# Introduction to Database Systems
# CSE 344

## Lecture 11:
## Basics of Query Optimization and
## Query Cost Estimation

# Motivation

- My database application is too slow… why?
- One of the queries is very slow… why?
- …


- To understand performance, need to understand how a DBMS works

# Recap

- What is a disk block? (A.k.a. page)

- What is an index?

- What are clustered/unclustered indexes?

# Recap – Indexes

V(M, N, P);

SELECT *
FROM V
WHERE V.M = 33

SELECT *
FROM V
WHERE V.M = 33 and V.P = 55

Suppose we only had _one_ of these indexes. Can the optimizer use it?

INDEX I1 on V(M)

INDEX I2 on V(M,P)

INDEX I3 on V(P,M)

# Recap – Indexes

Movie(mid, title, year)

CLUSTERED INDEX I on Movie(id)
INDEX J on Movie(year)

```
SELECT *
FROM Movie
WHERE year = 2010
```

The system uses the index
J for one of the queries,
but not for the other.

```
SELECT *
FROM Movie
WHERE year = 1910
```

Which and why?

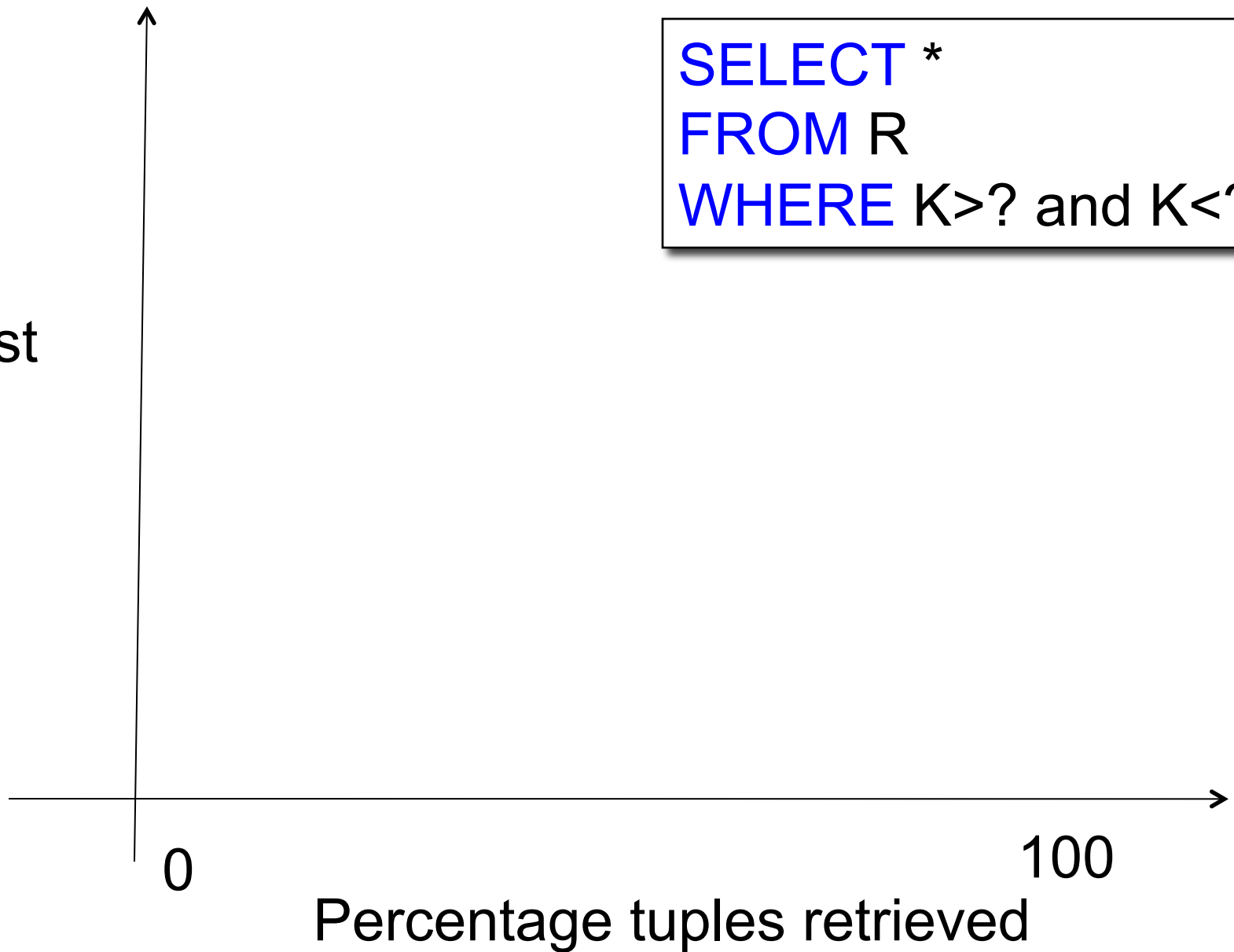# Basic Index Selection Guidelines

- 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

Cost

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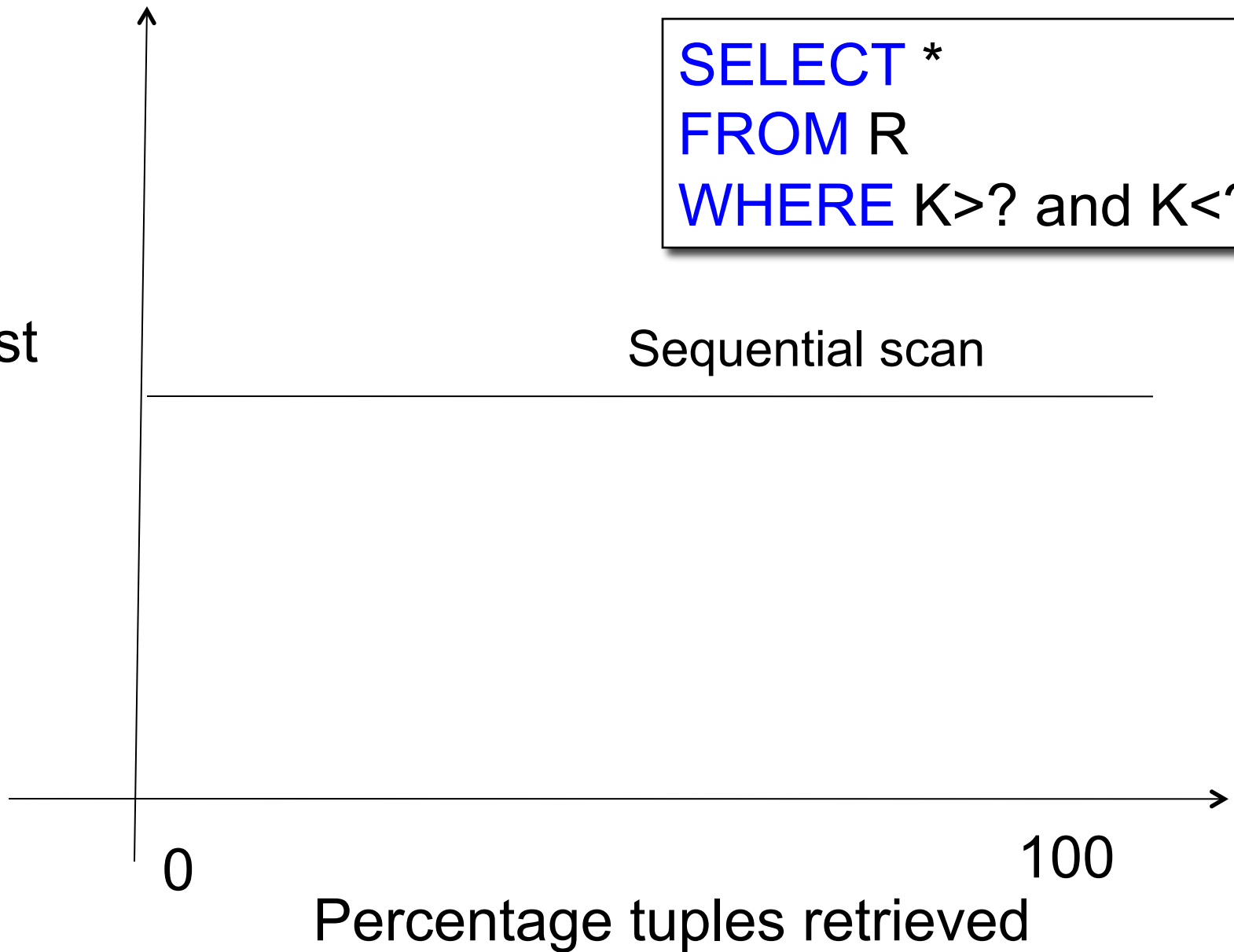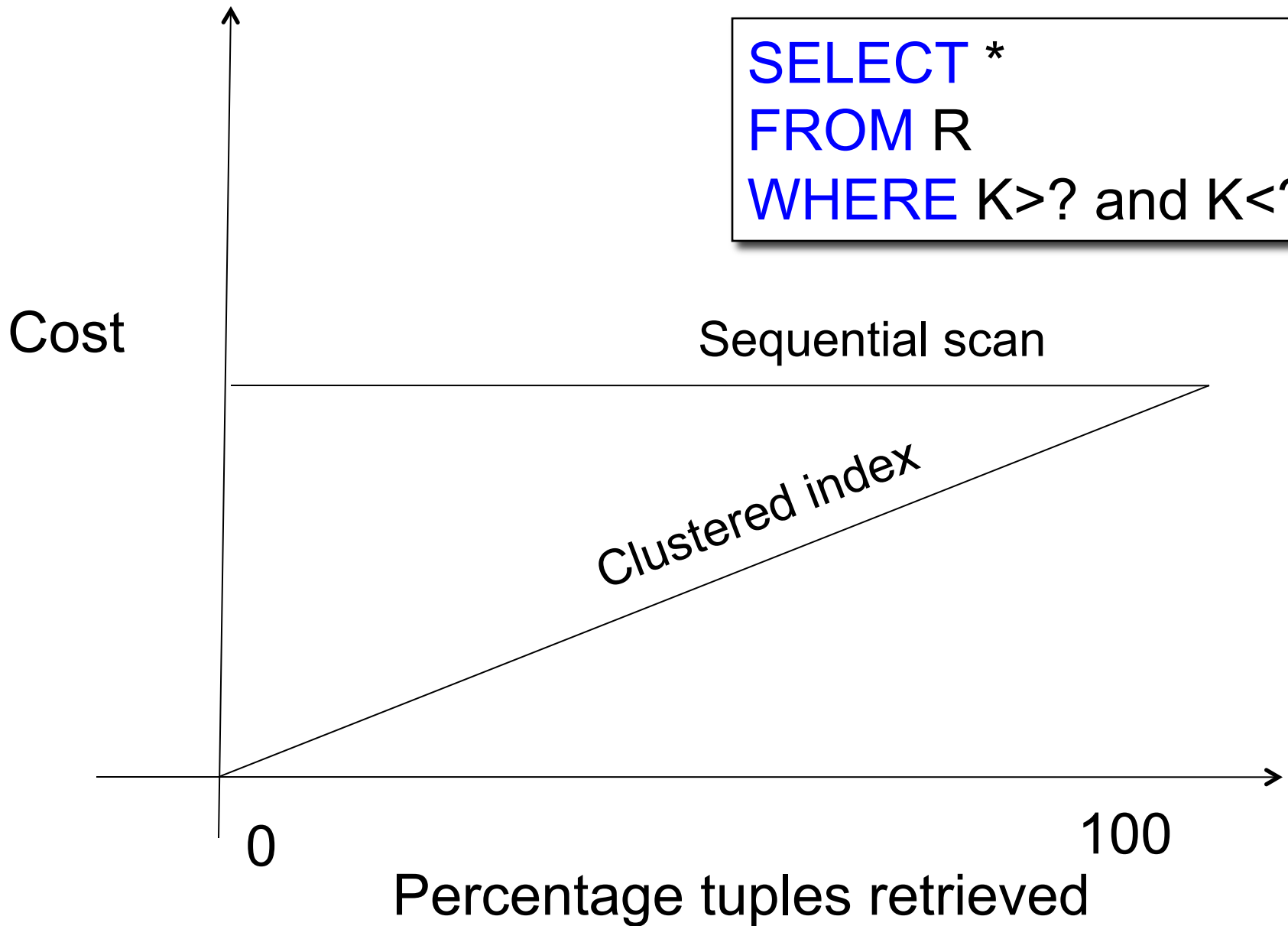WHERE K>? and K<?

0                                    100

Percentage tuples retrieved

Cost

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

Sequential scan

0

100

Percentage tuples retrieved

Cost

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

Sequential scan

Clustered index

0

100

Percentage tuples retrieved

Cost

Unclustered index

Sequential scan

Clustered index

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

0

100

Percentage tuples retrieved

# Today

- Cost of reading from disk

- Cost of single operators

- Cost of query plans

# Cost of Reading Data From Disk

# Cost Parameters

- Cost = I/O + CPU + Network BW
  - We will focus on I/O

- Parameters:
  - **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 < **T(R)**

- Where do these values come from?
  - DBMS collects statistics about data on disk

# Selectivity Factors for Conditions

- A = c                    /* $\sigma_{A=c}(R)$ */
  - Selectivity = $1/V(R,A)$

- A < c                    /* $\sigma_{A<c}(R)$ */
  - Selectivity = $(c - Low(R, A))/(High(R,A) - Low(R,A))$

- c1 < A < c2                    /* $\sigma_{c1<A<c2}(R)$ */
  - Selectivity = $(c2 - c1)/(High(R,A) - Low(R,A))$

# Cost of Reading Data From Disk

- Sequential scan for relation R costs **B(R)**

- Index-based selection
  - Estimate selectivity factor **X** (see previous slide)
  - Clustered index: **X*B(R)**
  - Unclustered index **X*T(R)**

Note: we ignore I/O cost for index pages

# Index Based Selection

- Example:

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

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:

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

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

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

- Index based selection:
  - If index is clustered:
  - If index is unclustered:

# Index Based Selection

- Example:

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

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

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

- Index based selection:
  - If index is clustered: B(R)/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: $B(R)/V(R,a) = 100$ I/Os
  - If index is unclustered: $T(R)/V(R,a) = 5{,}000$ 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)/V(R,a) = 100 I/Os
  - If index is unclustered: T(R)/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)

- One-pass algorithm when B(R) ≤ M

# Hash Join Example

Patient(pid, name, address)

Insurance(pid, provider, policy_nb)

Patient ⋈ 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 ⋈ Insurance

Some large-enough nb

Memory M = 21 pages

Showing pid only

Disk

Patient  Insurance

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

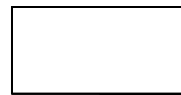| 2 | 4 |
| 4 | 3 |
| 2 | 8 |
| 8 | 9 |

| 6 | 6 |
| 1 | 3 |

This is one page with two tuples

# Hash Join Example

Step 1: Scan Patient and build hash table in memory

Can be done in method open()

Memory M = 21 pages

Hash h: pid % 5

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

Input buffer

Disk

Patient  Insurance

| 1 | 2 |

| 3 | 4 |

| 9 | 6 |

| 8 | 5 |

| 2 | 4 |

| 4 | 3 |

| 2 | 8 |

| 8 | 9 |

| 6 | 6 |

| 1 | 3 |

# 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 |

| 2 | 4 |
Input buffer

| 4 | 4 |
Output buffer

Disk

Patient Insurance

| 1 | 2 |

| 2 | 4 | | 6 | 6 |

| 3 | 4 |

| 4 | 3 | | 1 | 3 |

| 9 | 6 |

| 2 | 8 |

| 8 | 5 |

| 8 | 9 |

# 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 |

| 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 ⋈ 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 $t_1$ in r, $t_2$ in s
      if $t_1$ and $t_2$ join then output $(t_1, t_2)$
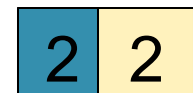
- Cost: $B(R) + B(R)B(S)$
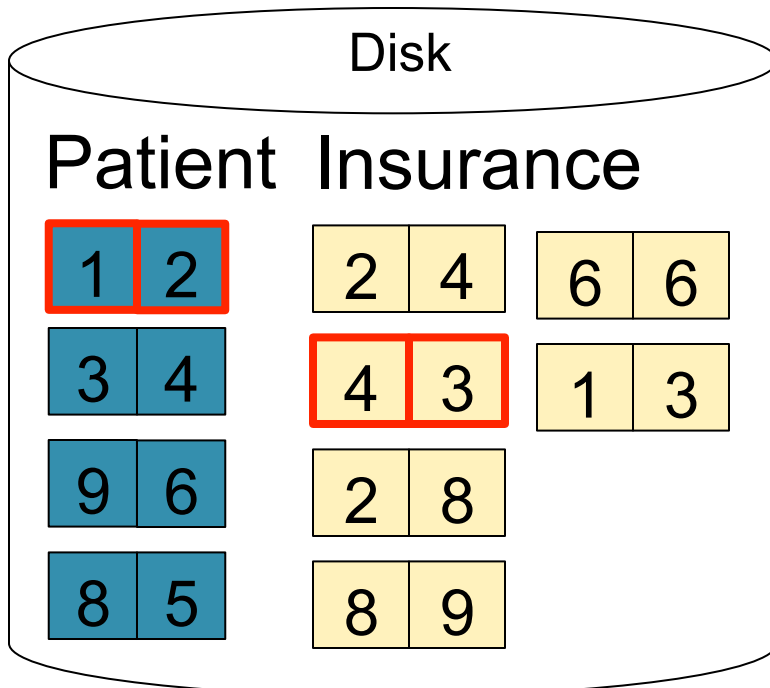
What is the Cost?

# Page-at-a-time Refinement

# Page-at-a-time Refinement

# Page-at-a-time Refinement



Disk

Patient    Insurance

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

| 2 | 4 |
| 4 | 3 |
| 2 | 8 |
| 8 | 9 |

| 6 | 6 |
| 1 | 3 |

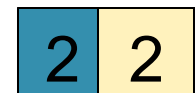| 1 | 2 |   Input buffer for Patient

| 2 | 8 |   Input buffer for Insurance

Keep going until read all of Insurance

| 2 | 2 |

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 $t_1$ in r, $t_2$ in s
         if $t_1$ and $t_2$ join then output $(t_1, t_2)$

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

What is the Cost?

# Sort-Merge Join

Sort-merge join:  R ⋈ S
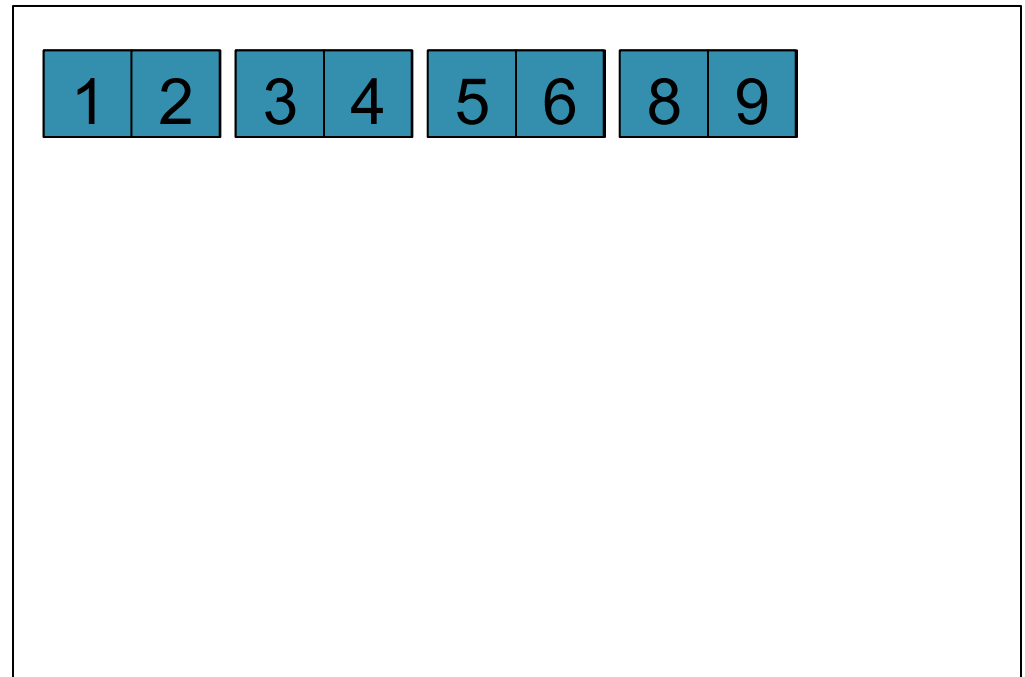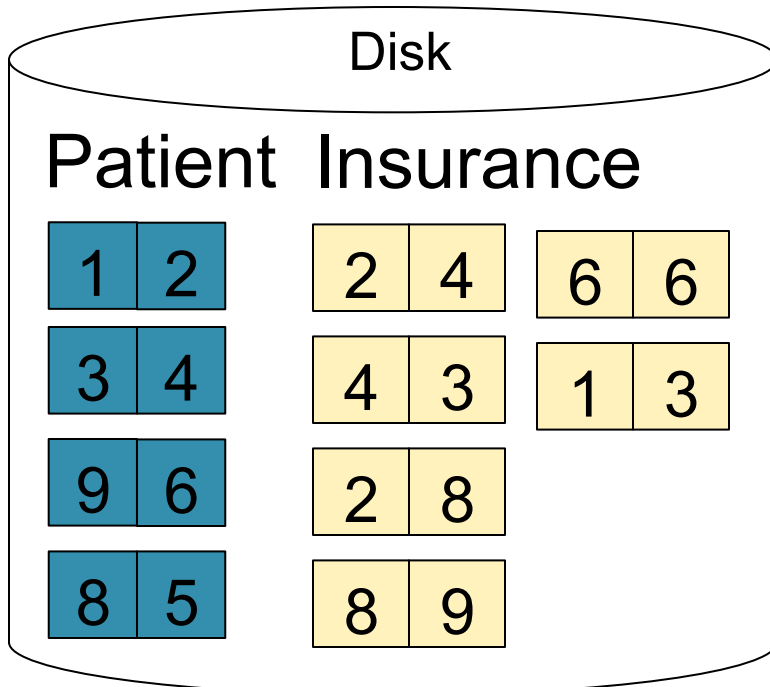
- Scan R and sort in main memory

- Scan S and sort in main memory

- Merge R and S


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

- One pass algorithm when B(S) + B(R) <= M

- Typically, this is NOT a one pass algorithm

# Sort-Merge Join Example

Step 1: Scan Patient and <span style="color:red">sort</span> in memory
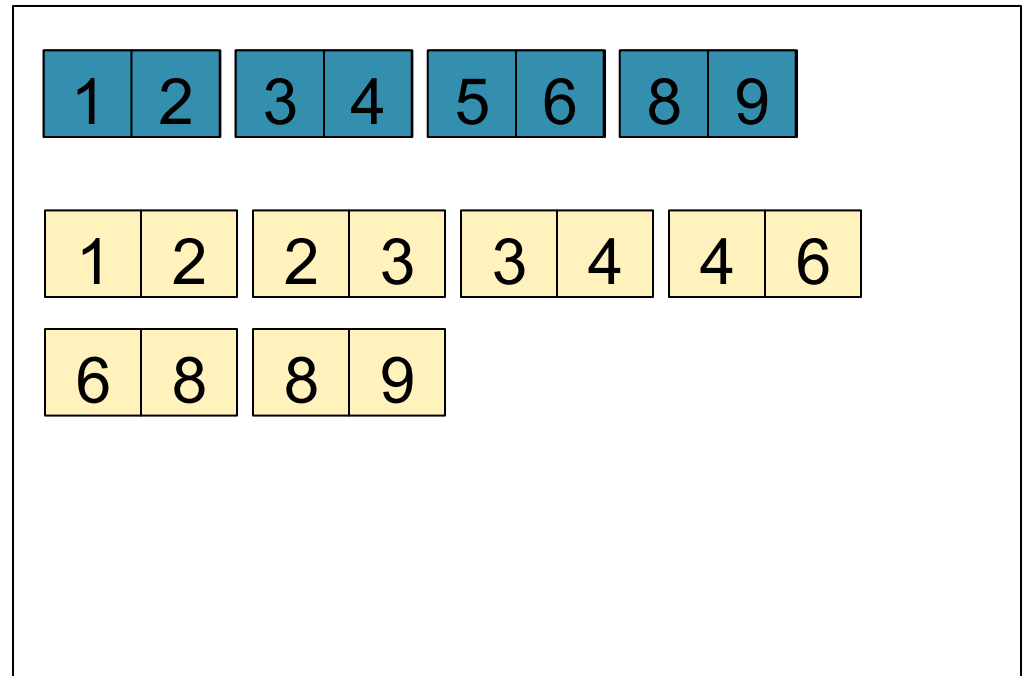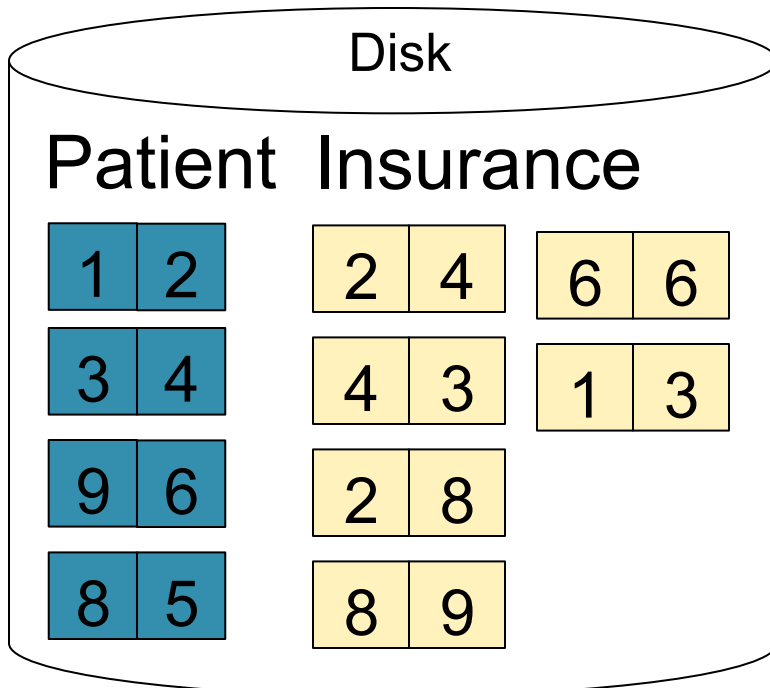
Memory M = 21 pages

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

Disk

## Patient  Insurance

| 1 | 2 |

| 3 | 4 |

| 9 | 6 |

| 8 | 5 |

| 2 | 4 |  | 6 | 6 |

| 4 | 3 |  | 1 | 3 |

| 2 | 8 |

| 8 | 9 |

41

# Sort-Merge Join Example

Step 2: Scan Insurance and sort in memory

Memory M = 21 pages

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

| 1 | 2 | 2 | 3 | 3 | 4 | 4 | 6 |

| 6 | 8 | 8 | 9 |

Disk

## Patient   Insurance

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

# Sort-Merge Join Example

Step 3: Merge Patient and Insurance
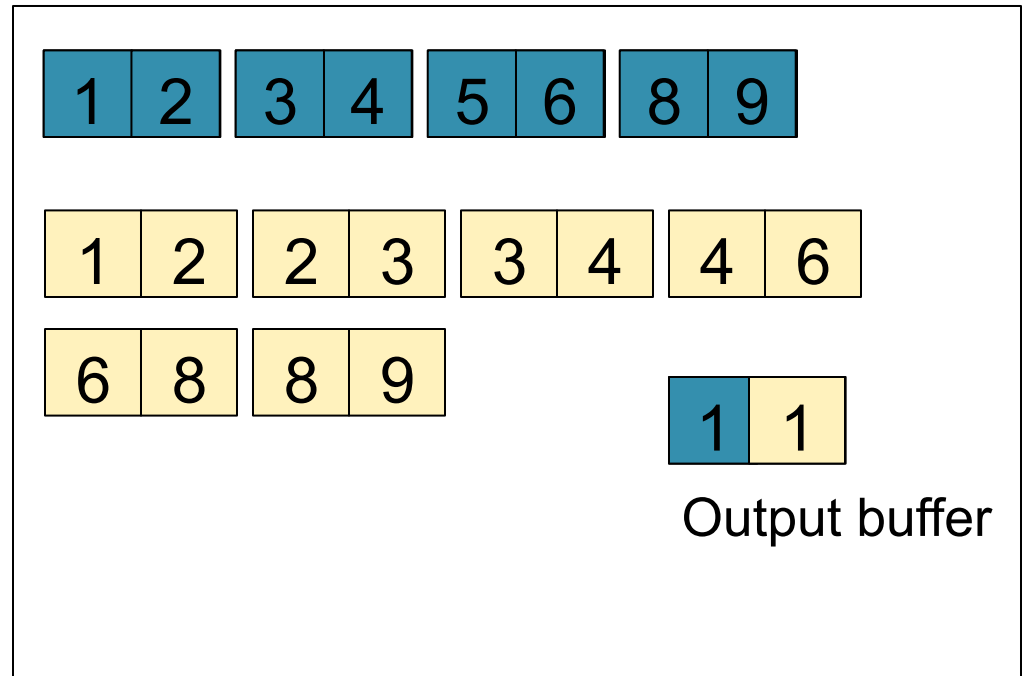
# Sort-Merge Join Example

Step 3: Merge Patient and Insurance

Memory M = 21 pages

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

| 1 | 2 | 2 | 3 | 3 | 4 | 4 | 6 |

| 6 | 8 | 8 | 9 |

| 2 | 2 |

Output buffer

Keep going until end of first relation

Disk

Patient    Insurance

| 1 | 2 |    | 2 | 4 |    | 6 | 6 |

| 3 | 4 |    | 4 | 3 |    | 1 | 3 |

| 9 | 6 |    | 2 | 8 |

| 8 | 5 |    | 8 | 9 |

# Index Nested Loop Join

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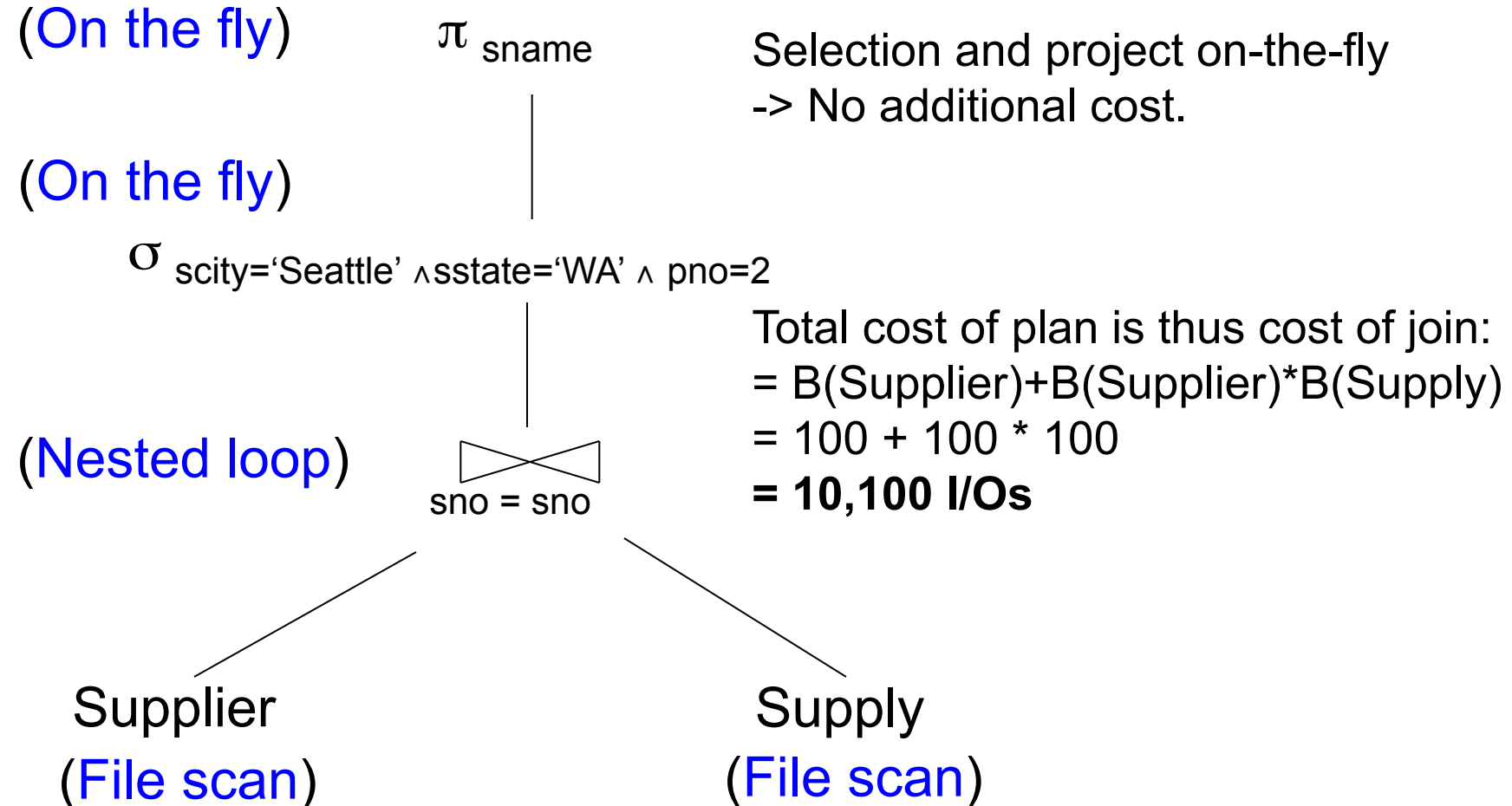  - If index on S is unclustered: $B(R) + T(R)T(S)/V(S,a)$

# Cost of Query Plans

# Physical Query Plan 1

(On the fly)          $\pi$ $_{sname}$          Selection and project on-the-fly
                                                -> No additional cost.

(On the fly)

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

                                                Total cost of plan is thus cost of join:
                                                = B(Supplier)+B(Supplier)*B(Supply)
                                                = 100 + 100 * 100
(Nested loop)          $\bowtie$                 **= 10,100 I/Os**
                       sno = sno

Supplier                              Supply
(File scan)                           (File scan)

CSE 344 - Winter 2016                                              47

# Physical Query Plan 2

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

(Sort-merge join)     $\bowtie$     (c)

sno = sno

(Scan
write to T1)

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

(Scan
write to T2)

Supplier
(File scan)

Supply
(File scan)

Total cost
= 100 + 100 * 1/20 * 1/10 (a)
+ 100 + 100 * 1/2500 (b)
+ 2 (c)
+ 0 (d)
Total cost $\approx$ **204 I/Os**

T(Supplier) = 1000          B(Supplier) = 100          V(Supplier,scity) = 20          M = 11
T(Supply) = 10,000          B(Supply) = 100          V(Supplier,state) = 10
                                                      V(Supply,pno) = 2,500

# Physical Query Plan 3

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

(On the fly)

(c)   $\sigma_{\text{scity='Seattle' } \wedge \text{sstate='WA'}}$

(b)   $\bowtie$
      sno = sno   (Index nested loop)

(Use hash index)   4 tuples

(a) $\sigma_{\text{pno=2}}$

Supply

(Index on pno )
Assume: clustered

Supplier

(Index on sno)
Clustering does not matter

Total cost
= 1 (a)
+ 4 (b)
+ 0 (c)
+ 0 (d)
Total cost $\approx$ **5 I/Os**

# Query Optimizer Overview

- **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
    - Compute number of I/Os
    - Optionally take into account other resources
  - Choose plan with lowest cost
  - This is called cost-based optimization