

# Data Management for Data Science

Lecture 4: Relational Algebra

Prof. Asoc. Endri Raço

# Announcements

- Assignment 1
- Hints and Grading

# Today's Lecture

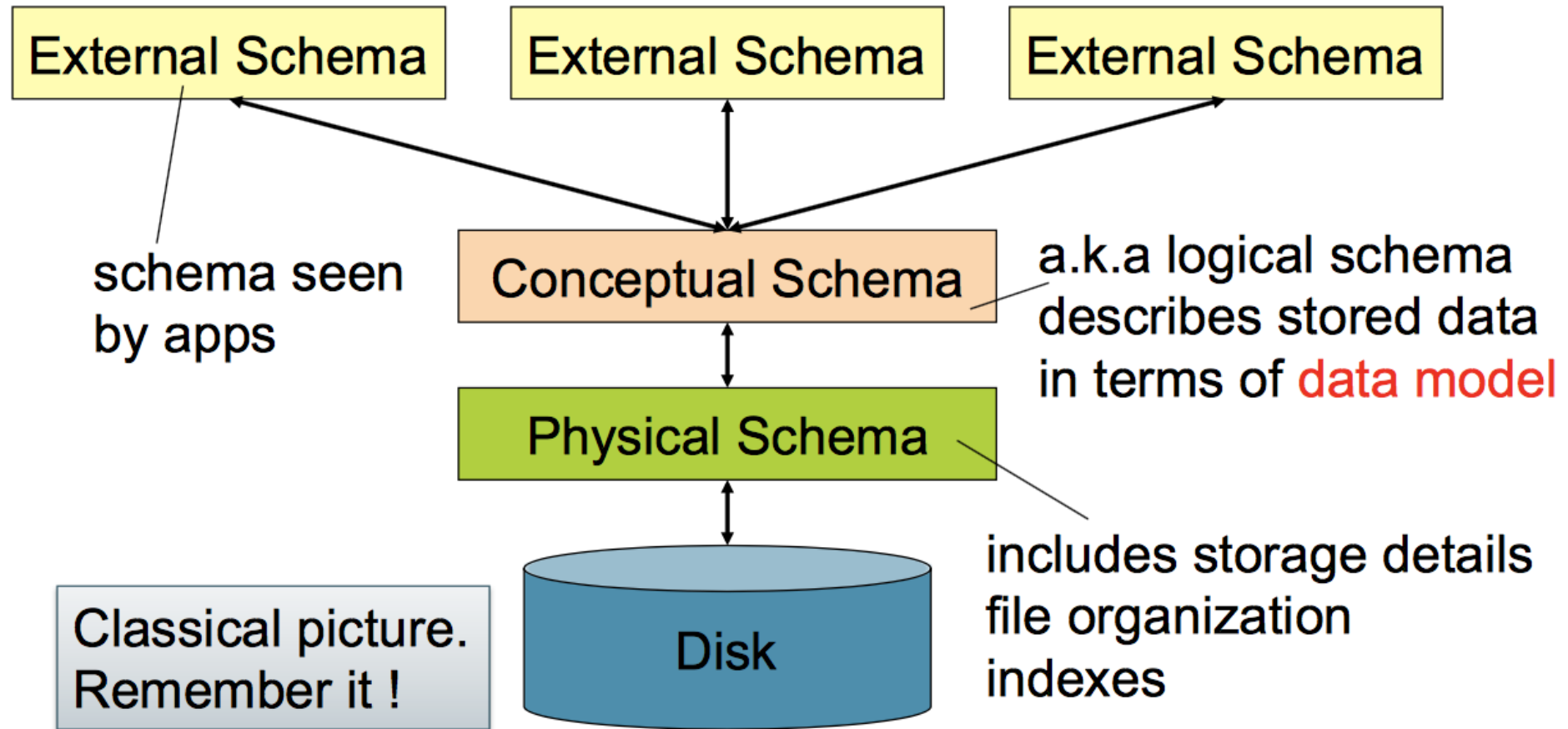
1. The Relational Model & Relational Algebra
2. Relational Algebra Pt. II

# 1. The Relational Model & Relational Algebra

# What you will learn about in this section

1. The Relational Model
2. Relational Algebra: Basic Operators

# Levels of abstraction



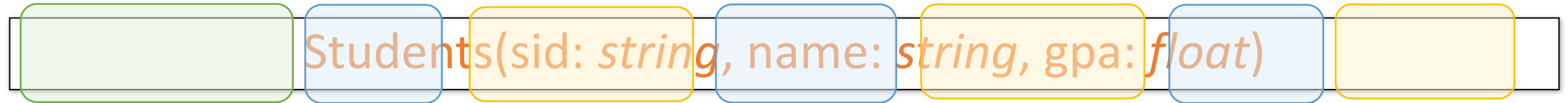
# Motivation

The Relational model is **precise**,  
**implementable**, and we can operate on it  
(query/update, etc.)

Database maps internally into this  
*procedural language*.

# The Relational Model: Schemata

- Relational Schema:



Relation name

*String, float, int, etc.*  
are the **domains** of  
the attributes

Attributes



# The Relational Model: Data

An attribute (or column) is a typed data entry present in each tuple in the relation

Student

sid	name	gpa
001	Bob	3.2
002	Joe	2.8
003	Mary	3.8
004	Alice	3.5

The number of attributes is the arity of the relation

# The Relational Model: Data

**Student**

sid	name	gpa
001	Bob	3.2
002	Joe	2.8
003	Mary	3.8
004	Alice	3.5

The number of tuples is the **cardinality** of the relation

A **tuple** or **row** (or *record*) is a single entry in the table having the attributes specified by the schema

# The Relational Model: Data

**Student**

sid	name	gpa
001	Bob	3.2
002	Joe	2.8
003	Mary	3.8
004	Alice	3.5

In practice DBMSs relax the set requirement, and use multisets.

A relational instance is a *set* of tuples all conforming to the same *schema*

# To Reiterate

- A relational schema describes the data that is contained in a relational instance

Let  $R(f_1:\text{Dom}_1, \dots, f_m:\text{Dom}_m)$  be a relational schema then, an instance of  $R$  is a subset of  $\text{Dom}_1 \times \text{Dom}_2 \times \dots \times \text{Dom}_n$

In this way, a relational schema  $R$  is a **total function from attribute names to types**

# One More Time

- A relational schema describes the data that is contained in a relational instance

A relation  $R$  of arity  $t$  is a function:  
 $R : \text{Dom}_1 \times \dots \times \text{Dom}_t \rightarrow \{0,1\}$

*i.e. returns whether or not a tuple of matching types is a member of it*

Then, the schema is simply the *signature* of the function

Note here that order matters, attribute name doesn't...  
We'll (mostly) work with the other model (last slide) in which **attribute name matters, order doesn't!**

# A relational database

- A relational database schema is a set of relational schemata, one for each relation
- A relational database instance is a set of relational instances, one for each relation

Two conventions:

1. We call relational database instances as simply *databases*
2. We assume all instances are valid, i.e., satisfy the domain constraints

# Remember the CMS

- *Relation DB Schema*

- Students(sid: *string*, name: *string*, gpa: *float*)
- Courses(cid: *string*, cname: *string*, credits: *int*)
- Enrolled(sid: *string*, cid: *string*, grade: *string*)

*Note that the schemas impose effective domain / type constraints, i.e. Gpa can't be "Apple"*

Sid	Name	Gpa
101	Bob	3.2
123	Mary	3.8

Students

## Relation Instances

sid	cid	Grade
123	564	A

Enrolled

cid	cname	credits
564	564-2	4
308	417	2

Courses

## 2<sup>nd</sup> Part of the Model: Querying

```
SELECT S.name  
FROM Students S  
WHERE S.gpa > 3.5;
```

*“Find names of all students  
with GPA > 3.5”*

We don't tell the system *how* or *where* to get the data- **just what we want**, i.e., Querying is *declarative*

To make this happen, we need to translate the *declarative* query into a series of operators... we'll see this next!



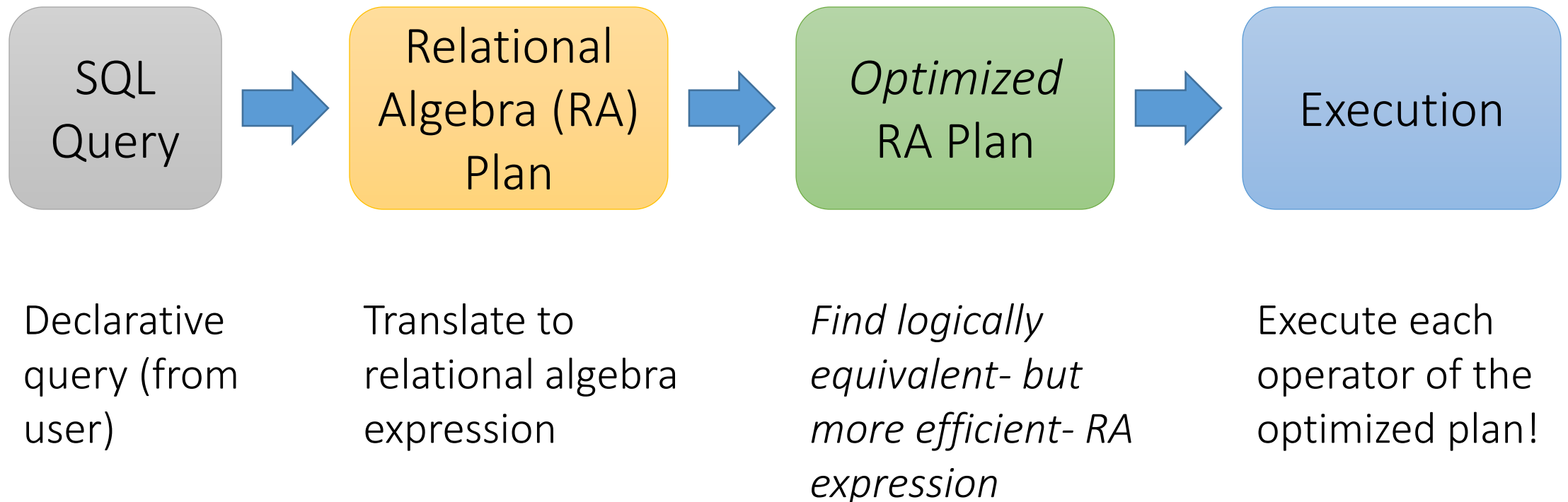
# Virtues of the model

- Physical independence (logical too), Declarative
- Simple, elegant clean: Everything is a relation

# Relational Algebra

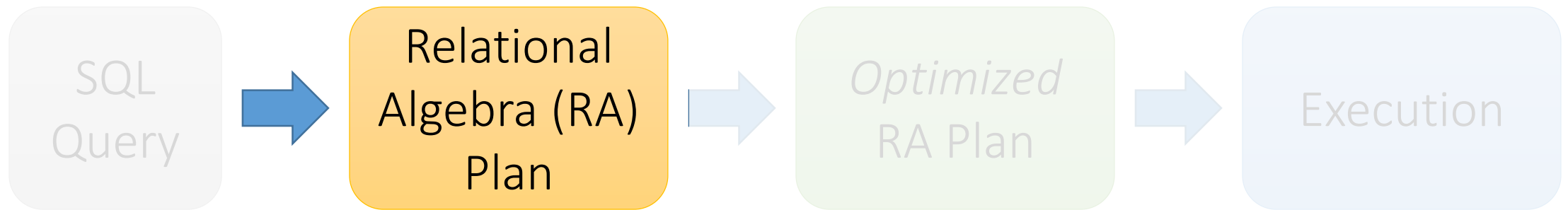
# RDBMS Architecture

How does a SQL engine work ?



# RDBMS Architecture

How does a SQL engine work ?



Relational Algebra allows us to translate declarative (SQL) queries into precise and optimizable expressions!

# Relational Algebra (RA)

- Five **basic** operators:

1. Selection:  $\sigma$
2. Projection:  $\Pi$
3. Cartesian Product:  $\times$
4. Union:  $\cup$
5. Difference:  $-$

- Derived or auxiliary operators:

- Intersection, complement
- Joins (natural, equi-join, theta join, semi-join)
- Renaming:  $\rho$
- Division

# Keep in mind: RA operates on sets!

- RDBMSs use *multisets*, however in relational algebra formalism we will consider **sets**!
- Also: we will consider the ***named perspective***, where every attribute must have a unique name
  - → attribute order does not matter...

Now on to the basic RA operators...

# 1. Selection ( $\sigma$ )

- Returns all tuples which satisfy a condition
- Notation:  $\sigma_c(R)$
- Examples
  - $\sigma_{\text{Salary} > 40000}(\text{Employee})$
  - $\sigma_{\text{name} = \text{"Smith"}}(\text{Employee})$
- The condition  $c$  can be  $=, <, \leq, >, \geq, <>$

Students(sid,sname,gpa)

SQL:

```
SELECT *  
FROM Students  
WHERE gpa > 3.5;
```



RA:

$\sigma_{gpa > 3.5}(\text{Students})$

Another example:

SSN	Name	Salary
1234545	John	20000
5423341	Smith	600000
4352342	Fred	500000

$\sigma_{\text{Salary} > 40000}$  (Employee)



SSN	Name	Salary
5423341	Smith	600000
4352342	Fred	500000



## 2. Projection ( $\Pi$ )

- Eliminates columns, then removes duplicates
- Notation:  $\Pi_{A1, \dots, An}(R)$
- Example: project social-security number and names:
  - $\Pi_{SSN, Name}(Employee)$
  - Output schema: Answer(SSN, Name)

Students(sid,sname,gpa)

SQL:

```
SELECT DISTINCT  
  sname,  
  gpa  
FROM Students;
```



RA:

$\Pi_{sname,gpa}(Students)$

Another example:

SSN	Name	Salary
1234545	John	200000
5423341	John	600000
4352342	John	200000

$\Pi_{\text{Name,Salary}}(\text{Employee})$



Name	Salary
John	200000
John	600000

# Note that RA Operators are Compositional!

Students(sid,sname,gpa)

```
SELECT DISTINCT  
  sname,  
  gpa  
FROM Students  
WHERE gpa > 3.5;
```

How do we represent  
this query in RA?


$$\Pi_{sname,gpa}(\sigma_{gpa>3.5}(Students))$$

$$\sigma_{gpa>3.5}(\Pi_{sname,gpa}(Students))$$

Are these logically equivalent?

### 3. Cross-Product ( $\times$ )

- Each tuple in R1 with each tuple in R2
- Notation:  $R1 \times R2$
- Example:
  - Employee  $\times$  Dependents
- Rare in practice; mainly used to express joins

```
Students(sid,sname,gpa)  
People(ssn,pname,address)
```

SQL:

```
SELECT *  
FROM Students, People;
```



RA:

*Students  $\times$  People*

Another example: People

ssn	pname	address
1234545	John	216 Rosse
5423341	Bob	217 Rosse

×

Students

sid	sname	gpa
001	John	3.4
002	Bob	1.3

*Students × People*



ssn	pname	address	sid	sname	gpa
1234545	John	216 Rosse	001	John	3.4
5423341	Bob	217 Rosse	001	John	3.4
1234545	John	216 Rosse	002	Bob	1.3
5423341	Bob	216 Rosse	002	Bob	1.3

# Renaming ( $\rho$ )

- Changes the schema, not the instance
- A 'special' operator- neither basic nor derived
- Notation:  $\rho_{B1,\dots,Bn}(R)$
- **Note: this is shorthand for the proper form (since names, not order matters!):**
  - $\rho_{A1 \rightarrow B1, \dots, An \rightarrow Bn}(R)$

Students(sid,sname,gpa)

SQL:

```
SELECT  
  sid AS studId,  
  sname AS name,  
  gpa AS gradePtAvg  
FROM Students;
```



RA:

$\rho_{studId,name,gradePtAvg}(Students)$

We care about this operator *because* we are working in a *named perspective*

Another example:

Students

sid	sname	gpa
001	John	3.4
002	Bob	1.3

$\rho_{studId,name,gradePtAvg}(Students)$



Students

studId	name	gradePtAvg
001	John	3.4
002	Bob	1.3

# Natural Join ( $\bowtie$ )

- Notation:  $R_1 \bowtie R_2$
- Joins  $R_1$  and  $R_2$  on *equality of all shared attributes*
  - If  $R_1$  has attribute set  $A$ , and  $R_2$  has attribute set  $B$ , and they share attributes  $A \cap B = C$ , can also be written:  $R_1 \bowtie_C R_2$
- Our first example of a *derived* RA operator:
  - Meaning:  $R_1 \bowtie R_2 = \Pi_{A \cup B}(\sigma_{C=D}(\rho_{C \rightarrow D}(R_1) \times R_2))$
  - Where:
    - The rename  $\rho_{C \rightarrow D}$  renames the shared attributes in one of the relations
    - The selection  $\sigma_{C=D}$  checks equality of the shared attributes
    - The projection  $\Pi_{A \cup B}$  eliminates the duplicate common attributes

Students(sid,name,gpa)  
People(ssn,name,address)

SQL:

```
SELECT DISTINCT
  ssid, S.name, gpa,
  ssn, address
FROM
  Students S,
  People P
WHERE S.name = P.name;
```



RA:

*Students*  $\bowtie$  *People*



Another example:

Students S

sid	S.name	gpa
001	John	3.4
002	Bob	1.3



People P

ssn	P.name	address
1234545	John	216 Rosse
5423341	Bob	217 Rosse

*Students ⋈ People*



sid	S.name	gpa	ssn	address
001	John	3.4	1234545	216 Rosse
002	Bob	1.3	5423341	216 Rosse

# Natural Join

- Given schemas  $R(A, B, C, D)$ ,  $S(A, C, E)$ , what is the schema of  $R \bowtie S$  ?
- Given  $R(A, B, C)$ ,  $S(D, E)$ , what is  $R \bowtie S$  ?
- Given  $R(A, B)$ ,  $S(A, B)$ , what is  $R \bowtie S$  ?

# Example: Converting SQL Query -> RA

Students(sid,sname,gpa)  
People(ssn,sname,address)

```
SELECT DISTINCT  
  gpa,  
  address  
FROM Students S,  
     People P  
WHERE gpa > 3.5 AND  
       sname = pname;
```

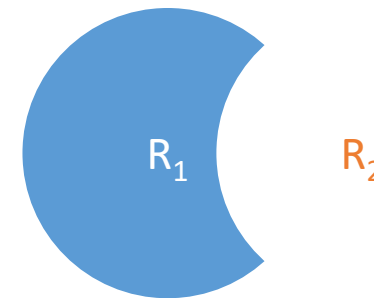
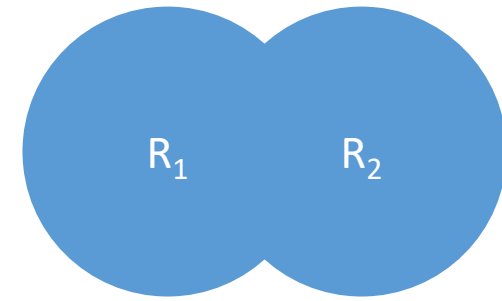

$$\Pi_{gpa,address}(\sigma_{gpa>3.5}(S \bowtie P))$$

# Logical Equivalence of RA Plans

- Given relations  $R(A,B)$  and  $S(B,C)$ :
  - Here, projection & selection commute:
    - $\sigma_{A=5}(\Pi_A(R)) = \Pi_A(\sigma_{A=5}(R))$
  - What about here?
    - $\sigma_{A=5}(\Pi_B(R)) \stackrel{?}{=} \Pi_B(\sigma_{A=5}(R))$

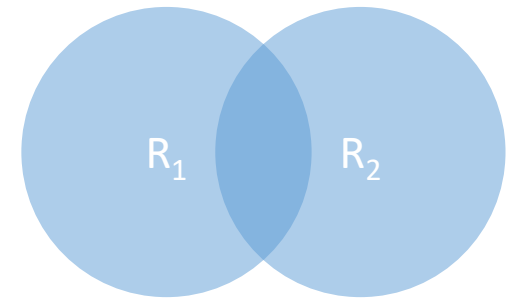
# 1. Union ( $\cup$ ) and 2. Difference ( $-$ )

- $R_1 \cup R_2$
- Example:
  - $\text{ActiveEmployees} \cup \text{RetiredEmployees}$
- $R_1 - R_2$
- Example:
  - $\text{AllEmployees} - \text{RetiredEmployees}$

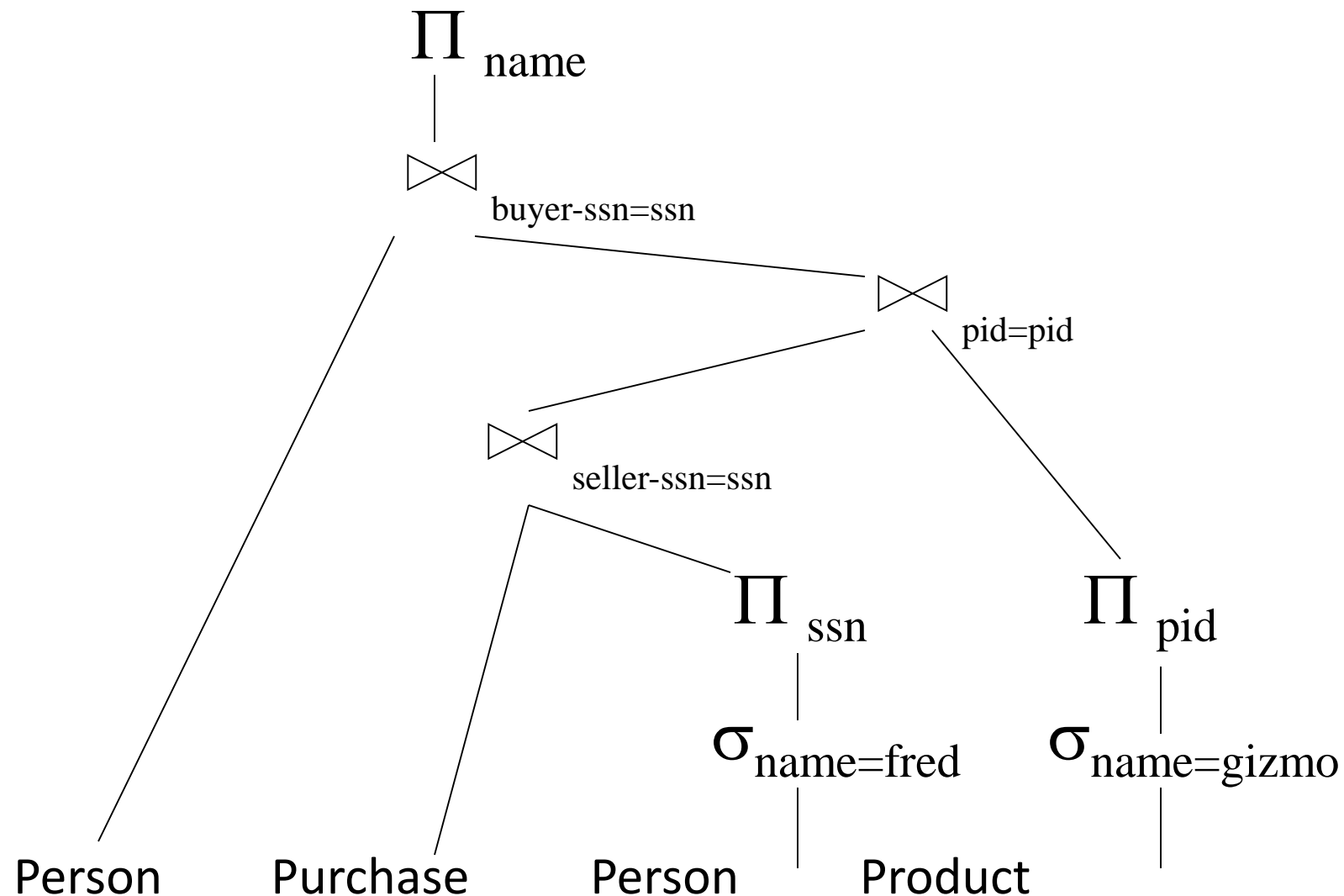


# What about Intersection ( $\cap$ ) ?

- It is a derived operator
- $R1 \cap R2 = R1 - (R1 - R2)$
- Also expressed as a join!
- Example
  - `UnionizedEmployees`  $\cap$  `RetiredEmployees`



# RA Expressions Can Get Complex!



# RA has Limitations !

- Cannot compute “transitive closure”

Name1	Name2	Relationship
Fred	Mary	Father
Mary	Joe	Cousin
Mary	Bill	Spouse
Nancy	Lou	Sister

- Find all direct and indirect relatives of Fred
- Cannot express in RA !!!
  - Need to write C program, use a graph engine, or modern SQL...