

Introduction to Data Management

CSE 344

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

(4 lectures)

Introduction to Data Management

CSE 344

Lecture 14: Introduction to Query Evaluation

Announcement

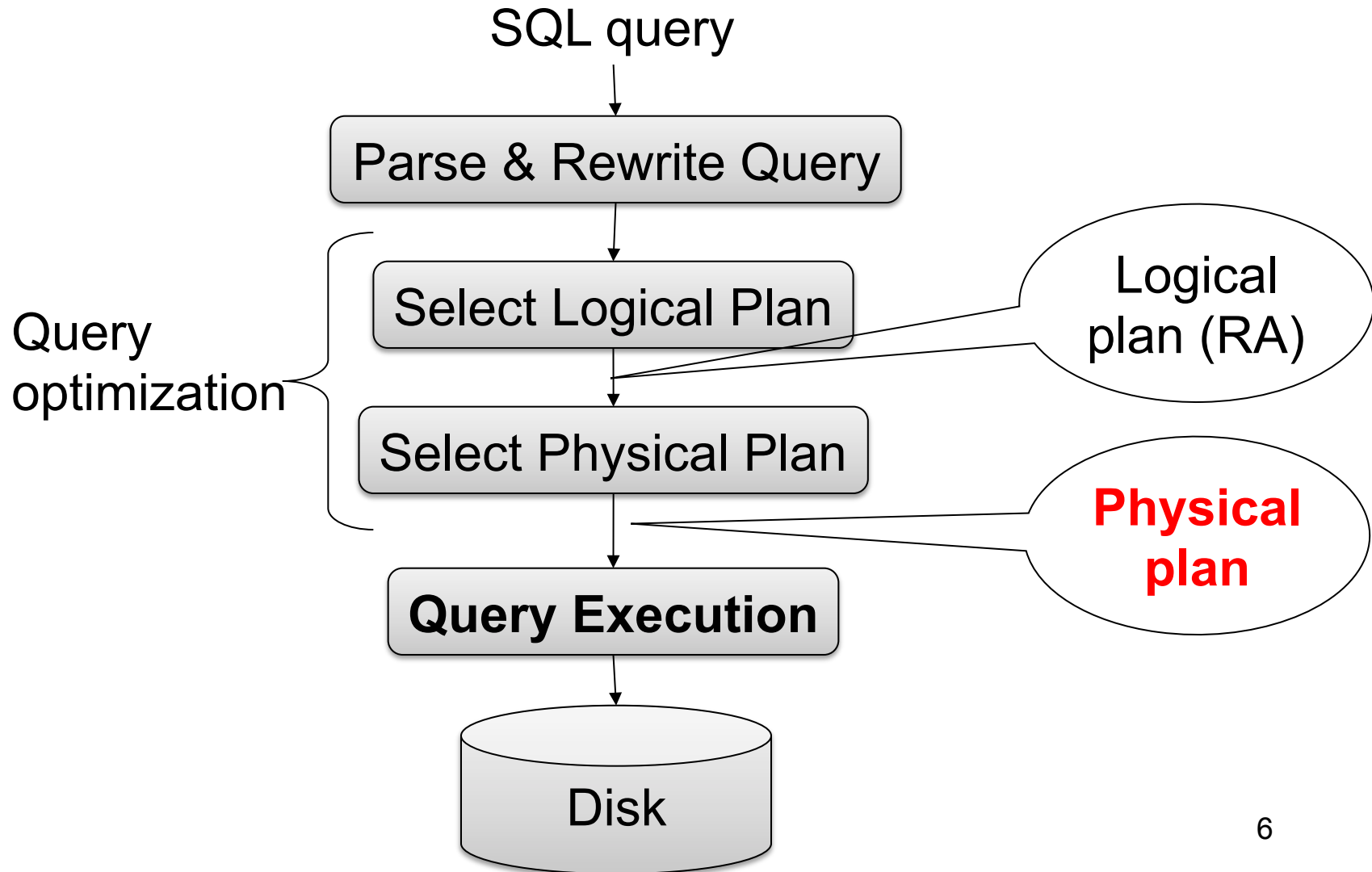
- Midterm review: 5:30pm in Smith 205
- WQ5 (datalog) due tonight
- HW4 (datalog) due on Tuesday
- Midterm: Wednesday, 1:30 in class

Class Overview

- Unit 1: Intro
- Unit 2: Relational Data Models and Query Languages
- Unit 3: Non-relational data
- Unit 4: RDMBS internals and query optimization
- Unit 5: Parallel query processing
- Unit 6: DBMS usability, conceptual design
- Unit 7: Transactions
- Unit 8: Advanced topics (time permitting)

From Logical RA Plans to Physical Plans

Query Evaluation Steps Review



Logical vs Physical Plans

- Logical plans:
 - Created by the parser from the input SQL text
 - Expressed as a relational algebra tree
 - Each SQL query has many possible logical plans
- Physical plans:
 - Goal is to choose an efficient implementation for each operator in the RA tree
 - Each logical plan has many possible physical plans

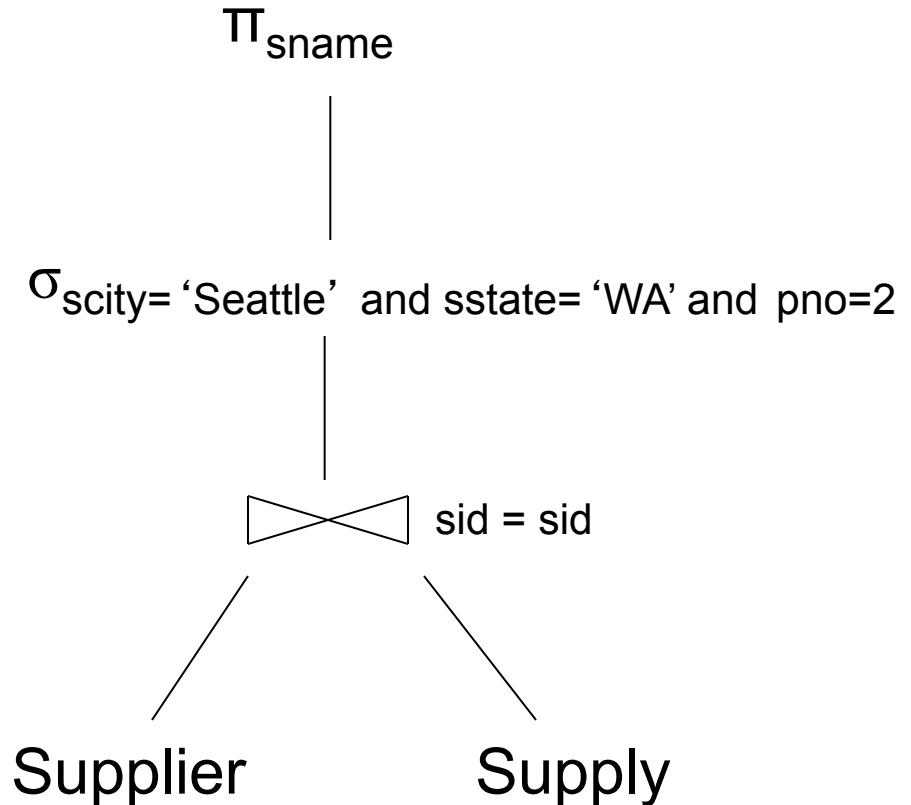
Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Review: 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”



Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Physical Query Plan 1

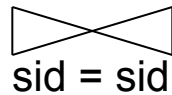
(On the fly)

Π_{sname}

(On the fly)

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

(Nested loop)



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

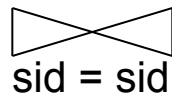
(On the fly)

Π_{sname}

(On the fly)

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

(Hash join)



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)

π_{sname} (d)

(Sort-merge join)

(c)
sid = sid

(Scan & write to T1)

(a) $\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA'}$

Supplier
(File scan)

(b) $\sigma_{\text{pno}=2}$ (Scan & write to T2)

Supply
(File scan)

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'
```

Query Optimization Problem

- For each SQL query... many logical plans
- For each logical plan... many physical plans
- Next: we will discuss physical operators;
how exactly are query executed?

Query Execution

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

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Main Memory Algorithms

Logical operator:

Supplier $\bowtie_{\text{sid}=\text{sid}}$ Supply

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

- 1.
- 2.
- 3.

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Main Memory Algorithms

Logical operator:

Supplier $\bowtie_{\text{sid}=\text{sid}}$ Supply

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 $O(??)$

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Main Memory Algorithms

Logical operator:

Supplier $\bowtie_{\text{sid}=\text{sid}}$ Supply

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

BRIEF Review of Hash Tables

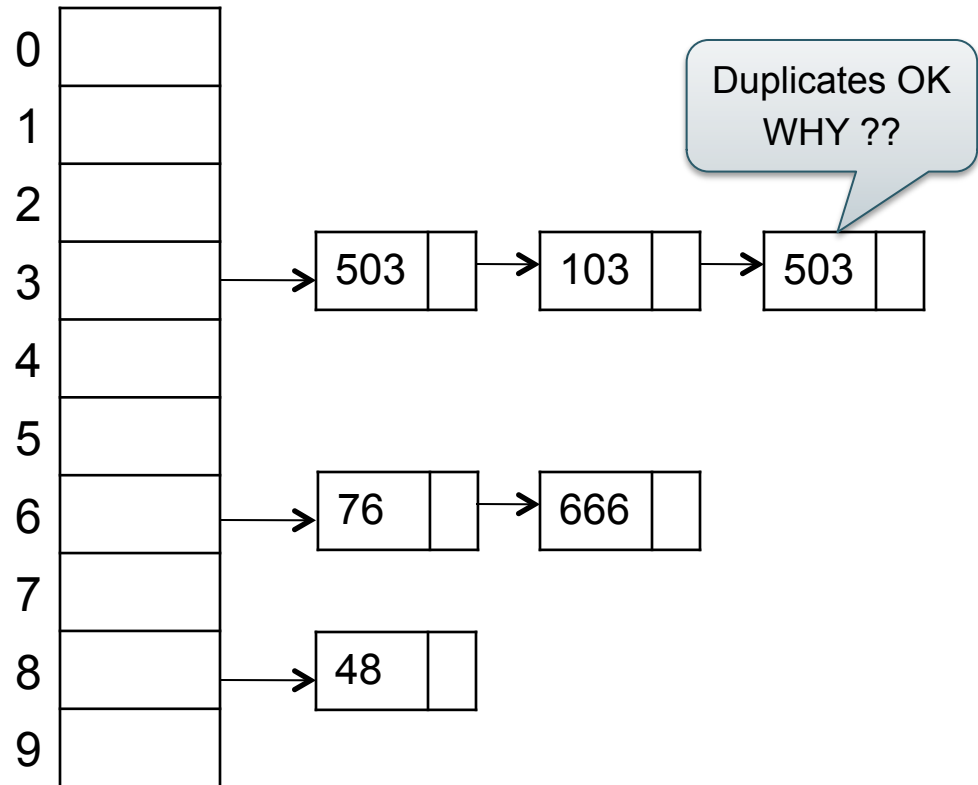
Separate chaining:

A (naïve) hash function:

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

Operations:

$$\begin{aligned} \text{find}(103) &= ?? \\ \text{insert}(488) &= ?? \end{aligned}$$



BRIEF Review of Hash Tables

- $\text{insert}(k, v)$ = inserts a key k with value v
- Many values for one key
 - Hence, duplicate k 's are OK
- $\text{find}(k)$ = returns the *list* of all values v associated to the key k

Implementing Query Operators with the Iterator Interface

Each operator implements three methods:

- `open()`
- `next()`
- `close()`

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {
```

```
}
```

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
}
```

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
}
```

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
    // cleans up (if any)  
    void close ();  
}
```

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
    // cleans up (if any)  
    void close ();  
}
```

```
class Select implements Operator {...  
    void open (Predicate p,  
               Operator child) {  
        this.p = p; this.child = child;  
    }  
}
```

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
    // cleans up (if any)  
    void close ();  
}
```

```
class Select implements Operator {...  
    void open (Predicate p,  
                Operator child) {  
        this.p = p; this.child = child;  
    }  
    Tuple next () {  
  
    }  
}
```

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
    // cleans up (if any)  
    void close ();  
}
```

```
class Select implements Operator {...  
    void open (Predicate p,  
               Operator child) {  
        this.p = p; this.child = child;  
    }  
    Tuple next () {  
        boolean found = false;  
        Tuple r = null;  
        while (!found) {  
            r = child.next();  
            if (r == null) break;  
            found = p(in);  
        }  
    }  
}
```

Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
    // cleans up (if any)  
    void close ();  
}
```

```
class Select implements Operator {...  
    void open (Predicate p,  
               Operator child) {  
        this.p = p; this.child = child;  
    }  
    Tuple next () {  
        boolean found = false;  
        Tuple r = null;  
        while (!found) {  
            r = child.next();  
            if (r == null) break;  
            found = p(in);  
        }  
        return r;  
    }  
}
```


Implementing Query Operators with the Iterator Interface

Example “on the fly” selection operator

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
    // cleans up (if any)  
    void close ();  
}
```

```
class Select implements Operator {...  
    void open (Predicate p,  
               Operator child) {  
        this.p = p; this.child = child;  
    }  
    Tuple next () {  
        boolean found = false;  
        Tuple r = null;  
        while (!found) {  
            r = child.next();  
            if (r == null) break;  
            found = p(in);  
        }  
        return r;  
    }  
    void close () { child.close(); }  
}
```

Implementing Query Operators with the Iterator Interface

```
interface Operator {  
  
    // initializes operator state  
    // and sets parameters  
    void open (...);  
  
    // calls next() on its inputs  
    // processes an input tuple  
    // produces output tuple(s)  
    // returns null when done  
    Tuple next ();  
  
    // cleans up (if any)  
    void close ();  
}
```

Query plan execution

```
Operator q = parse("SELECT ...");  
q = optimize(q);  
  
q.open();  
while (true) {  
    Tuple t = q.next();  
    if (t == null) break;  
    else printOnScreen(t);  
}  
q.close();
```

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

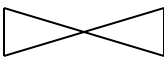
(On the fly)

Π_{sname}

(On the fly)

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

(Nested loop)


sno = sno

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

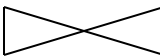
(On the fly)

Π_{sname} **open()**

(On the fly)

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

(Nested loop)


sno = sno

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

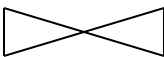
(On the fly)

Π_{sname} **open()**

(On the fly)

$\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA' \text{ and } \text{pno}=2}$ **open()**

(Nested loop)


sno = sno

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close for nested loop join

(On the fly)

Π_{sname} open()

(On the fly)

$\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA' \text{ and } \text{pno}=2}$ open()

(Nested loop)

sno = sno open()

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

(On the fly)

Π_{sname} **open()**

(On the fly)

$\sigma_{scity='Seattle' \text{ and } sstate='WA' \text{ and } pno=2}$ **open()**

(Nested loop)

open()
sno = sno

open()
Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

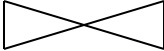
(On the fly)

Π_{sname} **open()**

(On the fly)

$\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA' \text{ and } \text{pno}=2}$ **open()**

(Nested loop)

open()

sno = sno

open()
Suppliers
(File scan)

open()
Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

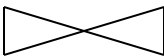
(On the fly)

Π_{sname} **next()**

(On the fly)

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

(Nested loop)


sno = sno

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

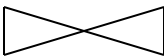
(On the fly)

Π_{sname} **next()**

(On the fly)

$\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA' \text{ and } \text{pno}=2}$ **next()**

(Nested loop)


sno = sno

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

(On the fly)

Π_{sname} **next()**

(On the fly)

$\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA' \text{ and } \text{pno}=2}$ **next()**

(Nested loop)

next()
sno = sno

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

(On the fly)

Π_{sname} **next()**

(On the fly)

$\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA' \text{ and } \text{pno}=2}$ **next()**

(Nested loop)

next()
sno = sno

next()
Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss: open/next/close
for nested loop join

(On the fly)

Π_{sname} **next()**

(On the fly)

$\sigma_{\text{scity}='Seattle' \text{ and } \text{sstate}='WA' \text{ and } \text{pno}=2}$ **next()**

(Nested loop)

next()
sno = sno

next()
Suppliers
(File scan)

next()
Supplies
(File scan)

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

Discuss: open/next/close
for nested loop join

(On the fly)

Π_{sname} **next()**

(On the fly)

$\sigma_{scity='Seattle' \text{ and } sstate='WA' \text{ and } pno=2}$ **next()**

(Nested loop)

next()
sno = sno

next()
Suppliers
(File scan)

next()
next()
Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss hash-join
in class

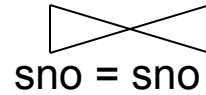
(On the fly)

Π_{sname}

(On the fly)

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

(Hash Join)



Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Pipelining

Discuss hash-join
in class

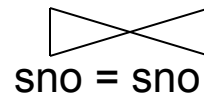
(On the fly)

Π_{sname}

(On the fly)

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

(Hash Join)



Suppliers
(File scan)

Supplies
(File scan)

Tuples from here are pipelined

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

Discuss hash-join
in class

(On the fly)

Π_{sname}

(On the fly)

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

(Hash Join)

sno = sno

Tuples from here are "blocked"

Tuples from here are pipelined

Suppliers
(File scan)



Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Blocked Execution

(On the fly)

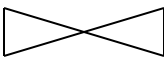
Π_{sname}

Discuss merge-join
in class

(On the fly)

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

(Merge Join)


sno = sno

Suppliers
(File scan)

Supplies
(File scan)

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Blocked Execution

(On the fly)

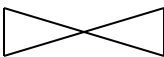
Π_{sname}

Discuss merge-join
in class


(On the fly)

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

(Merge Join)


sno = sno

Blocked


Suppliers
(File scan)


Supplies
(File scan)

Blocked

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

Query Execution Bottom Line

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

Recall: Physical Data Independence

- Applications are insulated from changes in physical storage details
- SQL and relational algebra facilitate physical data independence
 - Both languages input and output relations
 - Can choose different implementations for operators

Introduction to Database Systems

CSE 344

Lecture 15-16:
Basics of Data Storage and Indexes

Announcements

- HW4 (datalog) due tomorrow (Tuesday)
- Midterm: Wednesday
- No sections on Thursday!
- HW5 (SQL++) due next Tuesday

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

Data Storage

Student

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

- 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

10	Tom	Hanks	block 1
20	Amy	Hanks	
50	block 2
200	...		
220			block 3
240			
420			
800			

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

Data File Types

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

The data file can be one of:

- **Heap file**
 - Unsorted
- **Sequential file**
 - Sorted according to some attribute(s) called key

Student

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

Data File Types

The data file can be one of:

- **Heap file**
 - Unsorted
- **Sequential file**
 - Sorted according to some attribute(s) called key

Note: key here means something different from primary key: it just means that we order the file according to that attribute. In our example we ordered by **ID**. Might as well order by **fName**, if that seems a better idea for the applications running on our database.

Index

- An **additional** file, that allows fast access to records in the data file given a search 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

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

Key = means here search key

This



Is Not A Key

Different keys:

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



This is not a pipe.

CSE 344 - 2017au



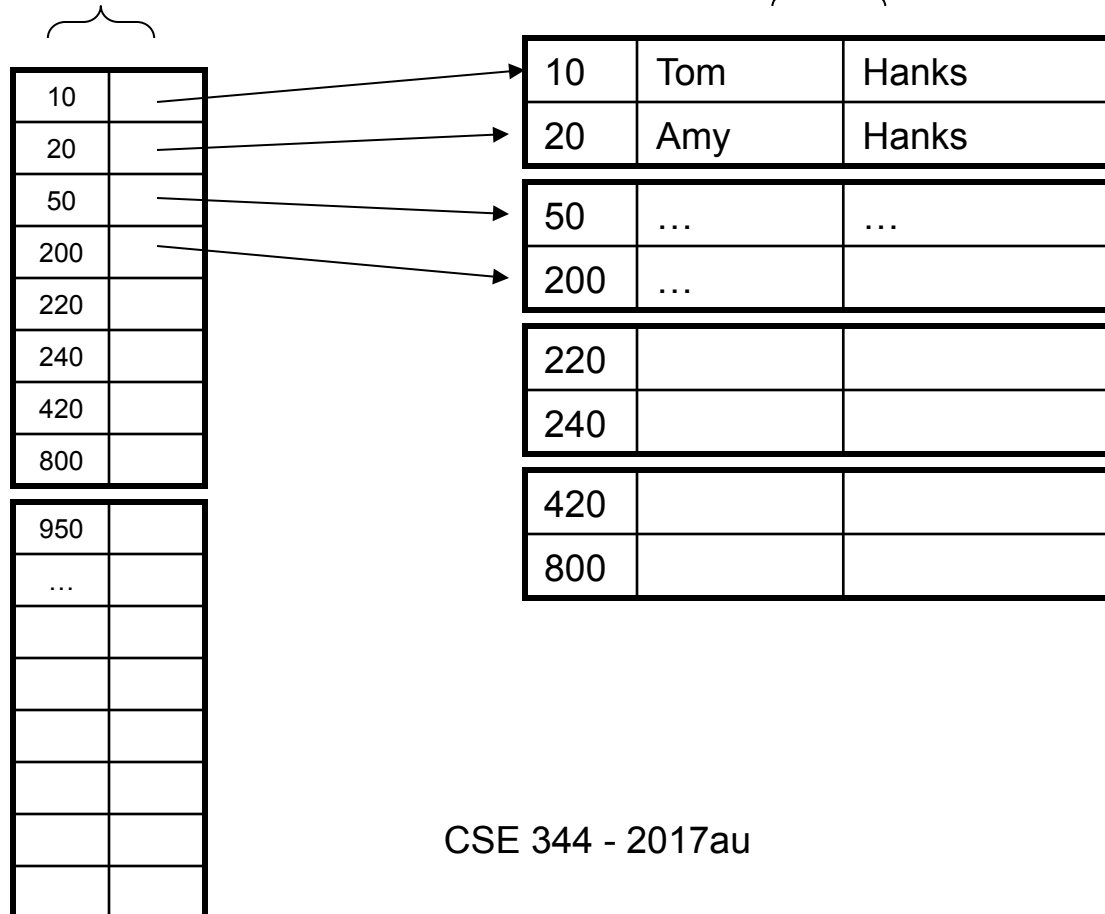
Student

Example 1: Index on ID

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

Index **Student_ID** on **Student.ID**

Data File **Student**

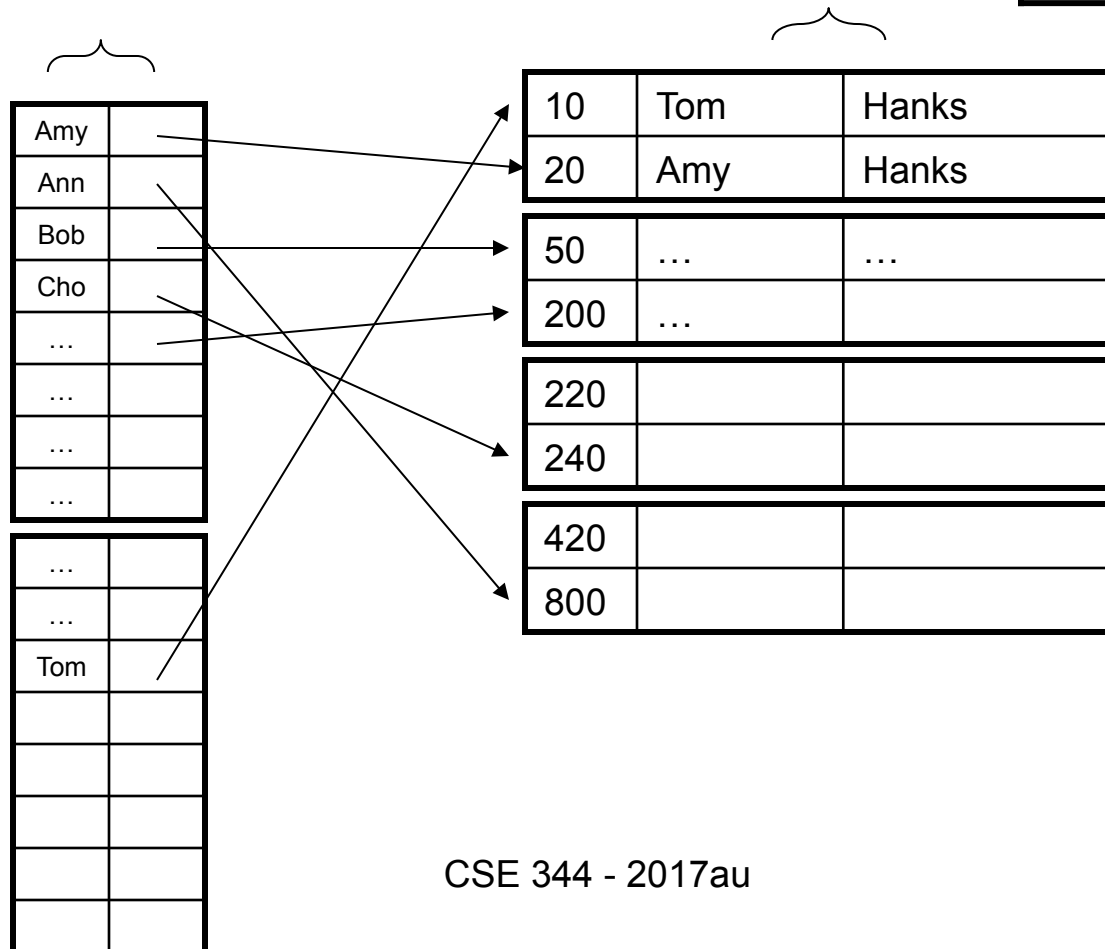


ID	fName	IName
10	Tom	Hanks
20	Amy	Hanks
...		

Example 2: Index on fName

Index **Student_fName**
on **Student.fName**

Data File **Student**



Index Organization

We need a way to represent indexes after loading into memory so that they can be used

Several ways to do this:

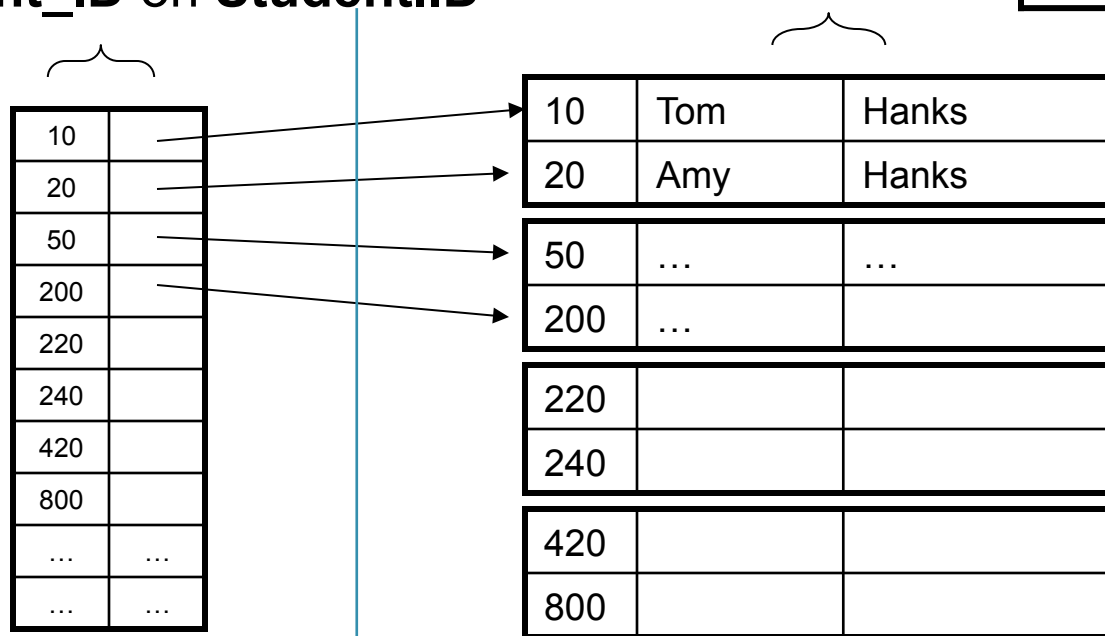
- Hash table
- B+ trees – most popular
 - They are search trees, but they are not binary instead have higher fanout
 - Will discuss them briefly next
- Specialized indexes: bit maps, R-trees, inverted index

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

Hash table example

Index **Student_ID** on **Student.ID**

Data File **Student**



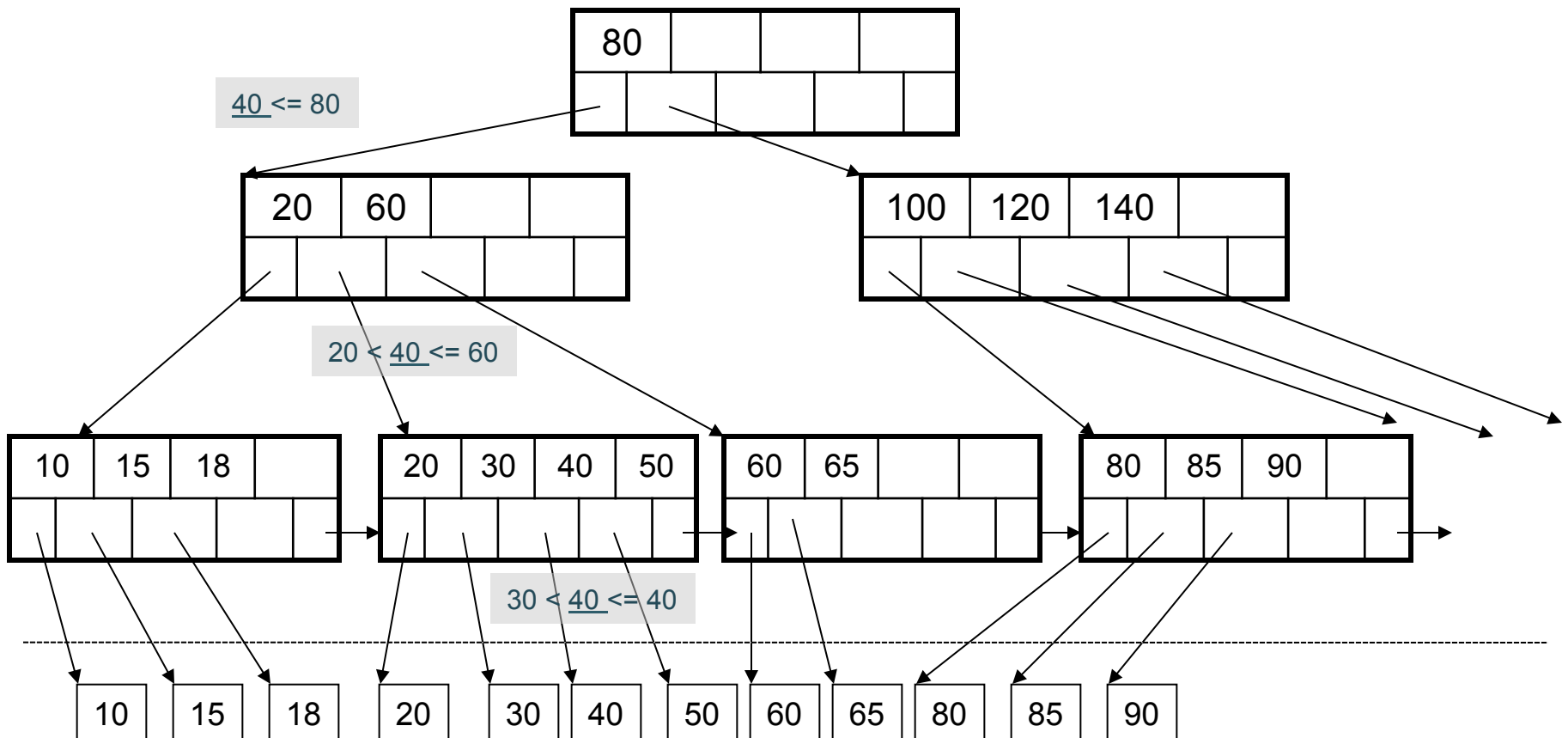
Index File
(preferably
in memory)

Data file
(on disk)

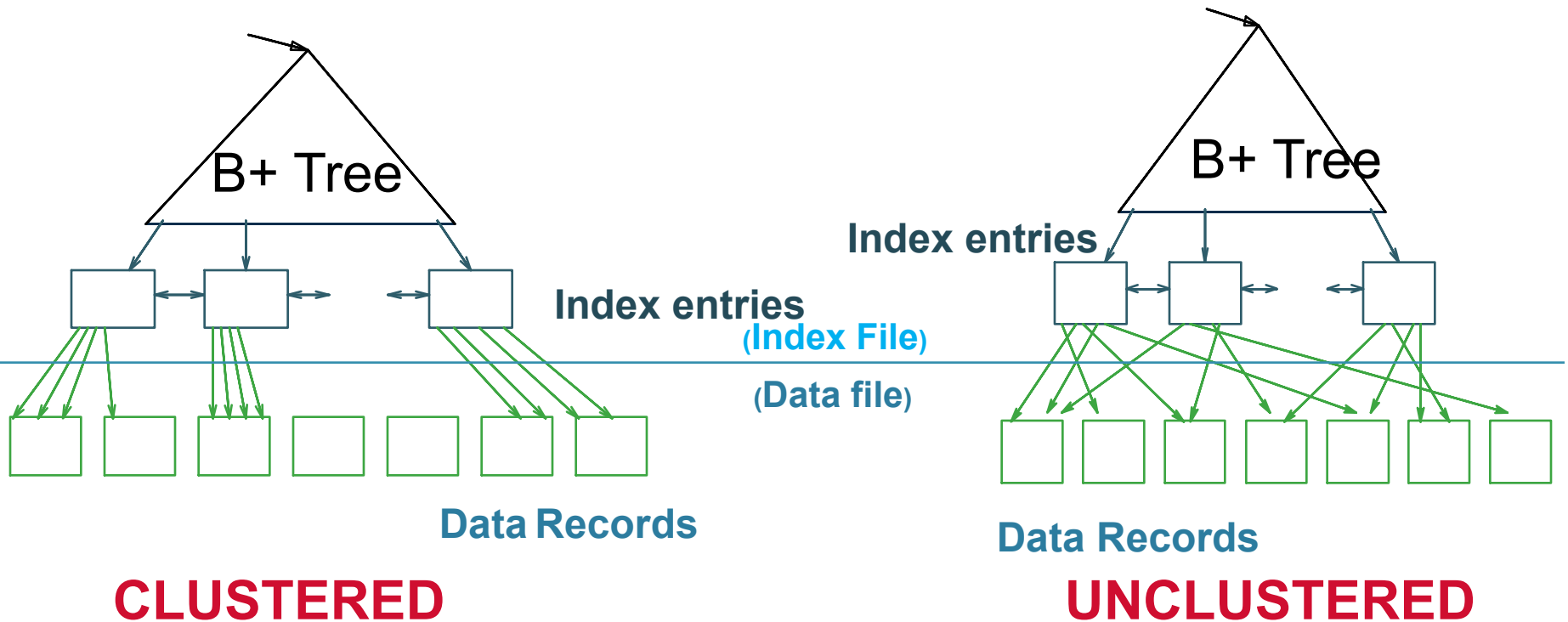
B+ Tree Index by Example

$d = 2$

Find the key 40



Clustered vs Unclustered



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

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

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:
 - Primary = is over attributes that include the primary key
 - Secondary = otherwise
 - Meaning 2: means the same as clustered/unclustered

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:
 - 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

Scanning a Data File

- Disks are mechanical devices!
 - Technology from the 60s; density much higher now
- Read only at the rotation speed!
- Consequence:
Sequential scan is MUCH FASTER than random reads
 - Good: read blocks 1,2,3,4,5,...
 - Bad: read blocks 2342, 11, 321,9, ...
- **Rule of thumb:**
 - Random reading 1-2% of the file \approx sequential scanning the entire file; this is decreasing over time (because of increased density of disks)
- Solid state (SSD): \$\$\$ expensive; put indexes, other “hot” data there, still too expensive for everything



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)

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

Example

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

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

Example

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

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

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

Example

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

Assume the database has indexes on these attributes:

- **Takes_courseID** = index on Takes.courseID
- **Student_ID** = index on Student.ID

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

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

Example

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

Assume the database has indexes on these attributes:

- **Takes_courseID** = index on Takes.courseID
- **Student_ID** = index on Student.ID

```
for y' in Takes_courseID where y'.courseID > 300
```

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

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

Example

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

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'
```


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

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

Example

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

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'
```

Index join

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

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

Example

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

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 *
```

Index join

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

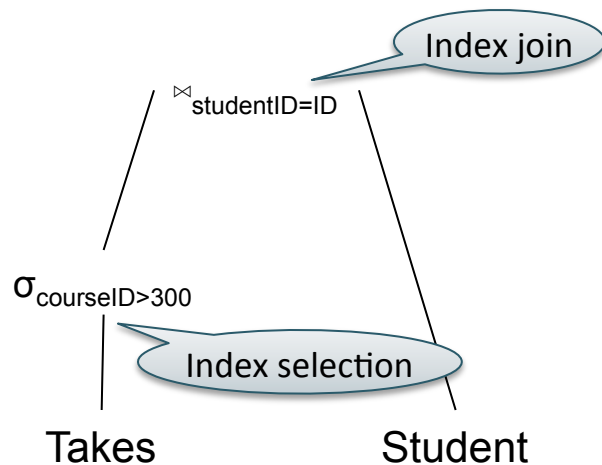
```
SELECT *  
FROM Student x, Takes y  
WHERE x.ID=y.studentID AND y.courseID > 300
```

Example

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

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 V2 ON V(P, M)
```

```
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
```

yes

```
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
```

yes

```
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
```

yes

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

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

no

```
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
```

yes

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

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

no

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

Not supported
in SQLite

Which Indexes?

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

- How many indexes **could** we create?
- Which indexes **should** we create?

Which Indexes?

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

- How many indexes **could** we create?
- Which indexes **should** we create?

In general this is a very hard problem

Which Indexes?

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

- 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

Which Indexes?

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

- 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)

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 (?, ?, ?)
```

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:

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

1000000 queries:

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

100000 queries:

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

What indexes ?

The Index Selection Problem 3

V(M, N, P);

Your workload is this

100000 queries:

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

1000000 queries:

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

100000 queries:

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

A: V(N, P)

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

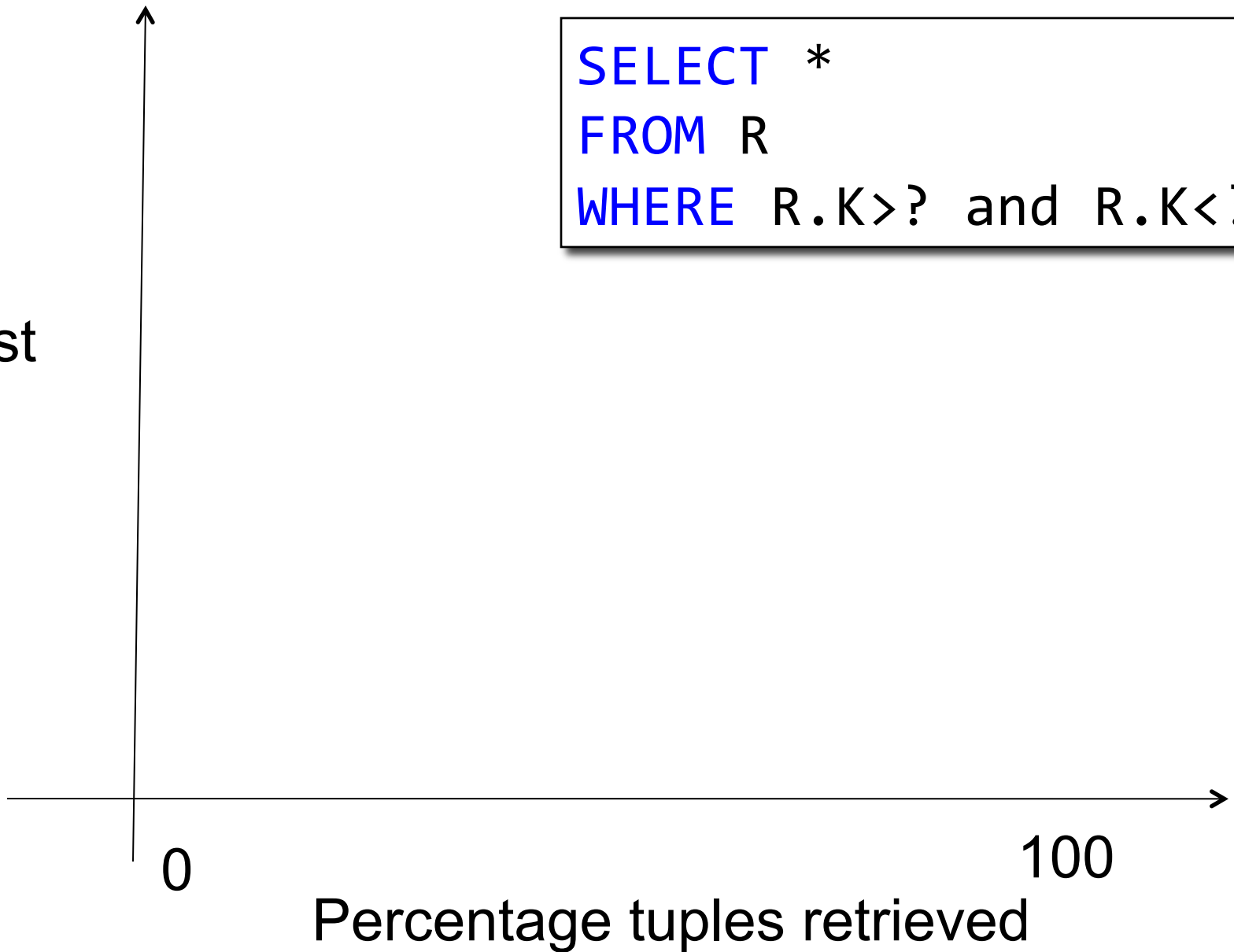
- Consider queries in workload in order of importance
- Consider relations accessed by query
 - No point indexing other relations
- Look at WHERE clause for possible search key
- Try to choose indexes that speed-up multiple queries

To Cluster or Not

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

```
SELECT *  
FROM R  
WHERE R.K>? and R.K<?
```

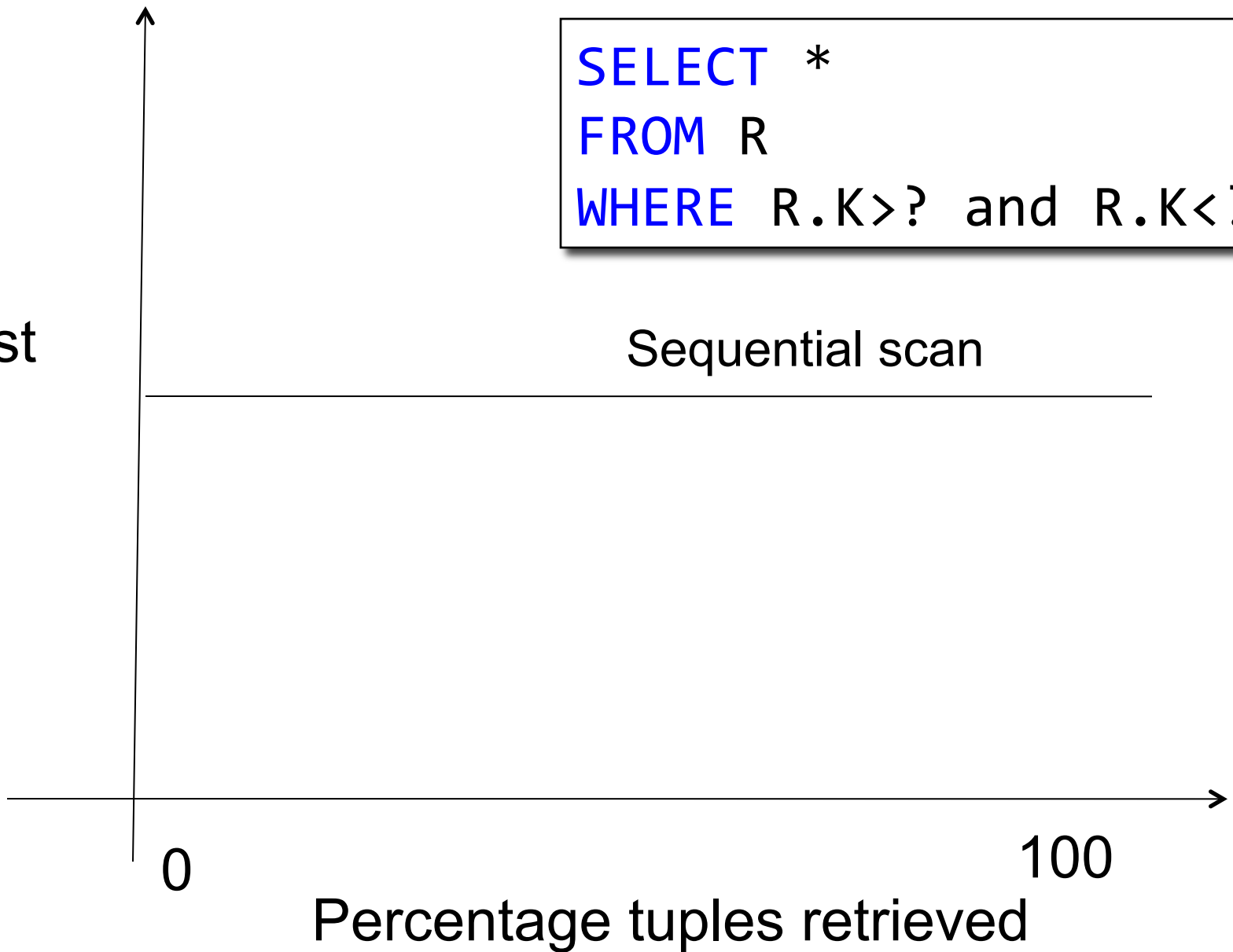
Cost



```
SELECT *  
FROM R  
WHERE R.K>? and R.K<?
```

Cost

Sequential scan



```
SELECT *  
FROM R  
WHERE R.K>? and R.K<?
```

Cost

Sequential scan

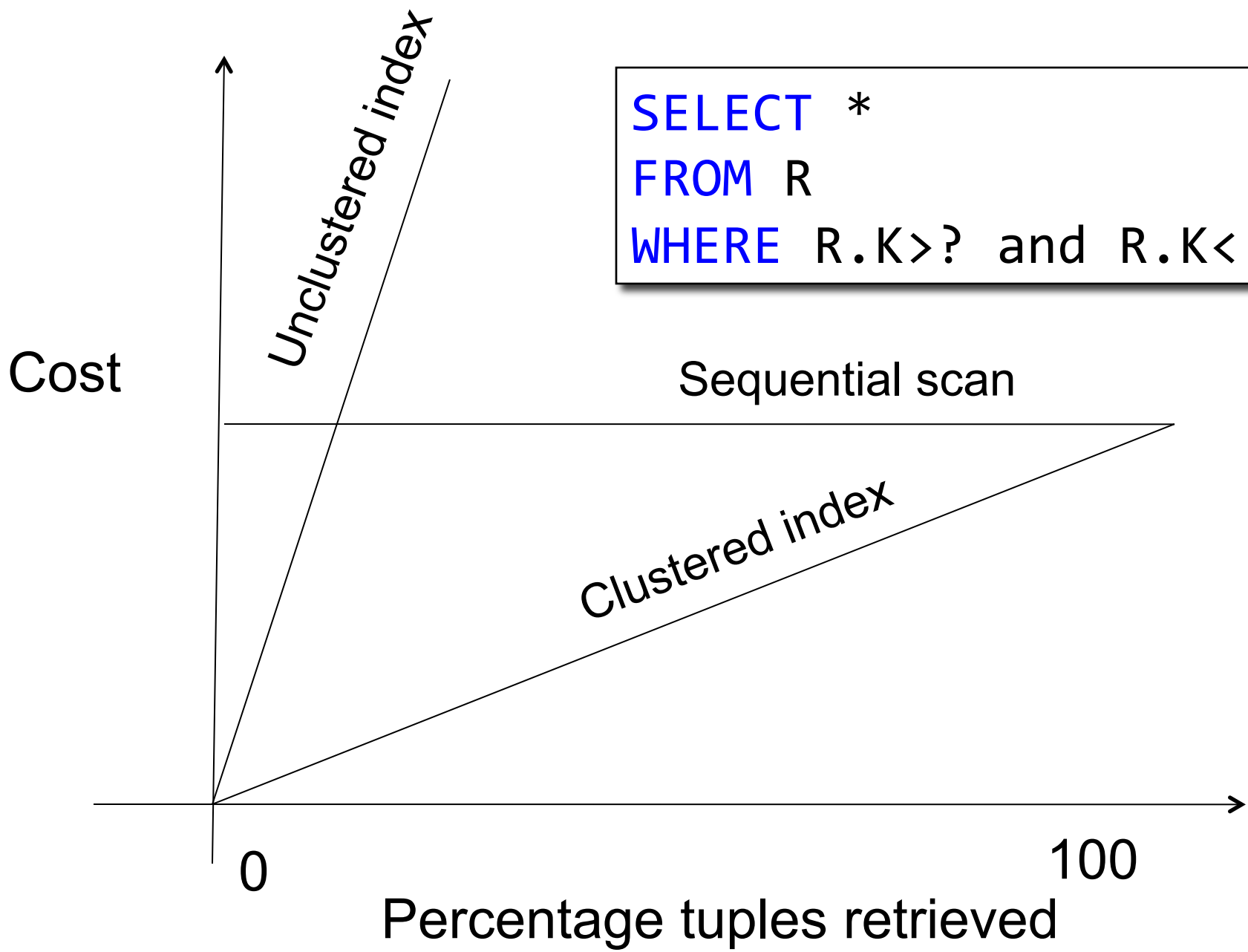
Clustered index

0

100

Percentage tuples retrieved

```
SELECT *  
FROM R  
WHERE R.K>? and R.K<?
```



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 = $(c - \min(R, A)) / (\max(R, A) - \min(R, A))$
- $c1 < A < c2$ $/* \sigma_{c1 < A < c2}(R) */$
 - Selectivity = $(c2 - c1) / (\max(R, A) - \min(R, A))$

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: $f \cdot B(R)$
 - Unclustered index $f \cdot T(R)$

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:

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

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

- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:

Index Based Selection

- Example:

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

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:

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

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:

Index Based Selection

- Example:

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

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

Index Based Selection

- Example:

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

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

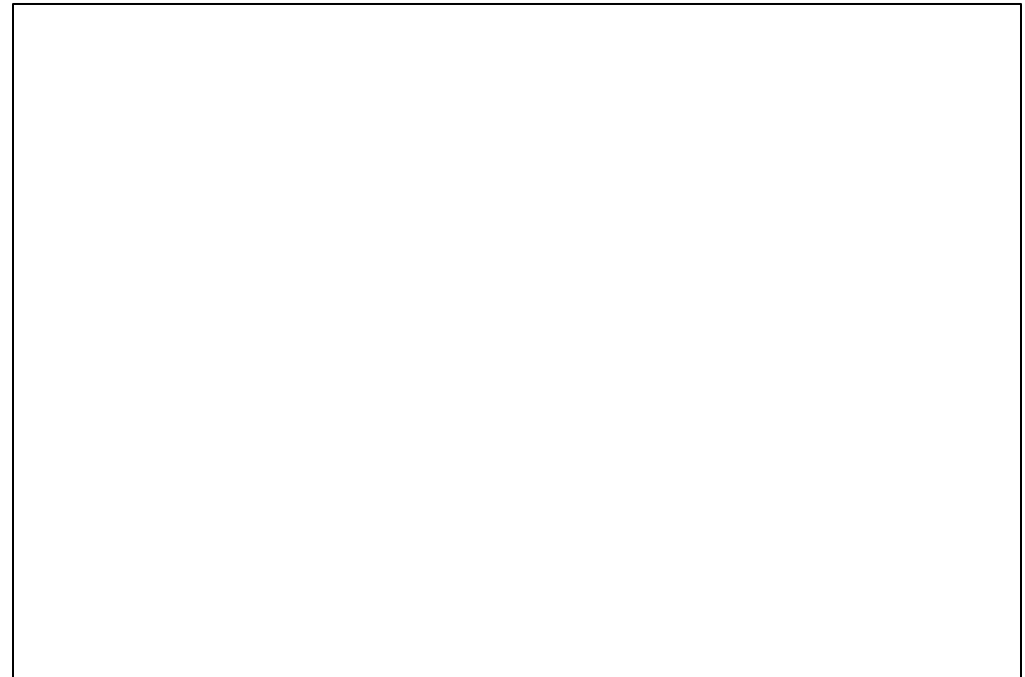
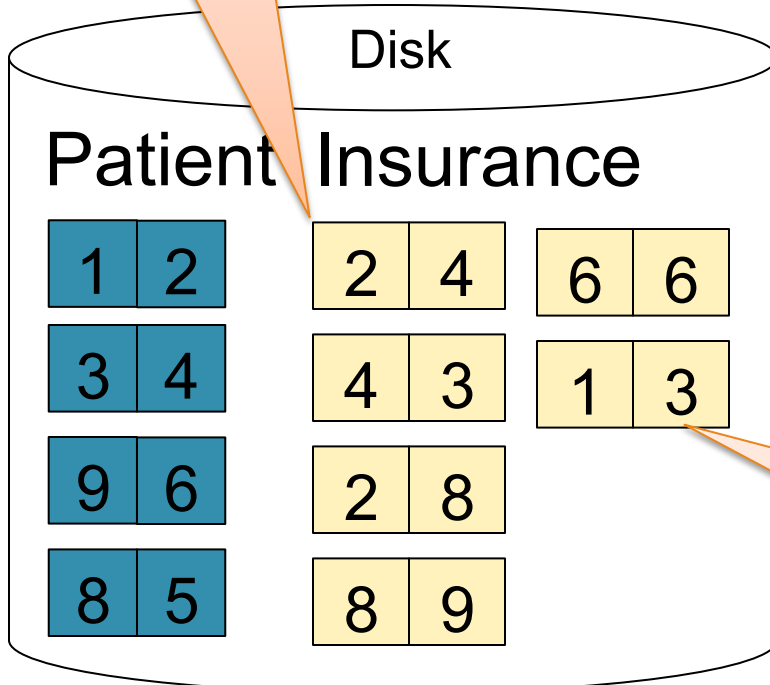
Hash Join Example

Patient \bowtie Insurance

Some large-enough #

Memory M = 21 pages

Showing pid only

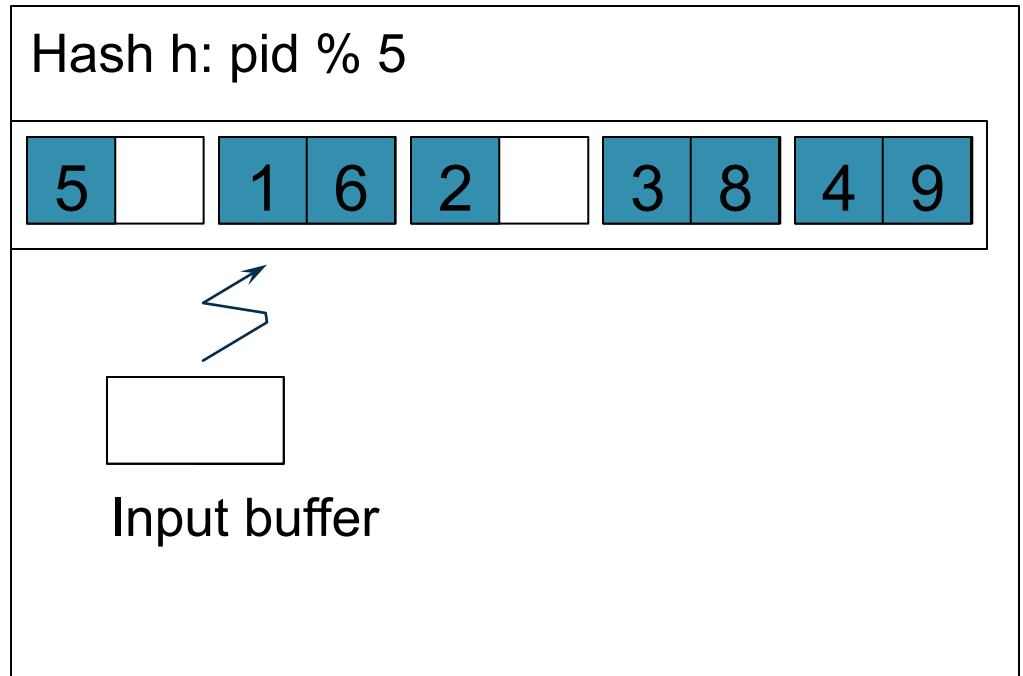
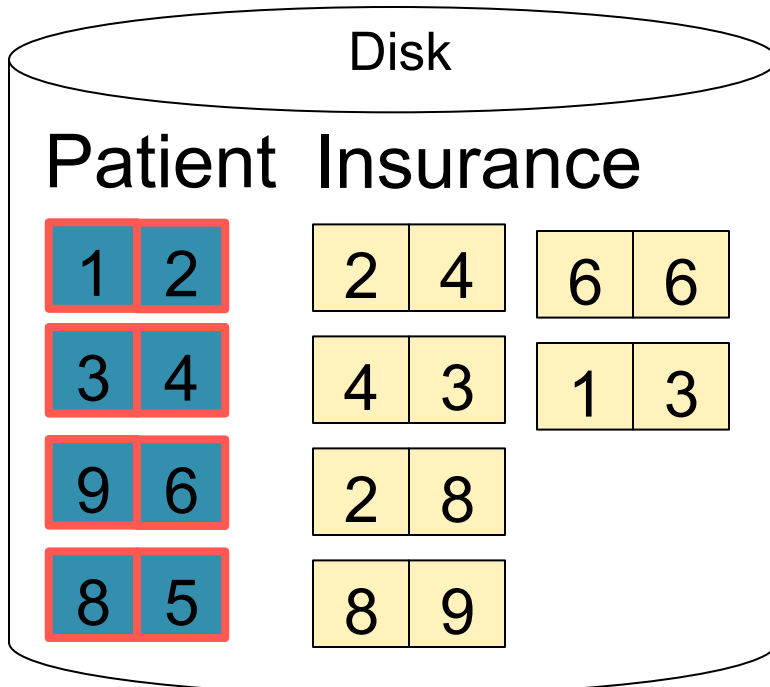


This is one page with two tuples

Hash Join Example

Step 1: Scan Patient and **build** hash table in memory
Can be done in method open()

Memory M = 21 pages



Hash 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
---	--	---	---	---	--	---	---	---	---

Disk

Patient Insurance

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

Disk

Patient Insurance

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

2	4
---	---

Input buffer

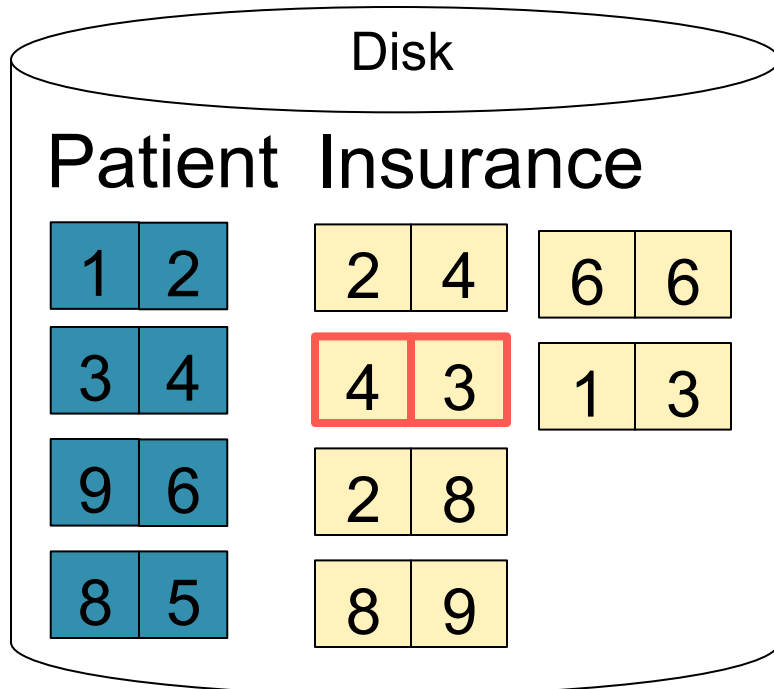
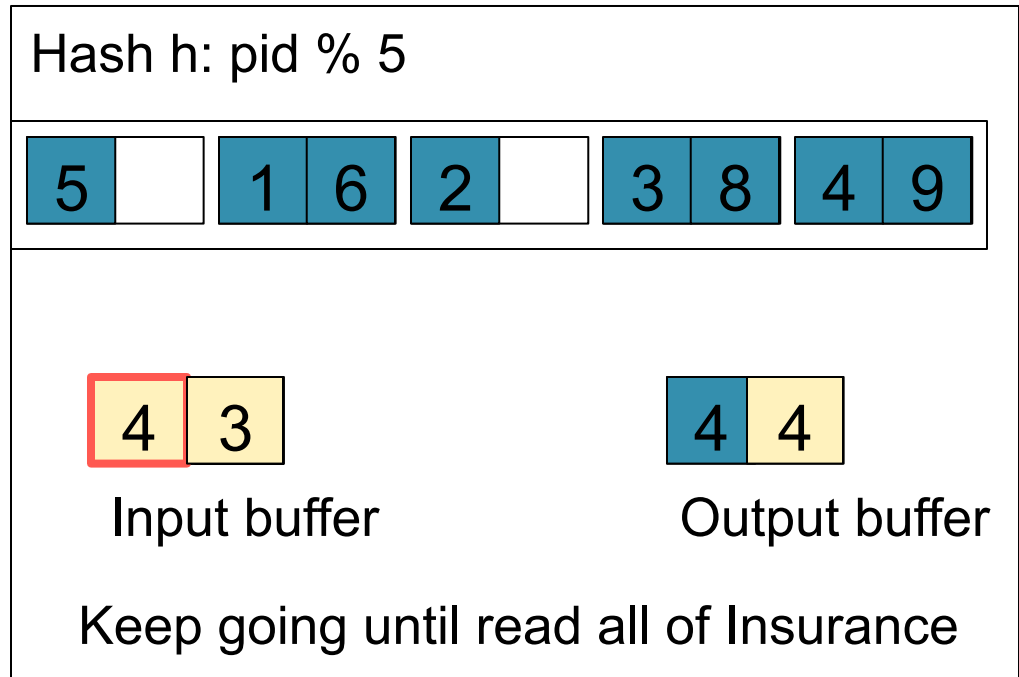
4	4
---	---

Output buffer

Hash Join Example

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

Memory M = 21 pages



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  $(t_1, t_2)$ 
```

What is the Cost?

- Cost: $B(R) + T(R) B(S)$
- Multiple-pass since S is read many times

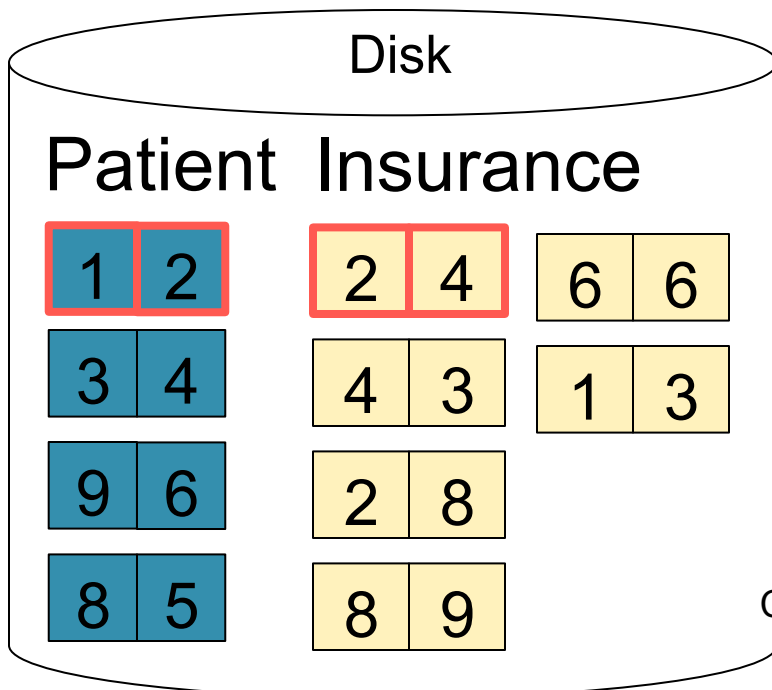
Page-at-a-time Refinement

```
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)
```

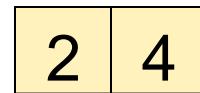
- Cost: $B(R) + B(R)B(S)$

What is the Cost?

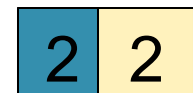
Page-at-a-time Refinement



Input buffer for Patient

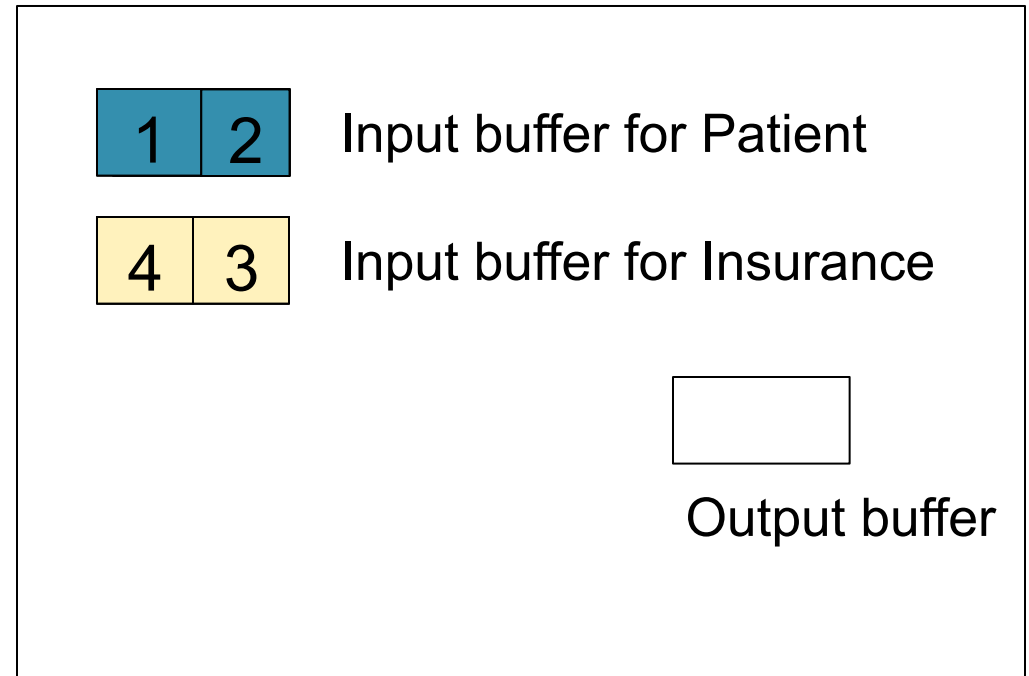
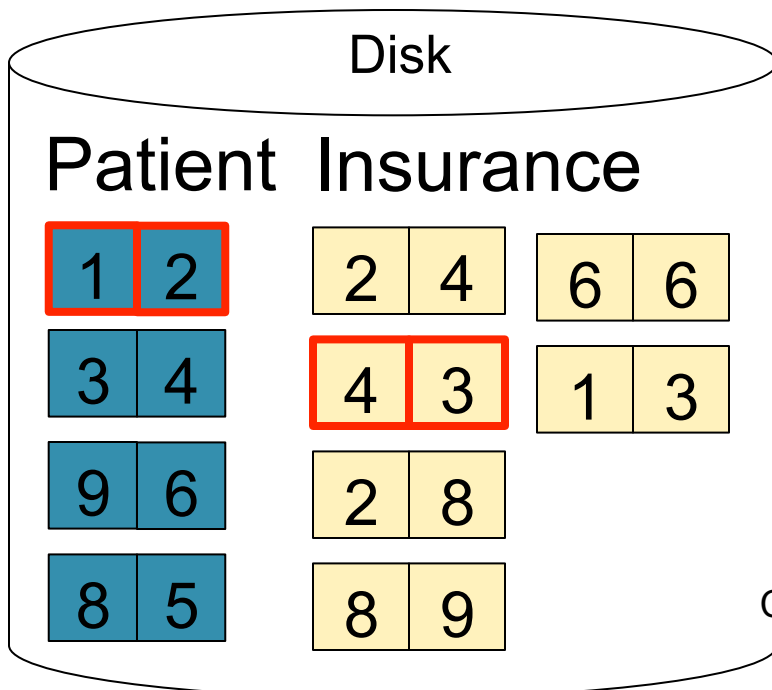


Input buffer for Insurance

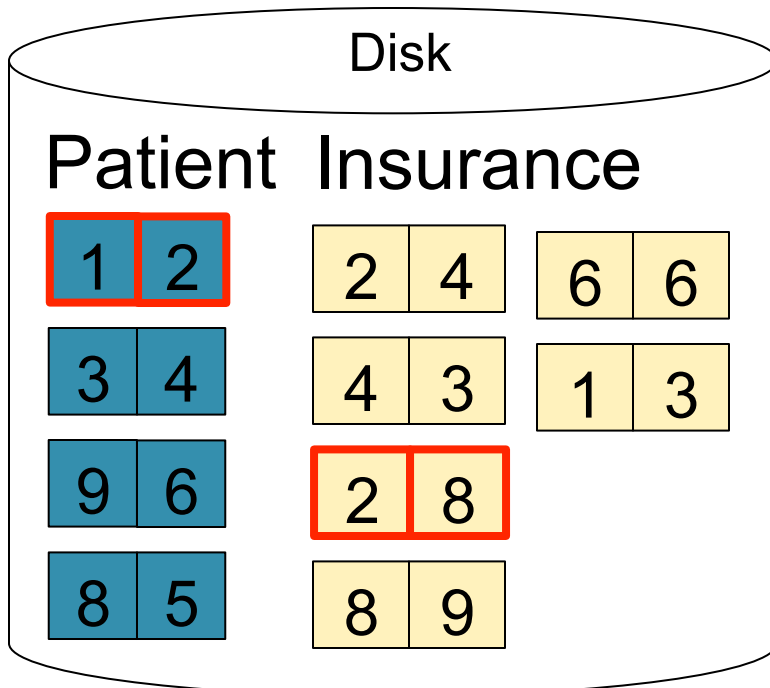


Output buffer

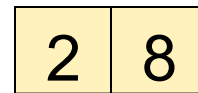
Page-at-a-time Refinement



Page-at-a-time Refinement

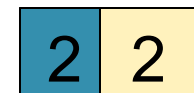


Input buffer for Patient



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: $R \bowtie S$

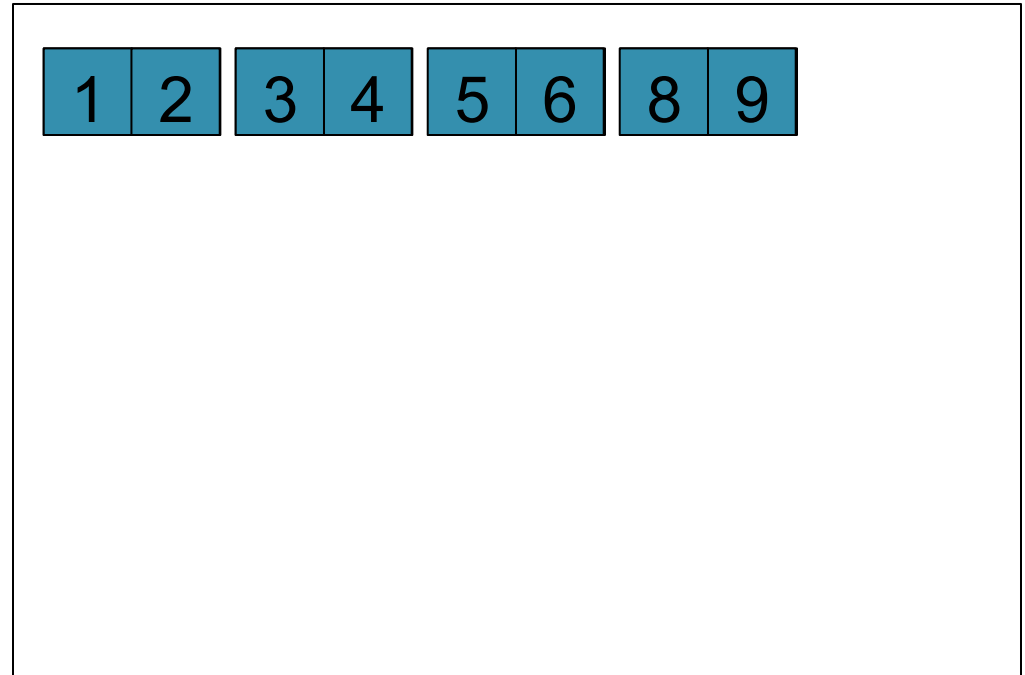
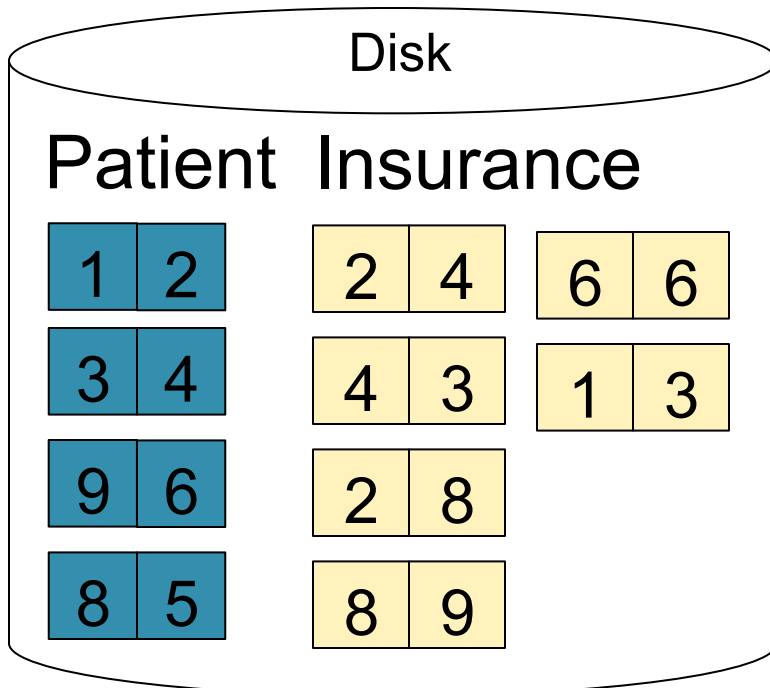
- Scan R and sort in main memory
- Scan S and sort in main memory
- Merge R and S

- 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

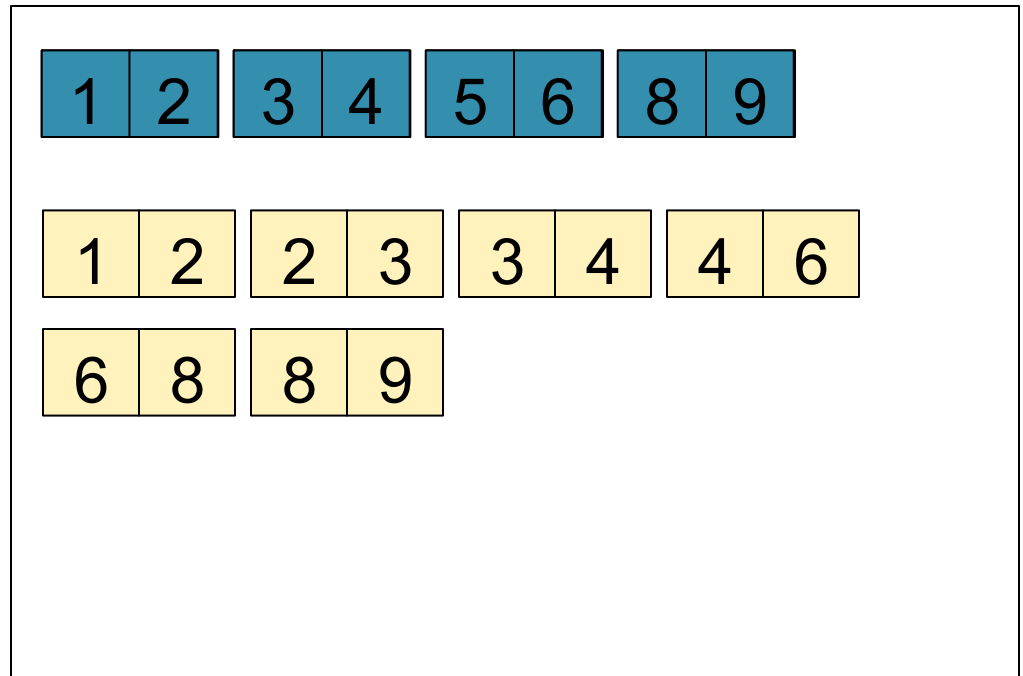
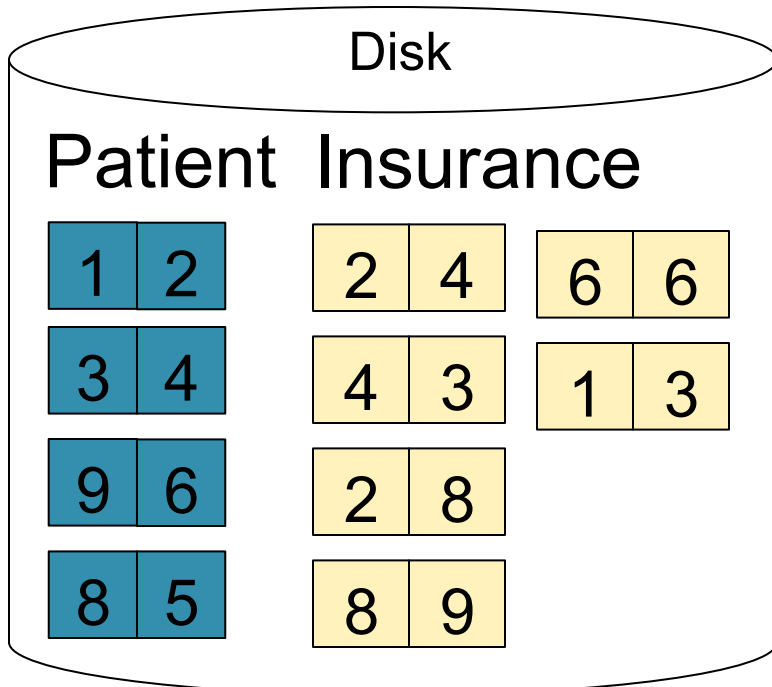
Memory M = 21 pages



Sort-Merge Join Example

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

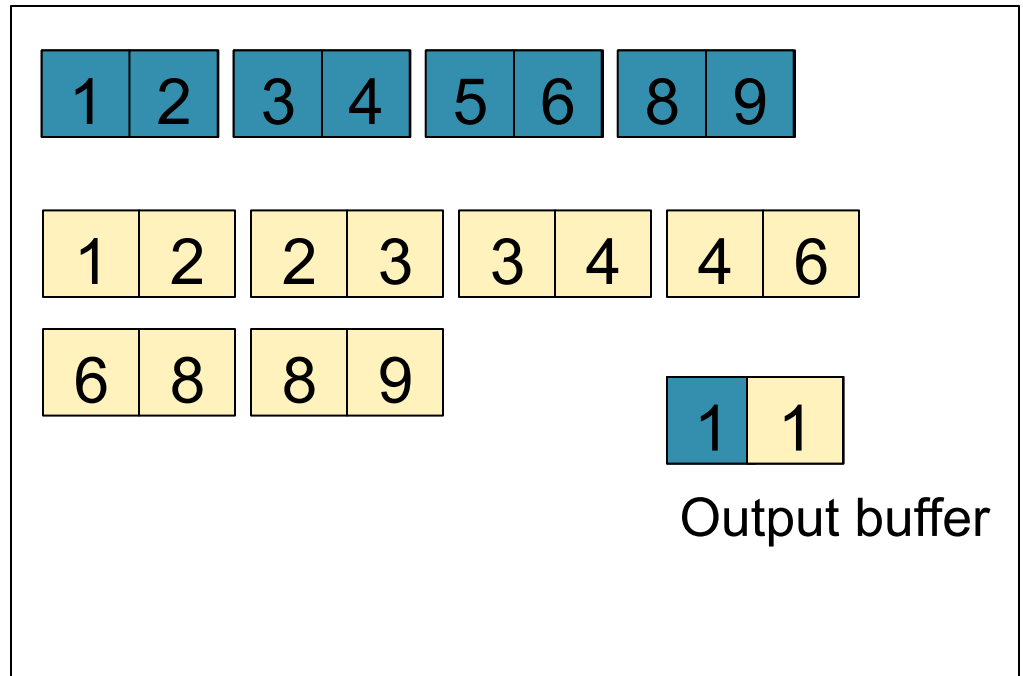
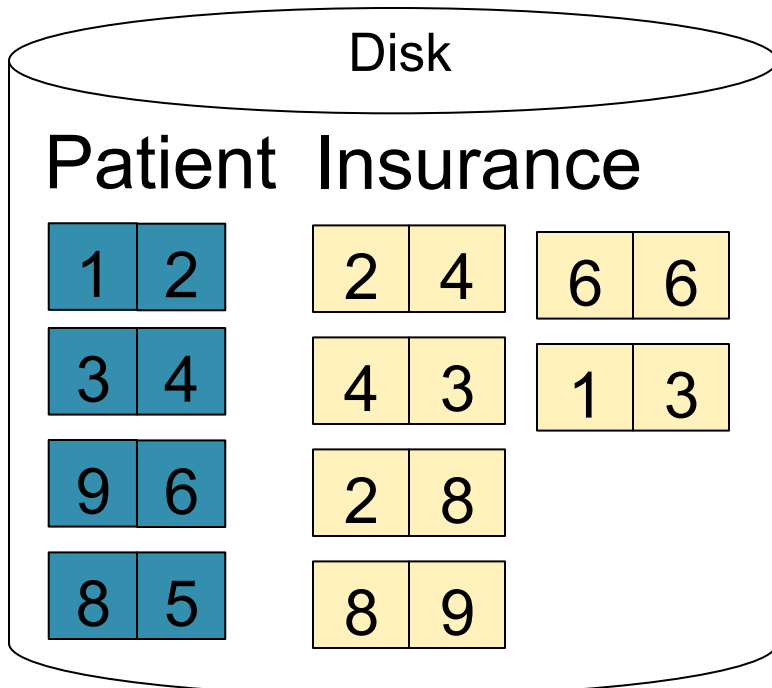
Memory M = 21 pages



Sort-Merge Join Example

Step 3: Merge Patient and Insurance

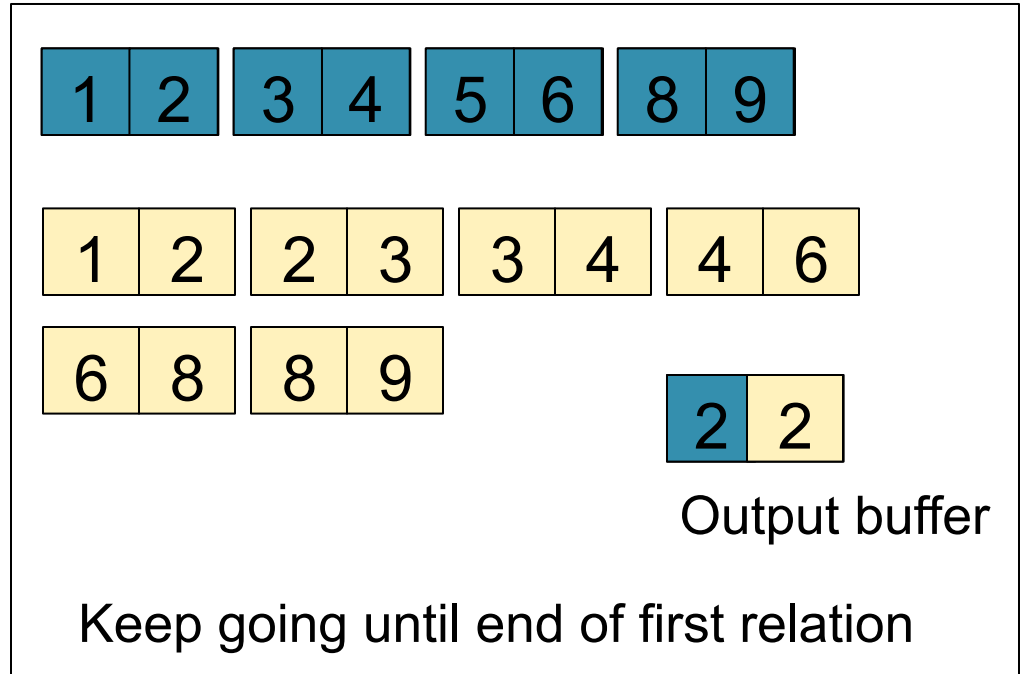
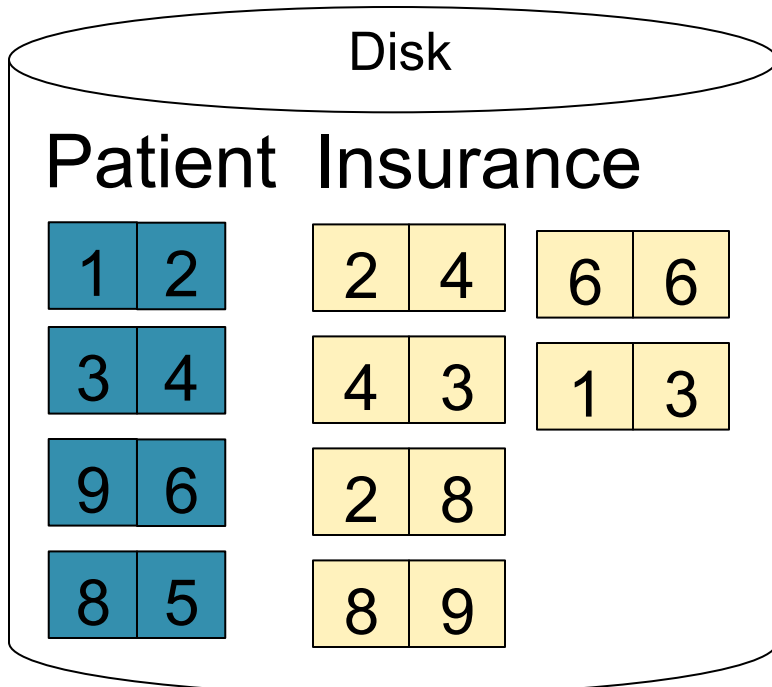
Memory M = 21 pages



Sort-Merge Join Example

Step 3: Merge Patient and Insurance

Memory M = 21 pages



Index Nested Loop Join

$R \bowtie S$

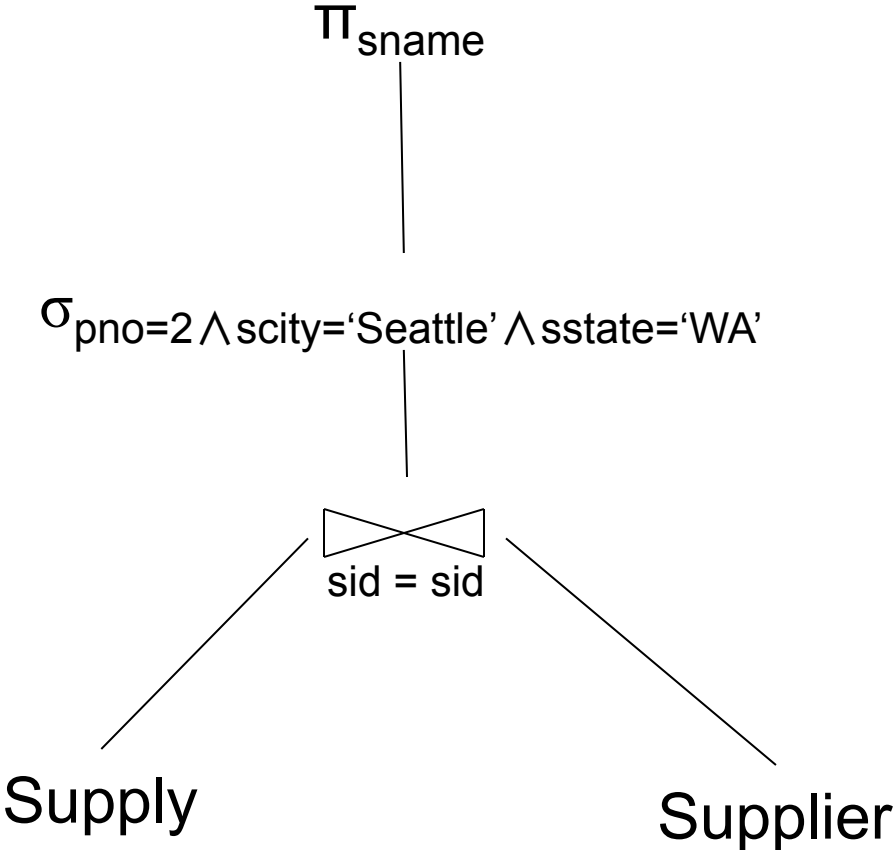
- Assume S has an index on the join attribute
- Iterate over R , for each tuple fetch corresponding tuple(s) from S
- **Cost:**
 - If index on S is clustered:
$$B(R) + T(R) * (B(S) * 1/V(S,a))$$
 - If index on S is unclustered:
$$B(R) + T(R) * (T(S) * 1/V(S,a))$$

Cost of Query Plans

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

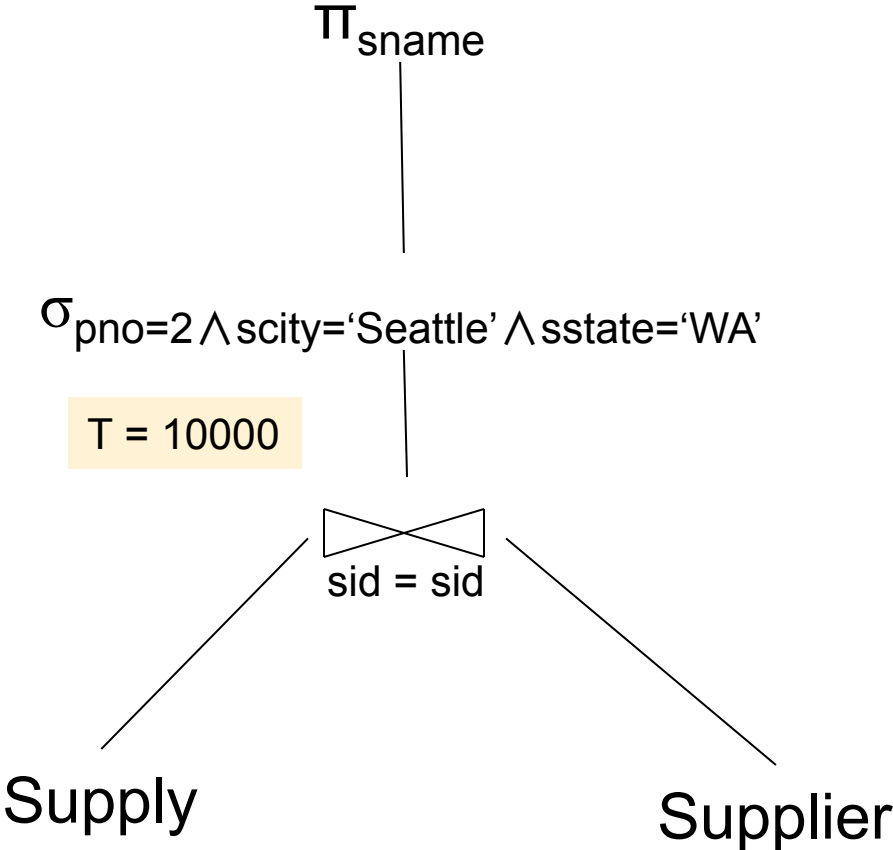
T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

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

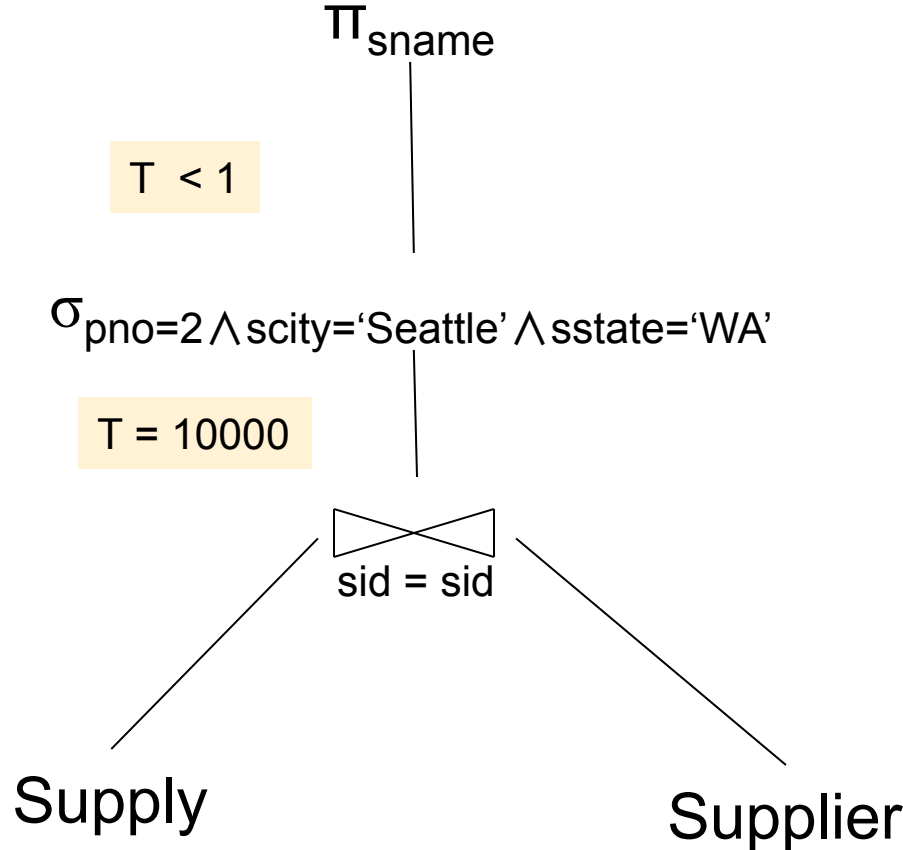
T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

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

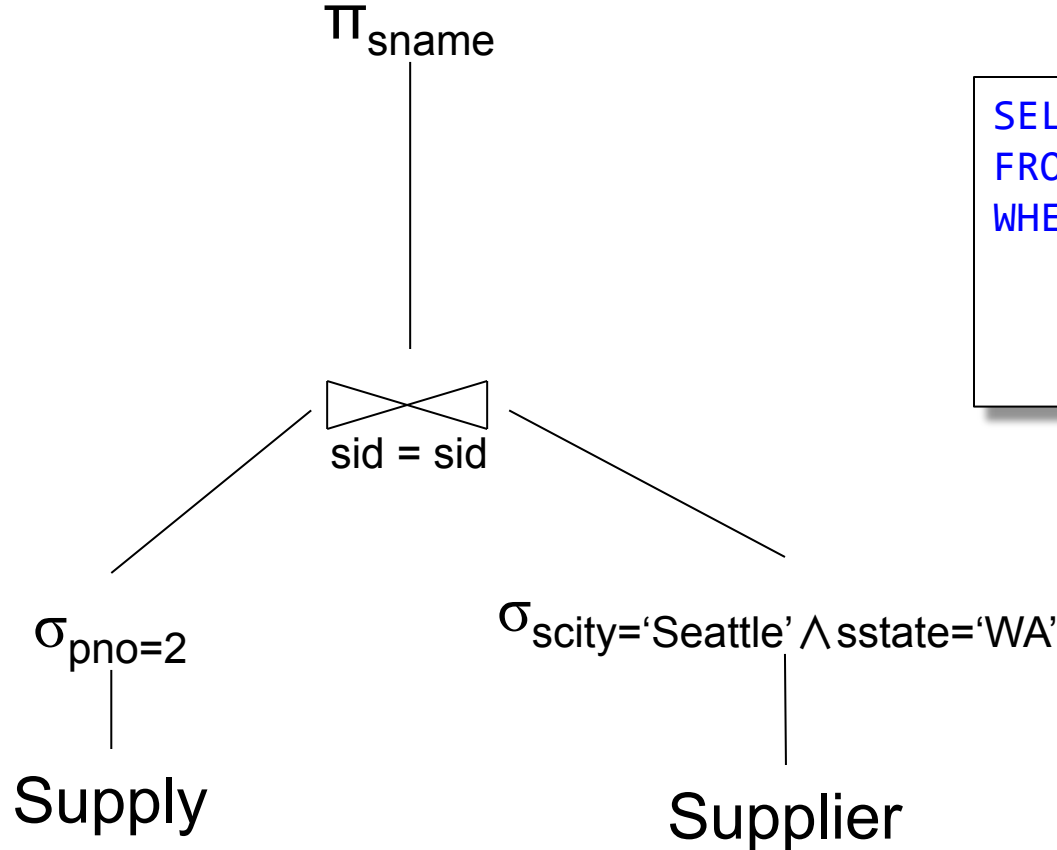
T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

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

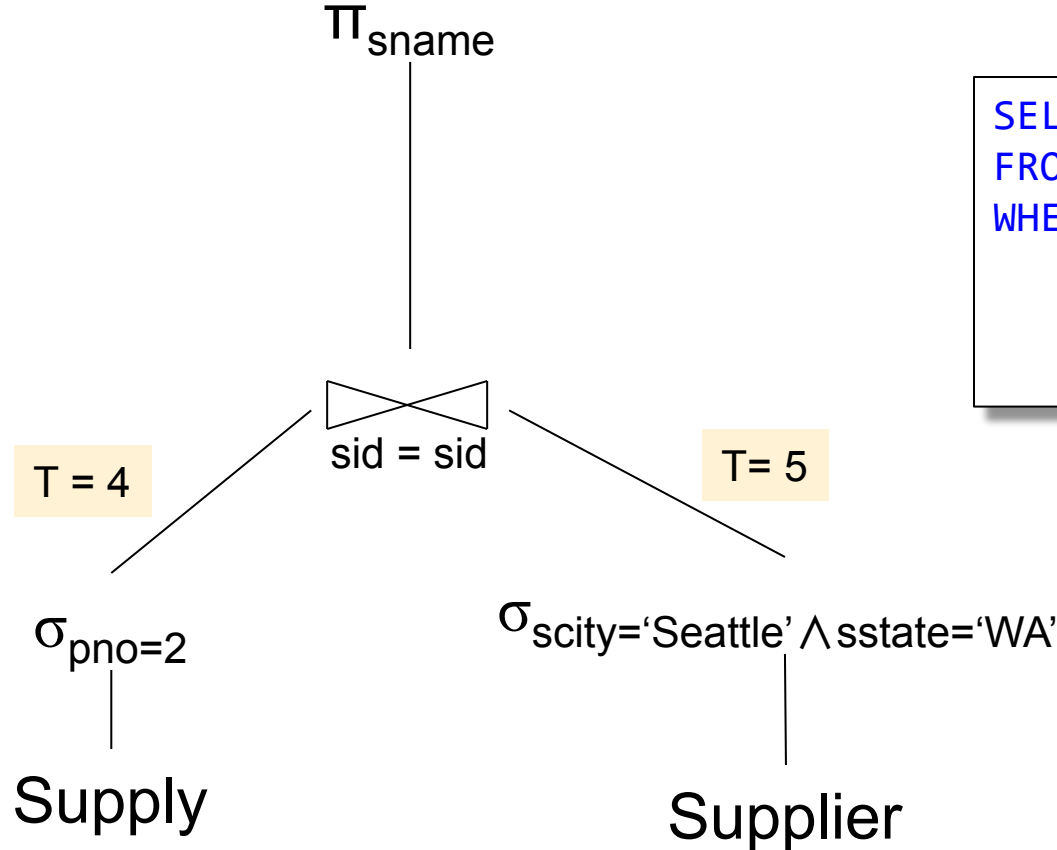
M=11

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

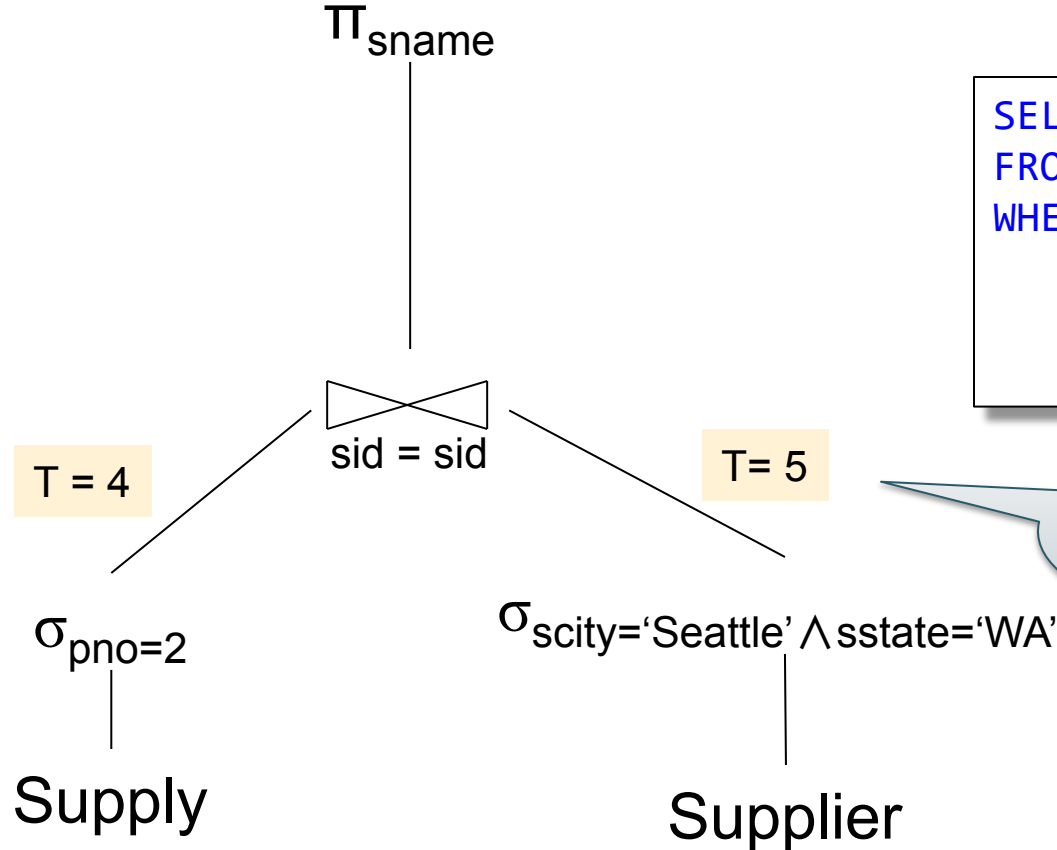
M=11

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

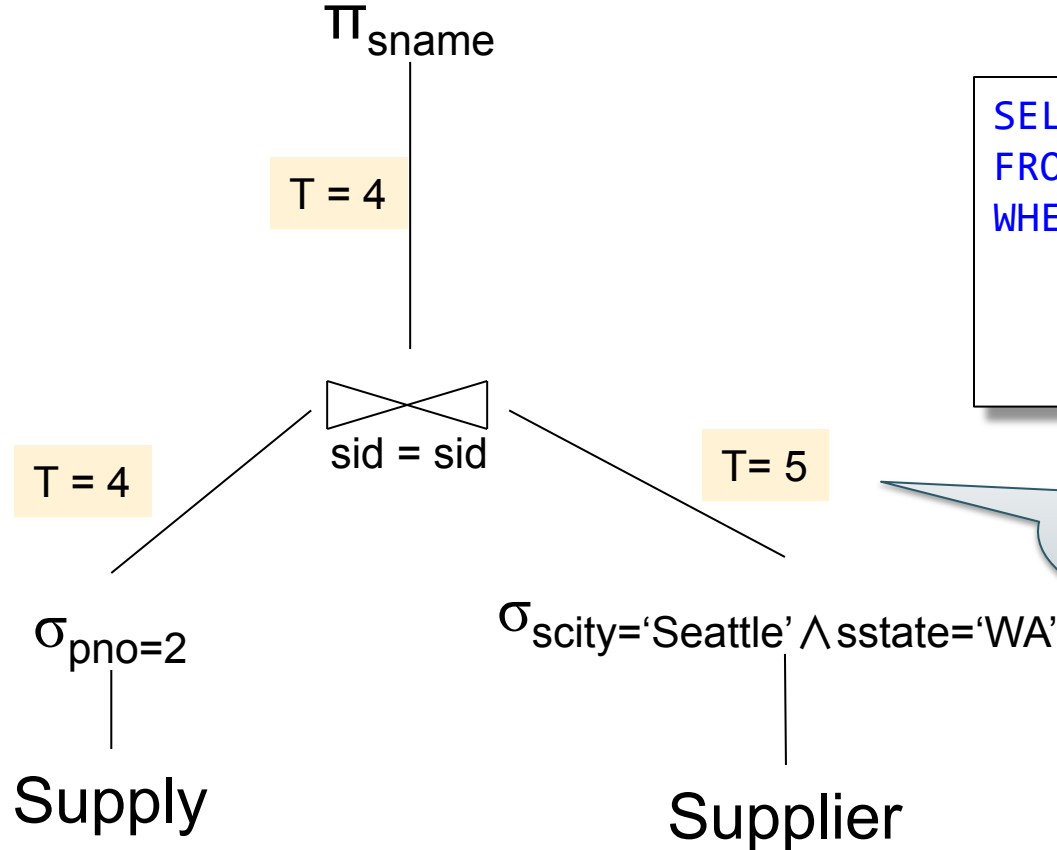
T(Supplier) = 1000
B(Supplier) = 100
V(Supplier, scity) = 20
V(Supplier, state) = 10

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

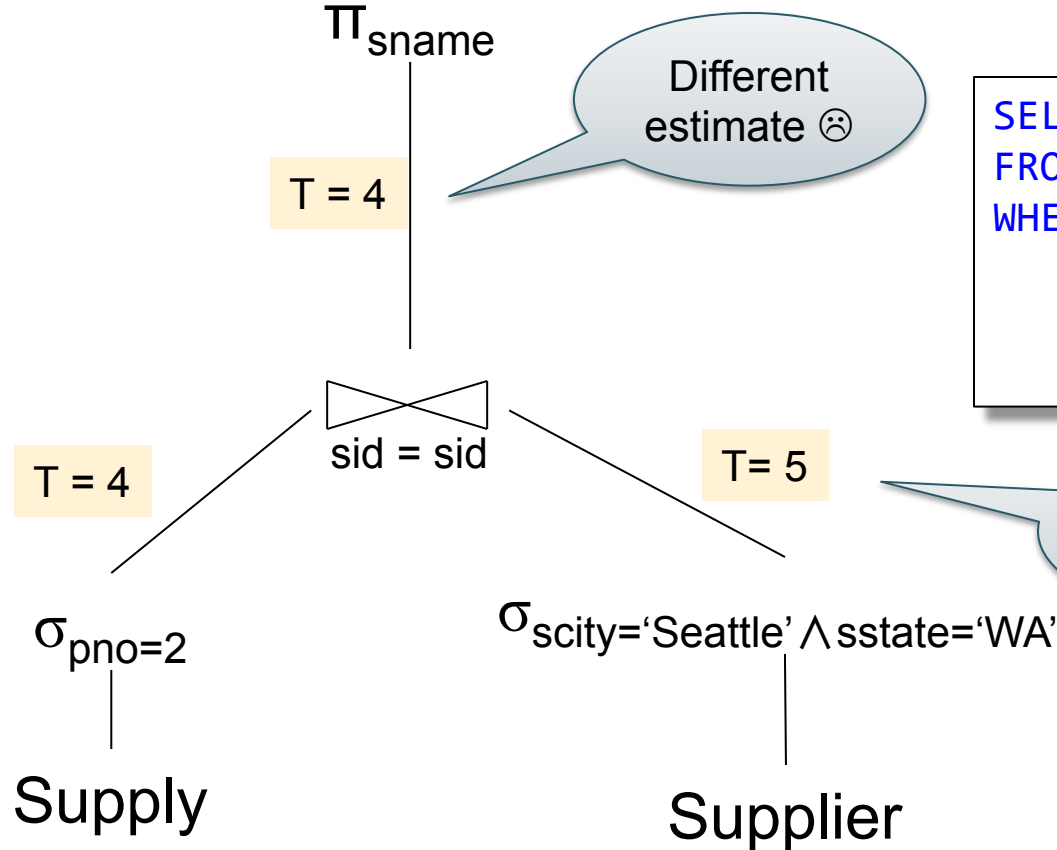
T(Supplier) = 1000
 B(Supplier) = 100
 V(Supplier, scity) = 20
 V(Supplier, state) = 10

M=11

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Logical Query Plan 2



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?

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

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

M=11

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Physical Plan 1

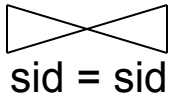
Π_{sname}

T < 1

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

T = 10000

Total cost:



sid = sid

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

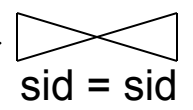
Π_{sname}

T < 1

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

T = 10000

Total cost: $100/10 * 100 = 1000$



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

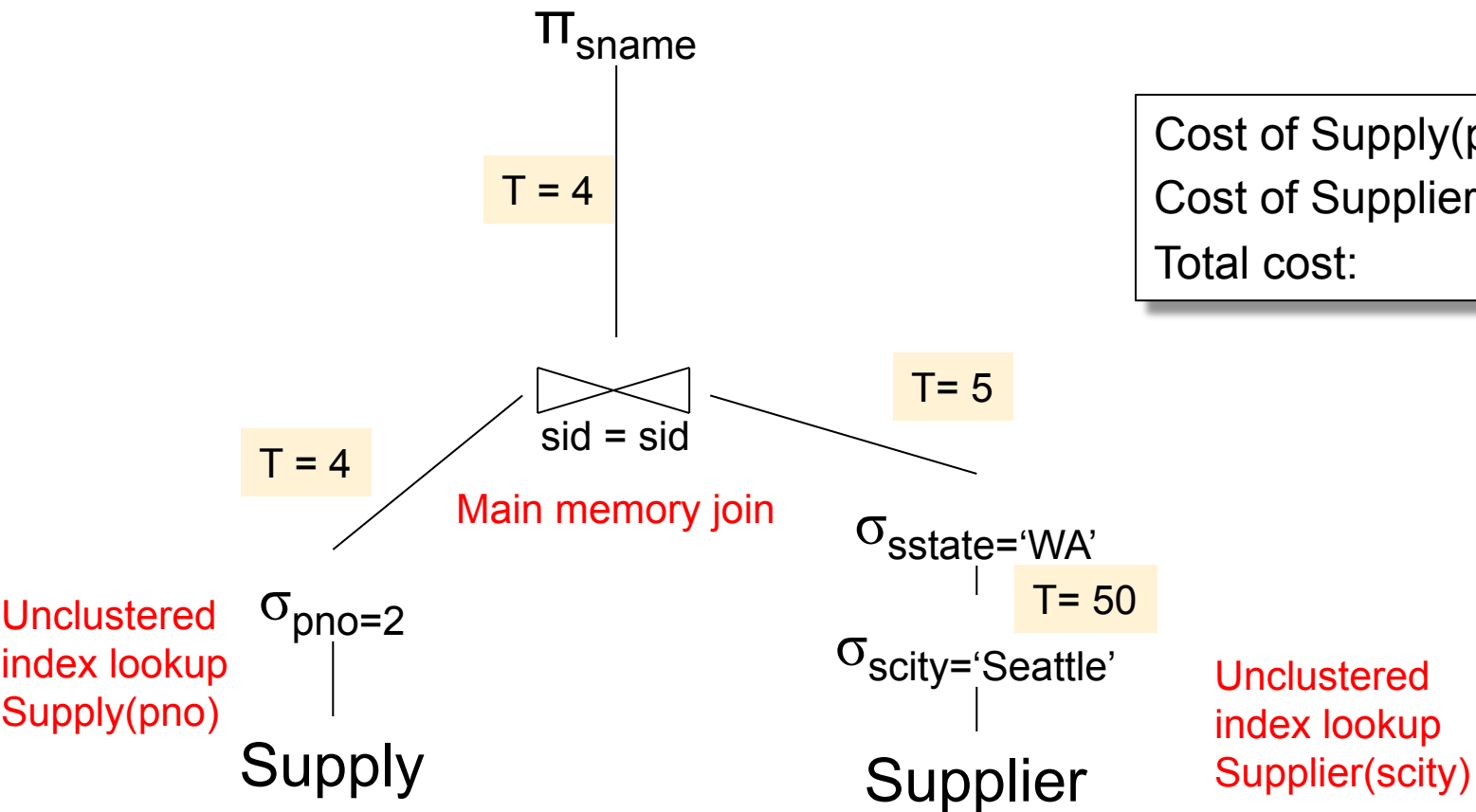
M=11

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Physical Plan 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

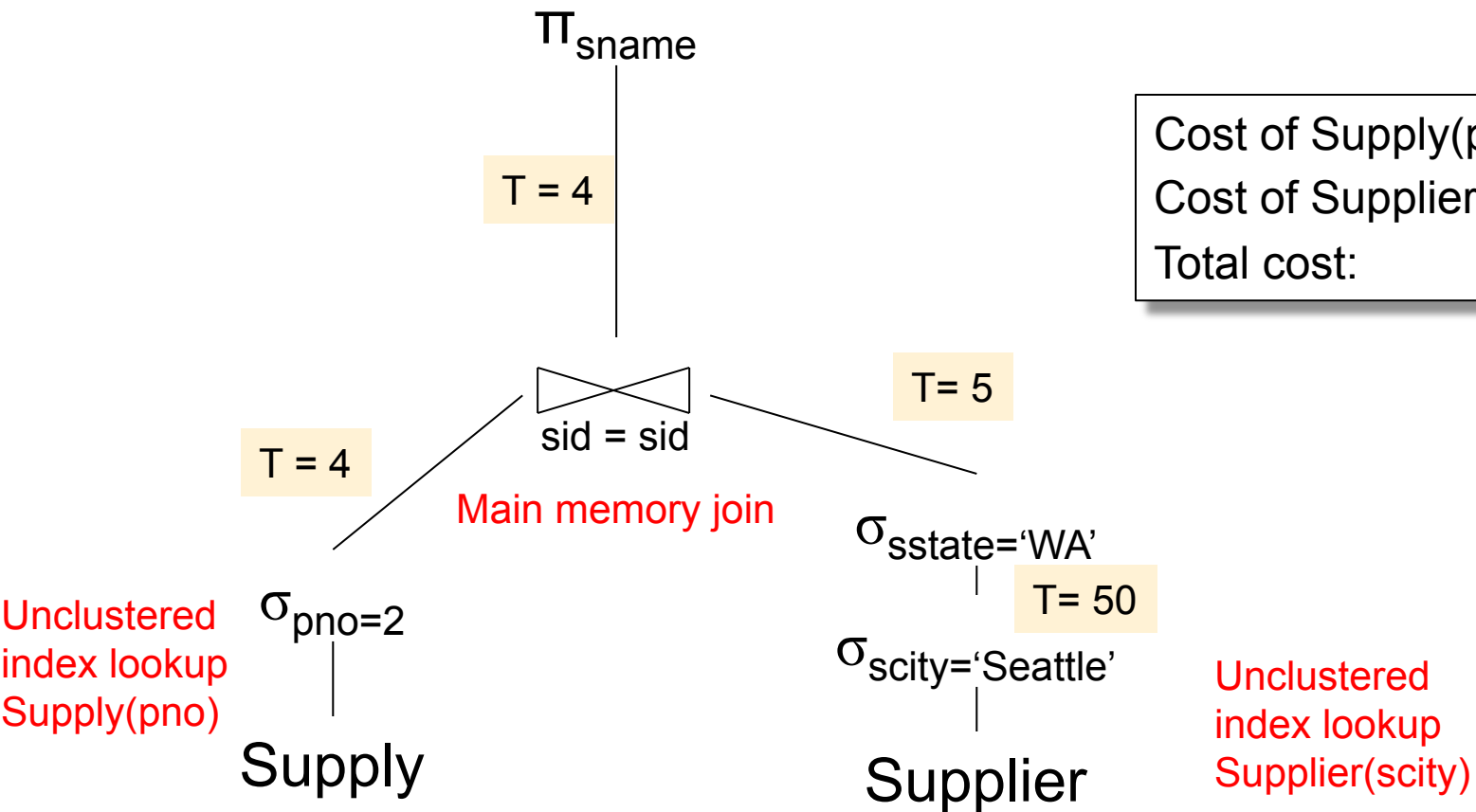
M=11

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Physical Plan 2

Cost of Supply(pno) = 4
Cost of Supplier(scity) =
Total cost:



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

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

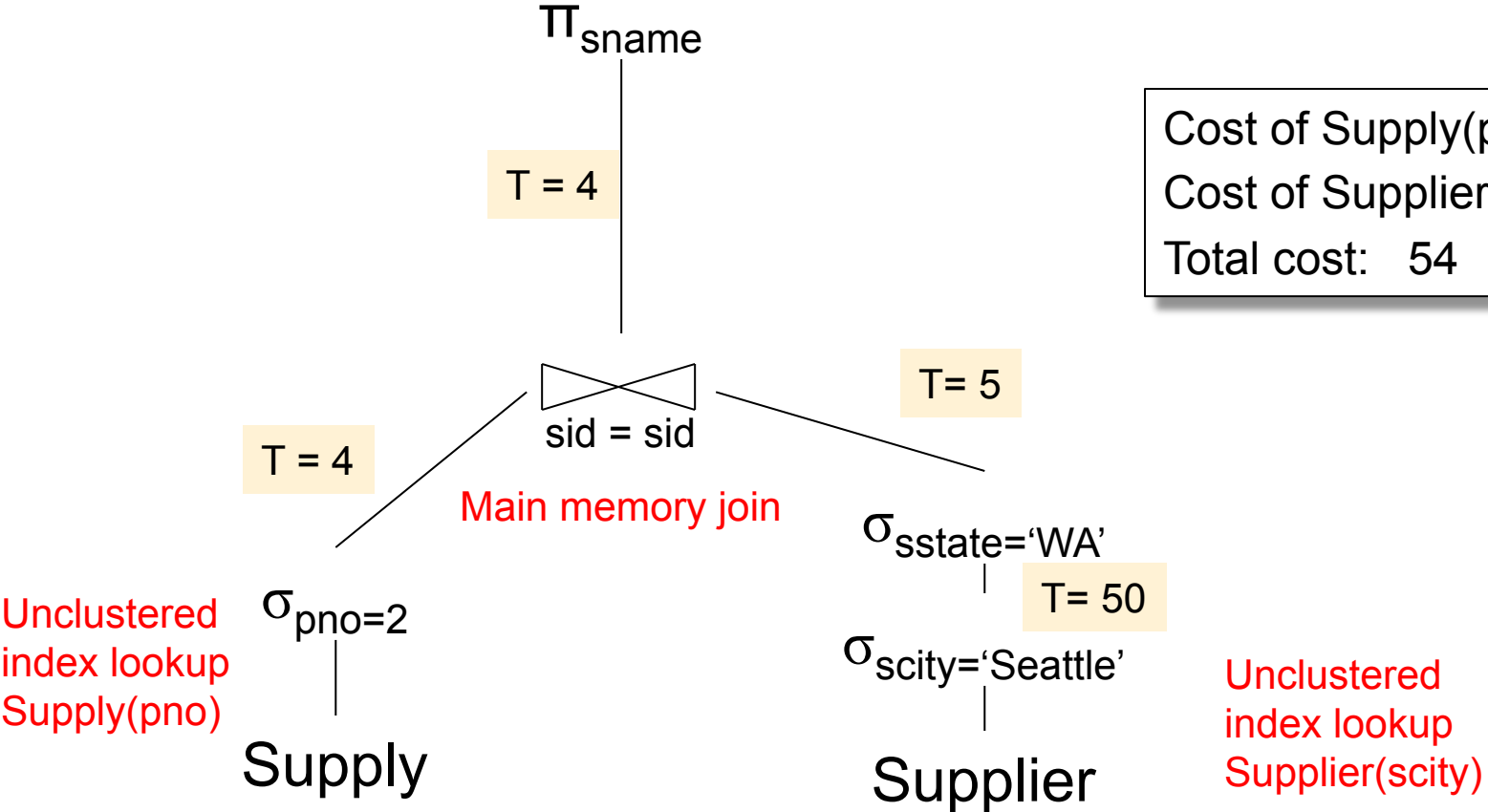
$M=11$

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Physical Plan 2

Cost of Supply(pno) = 4
Cost of Supplier(scity) = 50
Total cost: 54



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) =
Cost of Index join =
Total cost:

T = 4

sid = sid

Clustered
Index join

Unclustered
index lookup
Supply(pno)

$\sigma_{pno=2}$

Supply

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 =
Total cost:

T = 4

sid = sid

Clustered
Index join

Unclustered
index lookup
Supply(pno)

$\sigma_{pno=2}$

Supply

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

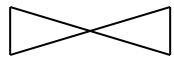
T = 4

Π_{sname}

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

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

T = 4



sid = sid

Clustered
Index join

Unclustered
index lookup
Supply(pno)

$\sigma_{pno=2}$

Supply

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

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