Introduction
○○○○○○○○○
○○○○○○○○○○

Relational Model
○○○○○○○○○○○

Query Languages
○○○○○○○○○○○○○○○
○○○○○○○○○○○○○○○○○○○

References
○○

# Databases

## Basics of Relational Database Management

Pierre Senellart

ENS
ÉCOLE NORMALE
SUPÉRIEURE

PSL★

7 February 2024

**Introduction**
○●○○○○○○○○○
Relational Model
○○○○○○○○○○○○
Query Languages
○○○○○○○○○○○○○○○
○○○○○○○○○○○○○○○
References
○○

# Plan

## Data management

Numerous applications (standalone software, Web sites, etc.)
need to manage data:

- Structure data useful to the application
- Store them in a persistent manner (data retained even
  when the application is not running)
- Efficiently query information within large data volumes
- Update data without violating some structural constraints
- Enable data access and updates by multiple users, possibly
  concurrently

Often, desirable to access the same data from several distinct
applications, from distinct computers.

# Example: Information system of a hotel

Access from an in-house software (front desk), a Website (guests), an accounting software suite. Requirements:

- Structured data representing rooms, customers, guests, bookings, rates, etc.

- No loss of data when these applications are unused, or when a power cut arises

- Find quasi-instantaneously which rooms are booked, by whom, a given day, within a history containing several years of bookings

- Easily add a booking by ensuring the same room is not booked twice the same day

- The guest, the front desk employee, the accountant, must not have the same view of the data (confidentiality, ease of use, etc.)

- If a room is available, it cannot be booked by different guests at the same time

# Naive implementation (1/2)

- Implementation in a classical programming language (C++, Java, Python, etc.) of data structures that can represent all useful data

- Definition of ad-hoc file formats to store data on disk, with regular synchronization and a mechanism for failure recovery

- In-memory storage of application data, with data structures (search trees, hash tables) and algorithms (search, sort, aggregation, graph navigation, etc.) for efficiently finding information

- Data update functions, with code ensuring no constraint is violated

# Naive implementation (2/2)

- Definition within the application of user access rights and an authentication process; use of parallel programming to answer different requests at the same time, locks/semaphores on critical data manipulation steps

- Definition and implementation of a network communication protocol to access this software component from the Web, from a Windows application, from an accounting suite, etc.

# Naive implementation (2/2)

- Definition within the application of user access rights and an authentication process; use of parallel programming to answer different requests at the same time, locks/semaphores on critical data manipulation steps

- Definition and implementation of a network communication protocol to access this software component from the Web, from a Windows application, from an accounting suite, etc.

Lots of work!

# Naive implementation (2/2)

- Definition within the application of user access rights and an authentication process; use of parallel programming to answer different requests at the same time, locks/semaphores on critical data manipulation steps
- Definition and implementation of a network communication protocol to access this software component from the Web, from a Windows application, from an accounting suite, etc.

Lots of work! Requires a programmer that masters OOP, serialization, failure recovery, data structures, algorithms, integrity constraint verification, role management, parallel programming, concurrency control, network programming, etc.

# Naive implementation (2/2)

- Definition within the application of user access rights and an authentication process; use of parallel programming to answer different requests at the same time, locks/semaphores on critical data manipulation steps

- Definition and implementation of a network communication protocol to access this software component from the Web, from a Windows application, from an accounting suite, etc.

Lots of work! Requires a programmer that masters OOP, serialization, failure recovery, data structures, algorithms, integrity constraint verification, role management, parallel programming, concurrency control, network programming, etc. Needs to be done again for every new application that manages data!

# Role of a DBMS

### Database Management System

Software that simplifies the design of applications that handle data, by providing a unified access to the functionalities required for data management, whatever the application.

### Database

Collection of data (specific to a given application) managed by a DBMS

# Features of DBMSs (1/2)

Physical independence. The user of a DBMS does not need to know how data are stored (in a file, on a raw partition, in a distributed filesystem, etc.); storage can be modified without impacting data access

Logical independence. It is possible to provide the user with a partial view of the data, corresponding to what he needs and is allowed to access

Ease of data access. Use of a declarative language describing queries and updates on the data, specifying the intent of a user rather than the way this will be implemented

Query optimization. Queries are automatically optimized to be implemented as efficiently as possible on the database

# Features of DBMSs (2/2)

Logical integrity.   The DBMS imposes constraints on data structure; every modification violating these constraints is denied

Physical integrity.   The database remains in a coherent state, and data are durably preserved, even in case of software or hardware failure

Data sharing.   Data are accessible by multiple users, concurrently, and these multiple and concurrent accesses cannot violate logical or physical data integrity

Standardization.   The use of a DBMS is standardized, so that it may be possible to replace a DBMS with another without changing (in a major way) the code of the application

# Plan

## Introduction

Data Management

### Types of DBMSs

## Relational Model

## Query Languages

## References

**Introduction**
○○○○○○○○○
○●○○○○○○○○○

Relational Model
○○○○○○○○○○○

Query Languages
○○○○○○○○○○○○○○○
○○○○○○○○○○○○○○○○○○

References
○○

# Diversity of DBMSs

- <span style="color:red">Dozens</span> of existing DBMSs that are broadly used
- All DBMSs do not provide <span style="color:red">all these features</span>
- DBMSs can be <span style="color:red">differentiated</span> based on:
    - data model used
    - trade-offs made between performance and features
    - ease of use
    - scalability
    - internal architecture

# Major types of DBMSs

Relational (RDBMS). Tables, complex queries (SQL), rich features

XML. Trees, complex queries (XQuery), features similar to RDBMS

Graph/Triples. Graph data, complex queries expressing graph navigation

Objects. Complex data model, inspired by OOP

Documents. Complex data, organized in documents, relatively simple queries and features

Key–Value. Very basic data model, focus on performance

Column Stores. Data model in between key–value and RDBMS; focus on iteration and aggregation on columns

# Major types of DBMSs

Relational (RDBMS). Tables, complex queries (SQL), rich features

XML. Trees, complex queries (XQuery), features similar to RDBMS

Graph/Triples. Graph data, complex queries expressing graph navigation

Objects. Complex data model, inspired by OOP

Documents. Complex data, organized in documents, relatively simple queries and features

Key–Value. Very basic data model, focus on performance

Column Stores. Data model in between key–value and RDBMS; focus on iteration and aggregation on columns

NoSQL

# Classical relational DBMSs

- Based on the relational model: decomposition of data into relations (i.e., tables)
- A standard query language: SQL
- Data stored on disk
- Relations (tables) stored row by row
- Centralized system, with limited distribution possibilities

Introduction
○○○○○○○○○
○○○○●○○○○○○

Relational Model
○○○○○○○○○○○

Query Languages
○○○○○○○○○○○○○○
○○○○○○○○○○○○○○○

References
○○

# Strengths of classical relational DBMSs

- Independence between:
  - data model and storage structures
  - declarative queries and the way queries are executed

- Complex queries

- Fine optimization of queries, indexes allowing quick access to data

- Mature technology, stable, efficient, rich in features and interfaces

- Integrity constraints ensuring invariants on data

- Efficient management of very large volume of data (up to terabytes)

- Transactions (sequences of elementary operations) with guaranties on concurrency control, isolation between users, failure recovery

Introduction
○○○○○○○○○
○○○○○●○○○○

Relational Model
○○○○○○○○○○○

Query Languages
○○○○○○○○○○○○○○
○○○○○○○○○○○○○○○○○○○

References
○○

# ACID properties

Classical relational DBMS transactions satisfy ACID properties:

# ACID properties

Classical relational DBMS transactions satisfy ACID properties:

Atomicity: The set of operations within a transaction is either executed as a whole or canceled as a whole

Consistency: Transactions ensure integrity constraints on the base are respected

Isolation: Two concurrent executions of transactions result in a state equivalent to serial execution of the transactions

Durability: Once transactions are committed, corresponding data stay durably in the base, even in case of system failure

Introduction
ooooooooo
ooooooo●oooo

Relational Model
oooooooooooo

Query Languages
oooooooooooooo
ooooooooooooooooooo

References
oo

# Weaknesses of classical RDBMSs

- Incapable of managing extremely large data volume (of the order of a petabyte)

- Impossible to manage extreme query rates (beyond thousands of queries per second)

- The relational data model is sometimes poorly adapted to the storage and querying of some data types (hierarchical data, unstructured data, semi-structured data)

- ACID properties imply major costs in latency, disk accesses, processing time (locks, logging, etc.)

- Performances limited by disk accesses

**Introduction**
○○○○○○○○○
○○○○○○○●○○

Relational Model
○○○○○○○○○○○

Query Languages
○○○○○○○○○○○○○
○○○○○○○○○○○○○○○○○

References
○○

# NoSQL

- No SQL or Not Only SQL

- DBMSs with other trade-offs than those made by classical systems

- Very diverse ecosystem

- Desiderata: different data model, scalability, extreme performance

- Abandoned features: ACID, (sometimes) complex queries

# NewSQL

- Some applications require:
  - A rich (SQL) query language
  - conformity to ACID properties
  - but greater performance than classical RDBMSs

# NewSQL

- Some applications require:
  - A rich (SQL) query language
  - conformity to ACID properties
  - but greater performance than classical RDBMSs
- Possible solutions:
  - Getting rid of classical bottleneck of RDBMSs: locks, logging, cache management
  - In-memory databases, with asynchronous copy on disk
  - Lock-free concurrency control (MVCC)
  - Shared-nothing distributed architecture with transparent load balancing

# In this course

- We mostly cover classical relational DBMSs (the most used, the most mature, by far!)
- ... and the rich underlying theory
- Many aspects also relevant to other forms of DBMSs

# Plan

# Relational schema

We fix countably infinite sets:

- $\mathcal{L}$ of labels
- $\mathcal{V}$ of values
- $\mathcal{T}$ of types, s.t., $\forall \tau \in \mathcal{T}, \tau \subseteq \mathcal{V}$

## Definition

A relation schema (of arity $n$) is an $n$-tuple $(A_1, \ldots, A_n)$ where each $A_i$ (called an attribute) is a pair $(L_i, \tau_i)$ with $L_i \in \mathcal{L}$, $\tau_i \in \mathcal{T}$ and such that all $L_i$ are distinct

## Definition

A database schema is defined by a finite set of labels $L \subseteq \mathcal{L}$ (relation names), each label of $L$ being mapped to a relation schema.

# Example database schema

- Universe:
    - $\mathcal{L}$ the set of alphanumeric character strings starting with a letter
    - $\mathcal{V}$ the set of finite sequences of bits
    - $\mathcal{T}$ is formed of types such as INTEGER (representation as a sequence of bits of integers between $-2^{31}$ and $2^{31} - 1$), REAL (representation of floating-point numbers following IEEE 754), TEXT (UTF-8 representation of character strings), DATE (ISO 8601 representation of dates), etc.

- Database schema formed of 2 relation names, Guest and Reservation

- Guest: $((\text{id}, \text{INTEGER}), (\text{name}, \text{TEXT}), (\text{email}, \text{TEXT}))$

- Reservation:
  $((\text{id}, \text{INTEGER}), (\text{guest}, \text{INTEGER}), (\text{room}, \text{INTEGER}),$
  $(\text{arrival}, \text{DATE}), (\text{nights}, \text{INTEGER}))$

# Database

### Definition

An instance of a relation schema $((L_1, \tau_1), \ldots, (L_n, \tau_n))$ (also called a relation on this schema) is a finite set $\{t_1, \ldots, t_k\}$ of tuples of the form $t_j = (v_{j1}, \ldots, v_{jn})$ with $\forall j \forall i \, v_{ji} \in \tau_i$.

### Definition

An instance of a database schema (or, simply, a database on this schema) is a function that maps each relation name to an instance of the corresponding relation schema.

Note: Relation is used somewhat ambiguously to talk about a relation schema or an instance of a relation schema.

Introduction
0000000
000000000

**Relational Model**
0000●000000

Query Languages
0000000000000
000000000000000

References
00

# Example

### Guest

| id | name | email |
|----|------|-------|
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

### Reservation

| id | guest | room | arrival | nights |
|----|-------|------|---------|--------|
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Introduction          **Relational Model**          Query Languages          References
ooooooooo              ooooo●oooooo             oooooooooooooooo          oo
ooooooooooo             oooooo                   oooooooooooooooooo

## Some notation

- If $A = (L, \tau)$ is the $i$th attribute of a relation $R$, and $t$ an $n$-tuple of an instance of $R$, we note $t[A]$ (or $t[L]$) the value of the $i$th component of $t$.

- Similarly, if $\mathcal{A}$ is a $k$-tuple of attributes among the $n$ attributes of $R$, $t[\mathcal{A}]$ is the $k$-tuple formed from $t$ by concatenating the $t[A]$ for $A$ in $\mathcal{A}$.

- A tuple is an $n$-tuple for some $n$.

Introduction
0000000
0000000000

**Relational Model**
00000●00000

Query Languages
0000000000000
0000000000000000

References
00

# Simple integrity constraints

One can add to the relational schema some integrity constraints, of different nature, to define a notion of validity of an instance

- Key. A tuple of attribute $\mathcal{A}$ of a relation schema $R$ is a key if there cannot exist two distinct tuples $t_1$ and $t_2$ in an instance of $R$ such that $t_1[\mathcal{A}] = t_2[\mathcal{A}]$

- Foreign key. A $k$-tuple of attributes $\mathcal{A}$ of a relation schema $R$ is a foreign key referencing a $k$-tuple of attributes $\mathcal{B}$ of a relation schema $S$ if for all instances $I^R$ and $I^S$ of $R$ and $S$, for every tuple $t$ of $I^R$, there exists a unique tuple $t'$ of $I^S$ with $t[\mathcal{A}] = t'[\mathcal{B}]$

- Check constraint. Arbitrary condition on the values of the attributes of a relation (applying to each tuple of the instances of that relation)

# Examples of constraints

- id is a <span style="color:red">key</span> of Guest
- email is a <span style="color:red">key</span> of Guest
- id is a <span style="color:red">key</span> of Reservation
- (room, arrival) is a <span style="color:red">key</span> of Reservation
- (guest, arrival) is a <span style="color:red">key</span> of Reservation (?)
- guest is a <span style="color:red">foreign key</span> of Reservation referencing id of Guest
- In Guest, email <span style="color:red">must</span> contain a "@"
- In Reservation, room <span style="color:red">must</span> be between 1 and 650
- In Reservation, nights <span style="color:red">must</span> be positive

Impossible to express more complex constraints (e.g., a room cannot be occupied twice the same night, which depends on the date and the number of nights for multiple tuples of Reservation)

## Variants: named and unnamed perspectives

The version presented considers the attributes of a relation are
ordered and have a name. This is what best matches the way
RDBMSs work, but not necessarily the most pleasant to reason
on the relational model.

Named perspective. We forget the position of attributes, and
consider they are uniquely identified by their
names.

Unnamed perspective. We forget the name of attributes, and
consider they are uniquely identified by their
position. One uses notation such as $t[2]$ to access
the value of the second attribute of a tuple.

No major impact, one will use one or the other depending on
what is convenient.

# Variant: bag semantics

- A relation instance is defined as a (finite) set of tuples. One can also consider a bag semantics of the relational model, where a relation instance is a multiset of tuples.

- This is what best matches how RDBMSs work. . .

- . . . but most of relational database theory has been established for the set semantics, more convenient to work with

- We will mostly discuss the set semantics in this lecture, but explain where differences matter

# Variant: untyped version

- In implementations, attributes are always typed
- In models and theoretical results, one often abstracts attribute types away and considers each attribute has a universal type $\mathcal{V}$
- We will most often omit attribute types

# Plan

# The relational algebra

- Algebraic language to express queries
- A relational algebra expression produces a new relation from the database relations
- Each operator takes 0, 1, or 2 subexpressions
- Main operators:

| Op. | Arity | Description | Condition |
|------|-------|--------------|-----------------|
| $R$ | 0 | Relation name | $R \in \mathcal{L}$ |
| $\rho_{A \to B}$ | 1 | Renaming | $A, B \in \mathcal{L}$ |
| $\Pi_{A_1 \ldots A_n}$ | 1 | Projection | $A_1 \ldots A_n \in \mathcal{L}$ |
| $\sigma_\varphi$ | 1 | Selection | $\varphi$ formula |
| $\times$ | 2 | Cross product | |
| $\cup$ | 2 | Union | |
| $\setminus$ | 2 | Difference | |
| $\bowtie_\varphi$ | 2 | Join | $\varphi$ formula |

# Relation name

| Guest | | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| Reservation | | | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression:  `Guest`

Result:

| id | name | email |
|---|---|---|
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

# Renaming

| Guest | | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| Reservation | | | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression: $\rho_{\texttt{id}\rightarrow\texttt{guest}}(\texttt{Guest})$

Result:

| guest | name | email |
|---|---|---|
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

# Projection

|  | Guest | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

|  | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression:  $\Pi_{\text{email,id}}(\text{Guest})$

Result:

| email | id |
|---|---|
| john.smith@gmail.com | 1 |
| alice@black.name | 2 |
| john.smith@ens.fr | 3 |

## Selection

| | Guest | |
| --- | --- | --- |
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | | Reservation | | |
| --- | --- | --- | --- | --- |
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression:  $\sigma_{\texttt{arrival}>\texttt{2017-01-12}\wedge\texttt{guest}=2}(\texttt{Reservation})$

Result:

| id | guest | room | arrival | nights |
| --- | --- | --- | --- | --- |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

The formula used in the selection can be any Boolean combination of comparisons of attributes to attributes or constants.

# Cross product

| | | Guest | |
|---|---|---|---|
| id | name | email | |
| 1 | John Smith | john.smith@gmail.com | |
| 2 | Alice Black | alice@black.name | |
| 3 | John Smith | john.smith@ens.fr | |

| | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression: $\Pi_{\texttt{id}}(\texttt{Guest}) \times \Pi_{\texttt{name}}(\texttt{Guest})$

Result:

| id | name |
|---|---|
| 1 | Alice Black |
| 2 | Alice Black |
| 3 | Alice Black |
| 1 | John Smith |
| 2 | John Smith |
| 3 | John Smith |

# Union

| Guest | | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| Reservation | | | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression:    $\Pi_{\text{room}}(\sigma_{\text{guest}=2}(\text{Reservation})) \cup$
$\Pi_{\text{room}}(\sigma_{\text{arrival}=\text{2017-01-15}}(\text{Reservation}))$

Result:

| room |
|---|
| 107 |
| 302 |
| 504 |

# Union

| | Guest | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression: $\Pi_{\text{room}}(\sigma_{\text{guest}=2}(\text{Reservation})) \cup$
$\Pi_{\text{room}}(\sigma_{\text{arrival}=2017\text{-}01\text{-}15}(\text{Reservation}))$

Result:

| room |
|---|
| 107 |
| 302 |
| 504 |

This simple union could have been written
$\Pi_{\text{room}}(\sigma_{\text{guest}=2 \vee \text{arrival}=2017\text{-}01\text{-}15}(\text{Reservation}))$. Not always possible.

# Difference

| Guest | | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| Reservation | | | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression:    $\Pi_{\text{room}}(\sigma_{\text{guest}=2}(\text{Reservation})) \setminus$
                $\Pi_{\text{room}}(\sigma_{\text{arrival}=\text{2017-01-15}}(\text{Reservation}))$

Result:

| room |
|---|
| 107 |

## Difference

| | Guest | |
| --- | --- | --- |
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | Reservation | | | |
| --- | --- | --- | --- | --- |
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression: $\Pi_{\text{room}}(\sigma_{\text{guest}=2}(\text{Reservation})) \setminus$
$\Pi_{\text{room}}(\sigma_{\text{arrival}=2017\text{-}01\text{-}15}(\text{Reservation}))$

Result:

| room |
| --- |
| 107 |

This simple difference could have been written
$\Pi_{\text{room}}(\sigma_{\text{guest}=2 \wedge \text{arrival}\neq 2017\text{-}01\text{-}15}(\text{Reservation}))$. Not always
possible.

## Join

Guest

| id | name | email |
|----|------|-------|
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

Reservation

| id | guest | room | arrival | nights |
|----|-------|------|---------|--------|
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

Expression:    `Reservation` $\bowtie_{\text{guest}=\text{id}}$ `Guest`

Result:

| id | guest | room | arrival | nights | name | email |
|----|-------|------|---------|--------|------|-------|
| 1 | 1 | 504 | 2017-01-01 | 5 | John Smith | john.smith@gmail.com |
| 2 | 2 | 107 | 2017-01-10 | 3 | Alice Black | alice@black.name |
| 3 | 3 | 302 | 2017-01-15 | 6 | John Smith | john.smith@ens.fr |
| 4 | 2 | 504 | 2017-01-15 | 2 | Alice Black | alice@black.name |
| 5 | 2 | 107 | 2017-01-30 | 1 | Alice Black | alice@black.name |

The formula used in the join can be any Boolean combination of comparisons of attributes of the table on the left to attributes of the table on the right.

# Note on the join

- The join is not an <span style="color:red">elementary</span> operator of the relational algebra (but it is very useful)

- It can be seen as a <span style="color:red">combination</span> of renaming, cross product, selection, projection

- Thus:

$$\texttt{Reservation} \bowtie_{\texttt{guest=id}} \texttt{Guest}$$

$$\equiv \frac{\Pi_{\texttt{id,guest,room,arrival,nights,name,email}}(}{\sigma_{\texttt{guest=temp}}(\texttt{Reservation} \times \rho_{\texttt{id}\rightarrow\texttt{temp}}(\texttt{Guest})))}$$

- If $R$ and $S$ have for attributes $\mathcal{A}$ and $\mathcal{B}$, we note $R \bowtie S$ the <span style="color:red">natural join</span> of $R$ and $S$, where the join formula is $\bigwedge_{A \in \mathcal{A} \cap \mathcal{B}} A = A$.

# Illegal operations

- All expressions of the relational algebra are not <span style="color:red">valid</span>

- The validity of an expression generally depends on the database schema

- For example:
    - No reference to the name of a relation that doesn't exist in the database schema
    - One cannot reference (within a renaming, projection, selection, join) an attribute that does not exist in the result of a sub-expression
    - One cannot union two relations with different attributes
    - One cannot produce (cross product, join, renaming) a table with two attributes with the same name

- Systems implementing the relational algebra may do a static or dynamic verification of these rules, or sometimes ignore them

# Bag semantics

In bag semantics (what is actually used by RDBMS):

- All operations return multisets
- In particular, projection and union can introduce multisets even when initial relations are sets

## Extension: Aggregation

- Various extensions have been proposed to the relational algebra to add additional features
- In particular, aggregation and grouping [Klug, 1982, Libkin, 2003] of results
- With a syntax inspired from [Libkin, 2003]:

$$\sigma_{\mathrm{avg}>3}(\gamma_{\mathrm{room}}^{\mathrm{avg}}[\lambda x.\mathrm{avg}(x)](\Pi_{\mathrm{room,nights}}(\mathrm{Reservation})))$$

computes the average number of nights per reservation for each room having an average greater than 3

| room | avg |
|------|-----|
| 302  | 6   |
| 504  | 3.5 |

# Plan

Introduction

Relational Model

Query Languages
   Relational Algebra
   SQL

References

## SQL

- Structured Query Language, standard language (ISO/IEC 9075, several versions [ISO, 1987, 1999]) to interact with an RDBMS

- Unfortunately, implementation of the standard very variable from one RDBMS to the next

- Many little things (e.g., available types) vary between RDBMSs instead of following the standard

- Differences more syntactical than major

- Where it makes a difference, we use the PostgreSQL version

- Two main parts: DDL (Data Definition Language) to define the schema and DML (Data Manipulation Language) to query and update data

- Declarative language: express what you mean, the system will take care of translating this into an efficient execution plan

# Syntax of SQL

- Quite verbose, designed to be almost readable as English words [Chamberlin and Boyce, 1974]

- Keywords are case-insensitive, traditionally written in all uppercase

- Identifiers often case-insensitive (depends of the RDBMS), often written with an initial uppercase for table names, in all lower case for attribute names

- Comments introduced with --

- SQL statements end with a ";" in some contexts (e.g., command line client) but the ";" id not properly part of the statement

# NULL

- In SQL, NULL is a special value that an attribute can take within a tuple
- Denotes absence of value
- Different from 0, from an empty string, etc.
- Weird tri-valued logic: True, False, NULL
- A normal comparison (equality, inequality, etc.) with NULL always returns NULL
- **IS NULL**, **IS NOT NULL** can be used to test whether a value is NULL
- NULL est ultimately converted to False
- Weird consequences, poor integration with the formal relational model

# Data Definition Language

**CREATE TABLE** Guest(id INTEGER, name TEXT, email TEXT);
**CREATE TABLE** Reservation(id INTEGER, guest INTEGER,
room INTEGER, arrival DATE, nights INTEGER);

But also:

- **DROP TABLE** Guest; to destroy a table
- **ALTER TABLE** Guest **RENAME TO** Guest2; to rename
  a table
- **ALTER TABLE** Guest **ALTER COLUMN** id **TYPE** TEXT;
  to change the type of a column

## Constraints

Specified at the creation of a table, or added later on (with
**ALTER TABLE**)

**PRIMARY KEY** for the primary key; only one per table, it is
the key that will be used for physical organization
of data; implies **NOT NULL**

    **UNIQUE** for other keys

**REFERENCES** for foreign keys

     **CHECK** for Check constraints

**NOT NULL** to indicate that an attribute cannot take the
NULL value

## Constraints

```
CREATE TABLE Guest(
  id INTEGER PRIMARY KEY,
  name TEXT NOT NULL,
  email TEXT UNIQUE CHECK (email LIKE '%@%')
);

CREATE TABLE Reservation(
  id INTEGER PRIMARY KEY,
  guest INTEGER NOT NULL REFERENCES Guest(id),
  room INTEGER NOT NULL CHECK (room>0
    AND room<651),
  arrival DATE NOT NULL,
  nights INTEGER NOT NULL CHECK (nights>0),
  UNIQUE(room, arrival),
  UNIQUE(guest, arrival)
);
```

# Updates

- Insertions:
  **INSERT INTO** Guest(id,name) **VALUES** (5,'John');
- Deletions:
  **DELETE FROM** Reservation **WHERE** id>4;
- Modifications:
  **UPDATE** Reservation **SET** room=205 **WHERE** room=204;

Introduction      Relational Model      **Query Languages**      References
ooooooooo         ooooooooooo           ooooooooooooooo          oo
ooooooooooo                             oooooo●oooooooooo

# Updates

**INSERT INTO** Guest **VALUES**
  (1,'Jean Dupont','jean.dupont@gmail.com'),
  (2,'Alice Dupuis','alice@dupuis.name'),
  (3,'Jean Dupont','jean.dupont@ens.fr');

**INSERT INTO** Reservation **VALUES**
  (1,1,504,'2017-01-01',5),
  (2,2,107,'2017-01-10',3),
  (3,3,302,'2017-01-15',6),
  (4,2,504,'2017-01-15',2),
  (5,2,107,'2017-01-30',1);

## Queries

General following form:

**SELECT** ... **FROM** ... **WHERE** ...
**GROUP BY** ... **HAVING** ...
**UNION SELECT** ... **FROM** ...

SELECT projection, renaming, aggregation

FROM cross product

WHERE selection (optional)

GROUP BY grouping (optional)

HAVING selection on the grouping (optional)

UNION union (optional)

Other keywords: ORDER BY to reorder, LIMIT to limit to
first $k$ results, DISTINCT to impose set semantics, EXCEPT
for set difference, etc.

# Renaming

| | Guest | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\rho_{\text{id}\rightarrow\text{guest}}(\text{Guest})$$

**SELECT** id **AS** guest, name, email
**FROM** Guest;

Introduction
0000000
0000000000

Relational Model
00000000000

Query Languages
0000000000000
0000000000●000000

References
00

# Projection

| Guest | | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| Reservation | | | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\Pi_{email,id}(\text{Guest})$$

**SELECT DISTINCT** email, id
**FROM** Guest;

# Selection

<table>
<tr><td colspan="3" align="center">Guest</td></tr>
<tr><td>id</td><td>name</td><td>email</td></tr>
<tr><td>1</td><td>John Smith</td><td>john.smith@gmail.com</td></tr>
<tr><td>2</td><td>Alice Black</td><td>alice@black.name</td></tr>
<tr><td>3</td><td>John Smith</td><td>john.smith@ens.fr</td></tr>
</table>

| | Reservation | | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\sigma_{\texttt{arrival}>\text{2017-01-12}\wedge\texttt{guest}=2}(\texttt{Reservation})$$

**SELECT** $*$
**FROM** Reservation
**WHERE** arrival>'2017-01-12' **AND** guest=2;

## Cross product

| | Guest | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\Pi_{\texttt{id}}(\texttt{Guest}) \times \Pi_{\texttt{name}}(\texttt{Guest})$$

**SELECT** $*$
**FROM**
  (**SELECT DISTINCT** id **FROM** Guest) **AS** temp1,
  (**SELECT DISTINCT** name **FROM** Guest) **AS** temp2
**ORDER BY** name, id;

# Union

| | Guest | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\Pi_{\text{room}}(\sigma_{\text{guest}=2}(\text{Reservation}))$$
$$\cup \Pi_{\text{room}}(\sigma_{\text{arrival}=2017\text{-}01\text{-}15}(\text{Reservation}))$$

```
SELECT room
FROM Reservation
WHERE guest=2
UNION
SELECT room
FROM Reservation
WHERE arrival='2017-01-15';
```

# Difference

| Guest | | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| Reservation | | | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\Pi_{\text{room}}(\sigma_{\text{guest}=2}(\text{Reservation}))$$
$$\setminus \Pi_{\text{room}}(\sigma_{\text{arrival}=2017\text{-}01\text{-}15}(\text{Reservation}))$$

```
SELECT room
FROM Reservation
WHERE guest=2
EXCEPT
  SELECT room
  FROM Reservation
  WHERE arrival='2017-01-15';
```

# Join

| | Guest | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\text{Reservation} \bowtie_{\text{guest}=\text{id}} \text{Guest}$$

**SELECT** Reservation.*, name, email
**FROM** Reservation **JOIN** Guest **ON** guest=Guest.id;

**SELECT** Reservation.*, name, email
**FROM** Reservation, Guest
**WHERE** guest=Guest.id;

# Aggregation

| | Guest | |
|---|---|---|
| id | name | email |
| 1 | John Smith | john.smith@gmail.com |
| 2 | Alice Black | alice@black.name |
| 3 | John Smith | john.smith@ens.fr |

| | | Reservation | | |
|---|---|---|---|---|
| id | guest | room | arrival | nights |
| 1 | 1 | 504 | 2017-01-01 | 5 |
| 2 | 2 | 107 | 2017-01-10 | 3 |
| 3 | 3 | 302 | 2017-01-15 | 6 |
| 4 | 2 | 504 | 2017-01-15 | 2 |
| 5 | 2 | 107 | 2017-01-30 | 1 |

$$\sigma_{\mathrm{avg}>3}(\gamma_{\mathrm{room}}^{\mathrm{avg}}[\lambda x.\mathrm{avg}(x)](\Pi_{\mathrm{room,nights}}(\mathtt{Reservation})))$$

**SELECT** room, **AVG**(nights) **AS avg**
**FROM** Reservation
**GROUP BY** room
**HAVING AVG**(nights)>3
**ORDER BY** room;

# Plan

Introduction

Relational Model

Query Languages

References

# References

- Basics on what data management research is [Benedikt and Senellart, 2012]

- Classic textbook on the foundations of data management, database theory [Abiteboul et al., 1995]

- Modern textbook, being developed, on database theory [Arenas et al., 2022]

- Classic textbook on database systems [Garcia-Molina et al., 2008]

- Details of SQL: standards are not public documents (and not useful for a final user); use the RDBMS documentation instead

- For example, PostgreSQL has a comprehensive documentation at https://www.postgresql.org/docs/ (and using \h in the command line client)

# Bibliography I

Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995. http://webdam.inria.fr/Alice/.

Marcelo Arenas, Pablo Barceló, Leonid Libkin, Wim Martens, and Andreas Pieris. Database theory. https://github.com/pdm-book/community, 2022. Preliminary version.

Michael Benedikt and Pierre Senellart. Databases. In Edward K. Blum and Alfred V. Aho, editors, *Computer Science. The Hardware, Software and Heart of It*, pages 169–229. Springer-Verlag, 2012.

Donald D. Chamberlin and Raymond F. Boyce. SEQUEL: A structured English query language. In *Proc. SIGFIDET/SIGMOD Workshop*, volume 1, 1974.

# Bibliography II

Hector Garcia-Molina, Jeffrey D. Ulman, and Jennifer Widom.
   Pearson, 2008.

ISO. *ISO 9075:1987: SQL*. International Standards
   Organization, 1987.

ISO. *ISO 9075:1999: SQL*. International Standards
   Organization, 1999.

Anthony C. Klug. Equivalence of relational algebra and
   relational calculus query languages having aggregate
   functions. *J. ACM*, 29(3):699–717, 1982.

Leonid Libkin. Expressive power of SQL. *Theor. Comput.
   Sci.*, 296(3):379–404, 2003.