

# Topic 1: Data models and query languages

## Unit 3: Relational Algebra (RA)

### Lecture 7

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp23)

<https://northeastern-datalab.github.io/cs7240/sp23/>

1/31/2023

# Pre-class conversations

- Last class summary
- Please keep on pointing out any errors on the slides
- It is time to start to hand in your first scribes (some ideas today)
- Project discussions (in 2 weeks: Fri 2/17: project ideas)
  
- today:
  - Algebra: independence and Codd's theorem
- next time:
  - Recursion (Datalog)

# Algebra and the connection to logic and queries

- Algebra
- Relational Algebra
  - Operators
  - Independence
  - Power of algebra: optimizations
- Equivalence RA and safe RC (Codd's theorem)
  - $RA \rightarrow RC$
  - $RC \rightarrow RA$

# Write in RA

Employee(id, name, salary)



*Q: Find the ID and name of those employees who earn more than the employee whose ID is 123?*





*Q: Find the ID and name of those employees who earn more than the employee whose ID is 123?*

$$\pi_{e.id, e.name} \left( \sigma_{e.salary > o.salary} (\rho_e(\text{employee}) \times \sigma_{id=123}(\rho_o(\text{employee}))) \right)$$

$$\pi_{id, name} \left( \sigma_{salary > s} \left( \text{employee} \times (\rho_{salary \rightarrow s} (\pi_{salary} (\sigma_{id=123}(\text{employee})))) \right) \right)$$

$$\pi_{\$1, \$2} \left( \sigma_{\$4=123 \wedge \$3 > \$6} (\text{employee} \times \text{employee}) \right)$$

# Relational Algebra (RA) operators

- Five basic operators:

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

- Auxiliary (or special) operator

6. Renaming:  $\rho$  ("rho")

- Derived (or implied) operators

7. Joins  $\bowtie$  (natural, theta join, equi-join, [semi-join: moved to T3-U1])
8. Intersection / complement
9. Division

*Derived relational operators:*

- *can be expressed in basic RA; thus not needed*

*But enhancing the basic operator set with derived operators is a good idea:*

- *Queries become easier to write/understand/maintain*
- *Easier for DBMS to apply specialized optimizations (recall the conceptual evaluation strategy)*

*we discuss later in class in detail (SJs are at the heart of efficient algorithms)*

*most important*

# 7a. Natural Join ( $\bowtie$ )

Product(pname, price, category, cid)  
Company(cid, cname, stockprice, country)



- Notation:  $R \bowtie S$
- Joins R and S on equality of all shared attributes
  - Only makes sense in **named perspective!**
  - If R has attribute set **A**, and S has attribute set **B**, and they share attributes  $A \cap B = C$ , can also be written as  $R \bowtie_C S$
- Natural join in basic RA:
  - Meaning:  $R \bowtie S = \pi_{A \cup B}(\sigma_{R.C=S.C}(R \times S))$
  - Meaning:  $R \bowtie S = \pi_{A \cup B}(\sigma_{C=D}(\rho_{C \rightarrow D}(R) \times S))$ 
    - 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

SQL

```
SELECT pname, price, category,  
P.cid, cname, stockprice, country  
FROM Product P, Company C  
WHERE P.cid= C.cid
```

SQL (alternative syntax)



# 7a. Natural Join ( $\bowtie$ )

```
Product(pname, price, category, cid)
Company(cid, cname, stockprice, country)
```



- Notation:  $R \bowtie S$
- Joins R and S on equality of all shared attributes
  - Only makes sense in named perspective!
  - If R has attribute set A, and S has attribute set B, and they share attributes  $A \cap B = C$ , can also be written as  $R \bowtie_C S$
- Natural join in basic RA:
  - Meaning:  $R \bowtie S = \pi_{A \cup B}(\sigma_{R.C=S.C}(R \times S))$
  - Meaning:  $R \bowtie S = \pi_{A \cup B}(\sigma_{C=D}(\rho_{C \rightarrow D}(R) \times S))$ 
    - 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

SQL

```
SELECT pname, price, category,
P.cid, cname, stockprice, country
FROM Product P, Company C
WHERE P.cid= C.cid
```

SQL (alternative syntax)

```
SELECT *
FROM Product
NATURAL JOIN Company
```

RA:





# 7a. Natural Join ( $\bowtie$ )

Product(pname, price, category, cid)  
Company(cid, cname, stockprice, country)



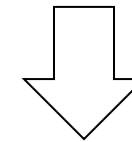
- Notation:  $R \bowtie S$
- Joins R and S on equality of all shared attributes
  - Only makes sense in named perspective!
  - If R has attribute set A, and S has attribute set B, and they share attributes  $A \cap B = C$ , can also be written as  $R \bowtie_C S$
- Natural join in basic RA:
  - Meaning:  $R \bowtie S = \pi_{A \cup B}(\sigma_{R.C=S.C}(R \times S))$
  - Meaning:  $R \bowtie S = \pi_{A \cup B}(\sigma_{C=D}(\rho_{C \rightarrow D}(R) \times S))$ 
    - 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

SQL

```
SELECT pname, price, category,  
P.cid, cname, stockprice, country  
FROM Product P, Company C  
WHERE P.cid= C.cid
```

SQL (alternative syntax)

```
SELECT *  
FROM Product  
NATURAL JOIN Company
```



RA:

Product  $\bowtie$  Company

# 7a. Natural Join ( $\bowtie$ ): an alternative perspective

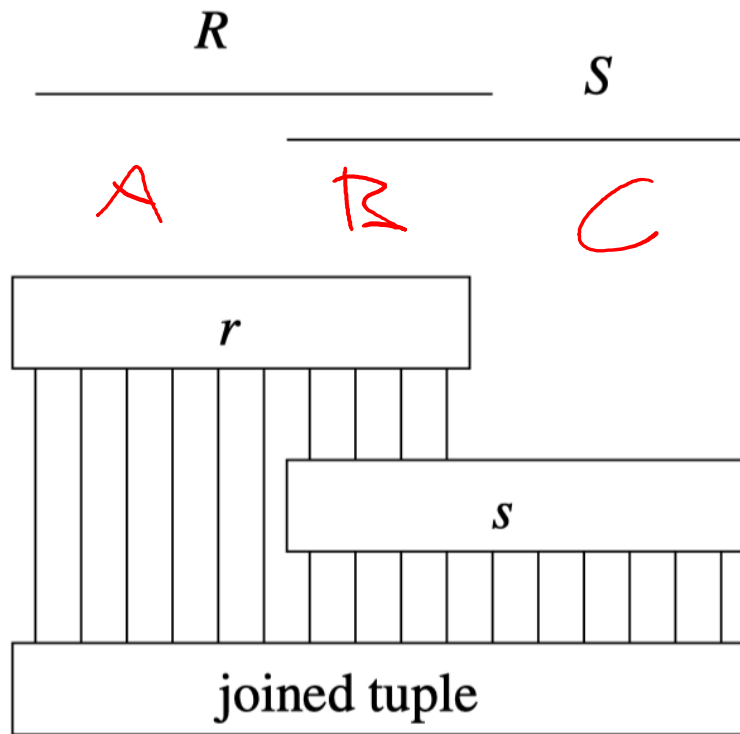


Figure 15: Joining tuples

We only want to pair those tuples that match in some way.

More formally the semantics of the natural join are defined as follows:

$$R \bowtie S = \{r \cup s \mid r \in R \wedge s \in S \wedge Fun(r \cup s)\} \quad (1)$$

where  $Fun(t)$  is a **predicate** that is true for a **relation**  $t$  (in the mathematical sense) **iff**  $t$  is a function. It is usually required that  $R$  and  $S$  must have at least one common attribute, but if this constraint is omitted, and  $R$  and  $S$  have no common attributes, then the natural join becomes exactly the Cartesian product.



# 7a. Natural Join ( $\bowtie$ ): An example

**R**

A	B
1	2
3	4

**S**

B	C	D
2	5	6
4	7	8
9	10	11

$$\rho_{B \rightarrow E}(R) \times S$$





# 7a. Natural Join ( $\bowtie$ ): An example

**R**

A	B
1	2
3	4

**S**

B	C	D
2	5	6
4	7	8
9	10	11

$R \bowtie S$



$\rho_{B \rightarrow E}(R) \times S$

A	E	B	C	D
1	2	2	5	6
1	2	4	7	8
1	2	9	10	11
3	4	2	5	6
3	4	4	7	8
3	4	9	10	11



# 7a. Natural Join ( $\bowtie$ ): An example

**R**

A	B
1	2
3	4

**S**

B	C	D
2	5	6
4	7	8
9	10	11

$R \bowtie S$

A	B	C	D
1	2	5	6
3	4	7	8

$R \bowtie S =$

*in basic RA*



$\rho_{B \rightarrow E}(R) \times S$

A	E	B	C	D
1	2	2	5	6
<del>1</del>	<del>2</del>	<del>4</del>	<del>7</del>	<del>8</del>
<del>1</del>	<del>2</del>	<del>9</del>	<del>10</del>	<del>11</del>
<del>3</del>	<del>4</del>	<del>2</del>	<del>5</del>	<del>6</del>
3	4	4	7	8
<del>3</del>	<del>4</del>	<del>9</del>	<del>10</del>	<del>11</del>



# 7a. Natural Join ( $\bowtie$ ): An example

**R**

A	B
1	2
3	4

**S**

B	C	D
2	5	6
4	7	8
9	10	11

$R \bowtie S$

A	B	C	D
1	2	5	6
3	4	7	8

$R \bowtie S =$

$$\Pi_{AR.BCD}(\sigma_{R.B=S.B}(R \times S)) =$$

$$\Pi_{ABCD}(\sigma_{B=E}(\rho_{B \rightarrow E}(R) \times S))$$

$\rho_{B \rightarrow E}(R) \times S$

A	E	B	C	D
1	2	2	5	6
1	2	4	7	8
1	2	9	10	11
3	4	2	5	6
3	4	4	7	8
3	4	9	10	11



## 7a. Natural Join ( $\bowtie$ ): practice

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





## 7a. Natural Join ( $\bowtie$ ): practice

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

Answer(A, B, C, D,E)

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

?





# 7a. Natural Join ( $\bowtie$ ): practice

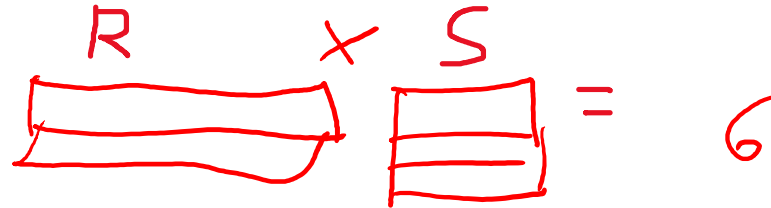
- Given schemas  $R(A, B, C, D)$ ,  $S(A, C, E)$ , what is the schema of  $R \bowtie S$  ?

Answer(A, B, C, D,E)

no condition in the selection that could be violated:

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

$R \times S$



S

A	B	C	D	E

- Given  $R(A, B)$ ,  $S(A, B)$ , what is  $R \bowtie S$  ?





# 7a. Natural Join ( $\bowtie$ ): practice

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

Answer(A, B, C, D,E)

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

$R \times S$

- Given  $R(A, B)$ ,  $S(\cancel{A}, \cancel{B})$ , what is  $R \bowtie S$  ?

$R \cap S$





# 7a. Natural Join ( $\bowtie$ ): practice

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

Answer(A, B, C, D,E)

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

$R \times S$

- Given  $R(A, B)$ ,  $S(A, B)$ , what is  $R \bowtie S$  ?

$R \cap S$





## 7b. Theta Join ( $\bowtie_{\theta}$ )

- A join that involves a predicate

$$R_1 \bowtie_{\theta} R_2 = \sigma_{\theta}(R_1 \times R_2)$$

- $\theta$  ("theta") can be any condition
- No projection: #attributes in output = sum #attributes in input *Note that natural join is a theta join + a selection*
- Example: **band-joins** for approx. matchings across tables

AnonPatient (age, zip, disease)  
Voters (name, age, zip)

*Assume relatively fresh data (within 1 year)*



# 7b. Theta Join ( $\bowtie_{\theta}$ )



- A join that involves a predicate

$$R_1 \bowtie_{\theta} R_2 = \sigma_{\theta}(R_1 \times R_2)$$

- $\theta$  ("theta") can be any condition
- No projection: #attributes in output = sum #attributes in input *Note that natural join is a theta join + a selection*
- Example: **band-joins** for approx. matchings across tables

AnonPatient (age, zip, disease)  
Voters (name, age, zip)

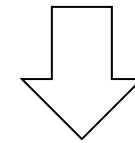
*Assume relatively fresh data (within 1 year)*

$$A \bowtie_{P.zip=V.zip \wedge (P.age \geq V.age - 1 \wedge P.age \leq V.age + 1)} V$$

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

SQL:

```
SELECT *  
FROM  
  Students, People  
WHERE  $\theta$ 
```



RA:



# 7b. Theta Join ( $\bowtie_{\theta}$ )



- A join that involves a predicate

$$R_1 \bowtie_{\theta} R_2 = \sigma_{\theta}(R_1 \times R_2)$$

- $\theta$  ("theta") can be any condition
- No projection: #attributes in output = sum #attributes in input *Note that natural join is a theta join + a selection*
- Example: **band-joins** for approx. matchings across tables

AnonPatient (age, zip, disease)  
Voters (name, age, zip)

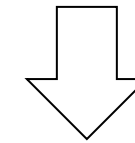
*Assume relatively fresh data (within 1 year)*

$$A \bowtie_{P.zip=V.zip \wedge P.age \geq V.age - 1 \wedge P.age \leq V.age + 1} V$$

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

SQL:

```
SELECT *  
FROM  
    Students, People  
WHERE  $\theta$ 
```



RA:

Students  $\bowtie_{\theta}$  People

## 7c. Equi-join ( $\bowtie_{A=B}$ )

- A theta join where q is an equality

$$R_1 \bowtie_{A=B} R_2 = \sigma_{A=B}(R_1 \times R_2)$$

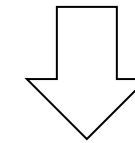
- Example over Gizmo DB:
  - Product  $\bowtie_{\text{manufacturer=cname}}$  Company
- Most common join in practice!



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

SQL:

```
SELECT *
FROM
  Students S, People P
WHERE sname = pname
```



RA:



## 7c. Equi-join ( $\bowtie_{A=B}$ )

- A theta join where q is an equality

$$R_1 \bowtie_{A=B} R_2 = \sigma_{A=B}(R_1 \times R_2)$$

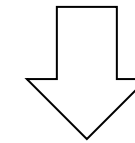
- Example over Gizmo DB:
  - Product  $\bowtie_{\text{manufacturer=cname}}$  Company
- Most common join in practice!



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

SQL:

```
SELECT *
FROM
  Students S, People P
WHERE sname = pname
```



RA:

$S \bowtie_{\text{sname=pname}} P$

What is the connection with a natural join?



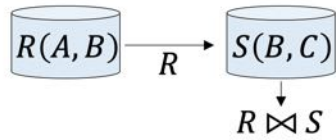


# 7d. Semi-join ( $\bowtie$ ) [moved to T3-U1]

- $R \bowtie S$ : Return tuples from  $R$  for which there is a matching tuple in  $S$  that is equal on their common attribute names.

## Semijoins as Message Passing

- **Semijoins** can reduce network use for equijoins in distributed databases

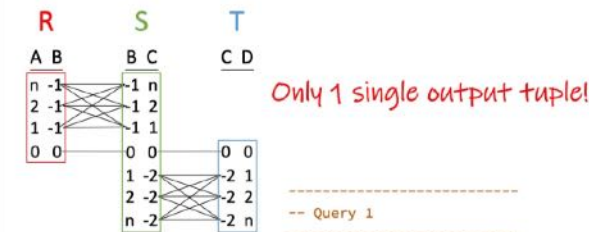


Effective if 1) the size of join attribute  $B$  (or  $n_{AB}$ ) is smaller than  $A$  and  $C$ , and 2) few tuples from  $R$

Towards Responsive DBMS, ICDE 2022 tutorial: <https://northeastern-datalab.github.io/responsive-dbms-tutorial>

## Semi-joins can also help if data are local

$P_3(x, y, z, x) : -R(x, y), S(y, z), T(z, w)$



select \*  
into record1  
from R natural join S natural join T;

$n=1,000$ :  $t_{Q1}=1.4 \text{ sec}$

$n=2,000$ :  $t_{Q1}=6.1 \text{ sec } O(n^2) \otimes$

SQL example 601 available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>  
Towards Responsive DBMS, ICDE 2022 tutorial: <https://northeastern-datalab.github.io/responsive-dbms-tutorial>

```
-- Query 2
With S2 as
(SELECT *
FROM S
WHERE S.B in
(SELECT R.B
FROM R)),
S3 as
(SELECT *
FROM S2
WHERE S2.C in
(SELECT T.C
FROM T))
select a, b, c, d
into record2
from R natural join S3
```

$t_{Q2}=5 \text{ msec}$

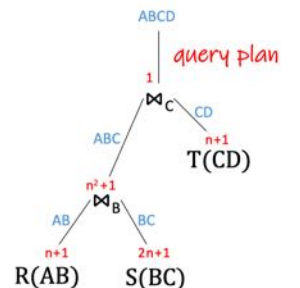
$t_{Q2}=8 \text{ msec}$

## The more general idea: "Sideways information passing"

### Sideways information passing:

- "sending information from one subexpression not simply to its parent expression, but also to some other correlated portion of the query computation, in order to prune irrelevant results" [Ives, Taylor 08]
- includes techniques like **two-way semijoins** [Bernstein, Goodman 81] and **magic sets** [Beeri, Ramakrishnan 91]

$$Q = (R \bowtie_B S) \bowtie_C T$$



[Bernstein, Goodman 81]. "Using Semi-Joins to Solve Relational Queries", JACM 1981. <https://doi.org/10.1145/322234.322238>  
[Beeri, Ramakrishnan 91]. "On the power of magic", Journal of Logic Programming, 1991. [https://doi.org/10.1016/0743-1066\(91\)90038-Q](https://doi.org/10.1016/0743-1066(91)90038-Q)  
Definition from: [Ives, Taylor 08]. "Sideways Information Passing for Push-Style Query Processing", ICDE 2008. <https://doi.org/10.1109/ICDE.2008.4497486>  
Towards Responsive DBMS, ICDE 2022 tutorial: <https://northeastern-datalab.github.io/responsive-dbms-tutorial>

See "Part 3: Acyclic queries & Enumeration": <https://northeastern-datalab.github.io/responsive-dbms-tutorial/slides/Responsive-DBMS-tutorial-part-3-AcyclicQueries-Enumeration.pdf>, [https://www.youtube.com/watch?list=PL\\_72ERKGF6DTInW\\_P3a9zTYPNSLwbqOAx&v=toi7ysuyRkw](https://www.youtube.com/watch?list=PL_72ERKGF6DTInW_P3a9zTYPNSLwbqOAx&v=toi7ysuyRkw) from ICDE'22 tutorial "Toward Responsive DBMS: Optimal Join Algorithms, Enumeration, Factorization, Ranking, and Dynamic Programming" by Tziavelis et al. <https://doi.org/10.1109/ICDE53745.2022.00299>, <https://northeastern-datalab.github.io/responsive-dbms-tutorial/> Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Join Summary

- **Theta-join:**  $R \bowtie_{\theta} S = \sigma_{\theta} (R \times S)$ 
  - Join of R and S with a join condition  $\theta$
  - Cross-product followed by selection  $\theta$
  - No projection
- **Equijoin:**  $R \bowtie_{\theta} S = \sigma_{\theta} (R \times S)$ 
  - Join condition  $\theta$  consists only of equalities
  - No projection
- **Natural join:**  $R \bowtie S = \pi_A (\sigma_{\theta} (R \times S))$ 
  - Equality on **all** fields with same name in R and in S
  - Projection  $\pi_A$  drops all redundant attributes

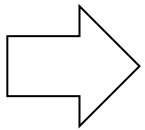
# Example: Converting SFW Query to RA



```
Student(sid,name,gpa)
People(ssn,name,address)
```

```
SELECT DISTINCT gpa, address
FROM Student S, People P
WHERE S.name = P.name
AND gpa > 3.5
```

*How do we represent this query in RA?*



# Example: Converting SFW Query to RA



```
Student(sid,name,gpa)
People(ssn,name,address)
```

```
SELECT DISTINCT gpa, address
FROM Student S, People P
WHERE S.name = P.name
AND gpa > 3.5
```

How do we represent this query in RA?

$$\begin{aligned} \Rightarrow & \Pi_{\text{gpa,address}}(\sigma_{\text{gpa}>3.5}(S \bowtie P)) \\ & \Pi_{\text{gpa,address}}(\sigma_{\text{gpa}>3.5 \wedge \text{S.name}=\text{P.name}}(S \times P)) \\ & \Pi_{\text{gpa,address}}(\sigma_{\text{gpa}>3.5 \wedge \text{name}=\text{name}_2}(S \times \rho_{\text{name} \rightarrow \text{name}_2} P)) \end{aligned}$$

# Some Examples

```
Supplier(sno,sname,scity,sstate)  
Part(pno,pname,psize,pcolor)  
Supply(sno,pno,qty,price)
```



Find names of suppliers of parts with size greater than 10



Find names of suppliers of red parts or parts with size greater than 10



# Some Examples

```
Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,qty,price)
```



Find names of suppliers of parts with size greater than 10

$$\pi_{\text{sname}}(\sigma_{\text{psize}>10}(\text{Supplier} \bowtie \text{Supply} \bowtie \text{Part}))$$
$$\pi_{\text{sname}}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{\text{psize}>10}(\text{Part})))$$

Find names of suppliers of red parts or parts with size greater than 10



# Some Examples

Supplier(sno,sname,scity,sstate)  
Part(pno,pname,psize,pcolor)  
Supply(sno,pno,qty,price)



Find names of suppliers of parts with size greater than 10

$\pi_{\text{sname}}(\sigma_{\text{psize}>10}(\text{Supplier} \bowtie \text{Supply} \bowtie \text{Part}))$

$\pi_{\text{sname}}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{\text{psize}>10}(\text{Part})))$

*Representation  
of RA as tree?*



Find names of suppliers of red parts or parts with size greater than 10

$\pi_{\text{sname}}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{\text{psize}>10}(\text{Part}) \cup \sigma_{\text{pcolor}='red'}(\text{Part})))$

$\pi_{\text{sname}}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{\text{psize}>10} \vee \text{pcolor}='red'(\text{Part})))$

# Some Examples



Supplier(sno,sname,scity,sstate)  
 Part(pno,pname,psize,pcolor)  
 Supply(sno,pno,qty,price)

Find names of suppliers of parts with size greater than 10

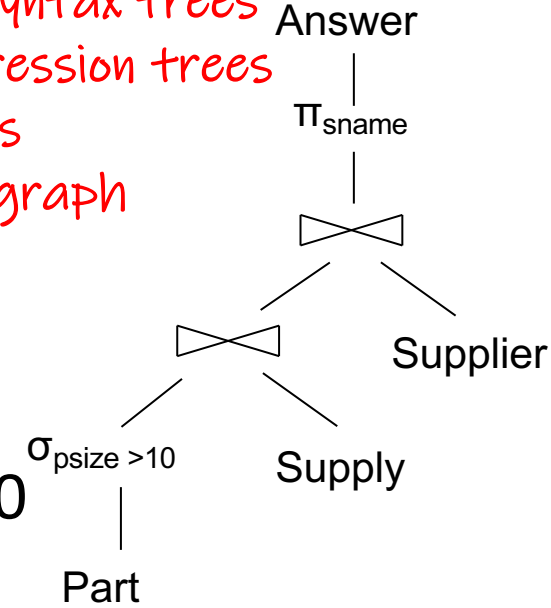
$\pi_{sname}(\sigma_{psize>10}(\text{Supplier} \bowtie (\text{Supply} \bowtie \text{Part})))$

$\pi_{sname}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{psize>10}(\text{Part})))$

*Representation of RA as tree?*

Usually unary or binary. Think of:

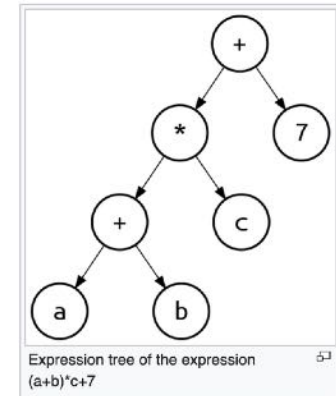
- abstract syntax trees
- binary expression trees
- parse trees
- data flow graph



Find names of suppliers of red parts or parts with size greater than 10

$\pi_{sname}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{psize>10}(\text{Part}) \cup \sigma_{pcolor='red'}(\text{Part})))$

$\pi_{sname}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{psize>10} \vee pcolor='red'(\text{Part})))$

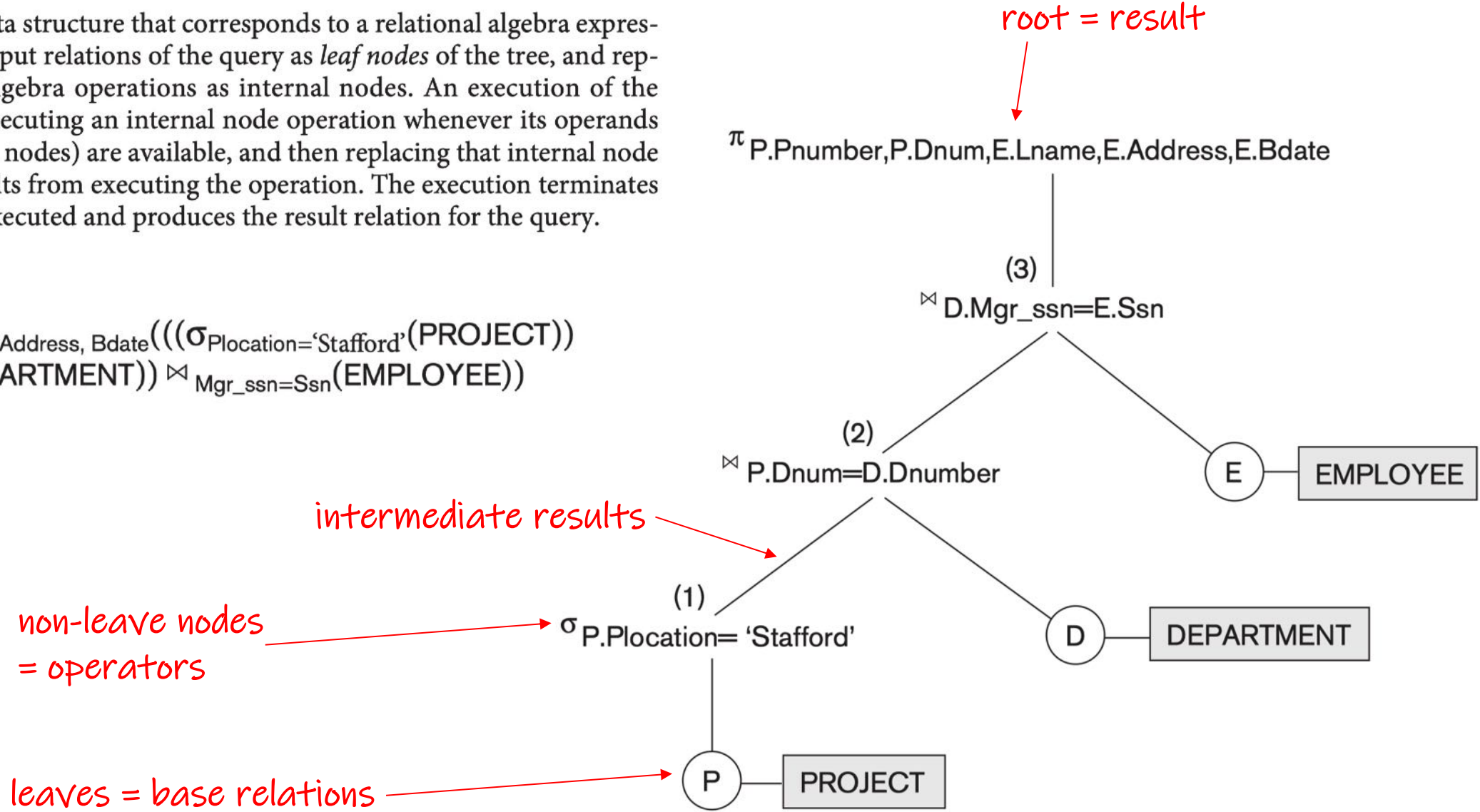




# Query (Evaluation / Execution) Tree, Data flow graph

A **query tree** is a tree data structure that corresponds to a relational algebra expression. It represents the input relations of the query as *leaf nodes* of the tree, and represents the relational algebra operations as internal nodes. An execution of the query tree consists of executing an internal node operation whenever its operands (represented by its child nodes) are available, and then replacing that internal node by the relation that results from executing the operation. The execution terminates when the root node is executed and produces the result relation for the query.

$\pi_{Pnumber, Dnum, Lname, Address, Bdate}(((\sigma_{Plocation='Stafford'}(PROJECT)))$   
 $\bowtie_{Dnum=Dnumber}(DEPARTMENT)) \bowtie_{Mgr\_ssn=Ssn}(EMPLOYEE))$



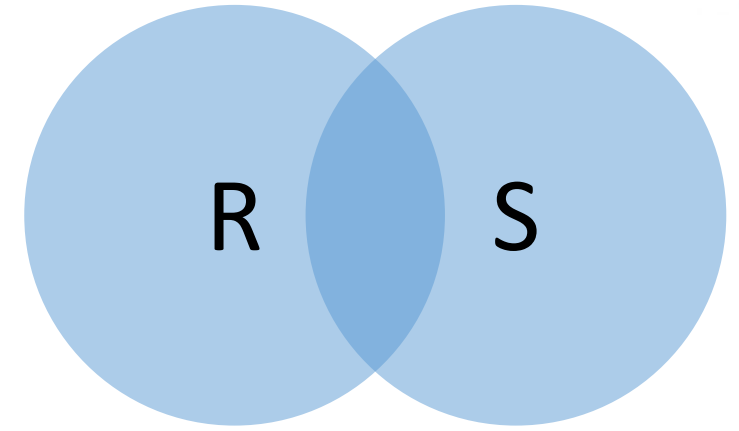
# Relational Algebra (RA) operators

- Five basic operators:
  1. Selection:  $\sigma$  ("sigma")
  2. Projection:  $\Pi$
  3. Cartesian Product:  $\times$
  4. Union:  $\cup$
  5. Difference:  $-$
- Auxiliary (or special) operator
  6. Renaming:  $\rho$  ("rho")
- Derived (or implied) operators
  7. Joins  $\bowtie$  (natural, theta join, equi-join, [semi-join: moved to T3-U1])
  8. Intersection / complement
  9. Division

## 8. What about Intersection $\cap$ ?

- As derived operator using **union** and **minus**

?



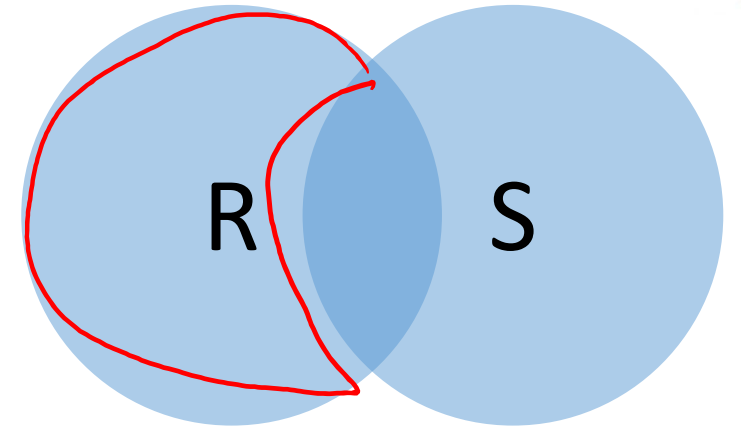
# 8. What about Intersection $\cap$ ?



- As derived operator using union and minus

?

$(R-S)$



$$S^c = D - S$$

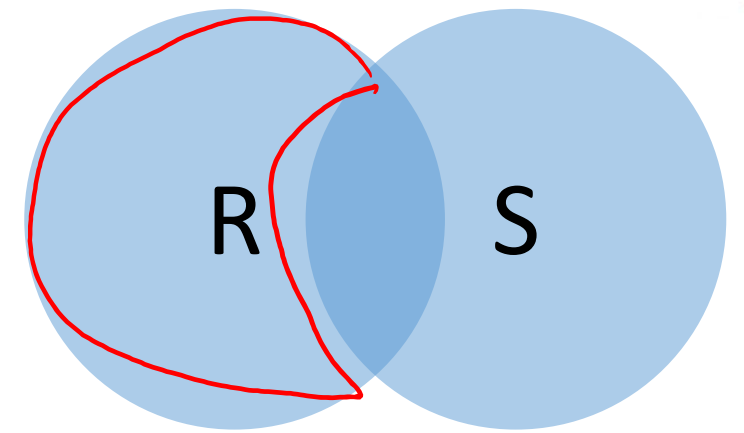
↓  
R

# 8. What about Intersection $\cap$ ?



- As derived operator using union and minus

?  $(R \cup S) - (R - S) - (S - R)$



$$S^c = D - S$$

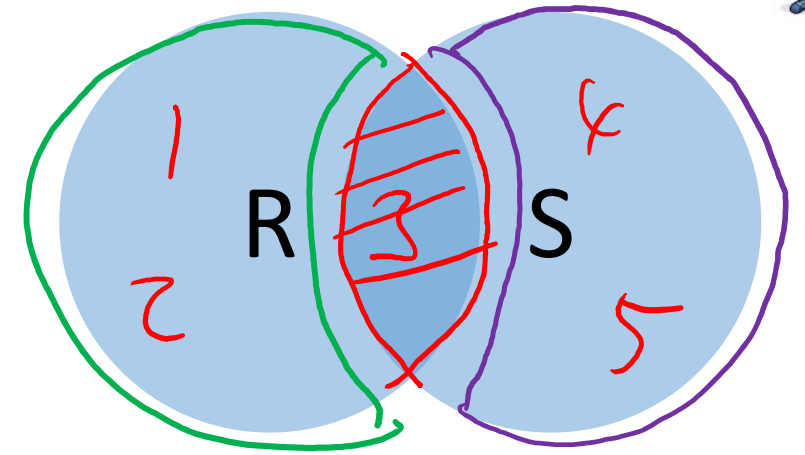
↓  
R

## 8. What about Intersection $\cap$ ?



- As derived operator using **union** and **minus**

$$(R \cup S) - ((R - S) \cup (S - R))$$



$$\{1, 2, 3\} \cap \{3, 4, 5\} = 3$$

## 8. What about Intersection $\cap$ ?



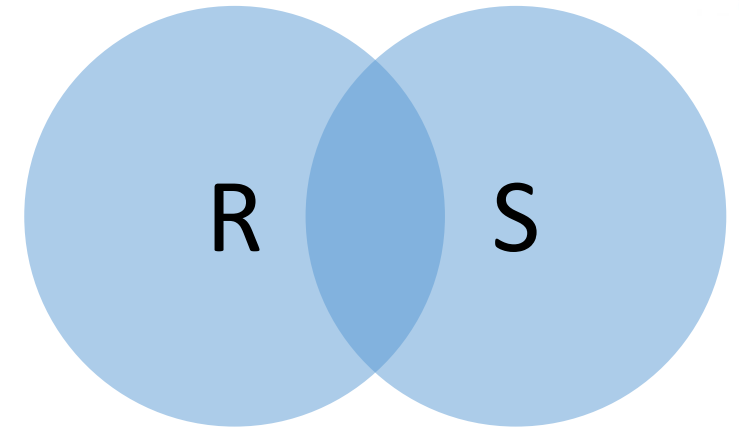
- As derived operator using **union** and **minus**

$$R \cap S = ((R \cup S) - (R - S)) - (S - R)$$

$$R \cap S = (R \cup S) - ((R - S) \cup (S - R))$$

- Derived operator using **minus** only!

?



## 8. What about Intersection $\cap$ ?



- As derived operator using **union** and **minus**

$$R \cap S = ((R \cup S) - (R - S)) - (S - R)$$

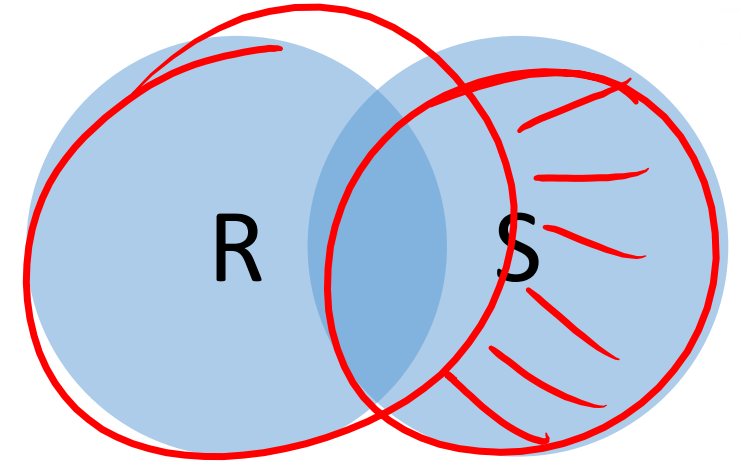
$$R \cap S = (R \cup S) - ((R - S) \cup (S - R))$$

- Derived operator using **minus** only!

$$R \cap S = S - (S - R)$$

- Derived using join

?



$$S - (S - R)$$



## 8. What about Intersection $\cap$ ?



- As derived operator using **union** and **minus**

$$R \cap S = ((R \cup S) - (R - S)) - (S - R)$$

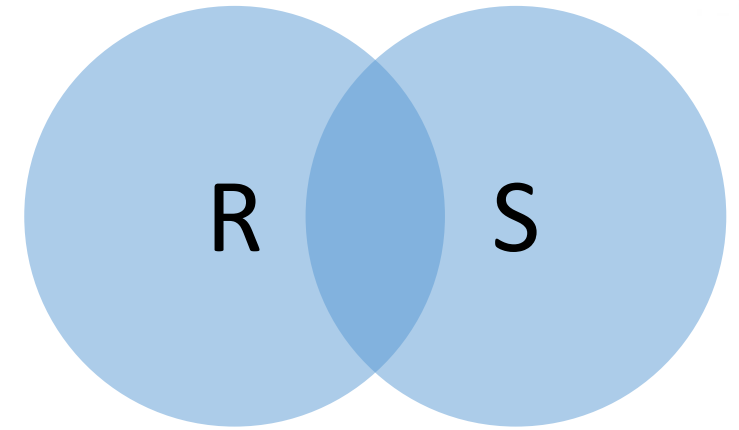
$$R \cap S = (R \cup S) - ((R - S) \cup (S - R))$$

- Derived operator using **minus** only!

$$R \cap S = S - (S - R)$$

- Derived using join

$$R \cap S = R \bowtie S$$



Legal input: schemas need to be union compatible (same schema). E.g. not:

$R(A,B,C)$

$S(A,B)$

If  $R$  and  $S$  have the same schema, then  $R \bowtie S$  and  $R \times S$  equal to  $R \cap S$

# Relational Algebra (RA) operators

- Five basic operators:
  1. Selection:  $\sigma$  ("sigma")
  2. Projection:  $\Pi$
  3. Cartesian Product:  $\times$
  4. Union:  $\cup$
  5. Difference:  $-$
- Auxiliary (or special) operator
  6. Renaming:  $\rho$  ("rho")
- Derived (or implied) operators
  7. Joins  $\bowtie$  (natural, theta join, equi-join, [semi-join: moved to T3-U1])
  8. Intersection / complement
  9. **Division**

## 9. Division ( $R \div S$ )

- Consider two relations  $R(X, Y)$  and  $S(Y)$
- Then  $R \div S$  is ...

*X, Y are sets of attributes  
Legal input:  $\text{att}(R) \supseteq \text{att}(S)$*

*What could be a meaningful definition of division*

*?*

*Compare to Integer division:  $7/2 = 3$*

*3 is the biggest integer that multiplied with 2 is smaller or equal to 7*

## 9. Division ( $R \div S$ )

- Consider two relations  $R(X, Y)$  and  $S(Y)$
- Then  $R \div S$  is ...
  - ... the largest relation  $T(X)$  s.t.  $S \times T \subseteq R$

*X, Y are sets of attributes  
Legal input:  $\text{att}(R) \supseteq \text{att}(S)$*

*(safety:  $T \subseteq \pi_X R$ )*

# 9. Division ( $R \div S$ )

- Consider two relations  $R(X, Y)$  and  $S(Y)$

*X, Y are sets of attributes*

- Then  $R \div S$  is ...

*Legal input:  $\text{att}(R) \supseteq \text{att}(S)$*

- ... the largest relation  $T(X)$  s.t.  $S \times T \subseteq R$ , or

*(safety:  $T \subseteq \pi_X R$ )*

- ... the relation  $T(X)$  that contains the  $X$ 's that occur with all  $Y$ 's in  $S$ , or

- ...  $\{t(X) \mid \forall s(Y) \in S. [\exists r(X, Y) \in R]\}$  (+ safety)

R Dividend

X	Y
Alice	1
Alice	2
Bob	1
Bob	2
Bob	3

S Divisor

Y
1
2
3

T

?

# 9. Division ( $R \div S$ )

- Consider two relations  $R(X, Y)$  and  $S(Y)$

*X, Y are sets of attributes*

- Then  $R \div S$  is ...

*Legal input:  $\text{att}(R) \supseteq \text{att}(S)$*

- ... the largest relation  $T(X)$  s.t.  $S \times T \subseteq R$ , or

*(safety:  $T \subseteq \pi_X R$ )*

- ... the relation  $T(X)$  that contains the  $X$ 's that occur with all  $Y$ 's in  $S$ , or

- ...  $\{t(X) \mid \forall s(Y) \in S. [\exists r(X, Y) \in R]\}$  (+ safety)

R Dividend

X	Y
Alice	1
Alice	2
Bob	1
Bob	2
Bob	3

S Divisor

Y
1
2
3

T

X
Bob

# Questions



sid	student	course
1	Alice	AI
1	Alice	DB
2	Bob	DB
2	Bob	ML
3	Charly	AI
3	Charly	DB
3	Charly	ML

$\div$

course
ML

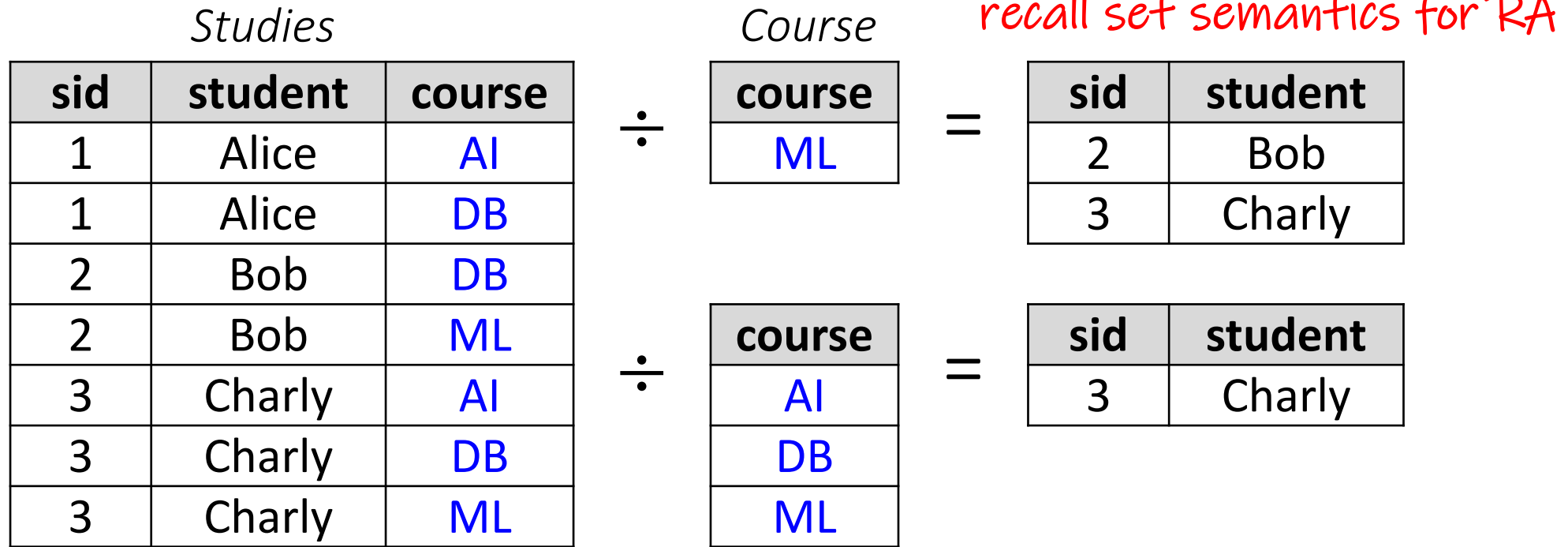
= ?

$\div$

course
AI
DB
ML

= ?

# Questions



*Assume R,S have disjoint attribute sets (possibly by renaming)*

$$(R \times S) \div S = ?$$

$$(R \times S) \div R = ?$$



# Questions



*Studies*

sid	student	course
1	Alice	AI
1	Alice	DB
2	Bob	DB
2	Bob	ML
3	Charly	AI
3	Charly	DB
3	Charly	ML

÷

*Course*

course
ML

=

sid	student
2	Bob
3	Charly

÷

*Course*

course
AI
DB
ML

=

sid	student
3	Charly

recall set semantics for RA

Assume R,S have disjoint attribute sets (possibly by renaming)

$$(R \times S) \div S = R$$

$$(R \times S) \div R = S$$

Q: If R has 1000 tuples and S has 100 tuples, how many tuples can be in  $R \div S$ ?

?

Q: If R has 1000 tuples and S has 1001 tuples, how many tuples can be in  $R \div S$ ?

?

# Questions



*Studies*

*Course*

recall set semantics for RA

sid	student	course
1	Alice	AI
1	Alice	DB
2	Bob	DB
2	Bob	ML
3	Charly	AI
3	Charly	DB
3	Charly	ML

÷

course
ML

=

sid	student
2	Bob
3	Charly

÷

course
AI
DB
ML

=

sid	student
3	Charly

1000

3

10

Assume R,S have disjoint attribute sets (possibly by renaming)

$$(R \times S) \div S = R$$

$$(R \times S) \div R = S$$

Q: If R has 1000 tuples and S has 100 tuples, how many tuples can be in  $R \div S$ ?

0-10

Q: If R has 1000 tuples and S has 1001 tuples, how many tuples can be in  $R \div S$ ?

0

# Questions



*Studies*

sid	student	course
1	Alice	AI
1	Alice	DB
2	Bob	DB
2	Bob	ML
3	Charly	AI
3	Charly	DB
3	Charly	ML

*CourseType*

course	type
AI	elective
DB	core
ML	core

Who took all **core** courses in RA with relational division?



# Questions

*Studies*

sid	student	course
1	Alice	AI
1	Alice	DB
2	Bob	DB
2	Bob	ML
3	Charly	AI
3	Charly	DB
3	Charly	ML

*CourseType*

course	type
AI	elective
DB	core
ML	core

Who took all **core** courses in RA with relational division?

$$\text{Studies} \div \pi_{\text{course}} \left( \sigma_{\text{type}='core'} \text{CourseType} \right)$$

# How to write $R \div S$ in Primitive RA? ( $\times, -, \pi$ )



$R(X, Y) \div S(Y)$

$R \div S = Q$

X	Y
a	0
a	1
a	2
b	1

Y
1
2

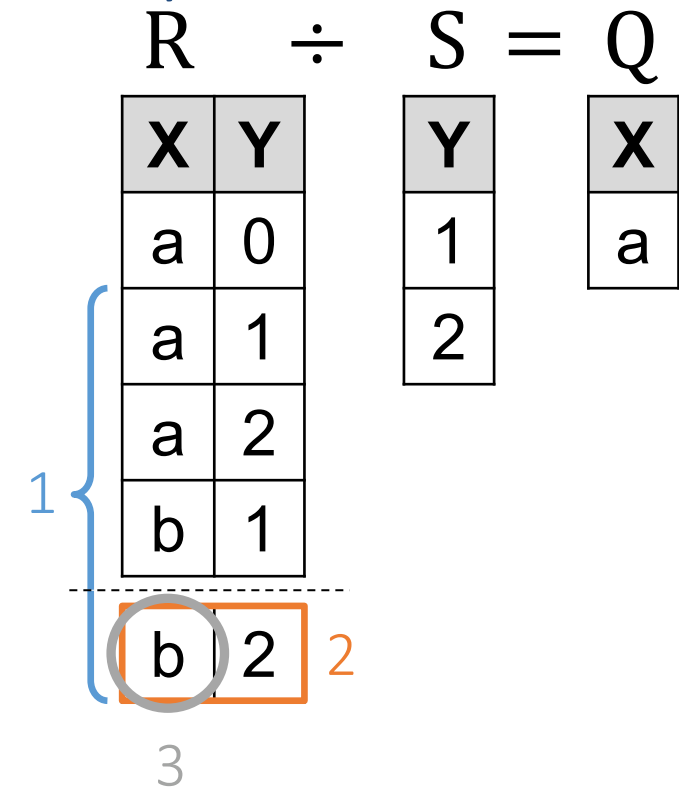
?

# How to write $R \div S$ in Primitive RA? ( $\times, -, \pi$ )



$R(X, Y) \div S(Y)$

?



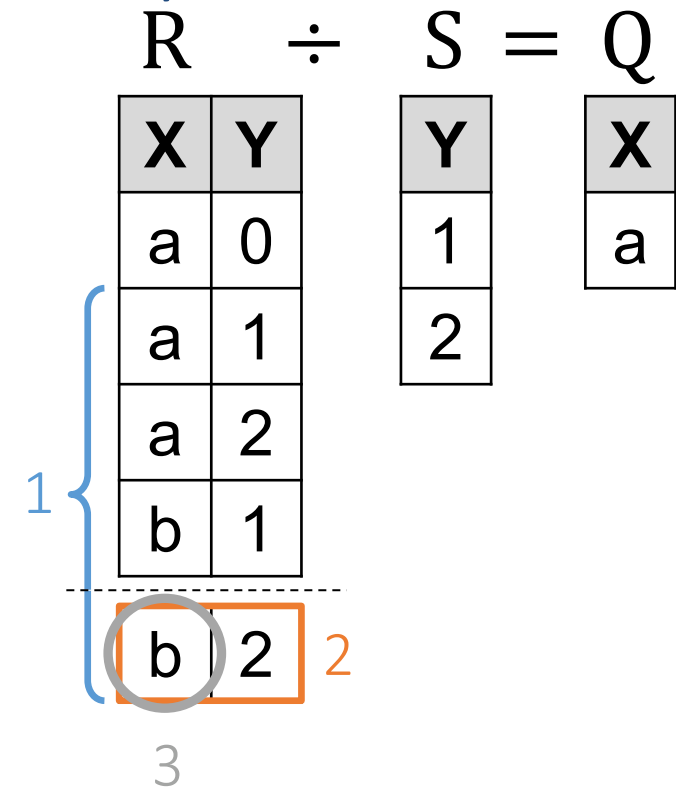
4:  $\{a\} = \{a, b\} - \{b\}$

# How to write $R \div S$ in Primitive RA? ( $\times, -, \pi$ )

$$R(X, Y) \div S(Y)$$

$$\pi_X R \times S$$

Each X of R w/ each Y of S



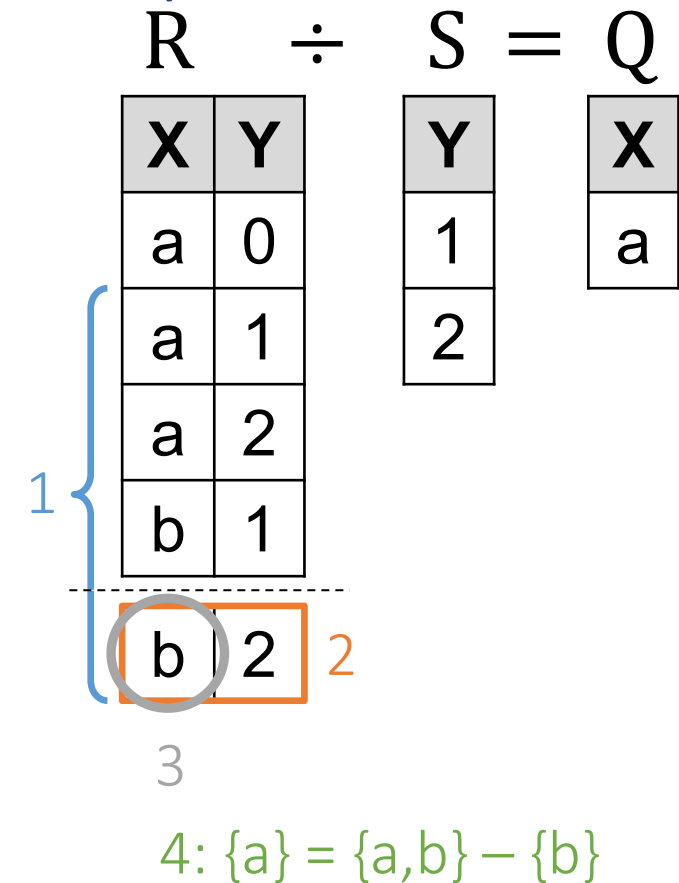
# How to write $R \div S$ in Primitive RA? ( $\times, -, \pi$ )

$$R(X, Y) \div S(Y)$$

$$(\pi_X R \times S) - R$$

Each X of R w/ each Y of S

(X,Y) s.t. X in R, Y in S, but (X,Y) not in R





# How to write $R \div S$ in Primitive RA? ( $\times, -, \pi$ )

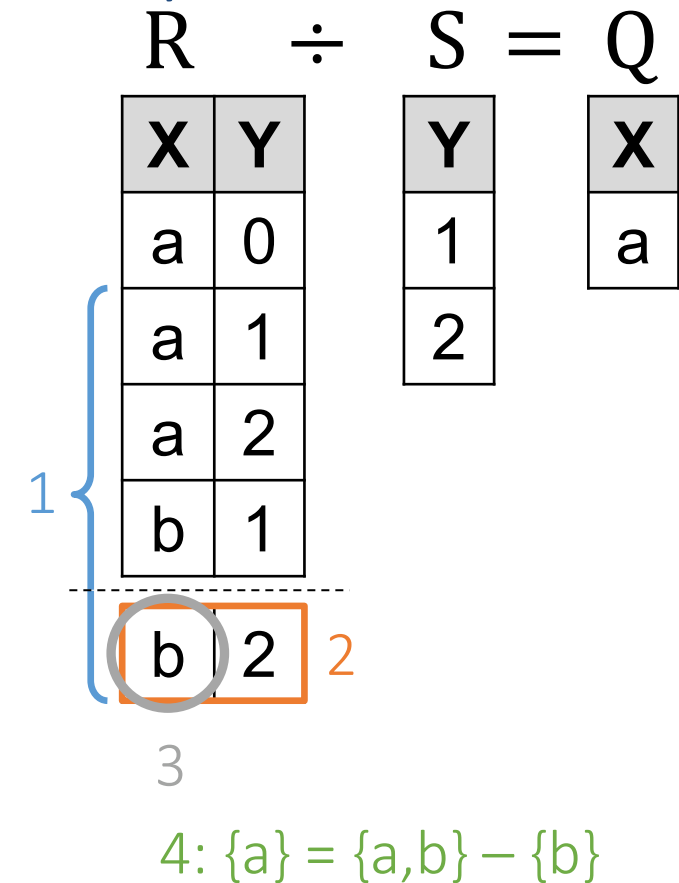
$$R(X, Y) \div S(Y)$$

$$\pi_X \left( \left( \pi_X R \times S \right) - R \right)$$

Each X of R w/ each Y of S

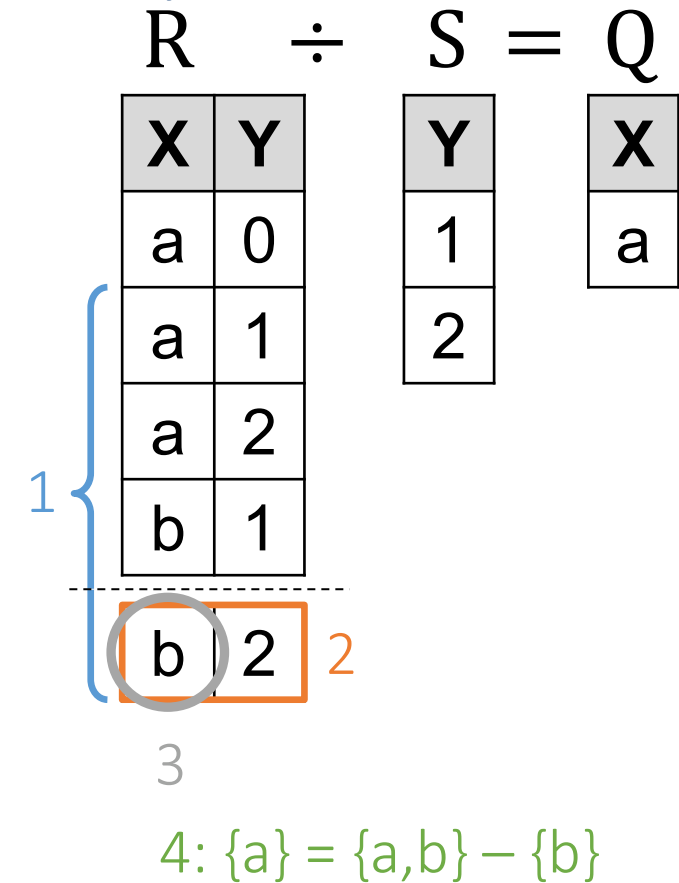
(X,Y) s.t. X in R, Y in S, but (X,Y) not in R

Xs in R where for some Y in S, (X,Y) is not in R



# How to write $R \div S$ in Primitive RA? ( $\times, -, \pi$ )

$$R(X, Y) \div S(Y)$$



$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

Each X of R w/ each Y of S

(X,Y) s.t. X in R, Y in S, but (X,Y) not in R

Xs in R where for some Y in S, (X,Y) is not in R

$R \div S$

What if  $S = \emptyset$ ?

$R(X, Y) \div S(Y)$

$R \div S = Q$

X	Y
a	0
a	1
a	2
b	1

$S$   $Y$   $Q$  ?

What if  $S = \emptyset$ ?

$$R(X, Y) \div S(Y)$$

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

	R	÷	S	=	Q														
	<table border="1" style="border-collapse: collapse;"><tr><th>X</th><th>Y</th></tr><tr><td>a</td><td>0</td></tr><tr><td>a</td><td>1</td></tr><tr><td>a</td><td>2</td></tr><tr><td>b</td><td>1</td></tr></table>	X	Y	a	0	a	1	a	2	b	1		<table border="1" style="border-collapse: collapse;"><tr><th>Y</th></tr></table>	Y		<table border="1" style="border-collapse: collapse;"><tr><th>X</th></tr><tr><td>a</td></tr><tr><td>b</td></tr></table>	X	a	b
X	Y																		
a	0																		
a	1																		
a	2																		
b	1																		
Y																			
X																			
a																			
b																			

Recall:  $\{t(X) \mid \forall s(Y) \in S. [\exists r(X, Y) \in R]\}$  (+ safety)

Now you see why we needed the safety condition " $T \subseteq \pi_X R$ " when defining " $R \div S$  as the largest relation  $T(X)$  s.t.  $S \times T \subseteq R$ "

# R ÷ S in Primitive RA vs. RC



$$R(X, Y) \div S(Y)$$

In RA:

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

In RC: ?

R ÷ S = Q

X	Y
a	0
a	1
a	2
b	1
b	2

Y
1
2

X
a

A blue bracket groups the first four rows of R. A dashed horizontal line is drawn between the fourth and fifth rows of R. The fifth row (b, 2) is circled in grey and has an orange box around it.

# R ÷ S in Primitive RA vs. RC



$$R(X, Y) \div S(Y)$$

In RA:

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

In DRC:

$$\{ X \mid \exists Z. [R(X, Z)] \wedge \quad ? \quad \}$$

*X is "guarded": safe and thus domain independent*

$$R \div S = Q$$

X	Y
a	0
a	1
a	2
b	1
b	2

Y
1
2

X
a

# R ÷ S in Primitive RA vs. RC



$$R(X, Y) \div S(Y)$$

In RA:

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

In DRC:

what if  $S(Y) = \emptyset$ ?

?

$$\{ X \mid \exists Z. [R(X, Z)] \wedge \forall Y. [S(Y) \rightarrow R(X, Y)] \}$$

R ÷ S = Q

X	Y
a	0
a	1
a	2
b	1
b	2

Y
1
2

X
a

A blue bracket groups the first four rows of R. A red circle highlights the cell (b, 2) in R, which is also highlighted with an orange box. A dashed horizontal line is drawn below the row (b, 1).

# R ÷ S in Primitive RA vs. RC



$$R(X, Y) \div S(Y)$$

In RA:

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

In DRC:

$$\{ X \mid \exists Z. [R(X, Z)] \wedge \forall Y. [S(Y) \rightarrow R(X, Y)] \}$$

R ÷ S = Q

X	Y
a	0
a	1
a	2
b	1
b	2

Y
1
2

X
a

A blue bracket groups the first four rows of R. A dashed line is below the fifth row. The cell (b, 2) in R is circled in grey and boxed in orange.

? without universal quantification



# R ÷ S in Primitive RA vs. RC



$$R(X, Y) \div S(Y)$$

In RA:

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

In DRC:

$$\left\{ X \mid \exists Z. [R(X, Z)] \wedge \forall Y. [S(Y) \rightarrow R(X, Y)] \right\}$$
$$\left\{ X \mid \exists Z. [R(X, Z)] \wedge \nexists Y. [S(Y) \wedge \neg R(X, Y)] \right\}$$

In TRC:

?

R ÷ S = Q

X	Y
a	0
a	1
a	2
b	1
b	2

Y
1
2

X
a

A blue bracket groups the first four rows of R. A dashed line is below the fifth row. The cell (b, 2) in R is circled in grey and boxed in orange.

# R ÷ S in Primitive RA vs. RC

$$R(X, Y) \div S(Y)$$

In RA:

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

In DRC:

$$\left\{ X \mid \exists Z. [R(X, Z)] \wedge \forall Y. [S(Y) \rightarrow R(X, Y)] \right\}$$

$$\left\{ X \mid \exists Z. [R(X, Z)] \wedge \nexists Y. [S(Y) \wedge \neg R(X, Y)] \right\}$$

In TRC:

$$\left\{ r.A \mid \exists r \in R. \left[ \nexists s \in S. \left[ \text{?} \right] \right] \right\}$$

R ÷ S = Q

X	Y
a	0
a	1
a	2
b	1
b	2

Y
1
2

X
a

# R ÷ S in Primitive RA vs. RC

$$R(X, Y) \div S(Y)$$

$R$		$\div$	$S = Q$	
X	Y		Y	X
a	0		1	a
a	1		2	
a	2			
b	1			
b	2			

In RA:

$$\pi_X R - \pi_X \left( (\pi_X R \times S) - R \right)$$

In DRC:

$$\{ X \mid \exists Z. [R(X, Z)] \wedge \forall Y. [S(Y) \rightarrow R(X, Y)] \}$$

$$\{ X \mid \exists Z. [R(X, Z)] \wedge \nexists Y. [S(Y) \wedge \neg R(X, Y)] \}$$

? in SQL

In TRC:

$$\{ r.A \mid \exists r \in R. [ \nexists s \in S. [ \nexists r_2 \in R. [ r_2.Y = s.Y \wedge r_2.X = r.X ] ] ] \}$$

# R ÷ S in Primitive RA vs. RC

In SQL

```
SELECT DISTINCT R.A
FROM R
WHERE not exists (
  SELECT *
  FROM S
  WHERE not exists (
    SELECT *
    FROM R AS R2
    WHERE R2.B=S.B
    AND R2.A=R.A))
```

R ÷ S = Q

X	Y
a	0
a	1
a	2
b	1
b	2

Y
1
2

X
a

A blue bracket groups the first four rows of R. A dashed horizontal line is drawn below the row (b, 1). The row (b, 2) is circled in grey and has an orange box around it.

In TRC:

$$\{ r.A \mid \exists r \in R. [ \nexists s \in S. [ \nexists r_2 \in R. [ r_2.B = s.B \wedge r_2.A = r.A ] ] ] \}$$

# Parentheses Convention

- We have defined 3 unary operators and 3 binary operators
- It is acceptable to omit the parentheses from  $o(R)$  when  $o$  is unary
  - Then, unary operators take precedence over binary ones
- Example:

$$(\sigma_{\text{course}='DB'}(\text{Course})) \times (\rho_{\text{cid} \rightarrow \text{cid1}}(\text{Studies}))$$

becomes

$$\sigma_{\text{course}='DB'}\text{Course} \times \rho_{\text{cid} \rightarrow \text{cid1}}\text{Studies}$$

# Algebra and the connection to logic and queries

- Algebra
- Relational Algebra
  - Operators
  - Independence
  - Power of algebra: optimizations
- Equivalence RA and safe RC (Codd's theorem)
  - $RA \rightarrow RC$
  - $RC \rightarrow RA$

# 5 Primitive Operators

1. Projection ( $\pi$ )
2. Selection ( $\sigma$ )
3. Union ( $\cup$ )
4. Set Difference ( $-$ )
5. Cross Product ( $\times$ )

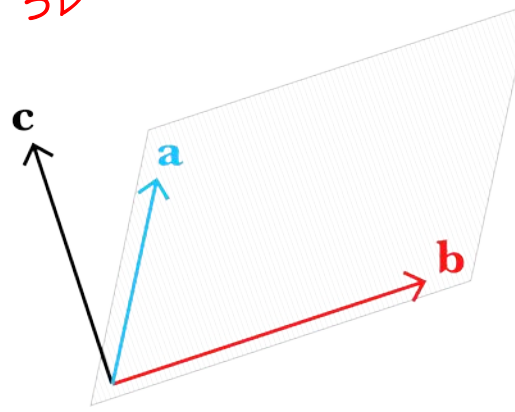
*Is this a well chosen set of primitives?*



# 5 Primitive Operators

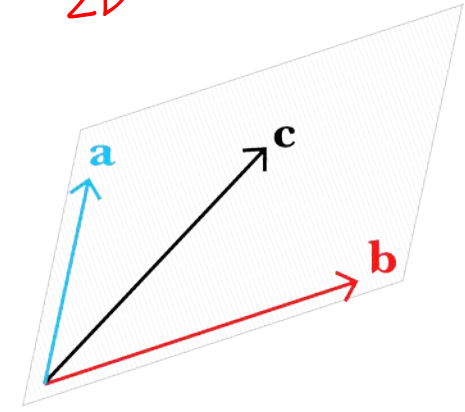
1. Projection ( $\pi$ )
2. Selection ( $\sigma$ )
3. Union ( $\cup$ )
4. Set Difference ( $-$ )
5. Cross Product ( $\times$ )

3D

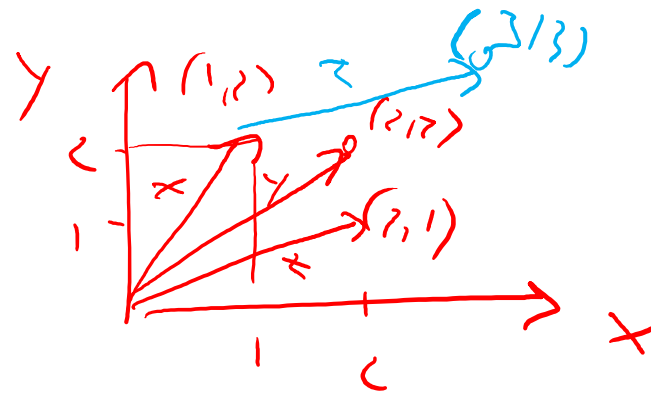


independent

2D



not independent



$$y = (x + z) \cdot \frac{1}{\sqrt{2}}$$
$$\frac{1}{\sqrt{2}}x + \frac{1}{\sqrt{2}}z = 0$$

Is this a well chosen set of primitives?

Could we drop an operator "without losing anything"?



# Independence among Primitives

- Let  $\circ$  be an RA operator, and let  $A$  be a set of RA operators
- We say that  $\circ$  is **independent** of  $A$  if  $\circ$  cannot be expressed in  $A$ ; that is, no expression in  $A$  is equivalent to  $\circ$

THEOREM: Each of the five primitives is independent of the other four

$\{\pi, \sigma, \times, \cup, -\}$

Proof:

- Separate argument for each of the 5 (For each operator, we need to discover a property that is possessed by that operator, but not by any RA expression that involves only the other 4 operations)
- Arguments follow a common pattern (union as example next slides)

# Recipe for Proving Independence of an operator $\circ$

1. Fix a schema  $S$  and an instance  $D$  over  $S$
2. Find some **property**  $P$  over relations
3. Prove: for every expression  $\varphi$  that does not use  $\circ$ , the relation  $\varphi(D)$  satisfies  $P$

*Such proofs are typically by induction on the size of the expression, since operators compose*

4. Find an expression  $\psi$  such that  $\psi$  uses  $\circ$  and  $\psi(D)$  violates  $P$

# Concrete Example: Proving Independence of Union $\cup$

1. Fix a schema  $S$  and an instance  $D$  over  $S$

$S: R(A), S(A)$        $D: \{R(0), S(1)\}$

$R$	$S$
<b>A</b>	<b>A</b>
0	1

2. Find some **property  $P$**  over relations

$\#tuples < 2$

3. Prove: for every expression  $\varphi$  that does not use  $\circ$ , the relation  $\varphi(D)$  satisfies  **$P$**

Induction base:  $R$  and  $S$  have  $\#tuples < 2$

Induction step: If  $\varphi_1(D)$  and  $\varphi_2(D)$  have  $\#tuples < 2$ , then so do:

$\sigma_c(\varphi_1(D)), \pi_A(\varphi_1(D)), \varphi_1(D) \times \varphi_2(D), \varphi_1(D) - \varphi_2(D), \rho_{A \rightarrow B}(\varphi_1(D))$

4. Find an expression  $\psi$  such that  $\psi$  uses  $\circ$  and  $\psi(D)$  violates  **$P$**

$\psi = R \cup S$

# Algebra and the connection to logic and queries

- Algebra
- Relational Algebra
  - Operators
  - Independence
  - Power of algebra: optimizations
- Equivalence RA and safe RC (Codd's theorem)
  - $RA \rightarrow RC$
  - $RC \rightarrow RA$

# Commutativity and distributivity of RA operators

- The basic commutators:

- Push projection through selection, join, union
- Push selection through projection, join, union
- Also: Joins can be re-ordered!

$$\pi_A(R \cup S) = \pi_A(R) \cup \pi_A(S)$$

$$\sigma_\theta(R \cup S) = \sigma_\theta(R) \cup \sigma_\theta(S)$$

$$(R \cup S) \times T = (R \times T) \cup (S \times T)$$

- Note that this is not an exhaustive set of operations

*What about sorting and joins?*

This simple set of tools allows us to greatly improve the execution time of queries by optimizing RA plans!

We next illustrate with an SFW (Select-From-Where) query

# An example: SQL to RA to Optimized RA

R(A,B) S(B,C) T(C,D)

```
SELECT R.A, T.D  
FROM   R, S, T  
WHERE  R.B = S.B  
       and S.C = T.C  
       and R.A < 10;
```

in RA ↓

?

# An example: SQL to RA to Optimized RA

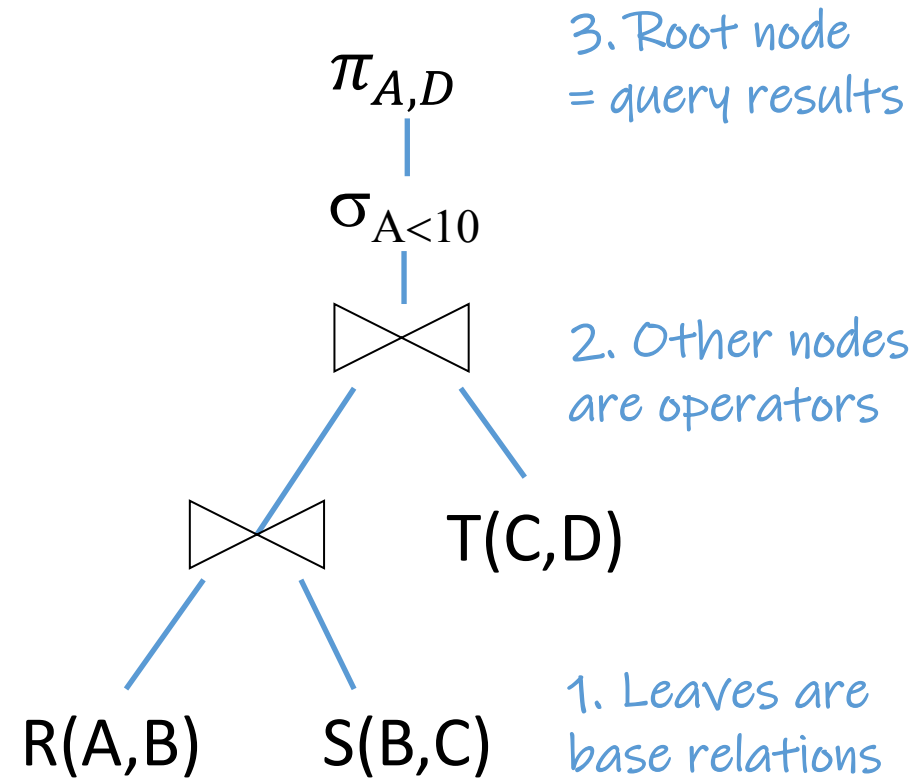
Heuristic: have selection and projection earlier to have fewer (or smaller) "intermediate" tuples

R(A,B) S(B,C) T(C,D)

```
SELECT R.A, T.D
FROM R, S, T
WHERE R.B = S.B
      and S.C = T.C
      and R.A < 10;
```

in RA ↓

$\pi_{A,D} \left( \sigma_{A < 10} \left( T \bowtie (R \bowtie S) \right) \right)$



Query tree / expression tree / computation tree / data flow graph

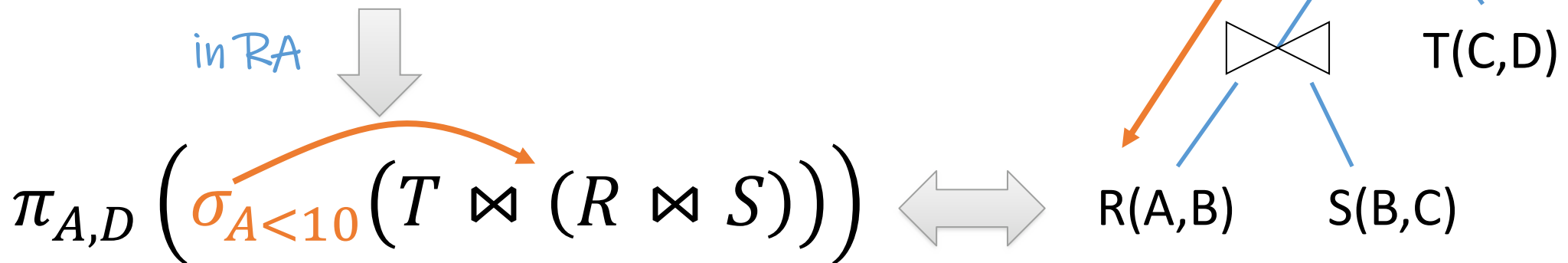
# An example: SQL to RA to Optimized RA

Heuristic: have selection and projection earlier to have fewer (or smaller) "intermediate" tuples

R(A,B) S(B,C) T(C,D)

```
SELECT R.A, T.D
FROM R, S, T
WHERE R.B = S.B
      and S.C = T.C
      and R.A < 10;
```

1. Push down selection on A



Pushing down may be suboptimal if selection condition is very expensive (e.g. running some image processing algorithm). Projection could be unnecessary effort (but more rarely).



# An example: SQL to RA to Optimized RA

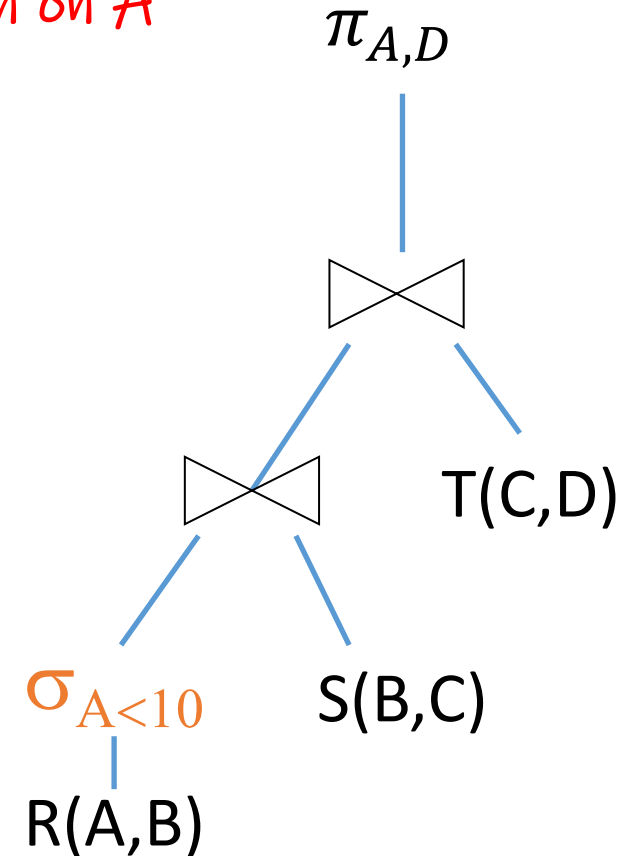
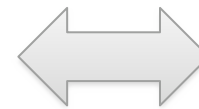
Heuristic: have selection and projection earlier to have fewer (or smaller) "intermediate" tuples

R(A,B) S(B,C) T(C,D)

```
SELECT R.A, T.D
FROM R, S, T
WHERE R.B = S.B
and S.C = T.C
and R.A < 10;
```

in RA

$\pi_{A,D} (T \bowtie (\sigma_{A < 10} R \bowtie S))$



# An example: SQL to RA to Optimized RA

Heuristic: have selection and projection earlier to have fewer (or smaller) "intermediate" tuples

R(A,B) S(B,C) T(C,D)

```
SELECT R.A, T.D
FROM   R, S, T
WHERE  R.B = S.B
      and S.C = T.C
      and R.A < 10;
```

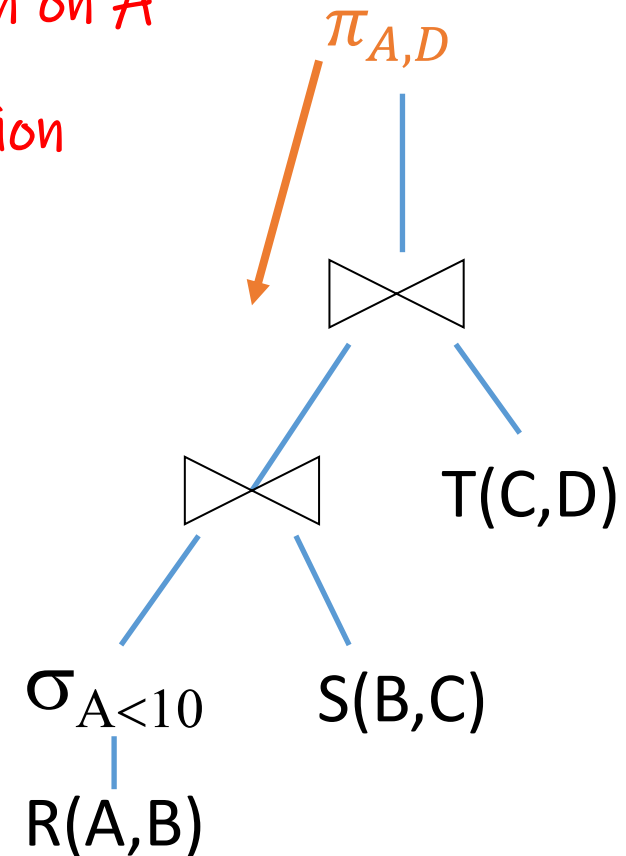
in RA

$\pi_{A,D} (T \bowtie (\sigma_{A < 10} R \bowtie S))$

$\pi_{-B,C}$

1. Push down selection on A

2. Push down projection



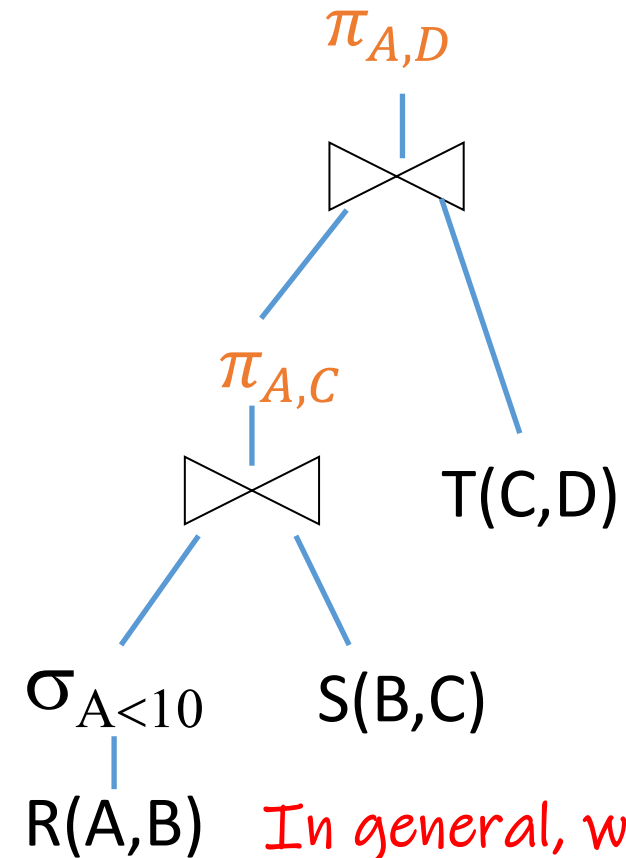
# An example: SQL to RA to Optimized RA

R(A,B) S(B,C) T(C,D)

```
SELECT R.A, T.D
FROM R, S, T
WHERE R.B = S.B
and S.C = T.C
and R.A < 10;
```

in RA ↓

$\pi_{A,D} \left( T \bowtie \pi_{A,C} (\sigma_{A < 10} R \bowtie S) \right)$



Variable Elimination!

$\pi_{-C}$

$\pi_{-B}$

$\pi_{-C}$

$\pi_{-B}$

We now eliminate B earlier

In general, when is an attribute not needed?

# Algebra and the connection to logic and queries

- Algebra
- Relational Algebra
  - Operators
  - Independence
  - Power of algebra: optimizations
- Equivalence RA and safe RC (Codd's theorem)
  - $RA \rightarrow RC$
  - $RC \rightarrow RA$



# "Clear" variables\*

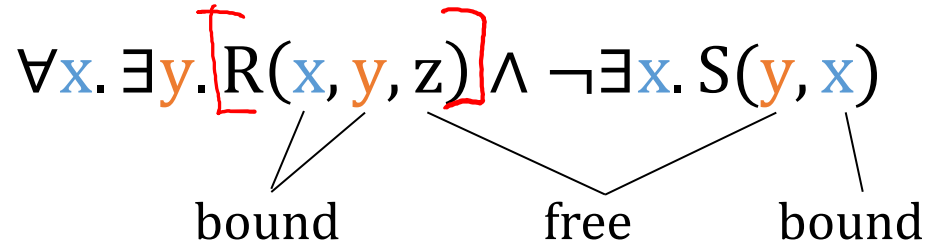
Formula with clear variables : each quantifier "has its own variables" & each variable has only free or only bound occurrences

$$\forall x. \exists y. R(x, y, z) \wedge \neg \exists x. S(y, x)$$

? *which variables are free or bound?*

# "Clear" variables\*

Formula with clear variables : each quantifier "has its own variables" & each variable has only free or only bound occurrences



notice operator precedence:  $\exists$  before  $\wedge$   
 $\forall x. \exists y. [R(x, y, z)] \wedge \neg \exists x. [S(y, x)]$

Not "clear": Two x's and y's are different variables.

? how to make it "clear"

# "Clear" variables\*

Formula with clear variables : each quantifier "has its own variables" & each variable has only free or only bound occurrences

$$\forall x. \exists y. R(x, y, z) \wedge \neg \exists x. S(y, x)$$

bound                      free                      bound

notice operator precedence:  $\exists$  before  $\wedge$   
 $\forall x. \exists y. [R(x, y, z)] \wedge \neg \exists x. [S(y, x)]$

Not "clear": Two x's and y's are different variables.

$$\forall x. \exists y. R(x, y, z) \wedge \neg \exists u. S(v, u)$$

now "clear"

$$\{(z, v) \mid \forall x. \exists y. R(x, y, z) \wedge \neg \exists u. S(v, u)\}$$

Now a query. But how to make it domain-independent



# "Clear" variables\*

Formula with clear variables : each quantifier "has its own variables" & each variable has only free or only bound occurrences

$$\forall x. \exists y. R(x, y, z) \wedge \neg \exists x. S(y, x)$$

bound                      free                      bound

notice operator precedence:  $\exists$  before  $\wedge$   
 $\forall x. \exists y. [R(x, y, z)] \wedge \neg \exists x. [S(y, x)]$

Not "clear": Two x's and y's are different variables.

$$\forall x. \exists y. R(x, y, z) \wedge \neg \exists u. S(v, u)$$

now "clear"

$$\{(z, v) \mid \forall x. \exists y. R(x, y, z) \wedge \neg \exists u. S(v, u)\}$$

$\forall x. [R(x, \_, \_) \rightarrow R(x, \_, z)]$                        $\forall u. [S(\_, v) \wedge \neg \exists x. S(x, v)]$

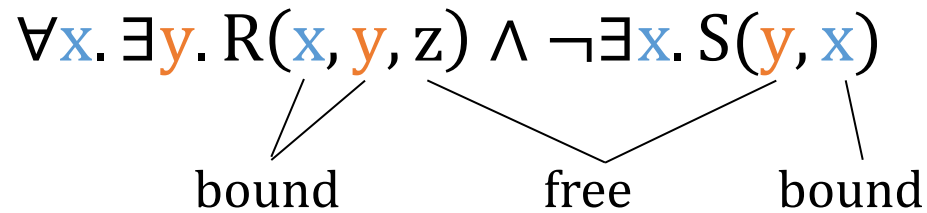
Now a query. But how to make it domain-independent

\* "Clear variable" is a non-standard term used in excellent slides by Marie Duzi: <http://www.cs.vsb.cz/duzi/>  
Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>



# "Clear" variables\*

Formula with clear variables : each quantifier "has its own variables" & each variable has only free or only bound occurrences



notice operator precedence:  $\exists$  before  $\wedge$   
 $\forall x. \exists y. [R(x, y, z)] \wedge \neg \exists x. [S(y, x)]$

Not "clear": Two x's and y's are different variables.

$$\forall x. \exists y. R(x, y, z) \wedge \neg \exists u. S(v, u)$$

now "clear"

$$\{ (z, v) \mid \exists s, t. R(s, t, z) \wedge \exists p. S(p, v) \wedge \forall x. \exists y. R(x, y, z) \wedge \neg \exists u. S(v, u) \}$$

$[(z, y, x). R(x, y, z) \rightarrow \exists y. R(x, y, z)]$

Now a query. But how to make it domain-independent

\* "Clear variable" is a non-standard term used in excellent slides by Marie Duzi: <http://www.cs.vsb.cz/duzi/>  
 Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Repeated variable names



In sentences with multiple quantifiers, distinct variables do not need to range over distinct objects! (cp. homomorphism vs. isomorphism)



*which of the following formulas imply each other?*

$$\forall x. \forall y. E(x, y)$$

$$\forall x. E(x, x)$$

$$\exists x. \exists y. E(x, y)$$

$$\exists x. E(x, x)$$

# Repeated variable names



In sentences with multiple quantifiers, distinct variables do not need to range over distinct objects! (cp. homomorphism vs. isomorphism)

Assume  $\text{DOM} = \{1, 2\}$ :

$$\forall x. \forall y. E(x, y)$$



$$\forall x. E(x, x)$$

E	
s	t
1	1
1	2
2	1
2	2

$$\exists x. \exists y. E(x, y)$$



$$\exists x. E(x, x)$$

E	
s	t
1	2

# Repeated variable names



In sentences with multiple quantifiers, distinct variables do not need to range over distinct objects! (cp. homomorphism vs. isomorphism)

Assume  $\text{DOM} = \emptyset$ :

$$\forall x. \forall y. E(x, y)$$

$\Rightarrow$

$$\forall x. E(x, x)$$

$\Downarrow$

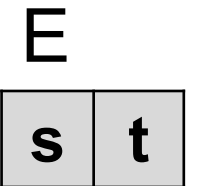
Only if domain is not empty!  $\text{Dom} \neq \emptyset$

$\Downarrow$

$$\exists x. \exists y. E(x, y)$$

$\Leftarrow$

$$\exists x. E(x, x)$$



# Example RC $\rightarrow$ RA

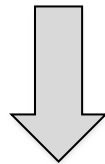
Person(id, name, country)  
Spouse(id1, id2)



In DRC:

$$\{ x \mid \exists z, w. \text{Person}(x, z, w) \wedge \forall y. [\neg \text{Spouse}(x, y)] \}$$

In RA:



# Example RC $\rightarrow$ RA

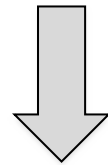
Person(id, name, country)  
Spouse(id1, id2)



In DRC:

$$\{ \mathbf{x} \mid \exists z, w. \text{Person}(\mathbf{x}, z, w) \wedge \forall y. [\neg \text{Spouse}(\mathbf{x}, y)] \}$$

$$\{ \mathbf{x} \mid \exists z, w. \text{Person}(\mathbf{x}, z, w) \wedge \neg \exists y. [\text{Spouse}(\mathbf{x}, y)] \}$$



In RA:

$$\pi_{\text{id}} \text{Person} - \pi_{\text{id1}} \text{Spouse}$$

$$\pi_{\text{id}} \text{Person} - \rho_{\text{id1} \rightarrow \text{id}} (\pi_{\text{id1}} \text{Spouse})$$

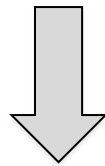
*Recall: named vs ordered perspective*

# Example $RA \rightarrow RC$

$$R(X, Y) \div S(Y)$$

In RA:

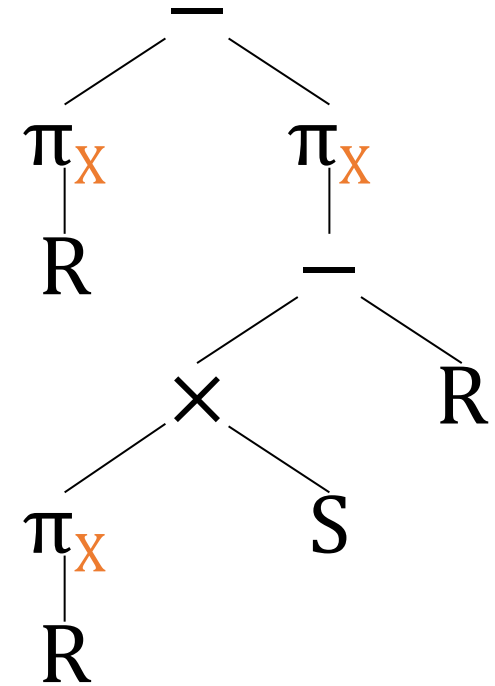
$$\pi_X R - \pi_X ((\pi_X R \times S) - R)$$



In DRC:

$$\{ X \mid \exists Z. [R(X, Z)] \wedge \forall Y. [S(Y) \rightarrow R(X, Y)] \}$$

$$\{ X \mid \exists Z. [R(X, Z)] \wedge \neg \exists Y. [S(Y) \wedge \neg R(X, Y)] \}$$



# Equivalence Between RA and Domain-Independent RC

CODD'S THEOREM:

RA and domain-independent RC  
have the same expressive power.

More formally, on every schema **S**:

1. For every RA expression **E**, there is a domain-independent RC query **Q** s.t.  $Q \equiv E$
2. For every domain-independent RC query **Q**, there is an RA expression **E** s.t.  $Q \equiv E$

The proof has two directions:

RA  $\rightarrow$  RC:

by induction on the size  
of the RA expression

RC  $\rightarrow$  RA:

more involved



# Algebra and the connection to logic and queries

- Algebra
- Relational Algebra
  - Operators
  - Independence
  - Power of algebra: optimizations
- Equivalence RA and safe RC (Codd's theorem)
  - $RA \rightarrow RC$
  - $RC \rightarrow RA$

# RA $\rightarrow$ DRC: Intuition

- Construction by induction
- Key technical detail: need to maintain a **mapping b/w attribute names and variables**

Intuition:  $\{x \mid \exists y. [R(x,y)] \wedge \exists y. [S(x,y)]\}$   
contrast with:  $\{x \mid \exists y. [R(x,y)] \wedge \exists z. [S(x,z)]\}$

$$Q(1) \leftarrow R(1, \underset{y=2}{2}), S(1, \underset{y=3}{3})$$

RA expression	DRC formula $\phi$ Here, $\phi_i$ is the formula constructed for expression $E_i$
$R$ (n columns)	$R(X_1, \dots, X_n)$
$E_1 \times E_2$	
$E_1 \cup E_2$	
$E_1 - E_2$	
$\pi_{A_1, \dots, A_k}(E_1)$	
$\sigma_c(E_1)$	

# RA $\rightarrow$ DRC: Intuition

- Construction by induction
- Key technical detail: need to maintain a **mapping b/w attribute names and variables**

Intuition:  $\{x \mid \exists y.[R(x,y)] \wedge \exists y.[S(x,y)]\}$   
contrast with:  $\{x \mid \exists y.[R(x,y)] \wedge \exists z.[S(x,z)]\}$

RA expression	DRC formula $\phi$
R (n columns)	$R(X_1, \dots, X_n)$
$E_1 \times E_2$	$\phi_1 \wedge \phi_2$ disjoint variables (rename)
$E_1 \cup E_2$	$\phi_1 \vee \phi_2$ use identical variables (rename) <span style="color: red;">UNION COMPATIBLE</span>
$E_1 - E_2$	$\phi_1 \wedge \neg \phi_2$ use identical variables (rename)
$\pi_{A_1, \dots, A_k}(E_1)$	
$\sigma_c(E_1)$	

# RA $\rightarrow$ DRC: Intuition

- Construction by induction

- Key technical detail: need to maintain a **mapping b/w attribute names and variables**

Intuition:  $\{x \mid \exists y.[R(x,y)] \wedge \exists y.[S(x,y)]\}$

contrast with:  $\{x \mid \exists y.[R(x,y)] \wedge \exists z.[S(x,z)]\}$

## RA expression

**DRC formula  $\phi$**  Here,  $\phi_i$  is the formula constructed for expression  $E_i$

R (n columns)

$R(X_1, \dots, X_n)$

$E_1 \times E_2$

$\phi_1 \wedge \phi_2$  **disjoint** variables (rename)

$E_1 \cup E_2$

$\phi_1 \vee \phi_2$  use **identical** variables (rename)

$E_1 - E_2$

$\phi_1 \wedge \neg \phi_2$  use **identical** variables (rename)

$\pi_{A_1, \dots, A_k}(E_1)$

$\exists X_1 \dots \exists X_m. \phi_1$  where  $X_1, \dots, X_m$  are the variables **not among**  $A_1, \dots, A_k$

$\sigma_c(E_1)$

$\phi_1 \wedge c$

*Correspondence more natural with  
project-away operator:  $\pi_{-A_1, \dots, A_m}(E_1)$*

# RA $\rightarrow$ DRC: Example $R \div S$

R(A,B) S(B)

RA	DRC	Mapping
R		
$\pi_A(R)$		
S		
$\pi_A(R) \times S$		
$(\pi_A(R) \times S) - R$		
$\pi_A((\pi_A(R) \times S) - R)$		
$\pi_A(R) -$ $\pi_A((\pi_A(R) \times S) - R)$		

# RA $\rightarrow$ DRC: Example $R \div S$

R(A,B) S(B)

RA	DRC	Mapping
R	$R(x, y)$	$x:R.A, y:R.B$
$\pi_A(R)$	$\exists y. R(x, y)$	$x:R.A$
S	$S(z)$	$z:S.B$
$\pi_A(R) \times S$		
$(\pi_A(R) \times S) - R$		
$\pi_A((\pi_A(R) \times S) - R)$		
$\pi_A(R) - \pi_A((\pi_A(R) \times S) - R)$		

# RA $\rightarrow$ DRC: Example $R \div S$

$R(A,B)$   $S(B)$

RA	DRC	Mapping
R	$R(x, y)$	$x:R.A, y:R.B$
$\pi_A(R)$	$\exists y. R(x, y)$	$x:R.A$
S	$S(z)$	$z:S.B$
$\pi_A(R) \times S$	$\exists y. \boxed{R(x, y)} \wedge S(z)$ <i>z needs to be different from y</i>	$x:R.A, z:S.B$
$(\pi_A(R) \times S) - R$	$(\exists y. R(x, y) \wedge S(z)) \wedge \neg R(x, z)$	$x:R.A, z:S.B$
$\pi_A((\pi_A(R) \times S) - R)$	$\exists z [ (\exists y. R(x, y) \wedge S(z)) \wedge \neg R(x, z) ]$	$x:R.A$
$\pi_A(R) -$	$\exists y. R(x, y) \wedge$ <i>x's need to be same variable</i>	$x:R.A$
$\pi_A((\pi_A(R) \times S) - R)$	$\neg \exists z [ (\exists y. R(x, y) \wedge S(z)) \wedge \neg R(x, z) ]$ <i>y's don't need to be same variable</i>	

# RA $\rightarrow$ DRC: Example $R \div S$

$R(A,B)$   $S(B)$

This is the DRC expression we got by translating from RA:

$$\{ x \mid \exists y [R(x, y)] \wedge \neg \exists z [\exists y [R(x, y) \wedge S(z)] \wedge \neg R(x, z)] \}$$

This is the DRC expression for relational division that we saw earlier.

$$\{ x \mid \exists y [R(x, y)] \wedge \neg \exists z [S(z) \wedge \neg R(x, z)] \}$$

Claim: there is no logically equivalent RA expression that uses the table R only twice.  
For details see: <https://arxiv.org/pdf/2203.07284>