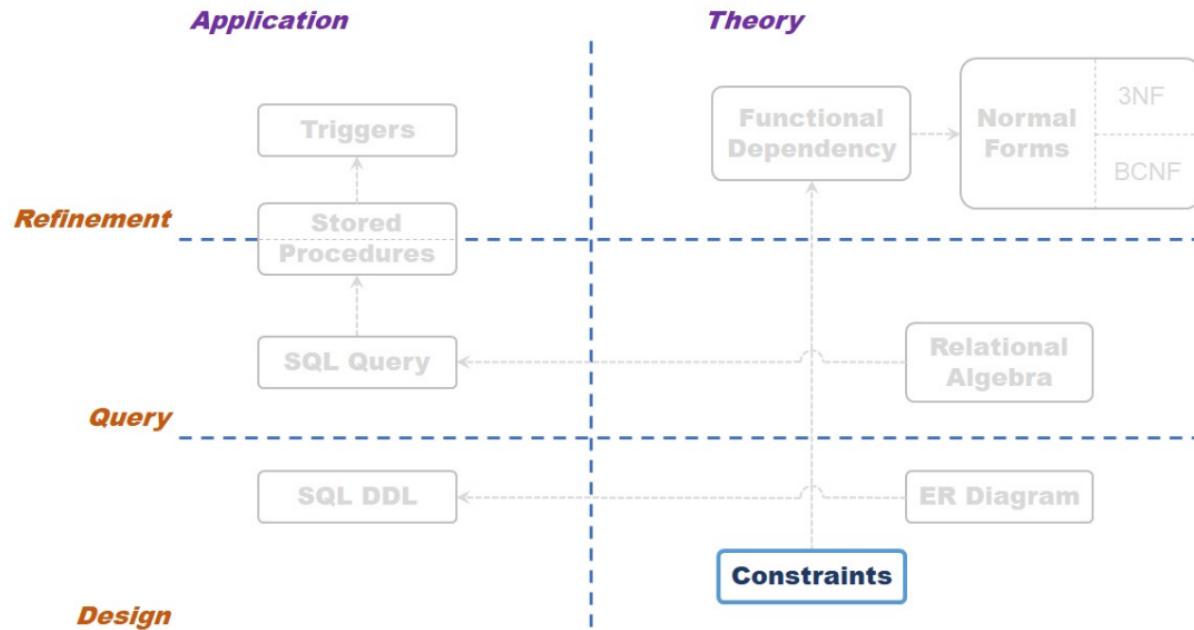


# CS2102: Database Systems

## Introduction

# Roadmap



# Roadmap

## Overview

## Overview

### *DBMS*

- Database
- Challenges
- File to DBMS
- Core Concepts

### *Relational Model*

- Motivation & History
- Core Concepts
- Constraints

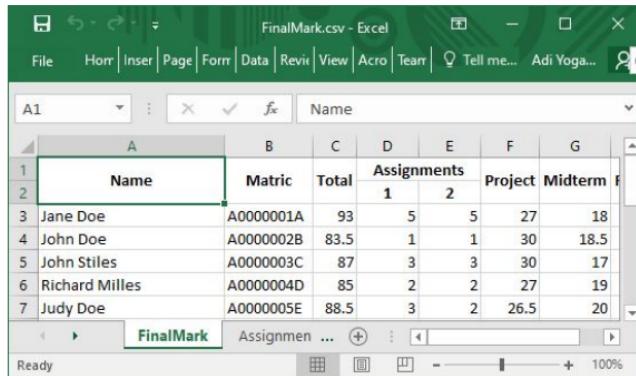
# DBMS

---

## What is a Database?

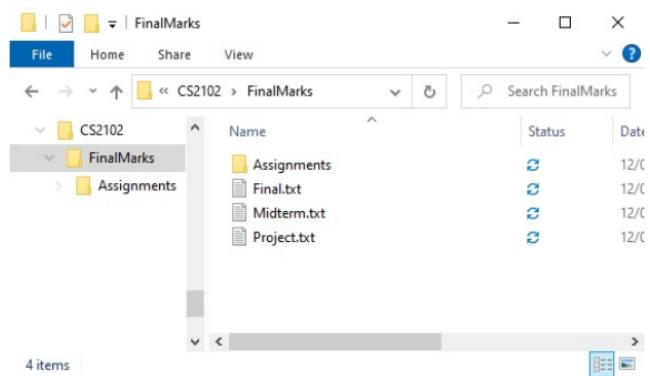
### Definition

A **database** is a collection of *organized* data, information, and records.



The screenshot shows a Microsoft Excel spreadsheet titled "FinalMark.csv - Excel". The table has columns labeled "Name", "Matric", "Total", "Assignments", "Project", and "Midterm". The "Assignments" column is further divided into "1" and "2". The data includes rows for five students: Jane Doe, John Doe, John Stiles, Richard Milles, and Judy Doe, with their respective scores.

	A	B	C	D	E	F	G
1	Name	Matric	Total	Assignments	Project	Midterm	
2				1	2		
3	Jane Doe	A0000001A	93	5	5	27	18
4	John Doe	A0000002B	83.5	1	1	30	18.5
5	John Stiles	A0000003C	87	3	3	30	17
6	Richard Milles	A0000004D	85	2	2	27	19
7	Judy Doe	A0000005E	88.5	3	2	26.5	20



# DBMS

Database  
Challenges  
Management  
Transactions  
Serialization  
Architecture  
Study

## Common Challenges

### Data-Intensive Applications

#### Efficiency

- | Fast access to information in **huge** volumes of data.



#### Transactions

- | "All-or-nothing" changes to data (*e.g., bank transfer*).

| 5,000 tps\*

#### Data Integrity

- | Parallel access and **changes** to data.

| 100,000 tps\*

#### Recovery

- | Fast and **reliable** handling of failures.



#### Security

- | Fine-grained data **access** rights.

| 544,000 tps\*

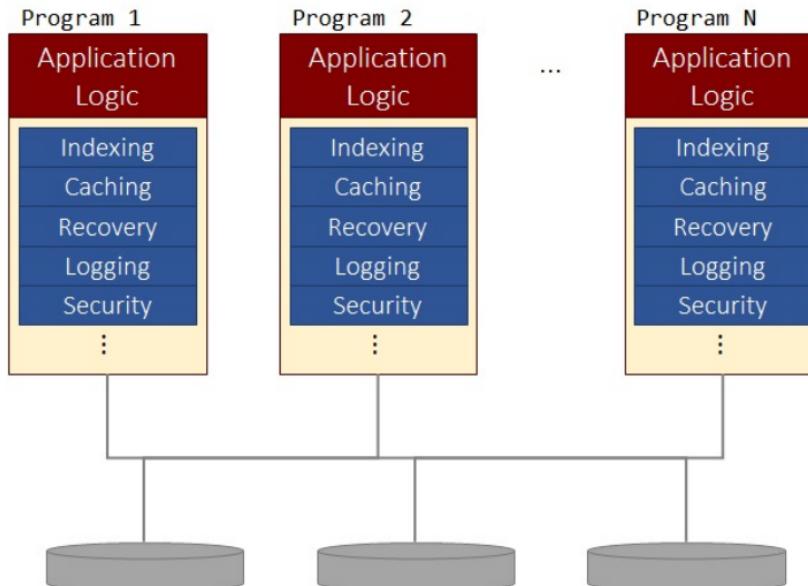
\***tps:** transactions per seconds (*self-reported peak values*).

# DBMS

Database Challenges  
**Management**  
- *File-Based*  
- *DBMS*  
- *Comparison*  
Transactions  
Serialization  
Architecture  
Study

# Data Management

# File-Based Data Management



## Note

- Complex, low-level code
  - Often similar requirements across different programs

## Problems

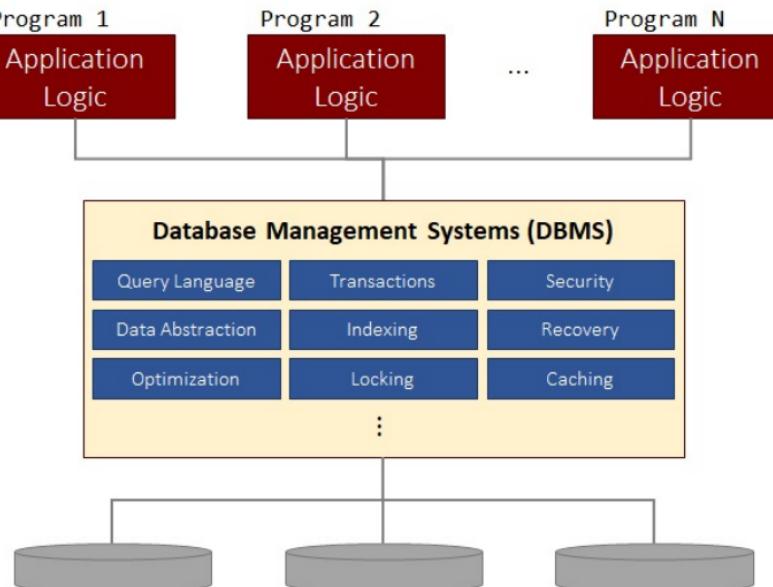
- High development effort
  - Long development times
  - Higher risk of errors

# DBMS

Database Challenges  
Management  
- File-Based  
- DBMS  
- Comparison  
Transactions  
Serialization  
Architecture Study

## Data Management

### Data Management with DBMS



#### Note

- Application logic moved to DBMS
- Other set of universal and powerful functionalities

#### Benefits

- Faster application development
- Increased productivity
- Higher stability

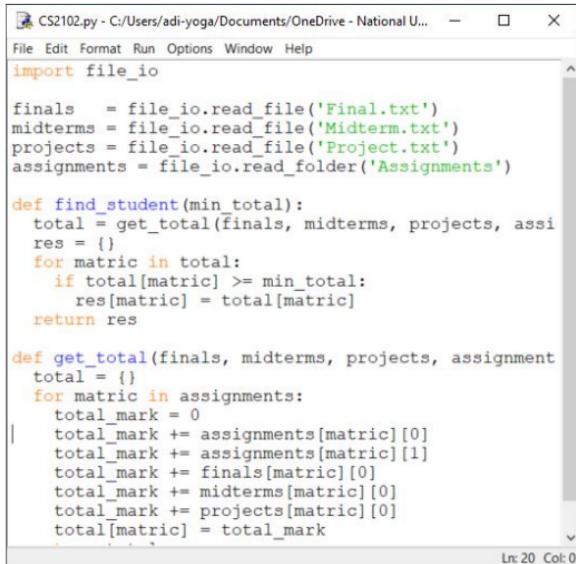
# DBMS

Database Challenges Management  
- File-Based  
- DBMS  
- Comparison  
Transactions  
Serialization  
Architecture Study

## Data Management

### Comparison

#### File-Based



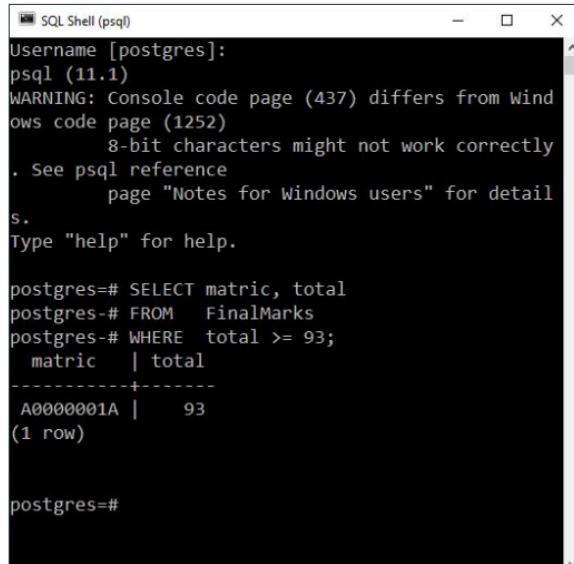
```
CS2102.py - C:/Users/adi-yoga/Documents/OneDrive - National U... - □ X
File Edit Format Run Options Window Help
import file_io

finals = file_io.read_file('Final.txt')
midterms = file_io.read_file('Midterm.txt')
projects = file_io.read_file('Project.txt')
assignments = file_io.read_folder('Assignments')

def find_student(min_total):
    total = get_total(finals, midterms, projects, assignments)
    res = {}
    for matric in total:
        if total[matric] >= min_total:
            res[matric] = total[matric]
    return res

def get_total(finals, midterms, projects, assignments):
    total = {}
    for matric in assignments:
        total_mark = 0
        total_mark += assignments[matric][0]
        total_mark += assignments[matric][1]
        total_mark += finals[matric][0]
        total_mark += midterms[matric][0]
        total_mark += projects[matric][0]
        total[matric] = total_mark
```

#### DBMS



```
SQL Shell (psql)
Username [postgres]: psql (11.1)
WARNING: Console code page (437) differs from Windows code page (1252)
         8-bit characters might not work correctly
See psql reference page "Notes for Windows users" for details.
Type "help" for help.

postgres=# SELECT matric, total
postgres-# FROM FinalMarks
postgres-# WHERE total >= 93;
      matric |      total
-----+-----
 A0000001A |      93
(1 row)

postgres=#

```

## Transactions

### Definition

A **finite sequence** of database operations (*reads and/or writes*) and constitutes the *smallest logical unit of work* from the application perspective.

### Properties

Each transaction  $T$  has the following properties:

- **Atomicity** either all effects of  $T$  are reflected in the database or none
- **Consistency**  $T$  guarantees to yield the correct state of the database
- **Isolation**  $T$  is isolated from the effect of concurrent transactions
- **Durability** after a commit of  $T$ , its effects are permanent

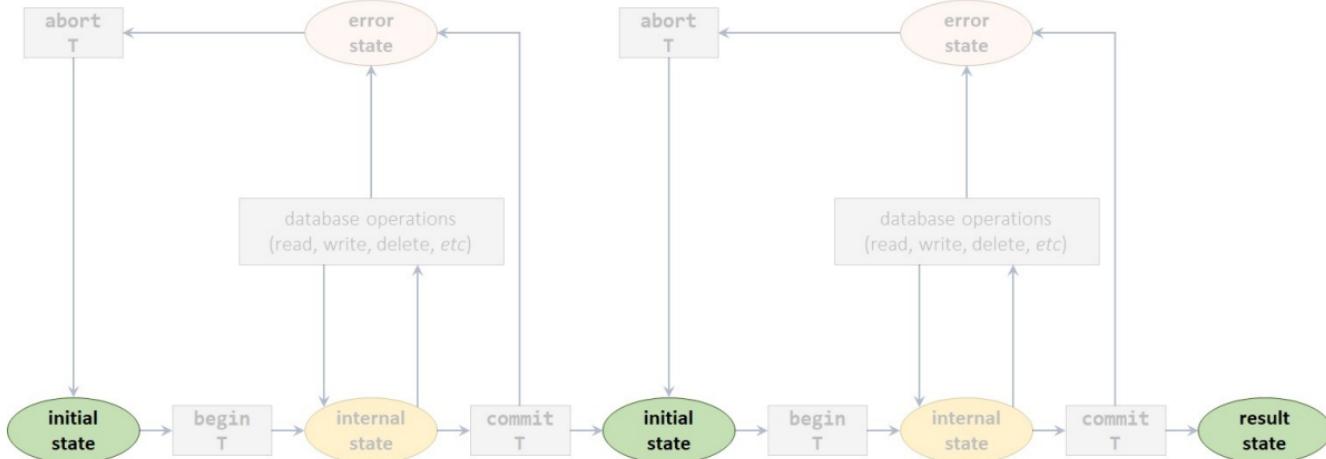
**ACID** properties of transactions

## Transactions

### Transition Graph

### Addendum

What you will see as a database user is the sequence of **consistent** state depicted in green. In between, the database may go into a **temporary inconsistent state**. Here, we define consistency as satisfying all user specified constraints.



# DBMS

Database  
Challenges  
Management  
**Transactions**  
- Properties  
- Transition  
- **Sequential**  
- Concurrent  
Serialization  
Architecture  
Study

## Transactions

### Sequential

```
define transaction T(U, Amount) {  
    BEGIN  
        Bal1 = READ(U)          -- read balance  
        Bal2 = Bal1 + Amount   -- update balance  
        WRITE(U, Bal2)         -- write balance  
    COMMIT  
}
```

#### Note

- Initial balance of U = 'A' to be \$1000
- Two transactions:
  - A. T1 = T('A', 500)
  - B. T2 = T('A', 100)

### Serial

T('A', 500)	T('A', 100)
BEGIN	
Bal1=READ('A')	
Bal2=Bal1+500	
WRITE('A',Bal2)	
COMMIT	
	BEGIN
	Bal1=READ('A')
	Bal2=Bal1+100
	WRITE('A',Bal2)
	COMMIT
:	:

**Expected** \$1600

**Actual** \$1600

# Transactions

## Lost Update

T('A', 500)	T('A', 100)
BEGIN	
Bal1=READ('A')	
Bal2=Bal1+500	
	BEGIN
	Bal1=READ('A')
	Bal2=Bal1+100
	WRITE('A',Bal2)
	COMMIT
WRITE('A',Bal2)	
COMMIT	
:	:

Expected \$1600  
 Actual \$1500

## Dirty Read

T('A', 500)	T('A', 100)
BEGIN	
Bal1=READ('A')	
Bal2=Bal1+500	
WRITE('A',Bal2)	
	BEGIN
	Bal1=READ('A')
	Bal2=Bal1+100
	WRITE('A',Bal2)
	COMMIT
ABORT	
:	:

Expected \$1100  
 Actual \$1600

## Unrepeatable Read

T('A', 500)	T('A', 100)
BEGIN	
Bal1=READ('A')	
	BEGIN
	Bal1=READ('A')
	Bal2=Bal1+100
	WRITE('A',Bal2)
	COMMIT
Bal1=READ('A')	
Bal2=Bal1+500	
WRITE('A',Bal2)	
ABORT	
:	:

1<sup>st</sup> read \$1000  
 2<sup>nd</sup> read \$1100

## Serialization

### Equivalent Transaction

#### Definition

Two executions are **equivalent** if they have the **same effect** on the data.

### Serializability

#### Definition

A concurrent execution of a set of transaction is **serializable** if this execution is **equivalent to some** serial execution of the same set of transactions.

### DBMS

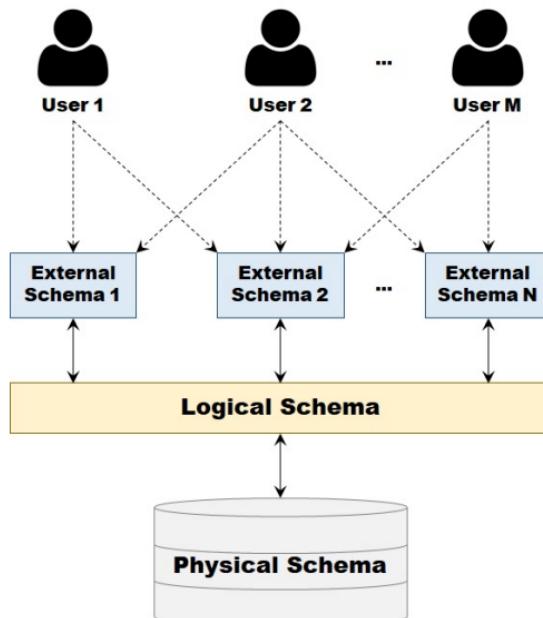
- Support concurrent executions of transactions to optimize performance
- Enforce serializability of concurrent executions to ensure integrity of data

# DBMS

Database  
Challenges  
Management  
Transactions  
Serialization  
**Architecture**  
- Abstraction  
- Independence  
Study

## 3-Tier Architecture of DBMS

### Levels of Data Abstraction



#### External Schema

- User or group-specific view on the data

#### Logical Schema

- Logical organization of data (*e.g., relations/tables, objects, graphs*)
- Unified representation of all data
- Support of **logical** data independence and **physical** data independence

#### Physical Schema

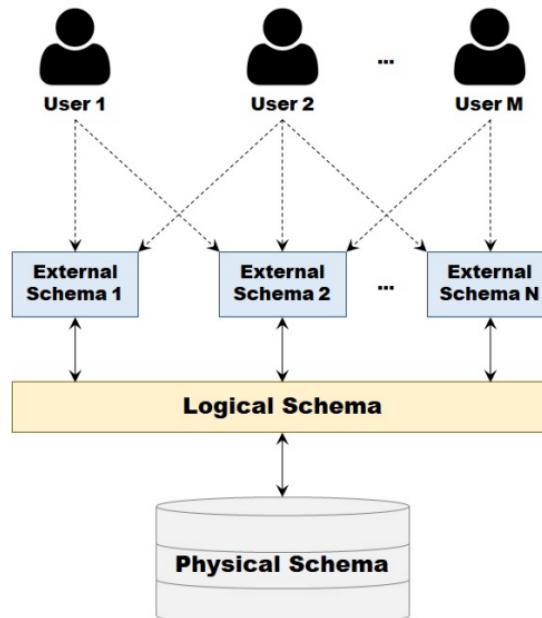
- Organization of data on disk and in memory
- Database as collections of fields, arrays, records, files, pages, etc

# DBMS

Database  
Challenges  
Management  
Transactions  
Serialization  
**Architecture**  
- Abstraction  
- Independence  
Study

## 3-Tier Architecture of DBMS

### Data Independence



### Logical Data Independence

- Ability to change *logical schema* without affecting *external schemas* (e.g., adding/deleting/updating attributes, changing data types, changing data model)

### Physical Data Independence

- Representation of data being independent from physical schema
- Physical schema can be changed without affecting logical schema (e.g., creating indexes, new caching strategies, different storage devices)

# DBMS

Database  
Challenges  
Management  
Transactions  
Serialization  
Architecture  
**Study**

## Study of DBMS

Scope

Topic

### Database Design

- How to model the data requirements
- How to organize data using a DBMS

Relational Model  
ER Model  
Schema Refinement

### Database Programming

- How to create, query, and update a database
- How to specify data constraints
- How to use SQL in applications

Relational Algebra  
SQL

### DBMS Implementation

- How to build a DBMS

CS3223

*out of scope*

# Relational Model

---

# Relational Model

History  
Popularity  
Model  
Integrity  
Key  
Foreign Key  
Considerations

## History of DBMS

### "Historical" Models

- Hierarchical Model
- Network Model

### *Relational Model*

- Commercial RDBMS
- Open-Source RDBMS

### Object-Oriented Model

- Native OO Model (*e.g., Objectstore, 1988*)
- Object-Relational Model\*

### More Recent Development

- NoSQL Model; In-Memory DBMS

### Commercial Systems



### Open-Source Systems



PostgreSQL



\*Supported by most DBMS

# Relational Model

History  
**Popularity**  
- Table  
- Graph  
- Job  
Model  
Integrity  
Key  
Foreign Key  
Considerations

## Popularity

### Table

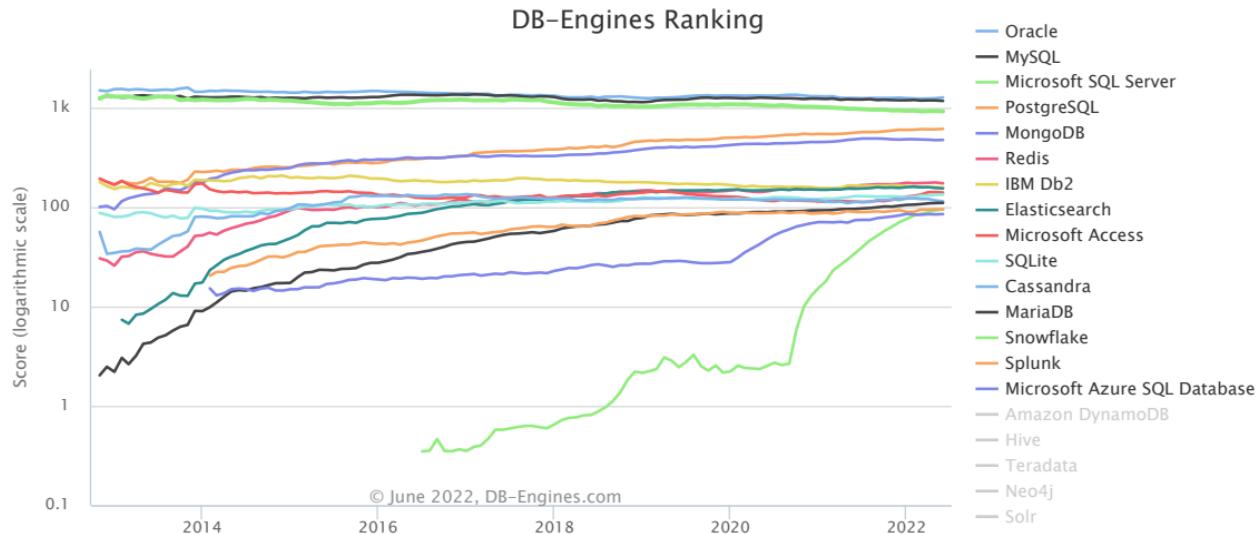
Jun 2022	Rank		DBMS	Database Model	Score		
	Jun 2022	May 2022			Jun 2022	May 2022	Jun 2021
1.	1.	1.	Oracle	Relational	1287.74	+24.92	+16.80
2.	2.	2.	MySQL	Relational	1189.21	-12.89	-38.65
3.	3.	3.	Microsoft SQL Server	Relational	933.83	-7.37	-57.25
4.	4.	4.	PostgreSQL	Relational	620.84	+5.55	+52.32
5.	5.	5.	MongoDB	Document	480.73	+2.49	-7.49
6.	6.	7. 	Redis	Key-value	175.31	-3.71	+10.06
7.	7.	6. 	IBM DB2	Relational	159.19	-1.14	-7.85
8.	8.	8.	Elasticsearch	Search engine	156.00	-1.70	+1.29
9.	9.	10. 	Microsoft Access	Relational	141.82	-1.62	+26.88
10.	10.	9. 	SQLite	Relational	135.44	+0.70	+4.90

#Source: [DB-Engines](#)

# Relational Model

History  
Popularity  
- Table  
- Graph  
- Job  
Model  
Integrity  
Key  
Foreign Key  
Considerations

## Popularity Graph



#Source: [DB-Engines](#)

# Relational Model

History  
**Popularity**

- Table  
- Graph  
- Job

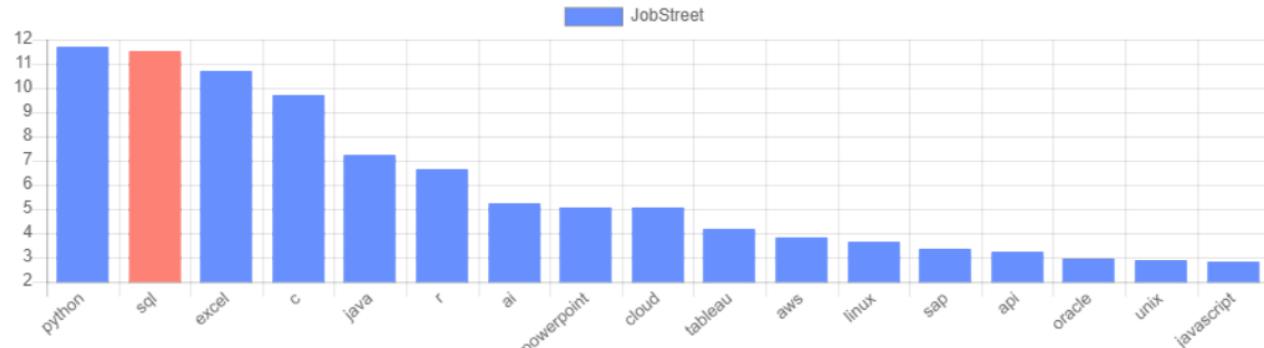
Model  
Integrity  
Key

Foreign Key  
Considerations

## Popularity

### Analysis of Job Descriptions

- 15k+ job offers from [JobStreet](#)  
(*data analyst, data engineer, data scientist*)
- Quick & dirty keyword extraction  
(*by Dr. Christian Von Der WETH*)



### Anecdotes

**Java, SQL and Python are the most in-demand digital skills**

**Key Skill: SQL, Because Companies are Obsessed with Data**

**Want a Job in Data? Learn SQL.**

# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- Relational Schema  
- Relational Database Integrity  
Key  
Foreign Key Considerations

## The Relational Model

### Preliminary

- Proposed by Edgar F. Codd in 1970
- **Basic Concept:** *Relations (tables with rows and columns)*

Table "Employees"

ID	Name	DOB	Salary
1	Alice	10-08-1988	7,500
2	Bob	06-11-2001	4,800
3	Carol	25-02-1995	5,500

### Addendum

When it is clear from the context, we often omit the data type (e.g., *INT* or *TEXT*). Additionally, we may omit the data type in favor of simplicity and assume that the data type is the most "*logical*" data type.

### A Relational Model of Data for Large Shared Data Banks

E. F. CODD  
IBM Research Laboratory, San Jose, California

Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation). A prompting service which supplies such information is not a satisfactory solution. Activities of users at terminals and most application programs should remain unaffected when the internal representation of data is changed and even when some aspects of the external representation are changed. Changes in data representation will often be needed as a result of changes in query, update, and report traffic and natural growth in the types of stored information.

- **Relation Schema:** *Definition of a Relation*
  - Specifies the attributes (*columns*) and data constraints (e.g., *domain constraints*)

`Employees(id:INT, name:TEXT, dob:DATE, salary:NUMERIC)`

# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- Relational Schema  
- Relational Database  
Integrity  
Key  
Foreign Key  
Considerations

## The Relational Model

### Domain

#### Definition

A **domain** is a set of *atomic* values (e.g., INT, NUMERIC, TEXT).

- We denote the domain of attribute  $A_i$  as  $\text{dom}(A_i)$  (i.e., the set of possible values of  $A_i$ )
- Each value  $v$  of attribute  $A_i$  is either:  $v \in \text{dom}(A_i)$  or  $v = \text{NULL}$ 
  - $\text{NULL}$  is a special value indicating that the value  $v$  is *unknown*, *invalid*, etc

### Relation

#### Definition

A **relation** is a set of tuple (i.e., records).

- Simplified notation of **relation schema** where the domain is clear from context is  $R(A_1, A_2, \dots, A_n)$
- Each **instance** of schema  $R$  is a relation which is a subset of  $\{(a_1, a_2, \dots, a_n) \mid a_i \in \text{dom}(A_i) \cup \{\text{NULL}\}\}$
- We write  $R.A_i$  to refer to the attribute  $A_i$  of the relation  $R$

# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- Relational Schema  
- Relational Database  
Integrity  
Key  
Foreign Key  
Considerations

## The Relational Model

### Domain & Relation

#### Relational Schema

- *Modules(course, mc, exam)*
  - $\text{dom}(\text{course}) = \{\text{cs2102, cs3223, cs4221}\}$
  - $\text{dom}(\text{mc}) = \{2, 4\}$
  - $\text{dom}(\text{exam}) = \{\text{yes, no}\}$

#### Instance

- Subset of:  
 $\{\text{cs2102, cs3223, cs4221, NULL}\}$   
 $\times \{2, 4, \text{NULL}\}$   
 $\times \{\text{yes, no, NULL}\}$

### Addendum

The symbol  $\times$  on set is a [Cartesian product](#). Do not forget to include the **NULL** value when computing the maximum possible cardinality.

course	mc	exam
cs2102	2	yes
cs2102	2	no
cs2102	2	NULL
cs2102	4	yes
cs2102	4	no
:	:	:
NULL	NULL	NULL

# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- Relational Schema  
- Relational Database  
Integrity  
Key  
Foreign Key  
Considerations

## The Relational Model

### Domain & Relation

#### Exercise #1

Assume a relation  $R(A, B)$  with

- $\text{dom}(A) = \{x, y, z\}$
- $\text{dom}(B) = \{1, 2, 3, 4\}$

Which tuples on the right **violate** the definition of relation  $R$ ?

A	B	Row
x	4	Row 1
z	4	Row 2
NULL	2	Row 3
NULL	0	Row 4

	Choice	Comment	
A	Row 1		?
B	Row 2		?
C	Row 3		?
D	Row 4		?



# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- Relational Schema  
- Relational Database  
Integrity  
Key  
Foreign Key  
Considerations

## The Relational Model

### Domain & Relation

#### Exercise #2

Assume a relation  $R(A, B)$  with

- $\text{dom}(A) = \{x, y, z\}$
- $\text{dom}(B) = \{1, 2, 3, 4\}$

Which tuples on the right **violate** the definition of relation  $R$ ?

A	B	Row
NULL	NULL	Row 1
x	4	Row 2
x	y	Row 3
z	1	Row 4

	Choice	Comment	
A	Row 1		?
B	Row 2		?
C	Row 3		?
D	Row 4		?



# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- Relational Schema  
- Relational Database  
Integrity  
Key  
Foreign Key  
Considerations

## The Relational Model

### Domain & Relation

#### Exercise #3

Assume a relation  $R(A, B)$  with

- $\text{dom}(A) = \{x, y, z\}$
- $\text{dom}(B) = \{1, 2, 3, 4\}$

Which tuples on the right **violate** the definition of relation  $R$ ?

A	B	Row
x	4	Row 1
y	1	Row 2
x	1	Row 3
x	4	Row 4

	Choice	Comment	
A	Row 1		?
B	Row 2		?
C	Row 3		?
D	Row 4		?



# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- **Relational Schema**  
- Relational Database Integrity  
Key  
Foreign Key Considerations

## The Relational Model

### Relational Database Schema

#### Definition

A **relational database schema** is a set of relation schemas + data constraints.

#### Example

Movies(

*id* : INT  
*title* : TEXT  
*genre* : TEXT  
*opened* : DATE

)

Casts(

*movie\_id* : INT  
*actor\_id* : INT  
*role* : TEXT

)

Actors(

*id* : INT  
*name* : TEXT  
*dob* : DATE

)

<sup>#</sup>When the data type is clear from context, we often omit. If it is omitted, we choose the most "logical" type.

# Relational Model

History  
Popularity  
**Model**  
- Preliminary  
- Domain & Relation  
- Relational Schema  
- **Relational Database**  
Integrity  
Key  
Foreign Key  
Considerations

## The Relational Model

### Relational Database

#### Definition

A **relational database** is a collections of tables.

#### Example

Table "Movies"

<b>id</b>	<b>title</b>	<b>genre</b>	<b>opened</b>
101	Aliens	action	1986
102	Logan	drama	2017
103	Heat	crime	1995
104	Terminator	action	1984
:	:	:	:

#### Addendum

Remember, this is a relation (*or relation instances*). The name "relation" comes from the fact that each row (*or tuple*), relate the values in the tuple. So 101 is related to Aliens, action, and 1986.

Table "Casts"

<b>movie_id</b>	<b>actor_id</b>	<b>role</b>
101	20	Ellen Ripley
101	23	Private Hudson
101	54	Corporal Hicks
102	21	Logan
104	23	Punk Leader
:	:	:

Table "Actors"

<b>id</b>	<b>name</b>	<b>dob</b>
20	Sigourney Weaver	08-10-1949
21	Hugh Jackman	12-10-1968
22	Tom Hanks	09-07-1956
23	Bill Paxton	17-05-1955
:	:	:

# Relational Model

History  
Popularity  
Model  
**Integrity**  
- Motivation  
- Constraints  
Key  
Foreign Key  
Considerations

## Data Integrity

### Motivation

#### Problem

Consider the *relational database* below. What are the problems with the tables?

### Relational Database

Table "Movies"

<b>id</b>	<b>title</b>	<b>genre</b>	<b>opened</b>
101	Aliens	action	1986
101	Logan	drama	2017
103	Heat	crime	1995
104	Terminator	action	1984

Table "Casts"

<b>movie_id</b>	<b>actor_id</b>	<b>role</b>
101	20	Ellen Ripley
101	23	Private Hudson
101	54	Corporal Hicks
102	21	Logan
104	23	Punk Leader

Table "Actors"

<b>id</b>	<b>name</b>	<b>dob</b>
20	Sigourney Weaver	08-10-2049
21	Hugh Jackman	12-10-1968
NULL	Tom Hanks	09-07-1956
23	Bill Paxton	17-05-1955

#In this example, what you see is the entire database.

## Addendum

1. The id 101 appears multiple times in Movies (*logically, the name id should be unique*).
2. The movie\_id 102 in Casts has no corresponding id in Movies.
3. Sigourney Weaver is born in the future?
4. Tom Hanks has no id?

# Relational Model

History  
Popularity  
Model  
**Integrity**  
- Motivation  
- Constraints  
Key  
Foreign Key  
Considerations

## Data Integrity

### Integrity Constraints

#### Definition

**Integrity constraints** are conditions that restrict what constitutes valid data.

#### Structural

- Independent of application
- Domain constraints (*e.g., cannot store TEXT in INT column*)\*
- Key constraints
- Foreign key constraints

#### General

- Dependent of application
- Check constraints
- Triggers

Table "Movies"

id	title	genre	opened
101	Aliens	Action	1986
101	Logan	Drama	2017

Table "Actors"

id	name	dob
20	Sigourney Weaver	08-10-2049
NULL	Tom Hanks	09-07-1956

\*Some automatic type conversion may be performed, but for simplicity, we often ignore this.



# Relational Model

History  
Popularity

Model  
Integrity

Key

- Superkey & Key
- Candidate & Primary
- Exercise

Foreign Key  
Considerations

## Key Constraints

### Superkey

#### Definition

A **superkey** is a subset of attributes that uniquely identifies a tuple in a relation.

### Key

#### Definition

A **key** is a **superkey** that is also minimal (i.e., no proper subset of the key is a superkey).

#### NOTE:

- Minimal (*cannot be made smaller*)  $\neq$  Minimum (*the smallest*)

#### Properties

- If  $(A, B, C)$  is definitely a **superkey** then  $(A, B, C, D)$  is a **superkey**
- If  $(A, B, C)$  is definitely a **key** then  $(A, B, C)$  is a **superkey**
- If  $(A, B)$  is definitely a **superkey** then  $(A, B, C)$  is **NOT** a **key**
- If  $(A, B)$  is definitely a **key** then it is possible that  $(B, C, D)$  is also a **key**
- Every relations have at least 1 superkey



Each table has exactly 1 primary key.



# Relational Model

History  
Popularity  
Model  
Integrity  
**Key**  
- Superkey & Key  
- Candidate & Primary  
- Exercise  
Foreign Key  
Considerations

## Key Constraints

### Candidate Keys

#### Definition

The **candidate keys** is a the set of all keys of a given relation.

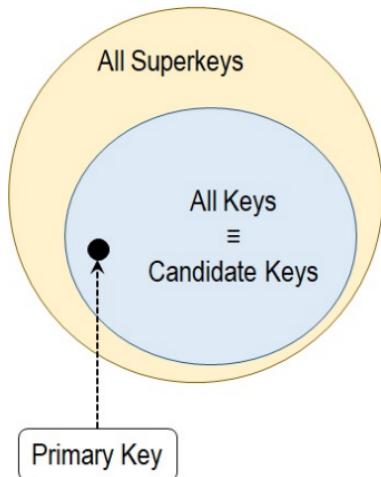
### Primary Key

#### Definition

A **primary key** is a selected candidate keys.

- We also add a constraint that primary key attributes cannot be **NULL**
- To simplify our notation, we underline the primary key in the schema notation

*Movies(id: **INT**, title: **TEXT**, genre: **TEXT**, opened: **DATE**)*



<sup>#</sup>We often say "a candidate key" to also mean "a key" and vice versa.

# Relational Model

History  
Popularity  
Model  
Integrity  
**Key**  
- Superkey & Key  
- Candidate & Primary  
- **Exercise**  
Foreign Key  
Considerations

## Key Constraints

### Exercise #4

Consider a forum database with the following relation filled with many thousands of users: Accounts(*email*: TEXT, *password*: TEXT, *name*: TEXT).

Using reasonable guess, which subsets of attributes are **superkeys** of relation Accounts?

Choice	Comment	
A {email}		?
B {pwd}		?
C {name}		?
D {email, pwd}		?
E {email, name}		?
F {pwd, name}		?
G {email, pwd, name}		?



# Relational Model

History  
Popularity  
Model  
Integrity  
**Key**  
- Superkey & Key  
- Candidate & Primary  
- Exercise  
Foreign Key  
Considerations

## Key Constraints

### Exercise #5

Consider a forum database with the following relation filled with many thousands of users: Accounts(*email*: TEXT, *password*: TEXT, *name*: TEXT).

Using reasonable guess, which subsets of attributes are **keys** of relation Accounts?

Choice	Comment	
A {email}		?
B {pwd}		?
C {name}		?
D {email, pwd}		?
E {email, name}		?
F {pwd, name}		?
G {email, pwd, name}		?



# Relational Model

History  
Popularity  
Model  
Integrity  
Key  
**Foreign Key**  
- Definition  
- Exercise  
Considerations

## Foreign Key Constraints

### Foreign Key

#### Definition

A **foreign key** is a subset of attributes of relation  $R_1$  that refers to the primary key\* of relation  $R_2$ .

- We say that  $R_1$  is the **referencing relation**
- We say that  $R_2$  is the **referenced relation**

Informally, we write  $(R_1.A_{j1}, R_1.A_{j1}, \dots) \rightsquigarrow (R_2.A_{j1}, R_2.A_{j2}, \dots)$ .

#### Requirement

Each foreign key in  $R_1$  must satisfy one of the following:

- Appear as **primary key** in  $R_2$
- Be a **NULL** value (*or a tuple containing at least one **NULL** value*)

## Addendum

In the tables below, we have:

- (Cast.movie\_id)  $\rightsquigarrow$  (Movies.id)
- (Cast.actor\_id)  $\rightsquigarrow$  (Actors.id)

Table "Movies"

<u>id</u>	title	genre	opened
101	Aliens	Action	1986
102	Logan	Drama	2017
:	:	:	:

Table "Cast"

movie_id	actor_id	role
101	20	Ellen Ripley
102	21	Logan
:	:	:

Table "Actors"

<u>id</u>	name	dob
20	Sigourney Weaver	08-10-1949
21	Hugh Jackman	12-10-1968
:	:	:

\*In most DBMS, it can actually refer to any **candidate keys**.



# Relational Model

History  
Popularity  
Model  
Integrity  
Key  
**Foreign Key**  
- Definition  
- Exercise  
Considerations

## Foreign Key Constraints

### Foreign Key

#### Definition

A **foreign key** is a subset of attributes of relation  $R_1$  that refers to the primary key\* of relation  $R_2$ .

- We say that  $R_1$  is the **referencing relation**
- We say that  $R_2$  is the **referenced relation**

Informally, we write  $(R_1.A_{j1}, R_1.A_{j1}, \dots) \rightsquigarrow (R_2.A_{j1}, R_2.A_{j2}, \dots)$ .

#### Notes

- A **referencing relation** can be a **referenced relation** for different foreign key
- The **referencing relation** and the **referenced relation** can be the same relation

#### Requirement

Each foreign key in  $R_1$  must satisfy one of the following:

- Appear as **primary key** in  $R_2$
- Be a **NULL** value (*or a tuple containing at least one NULL value*)

Table "Employees"

<u>id</u>	name	dob	salary	manager
1	Alice	10-08-1988	7500	NULL
2	Bob	06-11-2001	4800	3
3	Carol	25-02-1995	5500	1

\*In most DBMS, it can actually refer to any **candidate keys**.

#  $(\text{Employees}.\text{manager}) \rightsquigarrow (\text{Employees}.\text{id})$

# Relational Model

History  
Popularity  
Model  
Integrity  
Key  
**Foreign Key**  
- Definition  
- Exercise  
Considerations

## Foreign Key Constraints

### Exercise #6

Consider two tables  $R(A, B)$  and  $S(C, D)$  with  $S.D \rightsquigarrow R.A$ .

Which tuples in Table "S" **violate** the **foreign key** constraint?

Table "R"

A	B
x	4
z	2
y	NULL

Table "S"

C	D	Row
a	w	Row 1
b	x	Row 2
c	y	Row 3
d	NULL	Row 4

	Choice	Comment	
A	Row 1		?
B	Row 2		?
C	Row 3		?
D	Row 4		?

# Relational Model

History  
Popularity  
Model  
Integrity  
Key  
Foreign Key  
**Considerations**  
- Limitations  
- Practical

## Considerations

### Limitations

- Structural integrity constraints covers *application-independent* constraints
  - Limiting domain to valid values
  - Checking existence of values in another relations
- **Not covered:** *application-dependent* constraints derived from deeper semantics of the data
  - Some checking can be done using `CHECK` keyword
  - For example, checking salary not being negative

Table "Employees"

<u>id</u>	<u>name</u>	<u>dob</u>	<u>salary</u>	<u>manager</u>
1	Alice	10-08-1988	7500	NULL
2	Bob	06-11-2001	-4800	3
3	Carol	25-02-1995	5500	1
4	Dave	18-06-1999	6000	NULL

# Relational Model

History  
Popularity  
Model  
Integrity  
Key  
Foreign Key  
**Considerations**  
- Limitations  
- Practical

## Considerations

### Practical Considerations

- Integrity constraints are *optional* and not mandatory
  - It simply allows the check to be done by DBMS instead of application
  - In other words, no matter who uses the database, its entries will always be consistent with respect to the constraints (*assuming they cannot remove the constraints*)
- Integrity constraints may affect *performance*
  - May decrease performance due to additional checking
  - May improve performance due to creation of indexes that speeds up retrieval

# Summary

---

# Summary

## Cheat Sheet

- *Relational Model*

- *Constraints*

Summary

## Cheat Sheet

### Relational Model

Table "Movies"

<b>id</b>	<b>title</b>	<b>genre</b>	<b>opened</b>	...
101	Aliens	action	1986	...
102	Logan	drama	2017	...
103	Heat	crime	1995	...
⋮	⋮	⋮	⋮	⋮

### Informal Description

<b>Term</b>	<b>Description</b>
<i>Attribute</i>	column of a table
<i>Domain</i>	possible values
<i>Attribute Value</i>	element of a domain
<i>Relation Schema</i>	set of attributes + constraints + relation name
<i>Tuple</i>	row of a table
<i>Relation</i>	set of tuples
<i>Cardinality</i>	number of rows
<i>Database Schema</i>	set of relation schema
<i>Database</i>	set of relations

# Summary

## Cheat Sheet

- Relational Model

- Constraints

Summary

## Cheat Sheet

### Relational Model

The diagram illustrates the structure of a relational database table. At the top, the word "Attribute" is written above four arrows pointing downwards towards the columns of a table. Below the arrows, the text "Table 'Movies'" is written. The table itself has a light gray border and contains five columns labeled "id", "title", "genre", "opened", and "...". The first three columns have data entries: "101" in "id", "Aliens" in "title", and "action" in "genre". The fourth column has data entries: "1986" in "opened" and "2017" in the next row. The fifth column has data entries: "..." in "id", "Logan" in "title", "drama" in "genre", and "1995" in "opened". Ellipses (...) are used to indicate that the table continues.

id	title	genre	opened	...
101	Aliens	action	1986	...
102	Logan	drama	2017	...
103	Heat	crime	1995	...
...	...	...	...	...

### Informal Description

**Attributes:** Column of a table

# Summary

## Cheat Sheet

- Relational Model

- Constraints

Summary

## Cheat Sheet

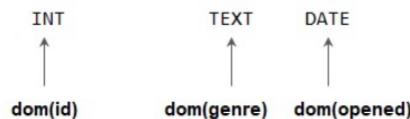
### Relational Model

#### Informal Description

| **Domain:** Possible values

Table "Movies"

<b>id</b>		<b>genre</b>	<b>opened</b>	
101	Aliens	action	1986	...
102	Logan	drama	2017	...
103	Heat	crime	1995	...
:	:	:	:	...



# Summary

## Cheat Sheet

- *Relational Model*

- *Constraints*

Summary

## Cheat Sheet

Relational Model

Informal Description

**Attribute Values:** Elements of a domain

Table "Movies"

101	Aliens	action	1986	...
102	Logan	drama	2017	...
103	Heat	crime	1995 ←	...
...	...	...	...	...

Attribute Value

DATE

# Summary

## Cheat Sheet

- Relational Model

- Constraints

Summary

## Cheat Sheet

### Relational Model

**Full:** Movies(*id*: INT, *title*: TEXT, *genre*: TEXT, *opened*: DATE)

Table "Movies"

**Simplified:** Movies(*id*, *title*, *genre*, *opened*)

<b>id</b>	<b>title</b>	<b>genre</b>	<b>opened</b>	...
101	Aliens	action	1986	...
102	Logan	drama	2017	...
103	Heat	crime	1995	...
...	...	...	...	...

} ← Relation Schema

### Informal Description

**Relation Schema:** Set of attributes + Constraints + Relation name

# Summary

## Cheat Sheet

- *Relational Model*

- *Constraints*

Summary

## Cheat Sheet

### Relational Model

#### Informal Description

| **Tuple:** Row of a table

Table "Movies"

101	Aliens	action	1986	...	...
102	Logan	drama	2017	...	...
103	Heat	crime	1995	...	...
...	...	...	...	...	...

← Tuple/Record

# Summary

## Cheat Sheet

- *Relational Model*

- *Constraints*

Summary

## Cheat Sheet

Relational Model

Informal Description

| **Relation:** Set of tuples

Table "Movies"

101	Aliens	action	1986	...
102	Logan	drama	2017	...
103	Heat	crime	1995	...
⋮	⋮	⋮	⋮	⋮

A curly brace on the right side of the table is labeled "Relation" with an arrow pointing to it.

# Summary

Cheat Sheet  
- Relational Model  
- Constraints  
Summary

## Cheat Sheet

### Relational Model

#### Informal Description

**Cardinality:** Number of rows

Table "Movies"

101	Aliens	action	1986	...
102	Logan	drama	2017	...
103	Heat	crime	1995	...
⋮	⋮	⋮	⋮	⋮

A diagram showing a table with five columns and four rows. The last cell in each row contains three dots (...). A curly brace on the right side of the table spans all five columns and all four rows. An arrow points from the text "Cardinality # of rows" to this brace.

# Summary

## Cheat Sheet

- Relational Model

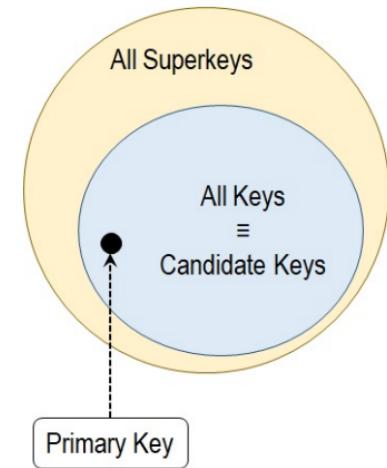
- **Constraints**

Summary

## Cheat Sheet

### Constraints

Term	Informal Description
<i>Superkey</i>	Set of attributes that <u>uniquely</u> identifies a tuple
<i>Key</i>	Minimal superkey ( <i>no subset is a superkey</i> )
<i>Candidate Keys</i>	Set of <b>ALL</b> keys
<i>Primary Key</i>	One selected key
<i>Prime Attribute</i>	Attributes of a candidate key



### DB vs DBS vs DBMS

$$\boxed{\text{DBS} = \text{DBMS} + (n \times \text{DB})} \quad (\text{for } n > 0)$$

<sup>#</sup>DB = Database ; DBS = Database Systems (not schema) ; DBMS = Database Management System

# Summary

Cheat Sheet  
**Summary**

## Summary

### Advantages of DBMS\*

- Transactions with *ACID* properties
- Level of abstraction for *data independence*

### Relational Model

- Unified representation of all data as relations (*i.e., tables*)
- (*Structural*) integrity constraints to specify restrictions on what constitutes correct/valid data

\*For large-scale data management compared to "files only".

```
postgres=# exit
```

```
Press any key to continue . . .
```

# Solutions

## Exercise #1

Exercise #2

Exercise #3

Exercise #4

Exercise #5

Exercise #6

## The Relational Model

### Domain & Relation

#### Exercise #1

Assume a relation  $R(A, B)$  with

- $\text{dom}(A) = \{x, y, z\}$
- $\text{dom}(B) = \{1, 2, 3, 4\}$

Which tuples on the right **violate** the definition of relation  $R$ ?

A	B	Row
x	4	Row 1
z	4	Row 2
NULL	2	Row 3
NULL	0	Row 4

	Choice	Comment	
<b>A</b>	Row 1	NO: $x \in \text{dom}(A) \& 4 \in \text{dom}(B)$	✗
<b>B</b>	Row 2	NO: $z \in \text{dom}(A) \& 4 \in \text{dom}(B)$	✗
<b>C</b>	Row 3	NO: $A = \text{NULL} \& 2 \in \text{dom}(B)$	✗
<b>D</b>	Row 4	YES: $A = \text{NULL} \& \emptyset \notin \text{dom}(B) \wedge \text{not } \text{NULL}$	✓

# Solutions

Exercise #1

Exercise #2

Exercise #3

Exercise #4

Exercise #5

Exercise #6

## The Relational Model

### Domain & Relation

#### Exercise #2

Assume a relation  $R(A, B)$  with

- $\text{dom}(A) = \{x, y, z\}$
- $\text{dom}(B) = \{1, 2, 3, 4\}$

Which tuples on the right **violate** the definition of relation  $R$ ?

A	B	Row
NULL	NULL	Row 1
x	4	Row 2
x	y	Row 3
z	1	Row 4

	Choice	Comment	
<b>A</b>	Row 1	NO: $A = \text{NULL} \wedge B = \text{NULL}$	✗
<b>B</b>	Row 2	NO: $x \in \text{dom}(A) \wedge 4 \in \text{dom}(B)$	✗
<b>C</b>	Row 3	YES: $x \in \text{dom}(A) \wedge y \notin \text{dom}(B) \wedge \text{not } \text{NULL}$	✓
<b>D</b>	Row 4	NO: $z \in \text{dom}(A) \wedge 1 \in \text{dom}(B)$	✗

# Solutions

Exercise #1

Exercise #2

**Exercise #3**

Exercise #4

Exercise #5

Exercise #6

## The Relational Model

### Domain & Relation

#### Exercise #3

Assume a relation  $R(A, B)$  with

- $\text{dom}(A) = \{x, y, z\}$
- $\text{dom}(B) = \{1, 2, 3, 4\}$

Which tuples on the right **violate** the definition of relation  $R$ ?

A	B	Row
x	4	Row 1
y	1	Row 2
x	1	Row 3
x	4	Row 4

Choice	Comment	
A	Row 1	LIKELY: Duplicate with Row 4
B	Row 2	NO: $y \in \text{dom}(A)$ & $1 \in \text{dom}(B)$
C	Row 3	NO: $x \in \text{dom}(A)$ & $1 \in \text{dom}(B)$
D	Row 4	LIKELY: Duplicate with Row 1

# Solutions

Exercise #1

Exercise #2

Exercise #3

**Exercise #4**

Exercise #5

Exercise #6

## Key Constraints

### Exercise #4

Consider a forum database with the following relation filled with many thousands of users: Accounts(*email*: **TEXT**, *password*: **TEXT**, *name*: **TEXT**).

Using reasonable guess, which subsets of attributes are **superkeys** of relation Accounts?

Choice	Comment	
A { <i>email</i> }	LIKELY: <i>email</i> is likely unique	✓
B { <i>pwd</i> }	UNLIKELY: two people may have the same password	✗
C { <i>name</i> }	UNLIKELY: unless this is username	✗
D { <i>email</i> , <i>pwd</i> }	LIKELY: it contains <i>email</i>	✓
E { <i>email</i> , <i>name</i> }	LIKELY: it contains <i>email</i>	✓
F { <i>pwd</i> , <i>name</i> }	UNLIKELY: does not contain <i>email</i>	✗
G { <i>email</i> , <i>pwd</i> , <i>name</i> }	LIKELY: it contains <i>email</i>	✓

# Solutions

Exercise #1

Exercise #2

Exercise #3

Exercise #4

**Exercise #5**

Exercise #6

## Key Constraints

### Exercise #5

Consider a forum database with the following relation filled with many thousands of users: Accounts(*email*: **TEXT**, *password*: **TEXT**, *name*: **TEXT**).

Using reasonable guess, which subsets of attributes are **keys** of relation Accounts?

Choice	Comment	
A { <i>email</i> }	LIKELY: <i>email</i> is likely unique	<input checked="" type="checkbox"/>
B { <i>pwd</i> }	UNLIKELY: two people may have the same password	<input type="checkbox"/>
C { <i>name</i> }	UNLIKELY: unless this is username	<input type="checkbox"/>
D { <i>email</i> , <i>pwd</i> }	UNLIKELY: it contains <i>email</i>	<input type="checkbox"/>
E { <i>email</i> , <i>name</i> }	UNLIKELY: it contains <i>email</i>	<input type="checkbox"/>
F { <i>pwd</i> , <i>name</i> }	UNLIKELY: does not contain <i>email</i>	<input type="checkbox"/>
G { <i>email</i> , <i>pwd</i> , <i>name</i> }	UNLIKELY: it contains <i>email</i>	<input type="checkbox"/>

# Solutions

Exercise #1

Exercise #2

Exercise #3

Exercise #4

Exercise #5

**Exercise #6**

## Foreign Key Constraints

### Exercise #6

Consider two tables  $R(A, B)$  and  $S(C, D)$  with  $S.D \rightsquigarrow R.A$ .

Which tuples in Table "S" **violate** the **foreign key** constraint?

Table "R"

A	B
x	4
z	2
y	NULL

Table "S"

C	D	Row
a	w	Row 1
b	x	Row 2
c	y	Row 3
d	NULL	Row 4

### Choice

### Comment

A	Row 1	YES: w not in Table "R"	✓
B	Row 2	NO: x is in Table "R"	✗
C	Row 3	NO: y is in Table "R"	✗
D	Row 4	NO: NULL is allowed	✗