

# DBM1

## IST Semester

Vincent Primault  
*[vincent.primault@insa-lyon.fr](mailto:vincent.primault@insa-lyon.fr)*

# Who am I?

Vincent Primault

- 3<sup>rd</sup> year PhD student at the LIRIS laboratory
- Office: 501.315, in CS department
- Email: [vincent.primault@insa-lyon.fr](mailto:vincent.primault@insa-lyon.fr)

# Course organization

- 20h in total, from now to October, 20th.
- Typical french organisation:
  - 6h of lectures (rather theoretical)
  - 4h of exercises (to apply abstract notions)
  - 8h of labs (practical experiments)
  - 2h of final exam
- Balance between theory and practice.
- Final mark: 50% exam, 50% labs.

# Course organization (cont'd)

- Take notes!
- Ask questions!

# Course outline

- Databases fundamentals
- Relational algebra
- SQL language
- Query processing & optimization
- Distributed databases & NoSQL

# Bibliography

- A huge amount of resources on the Web, but be careful...
- J.D. Ullman. **Database and Knowledge-base Systems**. 1988.
- J. Paredaens, P. De Bra, M. Gyssens, D. Van Gucht. **The Structure of the Relational Model**. 1989.
- H. Garcia-Molina, J. D. Ullman, J. Widom. **Database Systems**. 1999.
- P. O'Neil. **Database: Principles, Programming, Performance**. 1994.
- S. Abiteboul, R. Hull, V. Vianu. **Foundations of Databases**. 1995.

# DBM1

## Part 1: Introduction

Vincent Primault  
*[vincent.primault@insa-lyon.fr](mailto:vincent.primault@insa-lyon.fr)*

# From files...

- Data can always be stored within ordinary files
  - Various formats: plain text, Excel, CSV, JSON, YAML, XML...
  - More or less structured

Last name	First name	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01

Students.xls

Last name	First name	Course	Mark
Doe	John	DBM1	13
Doe	John	OPS	15
Smith	Jane	DBM1	12
Smith	Jane	OPS	17

Marks.xls



# From files... (cont'd)

- Data is duplicated (first/last name)
  - What happens if one of this information changes?
  - What happens if there is a mistake while typing a name?
- Nothing prevents to put inconsistent data (whether it is on purpose or not). E.g, a mark more than 20.
- How to uniquely identify student (thing about homonyms)?
- How to compute the average mark?

## ... to databases

- A database is a structured way to store data.
- A database management system (DBMS) is a software that acts as an interface between users or other applications and a database. It includes creating, modifying, deleting and querying data.
- Databases still store information into text and/or binary files!
  - But this is an implementation detail that differ from one DBMS to the other.
  - DBMS hides this complexity by providing a clean interface.
  - DBMS rely on standards (like SQL) to ensure interoperability.

# The promises of databases

- **Atomicity:** guarantee that a set of operations can be processed in a row.
- **Consistency:** possibility to put constraints on data.
- **Isolation:** it is possible for several users to interact with the same database.
- **Durability:** result of operations are stored permanently, and resilient to software and hardware errors.
- And more:
  - **Security:** access rights can be enforced with a fine granularity.
  - Powerful query language: **SQL**.

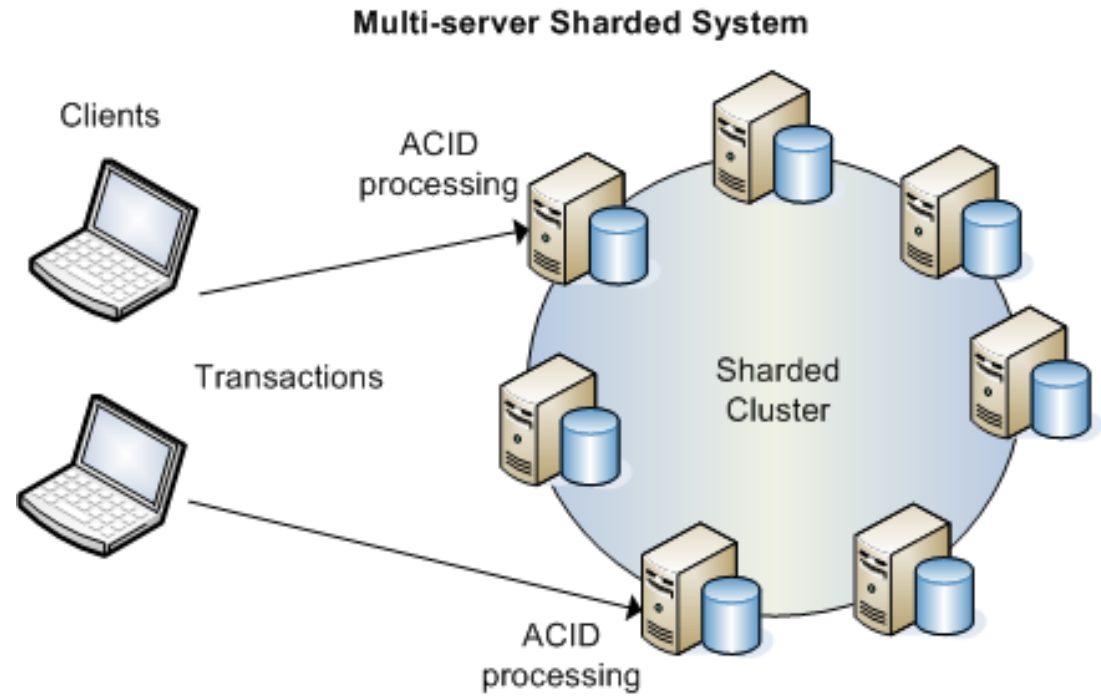
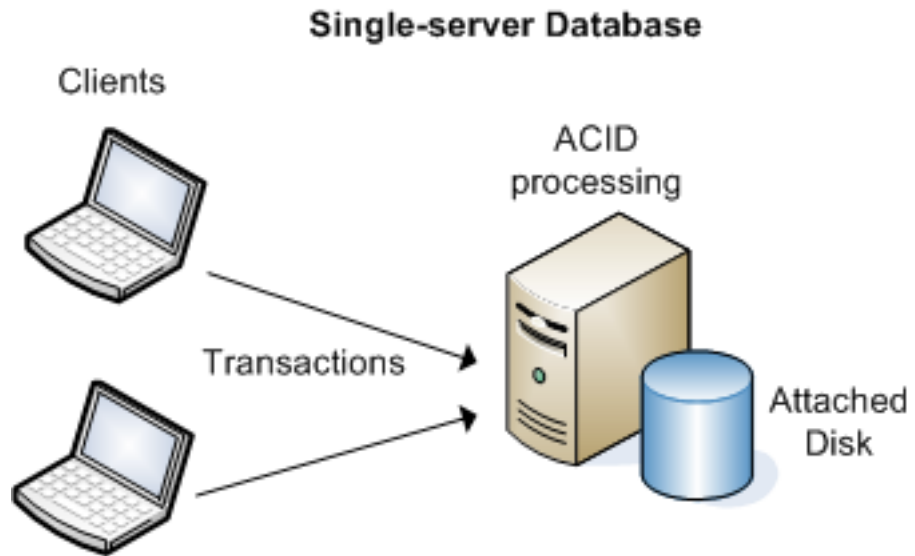
# A lot of players



Relational and centralized DBMS

"NoSQL" and distributed DBMS

# Centralized vs distributed

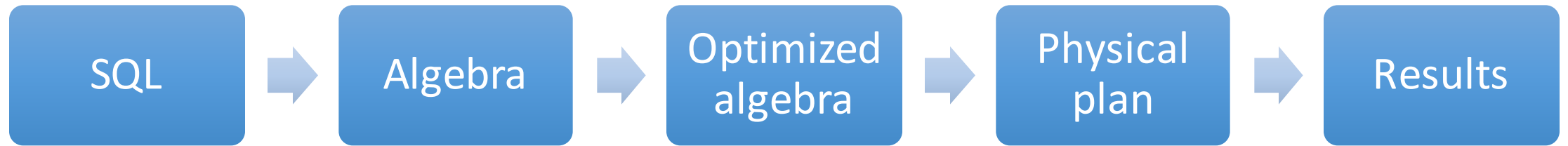


# DBM1

## Part 2: Relational algebra

Vincent Primault  
*[vincent.primault@insa-lyon.fr](mailto:vincent.primault@insa-lyon.fr)*

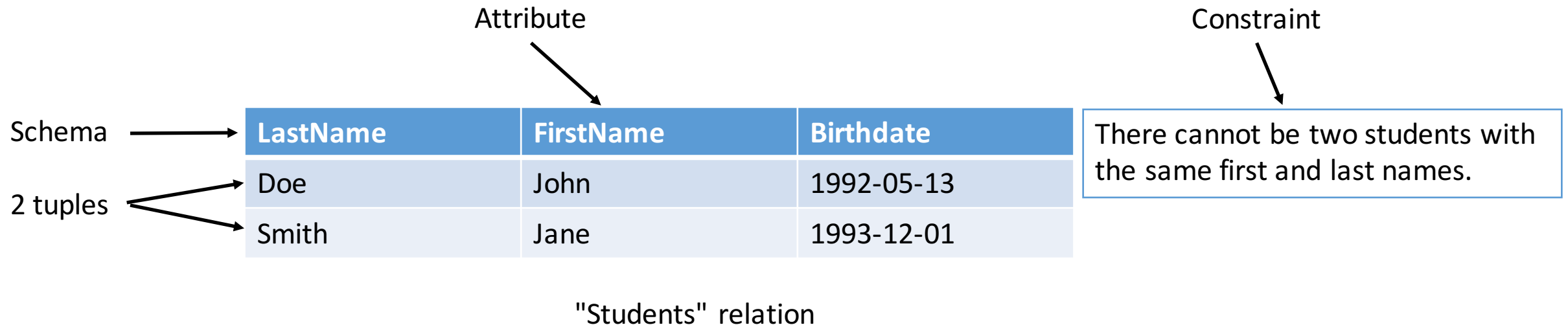
# Why learn relational algebra?



The life of a typical SQL query

# The relational model

- A *relation* is a set of *tuples*, bound to a *schema*.
- A schema contains *attributes* and constraints on them.
- Each attribute is associated with a *domain* of values.





# Some notations

Relation schema:  $(\Omega, \Delta, \text{dom}, \text{SC})$

- $\Omega$ : set of attributes
- $\Delta$ : set of domains
- $\text{dom}: \Omega \rightarrow \Omega$
- $\text{SC}$ : set of constraints

Tuple  $t$

- $t: \Omega \rightarrow \bigcup \Delta, t(A) \in \text{dom}(A)$
- It is a function!

# Relation instance

Schema of Students =  $(\Omega, \Delta, \text{dom}, \text{SC})$

$\Omega = \{\text{LastName}, \text{FirstName}, \text{Birthdate}\}$

$\Delta = \{\text{set of dates, set of strings}\}$

$\text{dom}(\text{LastName}) = \text{set of strings}$

$\text{dom}(\text{FirstName}) = \text{set of strings}$

$\text{dom}(\text{Birthdate}) = \text{set of dates}$

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01

An instance of Students

SC contains the following constraint:

1. There are no two different tuples with the same (FirstName, LastName) pair.

# A relational database

A relational database is a finite set of relations SR.

If two relations  $i$  and  $j$  both have an attribute  $A$  such that

$$A \in \Omega_i \cap \Omega_j$$

then

$$\text{dom}_i(A) = \text{dom}_j(A)$$

and they have the exact same meaning.

# What is a relational algebra?

- An algebra is a mathematical system consisting of:
  - Operands : variables or values from which new values can be constructed.
  - Operators: symbols denoting procedures that construct new values from given values.
- A relational algebra is an algebra whose operands are relations or variables that represent relations.
- Operators are designed to do the most common things that we need to do with relations in a database. The result is an algebra that can be used as a query language for relations.

# Relational algebra operators

- Set operations: **union**, **intersection** and **difference**.
- **Selection**: picking only certain tuples.
- **Projection**: picking only certain attributes.
- **Renaming** of attributes.
- **Products** and **joins**: composition of relations.
- **Division**.

# Intersection, union, difference

- These are the classical set operators, except that...
- ... **operands need to have the same schema.**
- Notations:
  - $R \cap S$
  - $R \cup S$
  - $R - S$

# Intersection, union, difference (cont'd)

Students2014

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01

Students2015

LastName	FirstName	Birthdate
Brown	Emily	1991-11-18
Doe	John	1992-05-13



Students2014  $\cap$  Students 2015

LastName	FirstName	Birthdate
Doe	John	1992-05-13

Students2014  $\cup$  Students 2015

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01
Brown	Emily	1991-11-18

Students2014  $-$  Students 2015

LastName	FirstName	Birthdate
Smith	Jane	1993-12-01

# Intersection, union, difference (cont'd)

- Schema of  $R$  is  $(\Omega, \Delta, \text{dom})$ .
- Schema of  $S$  is  $(\Omega, \Delta, \text{dom})$ .
  - Schema of  $R \cap S$  is  $(\Omega, \Delta, \text{dom})$
  - Schema of  $R \cup S$  is  $(\Omega, \Delta, \text{dom})$
  - Schema of  $R - S$  is  $(\Omega, \Delta, \text{dom})$
- Property:
  - $R \cap S = R - (R - S)$



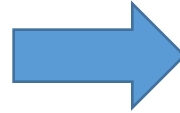
# Projection

- Eliminates columns, consequently removing duplicates across remaining attributes.
- Notation:  $\pi_{A_1, \dots, A_s}(R)$

# Projection (cont'd)

Students

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01
Doe	Anna	1994-03-01
Smith	Lara	1993-12-01
Jones	John	1993-01-16



$\pi_{\text{LastName, FirstName}}(\text{Students})$

LastName	FirstName
Doe	John
Smith	Jane
Doe	Anna
Smith	Lara
Jones	John

$\pi_{\text{LastName, Birthdate}}(\text{Students})$

LastName	Birthdate
Doe	1992-05-13
Smith	1993-12-01
Doe	1994-03-01
Jones	1993-01-16

$\pi_{\text{LastName}}(\text{Students})$

LastName
Doe
Smith
Jones

# Projection (cont'd)

- Schema of  $R$  is  $(\Omega, \Delta, \text{dom})$ .
- Schema of  $\pi_{A_1, \dots, A_s}(R)$  is  $(\{A_1, \dots, A_s\}, \text{dom}(\{A_1, \dots, A_s\}), \text{dom}_{|\{A_1, \dots, A_s\}})$ .
- Let  $r$  be an instance of  $R$ .
  - $\pi_{A_1, \dots, A_s}(r) = \{t_{|\{A_1, \dots, A_s\}} \mid t \in r\}$
- Properties:
  - $\text{card}(\pi_{A_1, \dots, A_s}(r)) \leq \text{card}(r)$
  - $\pi_{\Omega}(r) = r$
  - $\pi_{\emptyset}(r) = \text{singleton with empty tuple}$

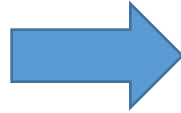
# Selection

- Filters tuples to keep only those satisfying a condition.
- Notation:  $\sigma_C(R)$ 
  - C is a logical condition.
  - The condition must be of the form  $A \text{ op } B$  or  $A \text{ op } b$ , with op being one of  $<$ ,  $>$ ,  $\leq$ ,  $\geq$ ,  $=$ ,  $\neq$ .
  - Multiple conditions can be combined with AND, OR and NOT conjunctions.

# Selection (cont'd)

Students

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01
Doe	Anna	1994-03-01
Smith	Lara	1993-12-01
Jones	John	1993-01-16



$\sigma_{\text{LastName} = \text{Doe}}(\text{Students})$

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Doe	Anna	1994-03-01

$\sigma_{\text{LastName} = \text{Doe AND FirstName} = \text{John}}(\text{Students})$

LastName	FirstName	Birthdate
Doe	John	1992-05-13

$\sigma_{\text{Birthdate} = 1992-05-14}(\text{Students})$

LastName	FirstName	Birthdate
----------	-----------	-----------

# Selection (cont'd)

- Schema of  $R$  is  $(\Omega, \Delta, \text{dom})$ .
- Schema of  $\sigma_C(R)$  is  $(\Omega, \Delta, \text{dom})$ .
- Let  $r$  be an instance of  $R$ .
  - $\sigma_C(r) = \{t \mid t \in r \text{ and } C(t)\}$
- Properties:
  - $\text{card}(\sigma_C(r)) \leq \text{card}(r)$
  - $\sigma_{A=A}(r) = r$
  - $\sigma_{\text{NOT } A=A}(r) = \emptyset$
  - $\sigma_{C \text{ AND } C'}(r) = \sigma_C(R) \cap \sigma_{C'}(R)$
  - $\sigma_{C \text{ OR } C'}(r) = \sigma_C(R) \cup \sigma_{C'}(R)$
  - $\sigma_{\text{NOT } C}(r) = r - \sigma_C(R)$

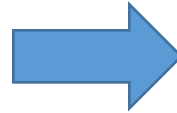
# Renaming

- Gives a new schema to a relation but...
- **It does not change the instance.**
- Notation:  $\rho_{\gamma}(R)$ 
  - $\gamma$  is a one-to-one function,  $\gamma : \Omega \rightarrow \Omega'$

# Renaming (cont'd)

Students

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01
Doe	Anna	1994-03-01
Smith	Lara	1993-12-01
Jones	John	1993-01-16



$\rho_\gamma(\text{Students})$

Nom	Prénom	DateDeNaissance
Doe	John	1992-05-13
Smith	Jane	1993-12-01
Doe	Anna	1994-03-01
Smith	Lara	1993-12-01
Jones	John	1993-01-16

$\gamma(\text{LastName}) = \text{Nom}$

$\gamma(\text{FirstName}) = \text{Prénom}$

$\gamma(\text{Birthdate}) = \text{DateDeNaissance}$



# Renaming (cont'd)

- Schema of  $R$  is  $(\Omega, \Delta, \text{dom})$ .
- Schema of  $\rho_\gamma(R)$  is  $(\gamma(\Omega), \Delta, \text{dom}.\gamma^{-1})$ .
- Let  $r$  be an instance of  $R$ .
  - $\rho_\gamma(r) = \{t.\gamma^{-1} \mid t \in r\}$
- Property:
  - $\text{card}(\Omega) = \text{card}(\Omega')$

# Cartesian product

- Associates every tuple in  $R$  with every tuple in  $S$ .
- Notation:  $R \times S$
- It is very rare in practice, mainly used to express joins.

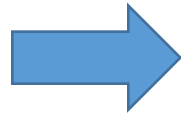
# Cartesian product (cont'd)

Students

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01

FirstNames

FirstName	Gender	Popularity
John	Male	60
Anna	Female	75
Lara	Female	30



Students x FirstNames

LastName <sub>1</sub>	FirstName <sub>1</sub>	Birthdate	FirstName <sub>2</sub>	Gender	Popularity
Doe	John	1992-05-13	John	Male	60
Doe	John	1992-05-13	Anna	Female	75
Doe	John	1992-05-13	Lara	Female	30
Smith	Jane	1993-12-01	John	Male	60
Smith	Jane	1993-12-01	Anna	Female	75
Smith	Jane	1993-12-01	Lara	Female	30

# Cartesian product (cont'd)

- Schema of R is  $(\Omega_R, \Delta_R, \text{dom}_R)$ .
- Schema of S is  $(\Omega_S, \Delta_S, \text{dom}_S)$ .
- Schema of  $R \times S$  is  $(\Omega', \Delta_R \cup \Delta_S, \text{dom}')$ .
  - $\Omega' = \{A_1 \mid A \in \Omega_R\} \cup \{A_2 \mid A \in \Omega_S\}$
  - $\text{dom}' = \text{dom}_R \cdot \gamma_R \cup \text{dom}_S \cdot \gamma_S, \gamma_R(A_1) = A, A \in \Omega_R, \gamma_S(A_2) = A, A \in \Omega_S$
- Let  $r, s$  be an instances of R and S respectively.
  - $r \times s = \{t \cdot \gamma_R \cup t' \cdot \gamma_S \mid t \in r \text{ and } t' \in s\}$
- Properties:
  - $\text{card}(r \times s) = \text{card}(r) \times \text{card}(s)$
  - $r \times \emptyset = \emptyset$

# Natural join

- Associates two relations by their common attributes.
- Notation:  $R \bowtie S$
- Meaning:  $R \bowtie S = \pi_A(\sigma_C(R \times S))$  where
  - C checks equality of common attributes
  - A is the set of common attributes
- This is a reason why it is important that two attributes with the same name in two different relations must represent the same notion!

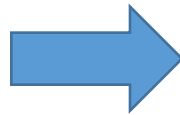
# Natural join (cont'd)

Students

LastName	FirstName	Birthdate
Doe	John	1992-05-13
Smith	Jane	1993-12-01
Doe	Anna	1994-03-01
Smith	Lara	1993-12-01
Jones	John	1993-01-16

FirstNames

FirstName	Gender	Popularity
John	Male	60
Jane	Female	25
Lara	Female	30
James	Male	28
Anna	Female	43



Students ⋈ FirstNames

LastName	FirstName	Birthdate	Gender	Popularity
Doe	John	1992-05-13	Male	60
Smith	Jane	1993-12-01	Female	25
Doe	Anna	1994-03-01	Female	43
Smith	Lara	1993-12-01	Female	30
Jones	John	1993-01-16	Male	60

FirstNames ⋈ Students

FirstName	Gender	Popularity	LastName	Birthdate
John	Male	60	Doe	1992-05-13
John	Male	60	Jones	1993-01-16
Jane	Female	25	Smith	1993-12-01
Lara	Female	30	Smith	1993-12-01

# Natural join (cont'd)

- Schema of R is  $(\Omega_R, \Delta_R, \text{dom}_R)$ .
- Schema of S is  $(\Omega_S, \Delta_S, \text{dom}_S)$ .
  - $\text{dom}_R|_{\Omega_R} \cap \Omega_S = \text{dom}_S|_{\Omega_R} \cap \Omega_S$
- Schema of  $R \bowtie S$  is  $(\Omega_R \cup \Omega_S, \Delta_R \cup \Delta_S, \text{dom}_R \cup \text{dom}_S)$ .
- Let r, s be an instances of R and S respectively.
  - $r \bowtie s = \{t \cup t' \mid t \in r \text{ and } t' \in s\}$
- Properties:
  - $\text{card}(r \bowtie s) \leq \text{card}(r) \times \text{card}(s)$
  - $r \bowtie \emptyset = \emptyset$
  - $r \times s = r \bowtie s$  if  $\Omega_R \cap \Omega_S = \emptyset$
  - $r \cap s = r \bowtie s$  if  $\Omega_R = \Omega_S$

# Theta join

- A join that involves a predicate.
- Notation:  $R \bowtie_{\theta} S$ 
  - $\theta$  is a logical condition.
  - The condition must be of the form  $A \text{ op } B$  or  $A \text{ op } b$ , with  $\text{op}$  being one of  $<$ ,  $>$ ,  $\leq$ ,  $\geq$ ,  $=$ ,  $\neq$ .
  - Multiple conditions can be combined with AND, OR and NOT conjunctions.



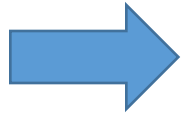
# Theta join (cont'd)

Students

LastName	FirstName	Birthdate
Smith	Jane	1993-12-01
Doe	Anna	1994-03-01
Smith	Lara	1993-12-01

Celebrities

LastName	FirstName	Birthdate
Bieber	Justin	1994-03-01
Smith	Will	1968-09-25
Firth	Colin	1960-09-10



Students  $\bowtie_{\text{LastName}_1=\text{LastName}_2}$  Celebrities

LastName <sub>1</sub>	FirstName <sub>1</sub>	Birthdate <sub>1</sub>	LastName <sub>2</sub>	FirstName <sub>2</sub>	Birthdate <sub>2</sub>
Smith	Jane	1993-12-01	Smith	Will	1968-09-25
Smith	Lara	1993-12-01	Smith	Will	1968-09-25

Students  $\bowtie_{\text{Birthdate}_2 \leq \text{Birthdate}_1}$  Celebrities

LastName <sub>1</sub>	FirstName <sub>1</sub>	Birthdate <sub>1</sub>	LastName <sub>2</sub>	FirstName <sub>2</sub>	Birthdate <sub>2</sub>
Smith	Jane	1993-12-01	Smith	Will	1968-09-25
Smith	Jane	1993-12-01	Firth	Colin	1960-09-10
Doe	Anna	1994-03-01	Bieber	Justin	1994-03-01
Doe	Anna	1994-03-01	Smith	Will	1968-09-25
Doe	Anna	1994-03-01	Firth	Colin	1960-09-10
Smith	Lara	1993-12-01	Smith	Will	1968-09-25
Smith	Lara	1993-12-01	Firth	Colin	1960-09-10

# Theta join (cont'd)

- Schema of R is  $(\Omega_R, \Delta_R, \text{dom}_R)$ .
- Schema of S is  $(\Omega_S, \Delta_S, \text{dom}_S)$ .
- Schema of  $R \bowtie_{\theta} S$  is  $(\Omega', \Delta_R \cup \Delta_S, \text{dom}')$ .
  - $\Omega' = \{A_1 \mid A \in \Omega_R\} \cup \{A_2 \mid A \in \Omega_S\}$
  - $\text{dom}' = \text{dom}_R \cdot \gamma_R \cup \text{dom}_S \cdot \gamma_S, \gamma_R(A_1) = A, A \in \Omega_R, \gamma_S(A_2) = A, A \in \Omega_S$
- Let r, s be an instances of R and S respectively.
  - $r \bowtie_{\theta} s = \{t \mid t \in r \times s \text{ and } \theta(t)\}$
- Properties:
  - $\text{card}(r \bowtie_{\theta} s) = \text{card}(\sigma_{\theta}(R \times S))$
  - $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$

# Division

- Not a primitive operator. Used to express queries such as "*find students that attended all lessons*".
- Notation:  $R \div S$

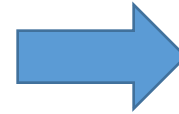
# Division (cont'd)

Attendance

LastName	FirstName	Teaching
Smith	Jane	DBM1
Smith	Jane	DBM2
Smith	Jane	OPS
Doe	John	DBM1
Doe	John	OPS
Jones	John	DBM1
Jones	John	DBM2
Jones	John	OPS

Teachings

Teaching
DBM1
DBM2
OPS



Attendance ÷ Teachings

LastName	FirstName
Smith	Jane
Jones	John

It gives all tuples over (LastName, FirstName) that in combination of every tuple of Teachings results in a tuple of Attendance.

# Division (cont'd)

- Schema of R is  $(\Omega_R, \Delta_R, \text{dom}_R)$ .
- Schema of S is  $(\Omega_S, \Delta_S, \text{dom}_S)$ .
- $\Omega_S \subseteq \Omega_R, \text{dom}_R|_{\Omega_S} = \text{dom}_S$ .
- Schema of  $R \div S$  is  $(\Omega_R - \Omega_S, \Delta_R, \text{dom}_R|_{\Omega_R - \Omega_S})$ .
- Let r, s be an instances of R and S respectively.
  - $r \div s = \{t \mid \forall t' \in s, t \cup t' \in r\}$

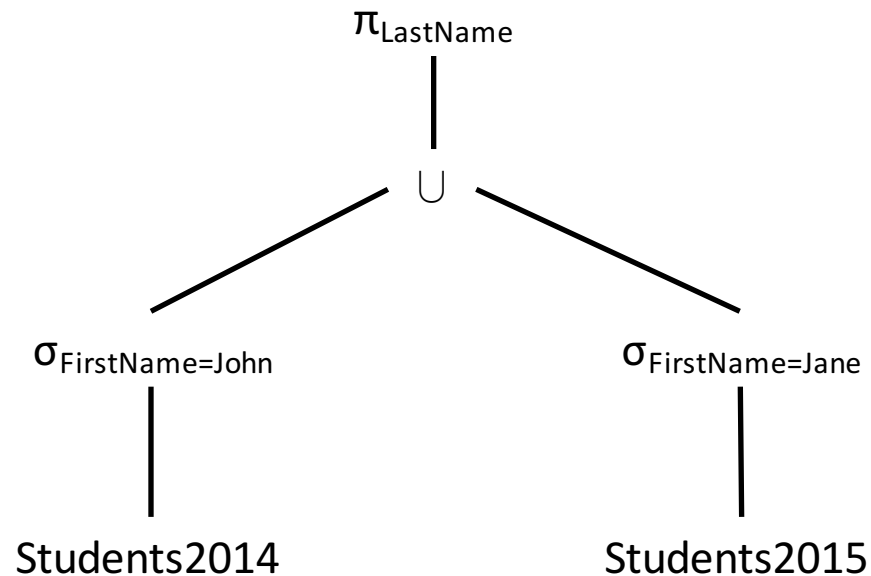
# Building complex expressions

- Combine operators w.r.t. precedence rules of operators:
  1. Projection, selection, renaming (highest)
  2. Cartesian product, joins
  3. Difference
  4. Union, intersection (lowest)
- Use parentheses to force a different order, or just to clarify.
- Example: find the last name of 2014 students named John and those of 2015 named Jane.

$$\pi_{\text{LastName}}(\sigma_{\text{FirstName}=\text{John}}(\text{Students2014}) \cup \sigma_{\text{FirstName}=\text{Jane}}(\text{Students2015}))$$

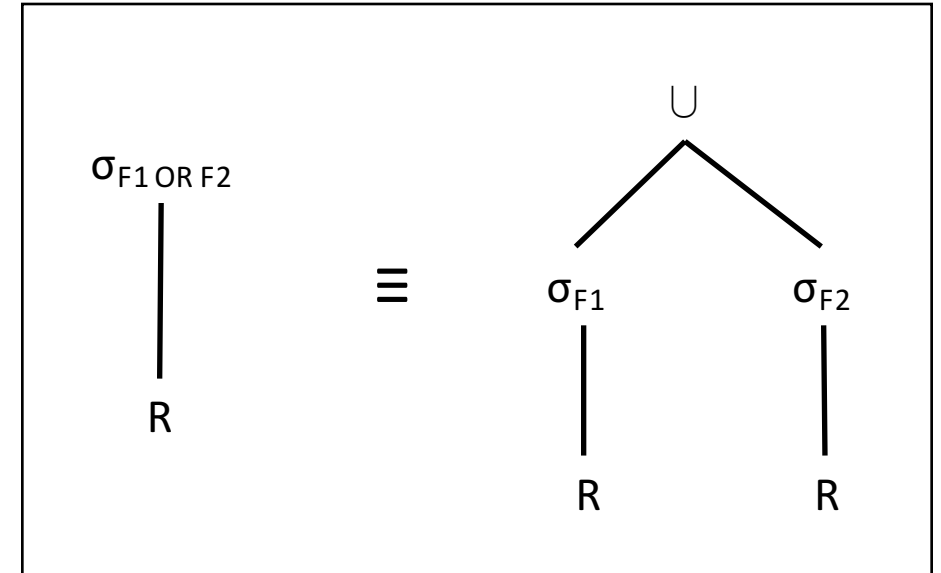
# Expression trees

- Another representation of any relational algebra expression.
- Leaves are operands, usually variables standing for relations.
- Nodes are operators, applied to their child or children.



# Equivalent expressions

- $E_1(R_1, \dots, R_k) \equiv E_2(R_1, \dots, R_k)$  if for all instances  $r_1, \dots, r_k$  of  $R_1, \dots, R_k$  we have  $E_1(r_1, \dots, r_k) = E_2(r_1, \dots, r_k)$
- It is independent of a given database.
- Example:  $\sigma_{F1 \text{ OR } F2}(R) = \sigma_{F1}(R) \cup \sigma_{F2}(R)$





# Exercise

