

# Introduction to Data Management

## CSE 414

11/2

Unit 4: RDBMS Internals  
Logical and Physical Plans  
Query Execution  
Query Optimization

(3 lectures)

# Introduction to Database Systems

## CSE 414

### Lecture 16:

### Basics of Data Storage and Indexes

# Query Performance

- My database application is too slow... why?
  - One of the queries is very slow... why?
  - To understand performance, we need to understand:
    - How is data organized on disk
    - How to estimate query costs
- In this course we will focus on **disk-based DBMSs**

extracting information from the disk is by far the slowest action, which is why we focus our look into efficiency by focusing how many disk reads are required.

why use a disk at all? cheapest option for storing large quantities of data

drawbacks

reader physically has to move to location on disk

read speed limited by rotation speed (faster to read things sequentially, or a block at a time)

CSE 414 - Autumn 2018

3

# Scanning a Data File

- Disks are mechanical devices!
  - Technology from the 60s;
  - Density increases over time
- Read only at the rotation speed!
- Consequence: sequential scan faster than random
  - **Good:** read blocks 1,2,3,4,5,...
  - **Bad:** read blocks 2342, 11, 321,9, ...
- **Rule of thumb:**
  - Random read 1-2% of file  $\approx$  sequential scan entire file;
  - 1-2% decreases over time, because of increased density
- Solid state (SSD): still too expensive today



why use a disk at all? cheapest option for storing large quantities of data

drawbacks

reader physically has to move to location on disk (MOST expensive operation)

read speed limited by rotation speed (

faster to read things sequentially, or a block at a time

good: read blocks 1,2,3,4,5

bad : read blocks 2342, 11, 321, 9 ...

to put into perspective:

Read 1 MB from memory (RAM) 250 us

Disk seek (i.e. move disk arm/reader) = 10,000 us

Read 1MB sequentially from disk (including disk seek) = 20,000 us

100 times slower!

# Data Storage

- DBMSs store data in **files**
- Most common organization is **row-wise storage**
- On disk, a file is split into **blocks**
- Each block contains a set of tuples

ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

10	Tom	Hanks
20	Amy	Hanks
50	...	...
200	...	
220		
240		
420		
800		

In the example, we have **4 blocks** with 2 tuples each



ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

# Data File Types

The data file can be one of:

- **Heap file**
  - Unsorted      put tuples completely random in file
- **Sequential file**
  - Sorted according to some attribute(s) called key

order tuples based on some sort of key.

# Index

- An **additional** file, that allows fast access to records in the data file given a search key
- The index contains (key, value) pairs:
  - The key = an attribute value (e.g., student ID or name)
  - The value = a pointer to the record
- Could have many indexes for one table

tuples; allows you to quickly access specific pieces of memory rather than scanning thru whole disk

Key = means here search key

# This Is Not A Key



Different keys:

- Primary key – uniquely identifies a tuple based on schema
- Key of the sequential file – how the data file is sorted, if at all
- Index key – how the index is organized



CSE 414 - Autumn 2018



9

ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

# Example 1: Index on ID

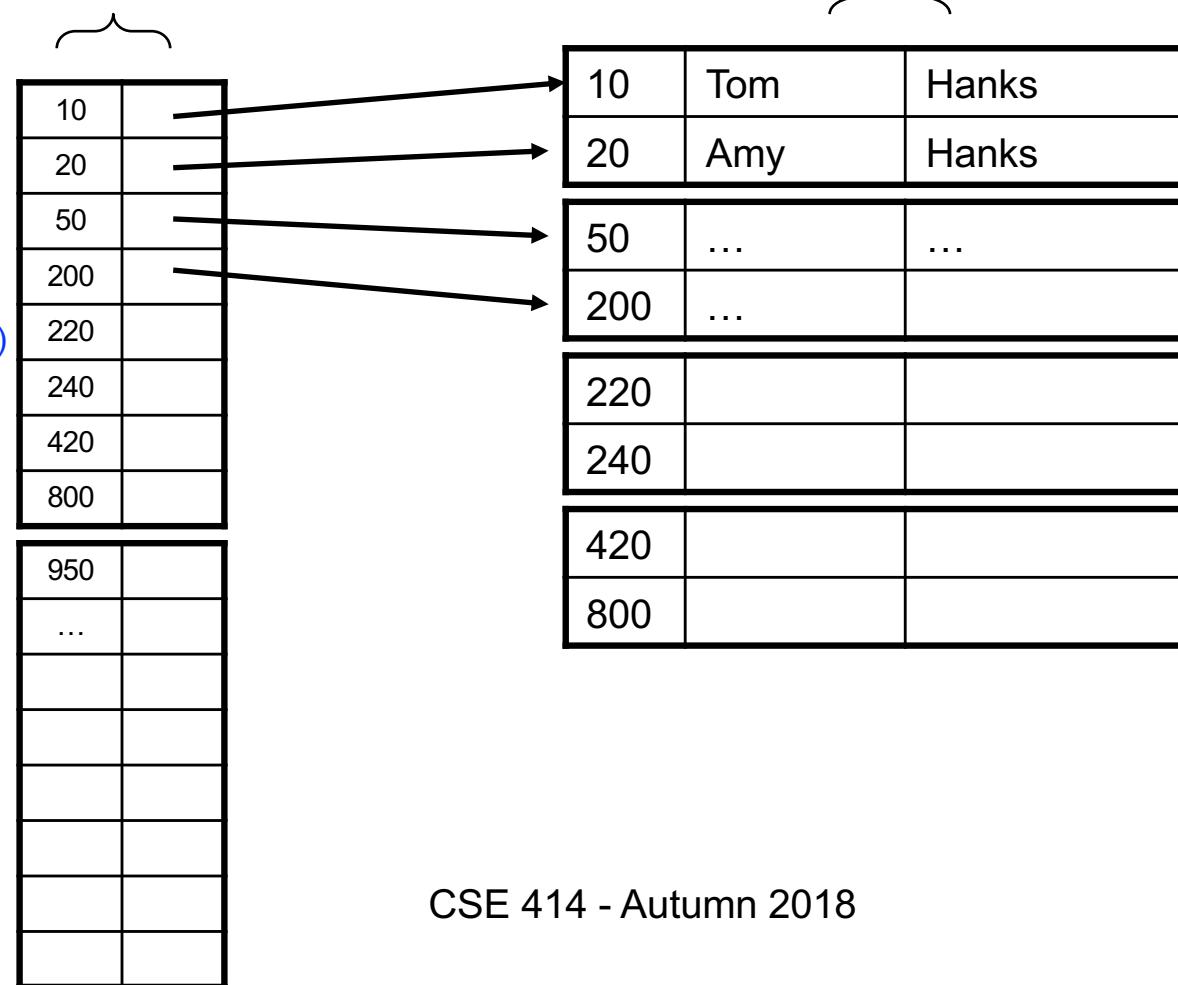
Index the student table on the Student.ID

Index **Student\_ID** on **Student.ID**

Data File **Student**

more indexs  
can fit in a  
block than full  
tuples; one  
primary advantage

can even store it in  
main memory as its  
so small (sometimes)  
giving fast acesses



effectiveness of index depends on how tuples are ordered on disk as well!

## Student

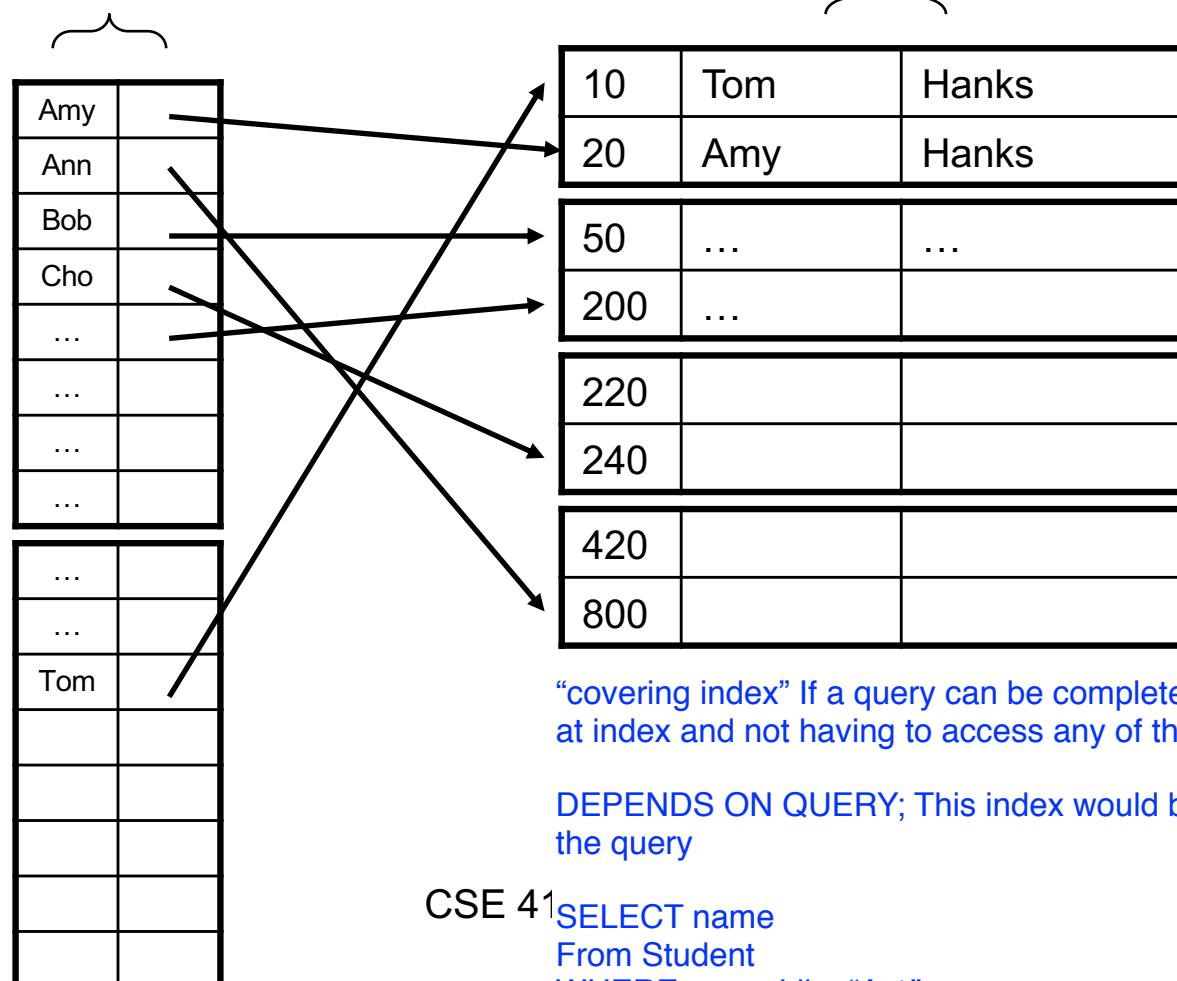
Unclustered index: tuples of data are not stored in file in an order based on the index;  
index less effective?

# Example 2: Index on fName

Index **Student\_fName**  
on **Student.fName**

Data File **Student**

ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		



"covering index" If a query can be completely answered just looking at index and not having to access any of the tuples in main memory

DEPENDS ON QUERY; This index would be a covering index for the query

CSE 41

SELECT name  
From Student  
WHERE name Like "A%"

# Index Organization

- Hash table
- B+ trees – most common
  - They are search trees, but they are not binary instead have higher fan-out
  - Will discuss them briefly next
- Specialized indexes: bit maps, R-trees, inverted index

## Student

ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

# Hash table example

Index **Student\_ID** on **Student.ID**



10	
20	
50	
200	
220	
240	
420	
800	
...	...
...	...

Data File **Student**



10	Tom	Hanks
20	Amy	Hanks
50	...	...
200	...	
220		
240		
420		
800		

Index File  
(preferably  
in memory)

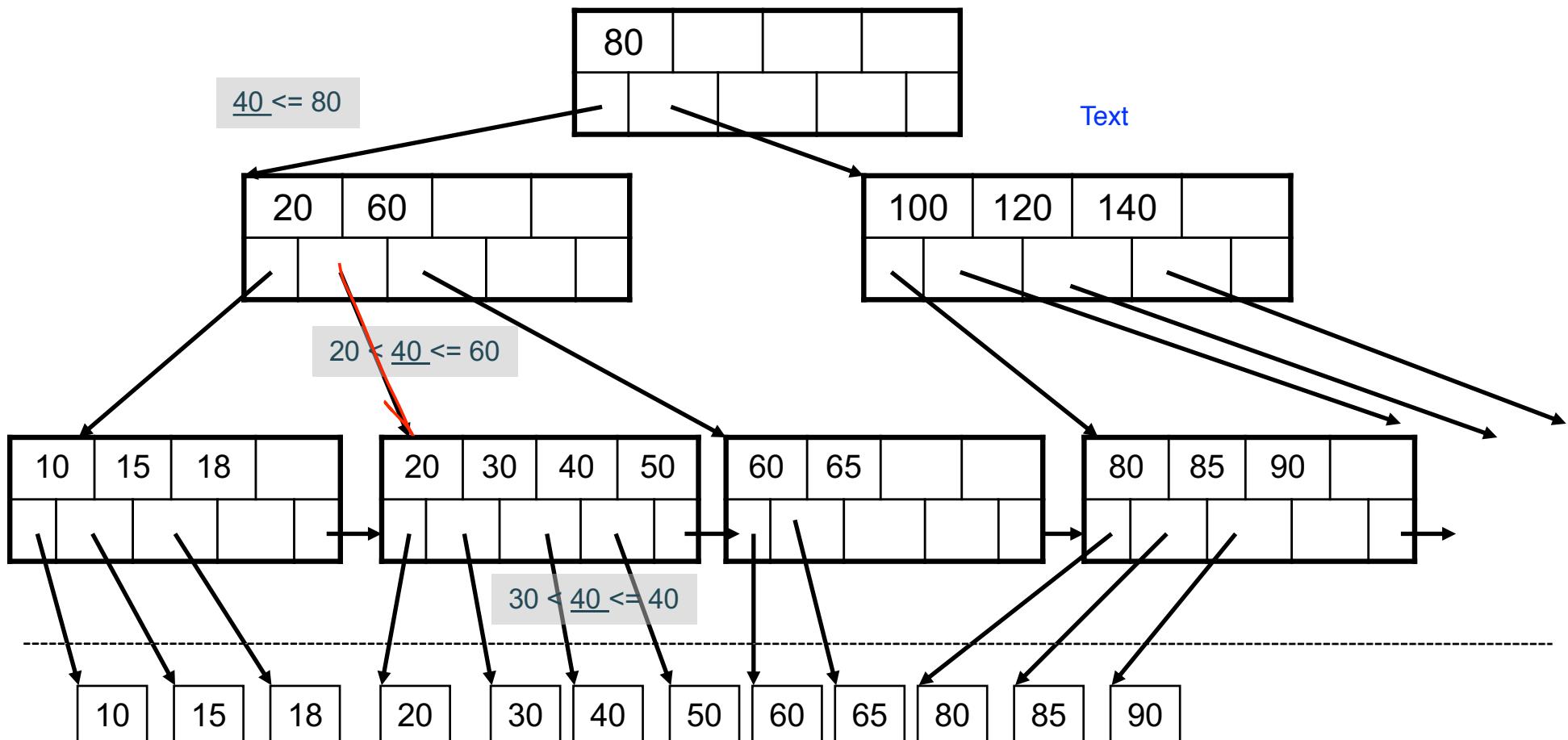
Data file  
(on disk)

# B+ Tree Index by Example

$d = 2$

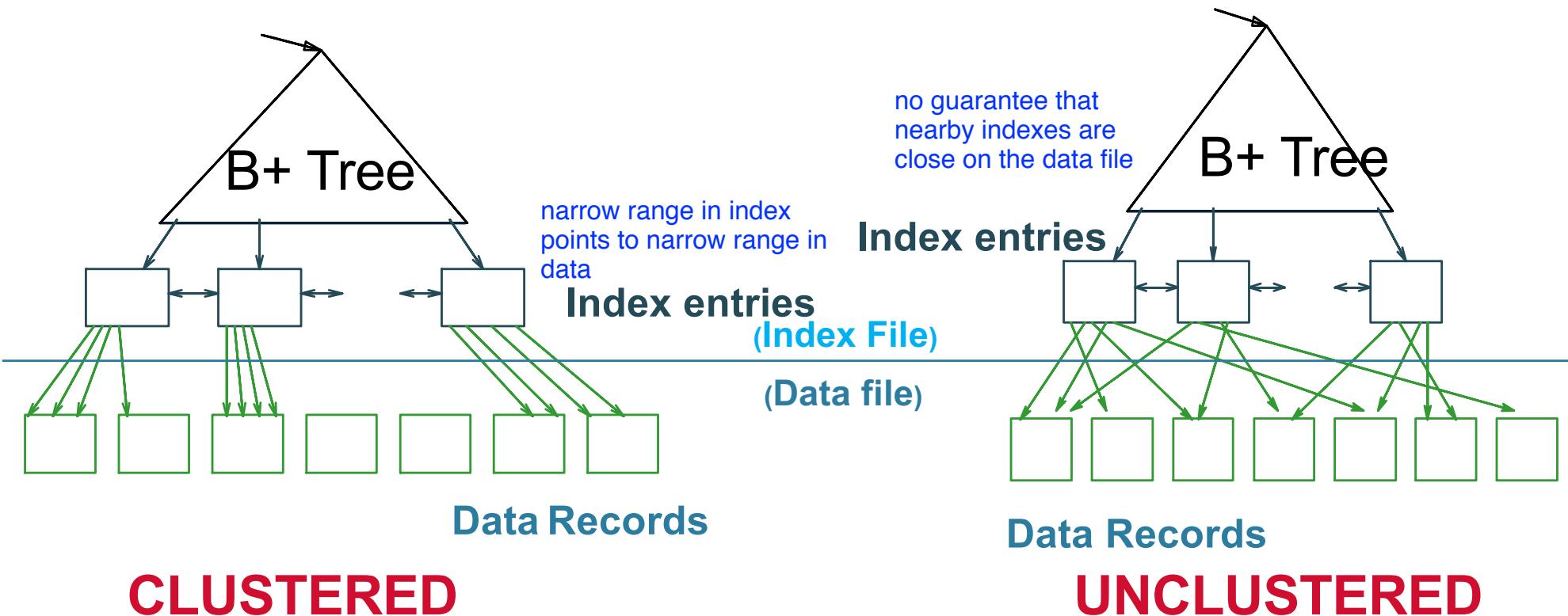
find the node that is in between the two values above it

Find the key 40



For clustered index, tuples (in data) are ordered based on the index

# Clustered vs Unclustered



Every table can have **only one** clustered and **many** unclustered indexes

because the data is ordered based on the clustered index; cant order on two separate qualities

Why?

unclustered indexes dont constrain organization of data at all, so we can have as many as we want

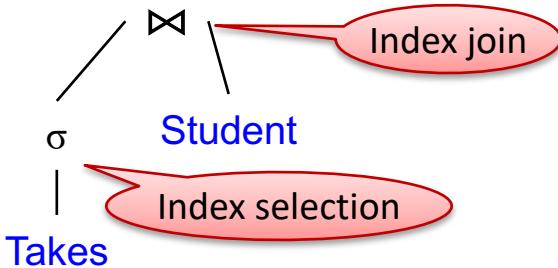
# Index Classification

- **Clustered/unclustered**
  - Clustered = records close in index are close in data
    - Option 1: Data inside data file is sorted on disk
    - Option 2: Store data directly inside the index (no separate files)
  - Unclustered = records close in index may be far in data
- **Primary/secondary**
  - Meaning 1: index over primary key; fastest look up by ordering data on disk by primary key
    - Primary = is over attributes that include the primary key
    - Secondary = otherwise
  - Meaning 2: means the same as clustered/unclustered
- **Organization** B+ tree or Hash table

# Summary So Far

- Index = a file that enables direct access to records in another data file
  - B+ tree / Hash table
  - Clustered/unclustered
- Data resides on disk
  - Organized in blocks
  - Sequential reads are efficient
  - Random access less efficient
  - Random read 1-2% of data worse than sequential

```
Student(ID, fname, lname)  
Takes(studentID, courseID)
```



```
for y in Takes  
  if courseID > 300 then  
    for x in Student  
      if x.ID=y.studentID  
        output *
```

faster to search through indexes and filter tuples as index files exist in main memory which can be accessed much faster;  
Limits the number of disk accesses we have to run

Index join

```
SELECT *  
FROM Student x, Takes y  
WHERE x.ID=y.studentID AND y.courseID > 300
```

# Example

indexes for speeding up queries!

Assume the database has indexes on these attributes:

- **Takes\_courseID** = index on Takes.courseID
- **Student\_ID** = index on Student.ID

Index selection

```
for y' in Takes_courseID where y'.courseID > 300  
  y = fetch the Takes record pointed to by y'  
  for x' in Student_ID where x'.ID = y.studentID  
    x = fetch the Student record pointed to by x'  
    output *
```

# Getting Practical: Creating Indexes in SQL

```
CREATE TABLE V(M int, N varchar(20), P int);
```

```
CREATE INDEX V1 ON V(N)
```

CREATE INDEX index name ON Relation(attribute)

```
CREATE INDEX V2 ON V(P, M)
```

With multiple attributes, order indices by first attribute, then for matching values of that order by second attribute

```
CREATE INDEX V3 ON V(M, N)
```

```
CREATE UNIQUE INDEX V4 ON V(N)
```

```
CREATE CLUSTERED INDEX V5 ON V(N)
```

# Getting Practical: Creating Indexes in SQL

```
CREATE TABLE V(M int, N varchar(20), P int);
```

```
CREATE INDEX V1 ON V(N)
```

```
CREATE INDEX V2 ON V(P, M)
```

What does this mean?

```
CREATE INDEX V3 ON V(M, N)
```

```
CREATE UNIQUE INDEX V4 ON V(N)
```

```
CREATE CLUSTERED INDEX V5 ON V(N)
```

# Getting Practical: Creating Indexes in SQL

```
CREATE TABLE V(M int, N varchar(20), P int);
```

yes

```
CREATE INDEX V1 ON V(N)
```

```
select *  
from V  
where P=55 and M=77
```

```
CREATE INDEX V2 ON V(P, M)
```

What does this mean?

```
CREATE INDEX V3 ON V(M, N)
```

```
CREATE UNIQUE INDEX V4 ON V(N)
```

```
CREATE CLUSTERED INDEX V5 ON V(N)
```

# Getting Practical: Creating Indexes in SQL

```
CREATE TABLE V(M int, N varchar(20), P int);
```

yes

```
CREATE INDEX V1 ON V(N)
```

```
select *  
from V  
where P=55 and M=77
```

```
CREATE INDEX V2 ON V(P, M)
```

What does this mean?

```
CREATE INDEX V3 ON V(M, N)
```

```
select *  
from V  
where P=55
```

```
CREATE UNIQUE INDEX V4 ON V(N)
```

```
CREATE CLUSTERED INDEX V5 ON V(N)
```

# Getting Practical: Creating Indexes in SQL

```
CREATE TABLE V(M int, N varchar(20), P int);
```

yes

```
CREATE INDEX V1 ON V(N)
```

```
select *  
from V  
where P=55 and M=77
```

```
CREATE INDEX V2 ON V(P, M)
```

What does this mean?

```
CREATE INDEX V3 ON V(M, N)
```

```
select *  
from V  
where P=55
```

```
CREATE UNIQUE INDEX V4 ON V(N)
```

yes

```
CREATE CLUSTERED INDEX V5 ON V(N)
```

# Getting Practical: Creating Indexes in SQL

```
CREATE TABLE V(M int, N varchar(20), P int);
```

yes

```
CREATE INDEX V1 ON V(N)
```

```
select *  
from V  
where P=55 and M=77
```

```
CREATE INDEX V2 ON V(P, M)
```

What does this mean?

```
CREATE INDEX V3 ON V(M, N)
```

```
select *  
from V  
where P=55
```

```
CREATE UNIQUE INDEX V4 ON V(N)
```

```
select *  
from V  
where M=77
```

```
CREATE CLUSTERED INDEX V5 ON V(N)
```

# Getting Practical: Creating Indexes in SQL

```
CREATE TABLE V(M int, N varchar(20), P int);
```

yes

```
CREATE INDEX V1 ON V(N)
```

```
select *  
from V  
where P=55 and M=77
```

```
CREATE INDEX V2 ON V(P, M)
```

What does this mean?

```
CREATE INDEX V3 ON V(M, N)
```

```
select *  
from V  
where P=55
```

```
CREATE UNIQUE INDEX V4 ON V(N)
```

yes

```
CREATE CLUSTERED INDEX V5 ON V(N)
```

no

forces data to be stored/sorted based on that given index

CSE 414 - Autumn 2018

M values are spread far  
throughout the indices  
file as the file is first  
sorted by P, then by M. So M = 77  
could occur in many places

35

wont speed up  
query

ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

# Which Indexes?

- How many indexes **could** we create?  
index files take up extra space...  
updating tuples and we also have to update index file
- Which indexes **should** we create?

In general this is a very hard problem

ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

# Which Indexes?

- The *index selection problem*
  - Given a table, and a “workload” (big Java application with lots of SQL queries), decide which indexes to create (and which ones NOT to create!)
- Who does index selection:
  - The database administrator DBA
    - person
  - Semi-automatically, using a database administration tool

ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

# Which Indexes?

- The *index selection problem*
  - Given a table, and a “workload” (big Java application with lots of SQL queries), decide which indexes to create (and which ones NOT to create!)
- Who does index selection:
  - The database administrator DBA
  - Semi-automatically, using a database administration tool



# Index Selection: Which Search Key

- Make some attribute K a search key if the WHERE clause contains:
  - An exact match on K
  - A range predicate on K
  - A join on K

# The Index Selection Problem 1

$V(M, N, P);$

Your workload is this

100000 queries:

```
SELECT *
FROM V
WHERE N=?
```

100 queries:

```
SELECT *
FROM V
WHERE P=?
```

# The Index Selection Problem 1

$V(M, N, P);$

Your workload is this

100000 queries:

```
SELECT *
FROM V
WHERE N=?
```

100 queries:

```
SELECT *
FROM V
WHERE P=?
```

What indexes ?

# The Index Selection Problem 1

$V(M, N, P);$

Your workload is this

100000 queries:

```
SELECT *
FROM V
WHERE N=?
```

100 queries:

```
SELECT *
FROM V
WHERE P=?
```

A:  $V(N)$  and  $V(P)$  (hash tables or B-trees)

make both

CSE 414 - Autumn 2018

44

# The Index Selection Problem 2

V(M, N, P);

Your workload is this

100000 queries:

```
SELECT *  
FROM V  
WHERE N>? and N<?
```

100 queries:

```
SELECT *  
FROM V  
WHERE P=?
```

100000 queries:

```
INSERT INTO V  
VALUES (?, ?, ?)
```

What indexes ?

# The Index Selection Problem 2

$V(M, N, P);$

Your workload is this

100000 queries:

**SELECT \***  
**FROM V**  
**WHERE N>? and N<?**

100 queries:

**SELECT \***  
**FROM V**  
**WHERE P=?**

100000 queries:

**INSERT INTO V**  
**VALUES (?, ?, ?)**

need to enforced order of a B-tree; easier to find starting and end points of a range (close index values are stored near to each other)

A: definitely  $V(N)$  (must B-tree); unsure about  $V(P)$

# The Index Selection Problem 3

V(M, N, P);

Your workload is this

100000 queries: 1000000 queries: 100000 queries:

**SELECT \***  
**FROM V**  
**WHERE N=?**

**SELECT \***  
**FROM V**  
**WHERE N=? and P>?**

**INSERT INTO V**  
**VALUES (?, ?, ?)**

What indexes ?

# The Index Selection Problem 3

$V(M, N, P);$

Your workload is this

100000 queries:    1000000 queries:    100000 queries:

**SELECT \***  
**FROM V**  
**WHERE N=?**

**SELECT \***  
**FROM V**  
**WHERE N=? and P>?**

**INSERT INTO V  
VALUES (?, ?, ?)**

**A:  $V(N, P)$**

Limit our query first to N values (as this is the more exclusive condition) and we also need to sort by P; also B+ tree  
Then all the P > conditions will be next to each other

By pairing N and P together, we do not have to check tuples twice.  
Narrowing down tuples for N, then further narrow down on P rather than narrow on N, narrow on P separately, then join

How does this index differ from:  
1. Two indexes  $V(N)$  and  $V(P)$ ?  
2. An index  $V(P, N)$ ?

# The Index Selection Problem 4

V(M, N, P);

Your workload is this  
1000 queries:

```
SELECT *  
FROM V  
WHERE N>? and N<?
```

100000 queries:

```
SELECT *  
FROM V  
WHERE P>? and P<?
```

What indexes ?

# The Index Selection Problem 4

$V(M, N, P);$

Your workload is this  
1000 queries:

```
SELECT *  
FROM V  
WHERE N>? and N<?
```

100000 queries:

```
SELECT *  
FROM V  
WHERE P>? and P<?
```

A:  $V(N)$  secondary,  $V(P)$  primary index

# Two typical kinds of queries

```
SELECT *
FROM Movie
WHERE year = ?
```

```
SELECT *
FROM Movie
WHERE year >= ? AND
      year <= ?
```

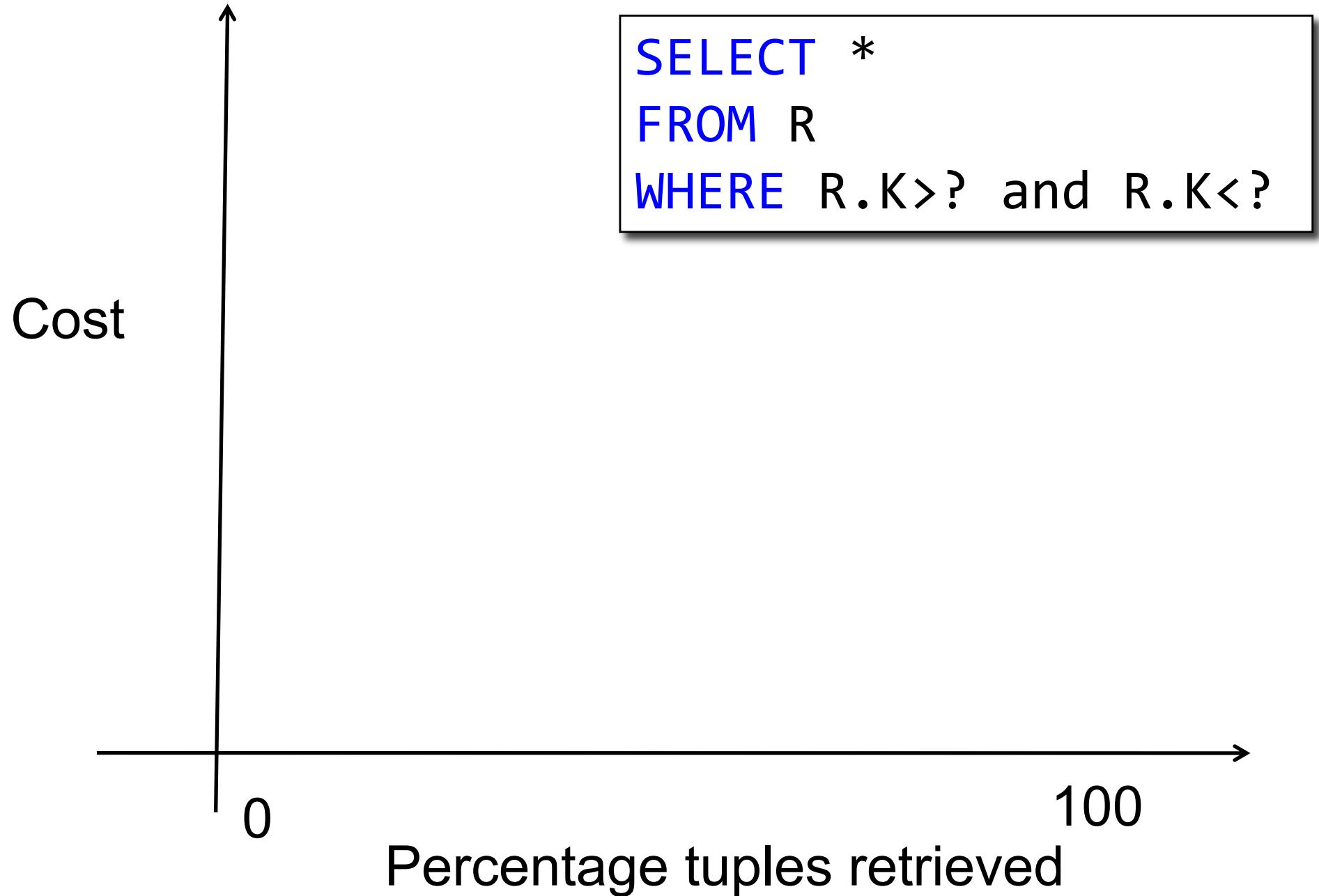
- Point queries
- What data structure should be used for index?
- Range queries
- What data structure should be used for index?

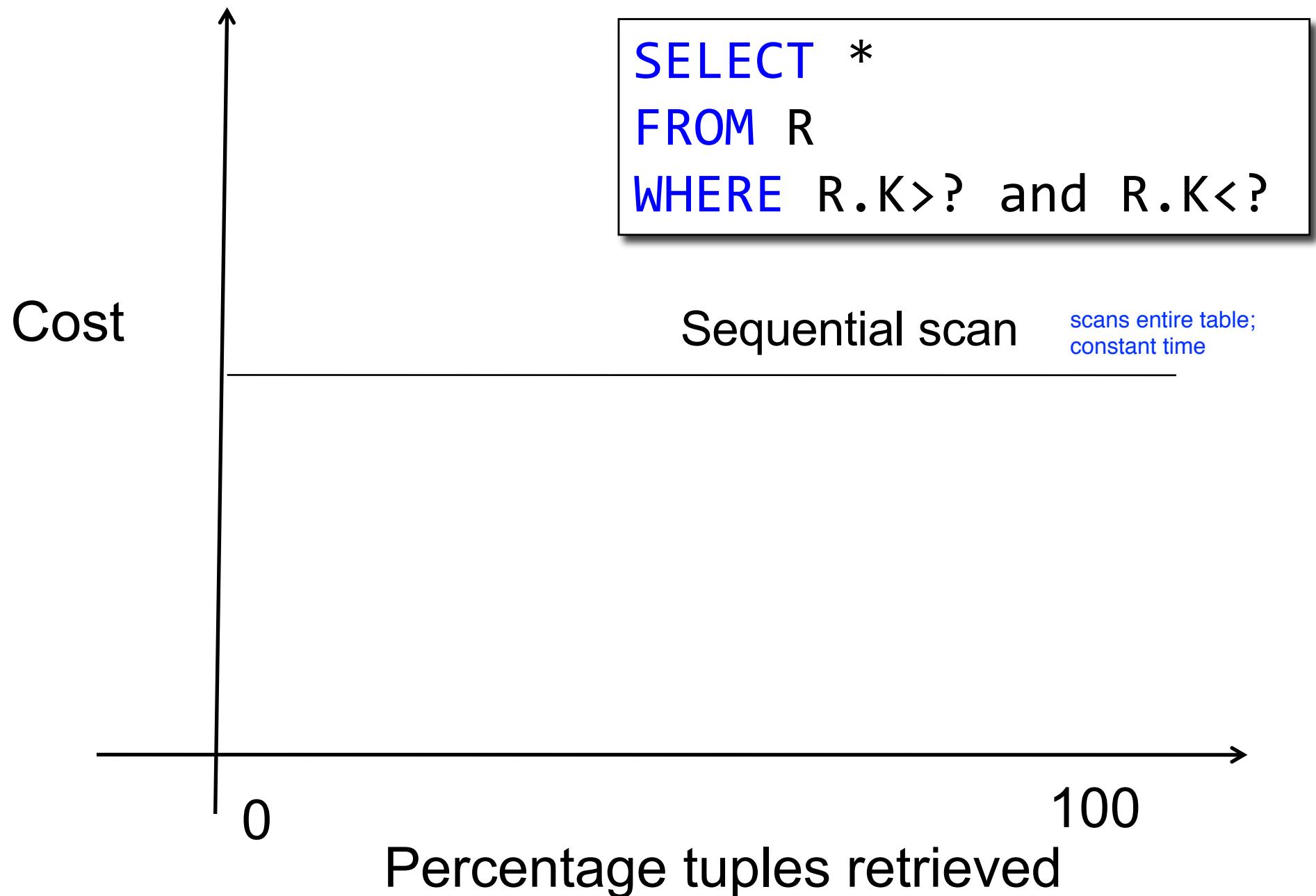
# Basic Index Selection Guidelines

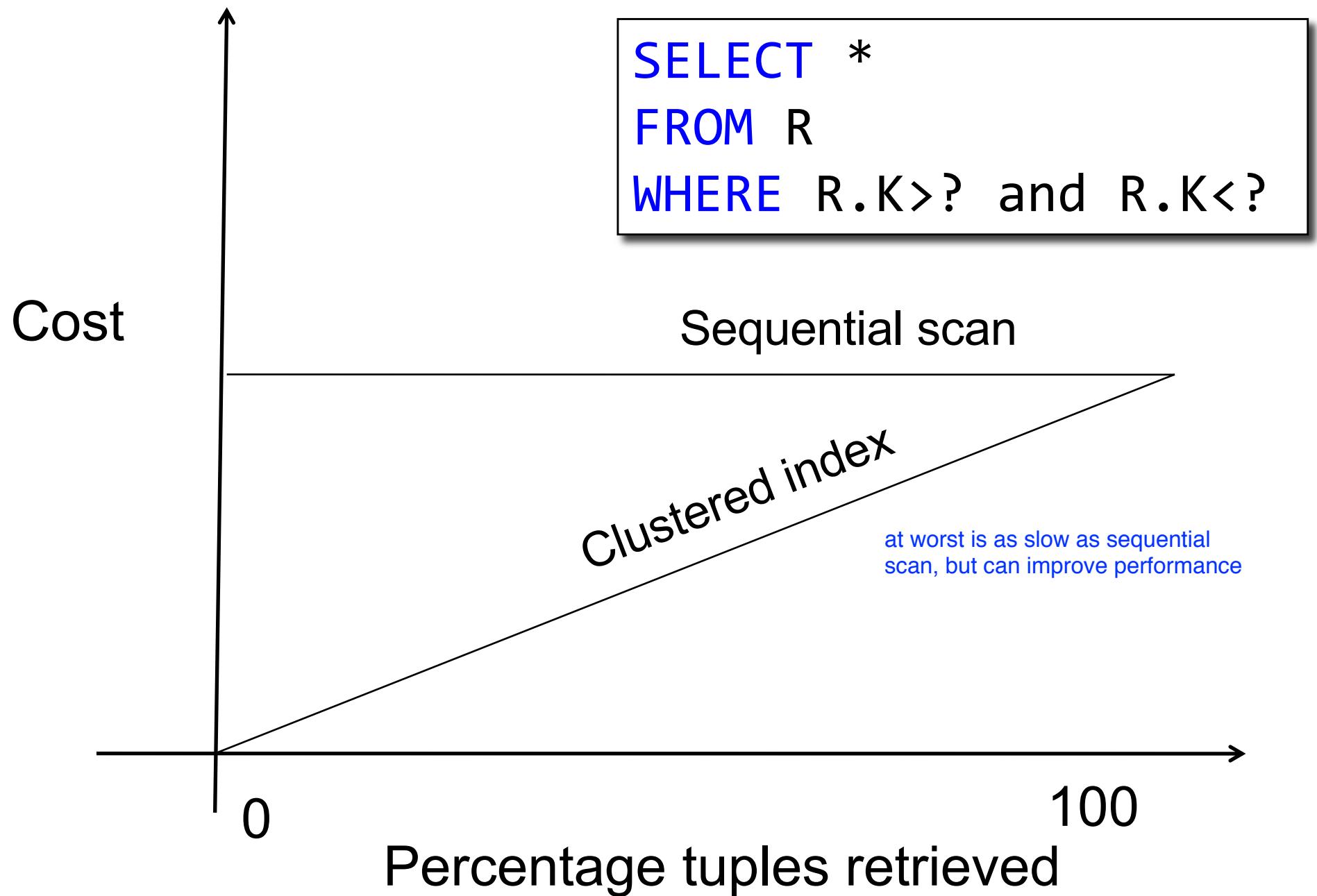
- Consider queries in workload in order of importance
- Consider relations accessed by query
  - No point indexing other relations *i.e. attributes not accessed by query.*
- Look at WHERE clause for possible search key
- Try to choose indexes that speed-up multiple queries

# To Cluster or Not

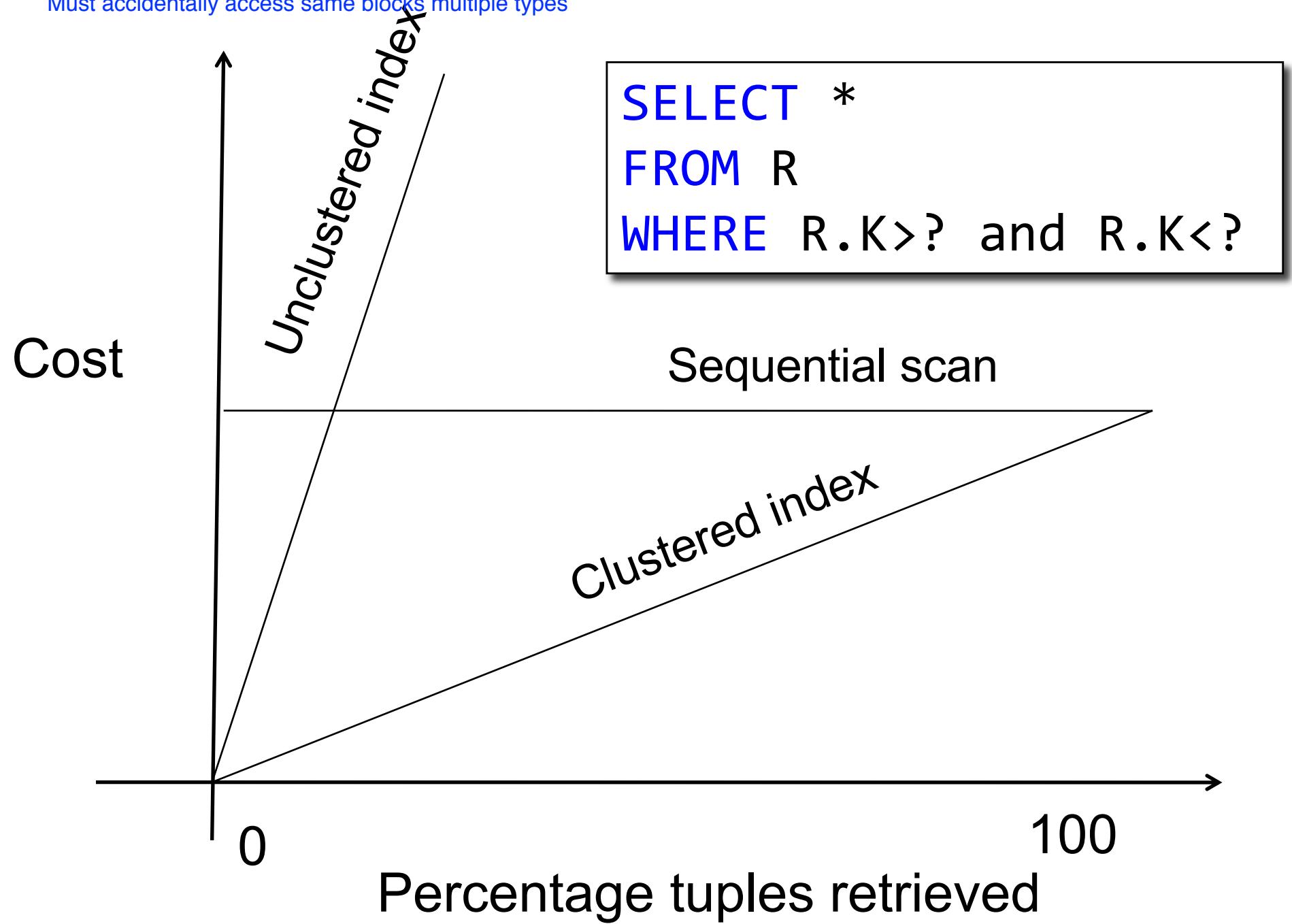
- Range queries benefit mostly from clustering
- Covering indexes do *not* need to be clustered: they work equally well unclustered







possibly worse than sequential  
Must accidentally access same blocks multiple types



# Introduction to Database Systems

## CSE 344

### Lecture 17:

#### Basics of Query Optimization and Query Cost Estimation

# Choosing Index is Not Enough

- To estimate the cost of a query plan, we still need to consider other factors:
  - How each operator is implemented
  - The cost of each operator
  - Let's start with the basics

# Cost of Reading Data From Disk

# Cost Parameters

- Cost = I/O + CPU + Network BW
  - We will focus on I/O in this class
- Parameters (a.k.a. statistics):
  - $B(R)$  = # of blocks (i.e., pages) for relation R
  - $T(R)$  = # of tuples in relation R
  - $V(R, a)$  = # of distinct values of attribute a

# Cost Parameters

- Cost = I/O + CPU + Network BW
  - We will focus on I/O in this class
- Parameters (a.k.a. statistics):
  - $B(R)$  = # of blocks (i.e., pages) for relation R
  - $T(R)$  = # of tuples in relation R
  - $V(R, a)$  = # of distinct values of attribute a

When  $a$  is a key,  $V(R,a) = T(R)$

When  $a$  is not a key,  $V(R,a)$  can be anything  $\leq T(R)$

# Cost Parameters

- Cost = I/O + CPU + Network BW
  - We will focus on I/O in this class
- Parameters (a.k.a. statistics):
  - $B(R)$  = # of blocks (i.e., pages) for relation R
  - $T(R)$  = # of tuples in relation R
  - $V(R, a)$  = # of distinct values of attribute a

When  $a$  is a key,  $V(R,a) = T(R)$

When  $a$  is not a key,  $V(R,a)$  can be anything  $\leq T(R)$

- DBMS collects **statistics** about base tables  
must infer them for intermediate results

# Selectivity Factors for Conditions

- $A = c$  /\*  $\sigma_{A=c}(R)$  \*/
  - Selectivity =  $1/V(R, A)$
- $A < c$  /\*  $\sigma_{A < c}(R)$  \*/
  - Selectivity =  $(c - \min(R, A)) / (\max(R, A) - \min(R, A))$
- $c_1 < A < c_2$  /\*  $\sigma_{c_1 < A < c_2}(R)$  \*/
  - Selectivity =  $(c_2 - c_1) / (\max(R, A) - \min(R, A))$

# Cost of Reading Data From Disk

- Sequential scan for relation R costs **B(R)**
- Index-based selection
  - Estimate selectivity factor  $f$  (see previous slide)
  - Clustered index:  $f^* \mathbf{B(R)}$
  - Unclustered index  $f^* \mathbf{T(R)}$

Note: we ignore I/O cost for index pages

# 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:
- Index based selection:

# 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:

# 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:
  - If index is unclustered:

# Index Based Selection

- Example:

$B(R) = 2000$
$T(R) = 100,000$
$V(R, a) = 20$
- 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:

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

# Index Based Selection

- Example:

$B(R) = 2000$
$T(R) = 100,000$
$V(R, a) = 20$
- 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

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

# 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 !

# Cost of Executing Operators (Focus on Joins)

# 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'

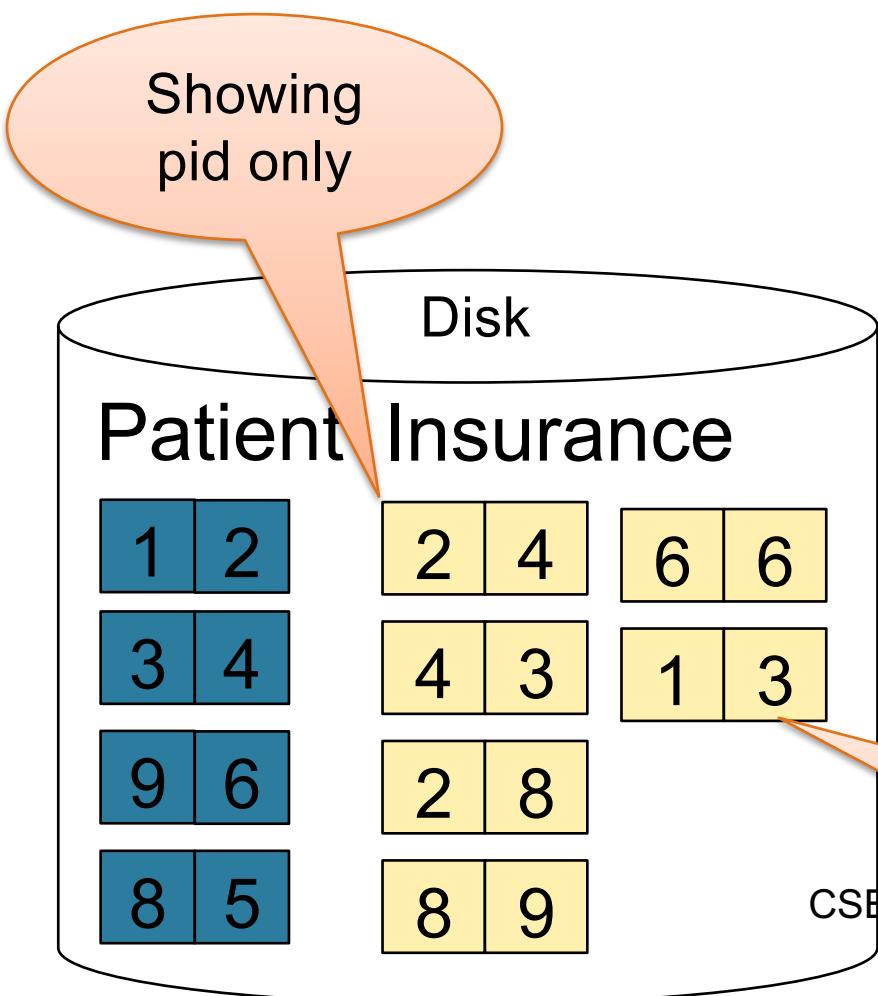
Insurance

2	'Blue'	123
4	'Prem'	432

4	'Prem'	343
3	'GrpH'	554

# Hash Join Example

Patient  $\bowtie$  Insurance

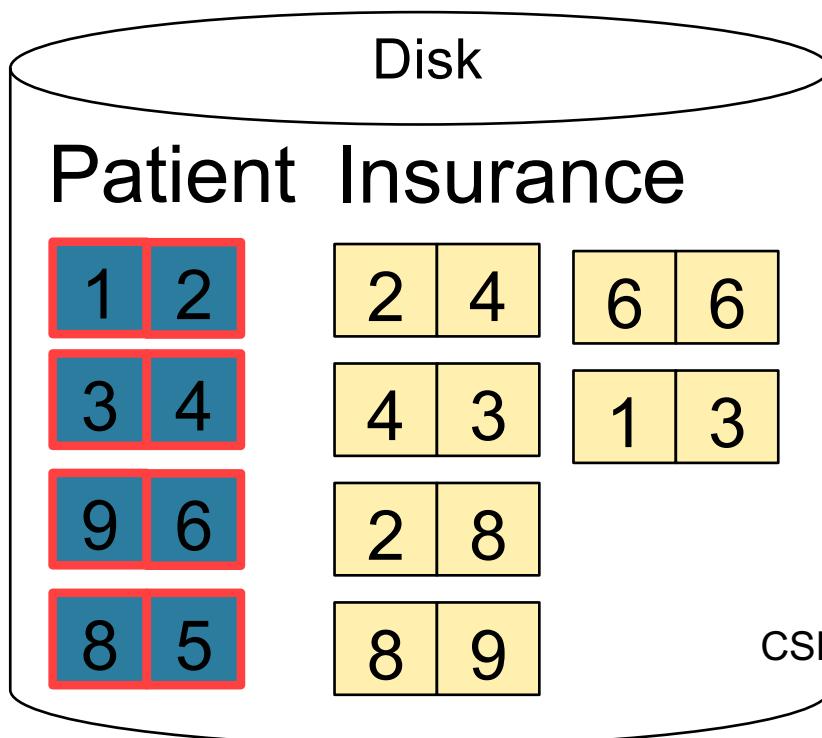


Memory M = 21 pages

Some large-enough #

# 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
---	--	---	---	---	--	---	---	---	---

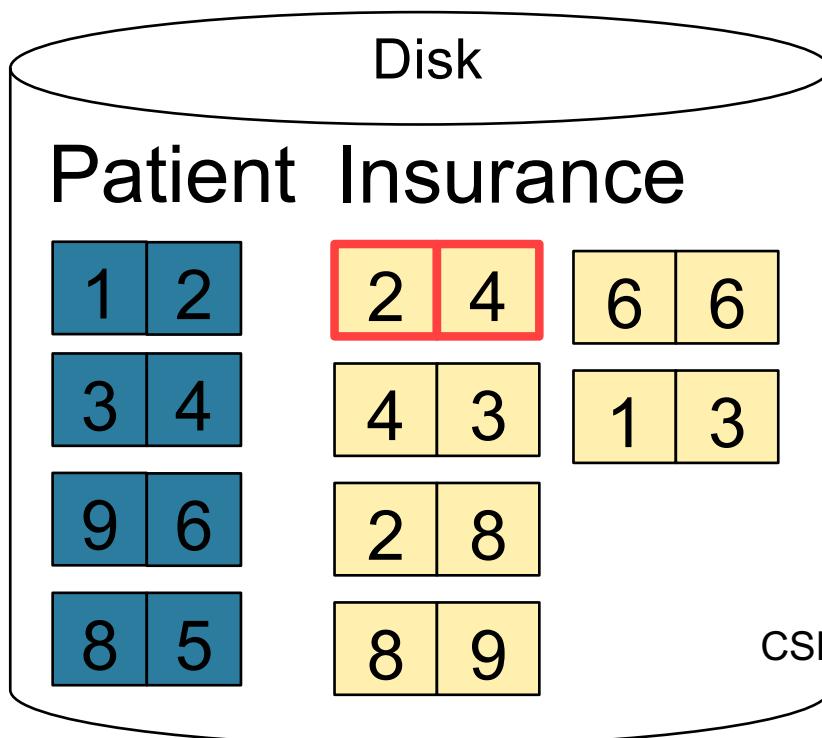


Input buffer

# 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
---	--	---	---	---	--	---	---	---	---

2	4
---	---

Input buffer

2	2
---	---

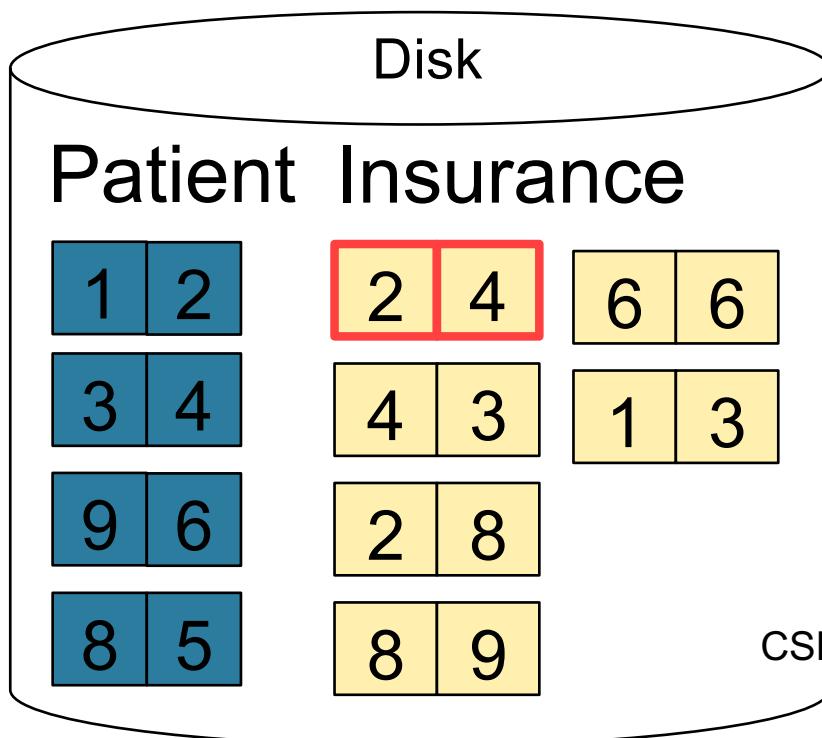
Output buffer

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
---	--	---	---	---	--	---	---	---	---

2	4
---	---

Input buffer

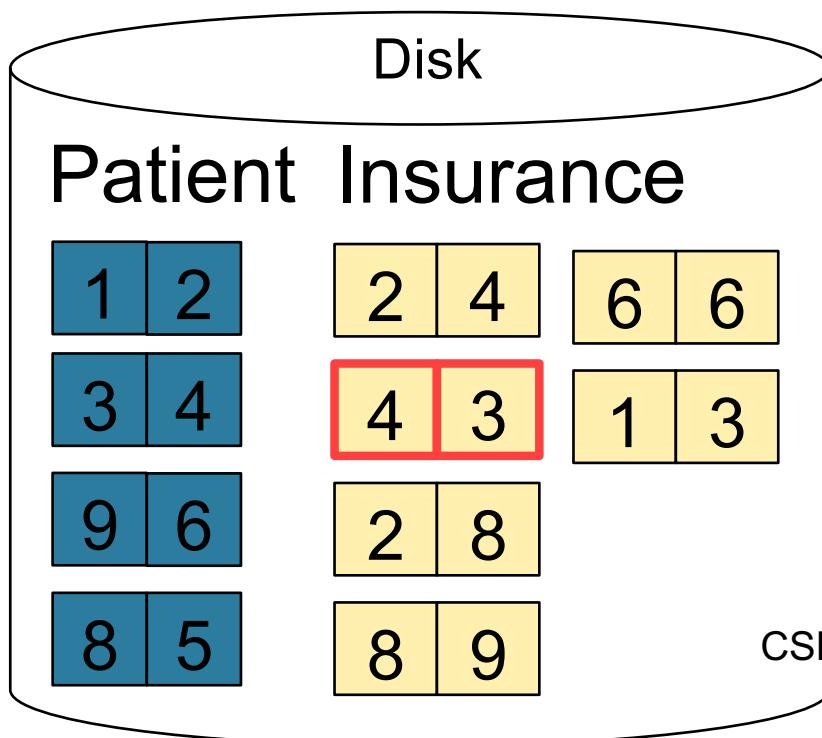
4	4
---	---

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
---	--	---	---	---	--	---	---	---	---

4	3
---	---

Input buffer

4	4
---	---

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 t1 in R do  
  for each tuple t2 in S do  
    if t1 and t2 join then output (t1,t2)
```

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 t1 in R do  
  for each tuple t2 in S do  
    if t1 and t2 join then output (t1,t2)
```

- 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

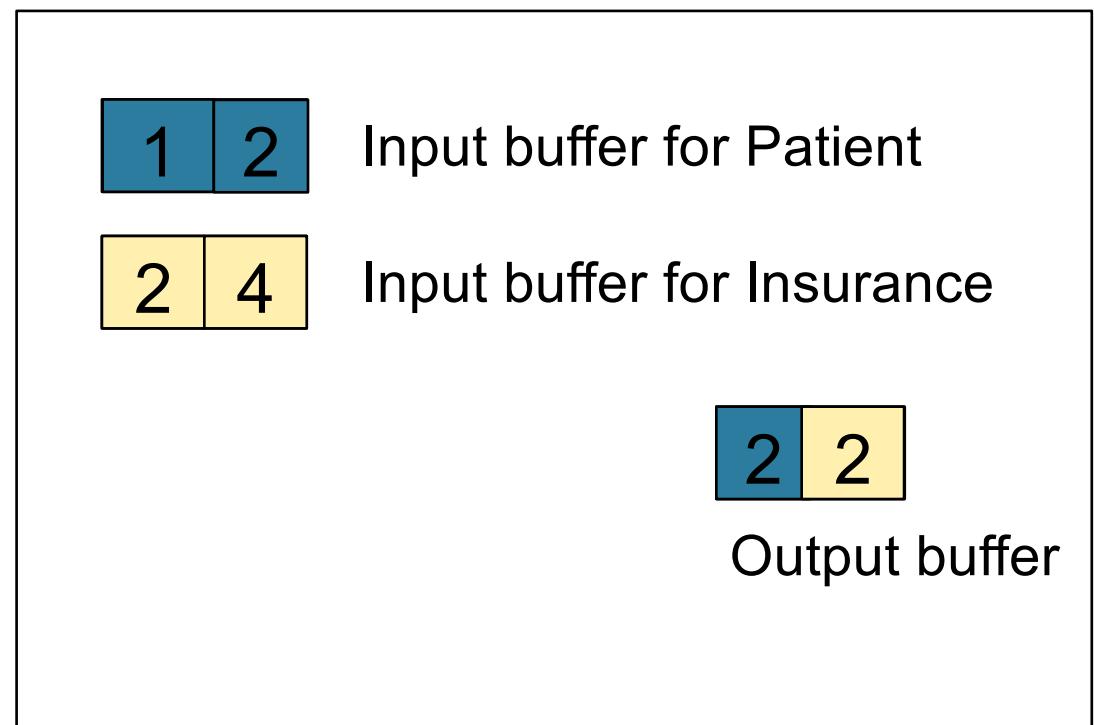
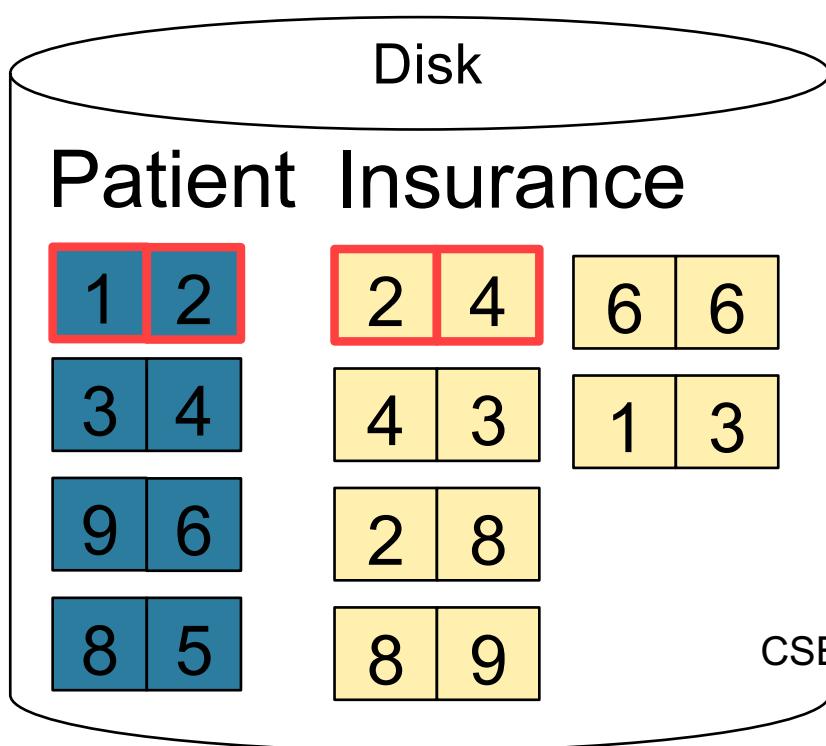
for all pairs of tuples  $t_1$  in r,  $t_2$  in s

if  $t_1$  and  $t_2$  join then output  $(t_1, t_2)$

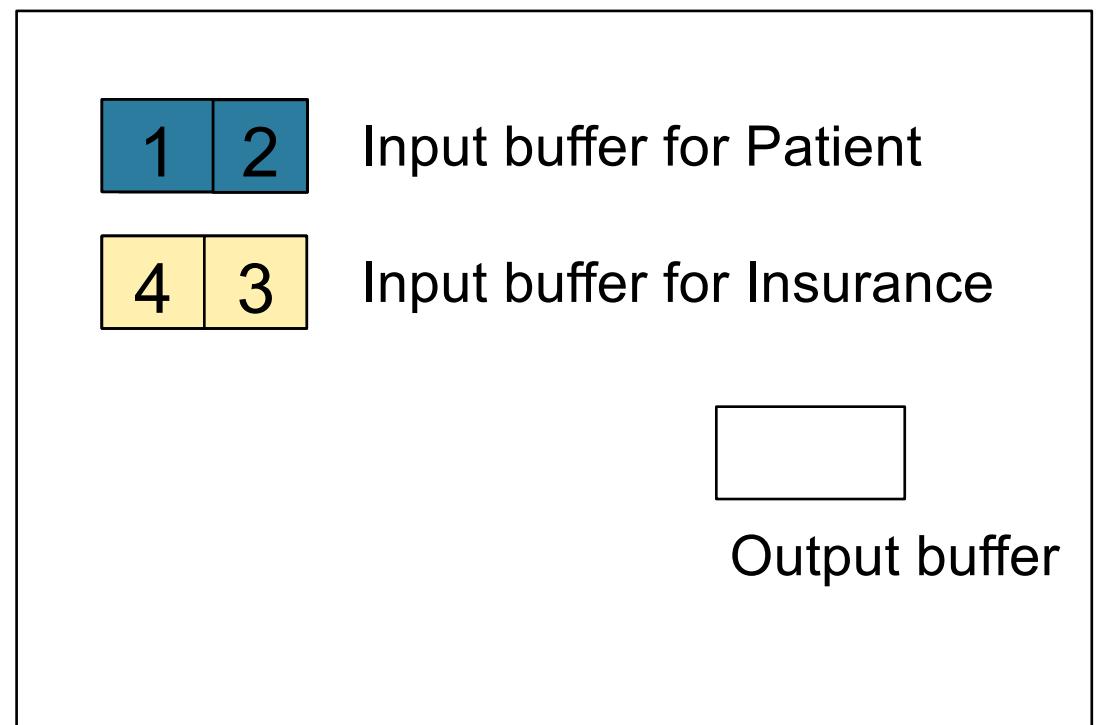
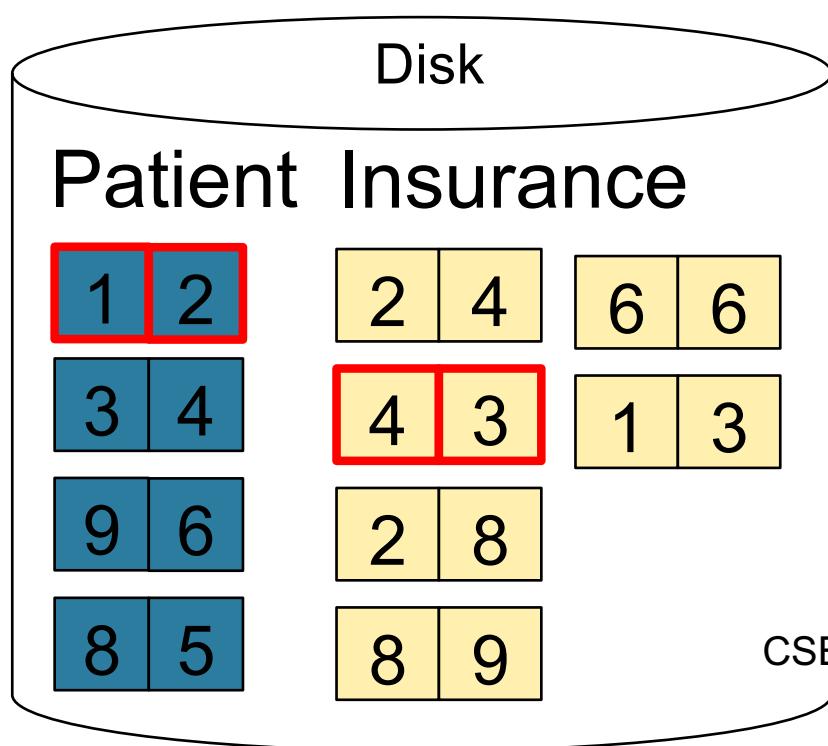
- Cost:  $B(R) + B(R)B(S)$

What is the Cost?

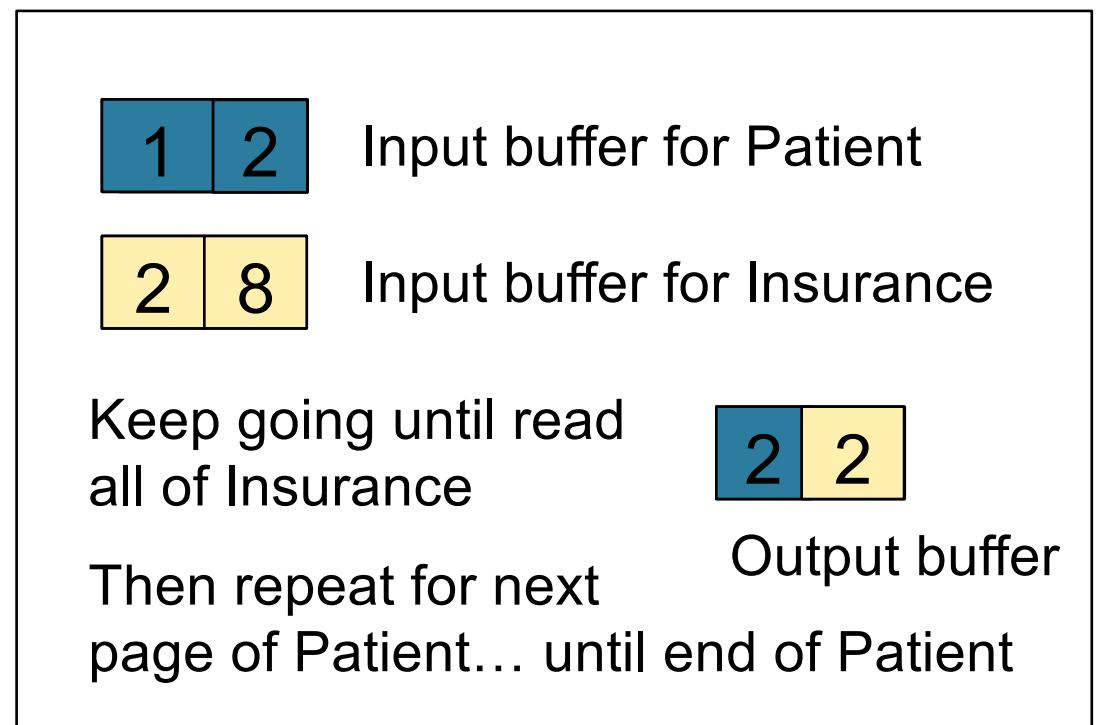
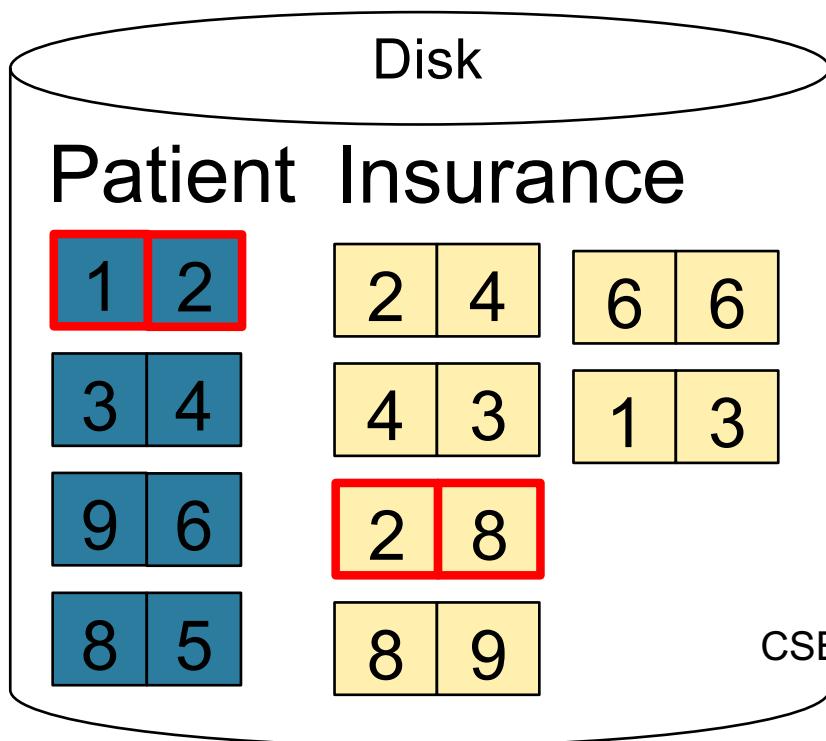
# Page-at-a-time Refinement



# Page-at-a-time Refinement



# Page-at-a-time Refinement



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:  $R \bowtie S$

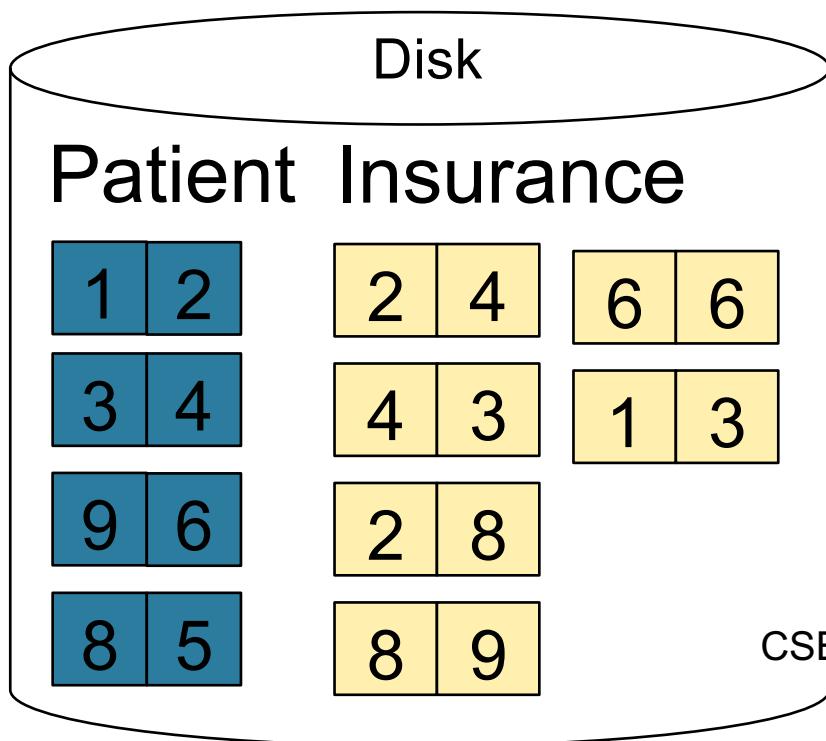
- Scan R and sort in main memory
  - Scan S and sort in main memory
  - Merge R and S
- 
- Cost:  $B(R) + B(S)$
  - One pass algorithm when  $B(S) + B(R) \leq M$
  - Typically, this is NOT a one pass algorithm

# Sort-Merge Join Example

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

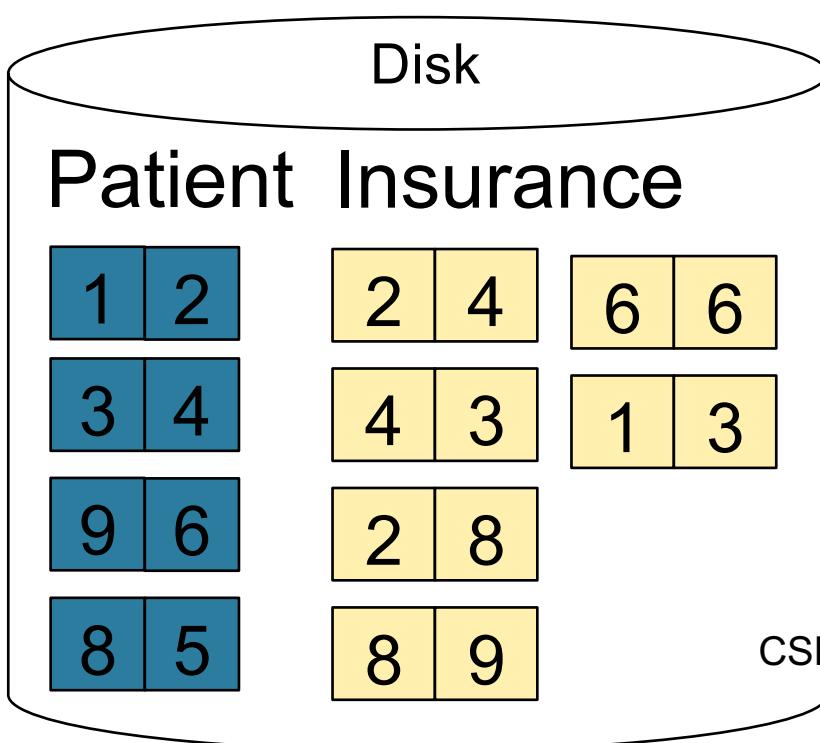
Memory M = 21 pages

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

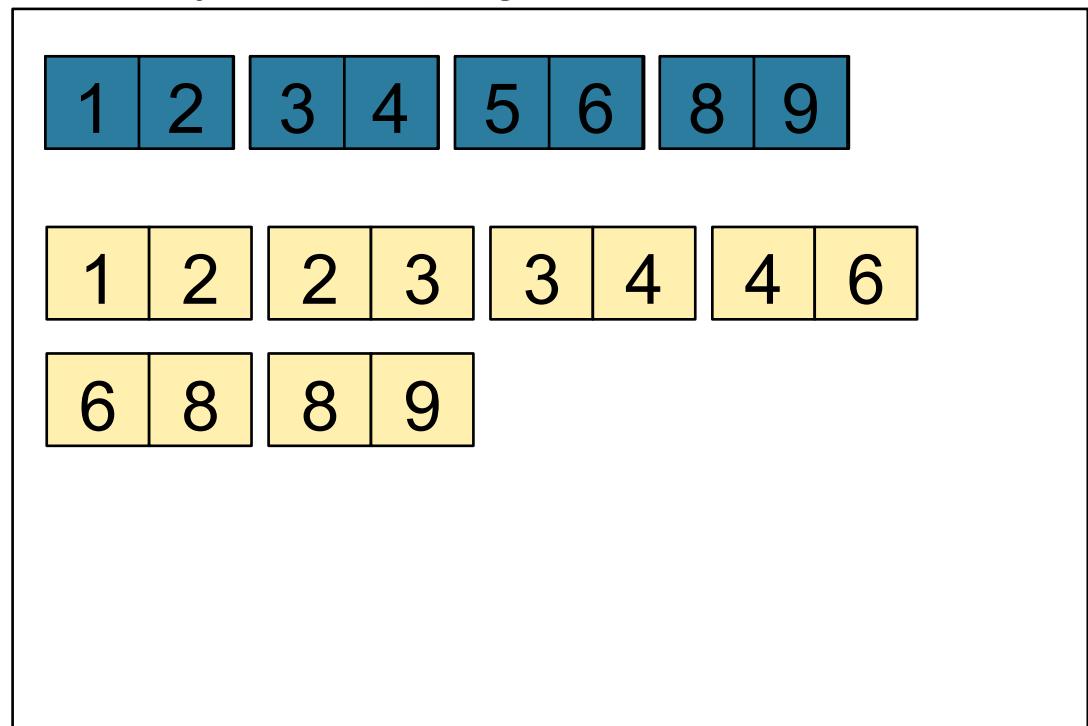


# Sort-Merge Join Example

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

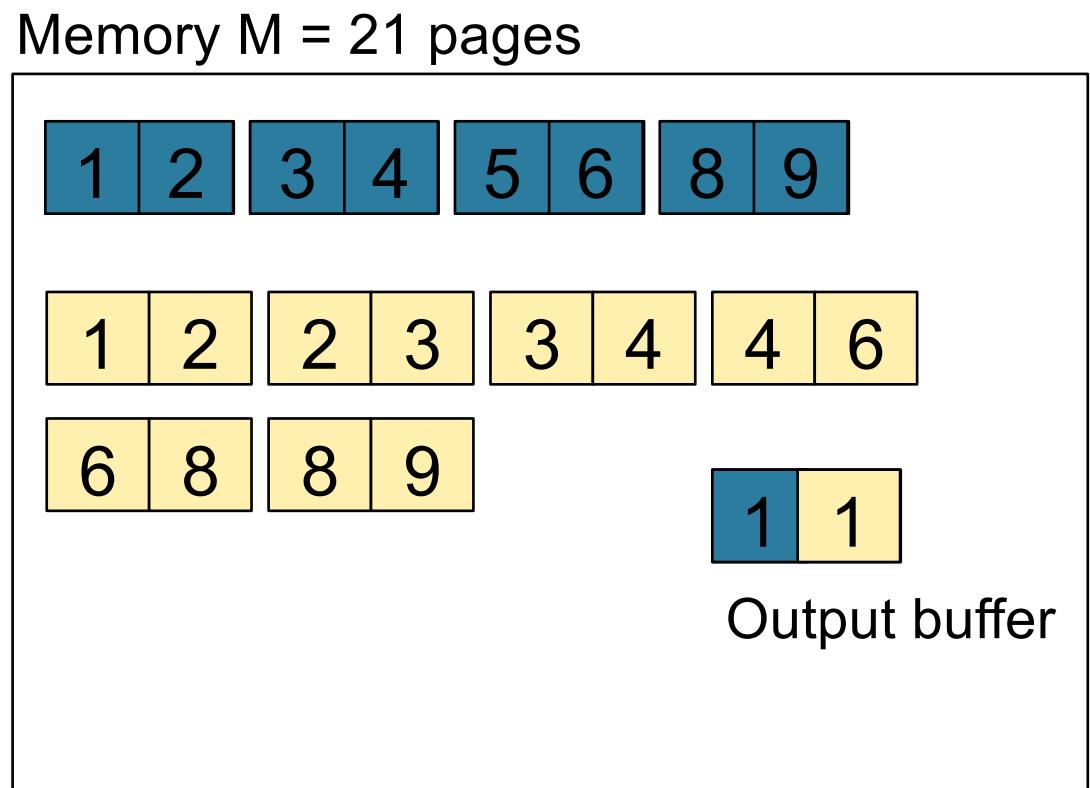
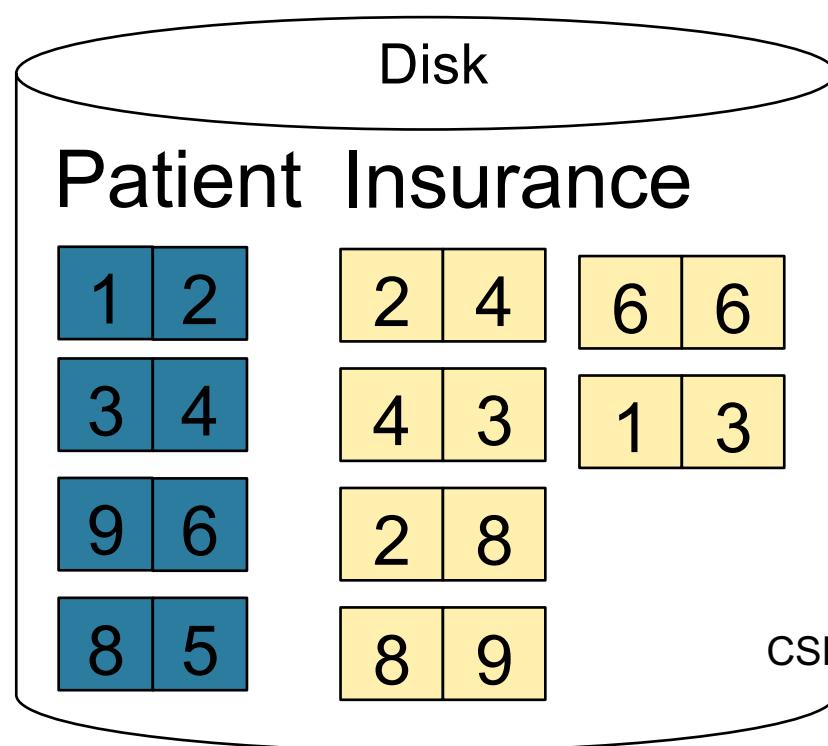


Memory M = 21 pages



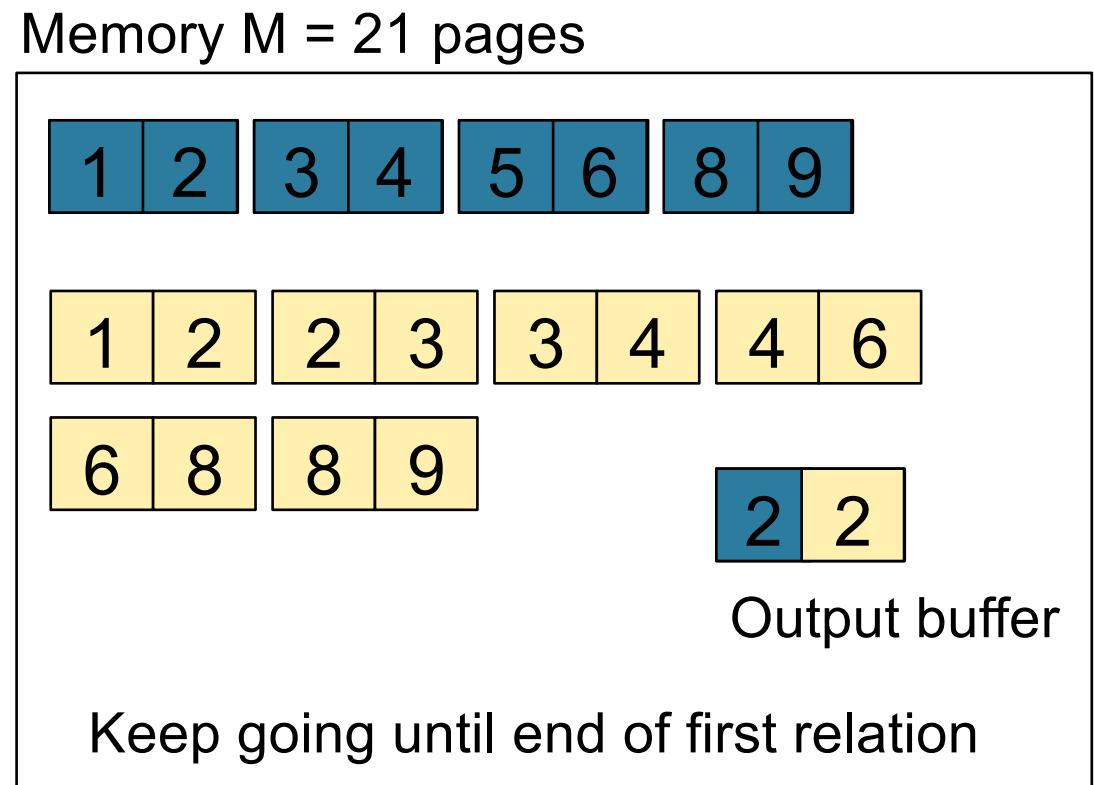
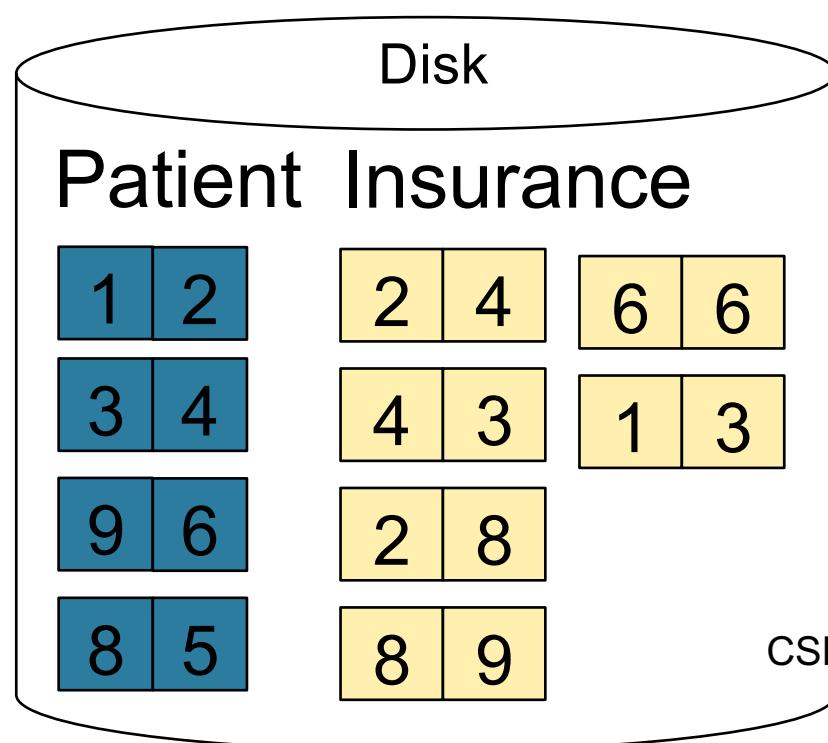
# Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance



# Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance



# 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:  
$$B(R) + T(R) * (T(S) * 1/V(S,a))$$

# Cost of Query Plans

`Supplier(sid, sname, scity, sstate)`

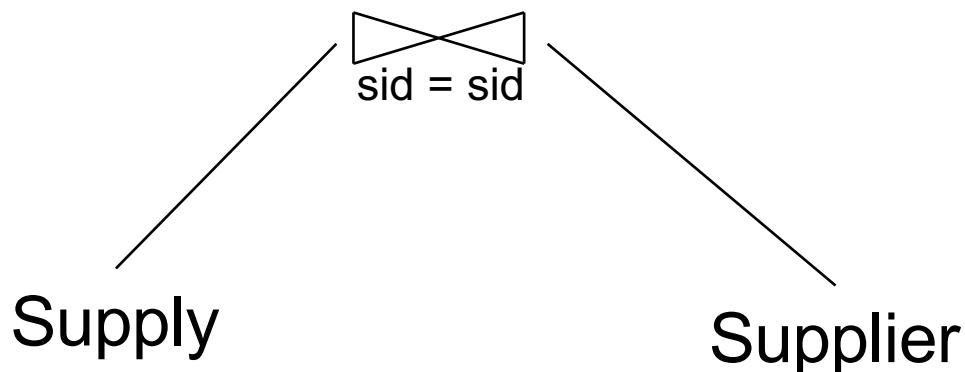
`Supply(sid, pno, quantity)`

# Logical Query Plan 1

$\Pi_{sname}$

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

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



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

CSE 414 Autumn 2018  
 $T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, scity) = 20$   
 $V(\text{Supplier}, state) = 10$

M=11 97

`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

# Logical Query Plan 1

$\Pi_{sname}$

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

$T = 10000$

$\bowtie$   
 $sid = sid$

Supply

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

Supplier

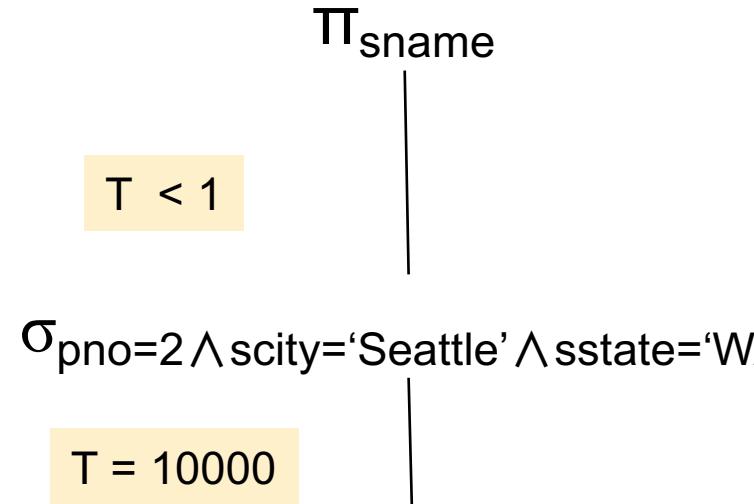
$T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, scity) = 20$   
 $V(\text{Supplier}, state) = 10$

```
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)`

# 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'
```

Supply

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

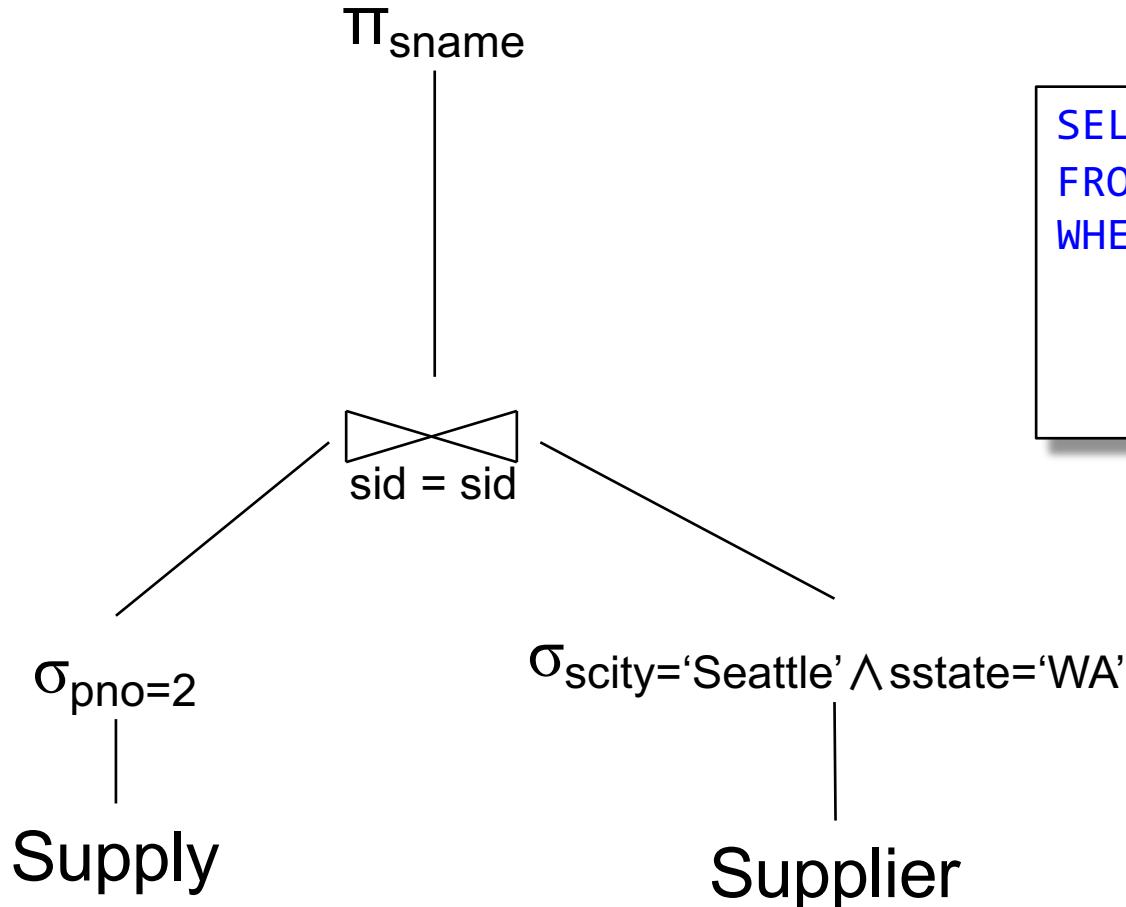
Supplier

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

`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'
```

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

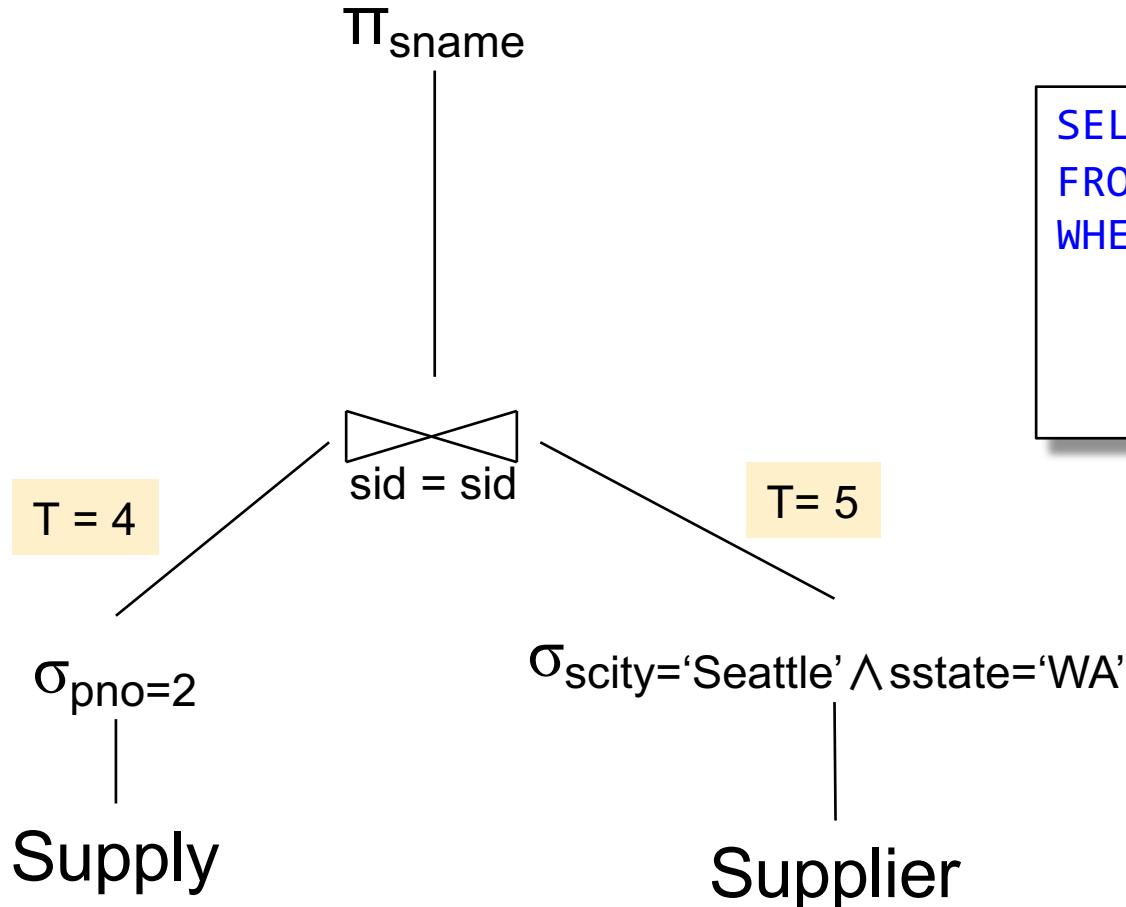
CSE 414 Autumn 2018  
 $T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, scity) = 20$   
 $V(\text{Supplier}, state) = 10$

$M=11_{100}$

`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

# Logical Query Plan 2



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

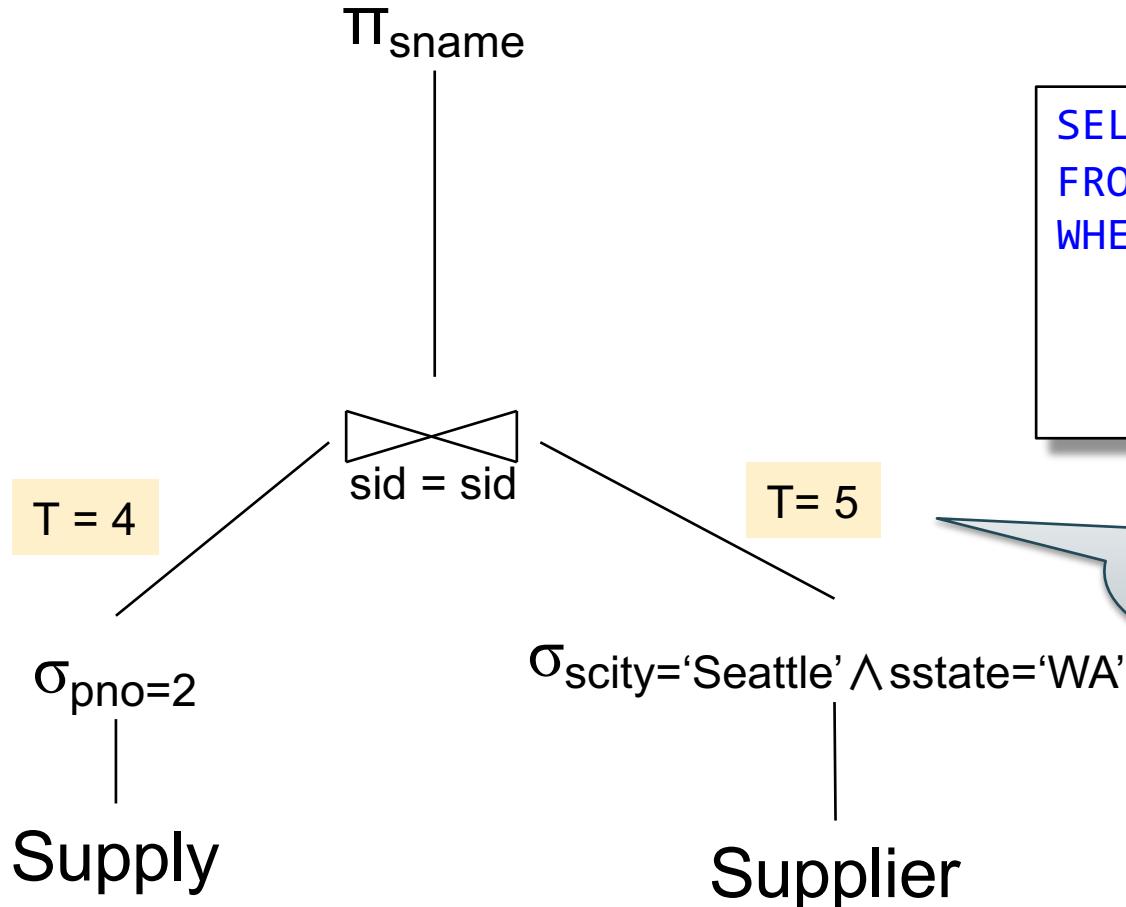
CSE 414 Autumn 2018  
 $T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, \text{scity}) = 20$   
 $V(\text{Supplier}, \text{sstate}) = 10$

$M = 11_{101}$

`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?

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

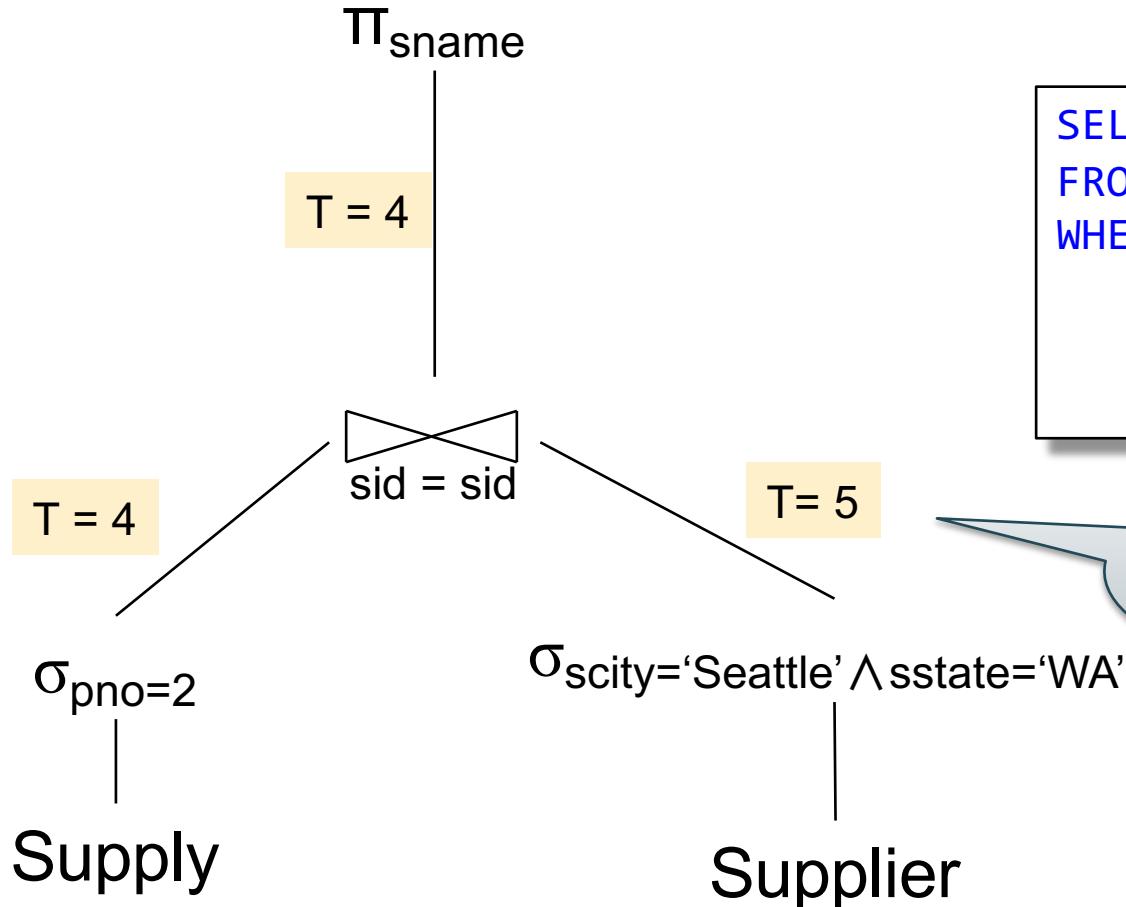
CSE 414 Autumn 2018  
 $T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, scity) = 20$   
 $V(\text{Supplier}, state) = 10$

$M=11_{102}$

`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?

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

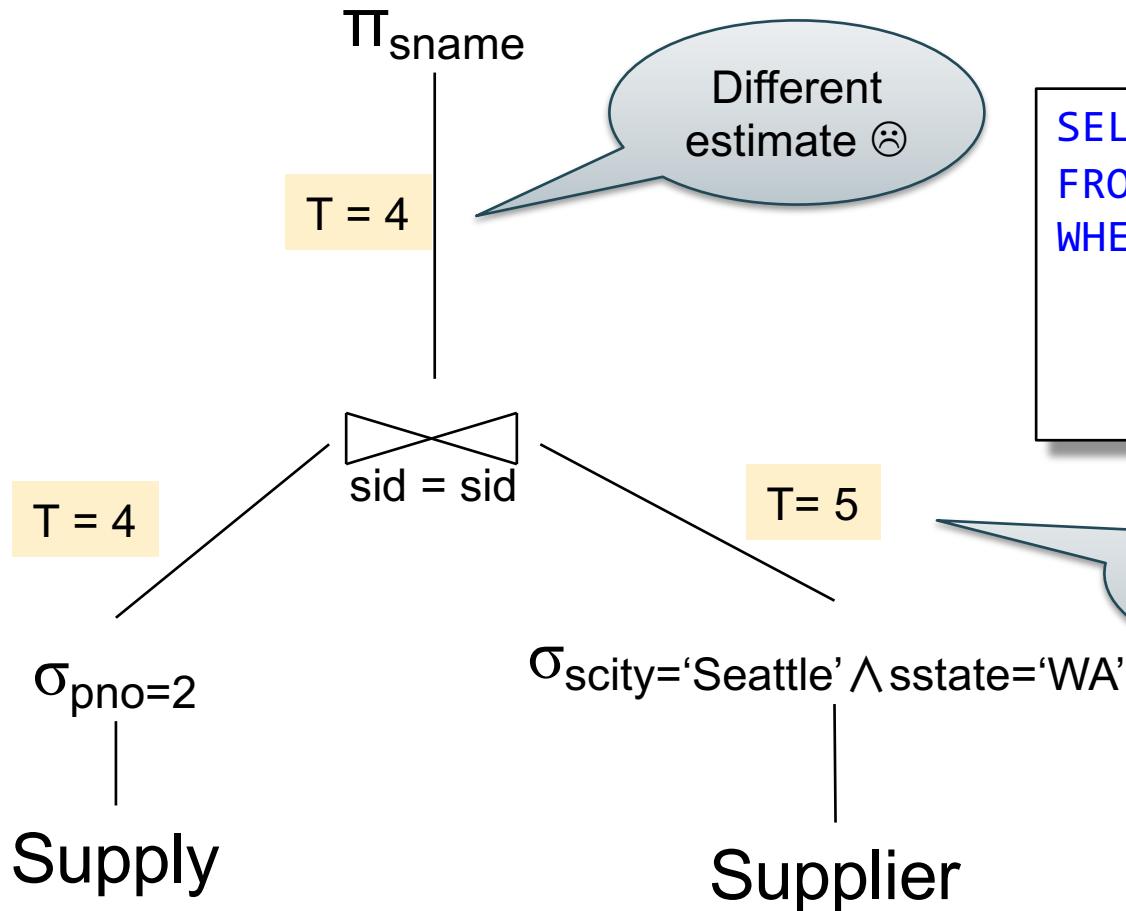
CSE 414 Autumn 2018  
 $T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, \text{scity}) = 20$   
 $V(\text{Supplier}, \text{sstate}) = 10$

$M=11_{103}$

`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

# Logical Query Plan 2



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

CSE 414 Autumn 2018  
 $T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, \text{scity}) = 20$   
 $V(\text{Supplier}, \text{sstate}) = 10$

$M=11_{104}$

Different estimate 😞

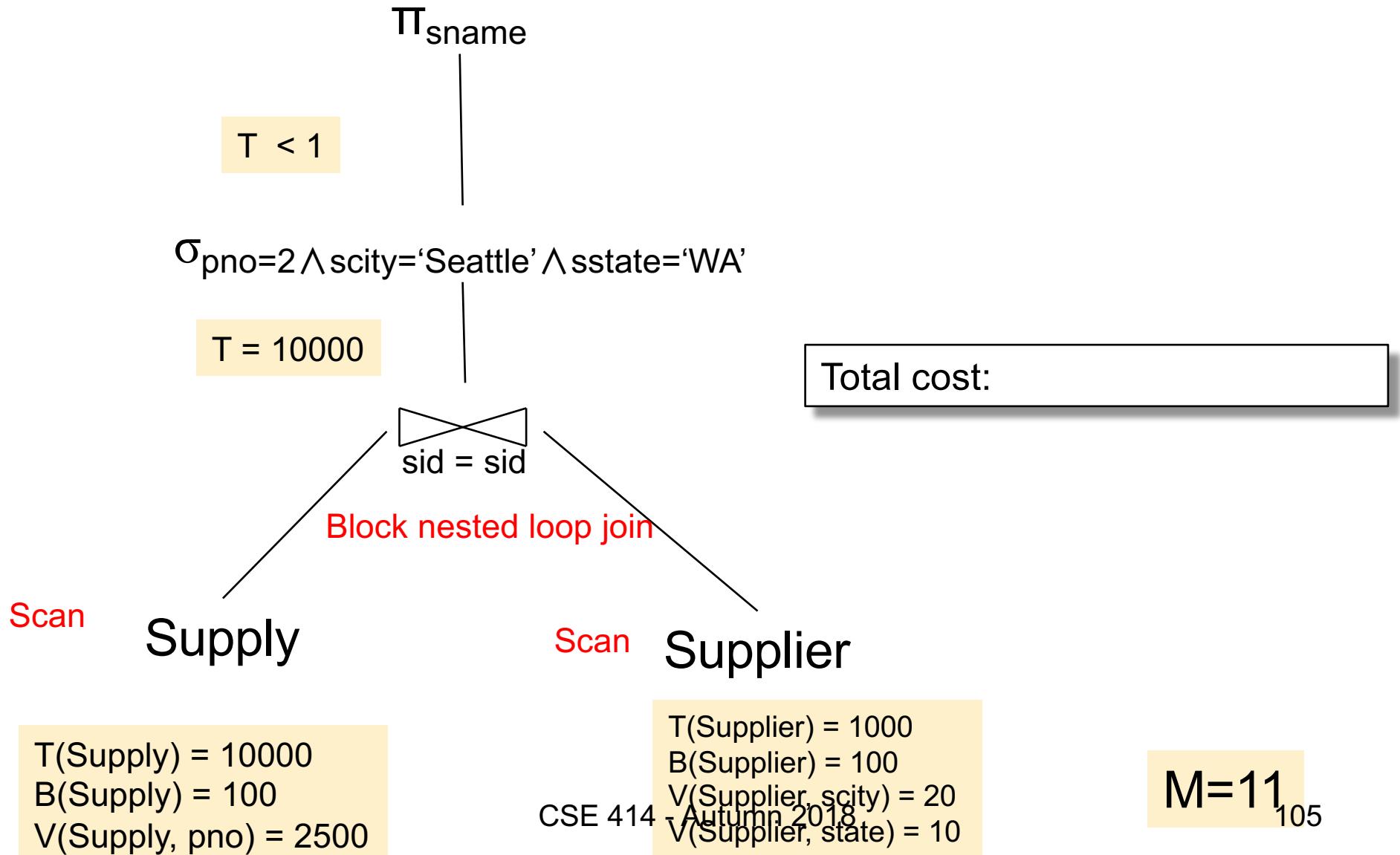
```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
and y.pno = 2
and x.scity = 'Seattle'
and x.sstate = 'WA'
```

Very wrong!  
Why?

`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

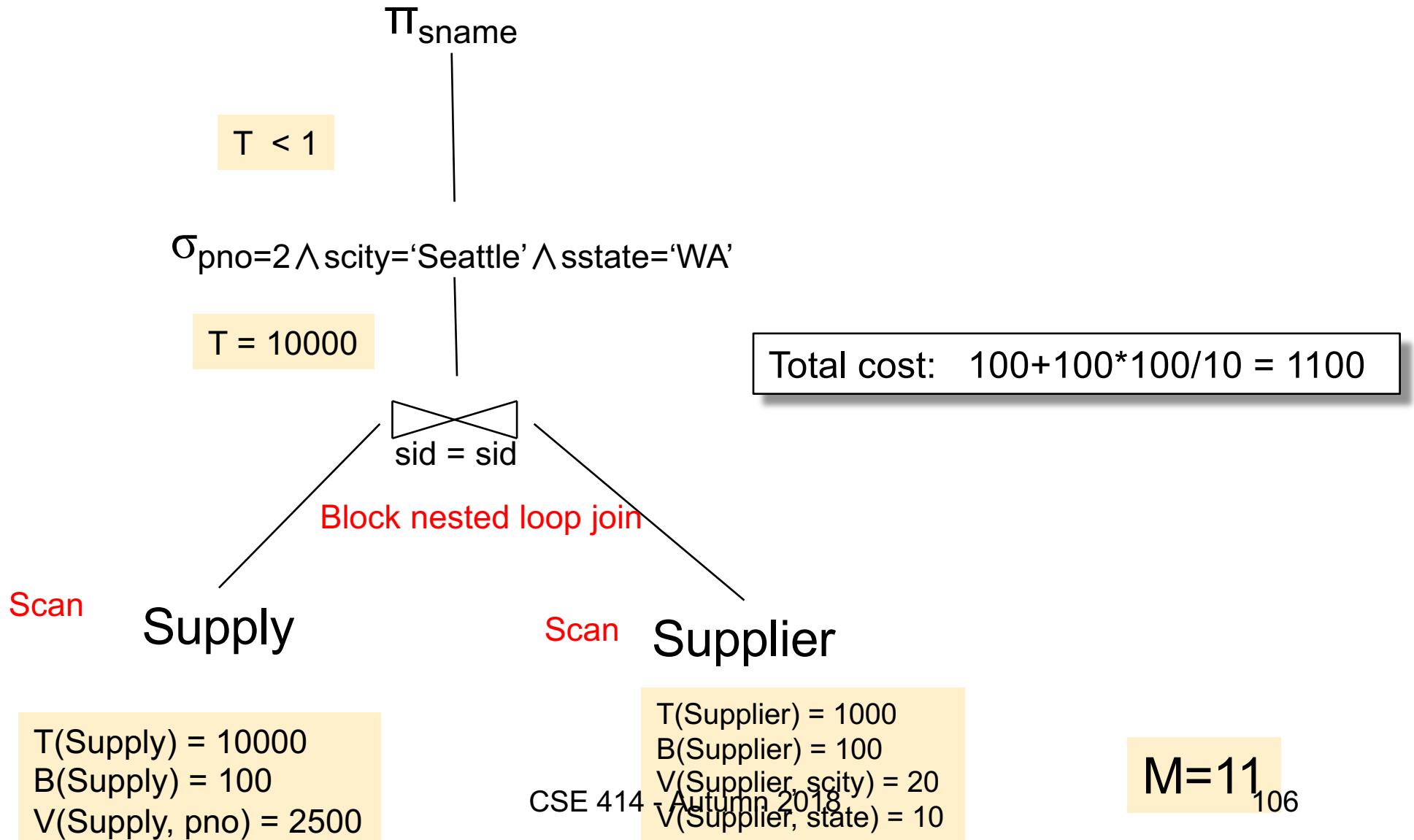
# Physical Plan 1



`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

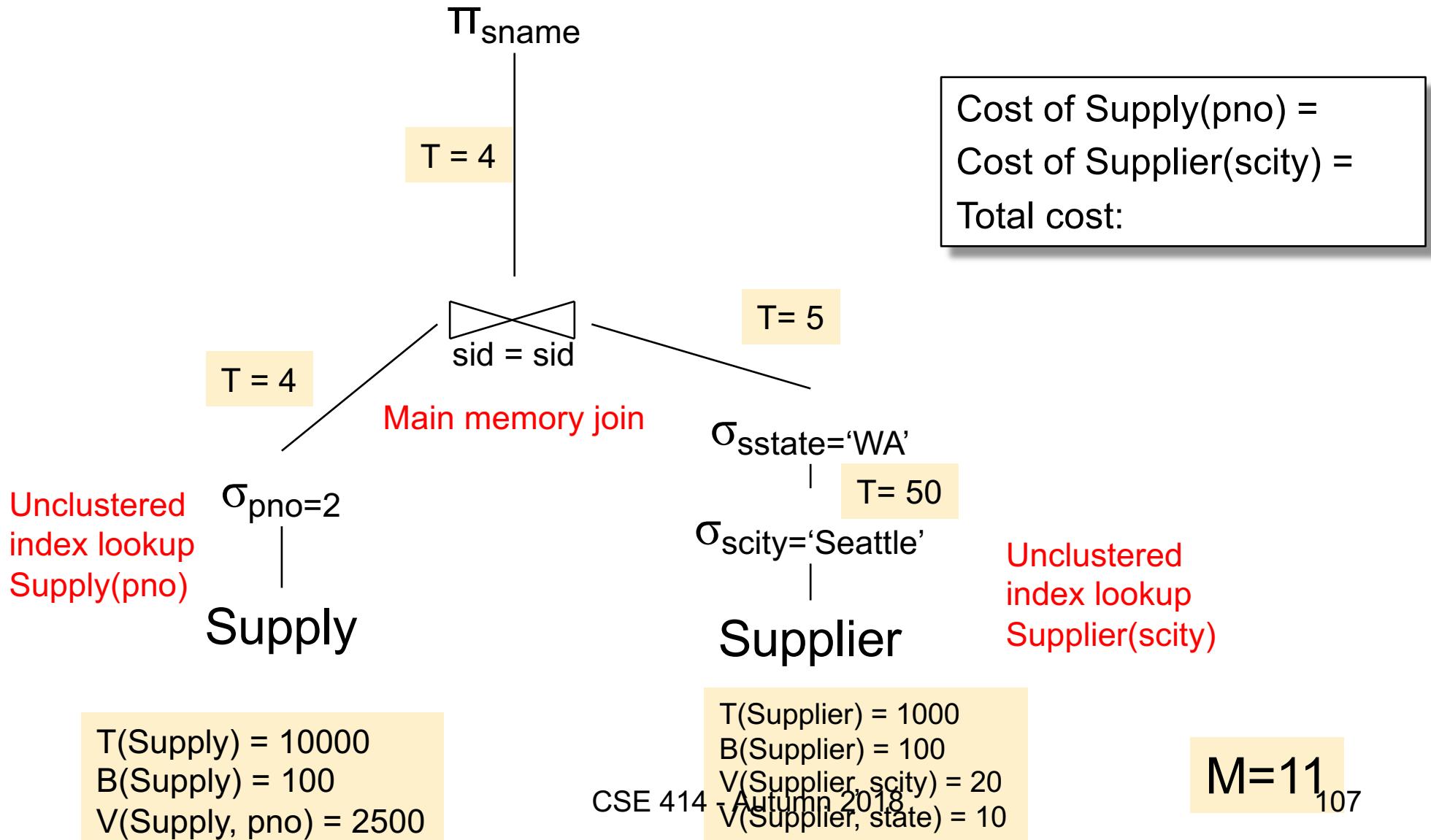
# Physical Plan 1



`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

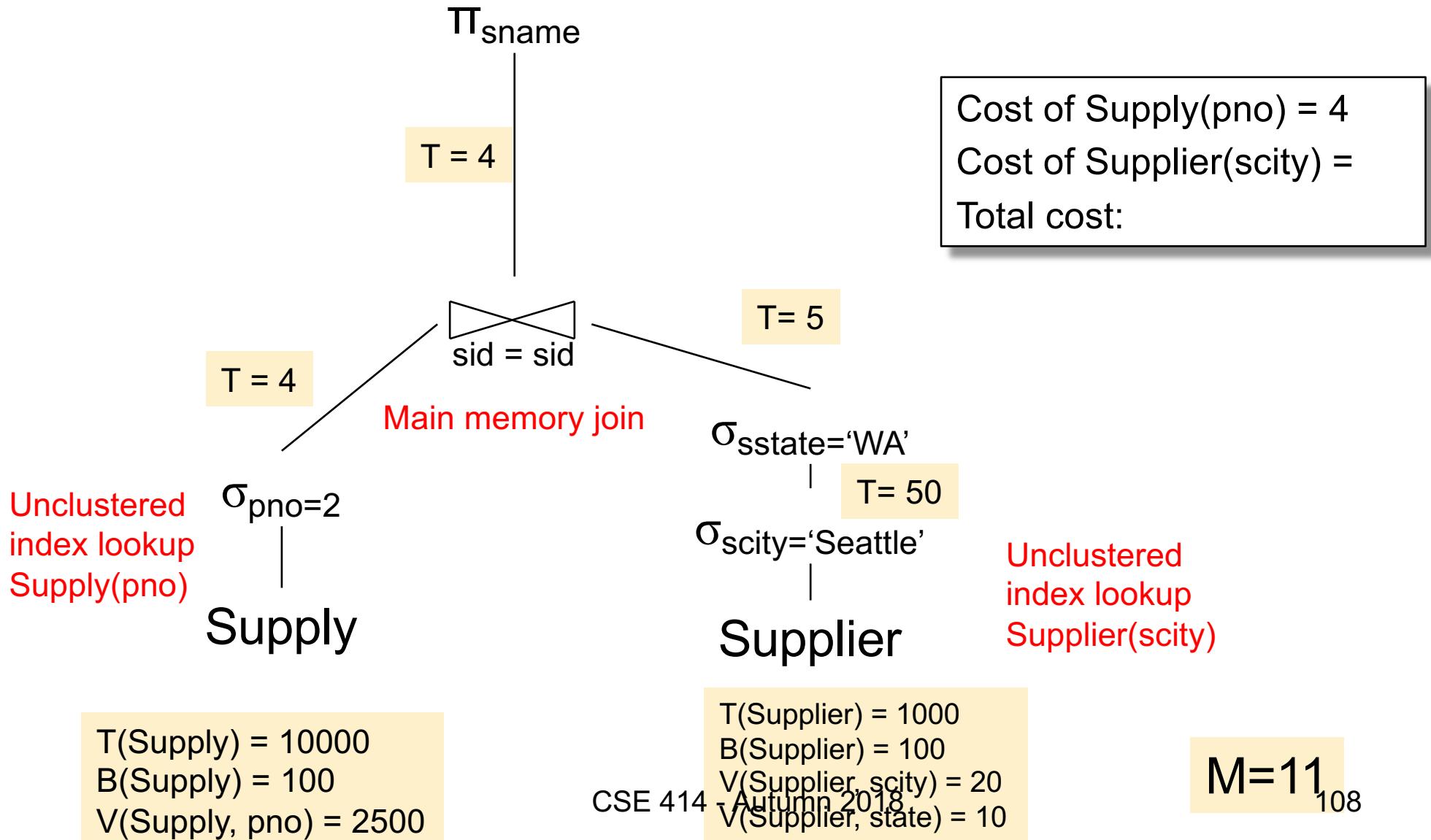
# Physical Plan 2



`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

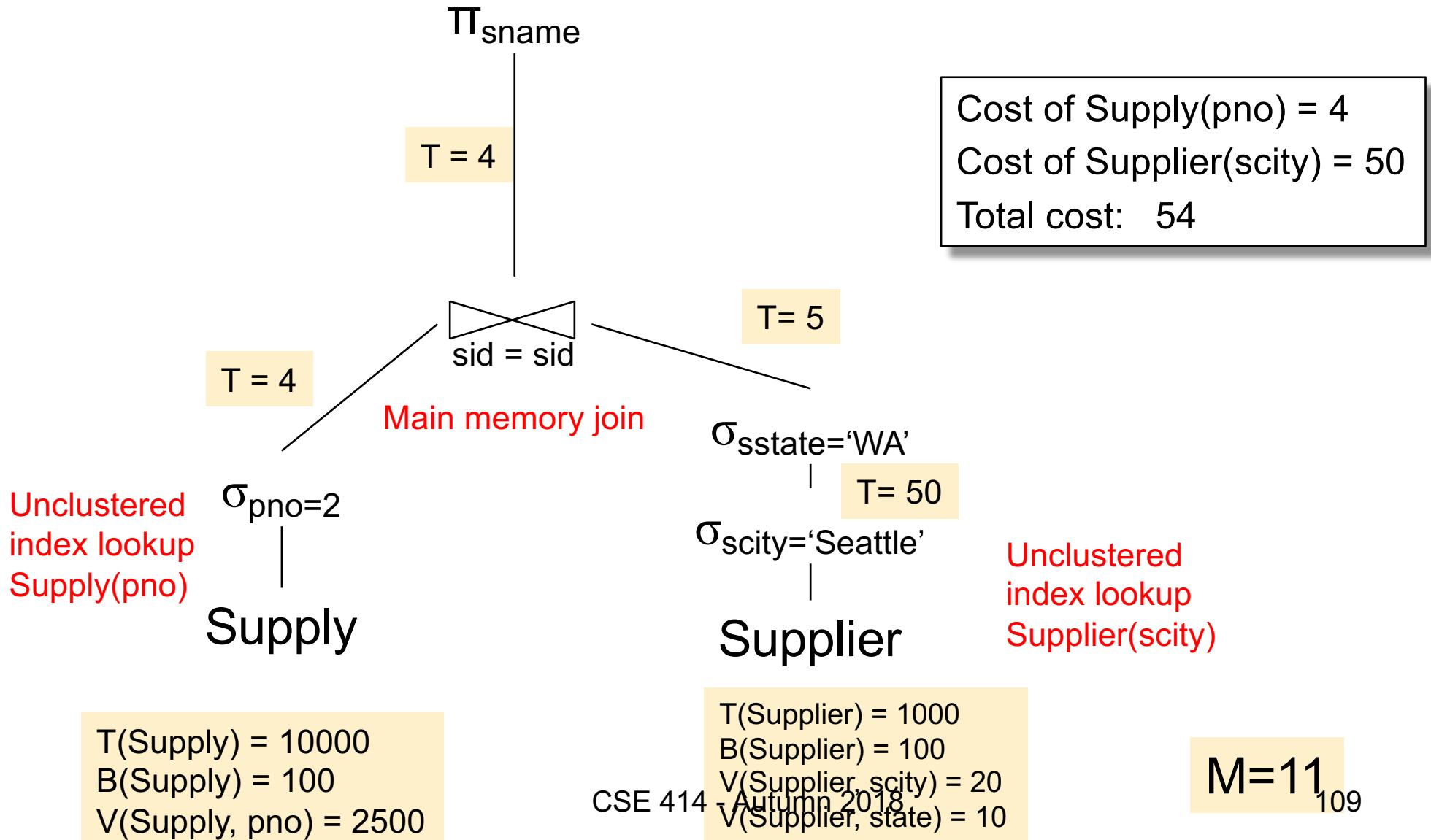
# Physical Plan 2



`Supplier(sid, sname, scity, sstate)`

`Supply(sid, pno, quantity)`

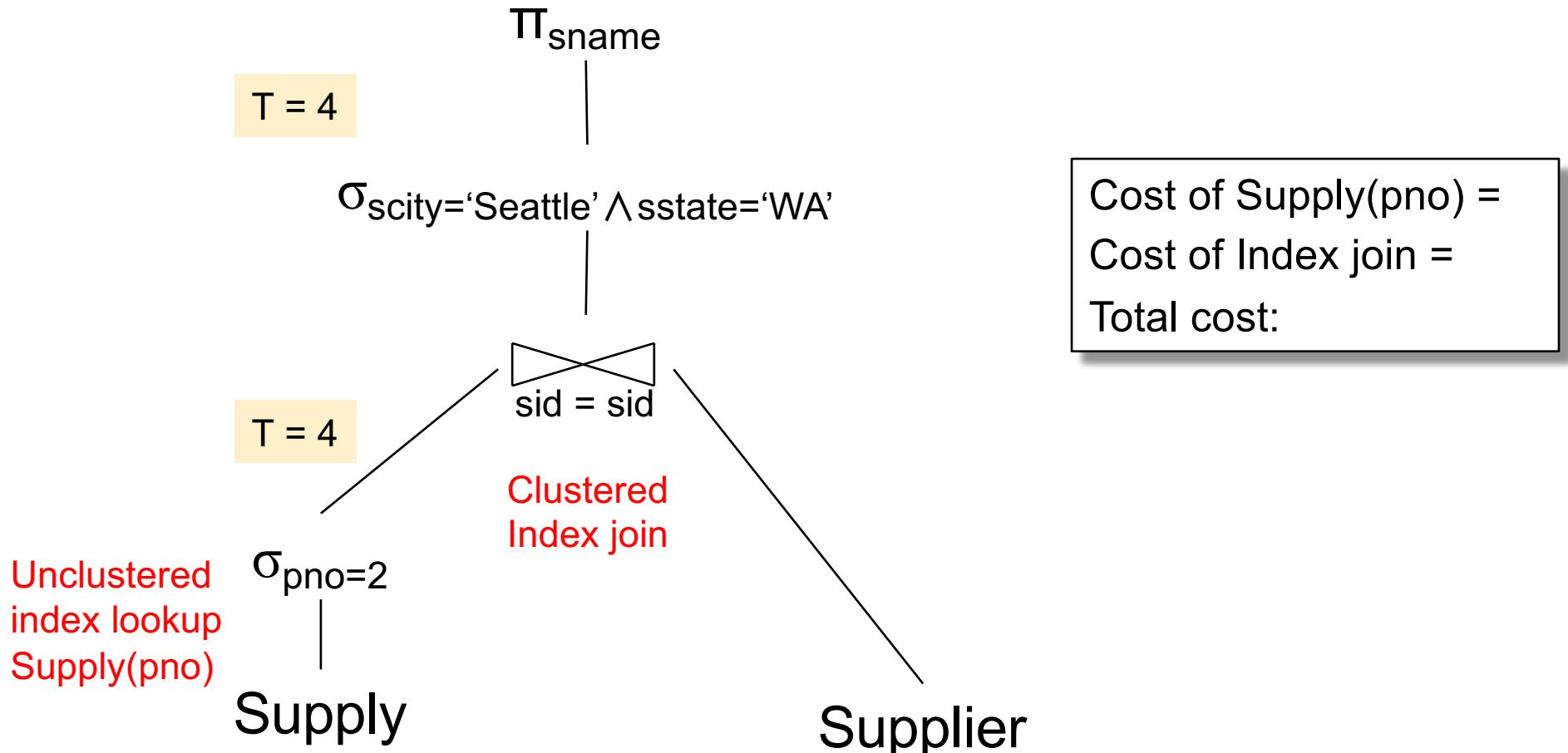
# Physical Plan 2



Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

# Physical Plan 3



$$T(\text{Supplier}) = 1000$$

$$B(\text{Supplier}) = 100$$

$$V(\text{Supplier}, \text{scity}) = 20$$

$$V(\text{Supplier}, \text{sstate}) = 10$$

CSE 414

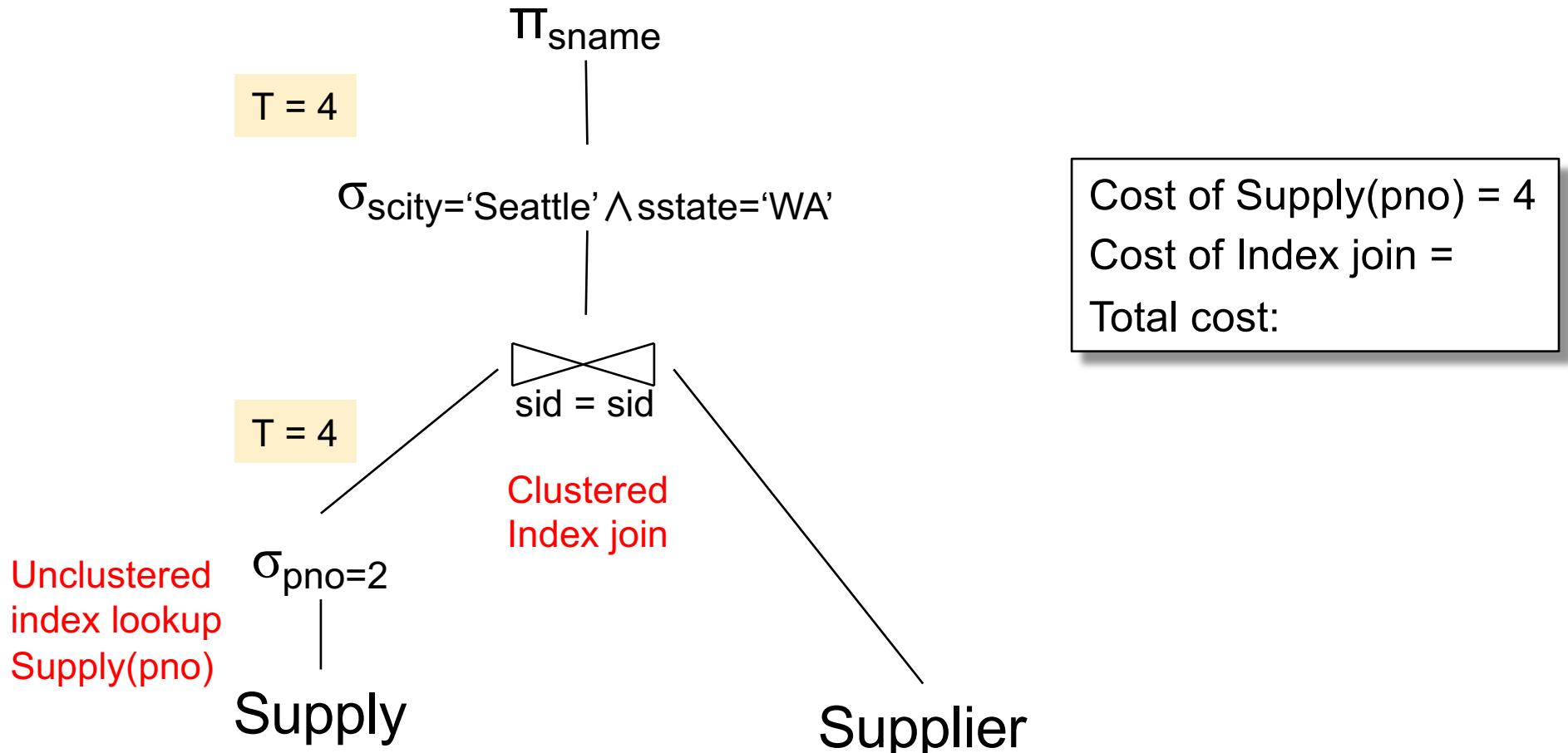
Autumn 2018

M=11<sub>110</sub>

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

# Physical Plan 3



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

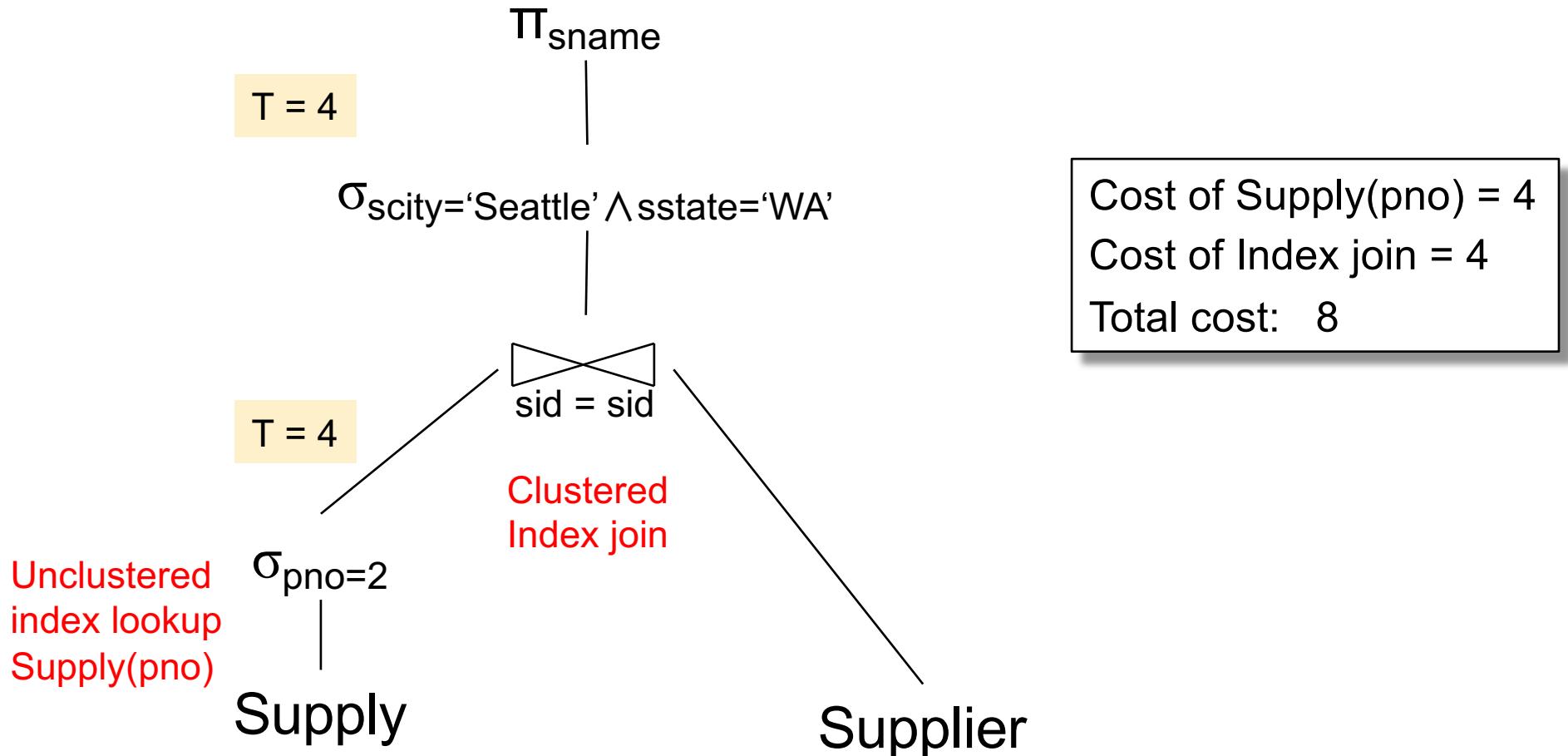
CSE 414 Autumn 2018  
 $T(\text{Supplier}) = 1000$   
 $B(\text{Supplier}) = 100$   
 $V(\text{Supplier}, \text{scity}) = 20$   
 $V(\text{Supplier}, \text{sstate}) = 10$

$M=11_{111}$

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

# Physical Plan 3



CSE 414 Autumn 2018

M=11<sub>112</sub>

# 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