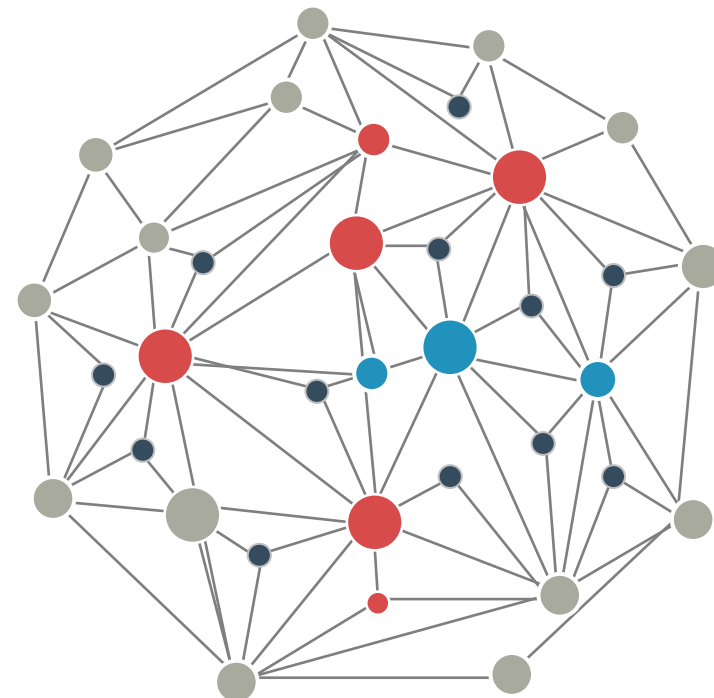2021-2022 ZJE Database and Software Technology 2

# Database
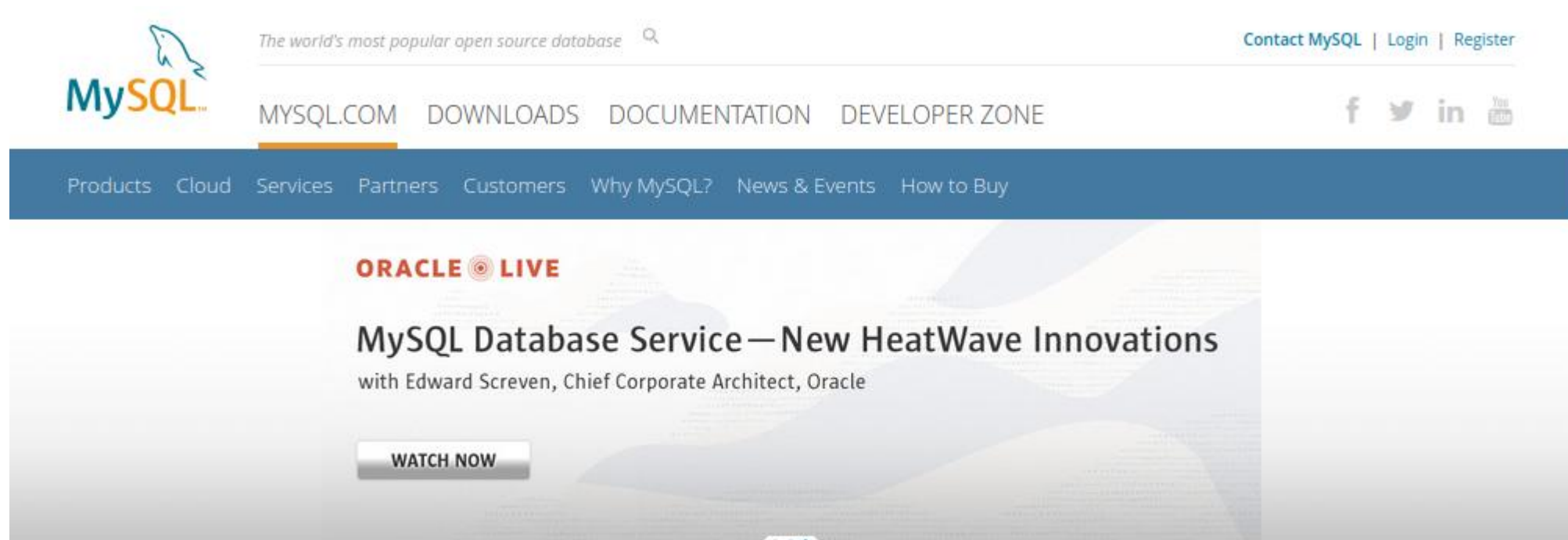
**Zhaoyuan Fang**

zhaoyuanfang@intl.zju.edu.cn

# Course set-up

- **Let's install MySQL!**
- **Go to : https://www.mysql.com**
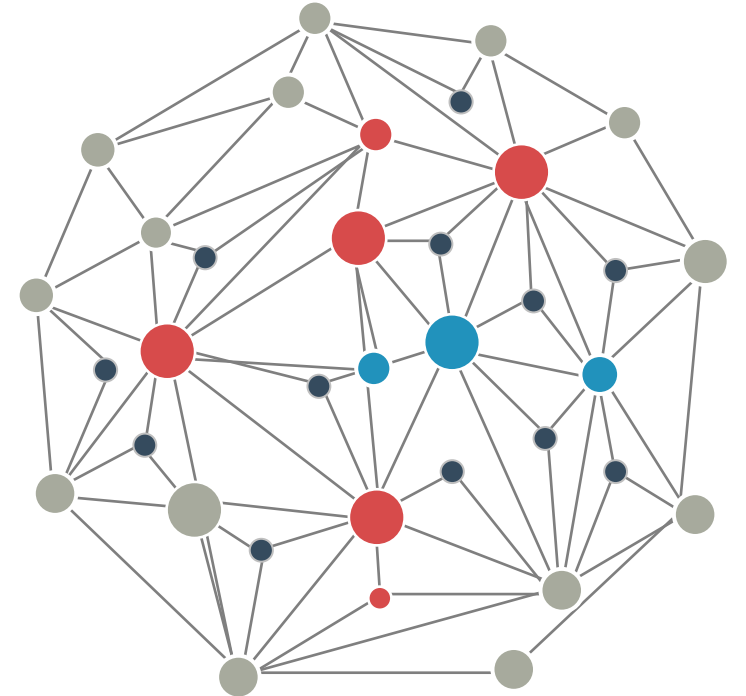
# Course set-up

**Tips:**

- **Follow the installation GUI**

- **Set a simple root password (e.g. 111)**

- **Type "show databases;"**

DST2 – Week 1
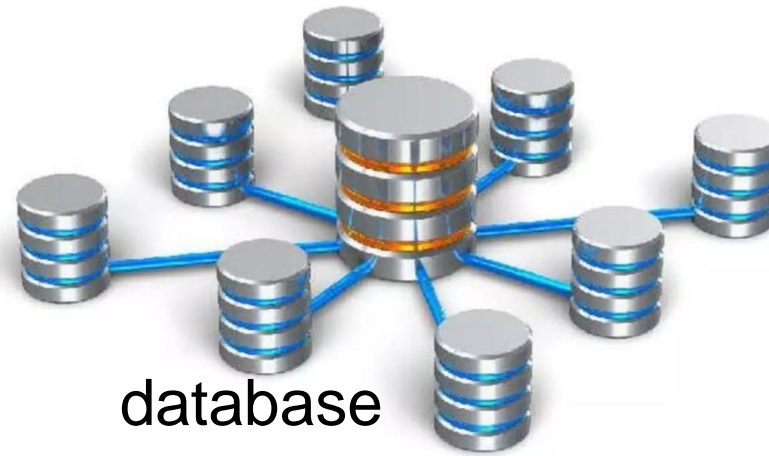
# Introduction to relational database

**Zhaoyuan Fang**

zhaoyuanfang@intl.zju.edu.cn

# General Info



database

Prof. Weng                  Prof. Fang                  Prof. Chen

# Course Info – Database part

- Goal:
  - Apply the concepts of relational data model in database design. Use Structured Query Language (SQL) to operate a relational database.
  - Apply the theory of database structure optimization. Balance non-redundancy and query efficiency by functional dependency analysis and design normalization.
  - Able to manage database transaction, data storage, indexing, and query optimization.

# Course Info – Database part

- Lecture – Database part

  - 6 weeks, every Tuesday 9:00 am – 11:50 am, total break at 15~20 mins

  - We will combine Lecture/tutorial/practical all in this 2hr30mins

- ICA

  - Mini-project 2: We will use Blackboard system, individual work is required.

    - Deadline: 10:00pm, Feb 13th, 2023 (Mon)

    - What to submit: your SQL code as well as a short report

    - Provisional marks and feedback date: Mar 6th, 2023

# Course Info – Database part

- Week 01: Introduction to relational databases

- Week 02: Fundamentals of Structured Query Language (SQL)

- Week 03: Advanced Structured Query Language (SQL)

- Week 04: Database design and E-R modelling

- Week 05: Database normalization, indexing and query optimization

- Week 06: Transaction management and introduction to MySQL

# Week 1 Learning Objectives – Part 1+2+3

- Understand the rationale for using database
- Differentiate data vs information vs knowledge
- See how modern databases evolved from file systems
- Describe the main functions of a database management system (DBMS)
- Structure of Relational Databases
- Fundamental Relational-Algebra-Operations
- Additional Relational-Algebra-Operations
- Extended Relational-Algebra-Operations
- Null Values
- Modification of the Database
- Install course set-up environment (MySQL)

# Week 1 Learning Objectives – Part 1

- **Understand the rationale for using database**

- **Differentiate data vs information vs knowledge**

- **See how modern databases evolved from file systems**

- **Describe the main functions of a database management system (DBMS)**

# Databases are everywhere: applications in daily life

A day in Susan's life
See how many databases she interacts with everyday

**In the morning, Susan checks her wechat, weibo accounts**

- Where is the data about the friends and groups stored?
- Where are the "likes" stored and what would they be used for?

Users
Friends
Posts

**At study, susan uses Baidu to search for useful information**

- Where is those information stored?
- How the information is queried upon a searching?

Keywords
Websites
Links

**After study, she plans for a trip and buys airline tickets and hotel reservations online**

- Where are the airline and hotel data from ?
- What customer data would be kept by the website?
- Where would the customer data be stored?

Flights
Hotels
Customers

**Then she makes a few online purchases**

- Where are the product and transaction data stored?
- Where does the system get the data to generate product "recommendations" to the customer?
- Where would credit card information be stored?

Products
Sales
Customers

# Databases are everywhere: applications in biomedical sciences

## Repositories

NCBI (GenBank)         http://www.ncbi.nlm.nih.gov/
EMBL-EBI  http://www.ebi.ac.uk/
DDBJ                   http://www.ddbj.nig.ac.jp/
INSDC                  http://www.insdc.org/
RefSeq                 http://www.ncbi.nlm.nih.gov/RefSeq/
VEGA                   http://vega.sanger.ac.uk/
CCDS                   http://www.ncbi.nlm.nih.gov/CCDS/

## Genome Browsers

Ensembl                    http://www.ensembl.org/
UCSC Genome Browser        http://genome.ucsc.edu/
NCBI MapViewer             http://www.ncbi.nlm.nih.gov/mapview/
1000 Genomes               http://www.1000genomes.org/

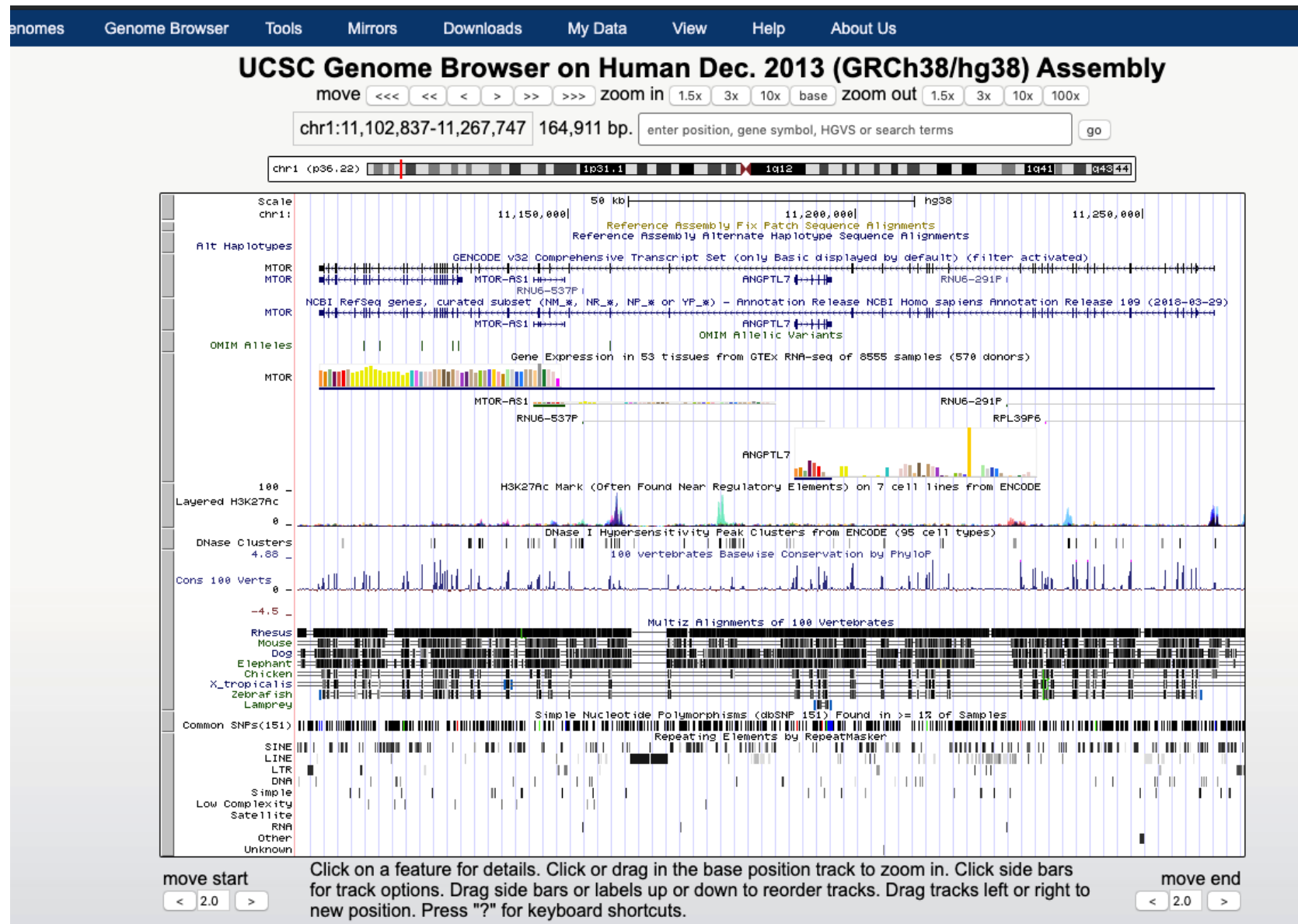## Subject-specific Database

PDB, Protein 3D structures        http://www.rcsb.org/
Pfam, protein families            http://pfam.sanger.ac.uk/
GEO, Gene Expression              http://www.ncbi.nlm.nih.gov/geo/
TCGA, cancer genome atlas         https://portal.gdc.cancer.gov/
ArrayExpress, transcriptomics     http://www.ebi.ac.uk/microarray-as/ae/
dbGaP                             http://www.ncbi.nlm.nih.gov/sites/entrez?db=gap
HuGE Navigator                    http://www.hugenavigator.net/

## Species & Taxa specific databases

Rat            http://rgd.mcw.edu/
Mouse          http://www.informatics.jax.org/
ZFIN, Zebrafish    http://zfin.org FlyBase,
Drosophila     http://flybase.org/ VectorBase,
human disease  http://www.vectorbase.org/
C. elegans     http://www.wormbase.org
crop grasses   http://www.gramene.org
Arabidopsis    http://www.arabidopsis.org/
Yeast          http://www.yeastgenome.org/
Microbial      http://img.jgi.doe.gov/ EcoliHub,
E. coli        http://www.ecolicommunity.org/ VBRC,
Viral          http://athena.bioc.uvic.ca/

## And many more!

# Example - UCSC Genome Browser (a huge database)



**Live demo:**
http://genome.ucsc.edu/

# What is Data? Data vs Information vs Knowledge

DATA

Processing →

INFORMATION

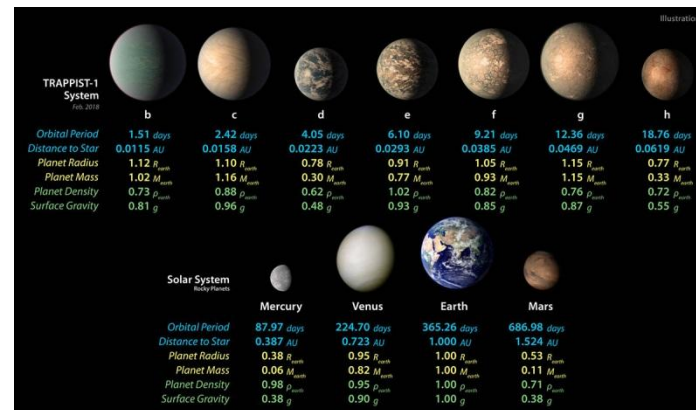Processing →

KNOWLEDGE

Raw facts
(number/text/figure/sequence…)

Reveal its meaning

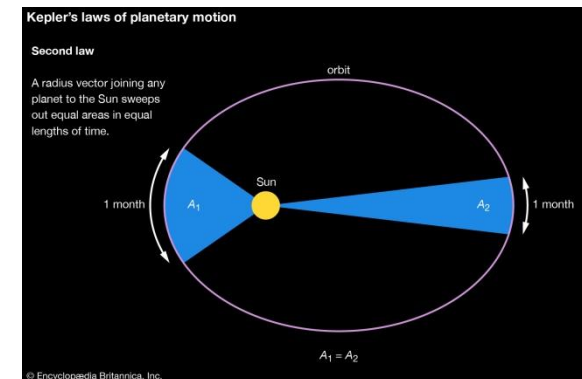information and facts
about a specific subject



Astronomical data of planets



Motion information for planets



Kepler's laws

# What is Data? Data vs Information vs Knowledge

**DATA**

@SN608:4:1101:268.60:93.50#0/1 :GCACTA
NACGCTGAATCAATGTTTCTCCAAAACCATTGATAACTAAATATCATAATA
+SN608:4:1101:268.60:93.50#0/1 :GCACTA
@PP['Q^`aa``ababaaaba^aa]```baaaaaaa]aa`aS^`aaaaaa]

Read name
sequence
Read name
Sequence quality

**Information**



**Knowledge is important!**

**Knowledge**      DNA methylation is decreased over *some genes* in *some cells*

# Data Management: Information explosion

# How is Data Managed? primitive vs. advanced

The basis of Data Management

- Let's say you need to develop a management system for the library, so the facts you need to figure out in the beginning are :
  - Who are the objects? What are the operations?

    Objects: books, students, administrator

    Operations: book borrowing, reservation, book returning etc…

- So your jobs are:

  How to define objects?
  How to store such data?
  How to implement those operations?
  A visible system that is user-friendly …

  We will answer these questions, through this course

keep records by tying knots

# Data Structure/Data Type/Data Management

**Data Structure:**

The **logical and mathematical** model of a particular organization of data and the relationship among its members is called a data structure. It is an organized collection of data which perform a **particular operation**.

**Data Type:**



**Data Management:**

- Data definition: logical + physical
- Data operation:
  - Query: type/characteristics/relationship
  - Update: insert, delete, edit
- Data constraints: consistency of data type, restriction on the data value etc.

# How is Data Managed? History

## Manual Management (before 1950s)
"Farming Age"

## File system (1950-1960)
"Industrial Age"

## Database system (After 1960-)
"Information Age"



A File System

| Personnel dept. | Sales dept. | Accounting dept. |
| --- | --- | --- |
| Employees | Customers · Sales · Inventory | Accounts |

A Database System

Personnel dept.
Sales dept.
Accounting dept.
→ DBMS →
Database
Employees
Customers
Sales
Inventory
Accounts

# Manual Data management

**Computers are not used or not widely used**
- **Small amount of data and simple structure, such as high-order equations, curve fitting, etc.**

**Data stored in sequential access "devices"**
- **Books, papers, tapes, cards; no disk**

**No operating system, no data management software**

**No computer – no 'light'**

Task 1 — access — data1

Task 2 — access — data2

Task n — access — data3

# Manual Data management

Without databases

- Each application manages its own data
- Data is stored multiple times (**redundancy**)
- Problems:
    - Waste of storage space
    - "Forgetting" of changes
    - No centralized, "standardized" data management

# Problems of Data Redundancy

Data Redundancy may cause:

- Other software systems cannot process large amounts of data efficiently
- Many users or applications cannot access the same data in parallel without interfering with each other
- Application developers / users cannot develop / use applications without knowing
  - internal representation of data
  - storage media or computers
  **(no data independence)**
- No data security; potential loss of data

# File System

- **After 1960s, Computer is used for computing and management**
- **Data is stored in disk operating system**

**Characteristics:**
Storage space management, Directory management, File read and write management, file protection, user operation interface

**Problem with file system:**
- Long development time
- Difficulty of getting quick answers
- Complex system administration
- Lack of security and limited data sharing
- Extensive programming
- Data redundancy

Program1

Program2

Program n

Access method

Data1

Data2

Data n

**Islands of information**

# Database motivation

- Database system are center piece of modern IT systems

- … ubiquitous

- Database scientists are in high demands in all area

Harvard Business Review

ANALYTICS

Data Scientist: The Sexiest Job of the 21st Century

**The skills Data Scientists need today**
(based on 300 job listings from tech companies in June 2019)

| Skill | Count |
|---|---|
| Python | 287 |
| R | 245 |
| SQL | 231 |
| Big Data | 221 |
| Apache Spark | 177 |
| Hadoop | 164 |
| Java | 125 |
| Scala | 104 |
| NLP | 83 |
| ETL | 81 |
| Deep Learning | 74 |
| TensorFlow | 38 |
| MATLAB | 30 |
| Pandas | 29 |
| scikit-learn | 27 |
| NumPY | 27 |
| Keras | 17 |
| Computer Vision | 11 |
| PyTorch | 6 |

# Database System background and characteristics

- **Computers become super powerful in managing huge data with complicated relationship**
- **Data need to be shareable (between different Apps, languages etc)**
- **The development of storage techniques**
- **Softwares become more expensive, but hardwares become cheaper**

➡️ Database management systems (DBMS)

Program2

Program1

Program n

Access

Data 1    Data 2

Data n

# School Database Management System

# DBMS is the core of a database system

```
Database
System
```

Application   Application   Application

**DBMS**

**Database Management System =**
Software for Managing Databases

Database

**Structured data,**
Which is managed by DBMS

**Advantages of DBMS:**
- Accessible for everyone/Apps/programming language in the group
- Data structure indicate the relationship between data
- Low data redundancy
- High independence from Apps
- Easy to expand
- Unified data control
  - Security
  - Integrity
  - Concurrency
  - Easier to recover

# How to manage data in the future?

Questions we need to consider:

- How to organize (model and use) data?
- How to store data safely and persistently?
- How to process huge amounts of data (>= terabytes) efficiently?
- How can many users (10,000 or even more) access data concurrently?

# Summary for Part 1

- **Recognize the importance of database**

- **Understand the relationship of data vs information vs knowledge**

- **Evolution of modern database (manual – file system – database)**

- **illustrate the functions of a database management system (DBMS)**

Break time

# Week 1 Learning Objectives – Part 2

- **Structure of Relational Databases**

- **Fundamental Relational-Algebra-Operations**

- **Additional Relational-Algebra-Operations**

- **Extended Relational-Algebra-Operations**

- **Null Values**

- **Modification of the Database**

# Group exercise/task rules

- Throughout this course, we will have some group exercises/tasks for you to complete

- Students sitting on same table naturally forming a group. You can discuss and solve the problems.

- Then show your answer as volunteer (otherwise I will pick one group randomly)

Before we start, give a fancy name for your team (1 mins).

☐ **Hello SQLer!**

# Relational database

- Some early models (immature)
- In 1970, Codd proposed a relational model of database. He won 1981 ACM Turing Award for this.
- A software system used to maintain relational databases is a relational database management system (RDBMS).

- System (early days):
  - System R : developed by IBM
  - INGRES: developed by UCBerkeley
- Current commercial RDBMS:
  - Oracle, SQL Server, DB2, MySQL, PostgreSQL
  - Access (Microsoft)

# What is a Relation? example

Relation:   **attribute 1**     **attribute 2**

| Name | Gender |
|------|--------|
| Smith | Male |
| Lisa | Female |
| Mohamed | Male |
| Ana | Female |
| Noah | Male |

(Lisa, Female)   **tuple**

**Domain D1**              **Domain D2**              **All possible tuples**
**{Smith, Lisa, …, Qian, …}**   **{Male, Female}**   **{ (Smith,Male), …}**

**X**                                    **=**

# Cartesian product

# Basic structure

- Formally, given sets $D_1$, $D_2$, …. $D_n$, a **relation** $r$ is a subset of Cartesian product $D_1$ x $D_2$ x … x $D_n$

$$r \in D_1 \text{ x } D_2 \text{ x } … \text{ x } D_n$$

- Thus, a **relation** is a **set of tuples** $(a_1, a_2, …, a_n)$ where each $a_i \in D_i$

| Name | Gender |
|------|--------|
| Smith | Male |
| Lisa | Female |
| Mohamed | Male |
| Ana | Female |
| Noah | Male |

(Lisa, Female)   **tuple**

**Domain D1**          **Domain D2**

# Basic structure

- Example: If

  - *customer_name* = {Jones, Smith, Curry, Lindsay, …}     /* Set of all customer names */

  - *customer_street* = {Main, North, Park, …}  /* set of all street names*/

  - *customer_city*    = {Harrison, Rye, Pittsfield, …}  /* set of all city names */

  Then *r* = {       (Jones,   Main,  Harrison),
                   (Smith,    North, Rye),
                   (Curry,    North, Rye),
                   (Lindsay, Park,  Pittsfield) }

    is a relation over
      *customer_name  x  customer_street  x  customer_city*

# Attribute and Domain

- Each attribute has its values from a **domain**

- Attribute values are (normally) required to be **atomic**

  - i.e. indivisible (not a set of values)

- The special value *null* is a member of every domain

  - null may cause complications in the definition of many operations

  - We shall ignore the effect of null values in our main presentation and consider their effect later

<div align="center">

attribute        domain

*customer_name*   <=   {Jones, Smith, Curry, Lindsay, …}

</div>

# Relational schema: an abstract form

- $A_1, A_2, \ldots, A_n$ are *attributes*

- $R = (A_1, A_2, \ldots, A_n)$ is a *relational schema*

  Example:

  *Customer_schema = (customer_name, customer_street, customer_city)*

- *r(R)* denotes a *relation r* on the *relation schema R*

  Example:

  *customer (Customer_schema)      -- a real table*

# Relation: an instance of relational schema

- The current values of a relational schema are specified by a table

**Relation Name**

**attributes**
(or columns)

| customer_name | customer_street | customer_city |
|---|---|---|
| Jones | Main | Harrison |
| Smith | North | Rye |
| Curry | North | Rye |
| Lindsay | Park | Pittsfield |

**R** ← Relation schema

**tuples**
(or rows)

*customer*

# Some synonyms

| Relation | = Relational Table | = Table |
|----------|--------------------|---------| 
| Column | = Attribute | = Field |
| Row | = Tuple | = Record |

# Characteristics of Relation

- Each attribute is **homogenous**
  - i.e. values are from **same domain**, and each column have the **same data type**
  - attribute values are **atomic** (a single value, but not a set)
  - *Null* can be in any data type
- Tuples are a set
  - so must be **unique**, but can be in any **order**

Attributes/column homogenous

| T# | S# | C# |
|----|----|----|
| t1 | s1 | c1 |
| t1 | t2 | c2 |

Task 1:        Yes or No?

Different attribute from same domain

| E# | ENAME | MGR |
|----|-------|-----|
| e1 | TOM   | e2  |
| e2 | jerry | null |

Task 2:        Yes or No?

# Relational Model – from relation to database



DEPT(<u>dno</u> , dname , **dean**)

S(<u>sno</u> , sname, sex, age, **dno**)

C(<u>cno</u> , cname , **pcno** , credit)

SC(<u>sno, cno</u> , grade)

PROF(<u>pno</u> , pname, sal, **dno**)

*DEPT: Department*
*S: student*
*C: Course*
*SC :  Student-Course*
*PROF :  Professor*

# Integrity Constraints - Keys

**Candidate Key:**

The **minimal set of attributes** which can **uniquely** identify a tuple is known as candidate key.

*dno, dname in DEPT relations.*

🎯 Task 3     What could be the candidate keys S relations?

🎯 Task 4     What could be the candidate keys C relations?

| Department | | |
|---|---|---|
| <u>dno</u> | dname | dean |

| Student | | | | |
|---|---|---|---|---|
| <u>sno</u> | sname | sex | age | dno |

| Professors | | | |
|---|---|---|---|
| <u>pno</u> | pname | sal | dno |

| Student - Course | | |
|---|---|---|
| <u>sno</u> | <u>cno</u> | grade |

| Course | | | |
|---|---|---|---|
| <u>cno</u> | cname | credit | pcno |

# Integrity Constraints - Keys

**Primary Key**

a candidate key chosen as the principal means of **identifying tuples within a relation** (unique & stable). Should choose an attribute that does not change with time.

e.g. *dno or dname in DEPT relations.*

🎯 Task 5: What could be the primary key S relations?

🎯 Task 6: What could be the primary key C relations?

| Department | | |
|---|---|---|
| <u>dno</u> | dname | dean |

| Student | | | | |
|---|---|---|---|---|
| <u>sno</u> | sname | sex | age | dno |

| Professors | | | |
|---|---|---|---|
| <u>pno</u> | pname | sal | dno |

| Student - Course | | |
|---|---|---|
| <u>sno</u> | <u>cno</u> | grade |

| Course | | | |
|---|---|---|---|
| <u>cno</u> | cname | credit | pcno |

# Integrity Constraints - Keys

**Foreign Key**
A relation schema may have an attribute that **corresponds** to the primary key of **another** relation.  The attribute is called a **foreign key**.

e.g. *dno in S relations.*

🎯 Task 7:    What are the other foreign keys in this database?

# Relational database

## Intension

The intension of a given relation **is independent of time**. It is the permanent part of the relation. It corresponds to what is specified in the relational schema.

Relation: Employee at time= t1

| EmpNo | EmpName | Age | Dept |
|-------|---------|-----|------|
|       |         |     |      |
|       |         |     |      |
|       |         |     |      |
|       |         |     |      |

The teaching building and each room in it is always there, but students inside each room at specific time may varies!

## Extension

The extension of a given relation is the set of tuples appearing in that relation at any given instance. The extension thus **varies with time**. It changes as tuples are created,

Relation: Employee at time= t1

| EmpNo | EmpName | Age | Dept |
|-------|---------|-----|------|
| 1001 | Jason | 23 | SD |
| 1002 | William | 24 | HR |
| 1003 | Jonathan | 28 | Fin |
| 1004 | Harry | 20 | Fin |

Relation: Employee at time= t2 after adding more records

| EmpNo | EmpName | Age | Dept |
|-------|---------|-----|------|
| 1001 | Jason | 23 | SD |
| 1002 | William | 24 | HR |
| 1003 | Jonathan | 28 | Fin |
| 1004 | Harry | 20 | Fin |
| 1005 | Smith | 22 | HR |
| 1006 | Mary | 19 | HR |
| 1007 | Sarah | 23 | SD |

Relation: Employee at time= t2 after adding more records

| EmpNo | EmpName | Age | Dept |
|-------|---------|-----|------|
| 1001 | Jason | 23 | SD |
| 1002 | William | 24 | HR |

# Query Language

- Query languages:
  - Language in which user requests information from the database
    - - Procedural: specify how to get the data wanted
    - - Non-procedural (declarative): only declares which data are wanted
- "Pure" languages:
  - Relational algebra, Tuple relational calculus, Domain relational calculus
  - Pure languages form underlying basis of query languages

- **Query languages:**
  - SQL: Developed by IBM for system R (Relational algebra)
  - QUEL: based on ALPHA by Codd, realized on INGRES
    (Tuple relational calculus)
  - QBE: developed by IBM (Domain relational calculus)

# Relational Algebra

- Six basic operators
  - unary operation
    - select: $\sigma$
    - project: $\Pi$
    - rename: $\rho$
  - Multivariate operation
    - union: $\cup$
    - set difference: $-$
    - Cartesian product: x
- Additional operations
  - Set intersection、Natural join、Division、Assignment
- Extended Relational-Algebra-Operations
  - Generalized projection、Aggregate functions、outer join
- The operators take one or two relations as inputs and produce a new relation.

**LET'S DO SOME HIGH SCHOOL MATH!** ☺

# Select Operation - Example

- Notation:

$$\sigma_p(r) = \{t \mid t \in r \text{ and } p(t)\}$$

  *i.e. "Choosing rows"*

- *p* is called the **selection predicate**, i.e. a formula consisting of: $\wedge$ (**and**), $\vee$ (**or**), $\neg$ (**not**)

- *p* examples:

  brand = "huawei" $\wedge$ price < 5000

  brand = "iphone" $\wedge$ price < 6000

- Relation r

| A | B | C | D |
|---|---|---|---|
| $\alpha$ | $\alpha$ | 1 | 7 |
| $\alpha$ | $\beta$ | 5 | 7 |
| $\beta$ | $\beta$ | 12 | 3 |
| $\beta$ | $\beta$ | 23 | 10 |

Task 8: $\sigma_{A=B \,\wedge\, D > 5}(r)$

# Project Operation - Example

- Notation:

$$\prod_{A1,A2,\ldots Ak} (r)$$

  where are attributes of relation $r$

  i.e. *"Choosing columns"* $A_1, A_2, \ldots A_k$

- duplicated rows removed (keep unique)

- Example:
  $$\prod_{account\_number, balance} (account)$$

- Relation $r$:

| A | B | C |
|---|---|---|
| $\alpha$ | 10 | 1 |
| $\alpha$ | 20 | 1 |
| $\beta$ | 30 | 1 |
| $\beta$ | 40 | 2 |

🎯 Task: $\prod_{A,C} (r)$

# Rename Operation - Example

- Allows us to **rename the attributes** in the results of relational-algebra expressions

- Example:

$$\rho_X(E)$$

returns the expression $E$ under the name $X$

- If a relational-algebra expression $E$ has arity $n$, then

$$\rho_{x(A_1, A_2, \dots, A_n,)}(E)$$

returns the result of expression $E$ under the name $X$, and with the attributes renamed to $A_1, A_2, \dots, A_n$.

# Set-Intersection Operation - Example

- Notation: $r \cap s$

- Defined as:

    $r \cap s = \{\, t \mid t \in r \textbf{ and } t \in s \,\}$

- Assume:

    - $r$, $s$ have the *same* **arity** (number of attributes)

    - attributes of $r$ and $s$ are **compatible**

- Relation $r$, $s$:

| A | B |
|---|---|
| α | 1 |
| α | 2 |
| β | 1 |

$r$

| A | B |
|---|---|
| α | 2 |
| β | 3 |

$s$

Task : $r \cap s$

# Set-Union Operation - Example

- Notation: $r \cup s$

- Defined as:

  $$r \cup s = \{t \mid t \in r \text{ or } t \in s\}$$

- For $r \cup s$ to be valid.

  1. $r, s$ must have the *same* **arity**

  2. The attribute domains must be **compatible**

- e.g. to find all customers with either an account or a loan

  $$\Pi_{customer}\ (depositor)\ \cup\ \Pi_{customer}\ (borrower)$$

- Relations $r, s$:

| A | B |
|---|---|
| $\alpha$ | 1 |
| $\alpha$ | 2 |
| $\beta$ | 1 |

$r$

| A | B |
|---|---|
| $\alpha$ | 2 |
| $\beta$ | 3 |

$s$

Task: $r \cup s$

# Set-Difference Operation - Example

- Notation $r - s$

- Defined as:

  $$r - s = \{t \mid t \in r \text{ and } t \notin s\}$$

- Set differences must be taken between **compatible** relations.

  - $r$ and $s$ must have the same arity

  - attribute domains of $r$ and $s$ must be compatible

- Relations $r$, $s$:

| A | B |
|---|---|
| $\alpha$ | 1 |
| $\alpha$ | 2 |
| $\beta$ | 1 |

$r$

| A | B |
|---|---|
| $\alpha$ | 2 |
| $\beta$ | 3 |

$s$

Task: $r - s$

# Cartesian-Product Operation - Example

- Notation *r* x *s*

- Defined as:

$$r \times s = \{t\ q \mid t \in r \textbf{ and } q \in s\}$$

- Assume that attributes of r(R) and s(S) are disjoint. (That is, $R \cap S = \varnothing$).

- If attributes of *r(R)* and *s(S)* are not disjoint, then need **rename** at first.

- Relations *r, s*:

| A | B |
|---|---|
| α | 1 |
| β | 2 |

*r*

| C | D | E |
|---|---|---|
| α | 10 | a |
| β | 10 | a |
| β | 20 | b |
| γ | 10 | b |

*s*

Task: *r* x *s*

# Composition-of Operation - Example

- Can build expressions using multiple operations

- Example: $\sigma_{A=C}(r \times s)$

- already know $r \times s$:

| A | B | C | D | E |
|---|---|---|---|---|
| $\alpha$ | 1 | $\alpha$ | 10 | a |
| $\alpha$ | 1 | $\beta$ | 10 | a |
| $\alpha$ | 1 | $\beta$ | 20 | b |
| $\alpha$ | 1 | $\gamma$ | 10 | b |
| $\beta$ | 2 | $\alpha$ | 10 | a |
| $\beta$ | 2 | $\beta$ | 10 | a |
| $\beta$ | 2 | $\beta$ | 20 | b |
| $\beta$ | 2 | $\gamma$ | 10 | b |

Task: $\sigma_{A=C}(r \times s)$

# Natural-join Operation - Example

- Notation: r ⋈ s

- r ⋈ s is a relation on schema $R \cup S$ obtained as follows:

  - Consider each pair of tuples $t_r$ from $r$ and $t_s$ from $s$ on **shared attributes**

  - If $t_r$ and $t_s$ have the same value on each of the attributes in $R \cap S$, add a tuple $t$ to the result, where

    - ▸ $t$ has the same value as $t_r$ on $r$

    - ▸ $t$ has the same value as $t_s$ on $s$

- Relations r, s:

| A | B | C | D |
|---|---|---|---|
| $\alpha$ | 1 | $\alpha$ | a |
| $\beta$ | 2 | $\gamma$ | a |
| $\gamma$ | 4 | $\beta$ | b |
| $\alpha$ | 1 | $\gamma$ | a |
| $\delta$ | 2 | $\beta$ | b |

r

| B | D | E |
|---|---|---|
| 1 | a | $\alpha$ |
| 3 | a | $\beta$ |
| 1 | a | $\gamma$ |
| 2 | b | $\delta$ |
| 3 | b | $\in$ |

s

🎯 Task: $r \bowtie s$

# Division Operation - Example

- Notation:    $r \div s$

- Suited to queries that include the phrase "for all".

- Let $r$ and $s$ be relations on schemas $R$ and $S$ respectively where

  - $R = (A_1, \ldots, A_m, B_1, \ldots, B_n)$

  - $S = (B_1, \ldots, B_n)$

  The result of $r \div s$ is a relation on schema

  $R - S = (A_1, \ldots, A_m)$

  $r \div s = \{\, t \mid t \in \prod_{R\text{-}S}(r) \wedge \forall\, u \in s\, (\, tu \in r\,)\, \}$

  Where $tu$ means the concatenation of tuples $t$ and $u$ to produce a single tuple

- Relations $r, s$:

| A | B |
|---|---|
| $\alpha$ | 1 |
| $\alpha$ | 2 |
| $\alpha$ | 3 |
| $\beta$ | 1 |
| $\gamma$ | 1 |
| $\delta$ | 1 |
| $\delta$ | 3 |
| $\delta$ | 4 |
| $\in$ | 6 |
| $\in$ | 1 |
| $\beta$ | 2 |

$r$

| B |
|---|
| 1 |
| 2 |

$s$

Task: $r \div s$

# Generalized Projection

- Extends the projection operation by **allowing arithmetic functions** to be used in the projection list.

$$\Pi_{F_1, F_2, \ldots, F_n}(E)$$

- $E$ is any relational-algebra expression

- Each of $F_1$, $F_2$, …, $F_n$ are are arithmetic expressions involving constants and attributes in the schema of $E$.

- e.g.

 Given relation *credit_info(customer_name, limit, credit_balance),* find how much more each person can spend:

$$\Pi_{customer\_name,\ limit\ -\ credit\_balance}(credit\_info)$$

# Aggregate Functions and Operations

■ **Aggregation function** takes a collection of values and returns a single value as a result.

> **avg**: average value
> **min**: minimum value
> **max**: maximum value
> **sum**: sum of values
> **count**: number of values

■ **Aggregate operation** in relational algebra

$$_{G_1,G_2,\dots,G_n} \mathcal{G}_{F_1(A_1),F_2(A_2),\dots,F_n(A_n)}(E)$$

*E* is any relational-algebra expression

- $G_1, G_2 \dots, G_n$ is a list of attributes on which to group (can be empty)
- Each $F_i$ is an aggregate function
- Each $A_i$ is an attribute name

■ Relation *r*:

| A | B | C |
|---|---|---|
| $\alpha$ | $\alpha$ | 7 |
| $\alpha$ | $\beta$ | 7 |
| $\beta$ | $\beta$ | 3 |
| $\beta$ | $\beta$ | 10 |

🎯 Task: $\mathcal{g}_{\textbf{sum(c)}}(\text{r})$

# Aggregate Functions and Operations

- Relation *account* **grouped by** *branch-name*:

| branch_name | account_number | balance |
|---|---|---|
| Perryridge | A-102 | 400 |
| Perryridge | A-201 | 900 |
| Brighton | A-217 | 750 |
| Brighton | A-215 | 750 |
| Redwood | A-222 | 700 |

Task: $_{branch\_name}\, g\, _{\textbf{sum}(balance)}\, (account)$

# Outer Join

- An extension of the join operation that avoids loss of information.

- Computes the join and then **adds tuples form one relation that does not match tuples in the other relation** to the result of the join.

- Uses *null* values:

  - *null* signifies that the value is unknown or does not exist

  - All comparisons involving *null* are (roughly speaking) **false** by definition.

    - We shall study precise meaning of comparisons with nulls later

# Outer Join

- Relation *loan*

| loan_number | branch_name | amount |
|---|---|---|
| L-170 | Downtown | 3000 |
| L-230 | Redwood | 4000 |
| L-260 | Perryridge | 1700 |

- Relation *borrower*

| customer_name | loan_number |
|---|---|
| Jones | L-170 |
| Smith | L-230 |
| Hayes | L-155 |

- Join

  *loan* ⋈ *borrower*

| loan_number | branch_name | amount | customer_name |
|---|---|---|---|
| L-170 | Downtown | 3000 | Jones |
| L-230 | Redwood | 4000 | Smith |

- Left Outer Join

  *loan* ⟕ *borrower*

| loan_number | branch_name | amount | customer_name |
|---|---|---|---|
| L-170 | Downtown | 3000 | Jones |
| L-230 | Redwood | 4000 | Smith |
| L-260 | Perryridge | 1700 | *null* |

# Null Values

- It is possible for tuples to have a null value, denoted by *null*, for some of their attributes

- *null* signifies an unknown value or that a value does not exist.

- *null* should best be avoided theoretically, though in practice they are used. *Null* could cause potential problems:

  - DBMSs (e.g. MySQL) may treat *null* in some way, e.g. arithmetic expression involving *null* as *null*, aggregate functions ignoring null values.

  - However, no such guarantee for all DBMSs

# Modification of the Database

- The content of the database may be modified using the following operations:
    - Deletion
    - Insertion
    - Updating
- All these operations are expressed using the assignment operator.

# Deletion

- A delete request is expressed similarly to a query, except instead of displaying tuples to the user, the selected tuples are removed from the database.

- Can delete only whole tuples; cannot delete values on only particular attributes

- A deletion is expressed in relational algebra by:

$$r \leftarrow r - E$$

where $r$ is a relation and $E$ is a relational algebra query.

# Insertion

- To insert data into a relation, we either:
    - specify a tuple to be inserted
    - write a query whose result is a set of tuples to be inserted
- in relational algebra, an insertion is expressed by:

$$r \leftarrow r \cup E$$

where $r$ is a relation and $E$ is a relational algebra expression.

- The insertion of a single tuple is expressed by letting $E$ be a constant relation containing one tuple.

# Updating

- A mechanism to **change a value in a tuple** without charging *all* values in the tuple

- Use the generalized projection operator to do this task

$$r \leftarrow \prod_{F_1, F_2, \ldots, F_l,} (r)$$

- Each $F_i$ is either

  - the $l^{th}$ attribute of $r$, if the $l^{th}$ attribute is not updated, or,

  - if the attribute is to be updated $F_i$ is an expression, involving only constants and the attributes of $r$, which gives the new value for the attribute

# Summary for Part 2

- **Structure of Relational Databases**
  - Relations, domain, attributes, Schema, tuples, Keys, intension/extension
- **Fundamental Relational-Algebra-Operations**
  - Select, project, union, set difference, cartesian product, rename
- **Additional Relational-Algebra-Operations**
  - Set intersection, natural join, division, assignment
- **Extended Relational-Algebra-Operations**
  - Generalized projection, aggregate functions, outer join
- **Null Values**
- **Modification of the Database**
  - Deletion, insertion, updating