# Introduction to Data Management CSE 344
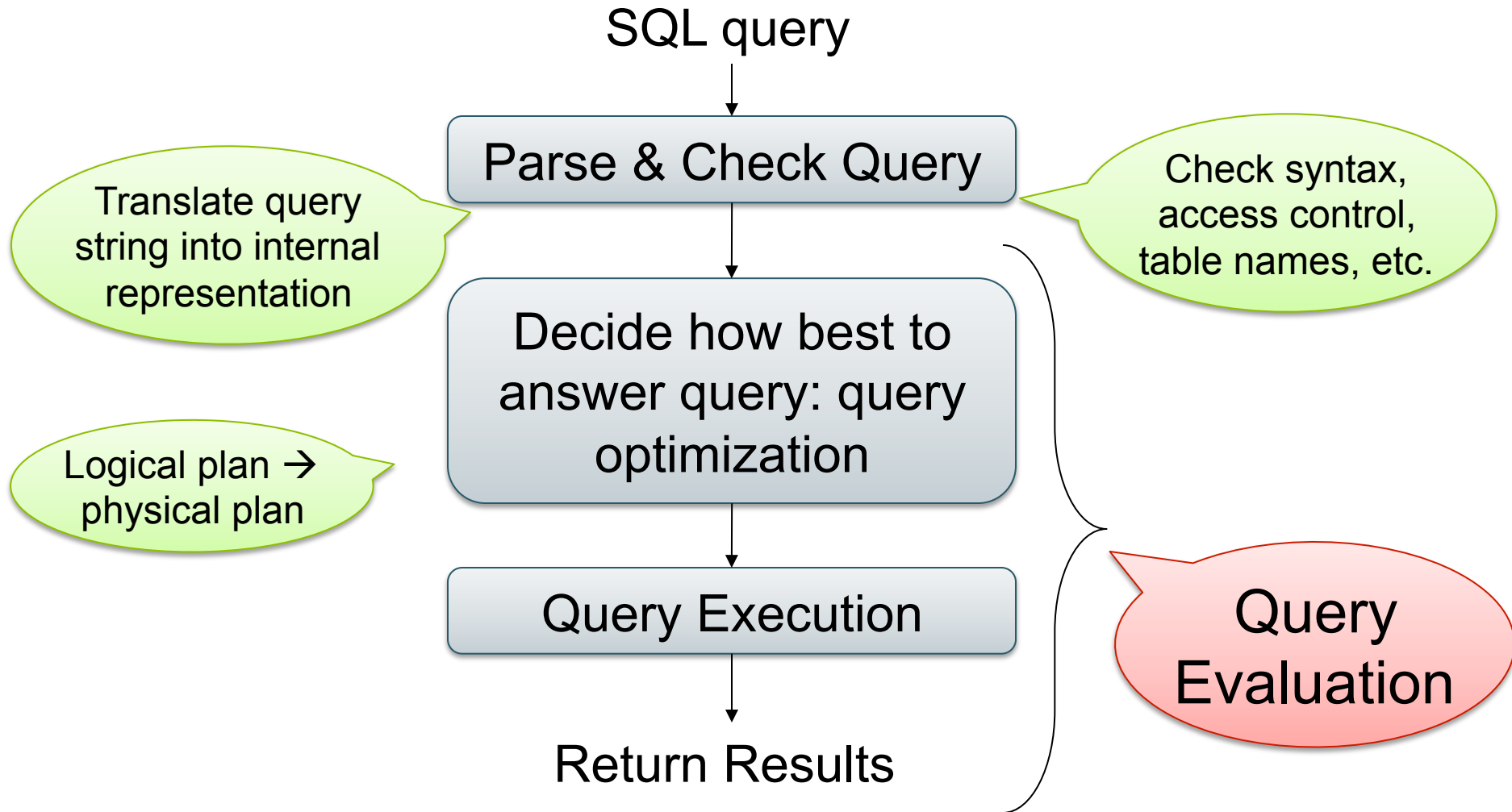
## Lectures 9: Relational Algebra (part 2) and Query Evaluation

Guest lecturer: Laurel Orr

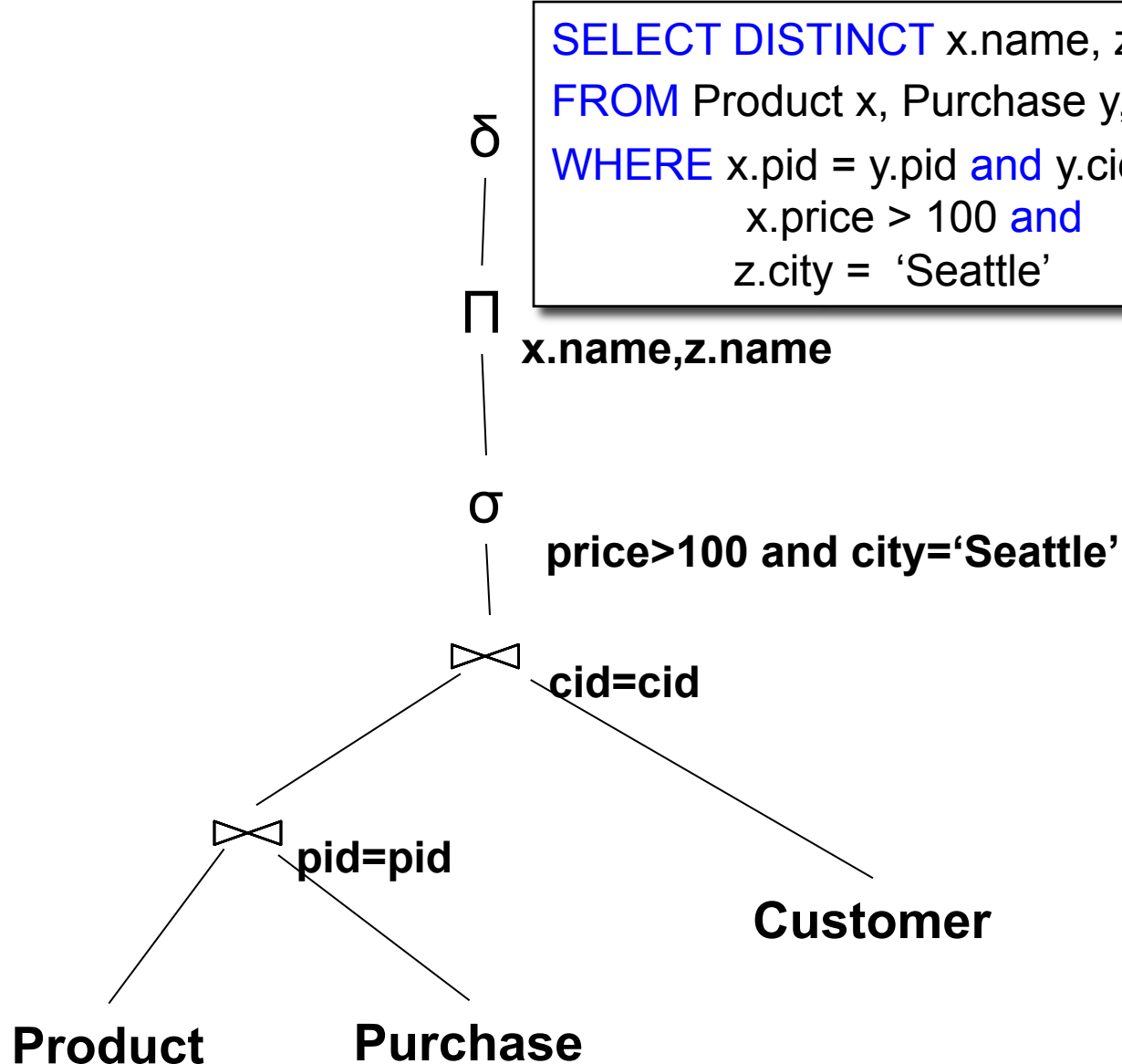# Announcements

- HW3 is due next Tuesday

# 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
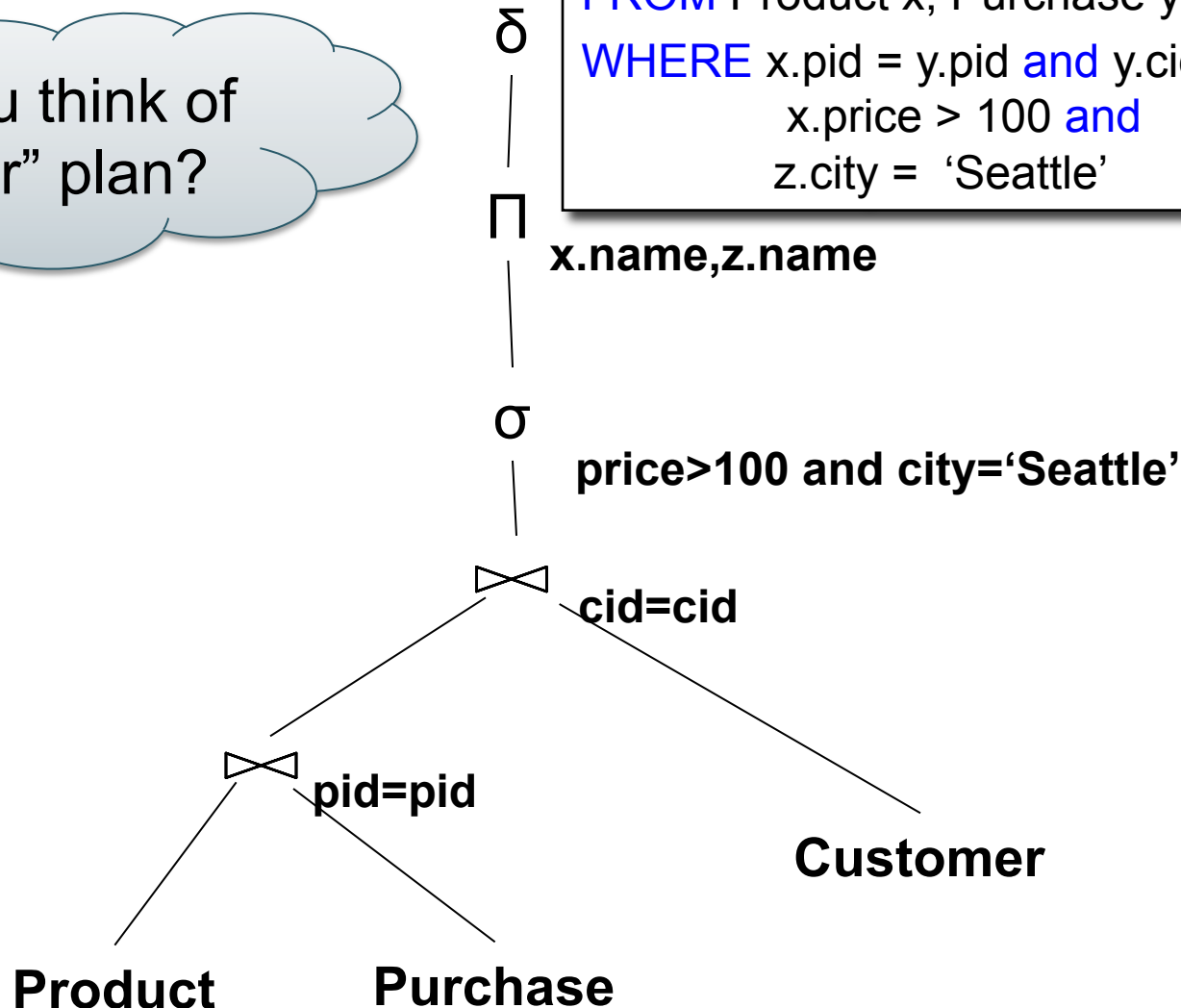FROM Product x, Purchase y, Customer z
WHERE x.pid = y.pid and y.cid = y.cid and
        x.price > 100 and
        z.city =  'Seattle'

δ

Π  **x.name,z.name**

σ  **price>100 and city='Seattle'**

⋈  **cid=cid**

⋈  **pid=pid**

**Product**      **Purchase**      **Customer**

4

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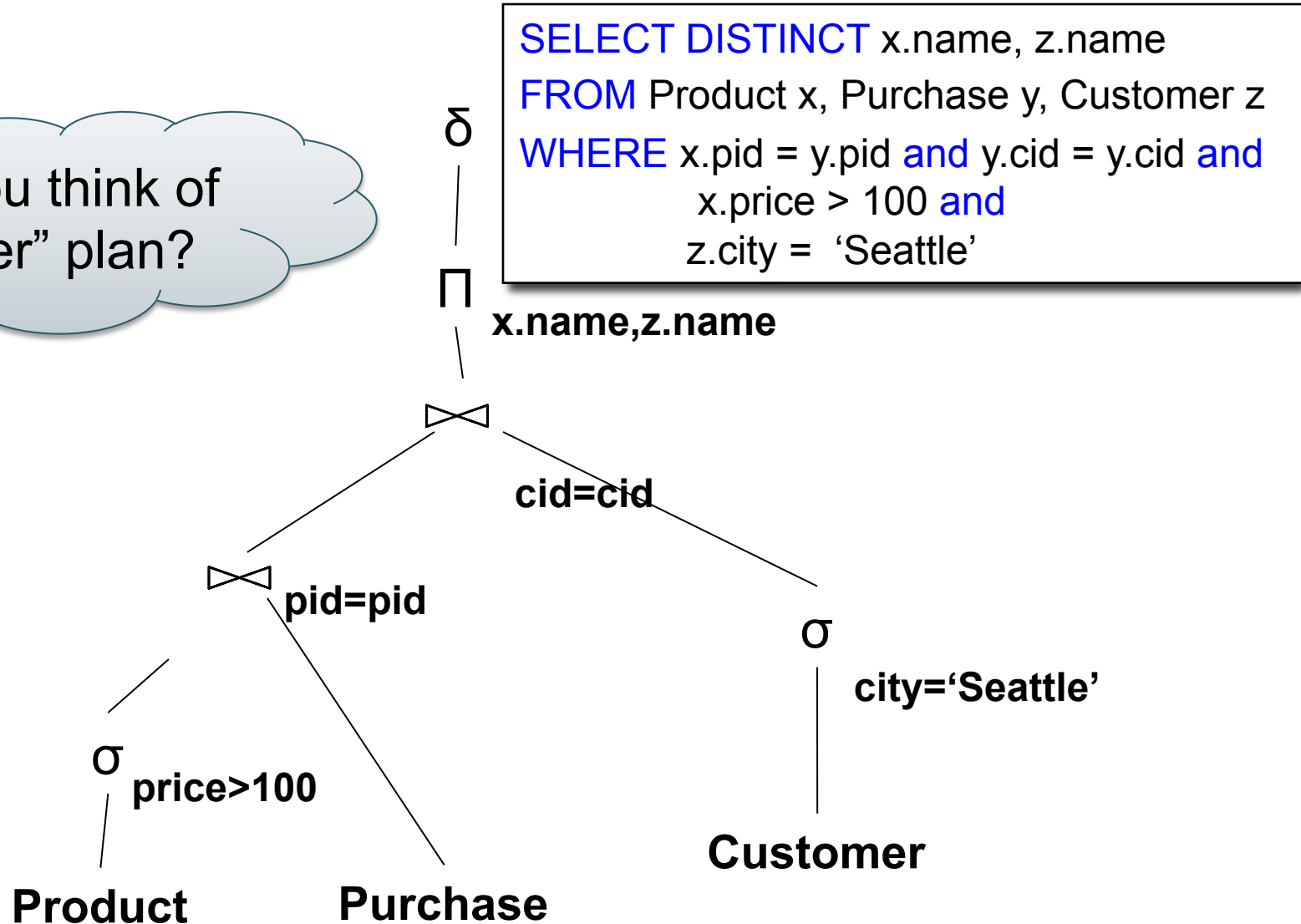Purchase(pid, cid, store)
Customer(cid, name, city)

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

Can you think of a "better" plan?

δ

Π **x.name,z.name**

σ **price>100 and city='Seattle'**

⋈ **cid=cid**

⋈ **pid=pid**

**Product**   **Purchase**

**Customer**

# Equivalent Expression

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

SELECT DISTINCT x.name, z.name
FROM Product x, Purchase y, Customer z
WHERE x.pid = y.pid and y.cid = y.cid and
        x.price > 100 and
        z.city = 'Seattle'

Can you think of
a "better" plan?

δ

Π  x.name,z.name

⋈  cid=cid

⋈  pid=pid

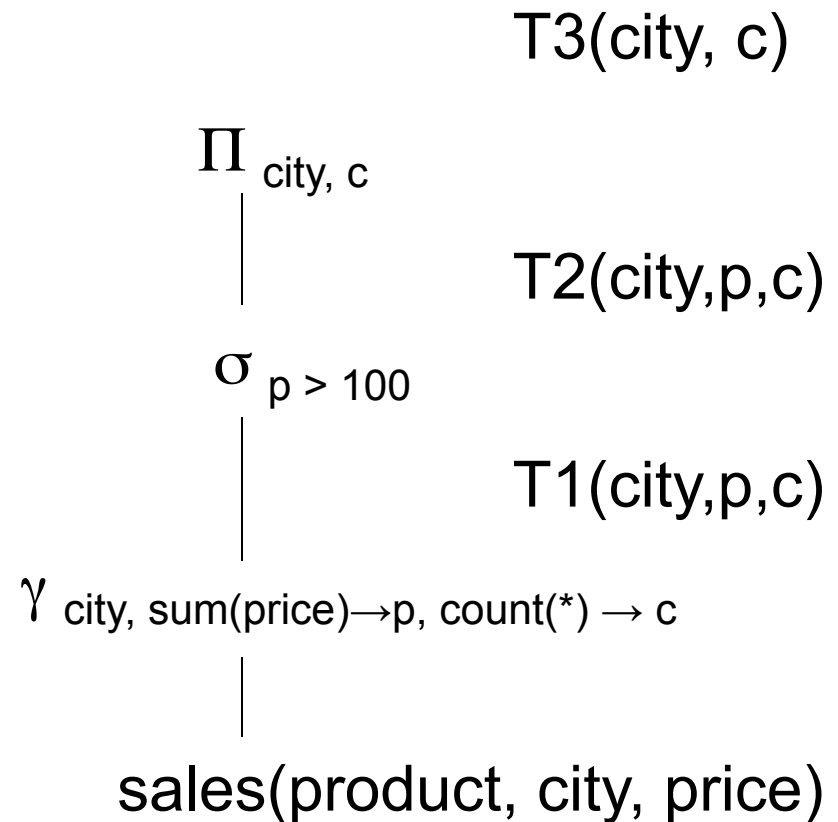σ  city='Seattle'

σ  price>100

Product

Purchase

Customer

6

Query optimization = finding cheaper, equivalent expressions

# Extended RA: Operators on Bags

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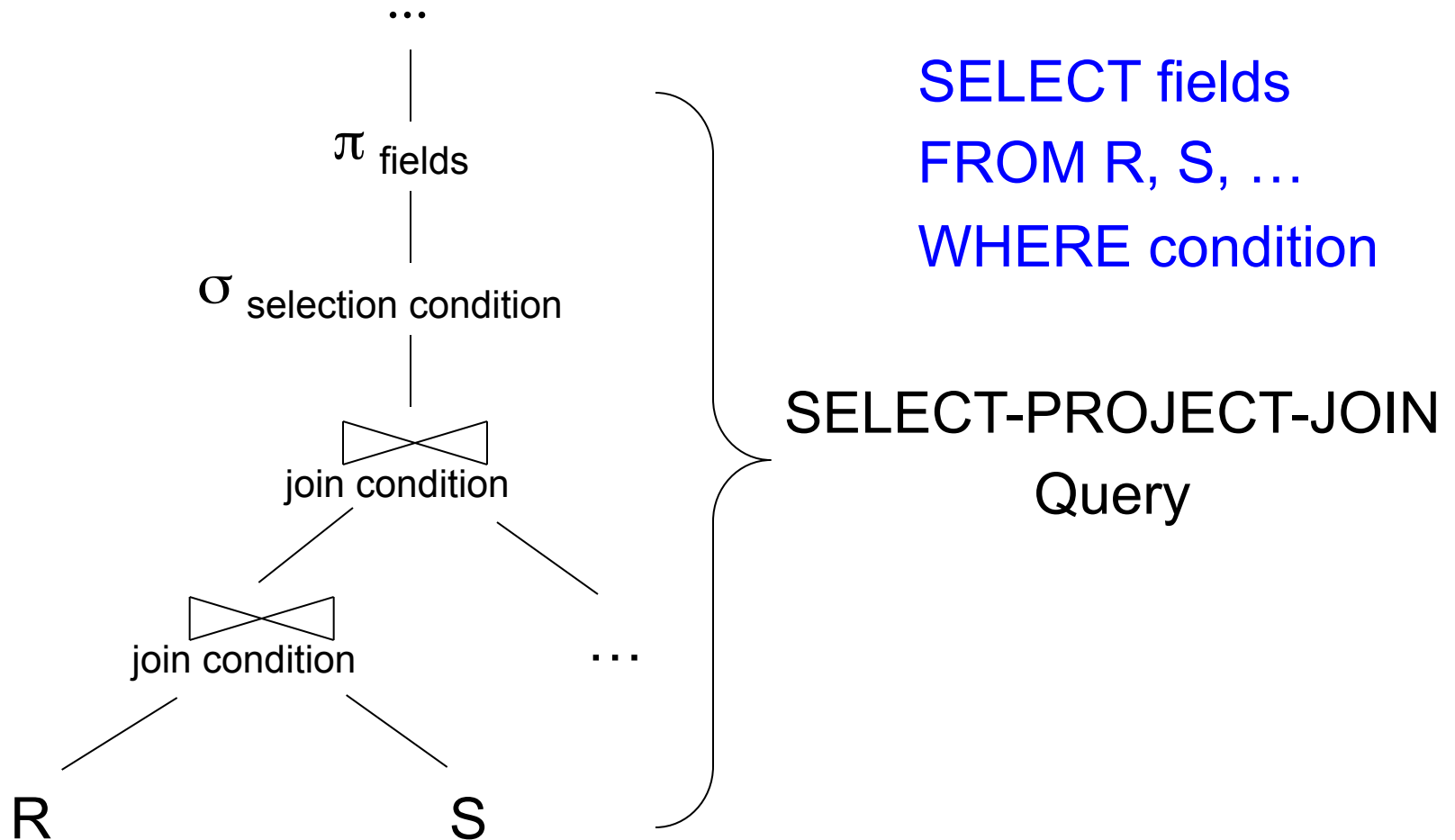- Sorting $\tau$

# 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)

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

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)

$\sigma_{\text{having condition}}$

|

$\gamma_{\text{fields, sum/count/min/max(fields)}}$

|

$\pi_{\text{fields}}$

|

$\sigma_{\text{where condition}}$

|

⋈
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)
```

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 !

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

# How about Subqueries?

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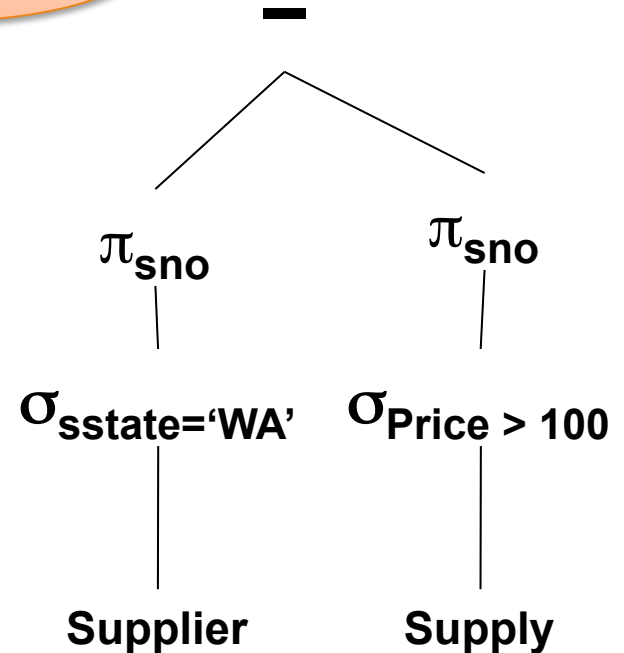
```
SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
    and Q.sno not in
        (SELECT P.sno
        FROM Supply P
        WHERE P.price > 100)
```

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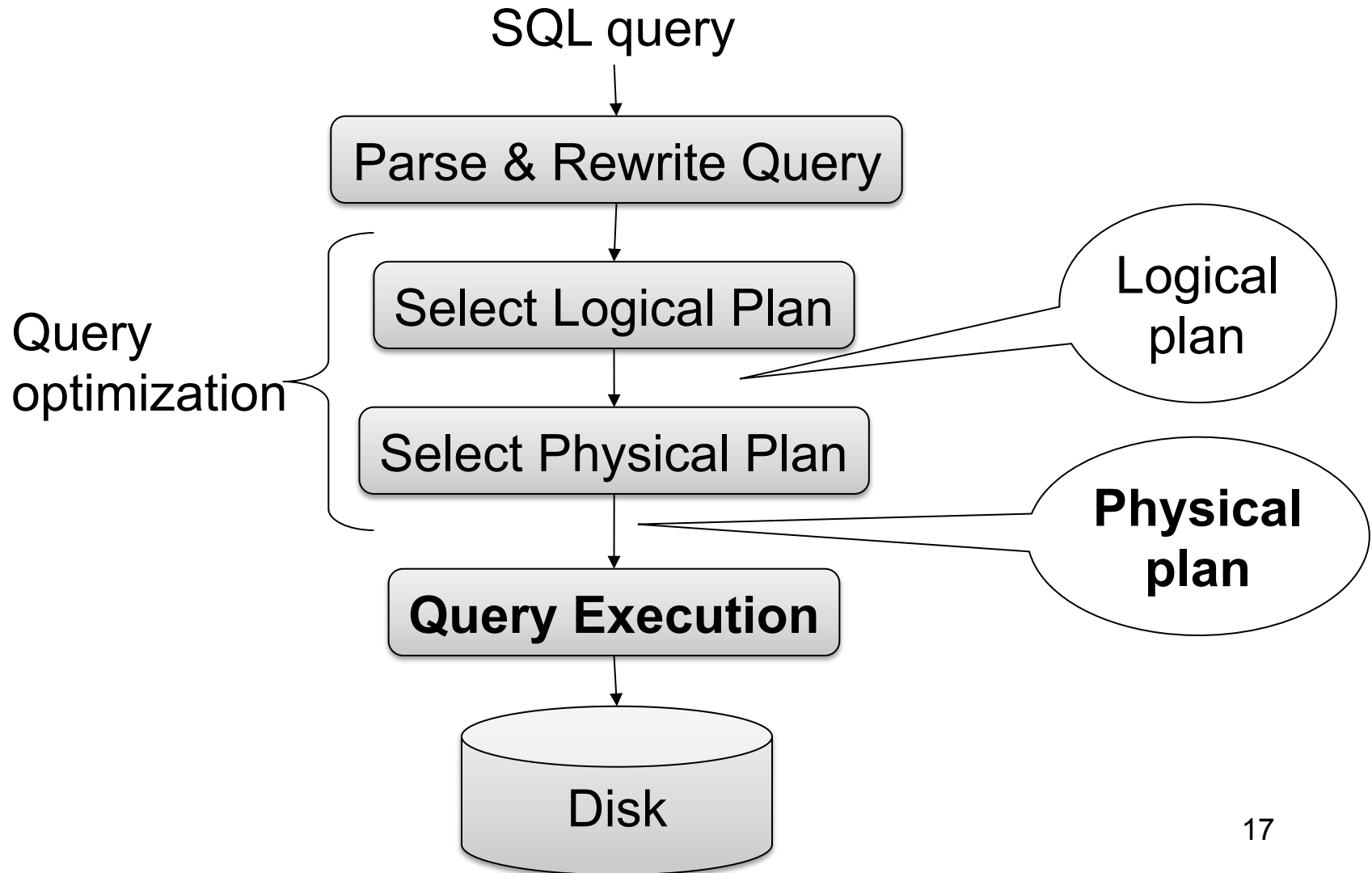Supply(sno,pno,price)

# 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)

# 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…



$$-$$

$\pi_{\textbf{sno}} \qquad \pi_{\textbf{sno}}$

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

**Supplier**      **Supply**

# From Logical Plans to Physical Plans

# Query Evaluation Steps Review

SQL query

↓

Parse & Rewrite Query

↓

Query optimization { 

Select Logical Plan → Logical plan

↓

Select Physical Plan → **Physical plan**

↓

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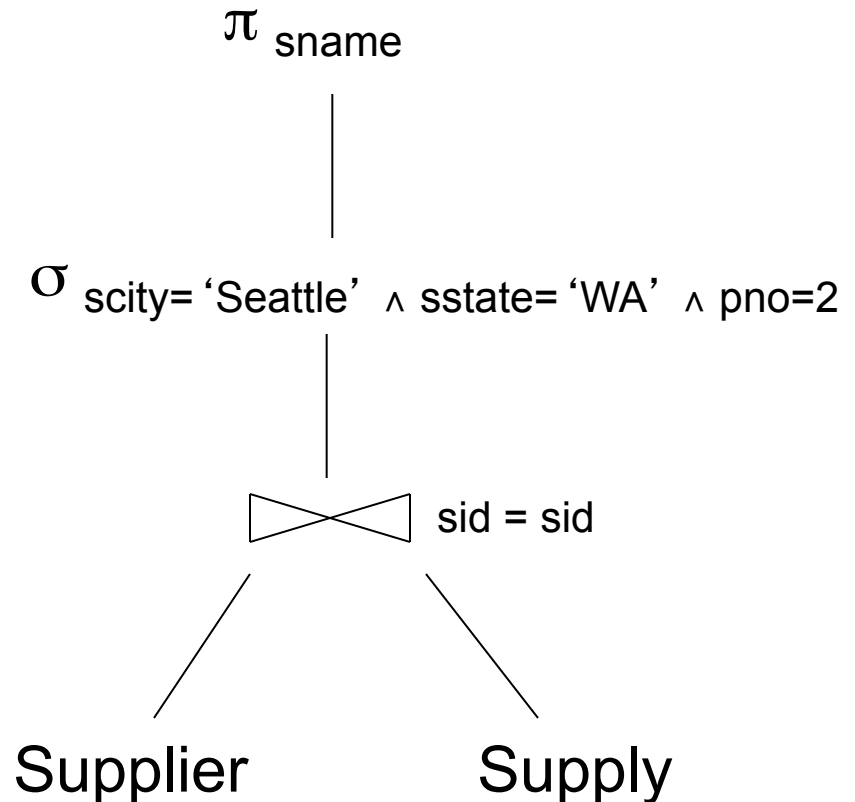
Supplier(<u>sid</u>, sname, scity, sstate)

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'

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

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"

$\pi_{sname}$

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

⋈ 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)

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

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

(Nested loop)

⋈ 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(<u>sid</u>, sname, scity, sstate)
Supply(<u>sid, pno</u>, quantity)
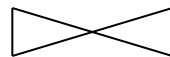
# Physical Query Plan 2

(On the fly)  $\pi$ ~sname~

(On the fly)

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

Same logical query plan
Different physical plan

(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

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

> Different but equivalent logical query plan; different physical plan

(Sort-merge join)     (c)
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'

(Scan & write to T1)

(Scan & write to T2)

(a) $\sigma_{scity='Seattle' \land sstate='WA'}$          (b) $\sigma_{pno=2}$

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?
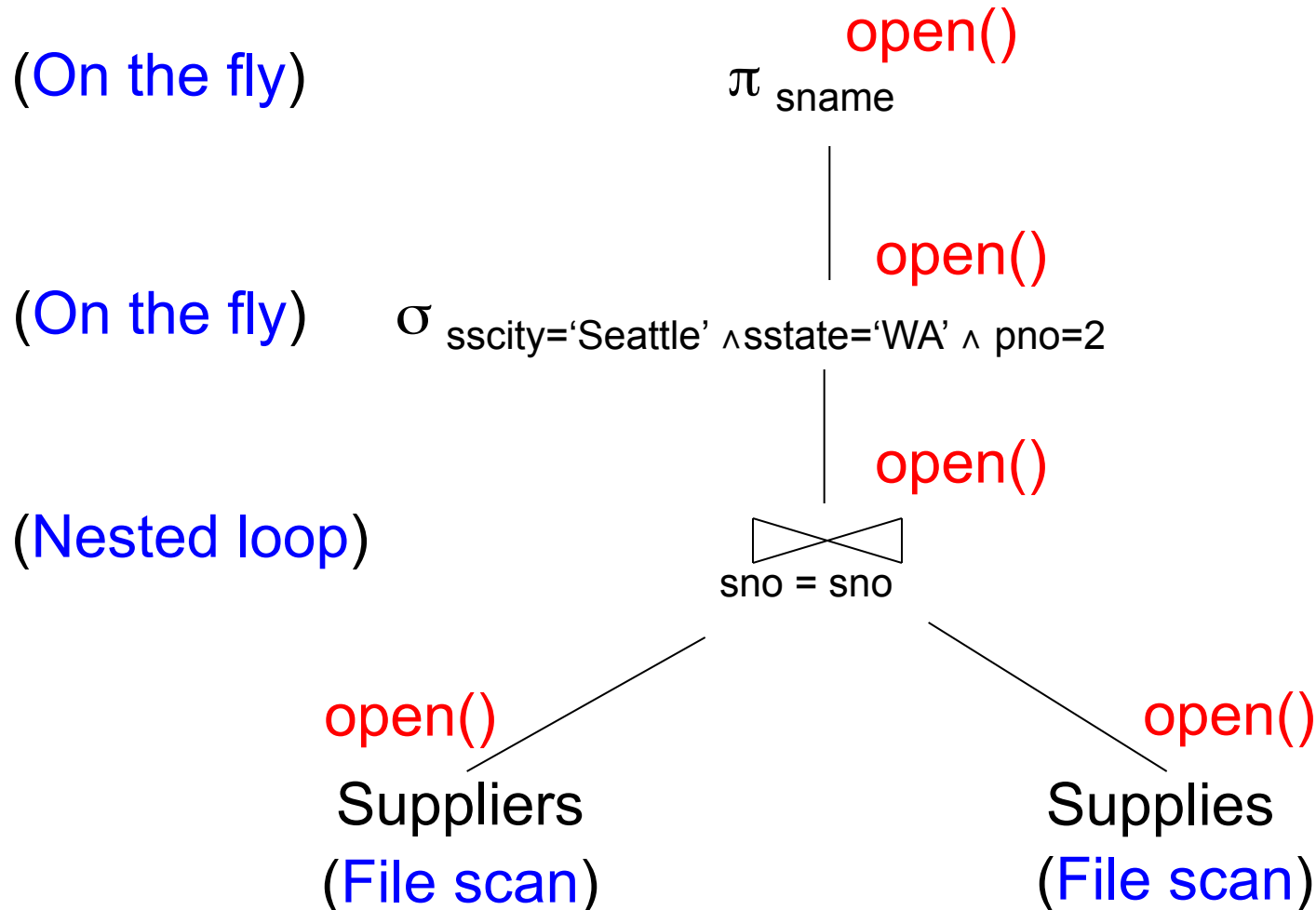  – Will discuss in a few lectures

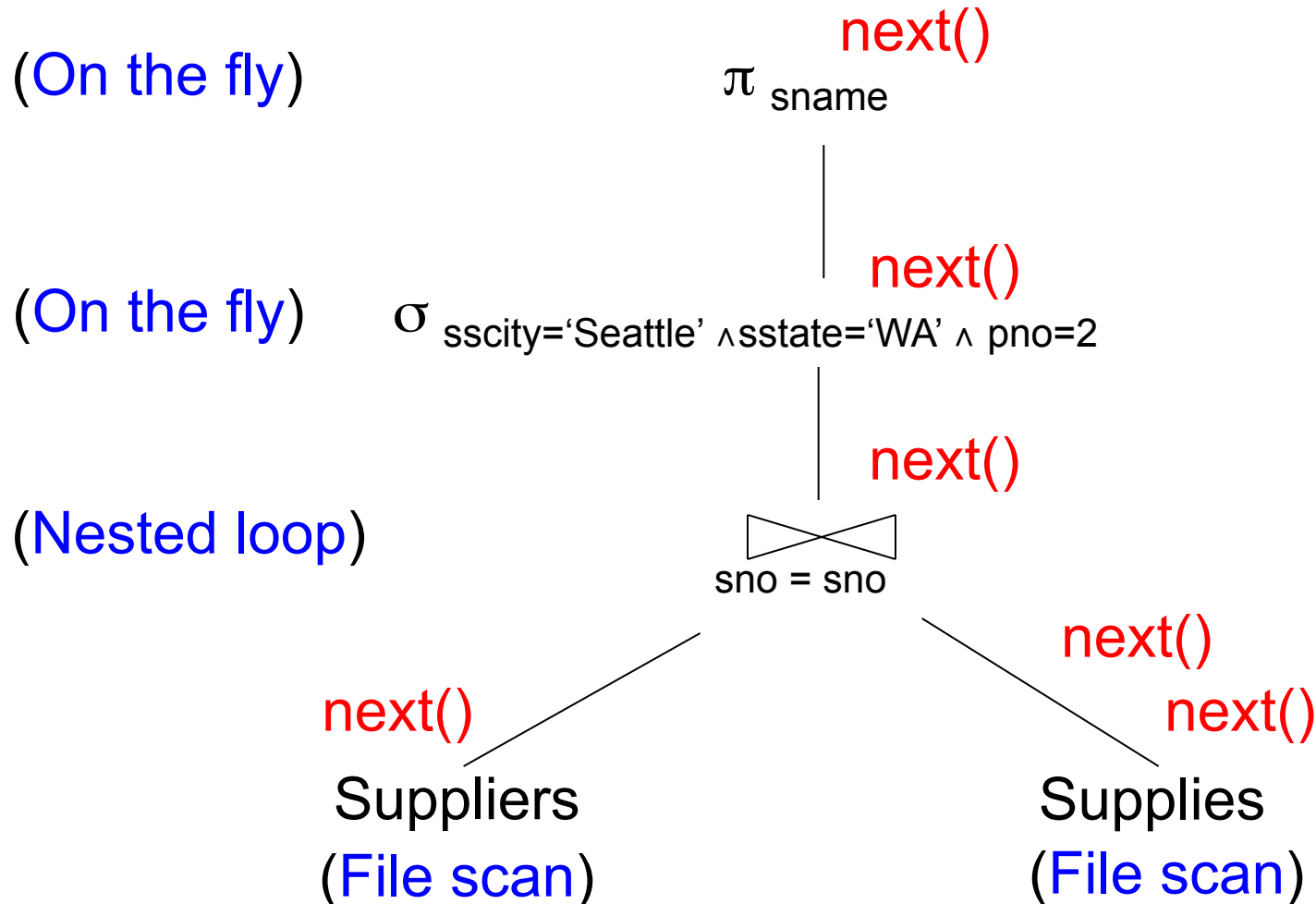# Demonstration with SQL Server Management Studio

# Query Execution

# Iterator Interface

- **open()**
  - Initializes operator state
  - Sets parameters such as selection condition
- **next()**
  - Operator invokes get_next() recursively on its inputs
  - Performs processing and produces an output tuple
- **close()**: clean-up state

# Pipelined Query Execution

(On the fly)

open()

$\pi_{\text{sname}}$

open()

(On the fly)

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

open()

(Nested loop)

⋈ sno = sno

open()

open()

Suppliers
(File scan)

Supplies
(File scan)

CSE 344 - Winter 2016

28

# Pipelined Query Execution

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

next()

(On the fly)    $\sigma_{\text{sscity='Seattle'} \wedge \text{sstate='WA'} \wedge \text{pno=2}}$

next()

(Nested loop)    ⋈
sno = sno

next()

next()

next()

Suppliers                Supplies
(File scan)              (File scan)

CSE 344 - Winter 2016                                    29

# Pipelined Execution

- Tuples generated by an operator are immediately sent to the parent

- Benefits:
  - No operator synchronization issues
  - No need to buffer tuples between operators
  - Saves cost of writing intermediate data to disk
  - Saves cost of reading intermediate data from disk

- This approach is used whenever possible
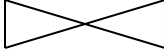
# Intermediate Tuple Materialization

- Tuples generated by an operator are written to disk an in intermediate table

- No direct benefit

- Necessary:
    - For certain operator implementations
    - When we don't have enough memory

# Intermediate Tuple Materialization

(On the fly)

$\pi$ sname

(Sort-merge join)

⋈ sno = sno

(Scan: write to T1)

(Scan: write to T2)

$\sigma$ sscity='Seattle' ∧sstate='WA'

$\sigma$ pno=2

Suppliers
(File scan)

Supplies
(File scan)

# Query Execution Bottom Line

- SQL query transformed into physical plan
  - **Access path selection** for each relation
    - Scan the relation or use an index (see next lecture)
  - **Implementation choice** for each operator
    - Nested loop join, hash join, etc.
  - **Scheduling decisions** for operators
    - Pipelined execution or intermediate materialization
- Execution of the physical plan is pull-based

# Physical Data Independence

- Means that applications are insulated from changes in physical storage details
  - E.g., can add/remove indexes without changing apps
  - Can do other physical tunings for performance

- SQL and relational algebra facilitate physical data independence because both languages are "set-at-a-time": Relations as input and output