consistency guaranteed by the ACID properties of relational systems.
- **Core NoSQL Models:** There are four primary NoSQL data models, each optimized for different use cases:
 1. **Key-Value Stores:** Simple, high-performance databases mapping a key to a value.
 2. **Document Stores:** Store flexible, semi-structured data in document formats like JSON/BSON.
 3. **Column-Family Stores:** Organize data by columns instead of rows, ideal for large-scale analytics.
 4. **Graph Databases:** Model data as nodes and edges, excelling at managing complex relationships.
- **MongoDB as a Case Study:** MongoDB is a leading document store that uses BSON (a binary version of JSON) for data storage. Its core concepts include databases, collections (analogous to tables), and documents (analogous to rows), offering a flexible schema where

documents within a single collection can have different structures. MongoDB supports a rich query language for full CRUD (Create, Read, Update, Delete) operations.

The Rationale for Non-Relational Databases

The Impedance Mismatch

A primary driver for the development of non-relational databases is the concept of “impedance mismatch.” This term describes the inherent difficulties that arise when trying to map the world of object-oriented programming (OOP) to the relational database model.

- **OOP Model:** In languages like Java, Python, or C++, data is represented as complex objects with properties, methods, and nested structures, including collections and inheritance hierarchies.
- **Relational Model:** Relational Database Management Systems (RDBMS) store data in flat tables composed of rows and columns containing simple scalar values (e.g., numbers, strings).

This fundamental difference forces developers to write a significant amount of translation code (Object-Relational Mapping, or ORM) to flatten objects into tables for storage and reconstruct them upon retrieval.

Early Solutions: Object-Oriented DBMS

Object-Oriented Database Management Systems (OODBMS) emerged as an early attempt to solve the impedance mismatch. These systems were designed to work directly with the objects used in OOP languages.

Core Features:

- **Object Identity (OID):** Each object is given a unique, immutable identifier, allowing it to be referenced regardless of its state or location.
- **Encapsulation and Inheritance:** OODBMS support core OOP principles, allowing for complex data types and class hierarchies to be persisted directly.
- **Persistence:** Objects in the programming language could be made persistent with minimal code changes.

- **Status:** Despite their conceptual advantages, OODBMS like GemStone/S and Objectivity/DB did not achieve widespread adoption and have largely been superseded in popularity by the NoSQL movement.

Introduction to the NoSQL Movement

NoSQL stands for “Not Only SQL” and represents a diverse range of database technologies that moved away from the relational model.

Core Characteristics

- **Non-relational:** Data is not stored in the rigid table-and-row format of RDBMS.
- **Distributed:** Systems are typically designed to run on clusters of machines.
- **Open-source:** Many popular NoSQL databases are developed and maintained by open-source communities.
- **Horizontal Scalability:** They are designed to “scale out” by adding more servers to a cluster, rather than “scaling up” by increasing the power of a single server. This makes them highly effective for handling large volumes of data and high traffic loads.

The BASE Principle

While traditional RDBMS adhere to ACID (Atomicity, Consistency, Isolation, Durability) properties to guarantee transaction reliability, many NoSQL systems favor the **BASE** principle, which is better suited for distributed environments where high availability is critical.

- **Basically Available:** The system guarantees availability, even in the event of partial server failure.
- **Soft state:** The state of the system may change over time, even without input, due to the “eventually consistent” model.
- **Eventually consistent:** The system will eventually become consistent once all inputs have propagated, but it does not guarantee immediate consistency across all nodes.

Primary NoSQL Data Models

The NoSQL landscape is categorized into four main data models, each with distinct strengths.

Data Model	Description	Strengths & Use Cases	Examples
Key-Value	A simple model using a dictionary or hash map where each unique key is associated with a single value.	Extremely high performance and scalability. Ideal for caching, session management, and real-time applications. Less suitable for complex queries involving multiple keys or relationships.	Redis, Memcached, Amazon DynamoDB
Document	Stores data in semi-structured documents, typically in formats like JSON or BSON. Documents can be nested and do not require a uniform schema.	Schema flexibility allows for easy evolution of applications. Data structure often maps naturally to objects in code. Ideal for content management, mobile apps, and e-commerce.	MongoDB, CouchDB
Column-Family	Stores data in tables, but data is organized by columns rather than rows. A “row” is a collection of columns that can vary from one row to the next.	Highly efficient for read/write operations on large datasets and for analytical queries that aggregate data over a subset of columns.	Apache Cassandra, HBase
Graph	Stores data as a network of nodes (entities) and edges (relationships).	Purpose-built for storing and navigating complex relationships. Ideal for social networks, fraud detection, and recommendation engines.	Neo4j

The Document Model in Detail: JSON and BSON

Document databases are one of the most popular types of NoSQL stores, and their data format is central to their functionality.

JSON (JavaScript Object Notation)

JSON is a lightweight, human-readable data-interchange format derived from a subset of JavaScript’s object literal notation. It is language-independent and has become a standard for data transfer on the web.

- **Core Structures:**

- **Object:** An unordered collection of `string: value` pairs enclosed in curly braces `{ }.`

This is analogous to a dictionary or hash map.

CS 340 Summary

- **Array:** An ordered list of values enclosed in square brackets []. This is analogous to a list or array.
- **Value Types:** A value in JSON can be one of seven types:
a string, number, object, array, true, false, or null.

BSON (Binary JSON)

BSON is a binary-encoded serialization format for JSON-like documents. It is the primary data storage format for MongoDB.

- **Purpose:** BSON is designed to be lightweight, traversable, and efficient for encoding and decoding by machines.
- **Enhancements over JSON:** While conceptually similar to JSON, BSON extends the specification to include additional data types that are not native to JSON, such as dates, binary data, and specific number types (e.g., 32/64-bit integers, doubles).

Case Study: MongoDB

MongoDB is a leading document-oriented NoSQL database that exemplifies the principles of the document model.

Overview and Features

MongoDB stores data in flexible, BSON-formatted documents. It is designed for high performance, high availability, and easy scalability.

- **Key Features:**
 - Ad hoc queries
 - Indexing on any field
 - Replication for high availability (via replica sets)
 - Load balancing and sharding for horizontal scaling

Core Concepts

MongoDB's architecture can be compared to that of a traditional RDBMS.

MongoDB Concept	Relational DB Equivalent	Description
Database	Database	A physical container for collections.
Collection	Table	A grouping of MongoDB documents. Collections do not enforce a schema; documents within a collection can have different fields.
Document	Row / Record	A single record stored in BSON format, consisting of field-and-value pairs.

CRUD Operations in MongoDB

MongoDB provides a rich set of operations to perform Create, Read, Update, and Delete actions on data.

Create: insertOne()

This operation inserts a single document into a collection.

```
db.users.insertOne(
{
  name: "sue",
  age: 26,
  status: "pending"
})
```

- **collection:** `users` is the target collection for the insertion.
- **document:** The BSON object `{ name: "sue", ... }` is the record being created, composed of `field: value` pairs.

Read: find()

This operation retrieves documents from a collection that match specified criteria. It returns a cursor, which is a pointer to the result set.

```
db.users.find(
  { age: { $gt: 18 } },
  { name: 1, address: 1 }
).limit(5)
```

- **collection:** `users` is the collection being queried.

CS 340 Summary

- **query criteria:** `{ age: { $gt: 18 } }` is the filter that selects documents where the `age` is greater than 18.
- **projection:** `{ name: 1, address: 1 }` specifies which fields to include in the returned documents (1 for include). The `_id` field is returned by default.
- **cursor modifier:** `.limit(5)` modifies the cursor to return a maximum of 5 documents.

Update: `updateMany()`

This operation modifies all documents that match a specified filter.

```
db.users.updateMany(  
  { age: { $lt: 18 } },  
  { $set: { status: "reject" } }  
)
```

- **collection:** `users` is the target collection.
- **update filter:** `{ age: { $lt: 18 } }` selects all documents where the `age` is less than 18.
- **update action:** `{ $set: { status: "reject" } }` is the modification to be applied. The `$set` operator updates the value of the `status` field to "reject".

Delete: `deleteMany()`

This operation removes all documents from a collection that match a filter.

```
db.users.deleteMany(  
  { status: "reject" }  
)
```

- **collection:** `users` is the target collection.
- **delete filter:** `{ status: "reject" }` specifies the criteria for deletion, removing all documents where the `status` field is "reject".

Study Guide: NoSQL and Object-Oriented Databases

This guide provides a comprehensive review of the concepts related to NoSQL and Object-Oriented Databases, based on the provided source materials. It includes a short-answer quiz to test your knowledge, a set of essay questions for deeper analysis, and a complete glossary of key terms.

Short-Answer Quiz

Instructions: Answer the following ten questions. Each answer should be approximately 2-3 sentences long and based entirely on the information provided in the source context.

1. What is NoSQL, and what was the primary driver for its creation?

2. Name the four main categories of NoSQL databases and briefly describe the data model for each.

3. Explain the concept of “impedance mismatch” and which type of database was developed to address it.

4. What is a “document” in the context of a document-oriented database like MongoDB?

5. Describe the difference between the ACID and BASE properties as they relate to databases.

6. In a MongoDB `find` query, what are the roles of the “query criteria” and “projection” parameters?

7. Define “object identity” (OID) as it applies to Object-Oriented Databases.

8. What is JSON, and what are its two primary structural components?

9. Explain the concept of “persistence” in an Object-Oriented Database Management System (OODBMS).

10. What is BSON and what is its relationship to JSON within MongoDB?

Essay Questions

Instructions: The following questions are designed for longer, more detailed responses. Formulate your answers by synthesizing information from the source material.

1. Compare and contrast the data modeling approach of a traditional relational database with that of a document-oriented NoSQL database. Discuss the advantages and disadvantages of each, particularly concerning schema flexibility, scalability, and handling complex, hierarchical data.

2. Explain the concept of “impedance mismatch” in detail. How do Object-Oriented Databases (OODBs) address this problem, and what are the key features of OODBs (like object identity, encapsulation, and persistence) that facilitate this solution?

3. Discuss the trade-offs between the ACID properties of traditional relational databases and the BASE properties common in NoSQL systems. In what types of applications would a BASE-compliant system be preferable to an ACID-compliant one, and why?

4. Describe the four main categories of NoSQL databases (Key-Value, Column-Family, Graph, Document). For each category, provide a hypothetical use case that illustrates its specific strengths for storing and querying data.

5. Using the provided MongoDB query examples (`insertOne`, `find`, `updateMany`, `deleteMany`), explain the fundamental CRUD (Create, Read, Update, Delete) operations in a document database. Deconstruct the syntax of each operation, explaining the role of collections, documents, filters, and other operators.

Answer Key (Short)

- 1. What is NoSQL, and what was the primary driver for its creation?** NoSQL, meaning “Not Only SQL,” refers to a class of non-relational databases designed for flexibility and scalability. The primary driver for its creation was the emergence of Big Data, as traditional relational databases struggled to handle the massive volume, velocity, and variety of data generated by modern applications.
- 2. Name the four main categories of NoSQL databases and briefly describe the data model for each.** The four main categories are Key-Value, which stores data as a simple key with an associated value; Column-Family, which stores data in tables with rows and dynamic columns; Graph, which uses nodes and edges to represent and store data relationships; and Document, which stores data in structured documents, such as JSON.
- 3. Explain the concept of “impedance mismatch” and which type of database was developed to address it.** Impedance mismatch is the difficulty and overhead that occurs when trying to store data from an object-oriented programming language into a relational database, due to their fundamentally different data models. Object-Oriented Databases (OODBs) were developed to solve this problem by allowing objects to be stored directly without conversion.
- 4. What is a “document” in the context of a document-oriented database like MongoDB?** In a document-oriented database, a document is the basic unit of data, roughly equivalent to a row in a relational database. It is a data structure composed of field and value pairs, similar in structure to a JSON object. For example, a user document might contain fields like “name”, “age”, and “status”, each with corresponding values.
- 5. Describe the difference between the ACID and BASE properties as they relate to databases.** ACID (Atomicity, Consistency, Isolation, Durability) is a set of properties guaranteeing that database transactions are processed reliably, which is characteristic of traditional SQL databases. BASE (Basically Available, Soft state, Eventual consistency), common in NoSQL systems, prioritizes availability over strict consistency, allowing for better performance and scalability in distributed environments.
- 6. In a MongoDB find query, what are the roles of the “query criteria” and “projection” parameters?** The “query criteria” parameter is a document used to specify the conditions for selecting which documents to return, such as finding users where `age` is greater than 18. The “projection” parameter specifies which fields from the selected documents should be included in the results, such as returning only the `name` and `address` fields.
- 7. Define “object identity” (OID) as it applies to Object-Oriented Databases.** In an OODB, object identity means that every object is assigned a unique, system-generated identifier (OID) that remains constant for the object’s entire lifetime, regardless of changes to its state or data. This allows the system to distinguish between two objects that might otherwise have identical attribute values.
- 8. What is JSON, and what are its two primary structural components?** JSON, or JavaScript Object Notation, is a lightweight format for storing and transporting data, often used in NoSQL databases. Its two primary structural components are objects, which are collections of key/value pairs enclosed in curly braces {}, and arrays, which are ordered lists of values enclosed in square brackets [].
- 9. Explain the concept of “persistence” in an Object-Oriented Database Management System (OODBMS).** Persistence in an OODBMS refers to the ability of an object to exist beyond the lifetime of the program that created it. This means the object’s state is saved in the database and can be retrieved and used later, even after the original application has terminated.

10. What is BSON and what is its relationship to JSON within MongoDB? BSON stands for “Binary JSON” and is a binary-encoded serialization of JSON-like documents. MongoDB uses BSON internally to store documents because it is lightweight, traversable, and efficient for encoding and decoding, while also extending the JSON model to include additional data types.

Glossary of Key Terms

Term	Definition
ACID	An acronym for Atomicity, Consistency, Isolation, Durability. A set of properties of database transactions intended to guarantee validity even in the event of errors, power failures, etc. It is typically associated with relational databases.
Array (JSON)	A structural component of JSON. An ordered collection of values, enclosed in square brackets [], with values separated by commas.
BASE	An acronym for Basically Available, Soft state, Eventual consistency. A model used by many NoSQL databases that prioritizes availability and scalability over the strict consistency offered by ACID.
Big Data	A term describing the large volume of data – both structured and unstructured – that inundates a business on a day-to-day basis. The emergence of Big Data was a key driver for the development of NoSQL databases.
BSON	An acronym for “Binary JSON.” A binary-encoded serialization of JSON-like documents used by MongoDB for data storage due to its lightweight, traversable, and efficient nature.
Collection	In MongoDB, a collection is a grouping of documents. It is the equivalent of a table in a relational database.
Column-Family Database	A type of NoSQL database that stores data in tables, rows, and dynamic columns. It is designed to store and process very large amounts of data distributed over many machines.
Complex Objects	Objects in an OODB that can store other objects, allowing for nested and intricate data structures that mirror real-world entities.
CRUD	An acronym for Create, Read, Update, and Delete, which are the four basic functions of persistent storage.
Document	The basic unit of data in a document-oriented database like MongoDB. It is a data structure composed of field and value pairs, similar to a JSON object.
Document Database	A type of NoSQL database that stores data in structured documents, such as JSON or BSON. MongoDB is a popular example.
Encapsulation	An object-oriented programming principle where an object’s data (attributes) and methods (behavior) are bundled together. In OODBs, this allows the database to manage both the data and its associated behavior.
Graph Database	A type of NoSQL database that uses graph structures with nodes, edges, and properties to represent and store data. It is optimized for analyzing interconnections between data points.

CS 340 Summary

Horizontal Scaling	The ability to increase capacity by connecting multiple hardware or software entities so that they work as a single logical unit. This is a key feature of many NoSQL databases.
Impedance Mismatch	The difficulties that arise when moving data between an object-oriented programming language and a relational database due to differences in their data models.
Inheritance	An object-oriented concept where a new class can be based on an existing class, inheriting its attributes and methods. OODBs support this, allowing for natural representation of “is-a” relationships.
JSON	An acronym for JavaScript Object Notation. A lightweight, human-readable text format for data interchange, built on objects (key/value pairs) and arrays (ordered lists).
Key-Value Database	The simplest type of NoSQL database where each data item is stored as a key together with its value.
MongoDB	A popular open-source, document-oriented NoSQL database that uses JSON-like documents with optional schemas.
NoSQL	Stands for “Not Only SQL.” A broad class of non-relational database management systems that differ from the classic relational model, often excelling in scalability, performance, and flexibility.
Object (JSON)	A structural component of JSON. An unordered set of key/value pairs, where a key is a string and a value can be one of the seven valid JSON data types. An object begins with { and ends with }.
Object Identity (OID)	A unique, system-generated identifier assigned to every object in an OODB, which persists for the object’s lifetime and distinguishes it from all other objects.
Object-Oriented Database (OODB)	A database system in which information is represented in the form of objects, as used in object-oriented programming, designed to resolve the impedance mismatch problem.
Persistence	The characteristic of an object in an OODBMS to continue to exist even after the program that created it has terminated.
Projection	A parameter in a MongoDB <code>find</code> query that specifies which fields of a document to return in the query results.
Relational Database	A database that stores and provides access to data points that are related to one another, typically organized into tables with rows and columns. It enforces a strict schema and uses SQL for queries.
Scalability	The capability of a system to handle a growing amount of work by adding resources. NoSQL databases are known for their horizontal scalability.
Schema-less	A characteristic of many NoSQL databases, like MongoDB, where a collection can hold documents with different fields and structures. This provides high flexibility.
SQL	Stands for Structured Query Language, the standard language used to communicate with and manage data in relational databases.

