

CSE 344

APRIL 27TH – COST ESTIMATION

ADMINISTRIVIA

- **HW5 Out**
 - Please verify that you can run queries
- **Midterm**
 - May 9th 9:30-10:20 – MLR 301
 - Review (in class) – May 7th
 - Practice exam – May 4th
 - Through parallelism: next week's material

INDEX BASED SELECTION

Example:

$$\begin{aligned} B(R) &= 2000 \\ T(R) &= 100,000 \\ V(R, a) &= 20 \end{aligned}$$

$$\text{cost of } \sigma_{a=v}(R) = ?$$

Table scan: $B(R) = 2,000$ I/Os

Index based selection:

- If index is clustered: $B(R) * 1/V(R,a) = 100$ I/Os
- If index is unclustered: $T(R) * 1/V(R,a) = 5,000$ I/Os

Lesson: Don't build unclustered indexes when $V(R,a)$ is small !

OUTLINE

Join operator algorithms

- One-pass algorithms (Sec. 15.2 and 15.3)
- Index-based algorithms (Sec 15.6)

Note about readings:

- In class, we discuss only algorithms for joins
- Other operators are easier: read the book

JOIN ALGORITHMS

Hash join

Nested loop join

Sort-merge join

HASH JOIN

Hash join: $R \bowtie S$

Scan R , build buckets in main memory

Then scan S and join

Cost: $B(R) + B(S)$

Which relation to build the hash table on?

HASH JOIN

Hash join: $R \bowtie S$

Scan R , build buckets in main memory

Then scan S and join

Cost: $B(R) + B(S)$

Which relation to build the hash table on?

One-pass algorithm when $B(R) \leq M$

- M = number of memory pages available

HASH JOIN EXAMPLE

Patient(pid, name, address)

Insurance(pid, provider, policy_nb)

Patient \bowtie Insurance

Two tuples
per page

Patient

1	'Bob'	'Seattle'
2	'Ela'	'Everett'

3	'Jill'	'Kent'
4	'Joe'	'Seattle'

Insurance

2	'Blue'	123
4	'Prem'	432

4	'Prem'	343
3	'GrpH'	554

HASH JOIN EXAMPLE

Patient \bowtie Insurance

Some large-enough #

Memory M = 21 pages

Showing pid only

Disk

Patient Insurance

1	2	2	4	6	6
3	4	4	3	1	3
9	6	2	8		
8	5	8	9		

This is one page with two tuples

HASH JOIN EXAMPLE

Step 1: Scan Patient and **build** hash table in memory

Can be done in
method open()

Memory M = 21 pages

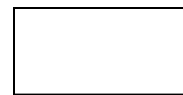
Hash h: pid % 5

5		1	6	2		3	8	4	9
---	--	---	---	---	--	---	---	---	---

Disk

Patient Insurance

1	2	2	4	6	6
3	4	4	3	1	3
9	6	2	8		
8	5	8	9		



Input buffer

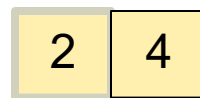
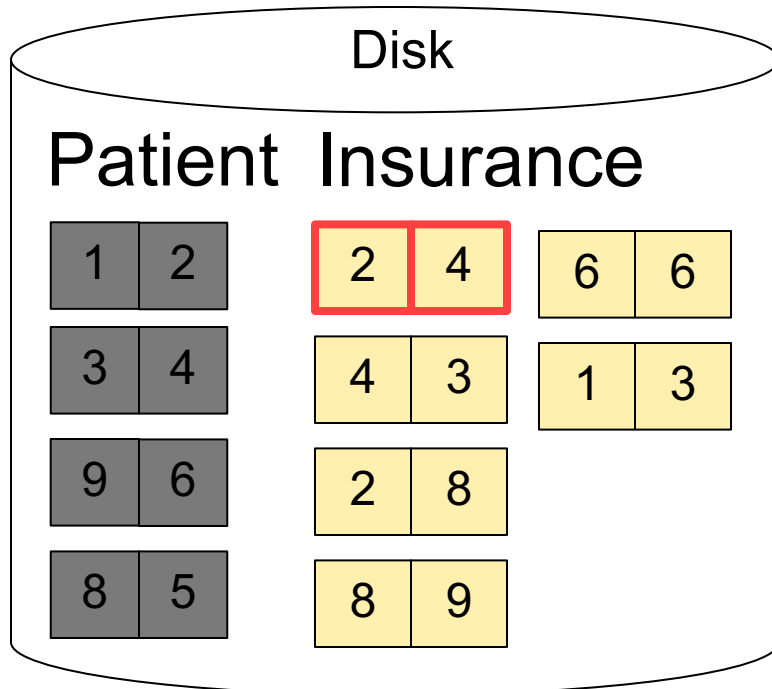
HASH JOIN EXAMPLE

Step 2: Scan Insurance and **probe** into hash table
Done during calls to next()

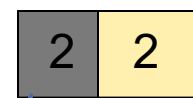
Memory M = 21 pages

Hash h: pid % 5

5		1	6	2		3	8	4	9
---	--	---	---	---	--	---	---	---	---



Input buffer



Output buffer

Write to disk or
pass to next
operator

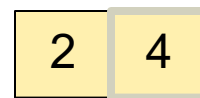
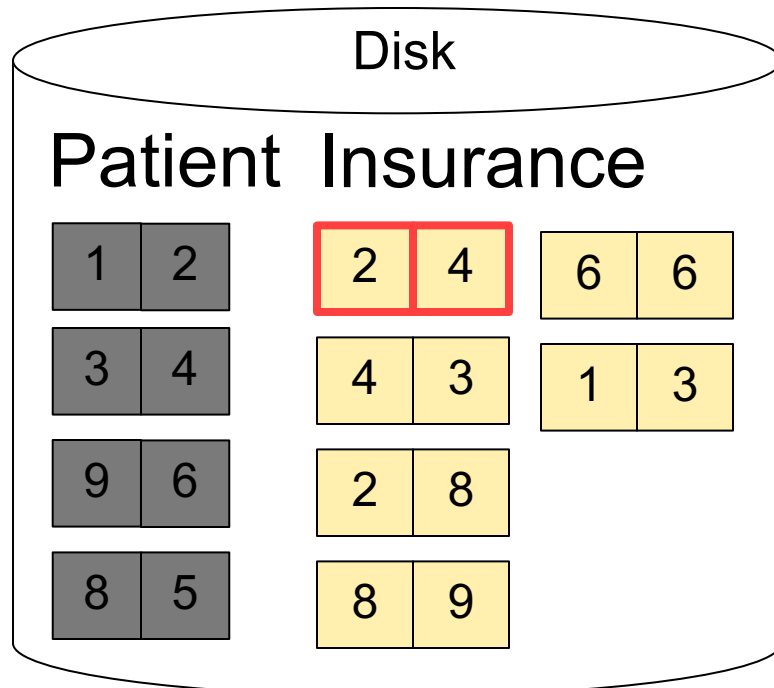
HASH JOIN EXAMPLE

Step 2: Scan Insurance and **probe** into hash table
Done during calls to next()

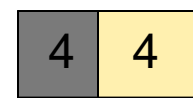
Memory M = 21 pages

Hash h: pid % 5

5		1	6	2		3	8	4	9
---	--	---	---	---	--	---	---	---	---



Input buffer



Output buffer

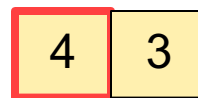
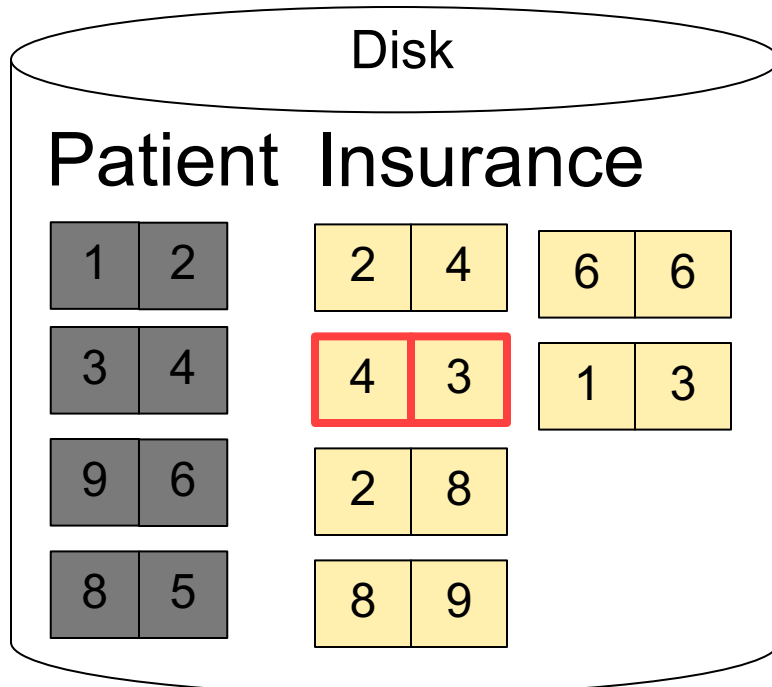
HASH JOIN EXAMPLE

Step 2: Scan Insurance and **probe** into hash table
 Done during calls to next()

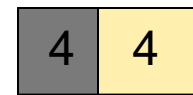
Memory M = 21 pages

Hash h: pid % 5

5		1	6	2		3	8	4	9
---	--	---	---	---	--	---	---	---	---



Input buffer



Output buffer

Keep going until read all of Insurance

Cost: $B(R) + B(S)$

NESTED LOOP JOINS

Tuple-based nested loop $R \bowtie S$

R is the outer relation, **S** is the inner relation

```
for each tuple  $t_1$  in R do  
  for each tuple  $t_2$  in S do  
    if  $t_1$  and  $t_2$  join then output  $(t_1, t_2)$ 
```

What is the **Cost**?

NESTED LOOP JOINS

Tuple-based nested loop $R \bowtie S$

R is the outer relation, S is the inner relation

```
for each tuple  $t_1$  in  $R$  do
  for each tuple  $t_2$  in  $S$  do
    if  $t_1$  and  $t_2$  join then output  $(t_1, t_2)$ 
```

Cost: $B(R) + T(R) B(S)$

Multiple-pass since S is read many times

What is the **Cost**?

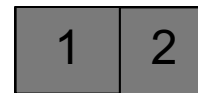
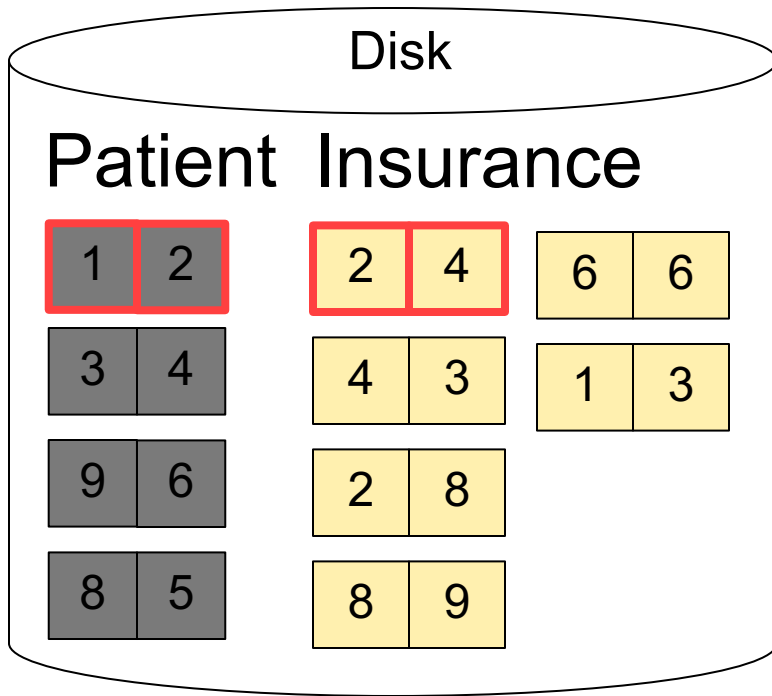
PAGE-AT-A-TIME REFINEMENT

```
for each page of tuples r in R do  
  for each page of tuples s in S do  
    for all pairs of tuples t1 in r, t2 in s  
      if t1 and t2 join then output (t1,t2)
```

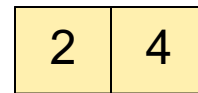
Cost: $B(R) + B(R)B(S)$

What is the **Cost**?

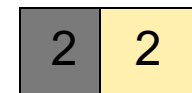
PAGE-AT-A-TIME REFINEMENT



Input buffer for Patient

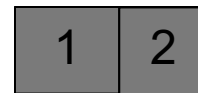
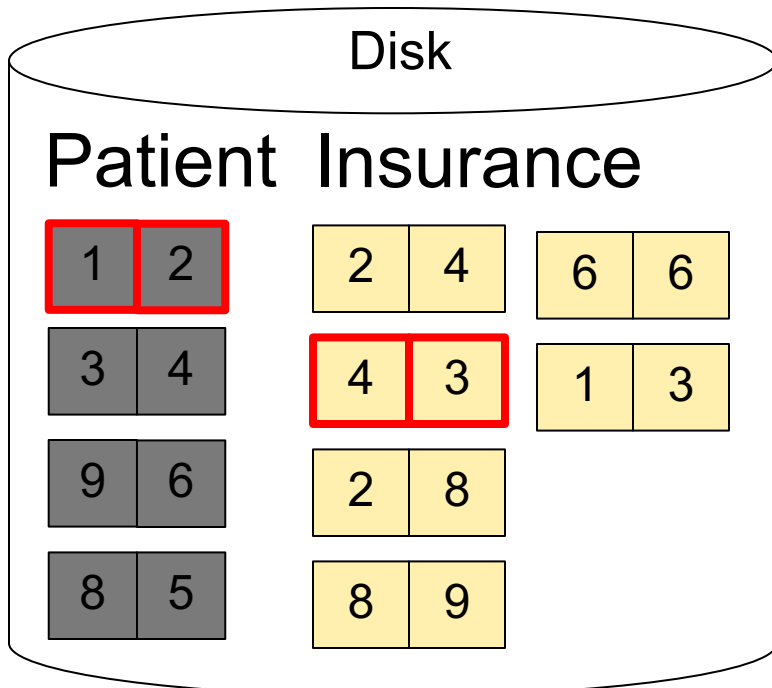


Input buffer for Insurance

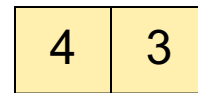


Output buffer

PAGE-AT-A-TIME REFINEMENT



Input buffer for Patient

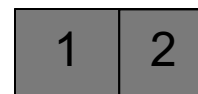
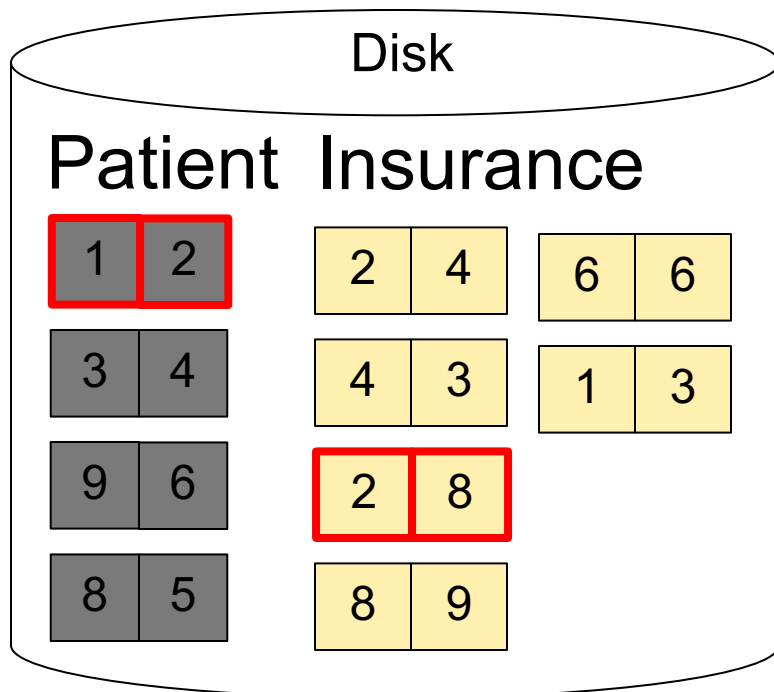


Input buffer for Insurance

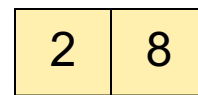


Output buffer

PAGE-AT-A-TIME REFINEMENT

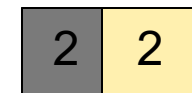


Input buffer for Patient



Input buffer for Insurance

Keep going until read
all of Insurance



Output buffer

Then repeat for next
page of Patient... until end of Patient

Cost: $B(R) + B(R)B(S)$

BLOCK-NESTED-LOOP REFINEMENT

```
for each group of M-1 pages r in R do  
  for each page of tuples s in S do  
    for all pairs of tuples t1 in r, t2 in s  
      if t1 and t2 join then output (t1,t2)
```

Cost: $B(R) + B(R)B(S)/(M-1)$

What is the **Cost**?

SORT-MERGE JOIN

Sort-merge join: $R \bowtie S$

Scan R and sort in main memory

Scan S and sort in main memory

Merge R and S

Cost: $B(R) + B(S)$

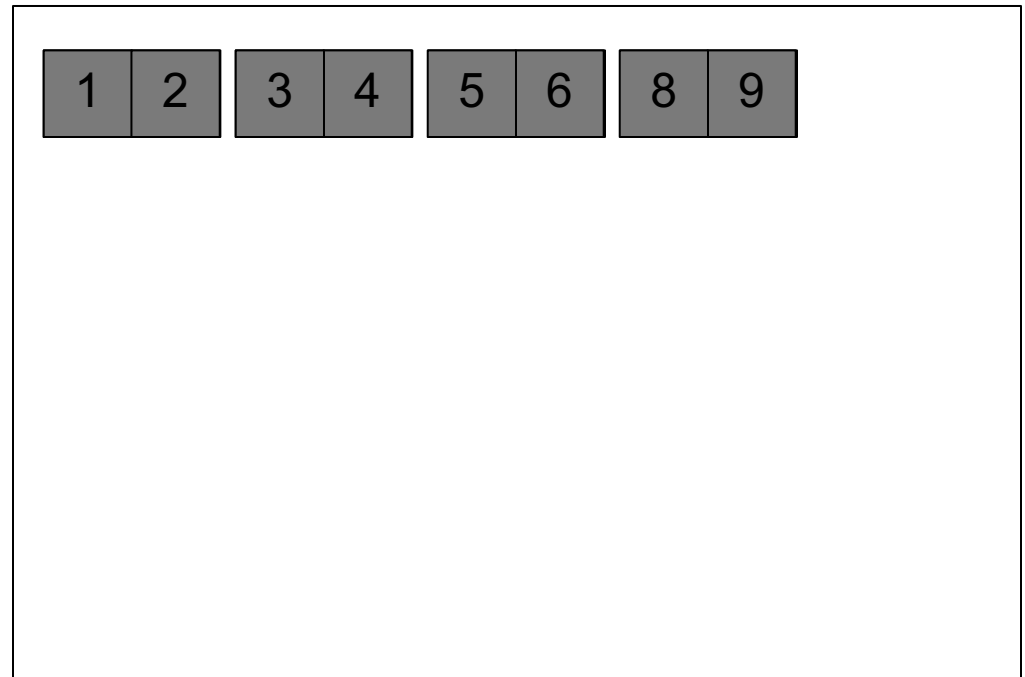
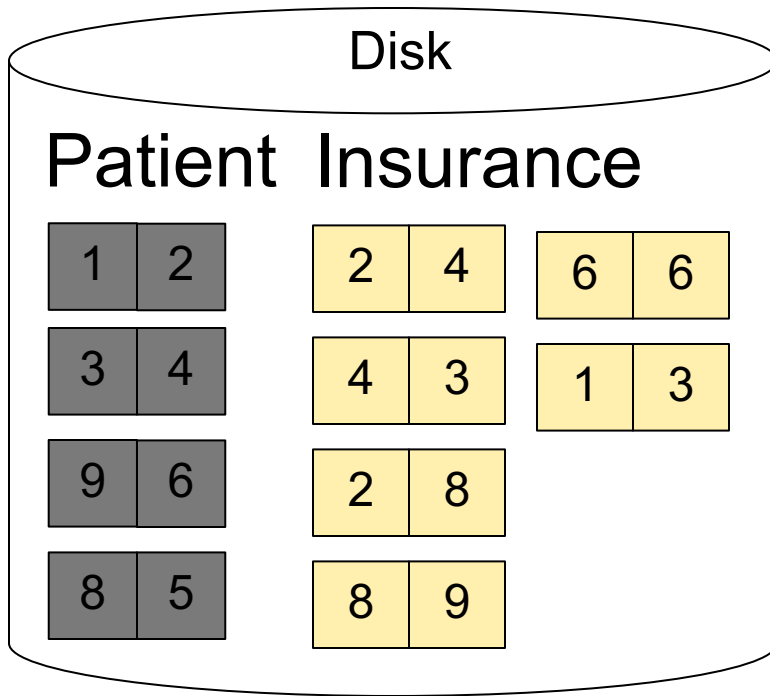
One pass algorithm when $B(S) + B(R) \leq M$

Typically, this is NOT a one pass algorithm

SORT-MERGE JOIN EXAMPLE

Step 1: Scan Patient and **sort** in memory

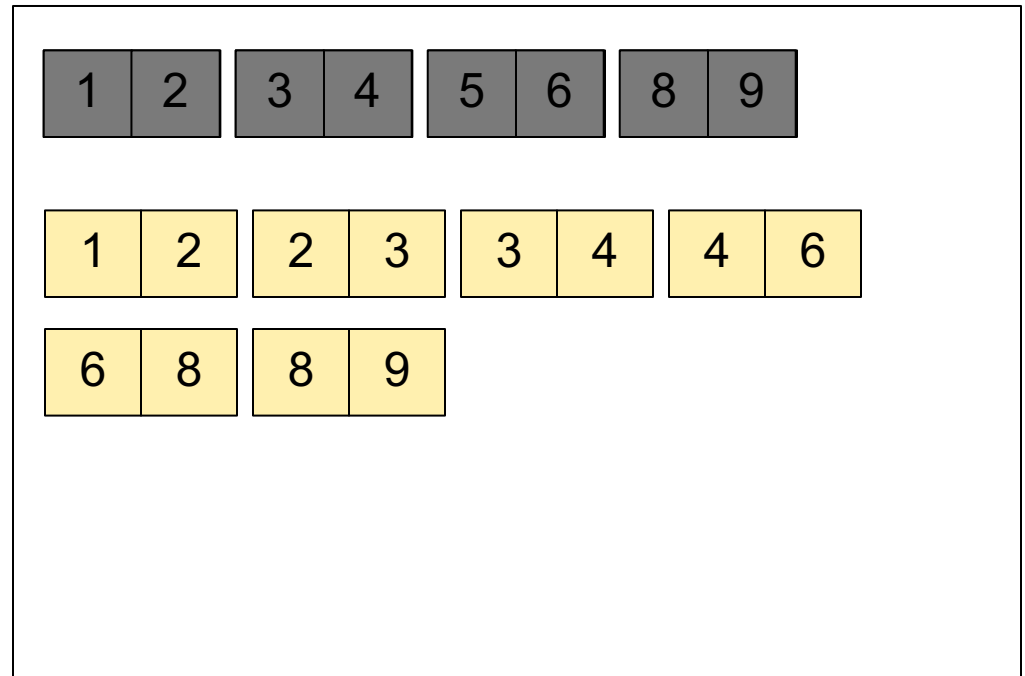
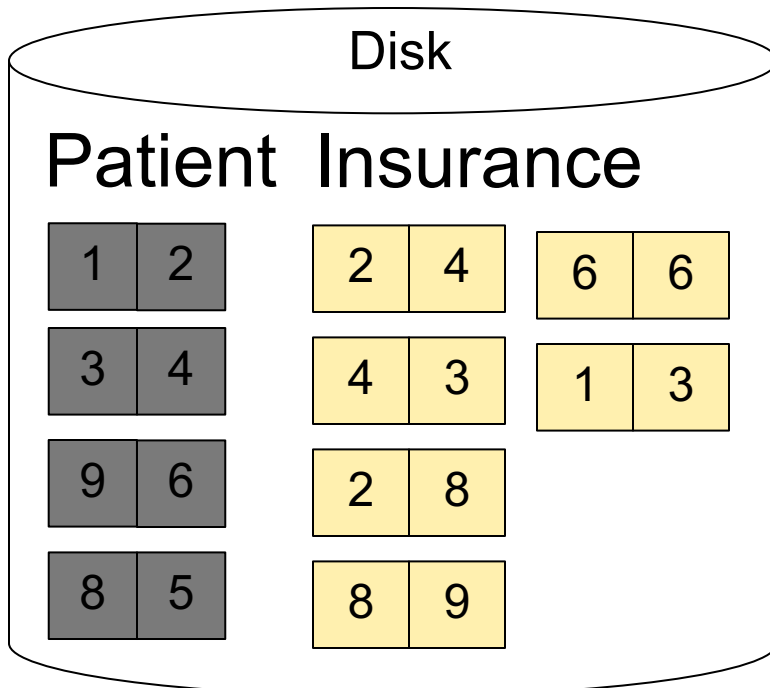
Memory M = 21 pages



SORT-MERGE JOIN EXAMPLE

Step 2: Scan Insurance and **sort** in memory

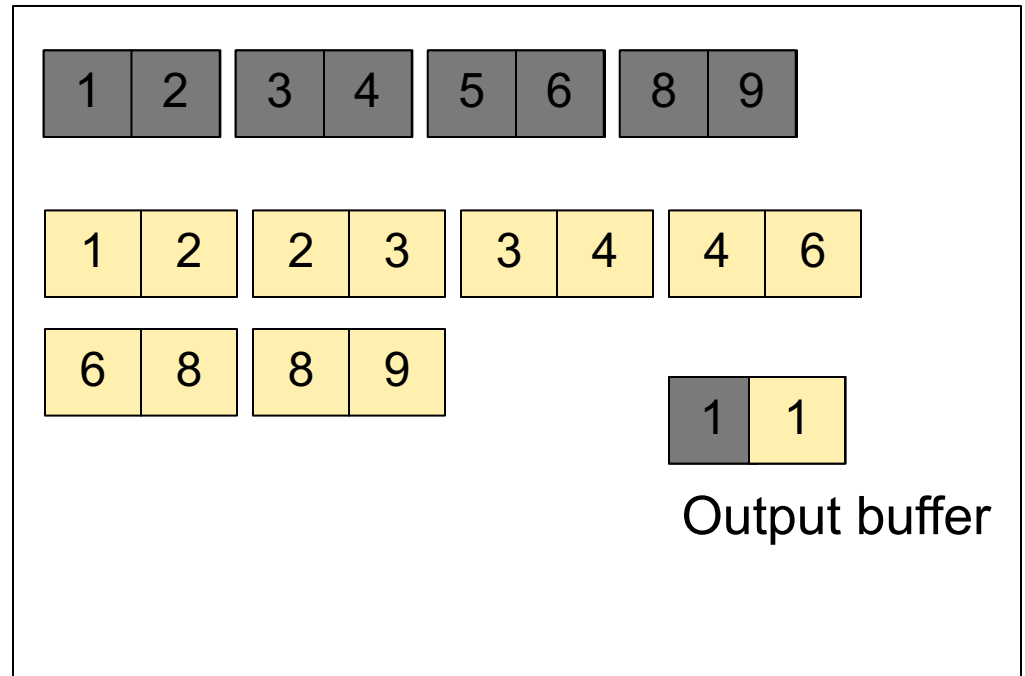
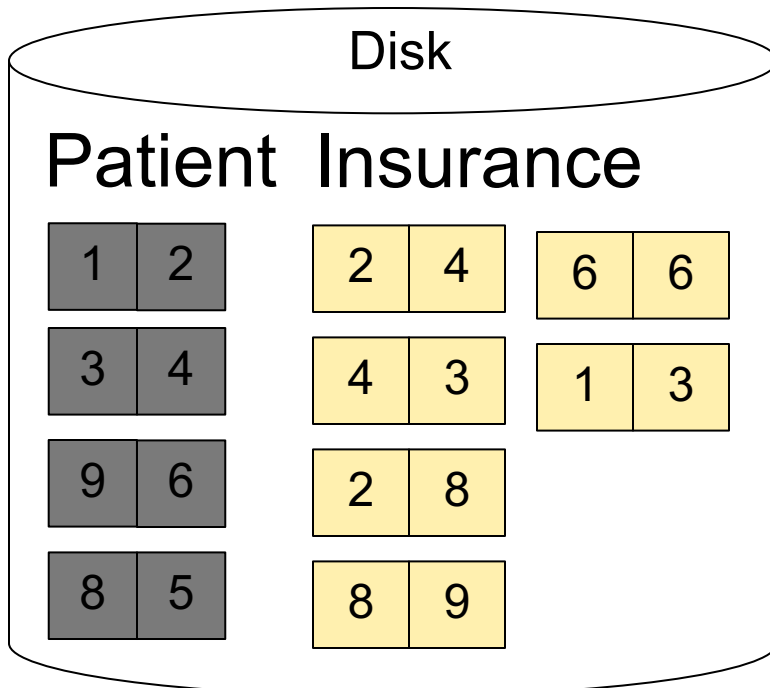
Memory M = 21 pages



SORT-MERGE JOIN EXAMPLE

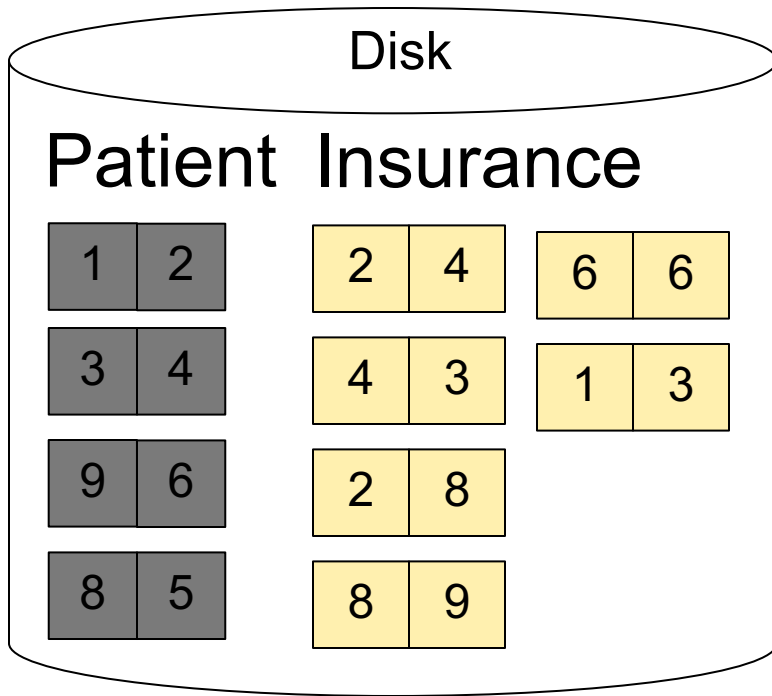
Step 3: **Merge** Patient and Insurance

Memory M = 21 pages

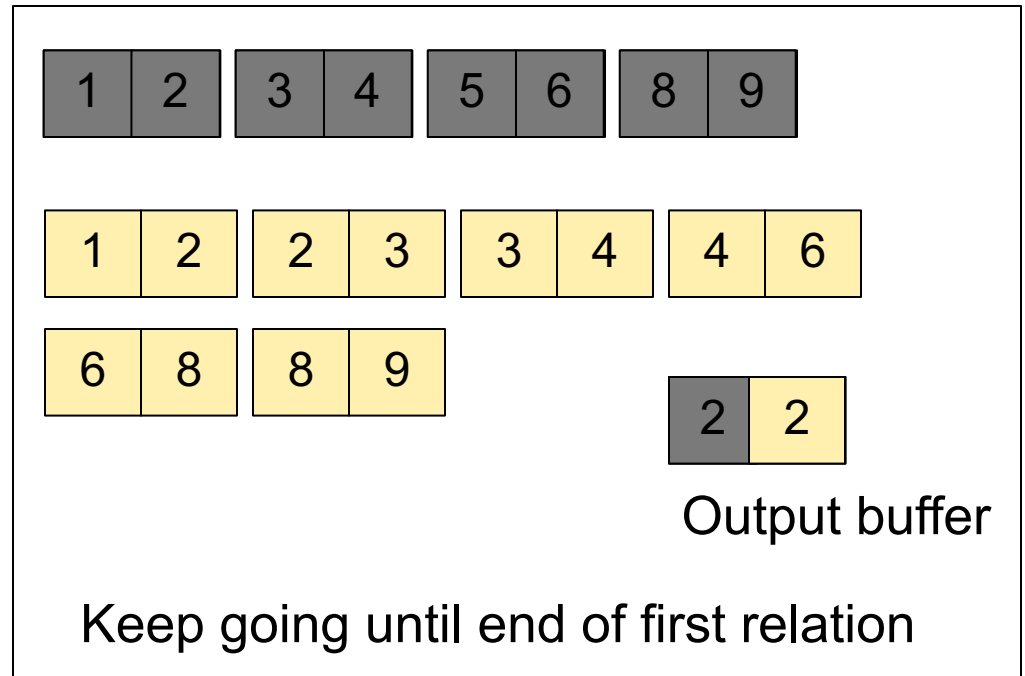


SORT-MERGE JOIN EXAMPLE

Step 3: Merge Patient and Insurance



Memory M = 21 pages



INDEX NESTED LOOP JOIN

$R \bowtie S$

Assume S has an index on the join attribute

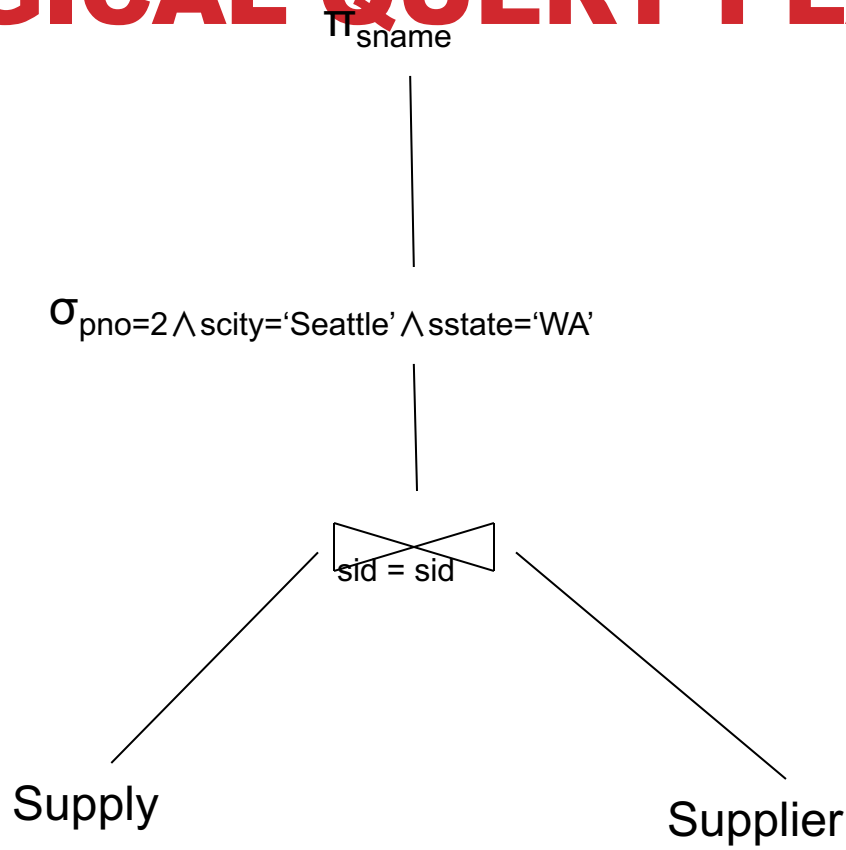
Iterate over R , for each tuple fetch corresponding tuple(s) from S

Cost:

- If index on S is clustered:
$$B(R) + T(R) * (B(S) * 1/V(S,a))$$
- If index on S is unclustered:
$$B(R) + T(R) * (T(S) * 1/V(S,a))$$

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

LOGICAL QUERY PLAN 1



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

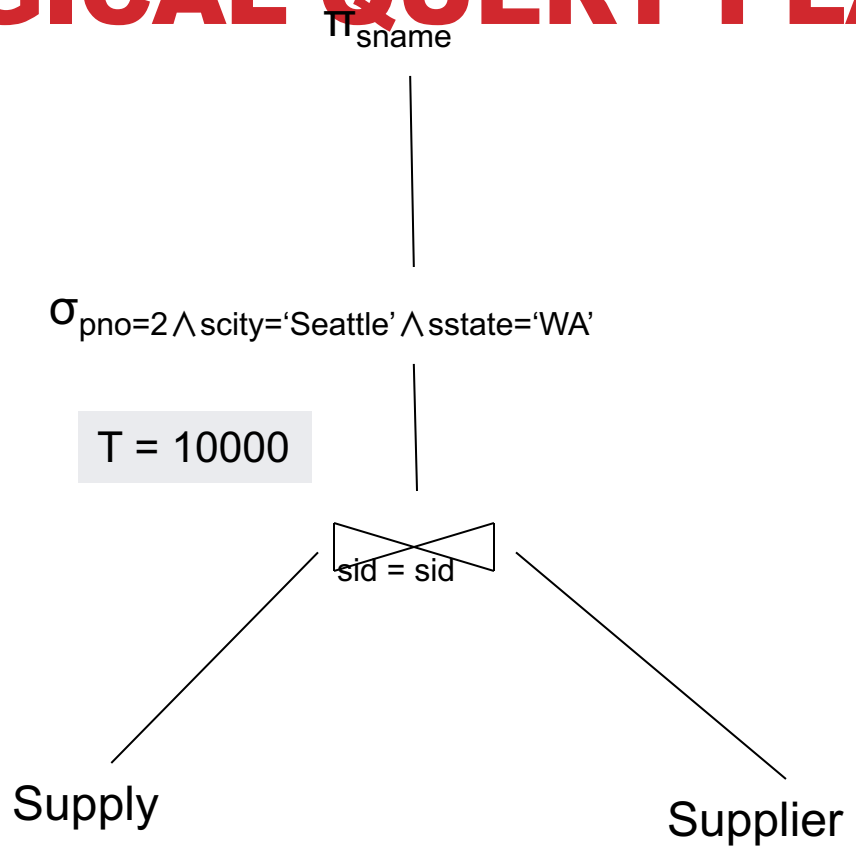
T(Supply) = 10000
B(Supply) = 100
V(Supply, pno) = 2500

T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

M=11

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

LOGICAL QUERY PLAN 1



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

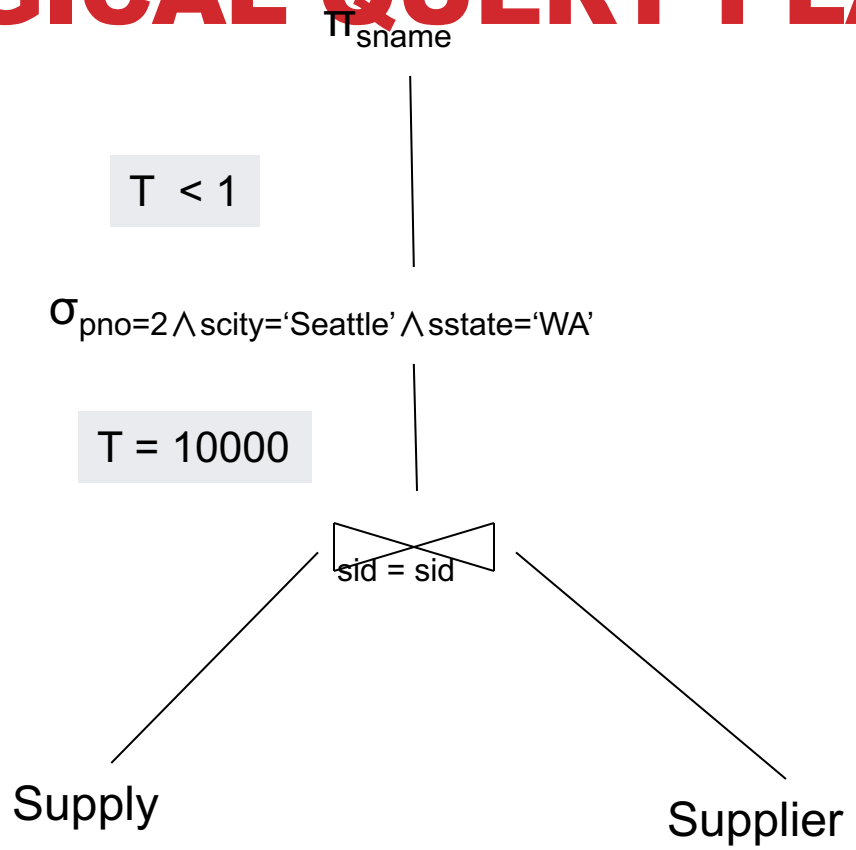
T(Supply) = 10000
B(Supply) = 100
V(Supply, pno) = 2500

T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

M=11

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

LOGICAL QUERY PLAN 1



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

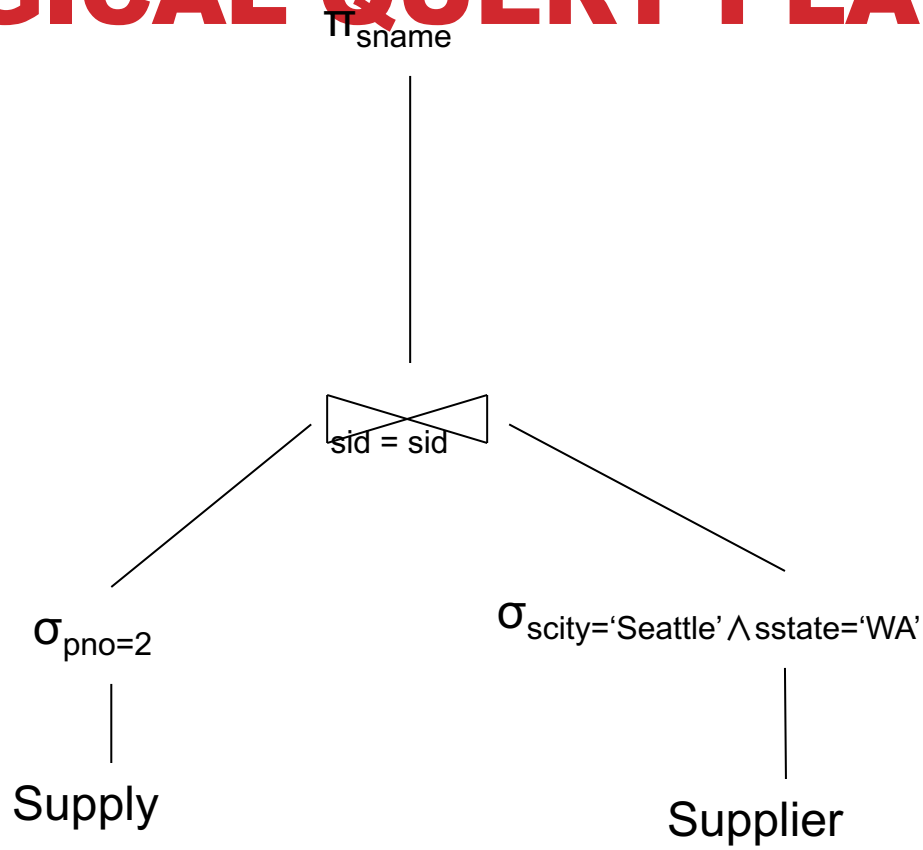
T(Supply) = 10000
B(Supply) = 100
V(Supply, pno) = 2500

T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

M=11

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

LOGICAL QUERY PLAN 2



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

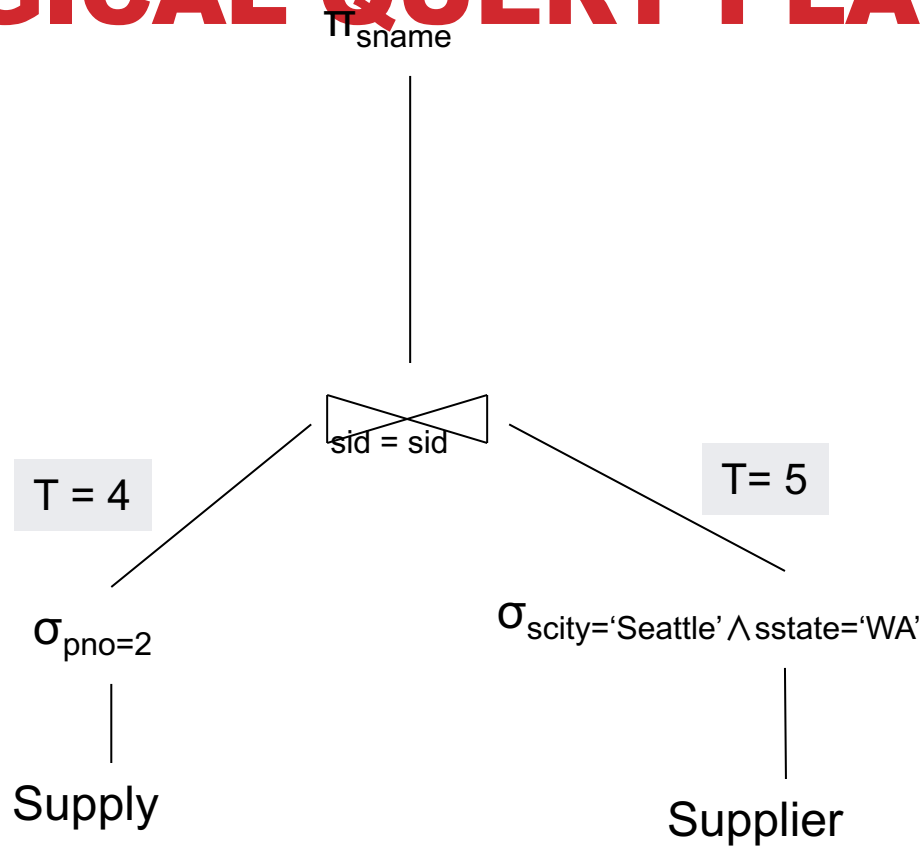
T(Supply) = 10000
B(Supply) = 100
V(Supply, pno) = 2500

T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

M=11

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

LOGICAL QUERY PLAN 2



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

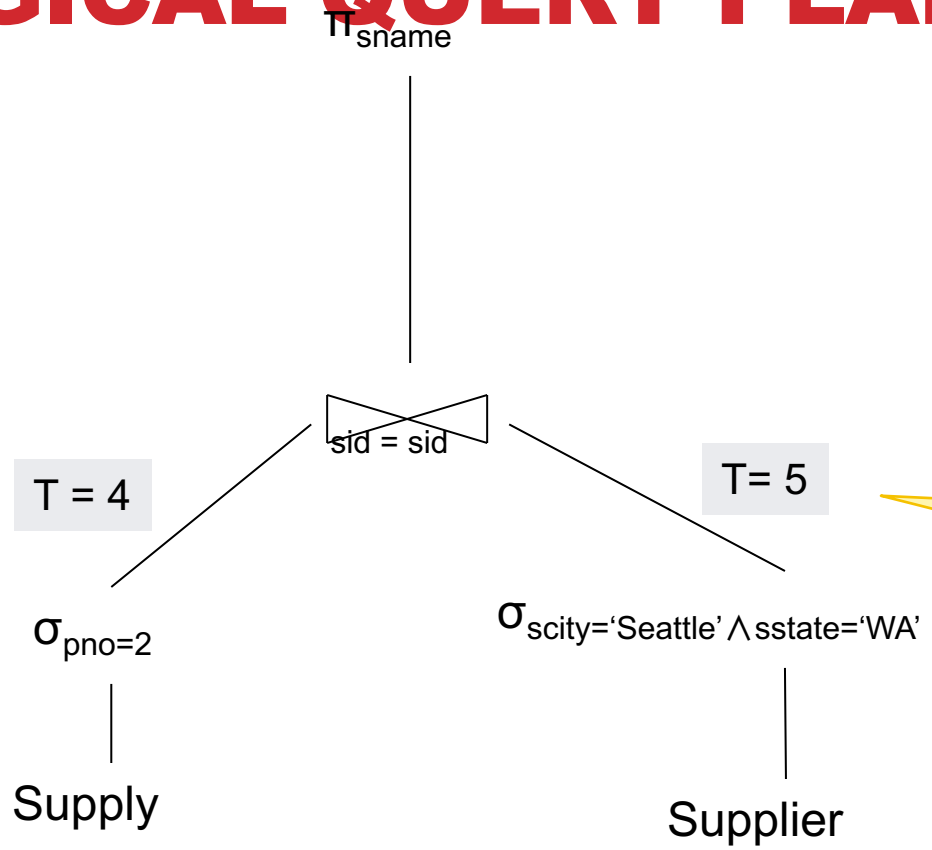
T(Supply) = 10000
B(Supply) = 100
V(Supply, pno) = 2500

T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

M=11

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

LOGICAL QUERY PLAN 2



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

Very wrong!
Why?

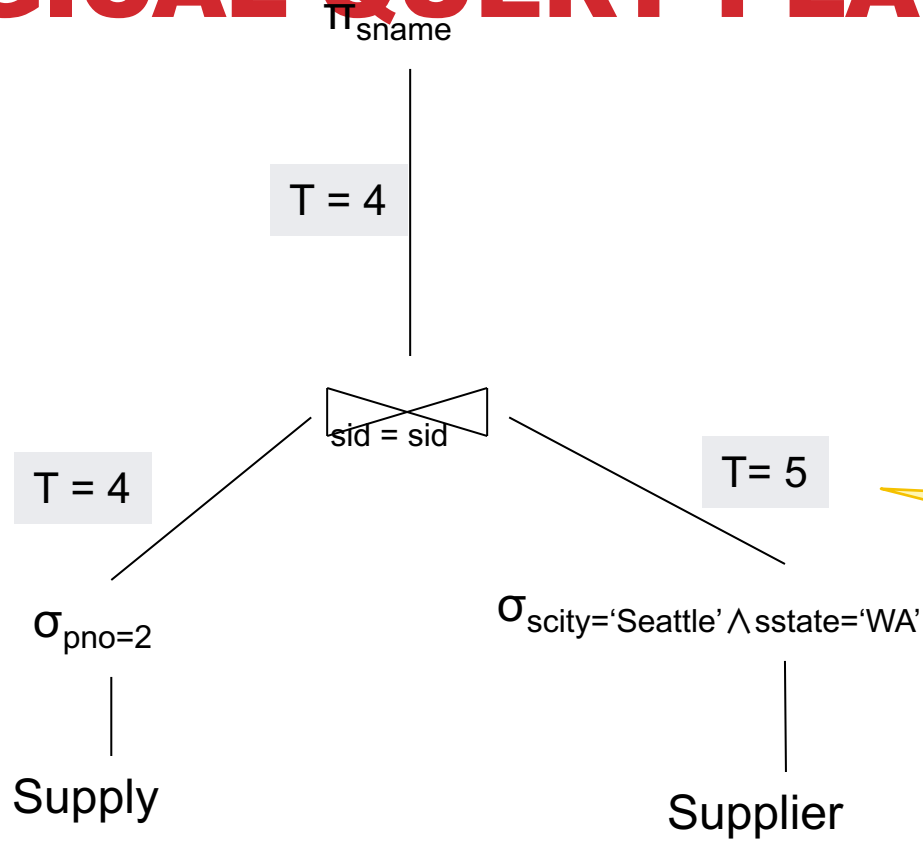
T(Supply) = 10000
B(Supply) = 100
V(Supply, pno) = 2500

T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

M=11

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

LOGICAL QUERY PLAN 2



```

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

Very wrong!
Why?

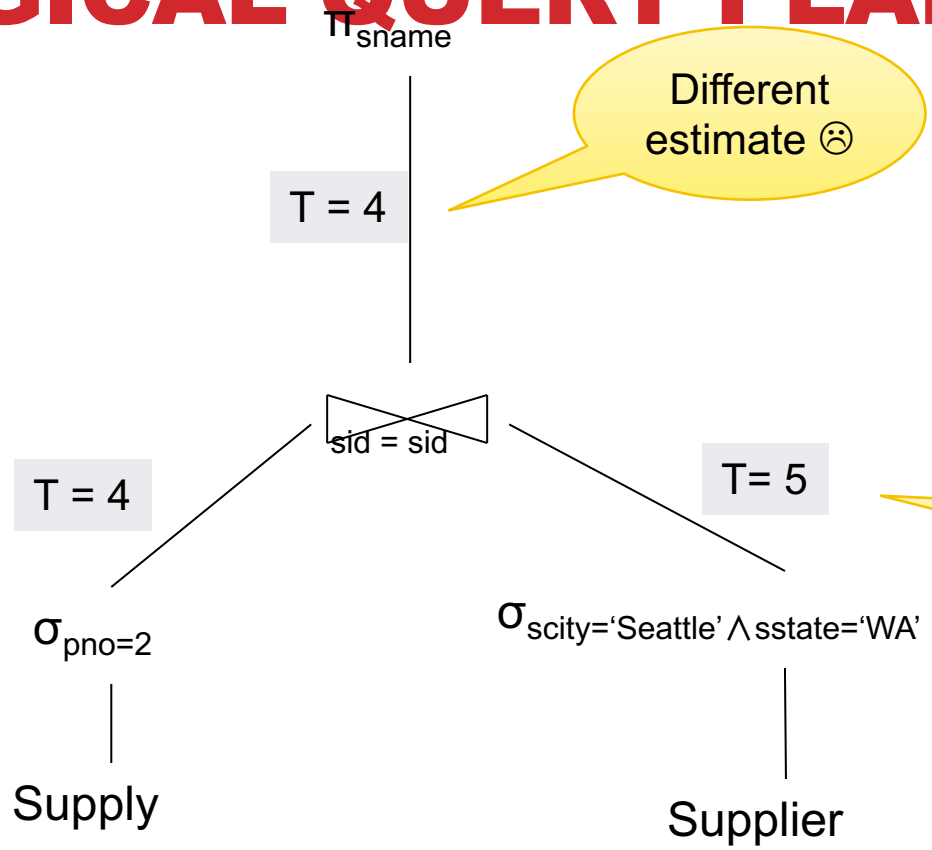
T(Supply) = 10000
 B(Supply) = 100
 V(Supply, pno) = 2500

T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

M=11

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

LOGICAL QUERY PLAN 2



```

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

Very wrong!
Why?

T(Supply) = 10000
 B(Supply) = 100
 V(Supply, pno) = 2500

T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

M=11

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

PHYSICAL PLAN 1

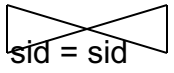
π_{sname}

T < 1

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

T = 10000

Total cost:



Block nested loop join

Scan

Supply

Scan

Supplier

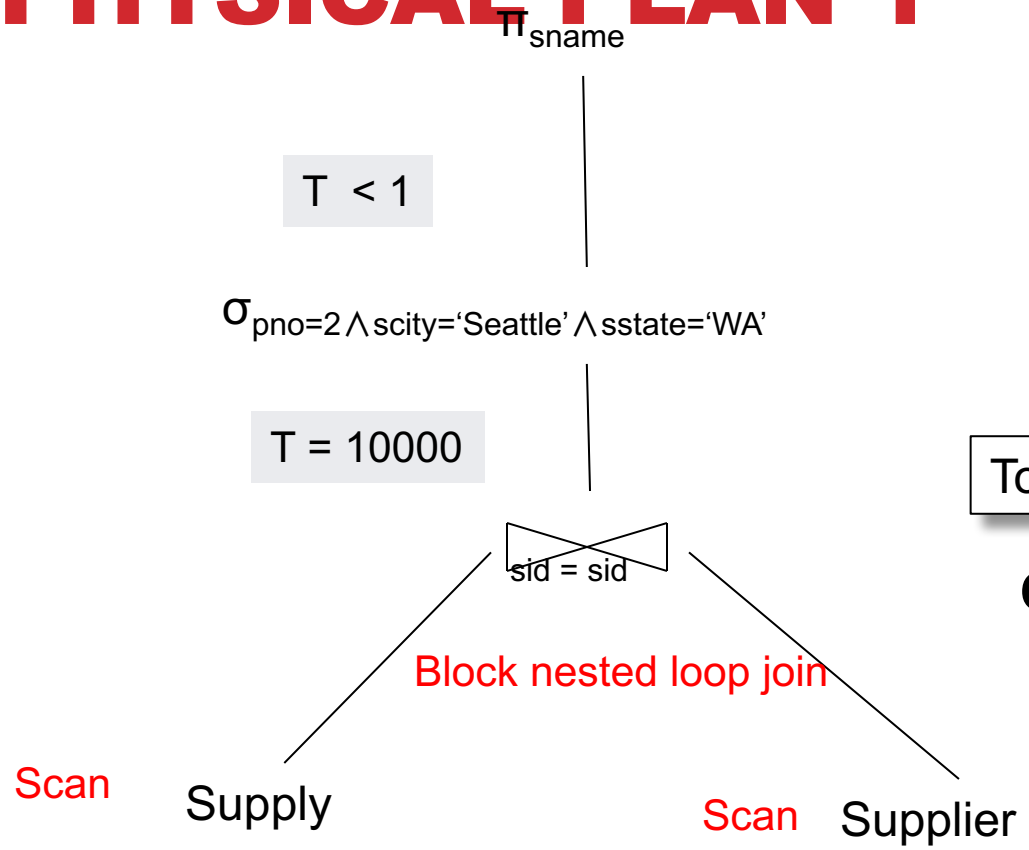
T(Supply) = 10000
B(Supply) = 100
V(Supply, pno) = 2500

T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

M=11

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

PHYSICAL PLAN 1



Total cost: $100 + 100 * 100 / 10 = 1100$

Cost: $B(R) + B(R)B(S)/(M-1)$

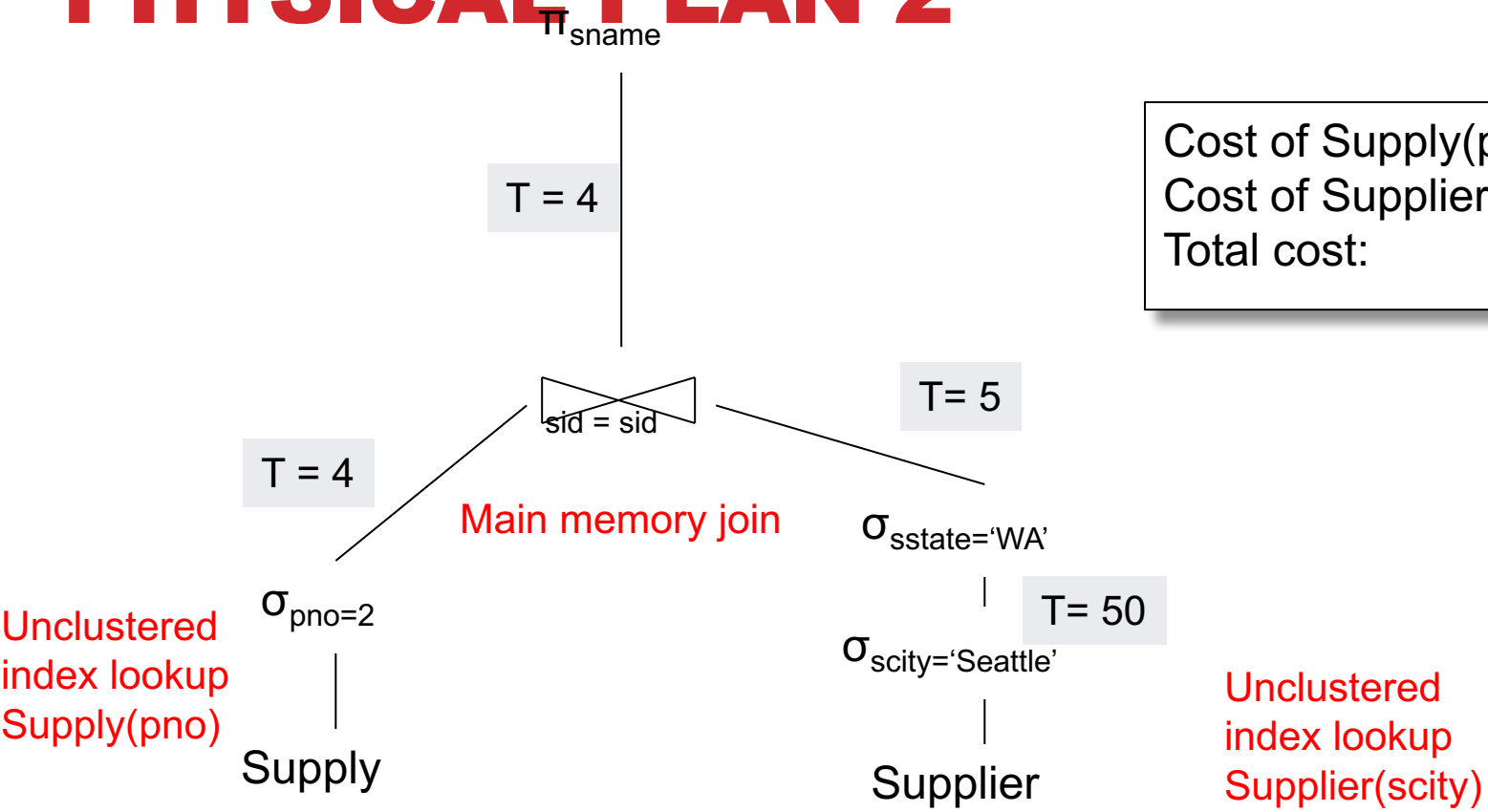
$T(\text{Supply}) = 10000$
 $B(\text{Supply}) = 100$
 $V(\text{Supply}, pno) = 2500$

$T(\text{Supplier}) = 1000$
 $B(\text{Supplier}) = 100$
 $V(\text{Supplier}, scity) = 20$
 $V(\text{Supplier}, state) = 10$

$M=11$

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

PHYSICAL PLAN 2



Cost of Supply(pno) =
 Cost of Supplier(scity) =
 Total cost:

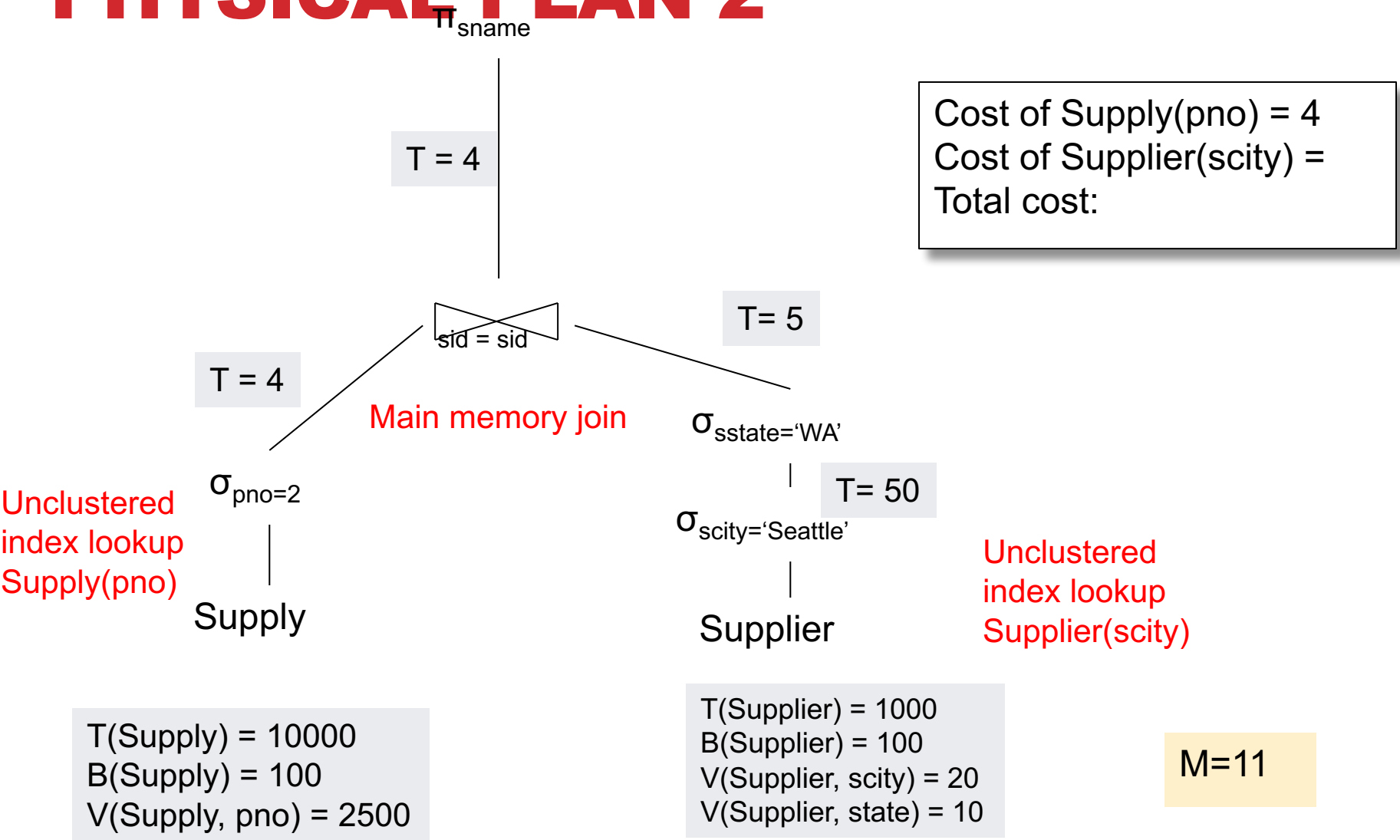
T(Supply) = 10000
 B(Supply) = 100
 V(Supply, pno) = 2500

T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

M=11

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

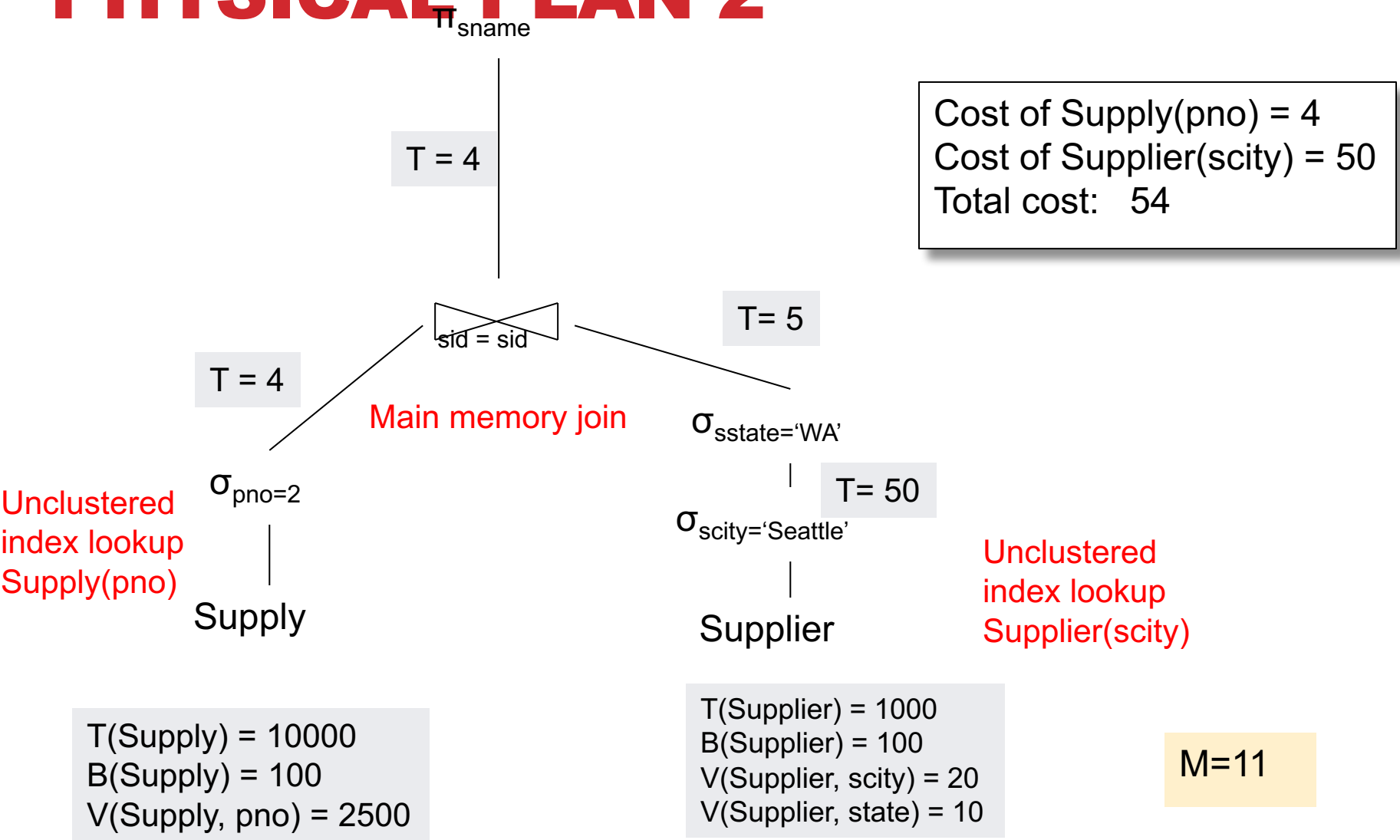
PHYSICAL PLAN 2



Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

PHYSICAL PLAN 2



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

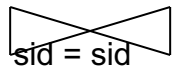
PHYSICAL PLAN 3

T = 4

Π_{sname}
 $\sigma_{scity='Seattle' \wedge sstate='WA'}$

T = 4

$\sigma_{pno=2}$
 Supply



Clustered
Index join

Supplier

Cost of Supply(pno) =
 Cost of Index join =
 Total cost:

Unclustered
index lookup
Supply(pno)

T(Supply) = 10000
 B(Supply) = 100
 V(Supply, pno) = 2500

T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

M=11

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

PHYSICAL PLAN 3

T = 4

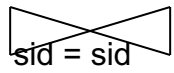
Π_{sname}
 $\sigma_{scity='Seattle' \wedge sstate='WA'}$

Cost of Supply(pno) = 4
 Cost of Index join =
 Total cost:

T = 4

Unclustered
 index lookup
 Supply(pno)

$\sigma_{pno=2}$
 Supply



Clustered
 Index join

Supplier

T(Supply) = 10000
 B(Supply) = 100
 V(Supply, pno) = 2500

T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

M=11

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

PHYSICAL PLAN 3

T = 4

Π_{sname}
 $\sigma_{scity='Seattle' \wedge sstate='WA'}$

Cost of Supply(pno) = 4
 Cost of Index join = 4
 Total cost: 8

T = 4

$\sigma_{pno=2}$
 Supply

sid = sid

Clustered
 Index join

Supplier

Unclustered
 index lookup
 Supply(pno)

T(Supply) = 10000
 B(Supply) = 100
 V(Supply, pno) = 2500

T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

M=11

QUERY OPTIMIZER SUMMARY

Input: A logical query plan

Output: A good physical query plan

Basic query optimization algorithm

- Enumerate alternative plans (logical and physical)
- Compute estimated cost of each plan
- Choose plan with lowest cost

This is called cost-based optimization

DISK SCHEDULING

- **Query optimization**
 - Good DB design
 - Good estimation
 - Hardware independent
- **All Disk I/Os are not created equal**
 - Sectors close to each other are more preferable to read