# Database Systems
# CSE 414

## Lectures 14: Relational Algebra (part 2) and Query Evaluation
## (Ch. 5.2 & 16.3 (skim 16.3.2))

# 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 ⋈ 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)
```
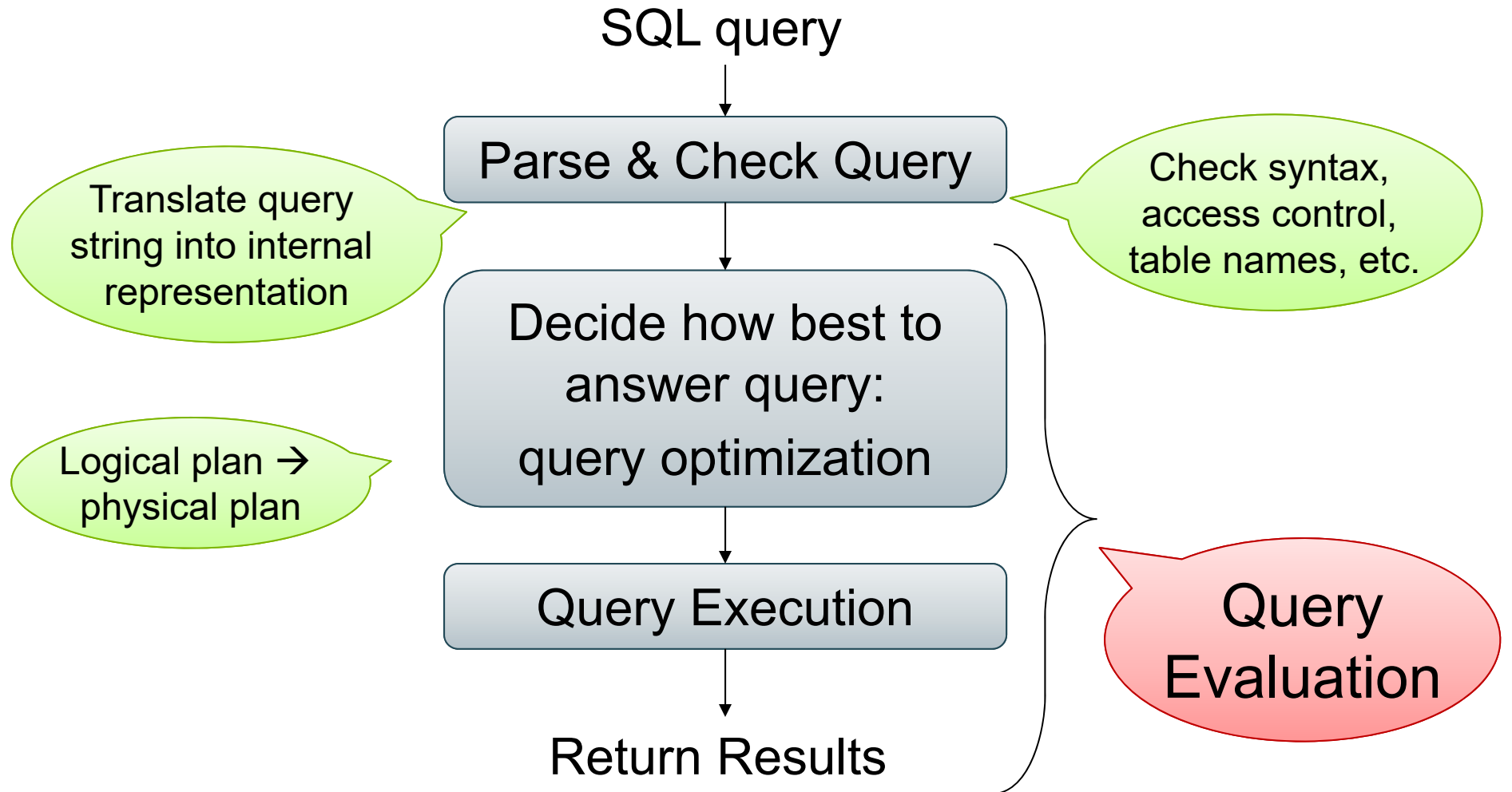
Name of supplier of parts with size greater than 10

$\pi_{sname}$(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
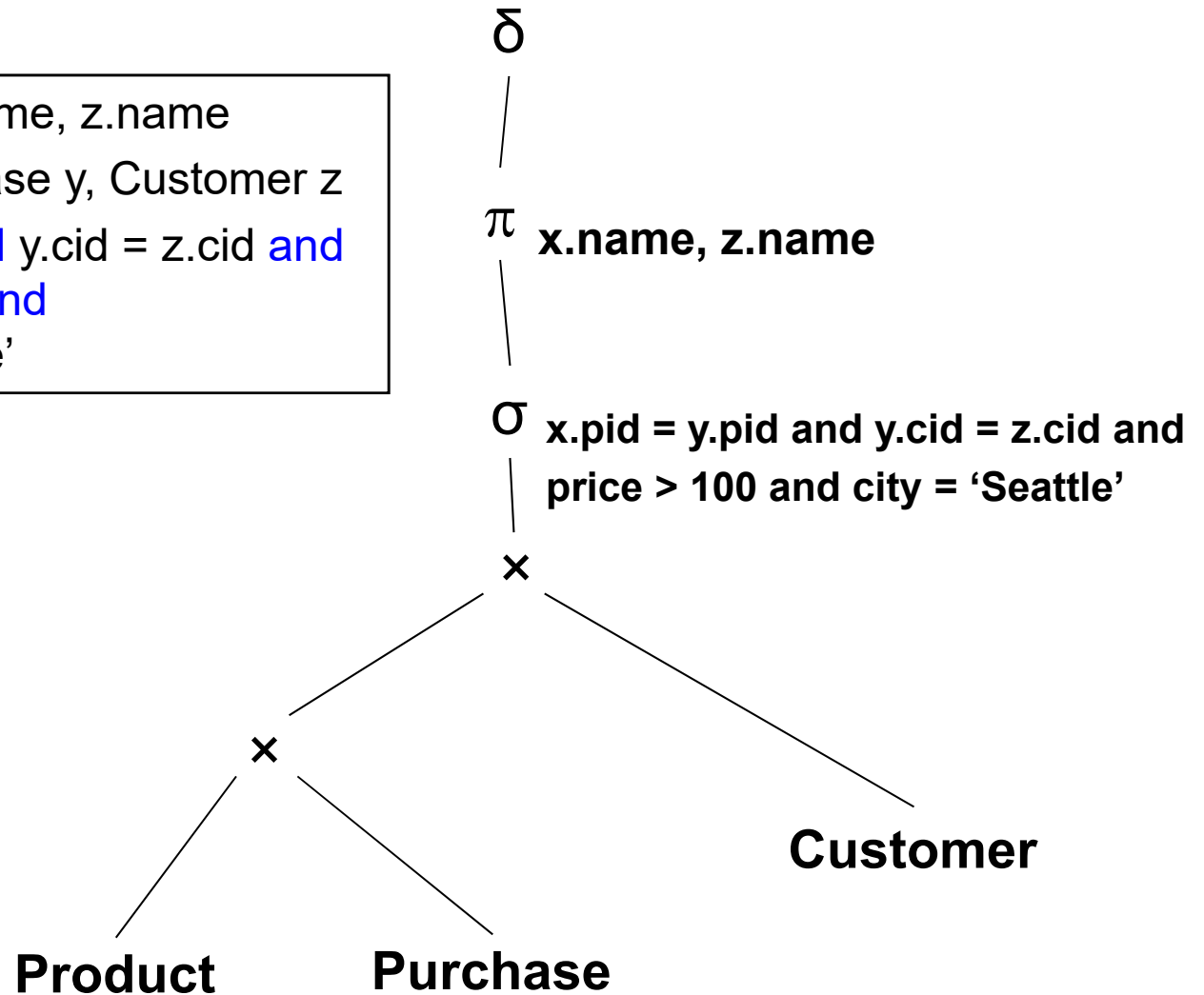
↓

Query Execution

↓

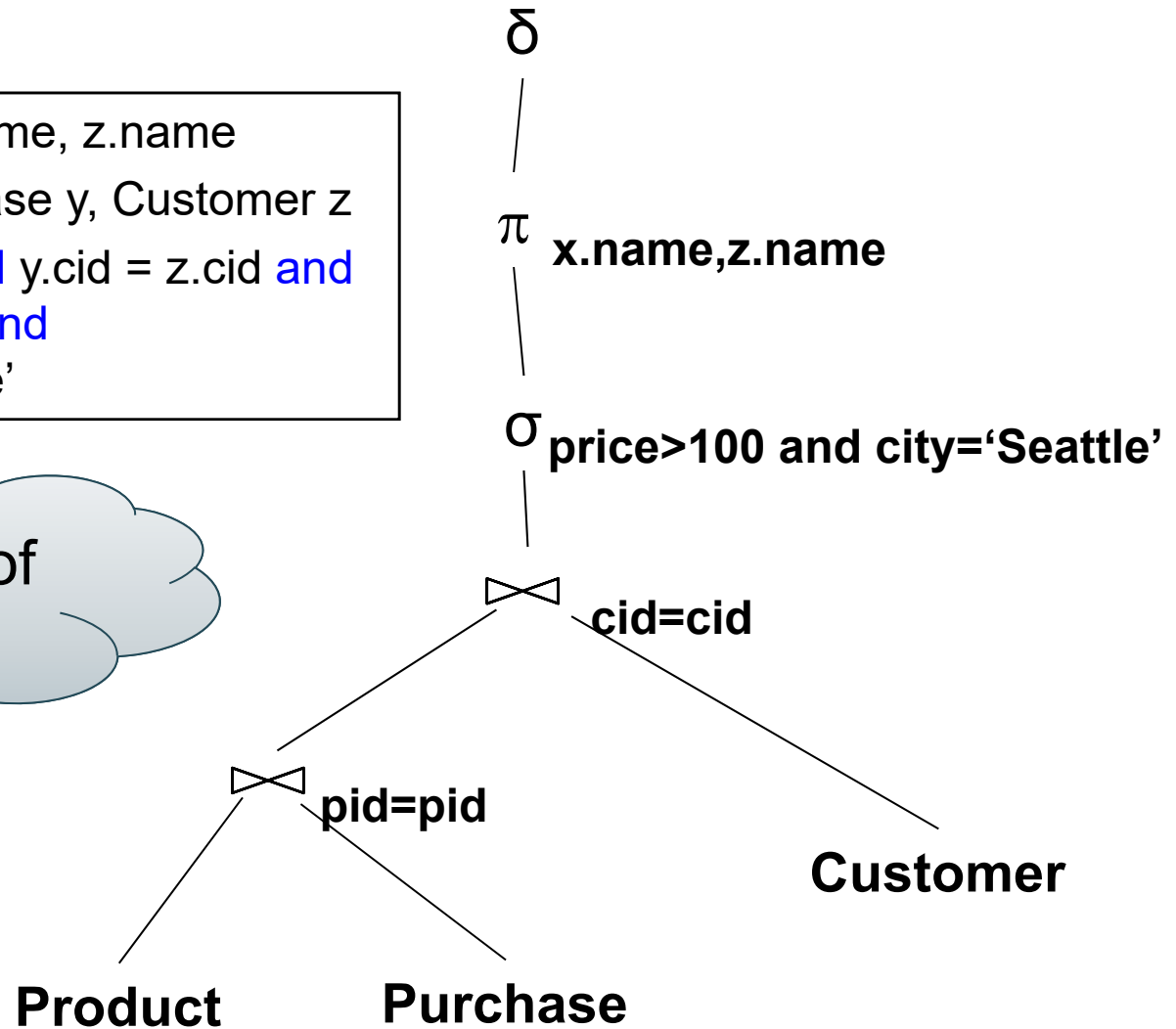Return Results

Query Evaluation

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**

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

×

×          **Customer**

**Product**      **Purchase**

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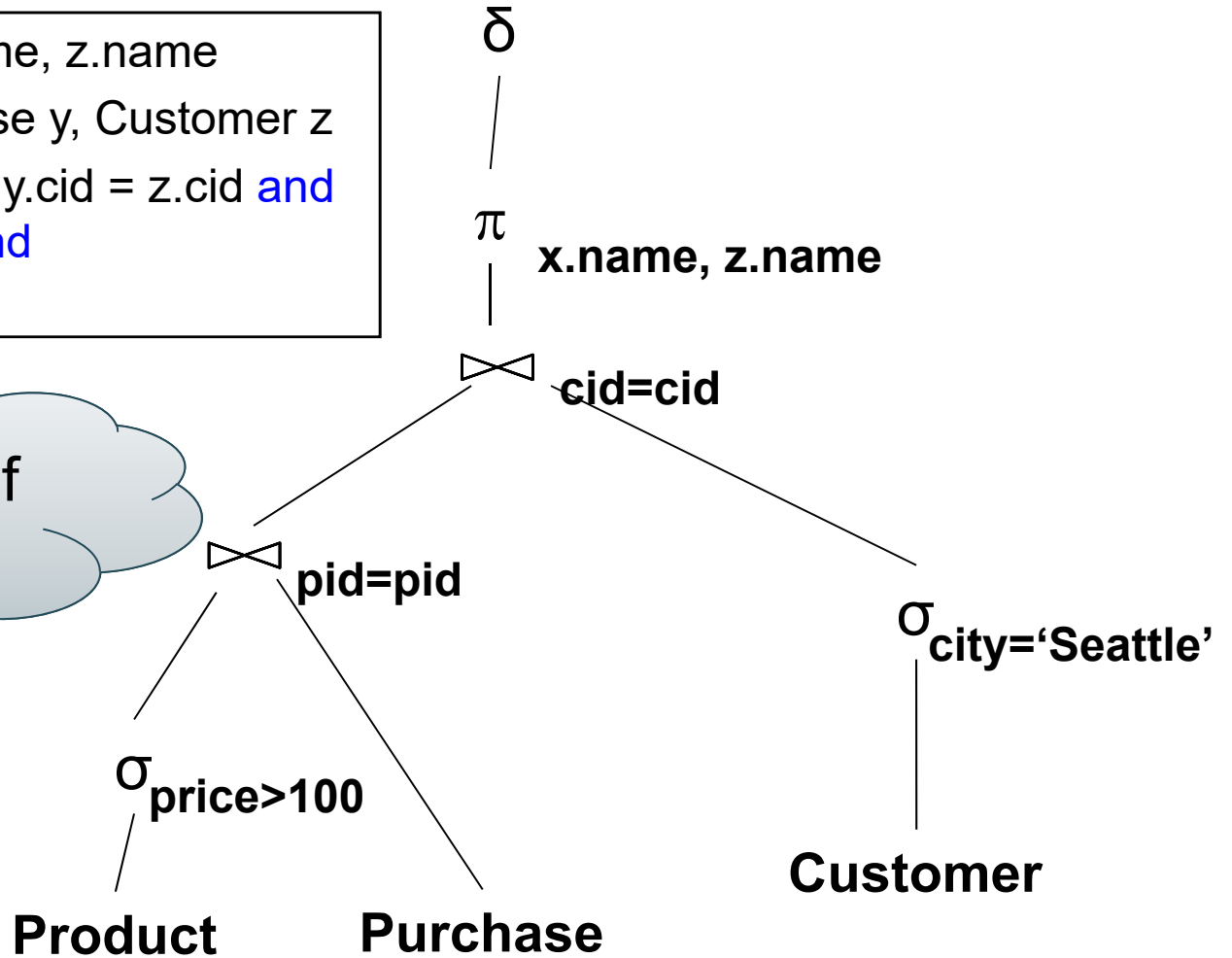Customer(cid, name, city)

# From SQL to RA

$\delta$

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'

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

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

Can you think of another plan?

⋈ **cid=cid**

⋈ **pid=pid**

**Customer**

**Product**   **Purchase**

CSE 414 - Fall 2017

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**

**Product**       **Purchase**

σ **city='Seattle'**

**Customer**

Can you think of
another plan?

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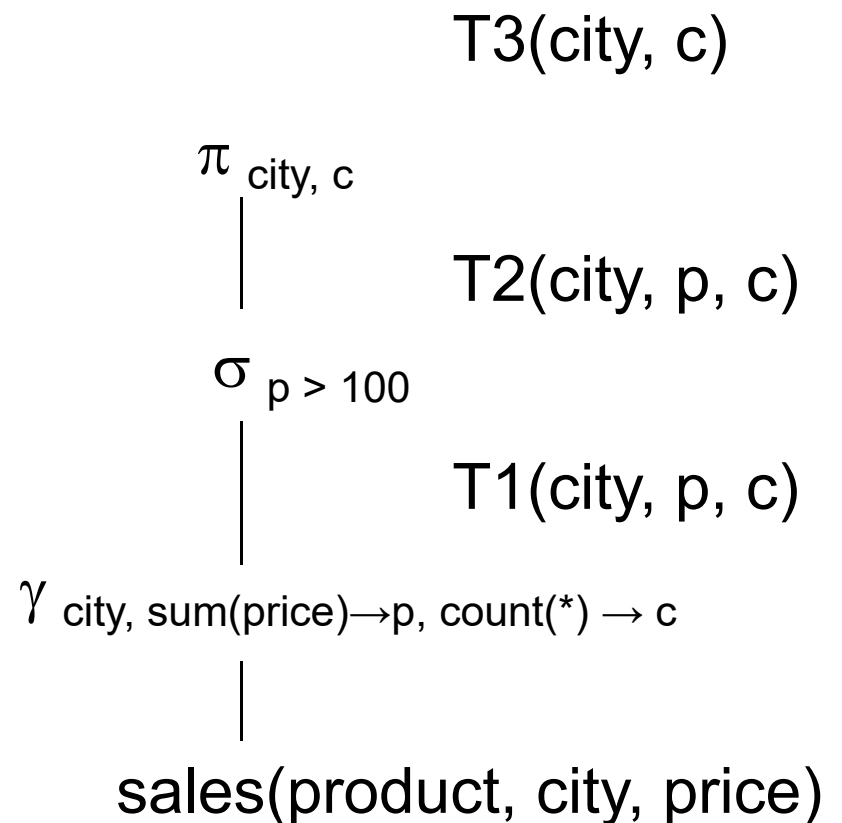
Query optimization: find
an equivalent optimal plan

CSE 414 - Fall 2017                                                    10

# Extended RA: Operators on Bags

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

# Logical Query Plan

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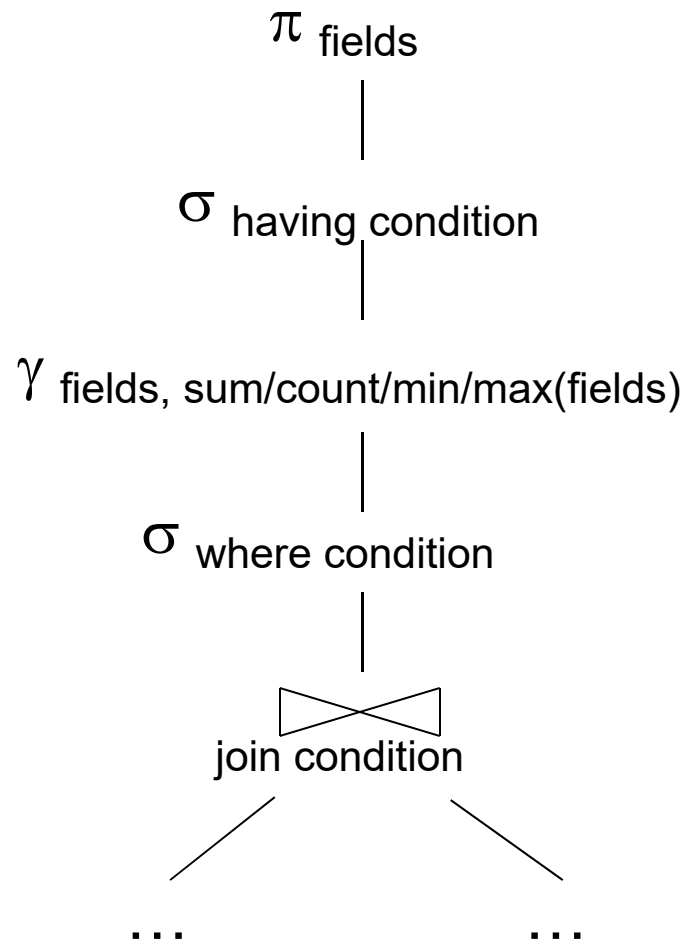GROUP BY city
HAVING sum(price) > 100

T3(city, c)

$\pi_{city, c}$

T2(city, p, c)

$\sigma_{p > 100}$

T1(city, p, c)

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

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

# Typical Plan for Block (1/2)

...

$\pi$ fields

$\sigma$ 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)

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)

De-Correlation

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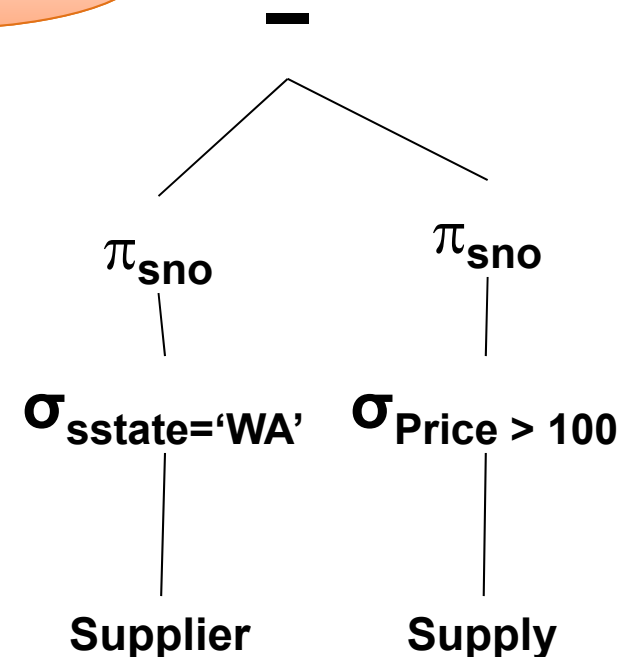        WHERE P.price > 100)

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

# How about Subqueries?

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)

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')
 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(pid, name, price)
Purchase(pid, cid, store)

# 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( ?? )
2.  Merge join                O( ?? )
3.  Hash join                 O( ?? )

(note that pid is a key)

# Main Memory Algorithms

Logical operator:

Product(pid, name, price) $\bowtie$ 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( ?? )

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

# 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 – $O(n)$

CSE 414 - Fall 2017                                      23

# 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                 $O(n) \dots O(n^2)$

add n to hash – O(n)?
lookup n in hash – O(n)?

# BRIEF Review of Hash Tables
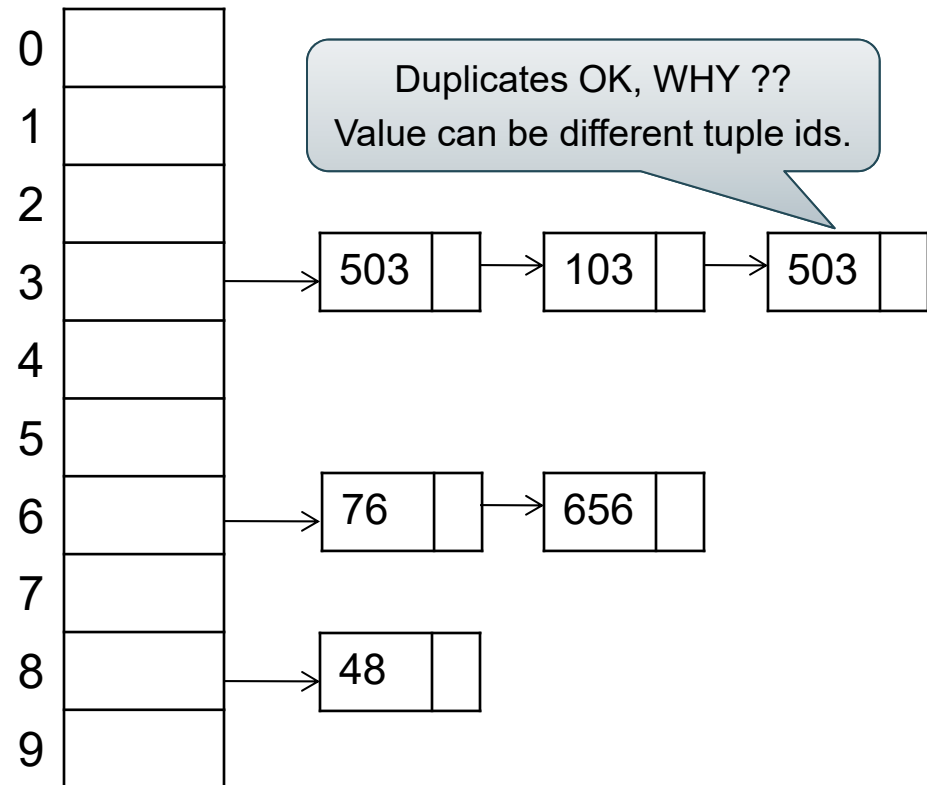
Separate chaining:

A (naïve) hash function:
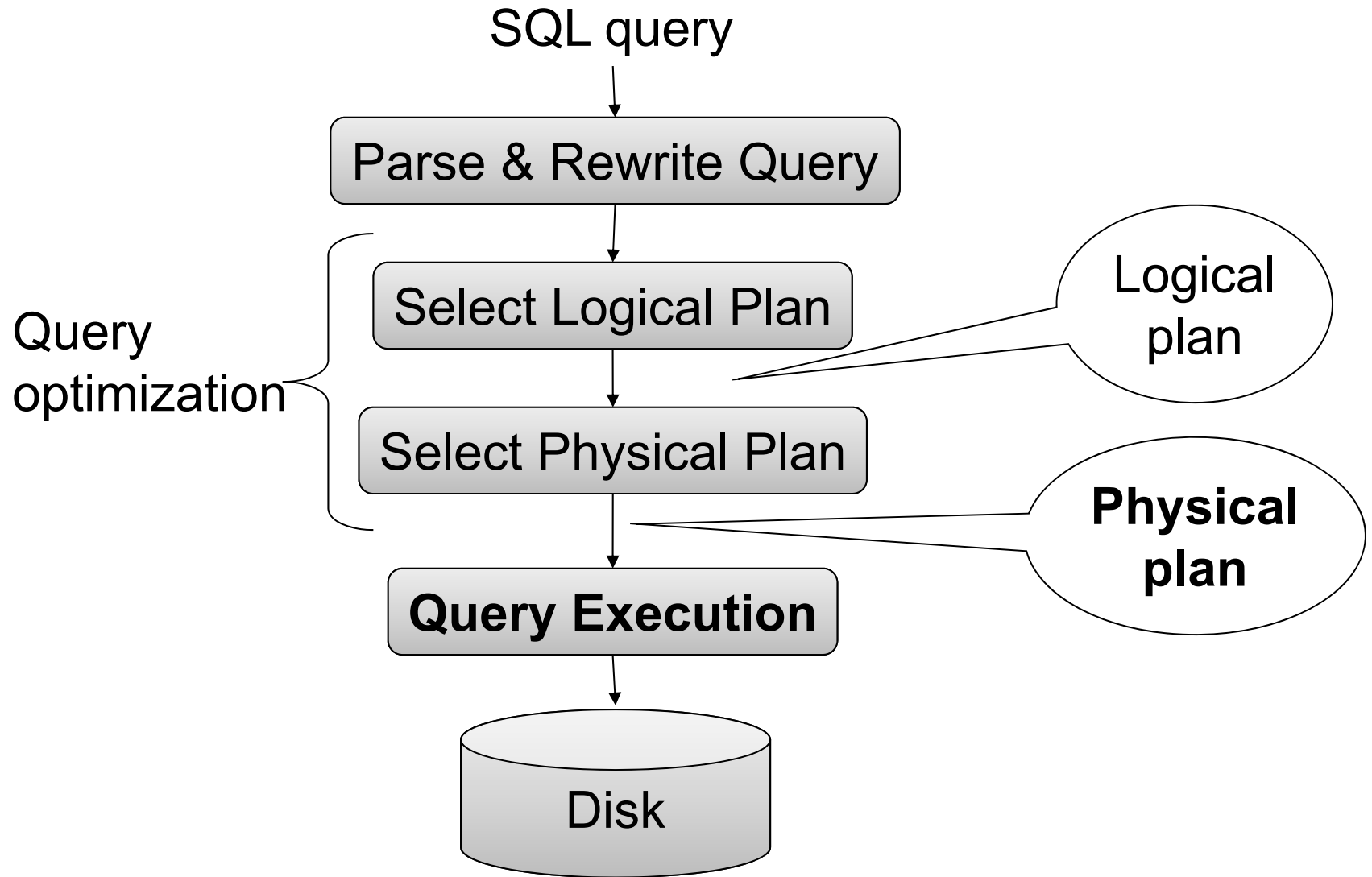
$$h(x) = x \bmod 10$$

Operations:

find(103) = ??

insert(488) = ??

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

Duplicates OK, WHY ??
Value can be different tuple ids.

# 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

↓

Select Logical Plan — Logical plan

↓

Select Physical Plan — **Physical plan**

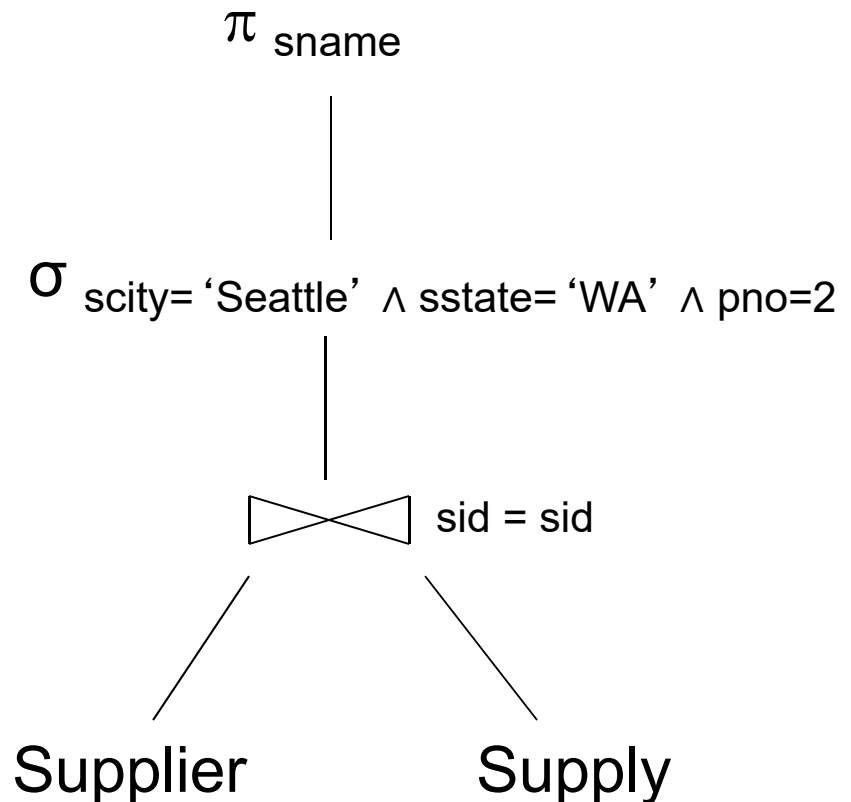Query optimization

↓

**Query Execution**

↓

Disk

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'

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

$\pi_{sname}$

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

$\bowtie$ sid = sid

Supplier        Supply

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

# Physical Query Plan 1

(On the fly)    $\pi_{sname}$

(On the fly)

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

(Nested loop)

⋈
sid = sid

Supplier
(File scan)

Supply
(File scan)

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

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)
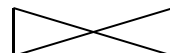
Supply(sid, pno, quantity)

# Physical Query Plan 2

(On the fly)  $\pi$ sname

Same logical query plan
Different physical plan

(On the fly)

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

(Hash join)

sid = sid

Supplier
(File scan)

Supply
(File scan)

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

Different but equivalent logical query plan; different physical plan

SELECT sname

FROM Supplier x, Supply y

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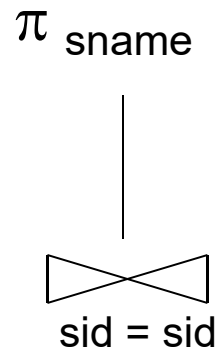  and  y.pno = 2
  and x.scity = 'Seattle'
  and x.sstate = 'WA'

(On the fly)  $\pi$ sname

(Sort-merge join)  $\bowtie$
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 to find a fast physical plan?
    - Will discuss in a few lectures