

Query Cost Estimation

CS 4750 Database Systems

[A. Silberschatz, H. F. Korth, S. Sudarshan, Database System Concepts, Ch.15]
[C.M. Ricardo, S.D. Urban, "Databases Illuminated, Ch.13]

Review 1: SQL and RA

Consider the following schema statements.

```
student(ID, name, dept_name, tot_cred)
takes(ID, course_ID, sec_id, semester, year, grade)
```

Find IDs and names of all students who have taken more than 3 courses

1. Write SQL query
2. Draw an RA plan

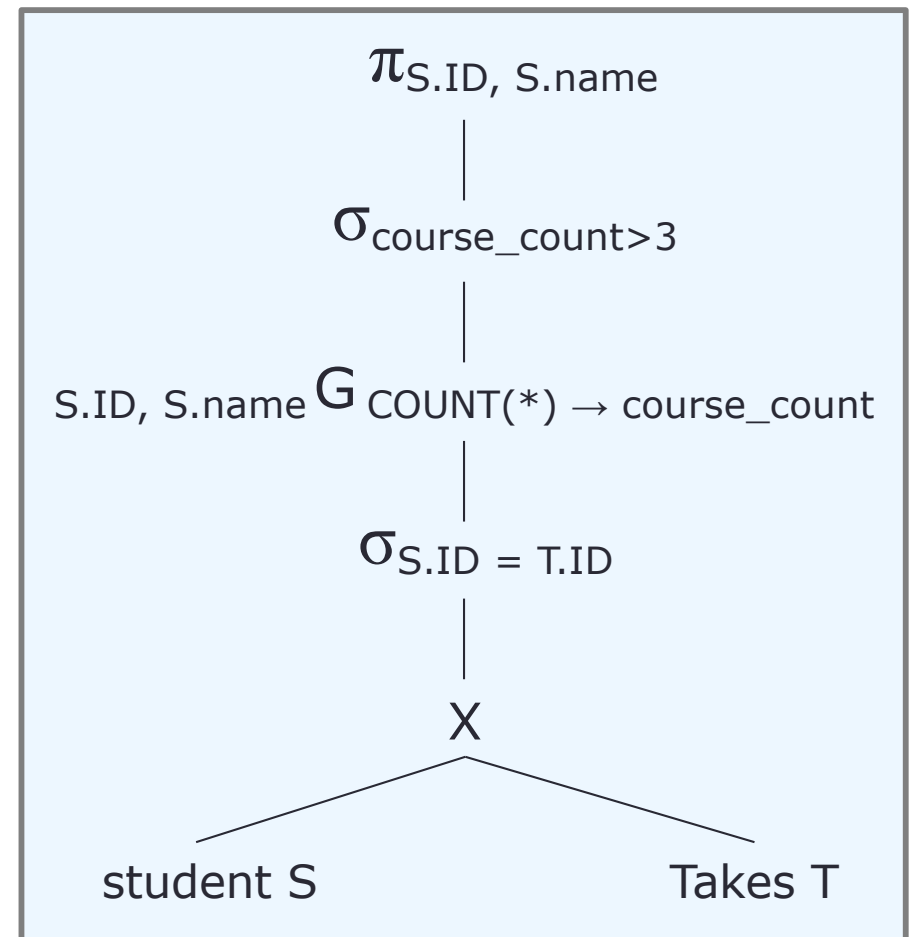
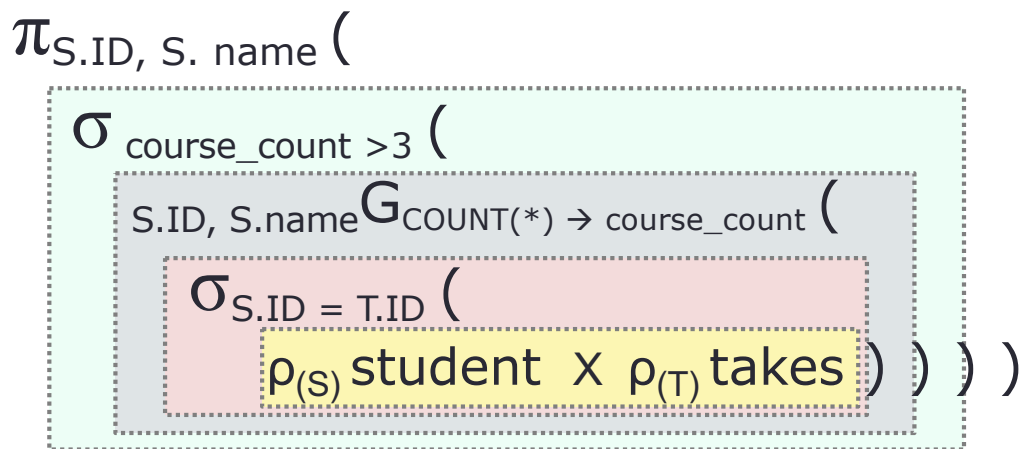
Review 1: SQL and RA (solution)

student(ID, name, dept_name, tot_cred)

takes(ID, course_ID, sec_id, semester, year, grade)

Find IDs and names of all students who have taken more than 3 courses

```
SELECT DISTINCT S.ID, S.name
FROM student S, takes T
WHERE S.ID = T.ID
GROUP BY S.ID, S.name
HAVING COUNT(*) > 3
```



Review 2: SQL and RA

Consider the following schema statements.

```
emp(empno, ename, job, mgr, hiredate, salary, comm, deptno)  
dept(deptno, dname, loc)
```

Find the names of departments where more than three employees are working

1. Write SQL query
2. Draw an RA plan

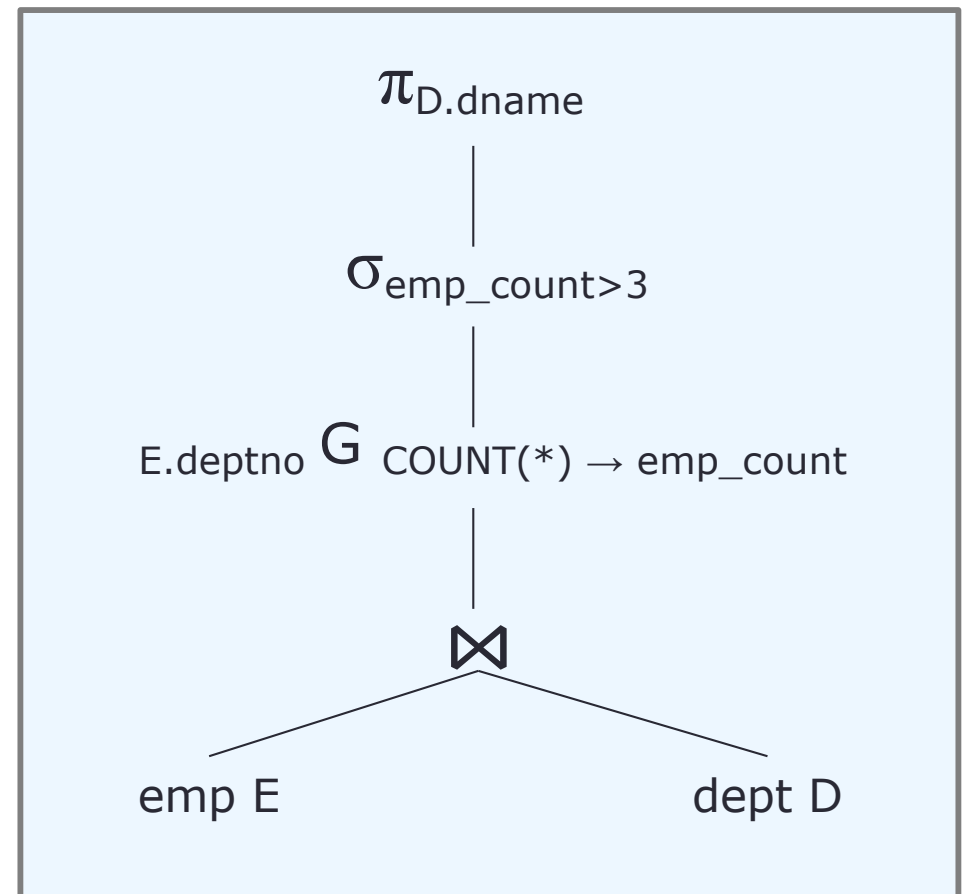
Review 2: SQL and RA (solution 1)

emp(empno, ename, job, mgr, hiredate, salary, comm, deptno)
dept(deptno, dname, loc)

Find the names of departments where more than three employees are working

```
SELECT D.dname
FROM emp E
      NATURAL JOIN dept D
GROUP BY E.deptno
HAVING COUNT(*) > 3;
```

$\pi_{D.dname} (\sigma_{\text{emp_count} > 3} ($
E.deptno \bowtie COUNT(*) \rightarrow emp_count ($\rho_{(E)} \text{ emp} \bowtie \rho_{(D)} \text{ dept}))$



Review 2: SQL and RA (solution 2)

```
emp(empno, ename, job, mgr, hiredate, salary, comm, deptno)  
dept(deptno, dname, loc)
```

Find the names of departments where more than three employees are working

Can you think of another solution?

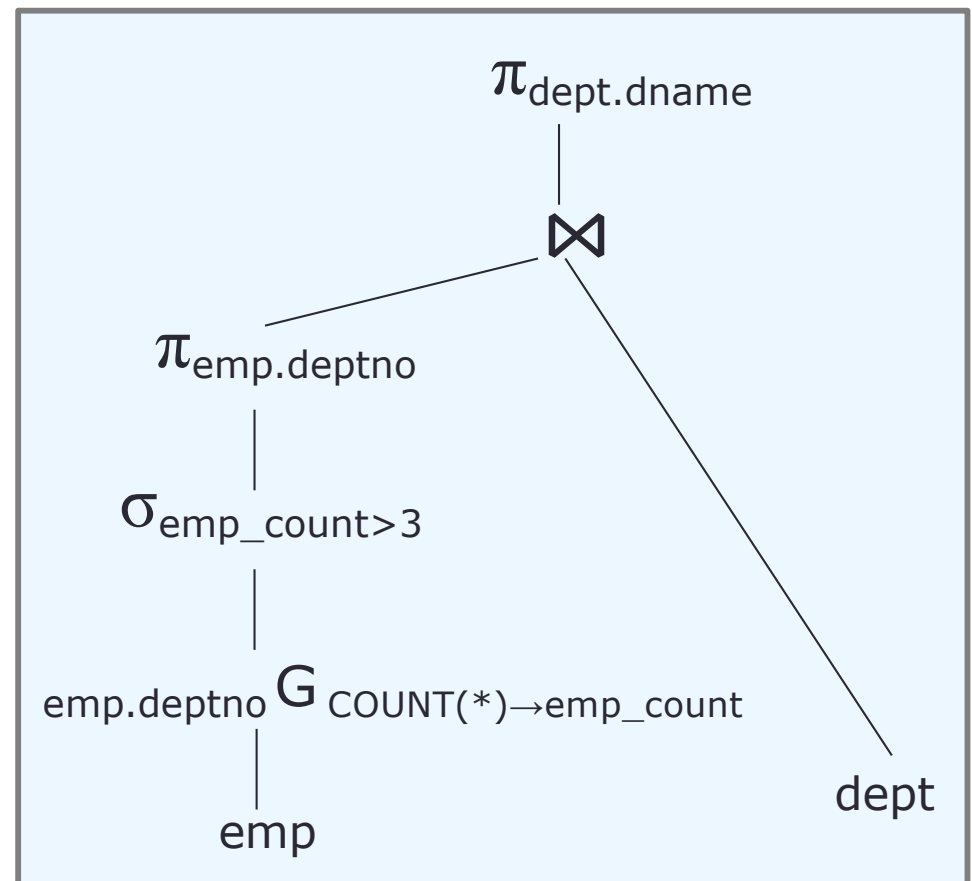
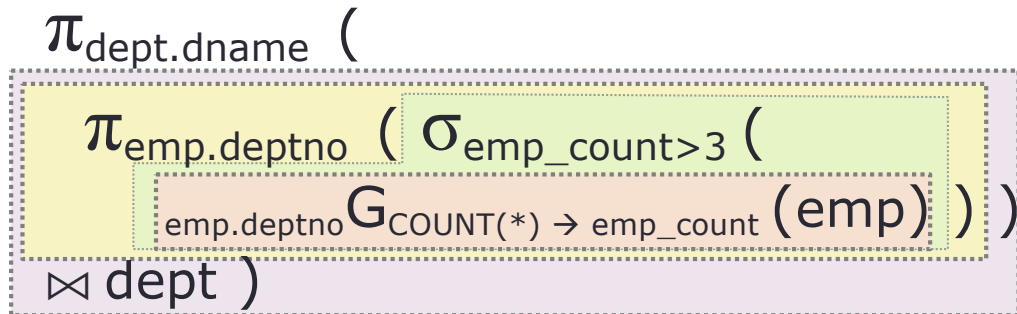
Review 2: SQL and RA (solution 2)

emp(empno, ename, job, mgr, hiredate, salary, comm, deptno)
 dept(deptno, dname, loc)

Find the names of departments where more than three employees are working

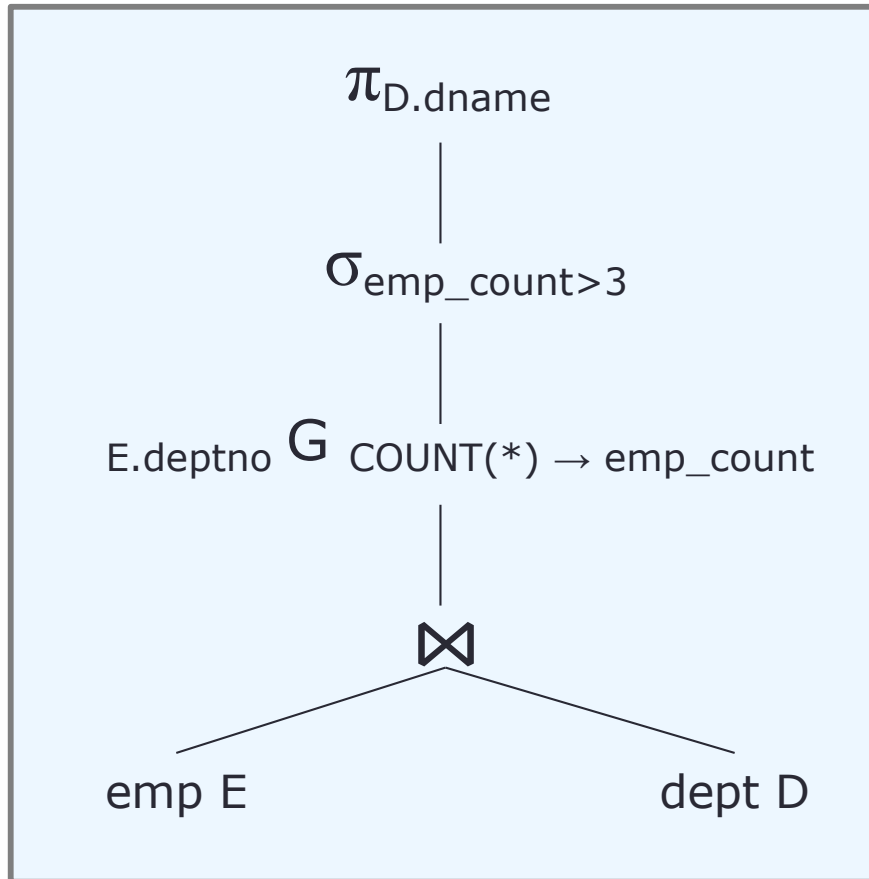
Another solution:

```
SELECT dept.dname
FROM (
    SELECT deptno
    FROM emp
    GROUP BY deptno
    HAVING COUNT(*) > 3 )
NATURAL JOIN dept
```

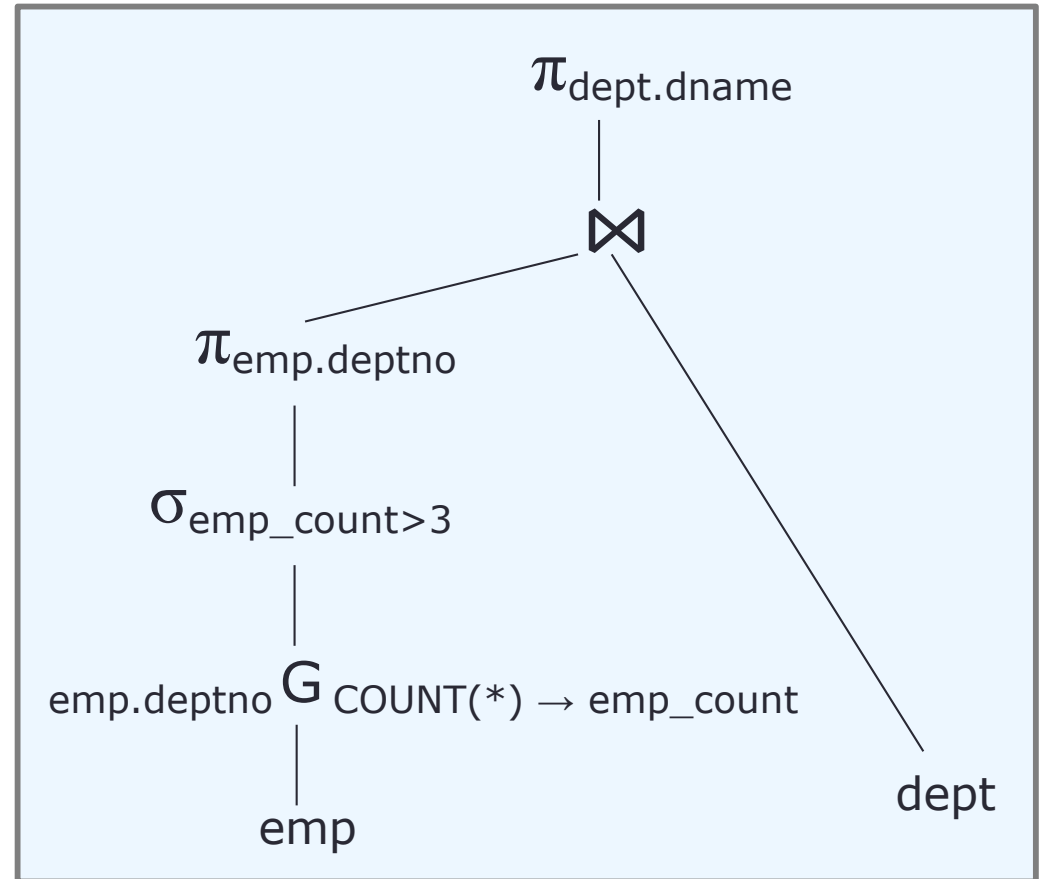


Review 2: SQL and RA (solutions 1 vs. 2)

Can you verify equivalence?



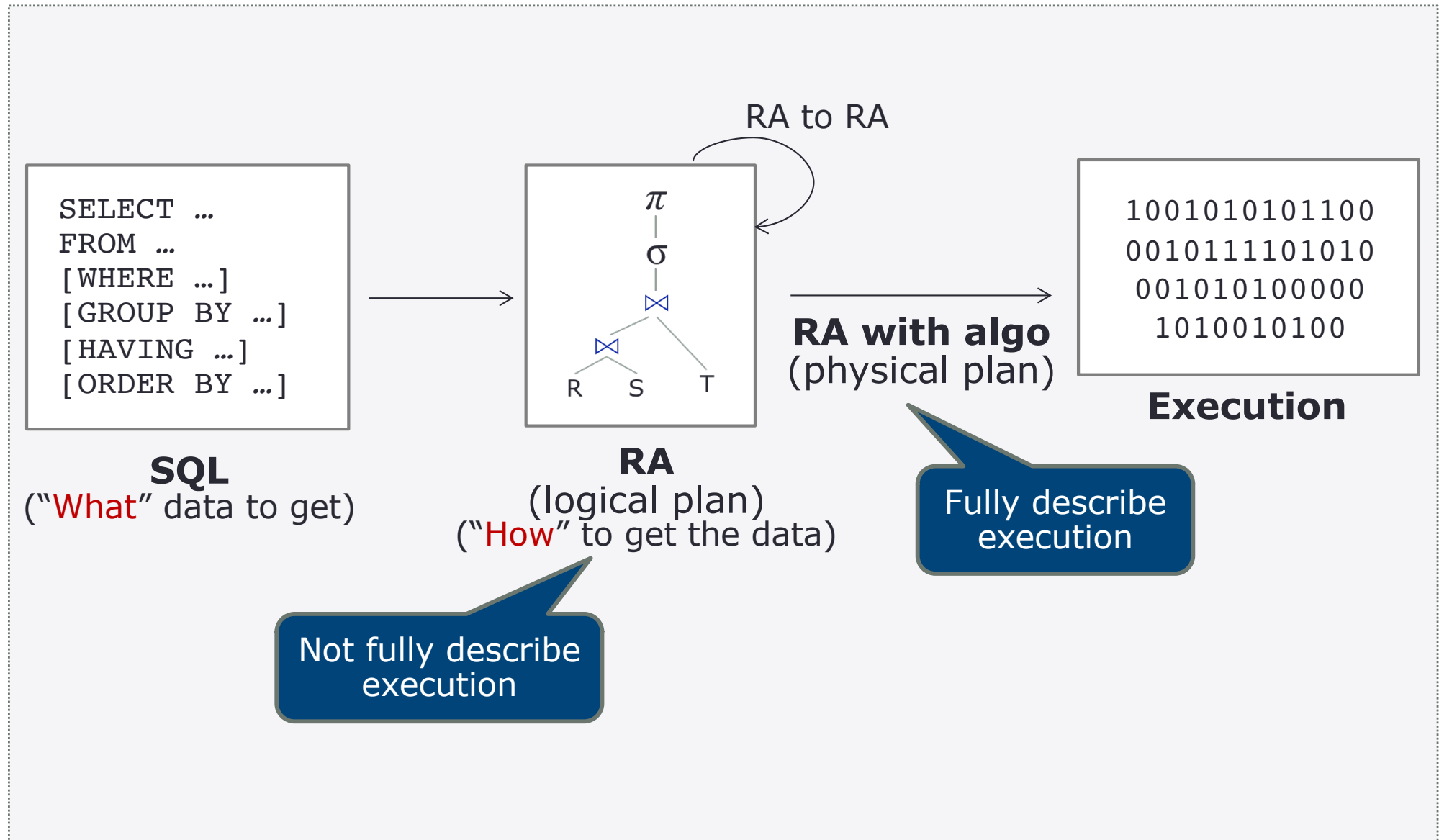
solution 1



solution 2

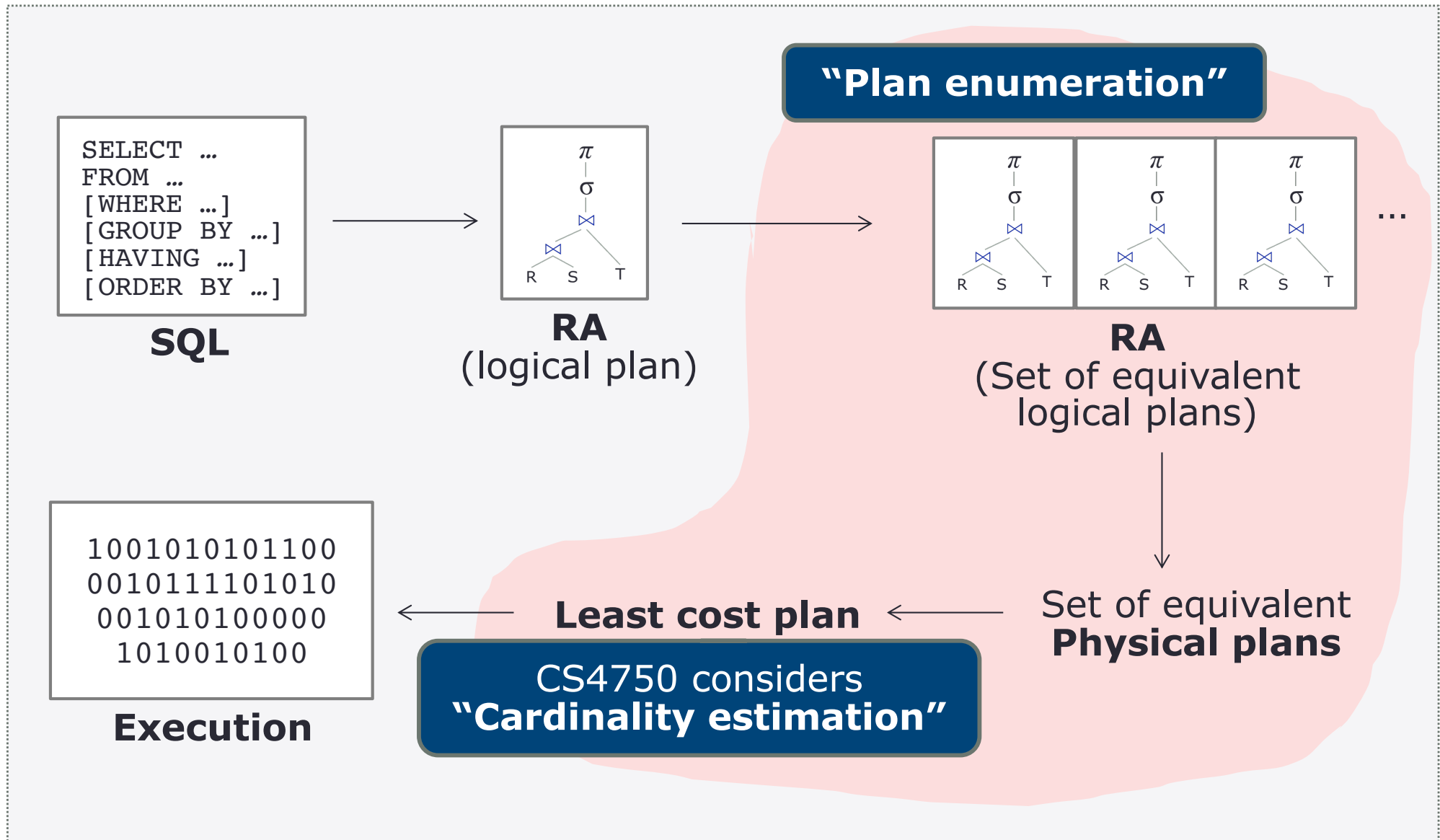
What's the Point of RA?

RDBMS



Overview: Query Processing

RDBMS

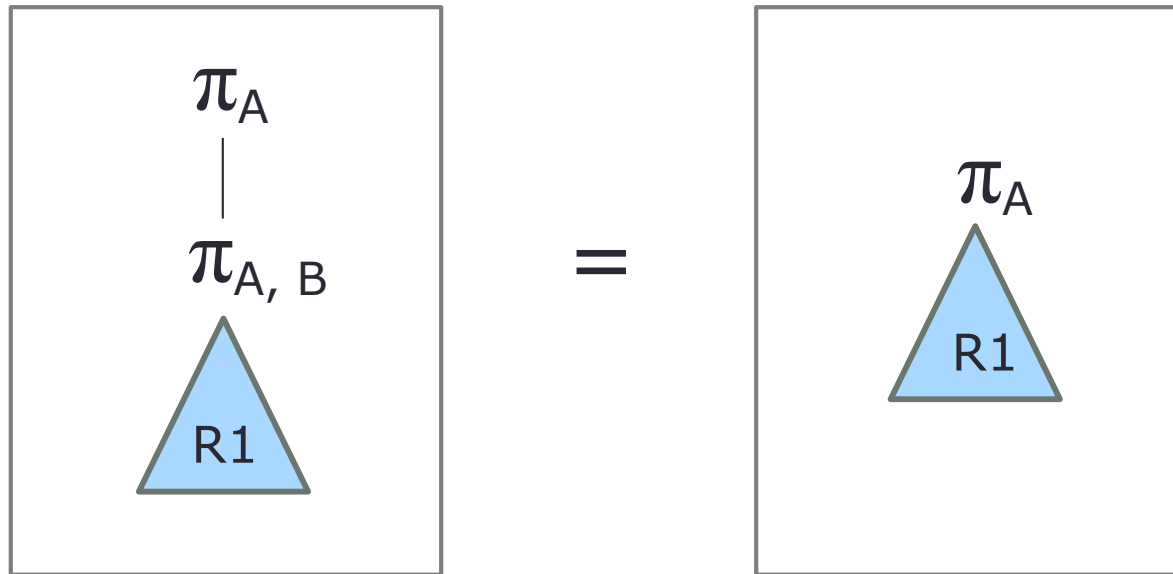


Plan Enumeration (RA to RA)

- Some queries can be expressed in different ways
- RA formulation of a query is important in query processing and optimization
 - RA specifies the order of operations
 - The order can largely determine how efficient the query plan will be
- Why RA equivalences?
 - Simplify queries
 - Make queries faster

Explore equivalent RA plans

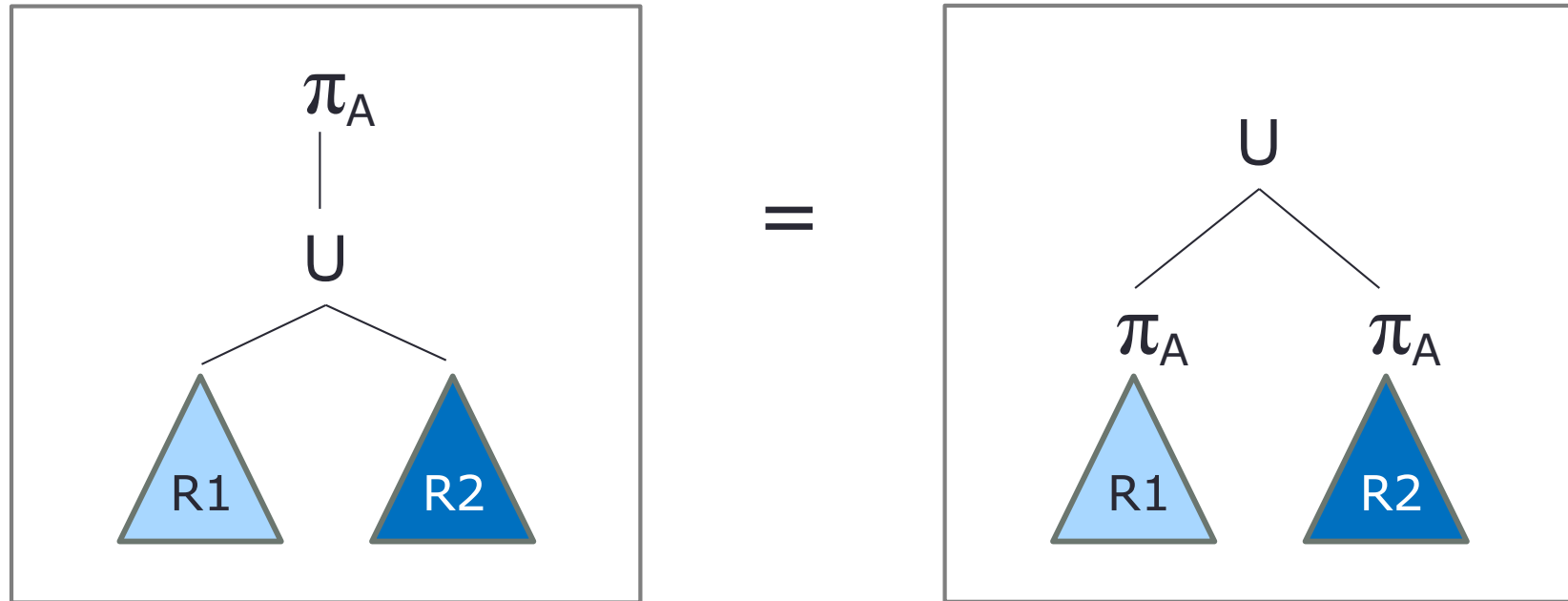
Simplify “Project”



Successive projects can be reduced to the final project
only the last project has to be executed

A and B are sets of attributes; $R1$ is a relation

“Project” and “Union”



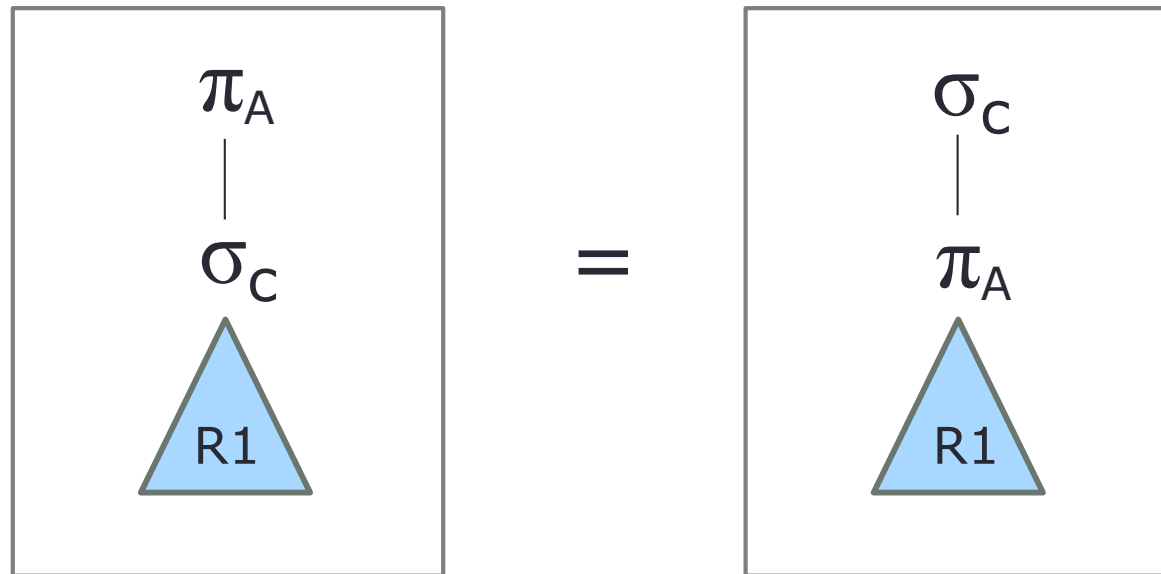
Make sure the **schema matches**

Note: this may not work with intersection or set difference

Project distributes over union

A is a set of attributes; $R1$ and $R2$ are relations

“Select” and “Project”



If C only references attributes in A

Select and project sometimes commute if the condition involves only the attributes in the project list

A is a set of attributes; C is a set of Boolean conditions; $R1$ is a relation

“Select” and “AND”

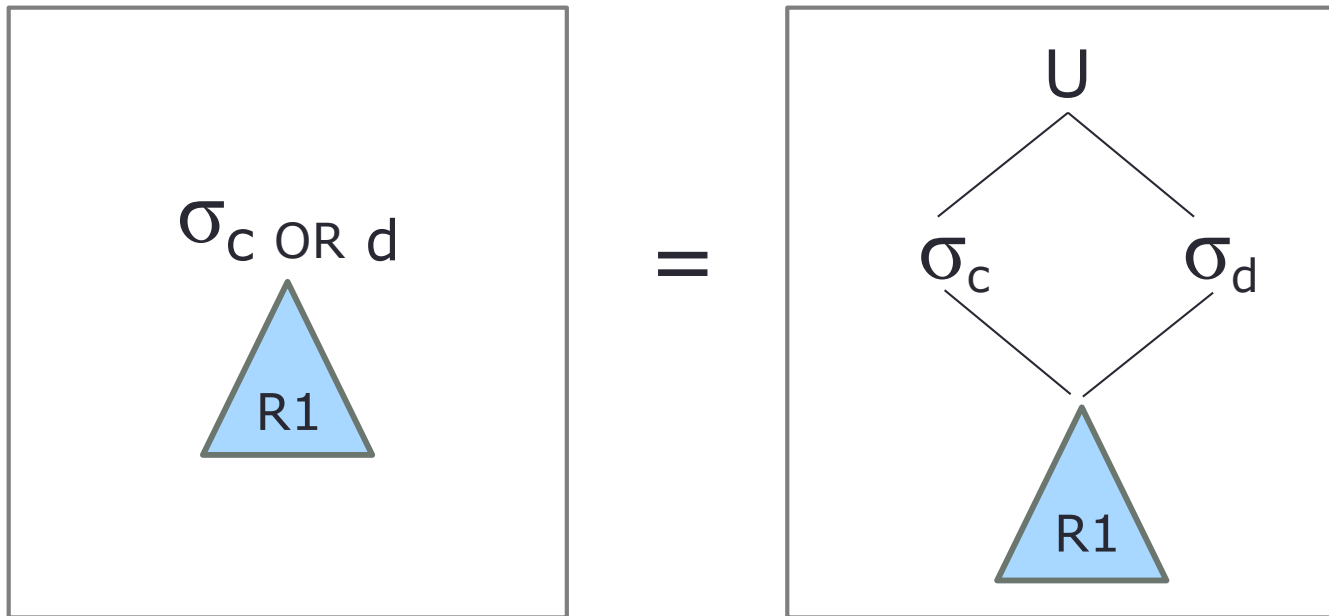


Conjunctive selects can cascade into individual selects

$$\sigma_{c \text{ AND } d} (R1)) = \sigma_d(\sigma_c(R1))$$

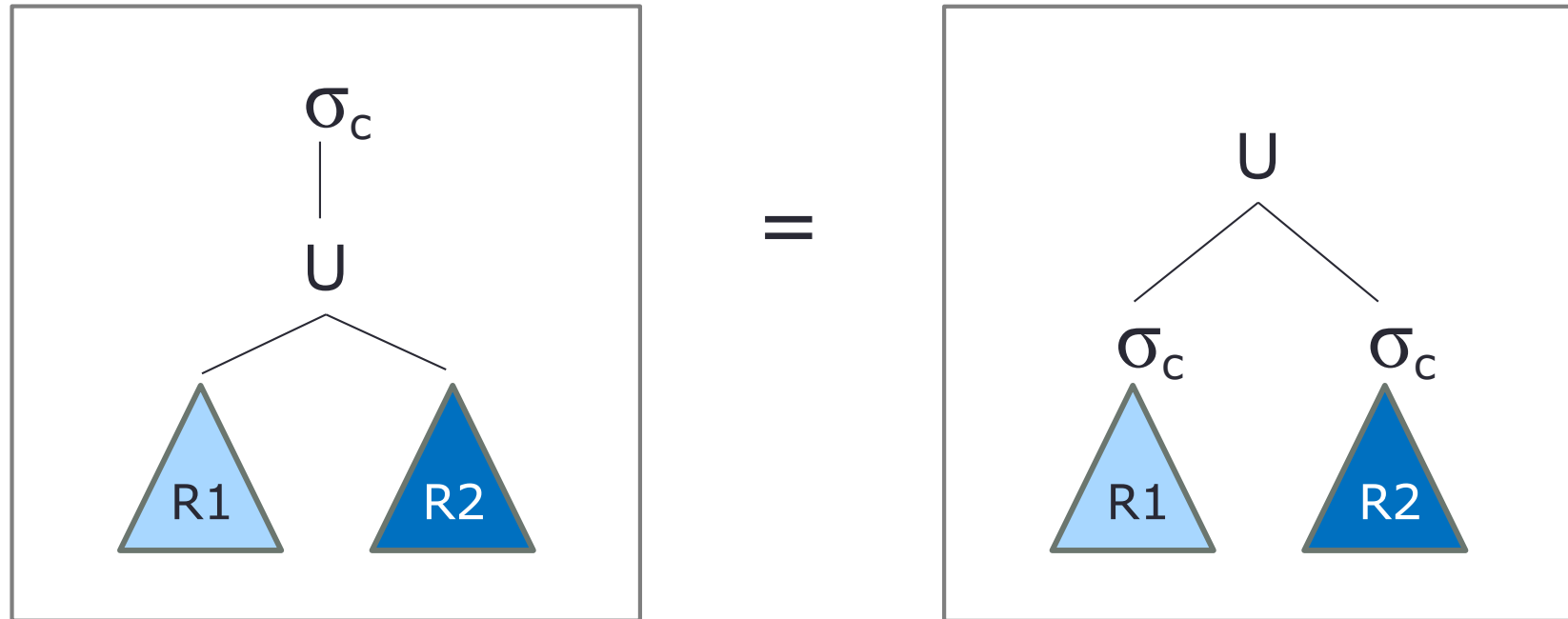
c and d are Boolean conditions; $R1$ is a relation

"Select" and "OR"



c and d are Boolean conditions; $R1$ is a relation

“Select” and “Union”

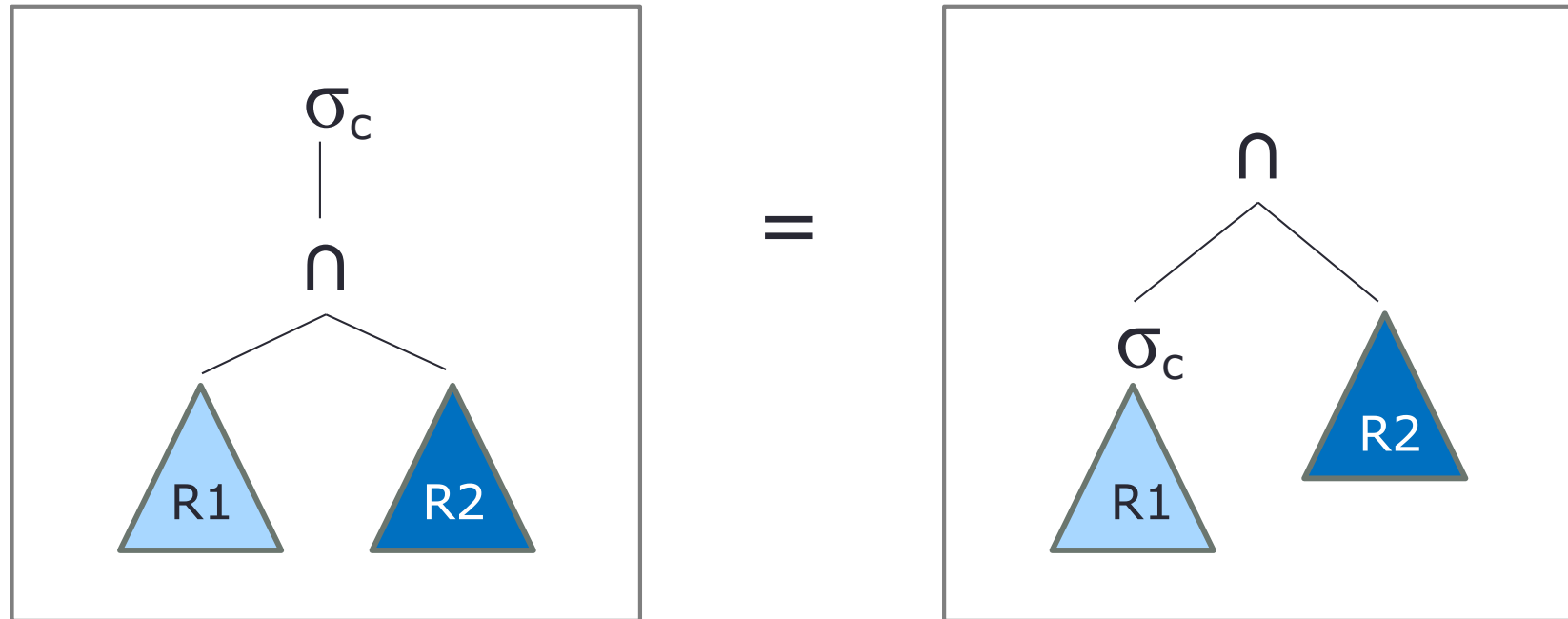


If C references attributes in $R1$ and $R2$

Perform “Select” operation as early as possible

C is a set of Boolean conditions; $R1$ and $R2$ are relations

“Select” and “Intersect”

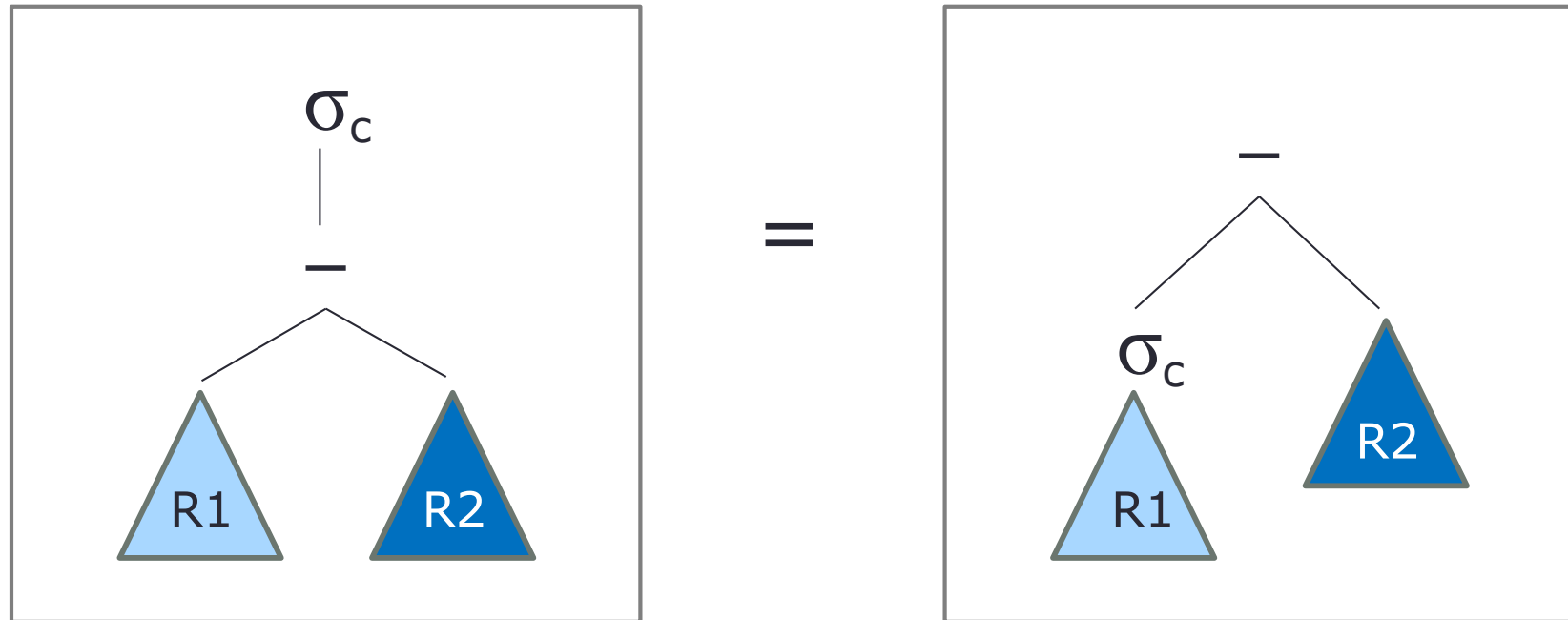


If C only references attributes in $R1$

Perform “Select” operation as early as possible

C is a set of Boolean conditions; $R1$ and $R2$ are relations

“Select” and “Difference”



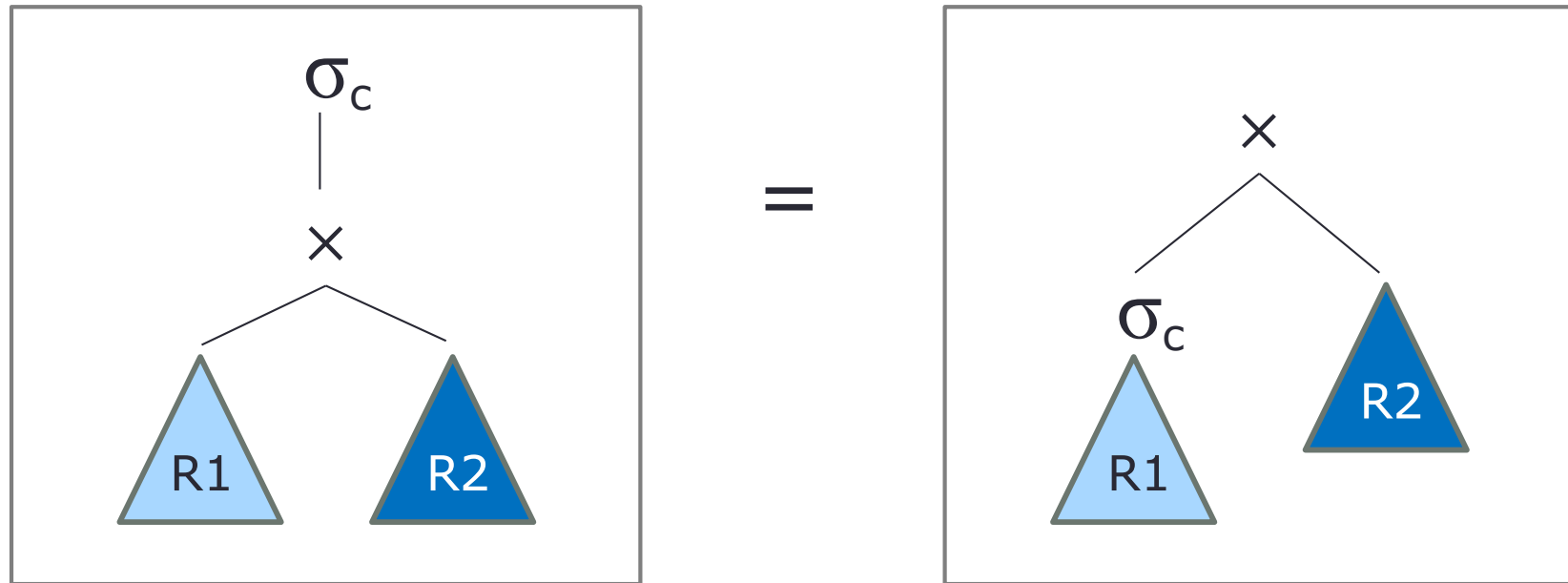
If C only references attributes in $R1$

Perform “Select” operation as early as possible

“EXCEPT” in SQL is equivalent to set “difference” in RA

C is a set of Boolean conditions; $R1$ and $R2$ are relations

“Select” and “Cartesian Product”



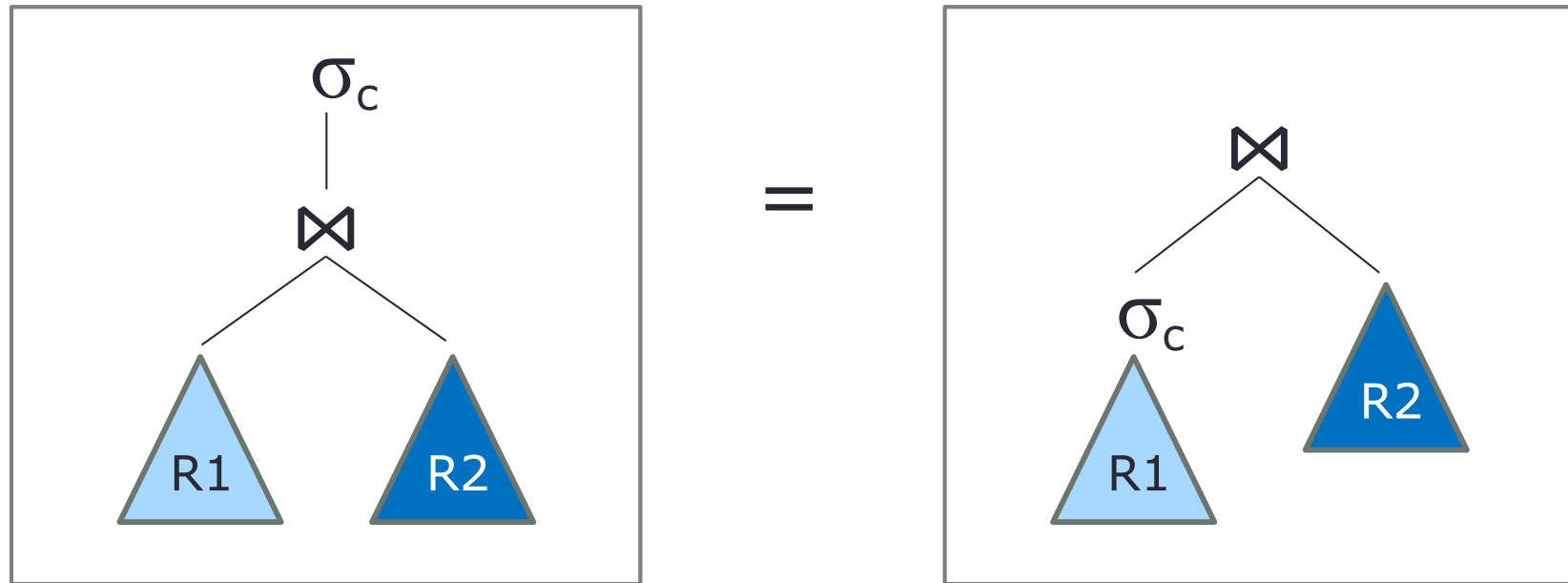
If C only references attributes in $R1$

Select and product
sometimes commute

Perform “Select” operation as
early as possible

C is a set of Boolean conditions; $R1$ and $R2$ are relations

“Select” and “Join”



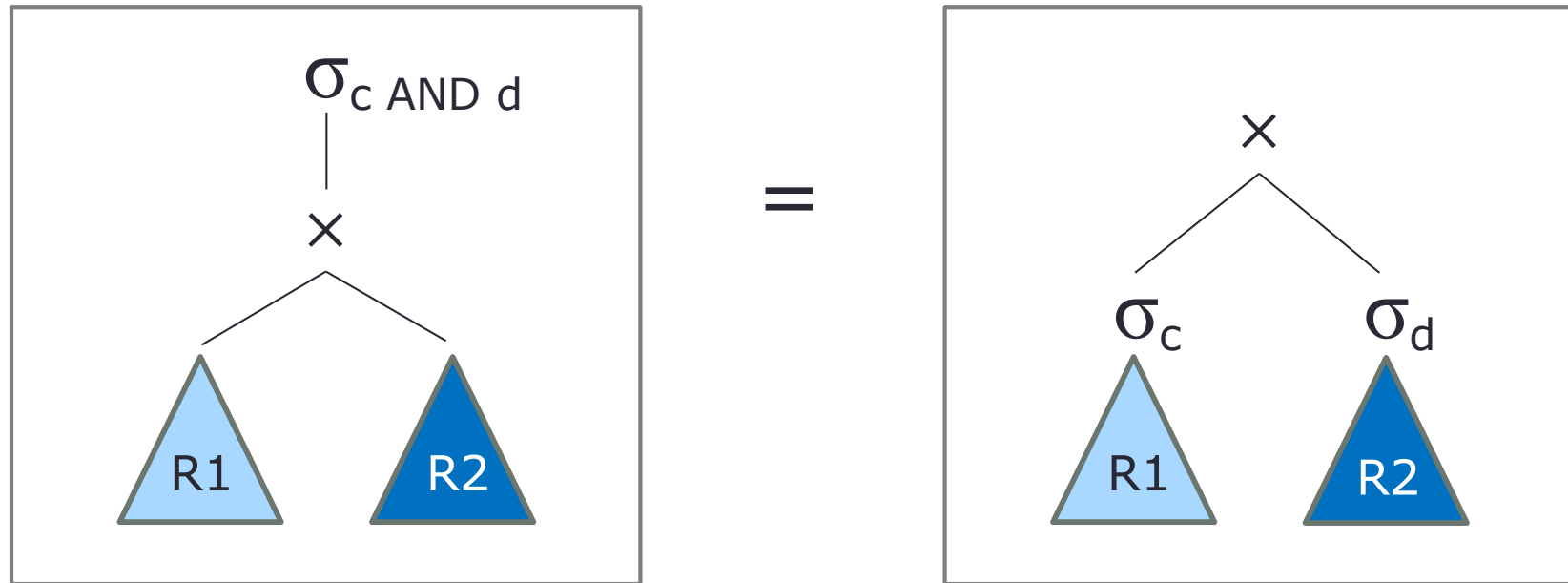
If C only references attributes in $R1$

Select and join
sometimes commute

Perform “Select” operation as
early as possible

C is a set of Boolean conditions; $R1$ and $R2$ are relations

"Select" and "Cartesian Product"



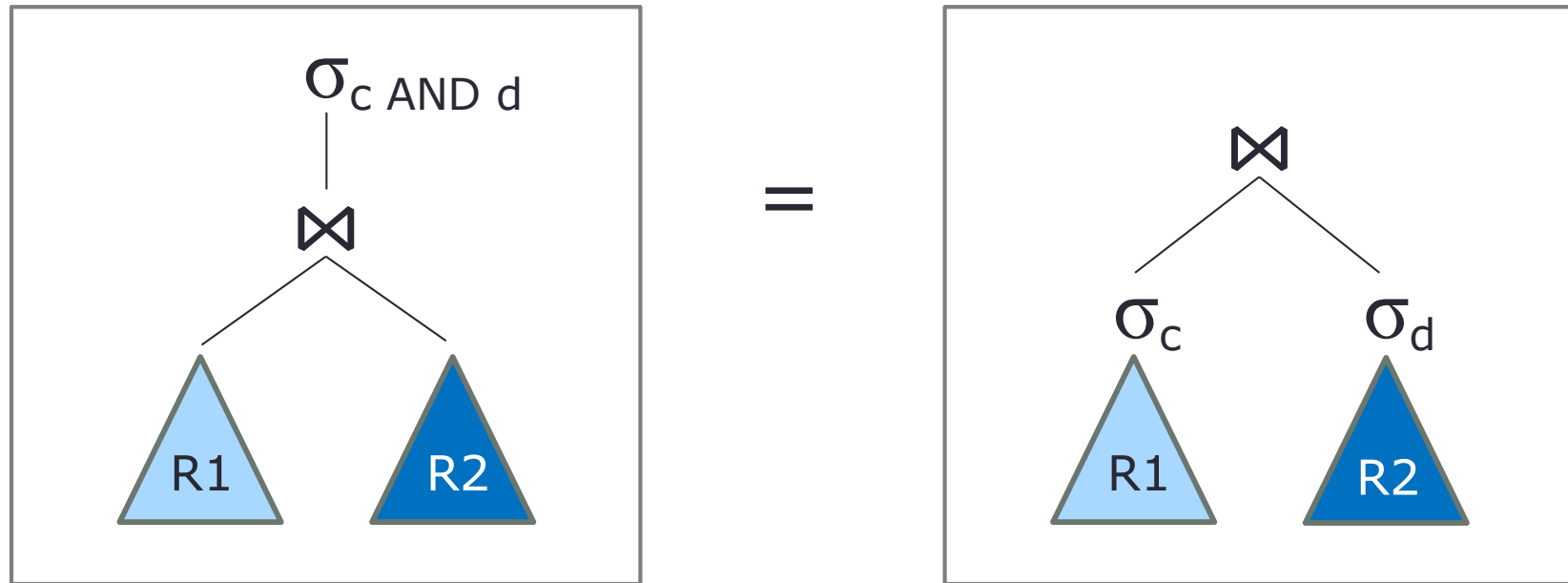
If c references attributes in $R1$ and d references attributes in $R2$

Select sometimes distribute
over product

Perform "Select" operation as
early as possible

C is a set of Boolean conditions; $R1$ and $R2$ are relations

“Select” and “Join”



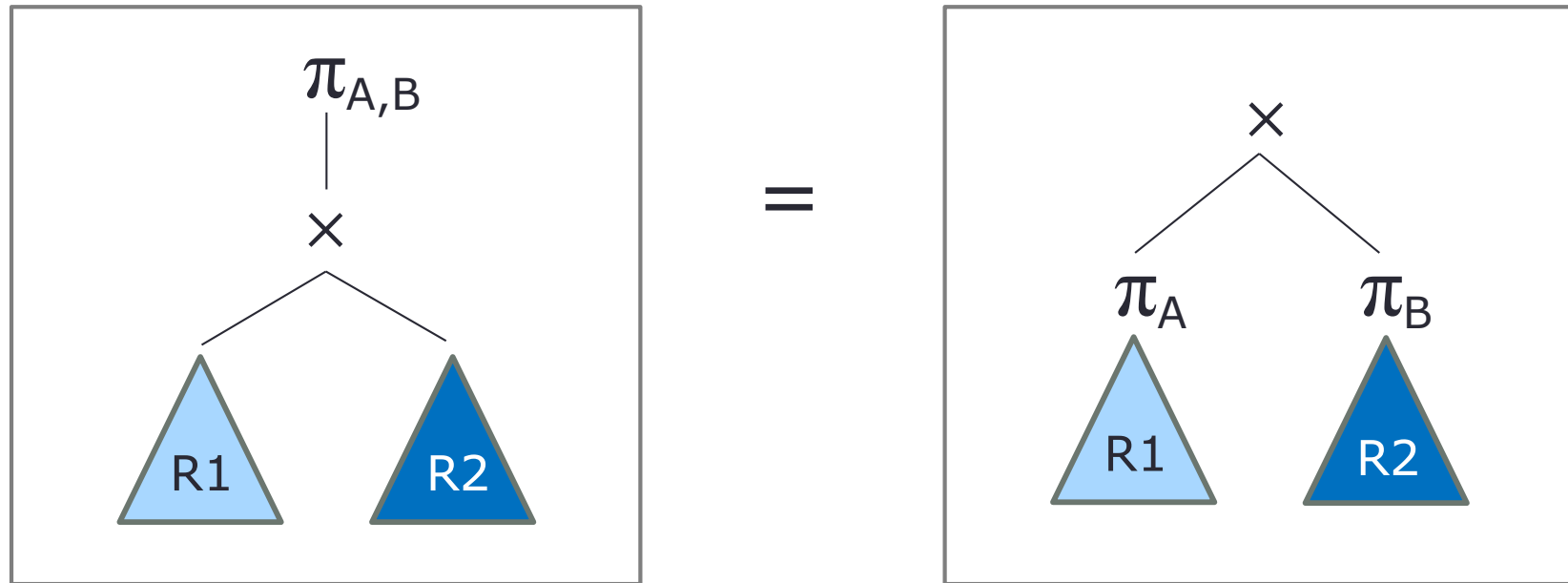
If c references attributes in $R1$ and d references attributes in $R2$

Select sometimes distribute
over join

Perform “Select” operation as
early as possible

C is a set of Boolean conditions; $R1$ and $R2$ are relations

“Project” and “Cartesian Product”

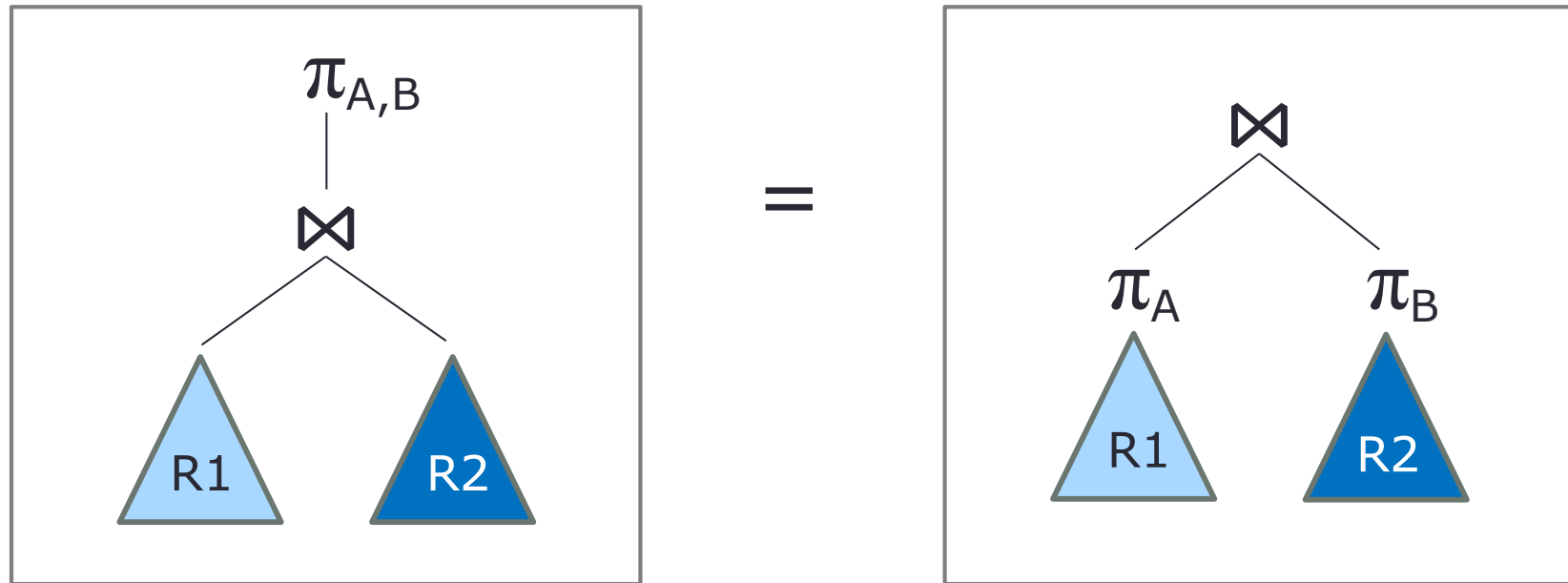


If A contains attributes in $R1$ and B contains attributes in $R2$

Project sometimes distribute
over product

A and B are sets of attributes; $R1$ and $R2$ are relations

“Project” and “Join”

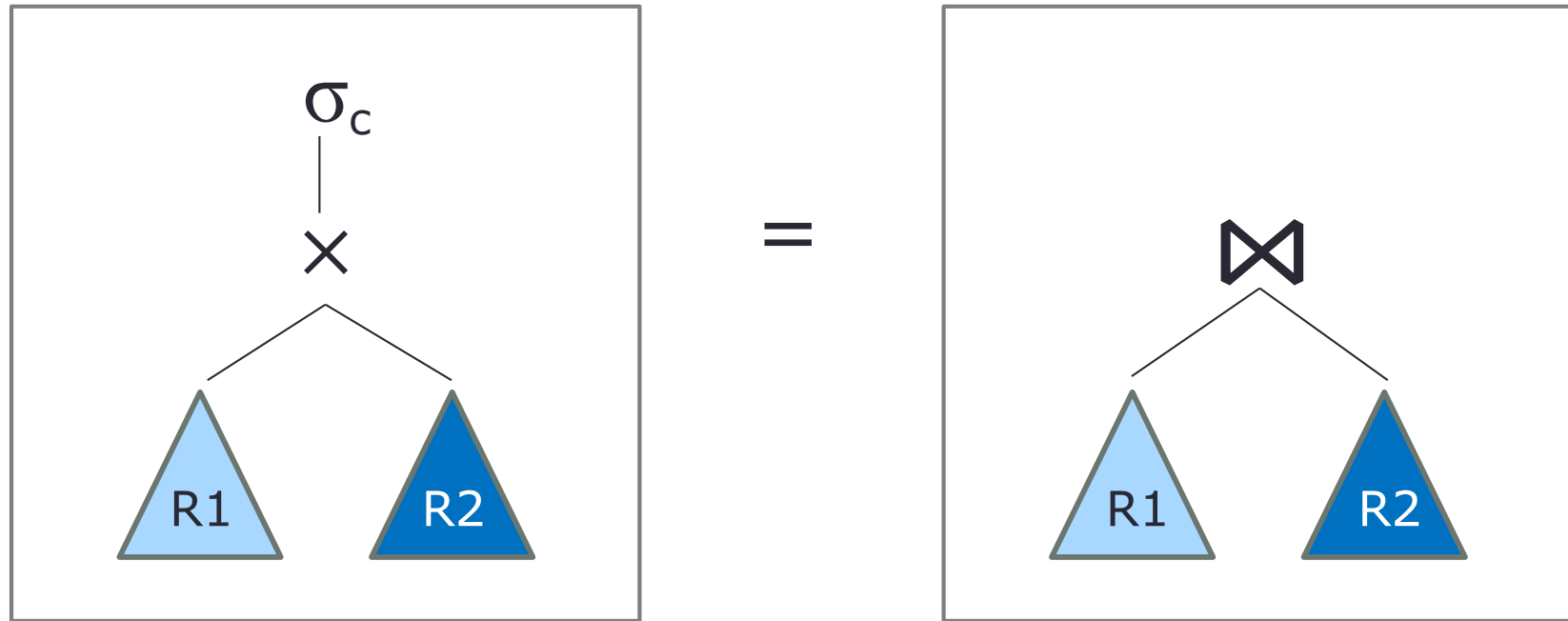


If A contains attributes in $R1$ and B contains attributes in $R2$

Project sometimes distribute
over join

A and B are sets of attributes; $R1$ and $R2$ are relations

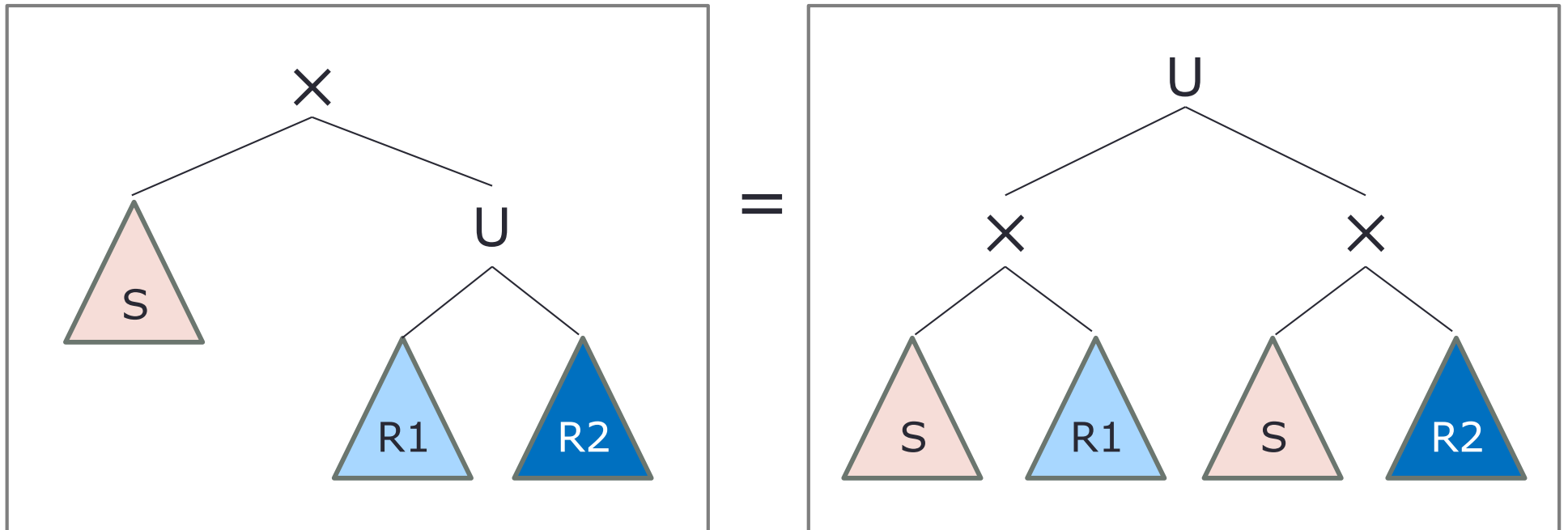
“Cartesian Product” and “Join”



Assume C contains attributes that are used to combine relations

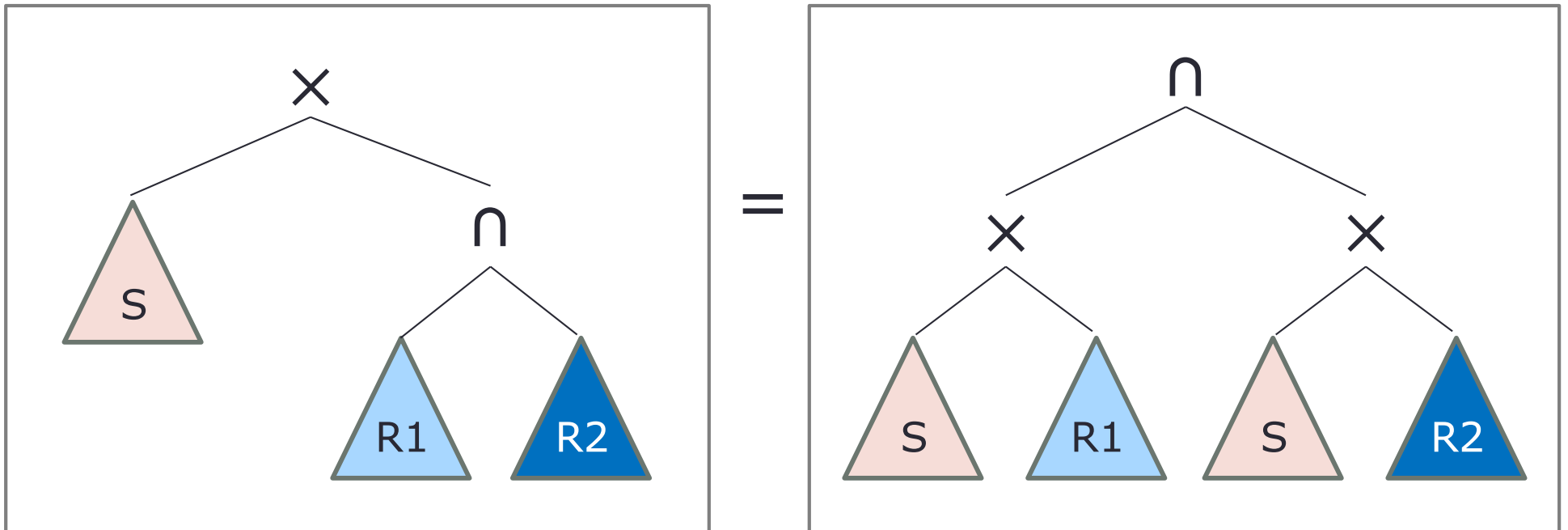
C is a set of Boolean conditions; $R1$ and $R2$ are relations

“Cartesian Product” and “Union”



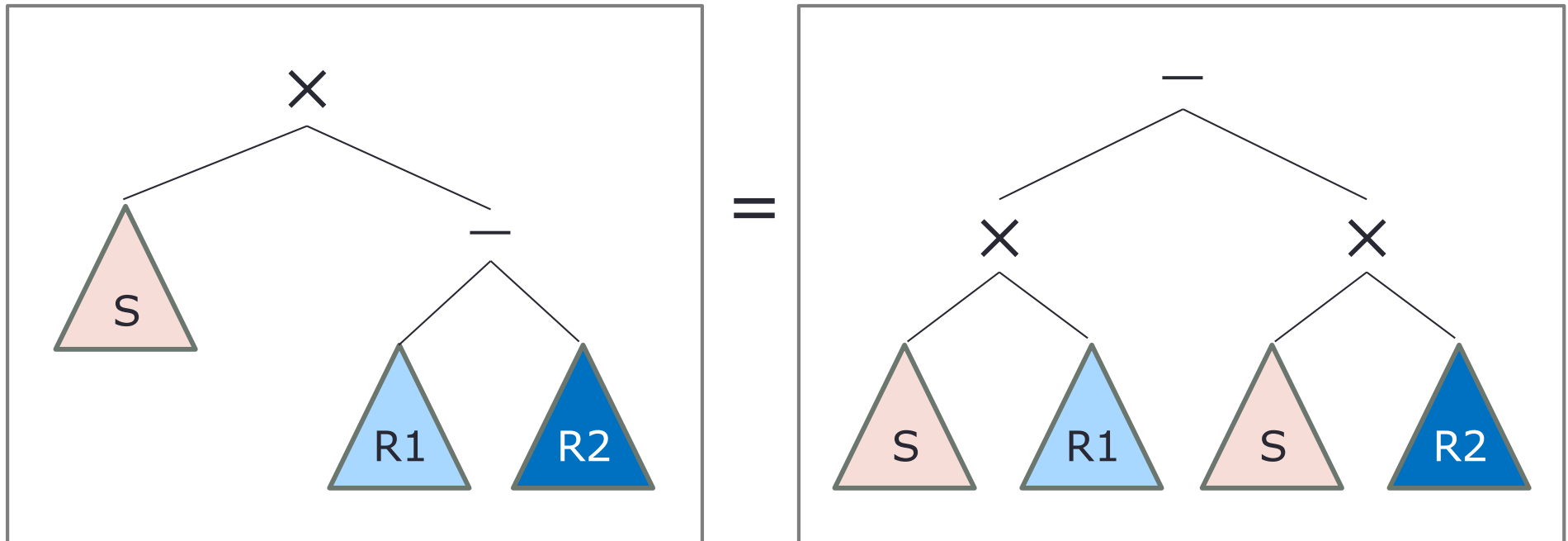
$R1$, $R2$, and S are relations

“Cartesian Product” and “Intersect”



$R1$, $R2$, and S are relations

“Cartesian Product” and “Difference”



“EXCEPT” in SQL is equivalent to set difference in RA

$R1$, $R2$, and S are relations

Let's Try: Equivalent RA (1)

Consider the Sailors database

```
Boats (bid, bname, color)
Sailors (sid, sname, rating, age)
Reserves (sid, bid, day)
```

Find the names of sailors who have reserved a red or green boat

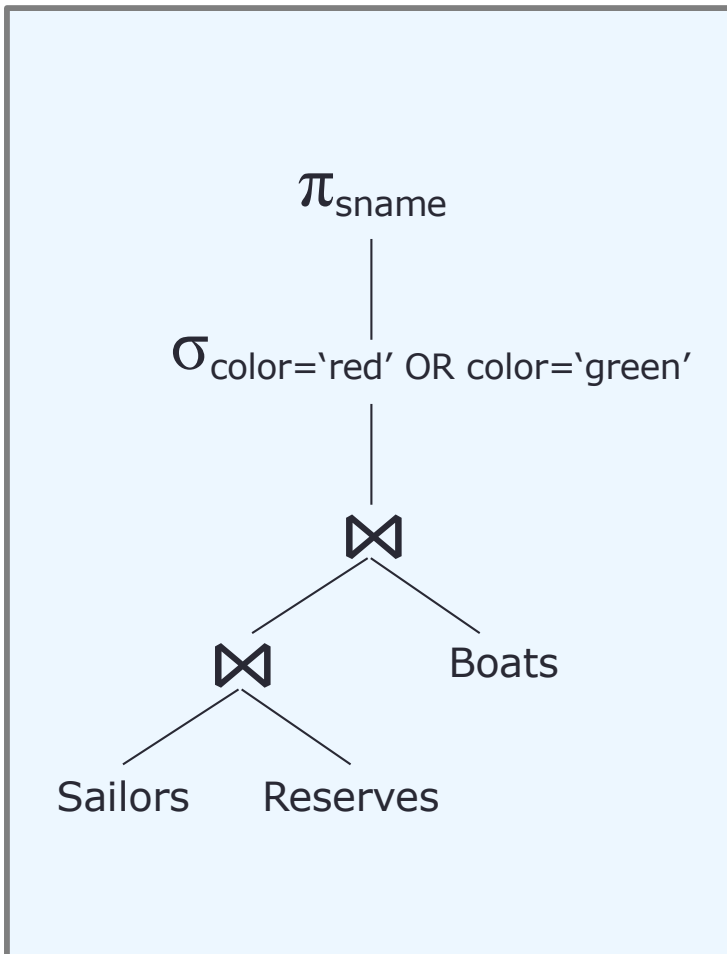
$$\pi_{\text{sname}}(\sigma_{\text{color}=\text{'red'} \text{ OR } \text{color}=\text{'green'}}(\text{Sailors} \bowtie \text{Reserves} \bowtie \text{Boats}))$$

Can you think of an equivalent RA?

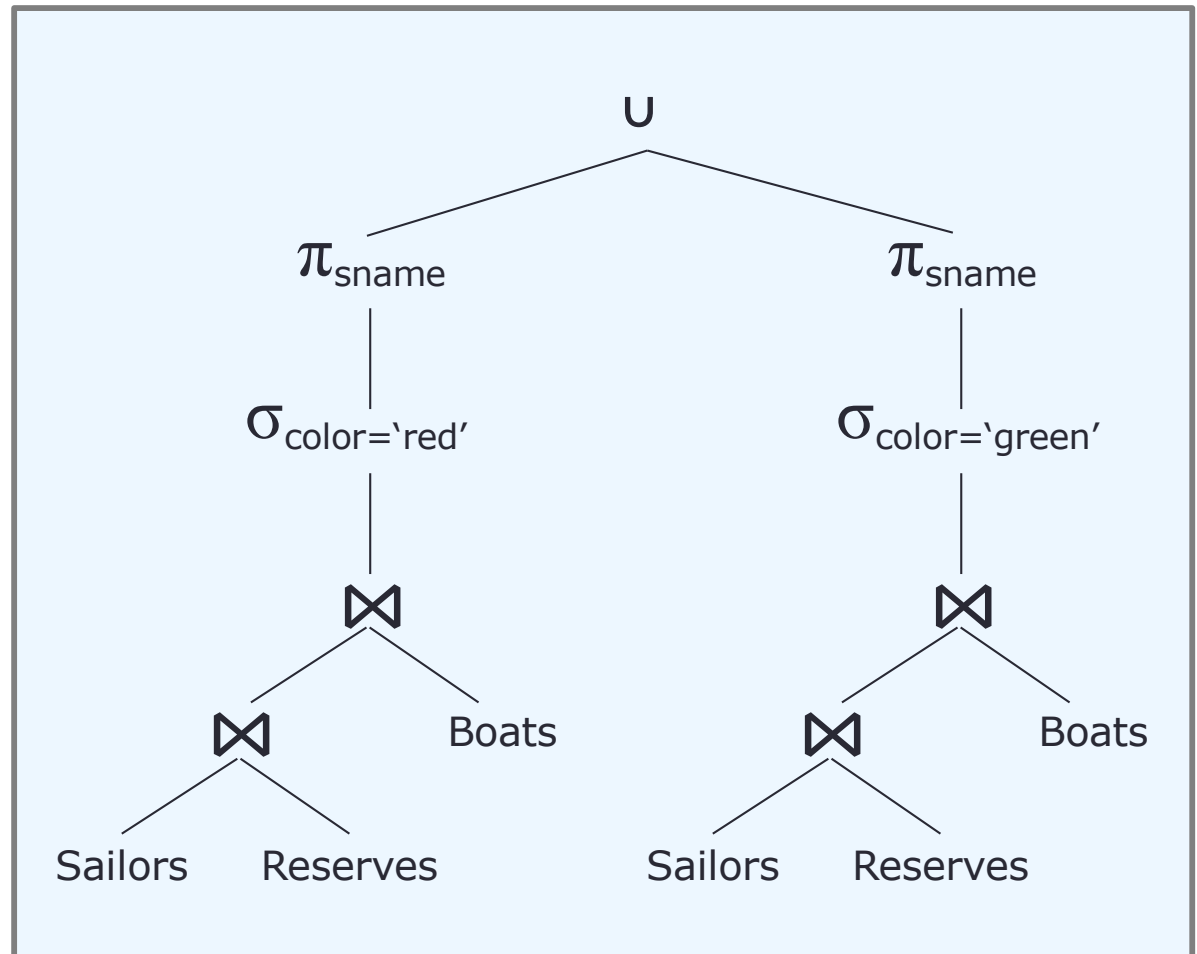
$$\pi_{\text{sname}}(\sigma_{\text{color}=\text{'red'}}(\text{Sailors} \bowtie \text{Reserves} \bowtie \text{Boats})) \cup \pi_{\text{sname}}(\sigma_{\text{color}=\text{'green'}}(\text{Sailors} \bowtie \text{Reserves} \bowtie \text{Boats}))$$

Let's Try: Equivalent RA (1)

$\pi_{\text{sname}}(\sigma_{\text{color}=\text{'red'}} \text{ OR } \text{color}=\text{'green'}} (\text{Sailors} \bowtie \text{Reserves} \bowtie \text{Boats}))$



$\pi_{\text{sname}}(\sigma_{\text{color}=\text{'red'}}(\text{Sailors} \bowtie \text{Reserves} \bowtie \text{Boats})) \cup \pi_{\text{sname}}(\sigma_{\text{color}=\text{'green'}}(\text{Sailors} \bowtie \text{Reserves} \bowtie \text{Boats}))$



Let's Try: Equivalent RA (2)

Consider the Sailors database

Boats (bid, bname, color)

Sailors (sid, sname, rating, age)

Reserves (sid, bid, day)

Find the names of sailors who have reserved boat 103

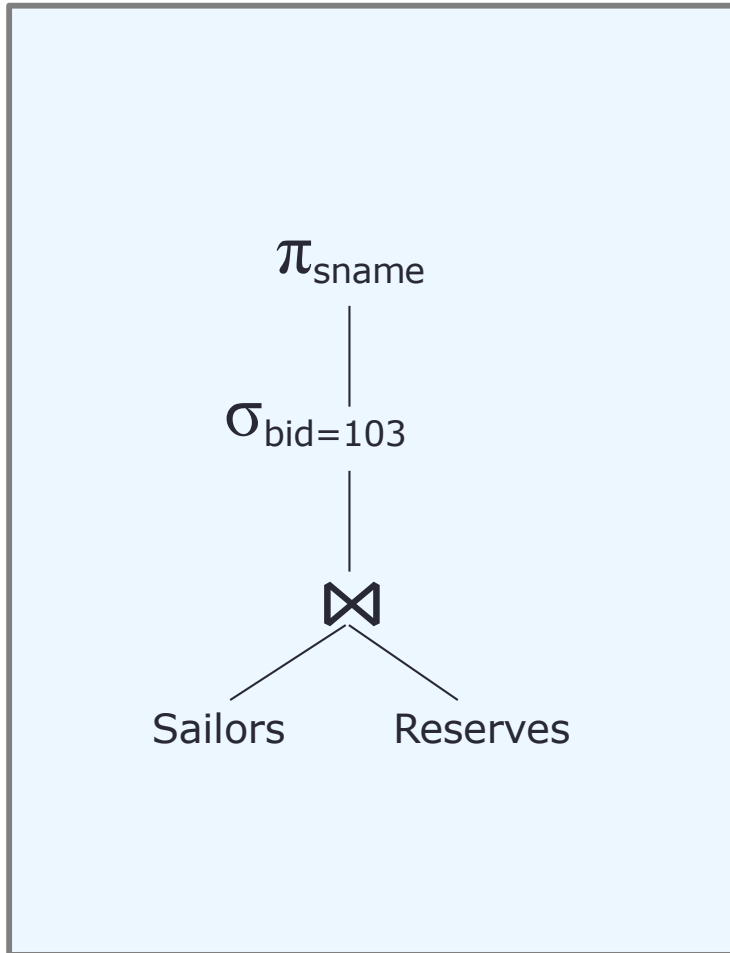
$$\pi_{\text{sname}}(\sigma_{\text{bid}=103}(\text{Sailors} \bowtie \text{Reserves}))$$

Can you think of an equivalent RA?

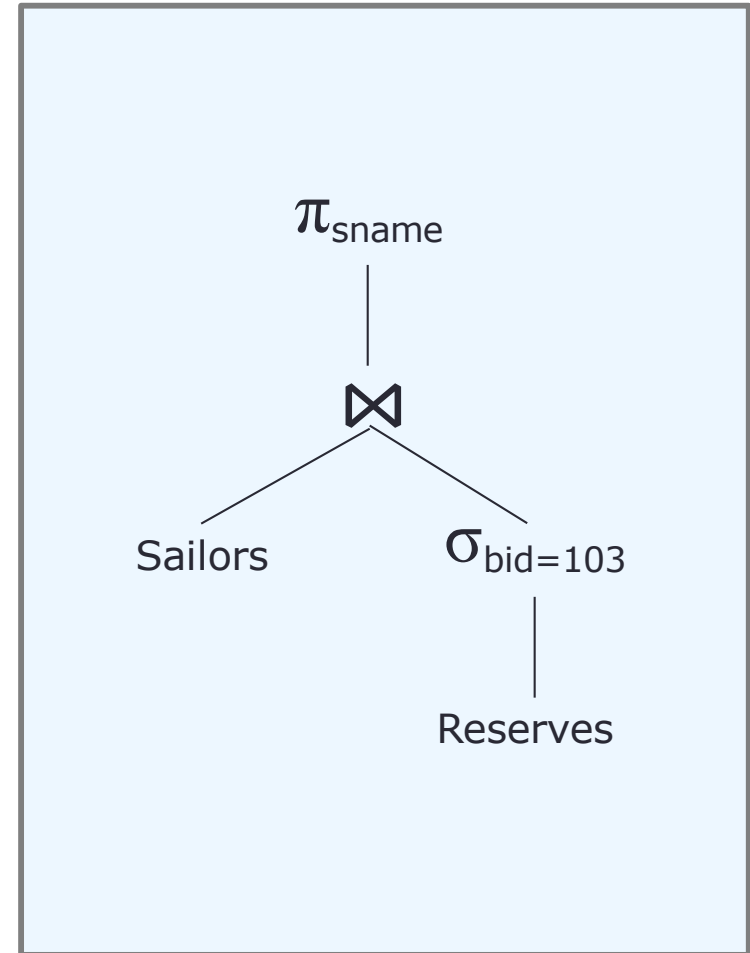
$$\pi_{\text{sname}}(\text{Sailors} \bowtie \sigma_{\text{bid}=103}(\text{Reserves}))$$

Let's Try: Equivalent RA (2)

$\pi_{\text{sname}}(\sigma_{\text{bid}=103}(\text{Sailors} \bowtie \text{Reserves}))$



$\pi_{\text{sname}}(\text{Sailors} \bowtie \sigma_{\text{bid}=103}(\text{Reserves}))$



Let's Try: Equivalent RA (3)

Consider the Sailors database

Boats (bid, bname, color)

Sailors (sid, sname, rating, age)

Reserves (sid, bid, day)

Find the IDs and names of sailors who have not reserved a boat

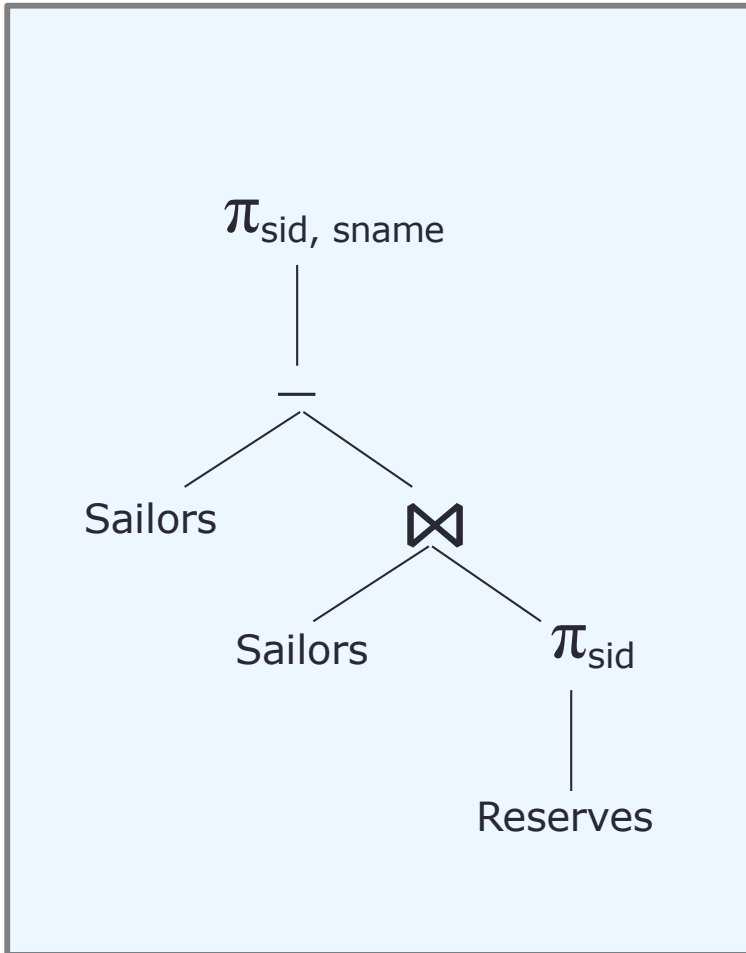
$$\pi_{\text{sid, sname}}(\text{Sailors} - (\text{Sailors} \bowtie \pi_{\text{sid}}(\text{Reserves})))$$

Can you think of an equivalent RA?

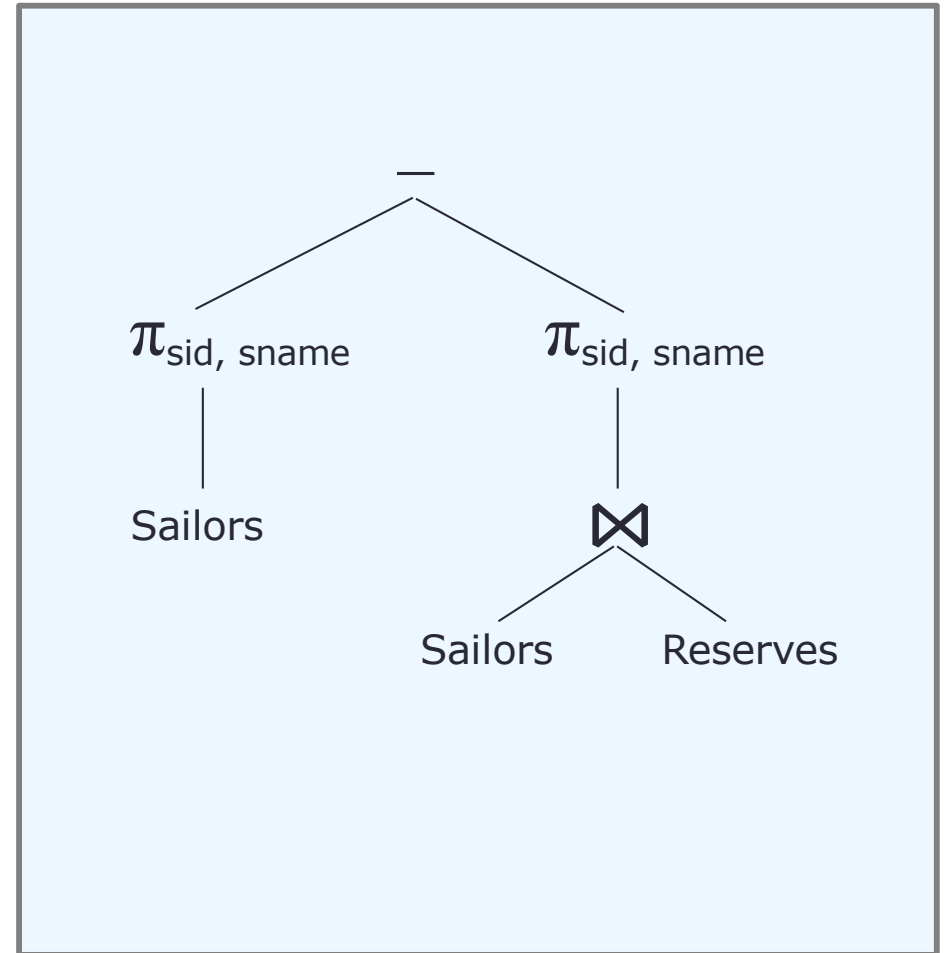
$$\pi_{\text{sid, sname}}(\text{Sailors}) - \pi_{\text{sid, sname}}(\text{Sailors} \bowtie \text{Reserves})$$

Let's Try: Equivalent RA (3)

$\pi_{\text{sid}, \text{sname}}(\text{Sailors} - (\text{Sailors} \bowtie \pi_{\text{sid}}(\text{Reserves})))$



$\pi_{\text{sid}, \text{sname}}(\text{Sailors}) - \pi_{\text{sid}, \text{sname}}(\text{Sailors} \bowtie \text{Reserves})$



More RA Equivalences

- All joins and Cartesian products are commutative

$$R \times S = S \times R \quad (\text{mostly})$$

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

- Joins and Cartesian products are associative

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

$$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$$

- Select is commutative

$$\sigma_c(\sigma_d(R)) = \sigma_d(\sigma_c(R))$$

- Union and intersection are commutative

$$R \cup S = S \cup R$$

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

- Union and intersection are associative

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

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

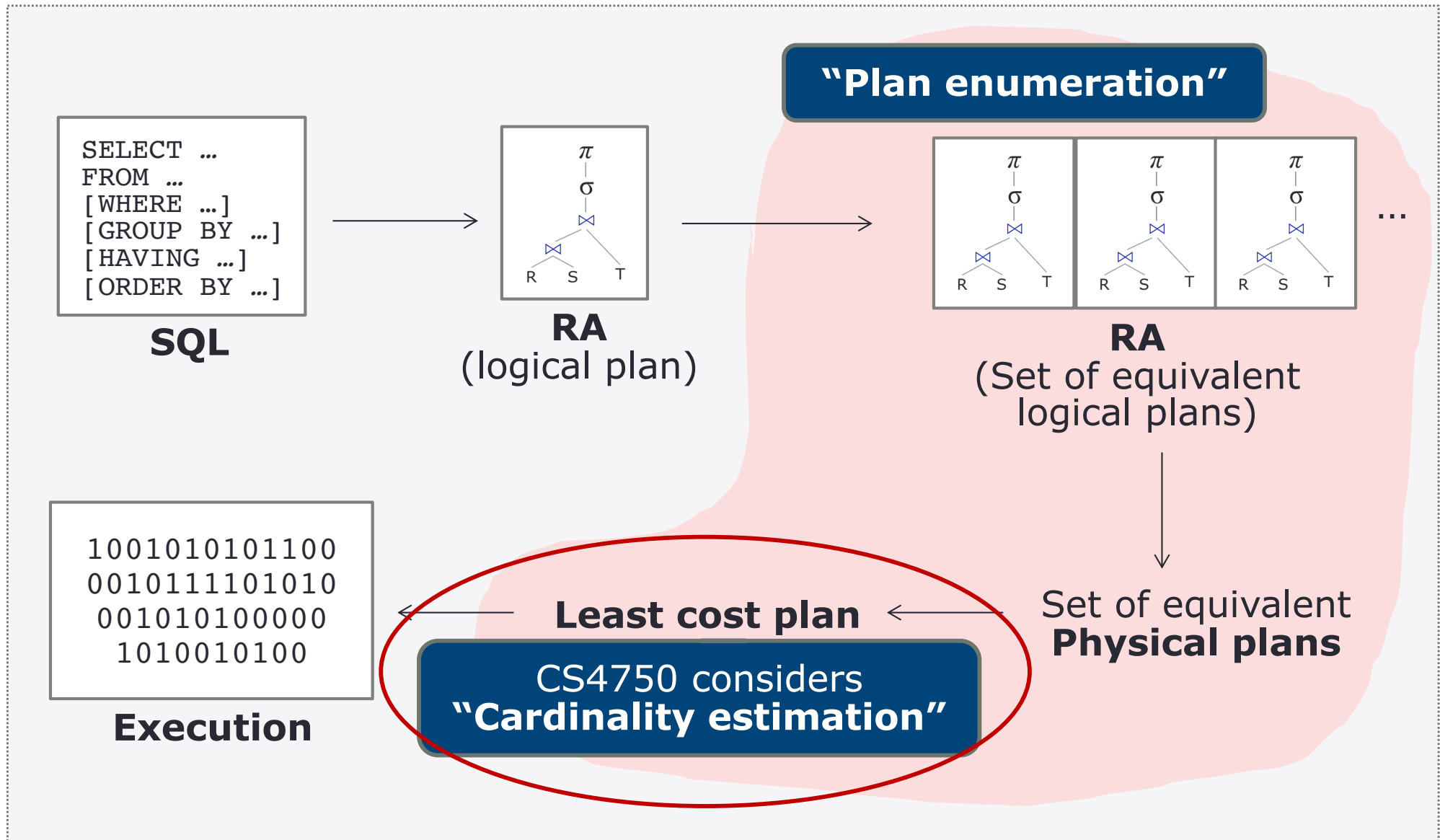
Plenty more equivalences

How to remember?
Use the definitions

R , S , and T are relations

Overview: Query Optimization

RDBMS



Disclaimer

Cost estimation is an active research topic.

Equations and methods discussed here form a foundation of concepts, but usually cannot compare to a commercialized solution.

General Idea on Plan Selection

- Which equivalent RA leads to the most efficient algorithm?
- What algorithm should we use to implement each operation?
- How should the operations pass data from one to the other?

Depends on info available to the query optimizer

- Size of each relation
- Statistics (#blocks, #tuples, #distinct values for an attribute)
- Indexes
- Layout of data on disk

For this class, we assume

- Disk-based storage – HDD
- Row-based storage – tuples are stored contiguously
- HDD I/O cost (reading from disk) only considered
- Sequential disk reads
- A block can be read at once
- No data preloaded

Cardinality Estimation

Estimate the **number of tuples** in the output of each RA operator

Let's use the University database schema as a running example

```
Student (studId, lastName, firstName, major, credits)
Class (classNumber, facId, schedule, room)
Faculty (facId, name, department, rank)
Enroll (studId, classNumber, grade)
```


Estimation: SELECT

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

$T(\text{Student}) = 10,000$

$V(\text{lastName}) = 9,500$

$V(\text{major}) = 10$

$\text{Range}(\text{credits}) = [1, 126)$

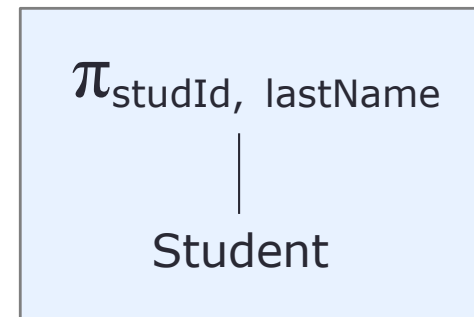
#of tuples

#of distinct values

#of distinct values

range of values

SELECT studId, lastName
FROM Student



How many tuples do we expect this query to output?

10,000

Let's Try: SELECT

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

T(Harris_Teeter) = 1,000

#of tuples

V(name) = 900

#of distinct values

V(category) = 10

#of distinct values

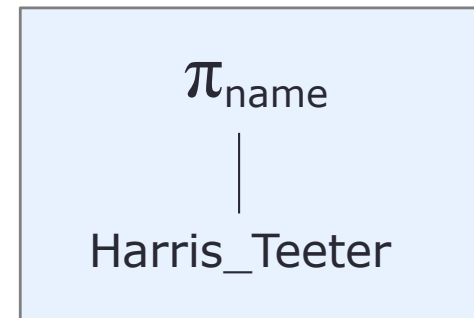
V(price) = 200

#of distinct values

Range(price) = [1,50)

range of values

```
SELECT name
FROM Harris_Teeter
```



How many tuples do we expect this query to output?

SQL (1000) vs. RA (900)

Estimation: DISTINCT

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

$T(\text{Student}) = 10,000$

$V(\text{lastName}) = 9,500$

$V(\text{major}) = 10$

$\text{Range}(\text{credits}) = [1, 126)$

#of tuples

#of distinct values

#of distinct values

range of values

```
SELECT DISTINCT lastName  
FROM Student
```

π_{lastName}

Student

How many tuples do we expect this query to output?

9,500

Let's Try: DISTINCT

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

T(Harris_Teeter) = 1,000

#of tuples

V(name) = 900

#of distinct values

V(category) = 10

#of distinct values

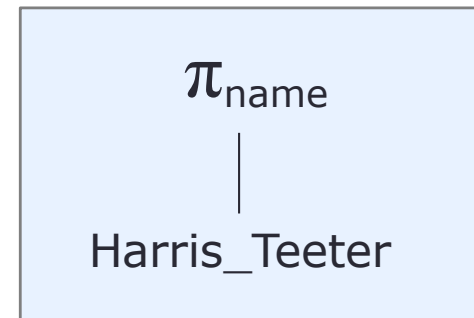
V(price) = 200

#of distinct values

Range(price) = [1,50)

range of values

```
SELECT DISTINCT name
FROM Harris_Teeter
```



How many tuples do we expect this query to output?

900

Estimation: AGGREGATE

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

T(Student) = 10,000

V(lastName) = 9,500

V(major) = 10

Range(credits) = [1, 126)

#of tuples

#of distinct values

#of distinct values

range of values

```
SELECT major, AVG(credits)
FROM Student
GROUP BY major
```

major^GAVG(credits)
|
Student

How many tuples do we expect this query to output?

10

Let's Try: AGGREGATE

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

T(Harris_Teeter) = 1,000

V(name) = 900

V(category) = 10

V(price) = 200

Range(price) = [1,50)

#of tuples

#of distinct values

#of distinct values

#of distinct values

range of values

```
SELECT category, COUNT(id)
FROM Harris_Teeter
GROUP BY category
```

category COUNT(id)
|
Harris_Teeter

How many tuples do we expect this query to output?

10

Estimation: WHERE Value

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

$T(\text{Student}) = 10,000$

$V(\text{lastName}) = 9,500$

$V(\text{major}) = 10$

$\text{Range}(\text{credits}) = [1, 126)$

#of tuples

#of distinct values

#of distinct values

range of values

```
SELECT *  
FROM Student  
WHERE studId = 1111
```

$\sigma_{\text{studId}=1111}$

Student

How many tuples do we expect this query to output?
(assume that 1111 exists)

1

Let's Try: WHERE Value

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

T(Harris_Teeter) = 1,000

V(name) = 900

V(category) = 10

V(price) = 200

Range(price) = [1,50)

#of tuples

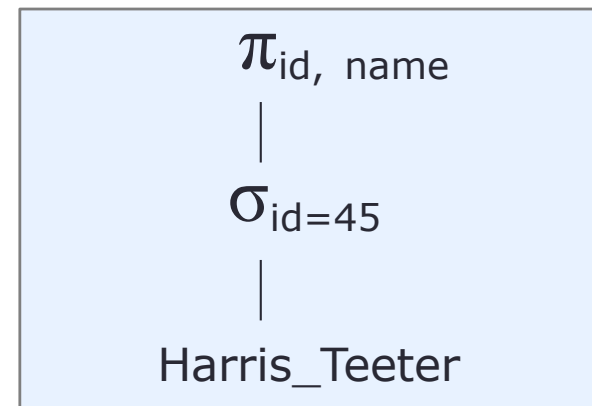
#of distinct values

#of distinct values

#of distinct values

range of values

```
SELECT id, name
FROM Harris_Teeter
WHERE id = 45
```



How many tuples do we expect this query to output?

1

Assume: '45' exists in the distinct values of id

Estimation: WHERE Value

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

$T(\text{Student}) = 10,000$

$V(\text{lastName}) = 9,500$

$V(\text{major}) = 10$

$\text{Range}(\text{credits}) = [1, 126)$

#of tuples

#of distinct values

#of distinct values

range of values

```
SELECT *  
FROM Student  
WHERE lastname = 'Happy'
```

$\sigma_{\text{lastname}='Happy'}$

Student

How many tuples do we expect this query to output?
(assume distinct values uniformly distribute; constant 'Happy' exists;
without assumption, estimate is impossible)

$$T(R) * \frac{1}{V(attr)} = 10000 * \frac{1}{9500} \approx 1.05 \text{ tuples}$$

Selectivity factor

Let's Try: WHERE Value

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

$T(\text{Harris_Teeter}) = 1,000$

#of tuples

$V(\text{name}) = 900$

#of distinct values

$V(\text{category}) = 10$

#of distinct values

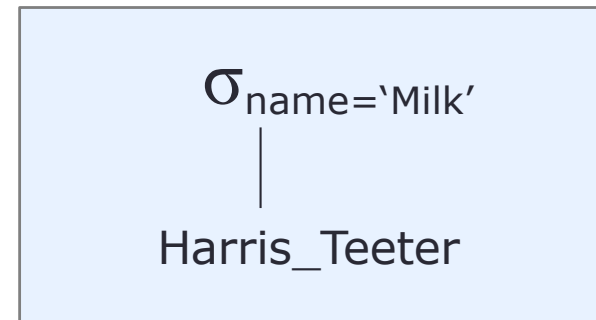
$V(\text{price}) = 200$

#of distinct values

$\text{Range}(\text{price}) = [1, 50]$

range of values

```
SELECT *  
FROM   Harris_Teeter  
WHERE  name = 'Milk'
```



How many tuples do we expect this query to output?

$$T(R) * \frac{1}{V(attr)} = 1000 * \frac{1}{900} \approx 1.11 \text{ tuples}$$

Selectivity factor

Assume: 'Milk' exists in name, and distinct values uniformly distributed

Estimation: WHERE Range

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

$T(\text{Student}) = 10,000$

$V(\text{lastName}) = 9,500$

$V(\text{major}) = 10$

$\text{Range}(\text{credits}) = [1, 126)$

#of tuples

#of distinct values

#of distinct values

range of values

```
SELECT *  
FROM Student  
WHERE credits < 30
```

$\sigma_{\text{credits} < 30}$

Student

How many tuples do we expect this query to output?
(assume uniformly distributed and continuous; without assumption,
estimate is impossible)

$$T(R) * \frac{(\text{val} - \text{min})}{(\text{max} - \text{min})} = 10000 * \frac{(30 - 1)}{(126 - 1)} \approx 2320 \text{ tuples}$$

Selectivity factor = (selection range) / (total range)

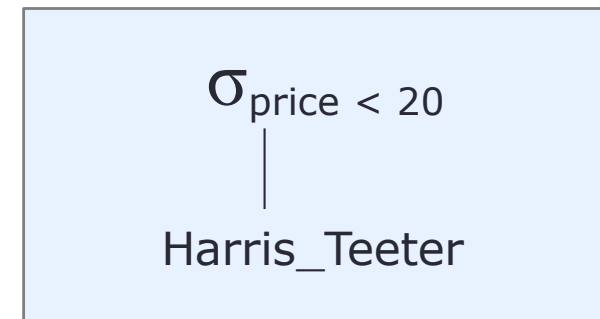
Let's Try: WHERE Range

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

T(Harris_Teeter)	= 1,000	#of tuples
V(name)	= 900	#of distinct values
V(category)	= 10	#of distinct values
V(price)	= 200	#of distinct values
Range(price)	= [1,50]	range of values

```
SELECT *  
FROM Harris_Teeter  
WHERE price < 20
```



How many tuples do we expect this query to output?

$$T(R) * \frac{(\text{val} - \text{min})}{(\text{max} - \text{min})} = 1000 * \frac{(20 - 1)}{(50 - 1)} \approx 387.8 \text{ tuples}$$

Selectivity factor

Assume: distinct values uniformly distributed and continuous

Estimation: AND

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

$T(\text{Student}) = 10,000$

#of tuples

$V(\text{lastName}) = 9,500$

#of distinct values

$V(\text{major}) = 10$

#of distinct values

$\text{Range}(\text{credits}) = [1, 126)$

range of values

```
SELECT *  
FROM Student  
WHERE credits < 30  
AND lastname = 'Happy'
```

$\sigma_{\text{credits} < 30 \text{ AND } \text{lastname} = \text{'Happy'}}$
|
Student

How many tuples do we expect this query to output?
(assume constants exist, distinct values uniformly distributed and continuous, 'Happy' exists; without assumption, estimate is impossible)

Estimation: AND (2)

```
SELECT *  
FROM Student  
WHERE credits < 30  
AND lastname = 'Happy'
```

$\sigma_{\text{credits} < 30 \text{ AND lastname} = \text{'Happy'}}$
|
Student

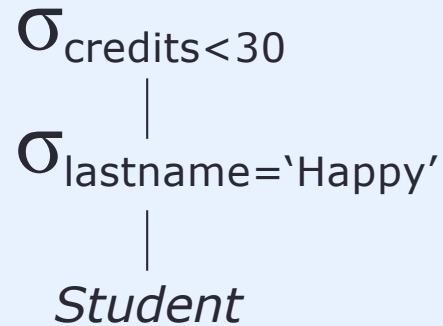
$\sigma_{\text{credits} < 30}$
|
 $\sigma_{\text{lastname} = \text{'Happy'}}$
|
Student



\cap
/ \
 $\sigma_{\text{lastname} = \text{'Happy'}}$ $\sigma_{\text{credits} < 30}$
| |
Student *Student*

Estimation: AND (3)

```
SELECT *  
FROM Student  
WHERE credits < 30  
AND lastname = 'Happy'
```



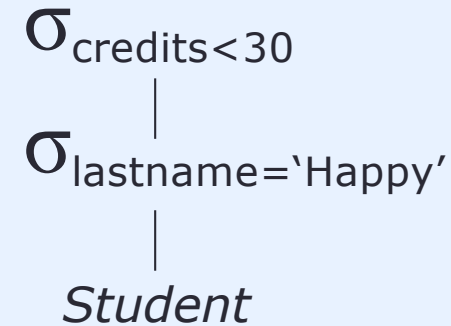
How many tuples do we expect this query to output? – hard to say

- If conditions disjoint, **0** tuple result
 - E.g. no student with lastname 'Happy' has credits <30
- If conditions independent, there will be **multiple** estimates
 - E.g. lastname and credits are independent
- If conditions fully overlap, take **minimum** of estimates
 - E.g. all students with lastname 'Happy' have credits <30

Assume independent unless you know for sure full overlap

Estimation: AND (4)

```
SELECT *  
FROM Student  
WHERE credits < 30  
AND lastname = 'Happy'
```



How many tuples do we expect this query to output? – hard to say

- If conditions disjoint, 0 tuple result

= 0

Selectivity factor

- If conditions independent, there will be multiple estimates

$$\approx 10000 * ((30-1) / (126-1)) * (1/9500) \approx 0.244 \text{ tuples}$$

- If conditions fully overlap, take minimum of estimates

$$\leq 10000 * \min\{ (30-1) / (126-1), (1/9500) \} \approx 1.053 \text{ tuples}$$

(assume independent unless otherwise specified --- answer: 0.244 tuples)

Let's Try: AND

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

T(Harris_Teeter) = 1,000

V(name) = 900

V(category) = 10

V(price) = 200

Range(price) = [1,50]

#of tuples

#of distinct values

#of distinct values

#of distinct values

range of values

```
SELECT *  
FROM Harris_Teeter  
WHERE name='Milk' AND  
category='meat'
```

$\sigma_{\text{name='Milk' AND category='meat'}}$
|
Harris_Teeter

Assume: 'Milk' exists in name, 'meat' exists in category, distinct values uniformly distributed, and conditions independent

$$T(R) * \frac{1}{V(\text{name})} * \frac{1}{V(\text{category})} = 1000 * \frac{1}{900} * \frac{1}{10} \approx 0.11 \text{ tuples}$$

Selectivity factor

Estimation: OR

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

T(Student) = 10,000

V(lastName) = 9,500

V(major) = 10

Range(credits) = [1, 126)

#of tuples

#of distinct values

#of distinct values

range of values

```
SELECT *  
FROM Student  
WHERE credits < 30  
OR lastname = 'Happy'
```

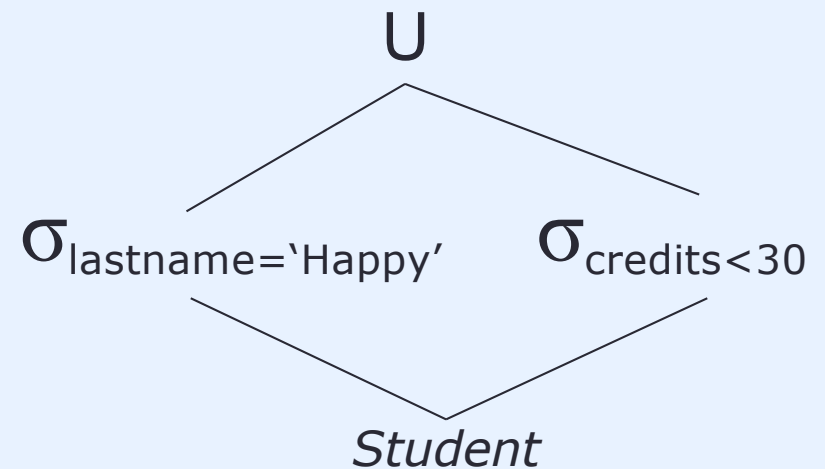
$\sigma_{\text{credits} < 30 \text{ OR } \text{lastname} = \text{'Happy'}}$
|
Student

How many tuples do we expect this query to output?
(assume constants exist, distinct values uniformly distributed and continuous; without assumption, estimate is impossible)

Estimation: OR (2)

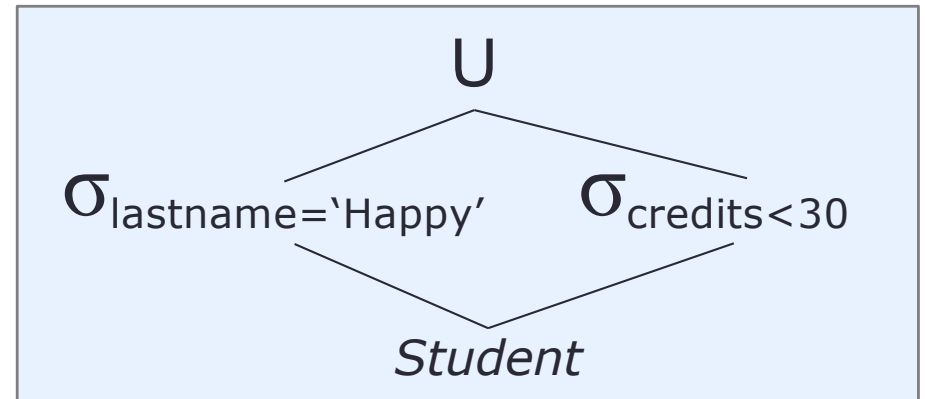
```
SELECT *  
FROM Student  
WHERE credits < 30  
OR lastname = 'Happy'
```

$\sigma_{\text{credits} < 30 \text{ OR lastname} = \text{'Happy'}}$
|
Student



Estimation: OR (3)

```
SELECT *  
FROM Student  
WHERE credits < 30  
OR lastname = 'Happy'
```



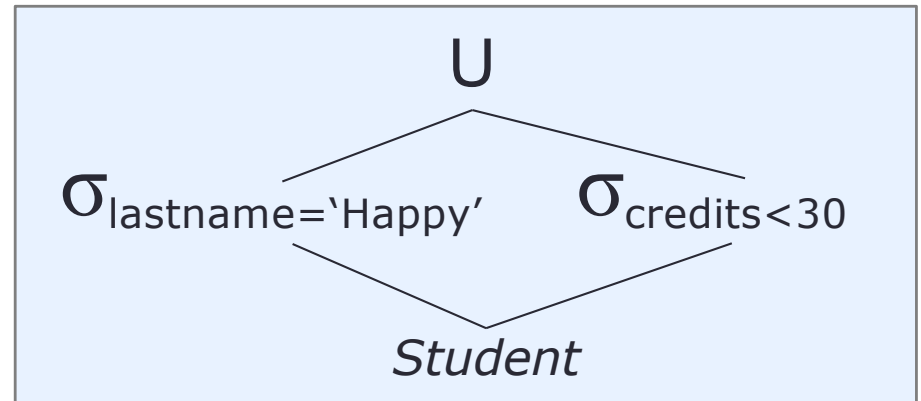
How many tuples do we expect this query to output? – hard to say

- If conditions disjoint, **add** estimates
 - E.g. no student with lastname 'Happy' has credits <30
- If conditions fully overlap, take **maximum** of estimates
 - E.g. all students with lastname 'Happy' have credits <30

Assume disjoint unless you know for sure full overlap

Estimation: OR (4)

```
SELECT *  
FROM Student  
WHERE credits < 30  
OR lastname = 'Happy'
```



How many tuples do we expect this query to output? – hard to say

- If conditions disjoint, **add** estimates

$$\leq 10000 * ((30-1) / (126-1)) + (10000 * 1/9500) \approx 2321 \text{ tuples}$$

- If conditions fully overlap, take **maximum** of estimates

$$\geq 10000 * \max\{ ((30-1) / (126-1)), (1/9500) \} \approx 2320 \text{ tuples}$$

(assume disjoint unless otherwise specified --- answer: 2321 tuples)

Selectivity factor

Let's Try: OR

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

T(Harris_Teeter) = 1,000
V(name) = 900
V(category) = 10
V(price) = 200
Range(price) = [1,50]

#of tuples
#of distinct values
#of distinct values
#of distinct values
range of values

SELECT *
FROM Harris_Teeter
WHERE name='Milk' **OR**
category='meat'

$\sigma_{\text{name='Milk' AND category='meat'}}$
|
Harris_Teeter

Assume: 'Milk' exists in name, 'meat' exists in category, distinct values uniformly distributed, and conditions disjoint

$$\leq \left[T(R) * \frac{1}{V(\text{name})} \right] + \left[T(R) * \frac{1}{V(\text{category})} \right] = \left[1000 * \frac{1}{900} \right] + \left[1000 * \frac{1}{10} \right] \approx 101.11 \text{ tuples}$$

Selectivity factor

Estimation: Cartesian Product

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

T(Student) = 10,000

#of tuples

V(lastName) = 9,500

#of distinct values

V(major) = 10

#of distinct values

Range(credits) = [1, 126)

range of values

Enroll (studId, classNumber, grade)

T(Enroll) = 50,000

#of tuples

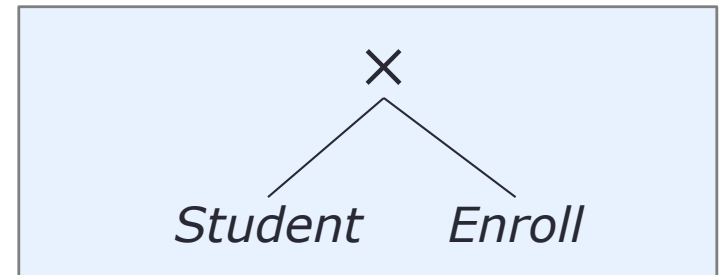
V(studId) = 10,000

#of distinct values

V(classNumber) = 200

#of distinct values

```
SELECT *  
FROM Student, Enroll
```



How many tuples do we expect this query to output?

$T(\text{Student}) * T(\text{Enroll}) = 10000 * 50000$ tuples

No selectivity factor
(because no WHERE clause
applied)

Let's Try: Cartesian Product

Let's go grocery shopping. Assume we know the following info:

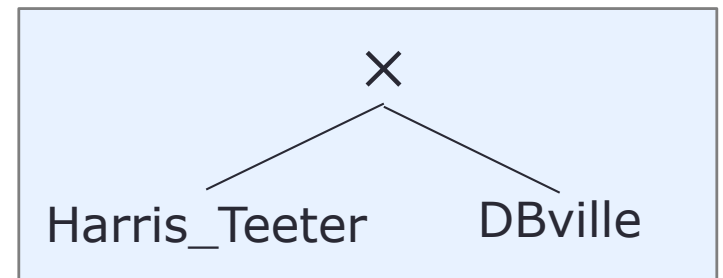
Harris_Teeter(id, name, category, price)

T(Harris_Teeter) = 1,000
V(name) = 900
V(category) = 10
V(price) = 200
Range(price) = [1,50)

DBville(id, dname, shelf, cost)

T(DBville) = 2,000
V(dname) = 1,900
V(shelf) = 12
V(cost) = 500

```
SELECT *  
FROM Harris_Teeter, DBville
```



How many tuples do we expect this query to output?

$T(\text{Harris_Teeter}) * T(\text{DBville}) = 1,000 * 2,000 = 2,000,000$ tuples

No selectivity factor
(because no WHERE clause
applied)

Estimation: JOIN

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

T(Student) = 10,000

#of tuples

V(lastName) = 9,500

#of distinct values

V(major) = 10

#of distinct values

Range(credits) = [1, 126)

range of values

Enroll (studId, classNumber, grade)

T(Enroll) = 50,000

#of tuples

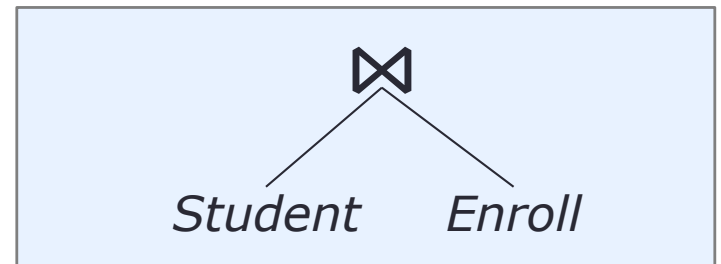
V(studId) = 10,000

#of distinct values

V(classNumber) = 200

#of distinct values

```
SELECT *  
FROM Student  
NATURAL JOIN Enroll
```



How many tuples do we expect this query to output?

$\leq T(\text{Student}) * T(\text{Enroll})$

$\leq 10000 * 50000$ tuples

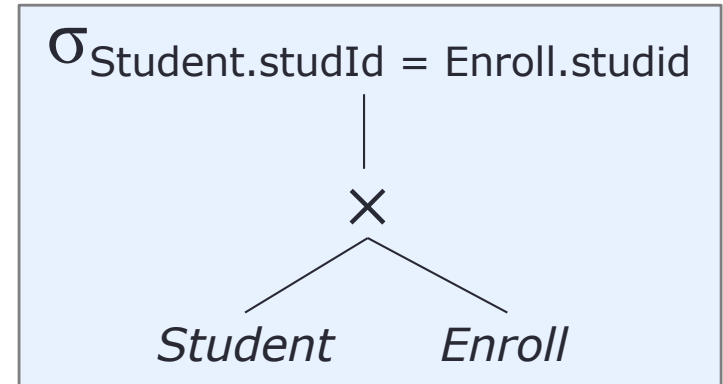
Can we do better?

Estimation: JOIN (2)

1. Start with Cartesian product

```
SELECT *  
FROM Student, Enroll  
WHERE Student.studId = Enroll.studId
```

$$T(\text{Student}) * T(\text{Enroll})$$



2. Suppose there are `studId` in both relations that match
3. How many times does `sid0` occur? (assume `sid0` is `studentId`)
How many tuples do we expect from

$\sigma_{\text{Student.studId}=\text{sid0} \text{ AND } \text{Enroll.studId}=\text{sid0}}$

$$\text{Selectivity factor} = \frac{1}{V(\text{Student}, \text{studId})} * \frac{1}{V(\text{Enroll}, \text{studId})}$$

Estimation: JOIN (3)

4. How many distinct values of `sid0s` exist in the join?

- If no overlap

0

- If full overlap

$$\leq \min\{ V(\text{Student}, \text{studId}), V(\text{Enroll}, \text{studId}) \}$$

Assume full overlap
(~ one is a subset of the other)

5. Multiply (1), (3), and (4)

$$\frac{T(\text{Student}) * T(\text{Enroll})}{V(\text{Student}, \text{studId}) * V(\text{Enroll}, \text{studId})} * \min\{V(\text{Student}, \text{studId}), V(\text{Enroll}, \text{studId})\}$$

Simplify to

$$\frac{T(\text{Student}) * T(\text{Enroll})}{\max\{V(\text{Student}, \text{studId}), V(\text{Enroll}, \text{studId})\}}$$

Estimation: JOIN (4)

Assume we know the following information:

Student (studId, lastName, firstName, major, credits)

T(Student) = 10,000

#of tuples

V(lastName) = 9,500

#of distinct values

V(major) = 10

#of distinct values

Range(credits) = [1, 126)

range of values

Enroll (studId, classNumber, grade)

T(Enroll) = 50,000

#of tuples

V(studId) = 10,000

#of distinct values

V(classNumber) = 200

#of distinct values

$$\frac{T(\text{Student}) * T(\text{Enroll})}{\max\{V(\text{Student}, \text{studId}), V(\text{Enroll}, \text{studId})\}} = \frac{10000 * 50000}{\max\{10000, 10000\}} = 50000 \text{ tuples}$$

Since we assume full overlap of studIds between Student and Enroll, we only need the studIds of the smaller relation

Let's Try: JOIN

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

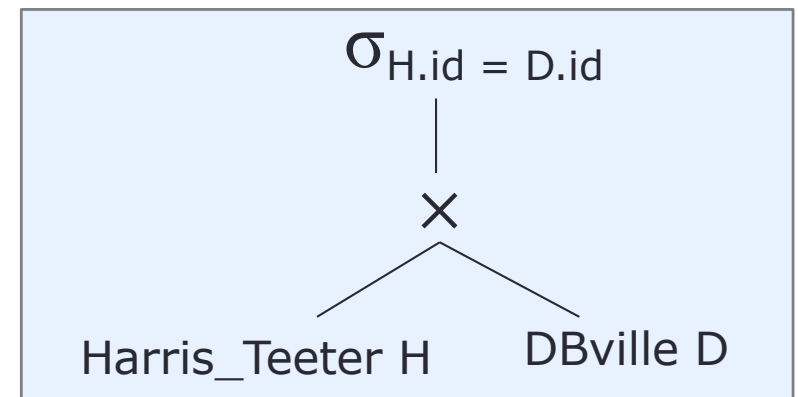
T(Harris_Teeter) = 1,000
V(name) = 900
V(category) = 10
V(price) = 200
Range(price) = [1,50]

DBville(id, dname, shelf, cost)

T(DBville) = 2,000
V(dname) = 1,900
V(shelf) = 12
V(cost) = 500

```
SELECT *  
FROM Harris_Teeter H  
NATURAL JOIN DBville D
```

Assume full overlap of id between the relations,
thus need the ids of the smaller relation



How many tuples do we expect this query to output?

$$\frac{T(\text{Harris_Teeter}) * T(\text{DBville})}{\max\{V(\text{Harris_Teeter}, \text{id}), V(\text{DBville}, \text{id})\}} = \frac{1000 * 2000}{\max\{1000, 2000\}} = 1000 \text{ tuples}$$

Let's Try: JOIN on Attr

Let's go grocery shopping. Assume we know the following info:

Harris_Teeter(id, name, category, price)

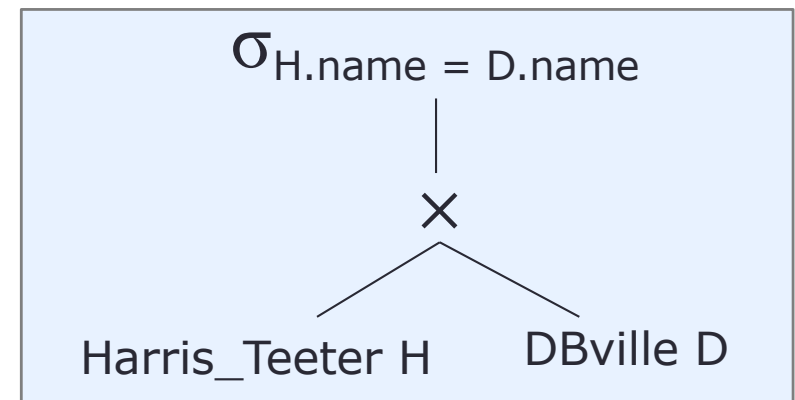
T(Harris_Teeter) = 1,000
V(name) = 900
V(category) = 10
V(price) = 200
Range(price) = [1,50)

DBville(id, dname, shelf, cost)

T(DBville) = 2,000
V(dname) = 1,900
V(shelf) = 12
V(cost) = 500

```
SELECT *  
FROM Harris_Teeter H, Dbville D  
WHERE H.name = D.dname
```

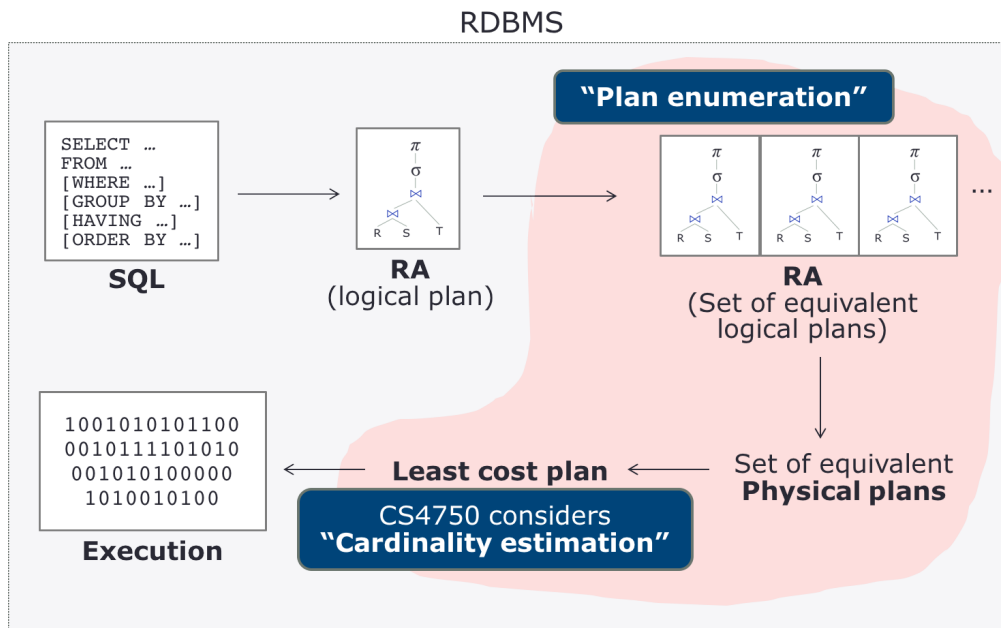
Assume full overlap of name and dname



How many tuples do we expect this query to output?

$$\frac{T(\text{Harris_Teeter}) * T(\text{DBville})}{\max\{V(\text{Harris_Teeter}, \text{name}), V(\text{DBville}, \text{dname})\}} = \frac{1000 * 2000}{\max\{900, 1900\}} \approx 1053 \text{ tuples}$$

Wrap-Up



What's next?

- Indexing

- Cardinality estimation
- Real RDBMS uses sophisticated cost model
- Making inappropriate assumptions to estimate cardinality may lead to:
 - Inaccurate estimates
 - Optimization selects a slow plan
 - Slow query execution
- Be careful and document your assumptions