# CS 143

Warren Kim

# Contents

## 0.1 Purpose of a Database

We will be studying (mostly) Relational DataBase Management Systems (RDBMS).

> **Definition: Database**
>
> A **database** abstracts how data is stored, maintained, and processed. It is a system that uses advanced data structures to store and index data.

A database abstracts away the data integrity and file management aspect of CRUD operations. Moreover, a database provides us with a single location for all of the data, even if the database itself is distributed.

## 0.2 Abstraction Layers

There are three layers of abstraction: physical, logical, and view.

> **Definition: Physical Abstraction**
>
> The **physical abstraction** defines the data and its relationships to other data within the database.

> **Definition: Logical Abstraction**
>
> The **logical abstraction** deals with how we interface with the database.

> **Definition: View Abstraction**
>
> The **view abstraction** refers to specific use cases and filters the data from the logical abstraction.

We start by learning the logical abstraction.

## 0.3 Instances and Schema

> **Definition: Schema and Instance**
>
> A **schema**[a] is the overall design of a database. It defines the structure of the data as well as how it is organized.
>
> An **instance** of a database is the actual set of data stored in the database at a particular moment in time.
> _____
> [a]Note: schema can also refer to a relation (table).

## 0.4 Data Models

Data models define how we design databases and interact with data. We want to answer the following:

 (i) How do we define data?

 (ii) How do we encode relationships among data?

 (iii) How do we impose constraints on data?

Data models are either an Implementation model or a Design mechanism. Implementation models build databases from the ground up while design mechanisms are implemented as features in a database. We discuss five major types (an several niche ones).

### 0.4.1 Relational

In a relational model, all data is stored as a **relation**[1]. Rows represent individual $n$-tuple units (**records**). Columns represent (typed) **attributes** common to all records in the relations.

### 0.4.2 Entity-Relationship (ER)

An entity-relationship model uses a collection of basic objects (**entities**) and define **relationships** among them.

### 0.4.3 Object-Oriented

The object-oriented model is similar to OOP with encapsulation, methods, adn object identity. It was originally an implementation model but is now a design mechanism.

### 0.4.4 Document (Semi-Structured)

A document model stores records as **documents**, which do **not** have an enforced schema. This allows for more versatility in the type of data stored in the database.

### 0.4.5 Network/Hierarchical/Graphical

A graph model is analogous to how we think. Records are stored as **nodes** and relationships between records as **edges**.

### 0.4.6 Vector

A vector model stores records as **vectors** in $\mathbb{R}^n$, and are stored in a way that enables efficient retrieval and comparison (e.g. nearest neighbor[s]).

### 0.4.7 Key-Value

A key-value model stores data as a key-value pair (typically using a hash function). In this model, data typically lives in RAM as opposed to disk.

## 0.5 Database Languages

There are two main semantic systems when working with databases:

(i) Data Manipulation Language (DML)

(ii) Data Definition Language (DDL)

Note that a relational model typically uses SQL for both DDL and DML.

### 0.5.1 Data Manipulation Language

DML's can either be procedural or declarative.

---

**Definition: Query**

A **query** is a written expression to retrieve or manipulate data.

---

[1]Note: tables are an implementation of relations.

### 0.5.2 Data Definition Language

DDL's specify a schema: a collection of attribute names and data types, consistency constraints, and optionally storage structure and access methods. There are four types of consistency constraints:

*(i)* Domain constraints define the domain of an attribute (e.g. `tinyint`, `enum`, etc.).

*(ii)* Assertions are business rules that must hold true (e.g. an enforced prerequisite for a class must be present in your transcript before you can add a class to your study list).

*(iii)* Authorization determines who can do what (e.g. full CRUD, read-only, etc.).

*(iv)* Referential integrity ensures that links from one table to another must be defined (Suppose we have two relations $R, R'$. If there is a link $f : R \to R'$, then $f$ is surjective).

## 0.6 Data Storage and Querying

> **Definition: Storage Manager**
>
> A **storage manager** that abstracts away how the data is laid out on disk.

A storage manager is helpful because reading data from disk to RAM is *slow*, and the storage manager handles swapping[2] and makes retrieval efficient.

> **Definition: Query Manager**
>
> A **query manager** takes the DML statements and organize them into a *query plan*[a] that "compiles" a query (using relational algebra) and executes the instruction(s).
>
> ---
> [a]Note: The query plan dictates the performance of a query.

## 0.7 Keys

---
[2]Swapping: Virtual memory in CS111!

### 0.7.1 Superkey

> **Definition: Superkey**
>
> A **superkey** is a set of one or more attributes that uniquely identifies a record (tuple) and distinguishes it from all other records in the relation.
>
> Formally, let $R$ be a relation with a set $S = \{a_1, a_2, \ldots, a_n : a$ is an attribute of $R\}$. A **superkey** is a subset $s \subseteq S$ such that $s$ uniquely identifies each $n$-tuple in $R$.

The superkey $s = S = \{a_1, a_2, \ldots, a_n\} = \bigcup_{i=1}^{n}\{a_i\}$ is called the ***trivial superkey***. Additionally, $\emptyset$ is not a superkey. Further note that for every relation $R$, there exists at most $2^n - 1$ superkeys where $n$ is the number of attributes.

### 0.7.2 Candidate Key

> **Definition: Candidate Key**
>
> A **candidate key** is a superkey such that no subset of the candidate key is a superkey; i.e. it is the minimal superkey.
>
> Formally, let $R$ be a relation with a set $S = \{a_1, a_2, \ldots, a_n : a$ is an attribute of $R\}$. A **candidate key** is a superkey $s \subseteq S$ such that for every propery subset $t \subsetneq s$, $t$ is not a superkey.

Candidate keys may vary in length, and the attributes of a candiate key may be `NULL` as long as it uniquely identifies an $n$-tuple in the relation.

### 0.7.3 Primary Key

> **Definition: Primary Key and Composite Key**
>
> A **primary key** is a candidate key (chosen by the database designer) to enforce uniqueness for a particular use case.

The primary key is typically chosen to be the minimal candidate key for simplicity. The attributes of a primary key may not be `NULL`.

### 0.7.4 Foreign Key

> **Definition: Foreign Key**
>
> A **foreign key** is a set of attributes that links tuples of two relations.
>
> Formally, let $R, R'$ be relations with sets $S = \{a_1, a_2, \ldots, a_n : a$ is an attribute of $R\}, S' = \{a'_1, a'_2, \ldots, a'_n : a'$ is an attribute of $R'\}$. A **foreign key** is a key $s \subseteq S$ of $R$ that maps to the primary key $p \subseteq S'$ of $R'$.

Foreign keys are used to enforce referential integrity constraints; i.e. foreign keys in a relation $R$ are used to protect data in $R$ from being orphaned and/or inconsistent. Given two relations $R, R'$ related via a foreign key, $R'$ is said to be the *referring* relation and $R$ the *referred* relation.

Let two relations $R, S$ be related via a foreign key, where $S$ is the *referring* relation and $R$ is the *referred* relation. Suppose we want to remove an $n$-tuple $r \in R$. Then there are two cases:

*Case 1* If there is no $s \in S$ such that $s \mapsto r$, we simply remove $r$.

*Case 2* If there is at least one $s \in S$ such that $s \mapsto r$, we can either throw an error to prevent the deletion of $r$ or *cascade*[3] the delete.

---

[3]Cascade: Delete $r$ and all $s \in S$ that refer to $r$.

## 0.8  Defining a Schema

A schema can be written as `relation(`<u>`attribute`</u>$_1$`, ..., attribute`$_n$`)` where underlined attributes represent the primary key.

## 0.9  Relational Algebra

describe relational algebra

### 0.9.1  Selection

> **Definition: Selection**
>
> **Selection** retrieves a subset of tuples from a *single* relation $R$ that satisfies some predicate $\psi$ and returns a new relation $R' \subseteq R$, and is defined by
>
> $$\sigma_\psi(R) = R' = \{t \in R : \psi(t)\}$$
>
> where $\psi$ is a boolean predicate on attributes and values with respect to a unary or binary operator[a]
>
> ---
> [a]We may use the following operators: $\{=, \neq, <, >, \leq, \geq, \neg\}$.

We can build complex predicates using conjunction $\wedge$ (and) or disjunction $\vee$ (or).

**Note: that selection $\sigma$ is implemented as `WHERE` in SQL.**

Below are a list of examples of selection, assuming all attributes and relations are well-defined:

*(i)* $\sigma_{(\texttt{dislikes}<\texttt{likes})}(\texttt{youtube\_video})$

*(ii)* $\sigma_{(\texttt{cat\_id}=17)}(\texttt{youtube\_video})$

*(iii)* $\sigma_{([\texttt{dislikes}<\texttt{likes}]\wedge[\texttt{views}>1000000]\wedge[\texttt{cat\_id}=24])}(\texttt{youtube\_video})$

*(iv)* $\sigma_{(\texttt{dislikes}<\texttt{likes})}(\sigma_{(\texttt{views}>1000000)}(\sigma_{(\texttt{cat\_id}=24)}(\texttt{youtube\_videos})))$

Note that *(iii)* and *(iv)* are equivalent.

### 0.9.2  Projection

> **Definition: Projection**
>
> **Projection** extracts attributes from a set of tuples and removes duplicates. Given a relation $R$, $n$-tuple $t$, and a set of attributes $a_1, \cdots, a_n$,
>
> $$\Pi_{a_1,\cdots,a_n}(R) = \{t[a_1, \cdots, a_n] : t \in R\}$$
>
> Projection is usually the last (outermost) operation done on a relation.

> **Aside: Projection?**
>
> We call it a projection because we are collapsing an $n$-tuple down to an $(n-k)$-tuple. That is, we take the $n$-tuples in a relation $R_n$ and collapse them into a set of $(n-k)$-tuples in a new relation $R'_{n-k}$.
>
> ---
> Here, $R_n$ is a relation with $n$ attributes.