# Database Systems
# CSE 414

Lectures 9: Relational Algebra
(part 2) and Query Evaluation

(Ch. 5.2 & 16.3 (skim 16.3.2))

# Announcements

- Should have used SQL / Azure now
  - let us know if you had any setup problems

- WQ3 is due on Sunday

- HW3 is due one week from Tuesday

- HW1 grades *should* be posted tonight

# 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 $\theta$

- **Equijoin**: $R \bowtie_\theta S = \sigma_\theta(R \times S)$

  - Join condition $\theta$ consists only of equalities

- **Natural join**: $R \bowtie S = \pi_A(\sigma_\theta(R \times S))$

  - Equijoin

  - Equality on **all** fields with same name in R and in S

  - Projection $\pi_A$ drops all redundant attributes

# So Which Join Is It ?

When we write R ⋈ S we usually mean an equijoin, but we often omit the equality predicate when it is clear from the context

# More Joins

- **Outer join**
  - Include tuples with no matches in the output
  - Use NULL values for missing attributes
  - Does not eliminate duplicate columns

- Variants
  - Left outer join
  - Right outer join
  - Full outer join

# Outer Join Example

AnonPatient P

| age | zip | disease |
|-----|-------|---------|
| 54 | 98125 | heart |
| 20 | 98120 | flu |
| 33 | 98120 | lung |

AnonJob J

| job | age | zip |
|---------|-----|-------|
| lawyer | 54 | 98125 |
| cashier | 20 | 98120 |

$$P \bowtie J$$

| P.age | P.zip | disease | job | J.age | J.zip |
|-------|-------|---------|---------|-------|-------|
| 54 | 98125 | heart | lawyer | 54 | 98125 |
| 20 | 98120 | flu | cashier | 20 | 98120 |
| 33 | 98120 | lung | null | null | null |

# More Examples

```
Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,qty,price)
```
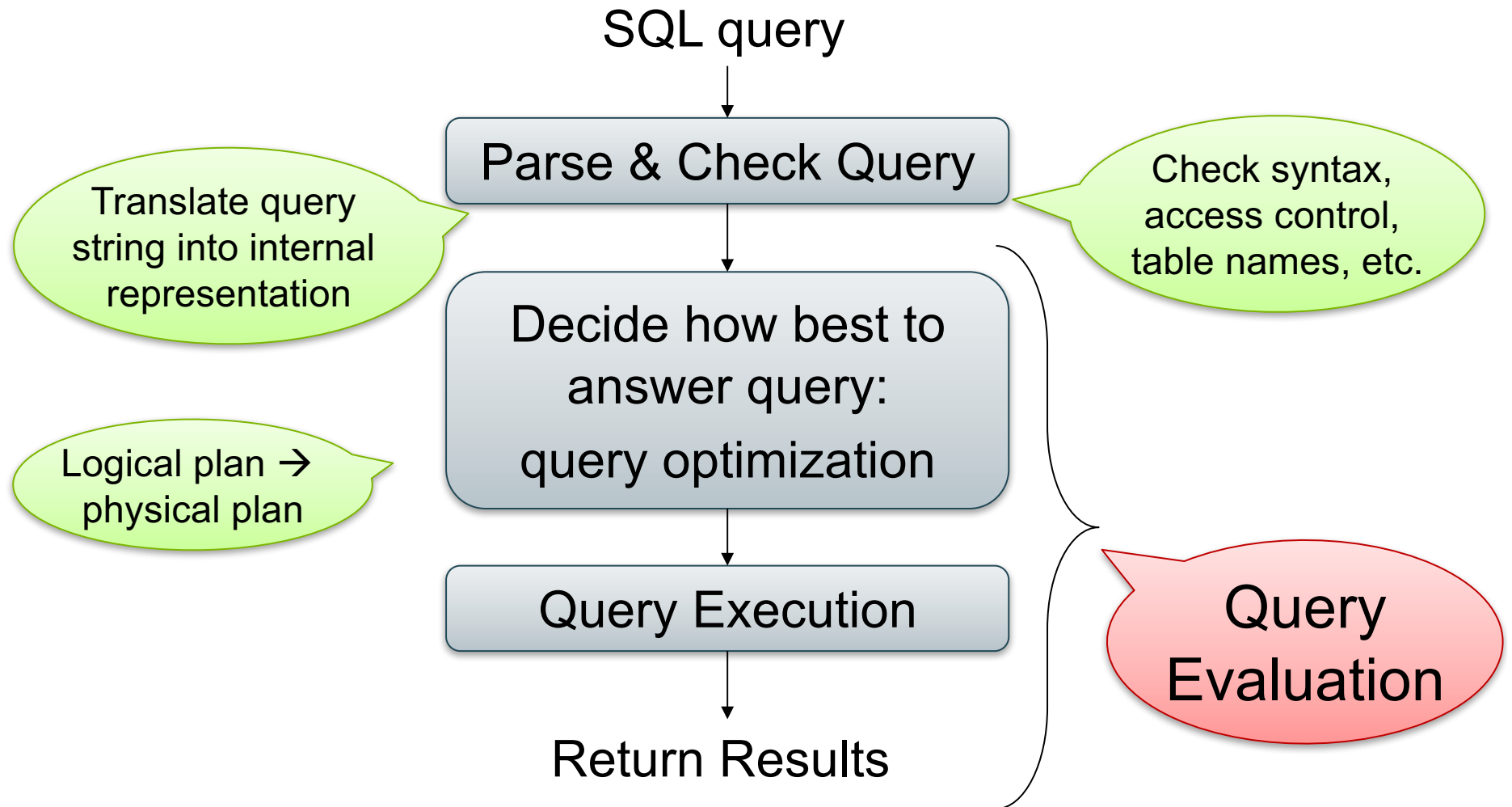
Name of supplier of parts with size greater than 10

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

Name of supplier of red parts or parts with size greater than 10

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

# Query Evaluation Steps

SQL query

↓

**Parse & Check Query**

*Translate query string into internal representation*

*Check syntax, access control, table names, etc.*

↓

**Decide how best to answer query: query optimization**

*Logical plan → physical plan*

↓

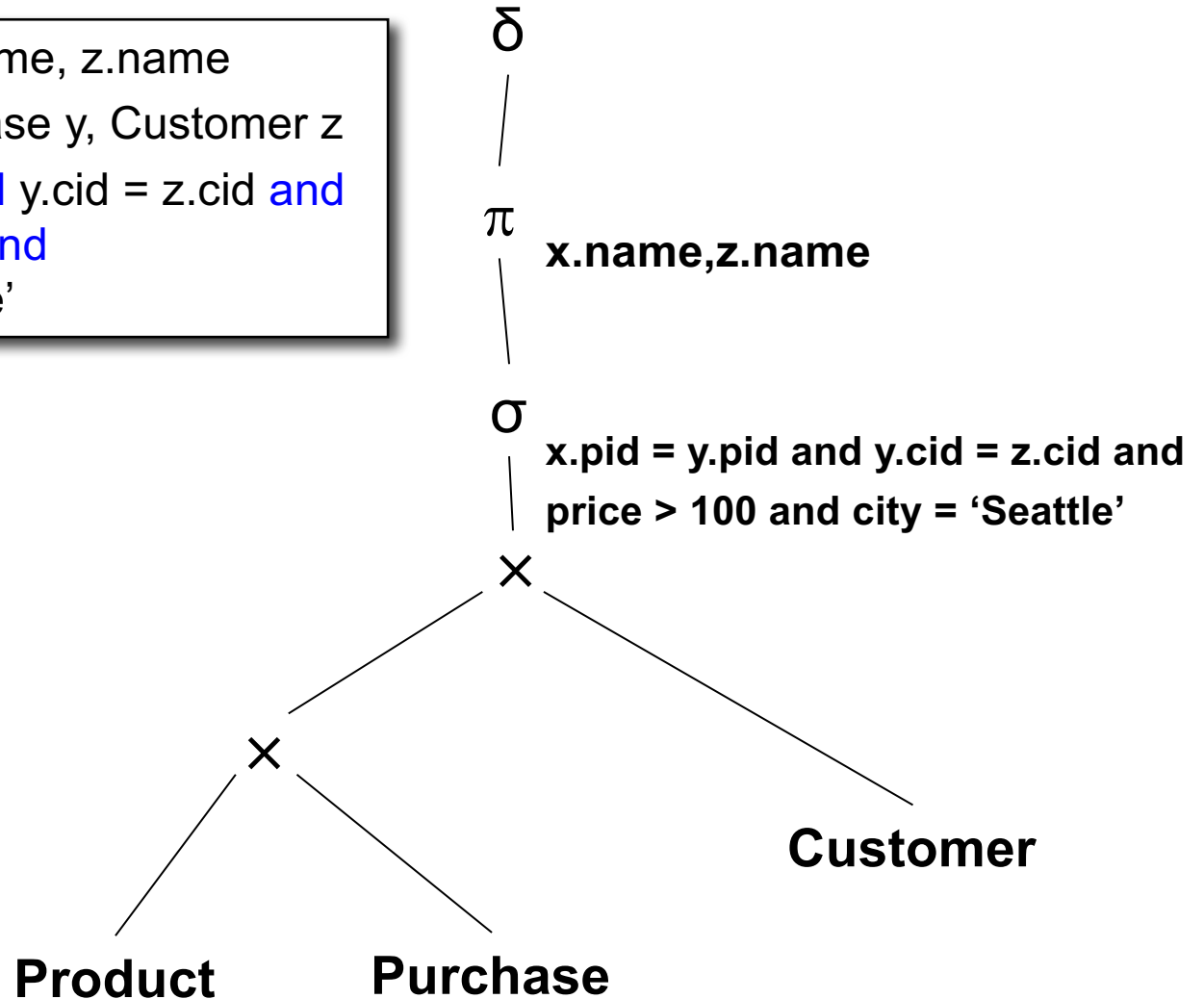**Query Execution**

**Query Evaluation**

↓

Return Results

Product(pid, name, price)
Purchase(pid, cid, store)
Customer(cid, name, city)

# From SQL to RA

SELECT DISTINCT x.name, z.name

FROM Product x, Purchase y, Customer z

WHERE x.pid = y.pid and y.cid = z.cid and
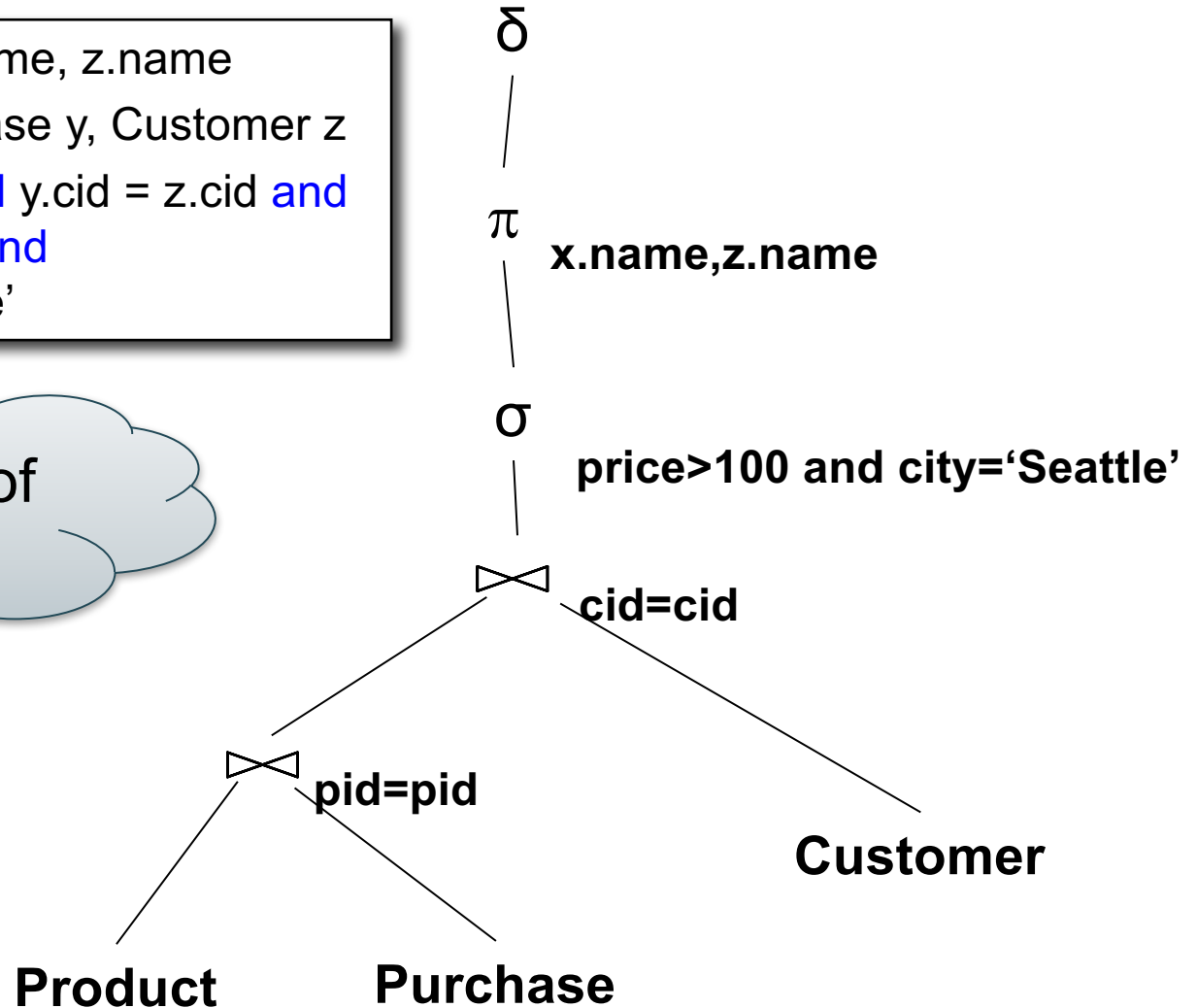       x.price > 100 and
       z.city = 'Seattle'

$\delta$

$\pi$ **x.name,z.name**

$\sigma$ **x.pid = y.pid and y.cid = z.cid and price > 100 and city = 'Seattle'**
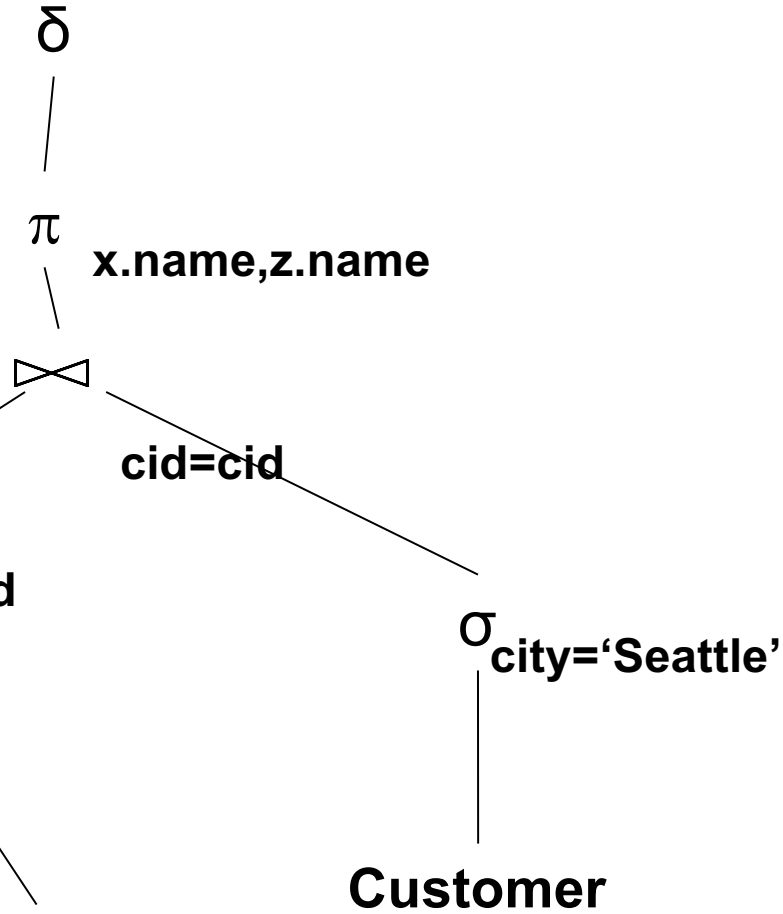
$\times$

$\times$ **Customer**

**Product** **Purchase**

Product(pid, name, price)
Purchase(pid, cid, store)
Customer(cid, name, city)

# From SQL to RA

SELECT DISTINCT x.name, z.name

FROM Product x, Purchase y, Customer z

WHERE x.pid = y.pid and y.cid = z.cid and
    x.price > 100 and
    z.city = 'Seattle'

$\delta$

$\pi$ **x.name,z.name**

$\sigma$ **price>100 and city='Seattle'**

Can you think of another plan?

⋈ **cid=cid**

⋈ **pid=pid**

**Product**      **Purchase**

**Customer**

Product(pid, name, price)
Purchase(pid, cid, store)
Customer(cid, name, city)

# From SQL to RA

SELECT DISTINCT x.name, z.name

FROM Product x, Purchase y, Customer z

WHERE x.pid = y.pid and y.cid = z.cid and
        x.price > 100 and
        z.city =  'Seattle'

δ

π  x.name,z.name

⋈  cid=cid

⋈  pid=pid

σ  price>100

σ  city='Seattle'

Product

Purchase

Customer

Can you think of another plan?

Push selections down the query plan!
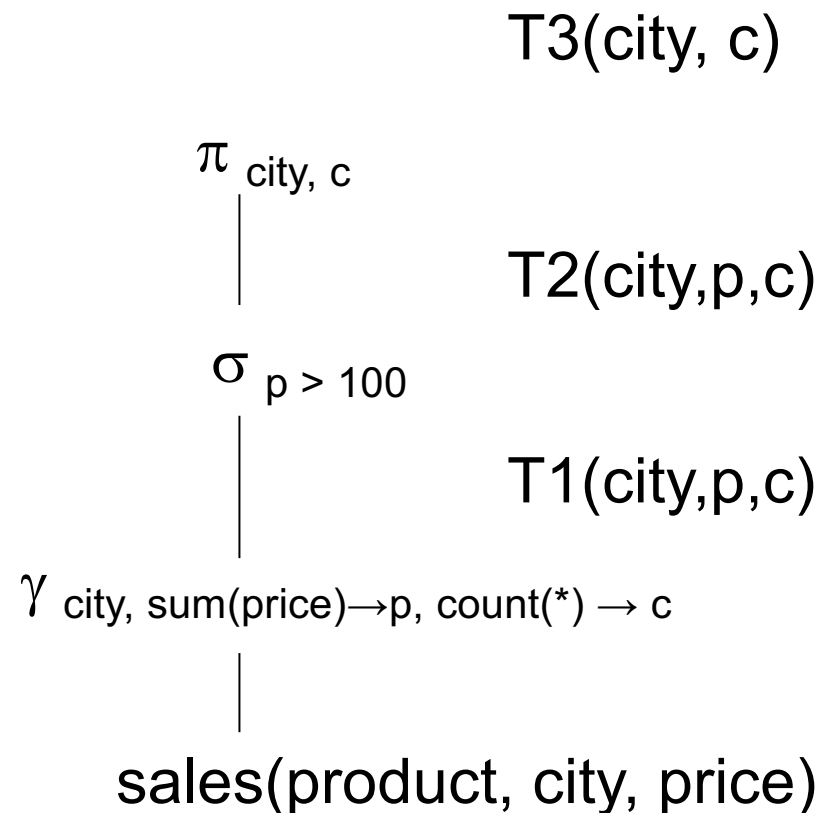
Query optimization: find an equivalent optimal plan

# Extended RA: Operators on Bags

- Duplicate elimination $\delta$
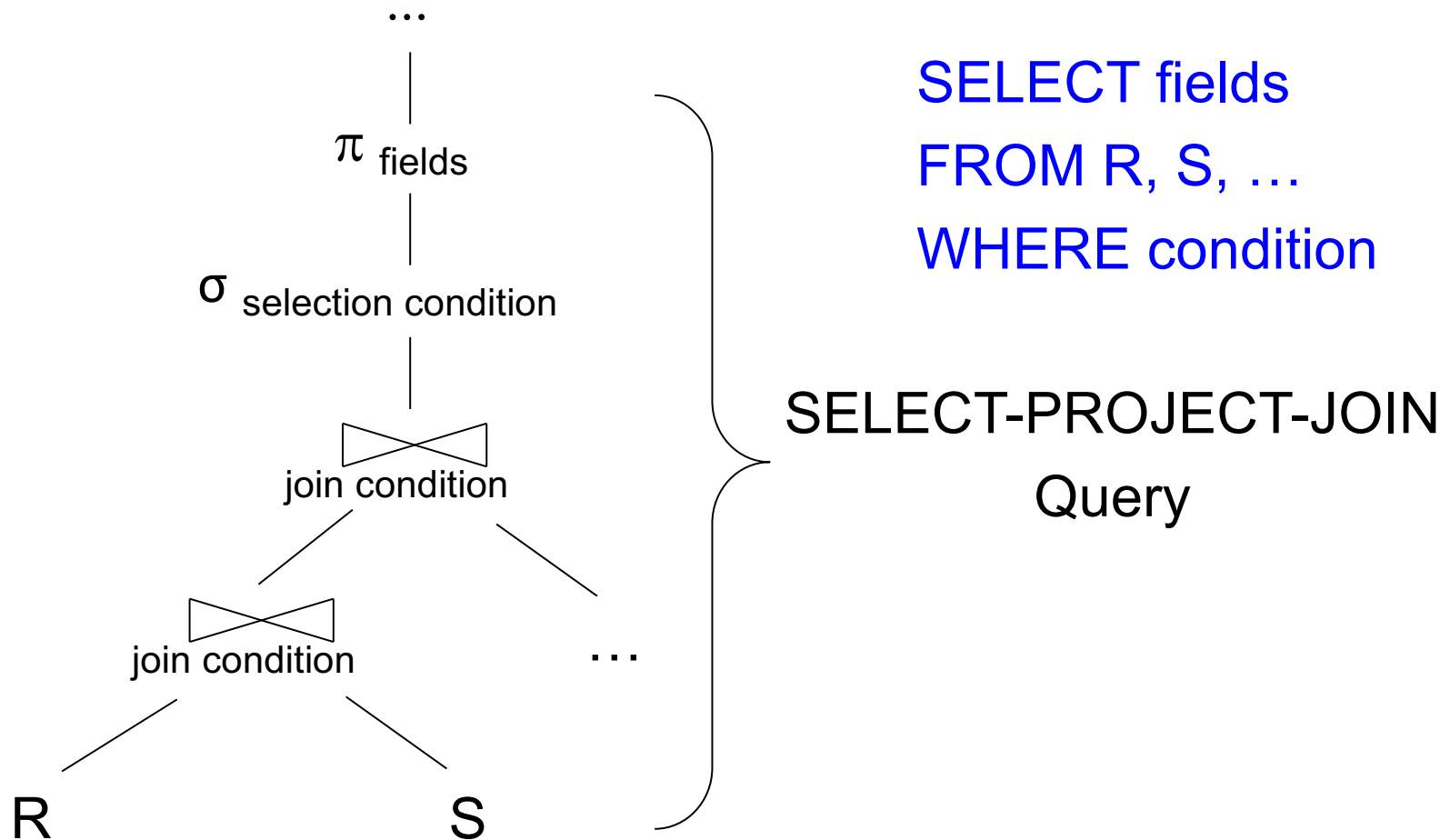- Grouping & aggregation $\gamma$
- Sorting $\tau$

# Logical Query Plan

SELECT city, count(*)
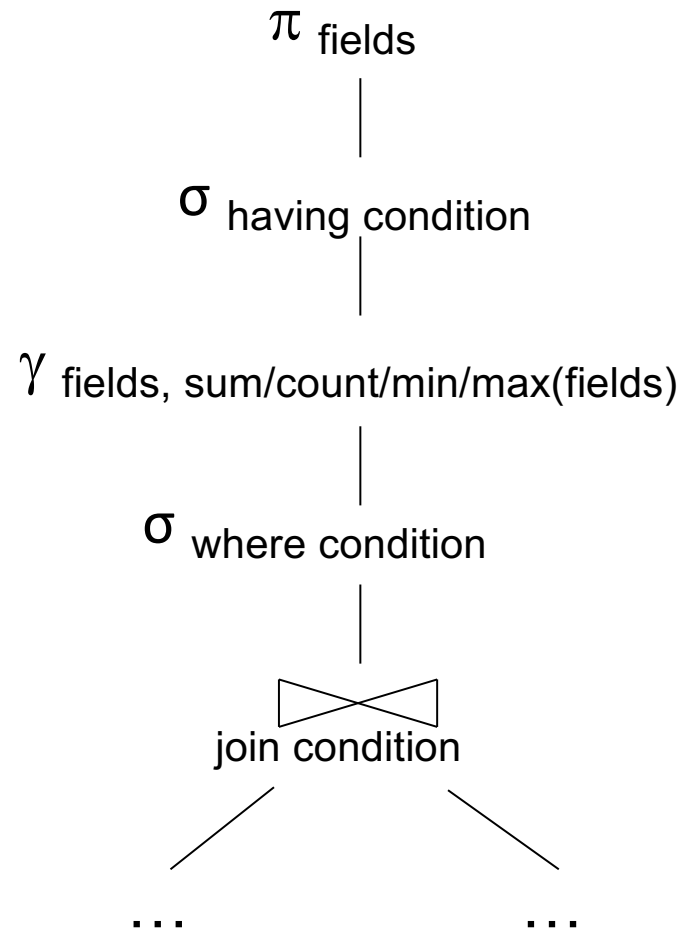FROM sales
GROUP BY city
HAVING sum(price) > 100

T3(city, c)

$\pi_{\text{city, c}}$

T2(city,p,c)

$\sigma_{\text{p > 100}}$

T1(city,p,c)

$\gamma_{\text{city, sum(price)} \rightarrow p, \text{count(*)} \rightarrow c}$

T1, T2, T3 = temporary tables           sales(product, city, price)

# Typical Plan for Block (1/2)

$\cdots$

$\pi_{\text{ fields}}$

$\sigma_{\text{ selection condition}}$

join condition

join condition

R          S

SELECT fields
FROM R, S, …
WHERE condition

SELECT-PROJECT-JOIN
Query

# Typical Plan for Block (2/2)

$\pi$ fields

|

$\sigma$ having condition

|

$\gamma$ fields, sum/count/min/max(fields)

|

$\sigma$ where condition

|

⋈ join condition

…          …

SELECT fields

FROM R, S, …

WHERE condition

GROUP BY fields

HAVING condition

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

# How about Subqueries?

SELECT  Q.sno

FROM Supplier Q

WHERE  Q.sstate = 'WA'
   and not exists

      (SELECT *

       FROM Supply P
       WHERE P.sno = Q.sno
          and P.price > 100)

Correlation !

# How about Subqueries?

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

De-Correlation

```
SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
   and not exists
       (SELECT *
       FROM Supply P
       WHERE P.sno = Q.sno
           and P.price > 100)
```

```
SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
   and Q.sno not in
       (SELECT P.sno
       FROM Supply P
       WHERE P.price > 100)
```

# How about Subqueries?

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

Un-nesting

(SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA')
  EXCEPT
(SELECT P.sno
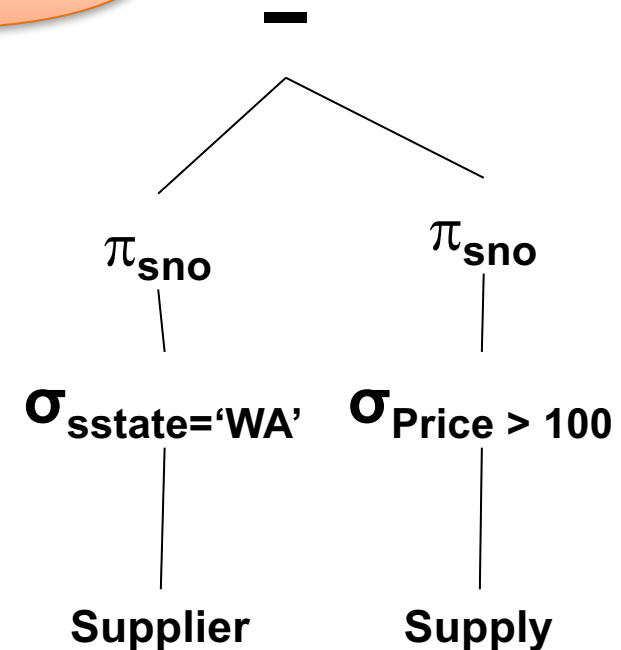FROM Supply P
WHERE P.price > 100)

EXCEPT = set difference

SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
  and Q.sno not in
    (SELECT P.sno
FROM Supply P
WHERE P.price > 100)

# How about Subqueries?

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

(SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA')
 EXCEPT
(SELECT P.sno
 FROM Supply P
 WHERE P.price > 100)

Finally…

$-$

$\pi_{sno}$       $\pi_{sno}$

$\sigma_{sstate='WA'}$    $\sigma_{Price > 100}$

**Supplier**      **Supply**

# From Logical Plans to Physical Plans

# Physical Operators

Each of the logical operators may have one or more implementations = physical operators

Will discuss several basic physical operators, with a focus on join

Product(<u>pid</u>, name, price)
Purchase(<u>pid, cid</u>, store)

# Main Memory Algorithms

Logical operator:

Product(<u>pid</u>, name, price) ⋈ <sub>pid=pid</sub> Purchase(<u>pid, cid</u>, store)

Propose three physical operators for the join, assuming the tables are in main memory:

1. Nested Loop Join          O( ?? )
2. Merge join                O( ?? )
3. Hash join                 O( ?? )

(note that pid is a key)

# Main Memory Algorithms

Logical operator:

Product(pid, name, price) ⋈ pid=pid Purchase(pid, cid, store)

Propose three physical operators for the join, assuming the tables are in main memory:

1. Nested Loop Join      $O(n^2)$ — two nested loops
2. Merge join      O( ?? )
3. Hash join      O( ?? )

# Main Memory Algorithms

Logical operator:

Product(pid, name, price) ⋈ _pid=pid_ Purchase(pid, cid, store)

Propose three physical operators for the join, assuming the tables are in main memory:

1. Nested Loop Join $O(n^2)$
2. Merge join $O(n \log n)$
3. Hash join $O(\ ??\ )$

sort both – $O(n \log n)$
merge – $O(n)$

# Main Memory Algorithms

Logical operator:

Product(<u>pid</u>, name, price) ⋈ <sub>pid=pid</sub> Purchase(<u>pid, cid</u>, store)

Propose three physical operators for the join, assuming the tables are in main memory:

1. Nested Loop Join $O(n^2)$
2. Merge join $O(n \log n)$
3. Hash join $O(n) \ldots O(n^2)$

> add n to hash – O(n)?
> lookup n in hash – O(n)?
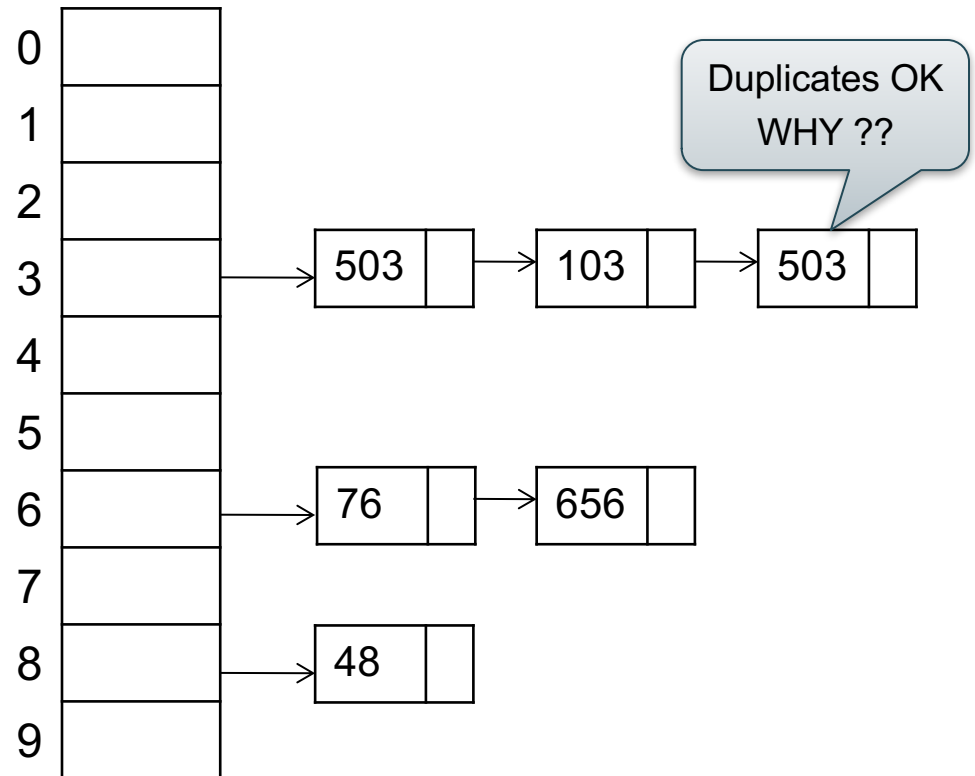
# BRIEF Review of Hash Tables

Separate chaining:

A (naïve) hash function:

$$h(x) = x \bmod 10$$

Operations:

find(103) = ??

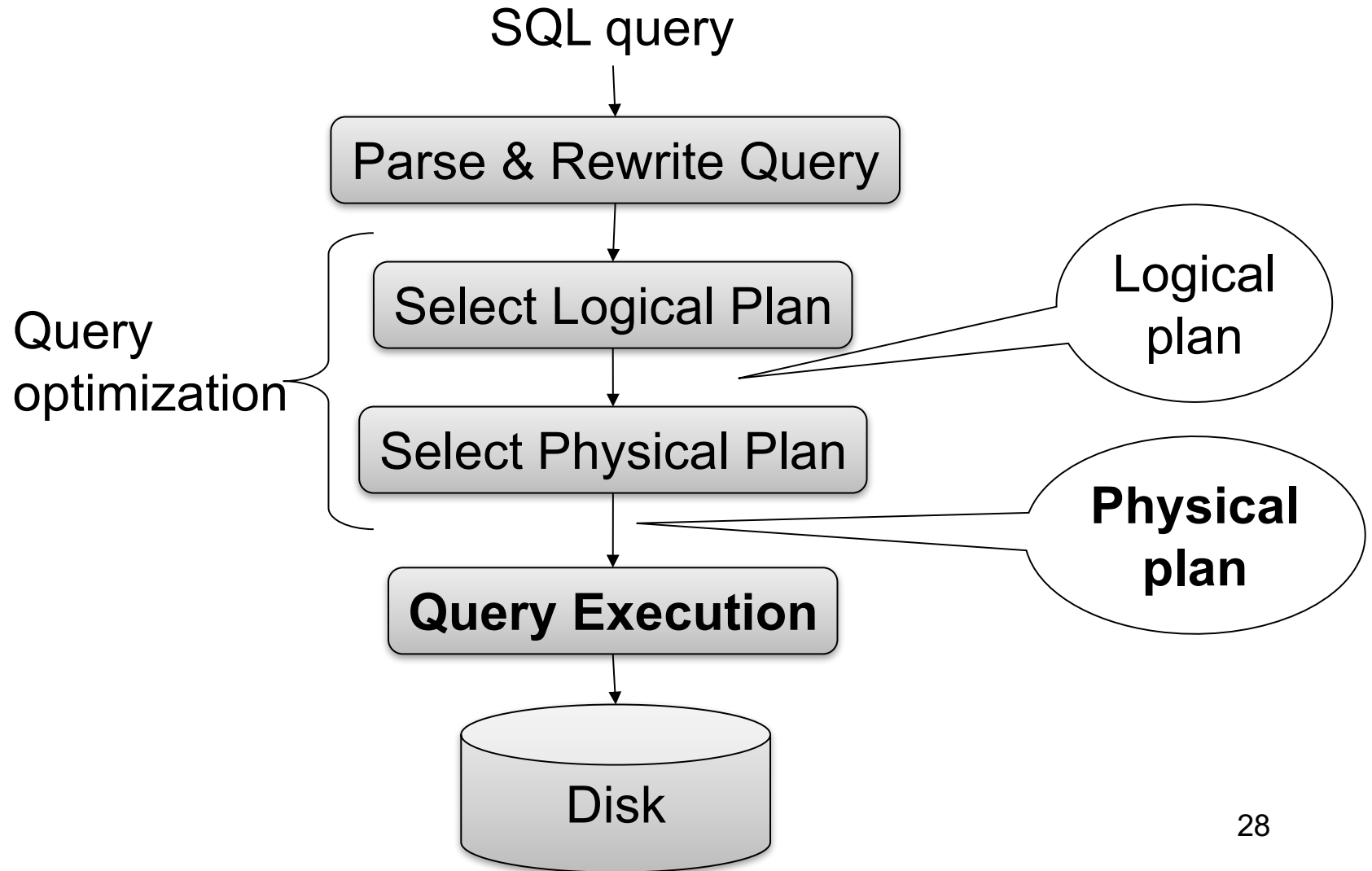insert(488) = ??

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | → 503 → 103 → 503 |
| 4 | |
| 5 | |
| 6 | → 76 → 656 |
| 7 | |
| 8 | → 48 |
| 9 | |

Duplicates OK
WHY ??

# BRIEF Review of Hash Tables

- insert(k, v) = inserts a key k with value v

- Many values for one key
  - Hence, duplicate k's are OK

- find(k) = returns the ***list*** of all values v associated to the key k

# Query Evaluation Steps Review

SQL query

↓

Parse & Rewrite Query

↓

Query optimization

Select Logical Plan → Logical plan

↓

Select Physical Plan

↓

Physical plan

Query Execution

↓

Disk

28

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)

# Relational Algebra

SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
    and  y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'

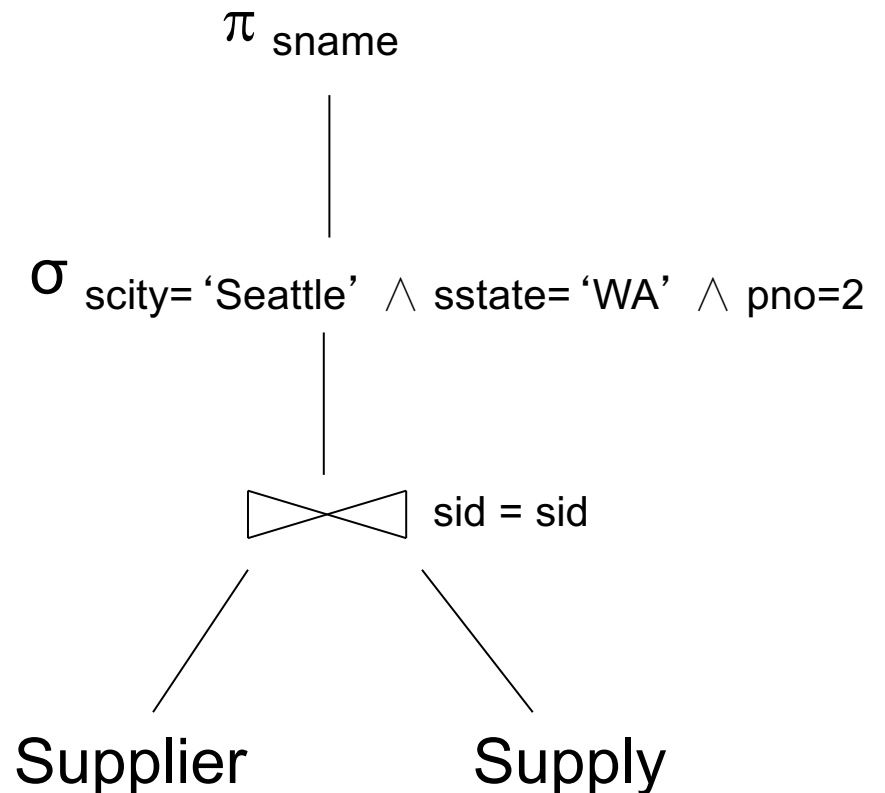Give a relational algebra expression for this query

Supplier(<u>sid</u>, sname, scity, sstate)

Supply(<u>sid, pno</u>, quantity)

# Relational Algebra

SELECT sname

FROM Supplier x, Supply y

WHERE x.sid = y.sid
   and  y.pno = 2
   and x.scity = 'Seattle'
   and x.sstate = 'WA'

$\pi_{sname}(\sigma_{scity= 'Seattle' \wedge sstate= 'WA' \wedge pno=2} (Supplier \bowtie_{sid = sid} Supply))$

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)

# Relational Algebra

SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
    and  y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'

$\pi_{sname}$

$\sigma_{scity=\text{'Seattle'} \land sstate=\text{'WA'} \land pno=2}$

⋈ sid = sid

Supplier      Supply

Relational algebra expression is also called the "logical query plan"

Supplier(<u>sid</u>, sname, scity, sstate)

Supply(<u>sid, pno</u>, quantity)

# Physical Query Plan 1

(On the fly)    $\pi$ sname

(On the fly)

$\sigma$ scity= 'Seattle' $\wedge$ sstate= 'WA' $\wedge$ pno=2

A physical query plan is a logical query plan annotated with physical implementation details

(Nested loop)

⋈
sid = sid

SELECT sname

FROM Supplier x, Supply y

WHERE x.sid = y.sid
    and  y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'

Supplier
(File scan)

Supply
(File scan)

Supplier(<u>sid</u>, sname, scity, sstate)
Supply(<u>sid, pno</u>, quantity)

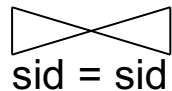# Physical Query Plan 2

(On the fly)        $\pi$ sname

(On the fly)

$\sigma$ scity= 'Seattle' $\wedge$ sstate= 'WA' $\wedge$ pno=2

(Hash join)

sid = sid

Supplier
(File scan)

Supply
(File scan)

Same logical query plan
Different physical plan

SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
    and  y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'

Supplier(sid, sname, scity, sstate)
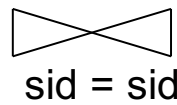
Supply(sid, pno, quantity)

# Physical Query Plan 3

(On the fly)

$\pi$ sname

SELECT sname

FROM Supplier x, Supply y

WHERE x.sid = y.sid

   and y.pno = 2

   and x.scity = 'Seattle'

   and x.sstate = 'WA'

(Sort-merge join)

sid = sid

(Scan & write to T1)

(Scan & write to T2)

$\sigma$ scity= 'Seattle' $\wedge$ sstate= 'WA'

$\sigma$ pno=2

Supplier
(File scan)

Supply
(File scan)

# Query Optimization Problem

- For each SQL query… many logical plans

- For each logical plan… many physical plans

- How do find a fast physical plan?
  - Will discuss in a few lectures