

Introduction to Database Systems

CSE 344

Lecture 11: Basics of Query Optimization and Query Cost Estimation

Announcements

- HW3 Azure V12 upgrade
- Section attendance

Review

- What is a disk block? (aka page)
- What is an index?
 - What data structures are used to represent indexes in memory?
- 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.
How can the optimizer use it?

INDEX I1 on V(M)

INDEX I2 on V(M,P)

INDEX I3 on V(P,M)

Two typical kinds of queries

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

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

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

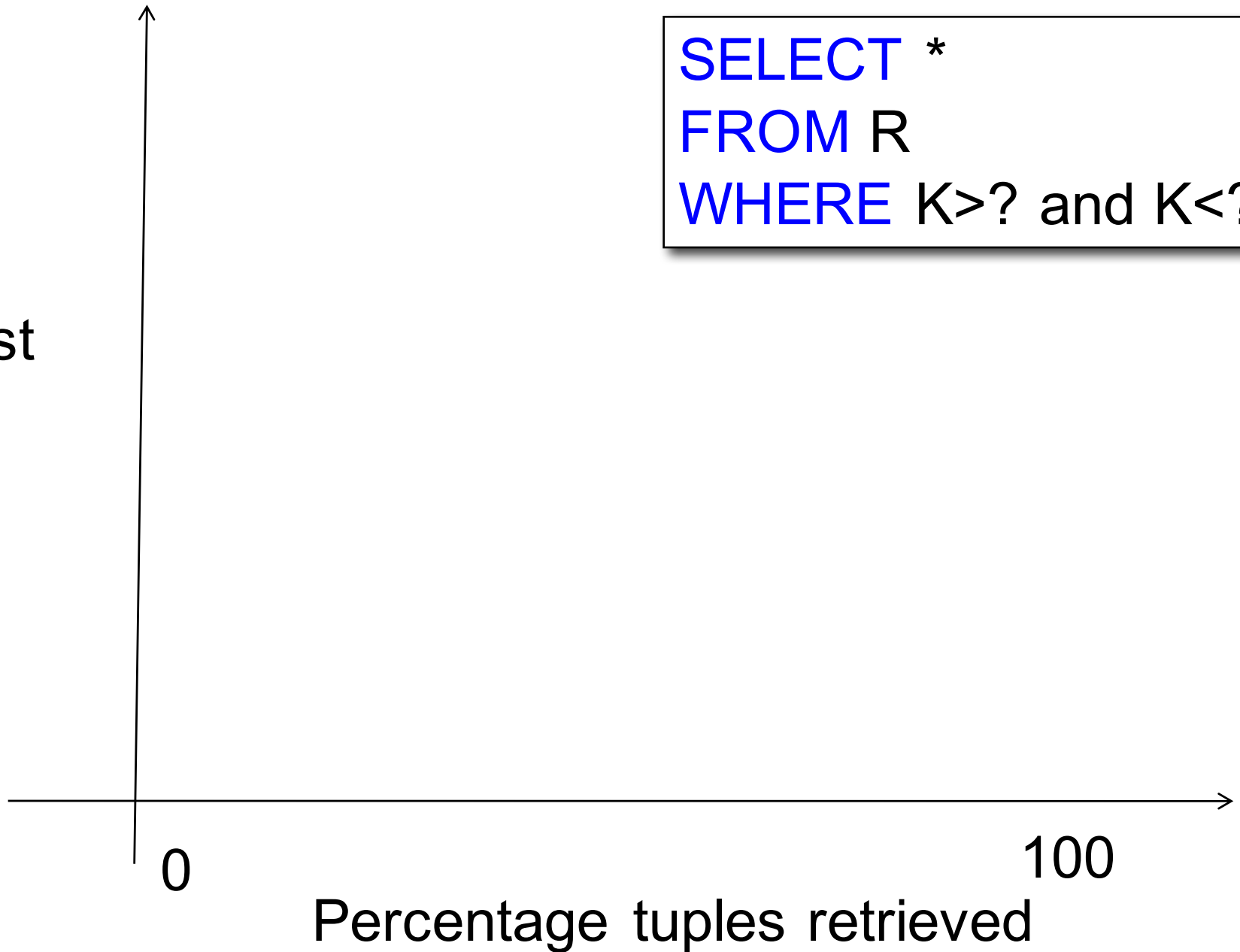
- Range queries
- What data structure should be used for index?

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
- Consider how each query will be processed
 - Which predicate will be processed first?
- Try to choose indexes that speed-up multiple queries

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

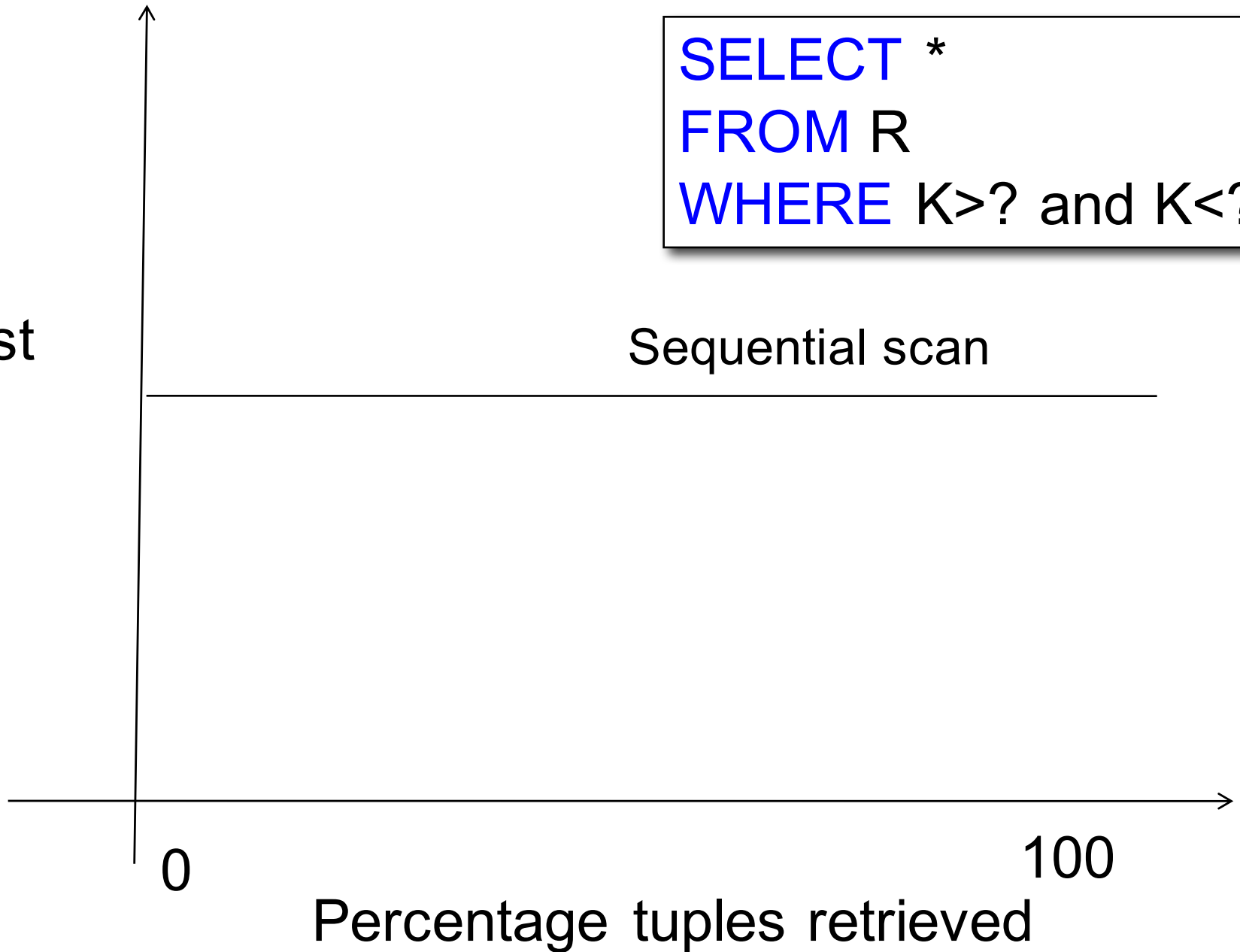
Cost



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

Cost

Sequential scan

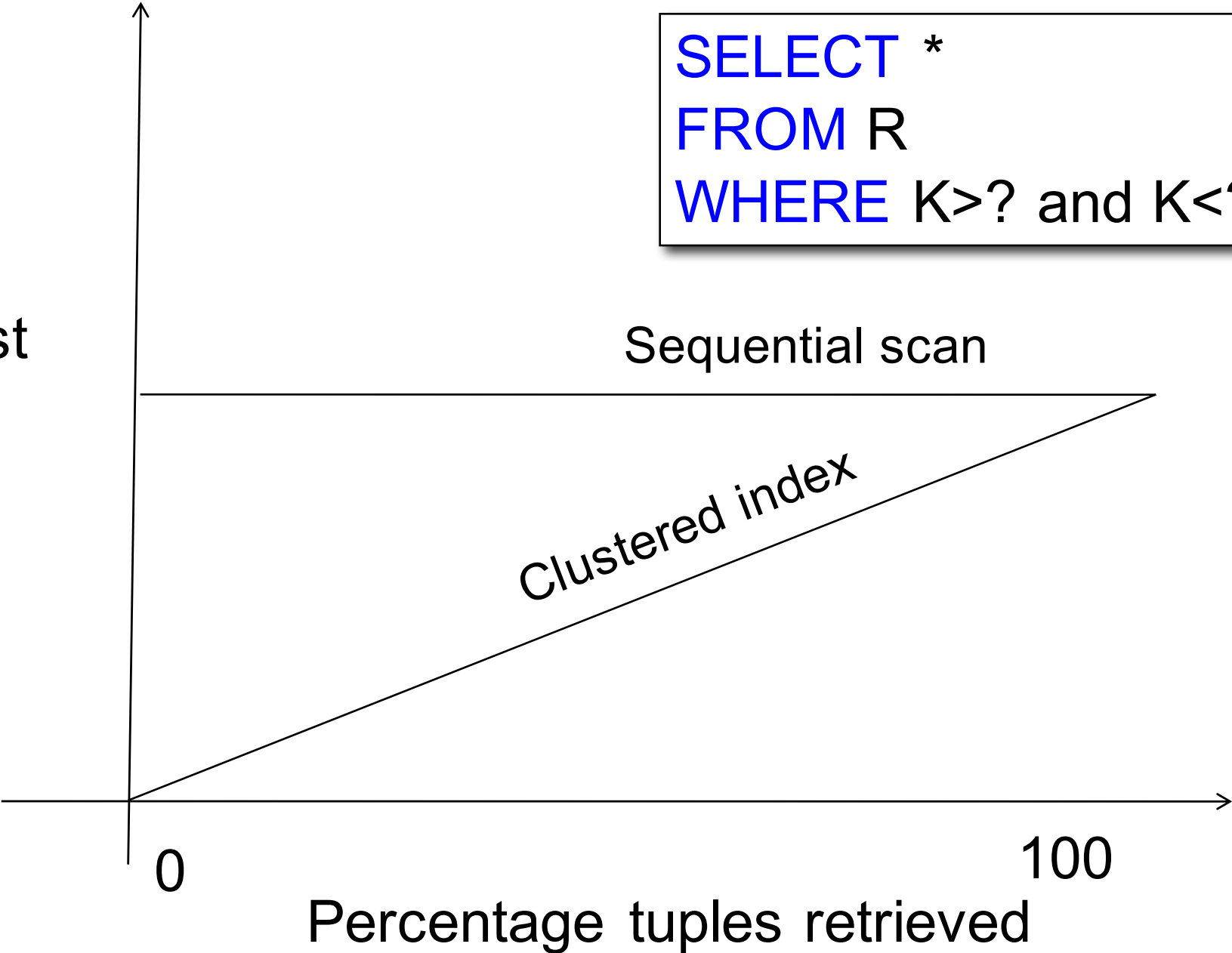



```
SELECT *  
FROM R  
WHERE K > ? and K < ?
```

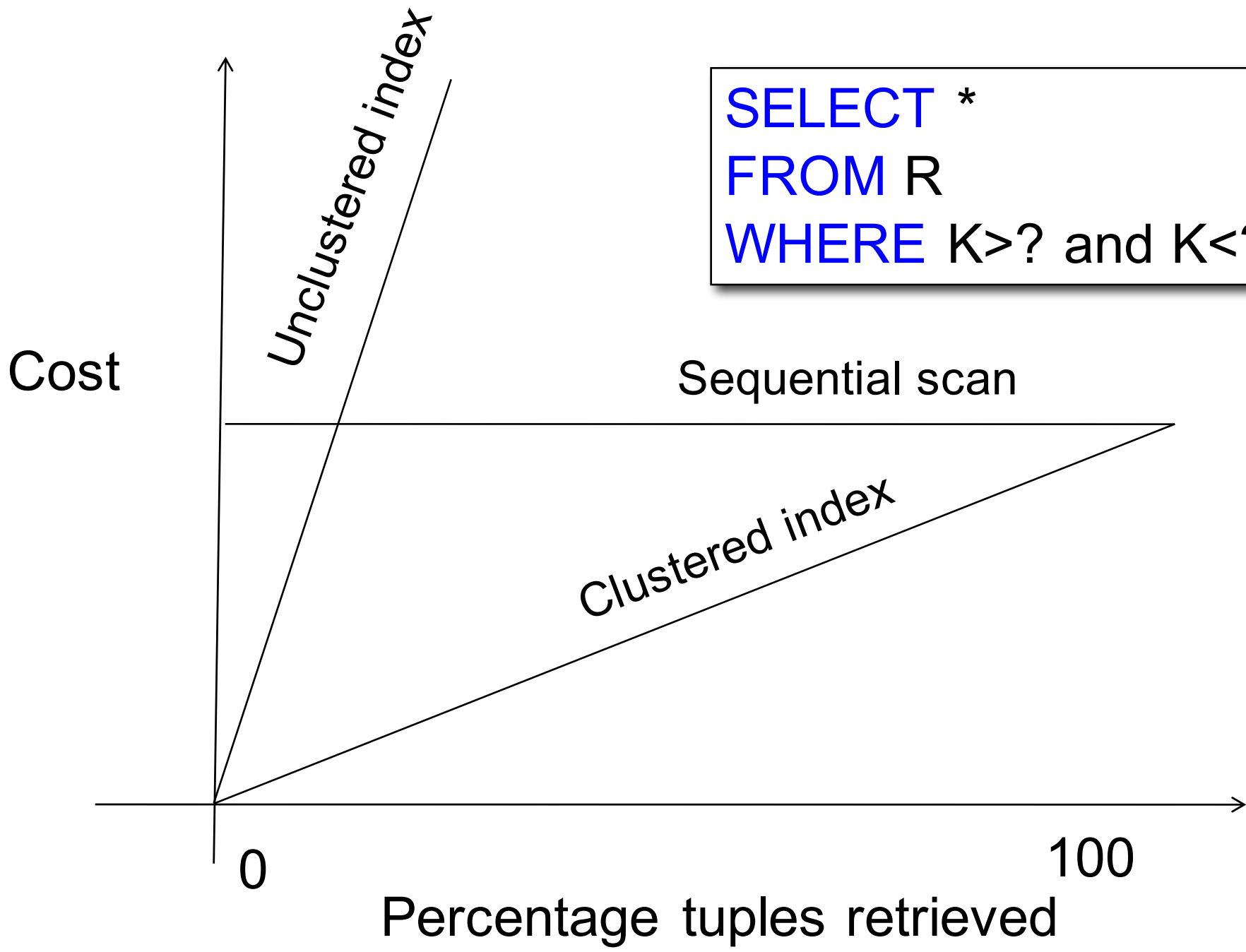
Cost

Sequential scan

Clustered index



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



Cost Models

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

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

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

- 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:

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

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

- 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:

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?

- 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 ⋈ Insurance

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

Two tuples
per page

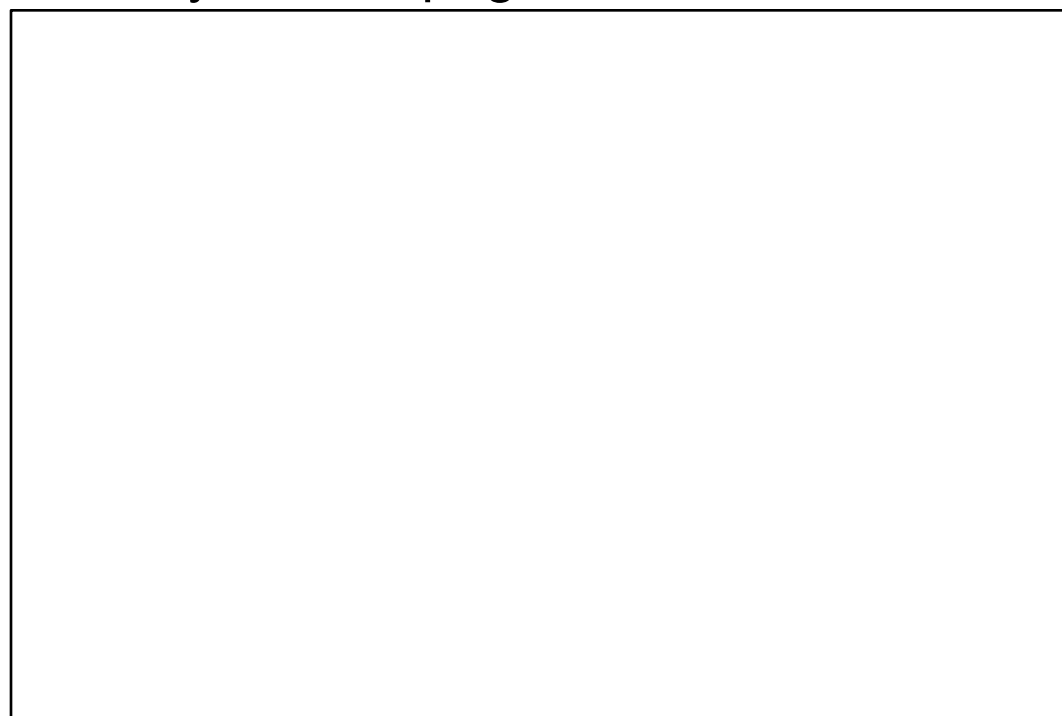
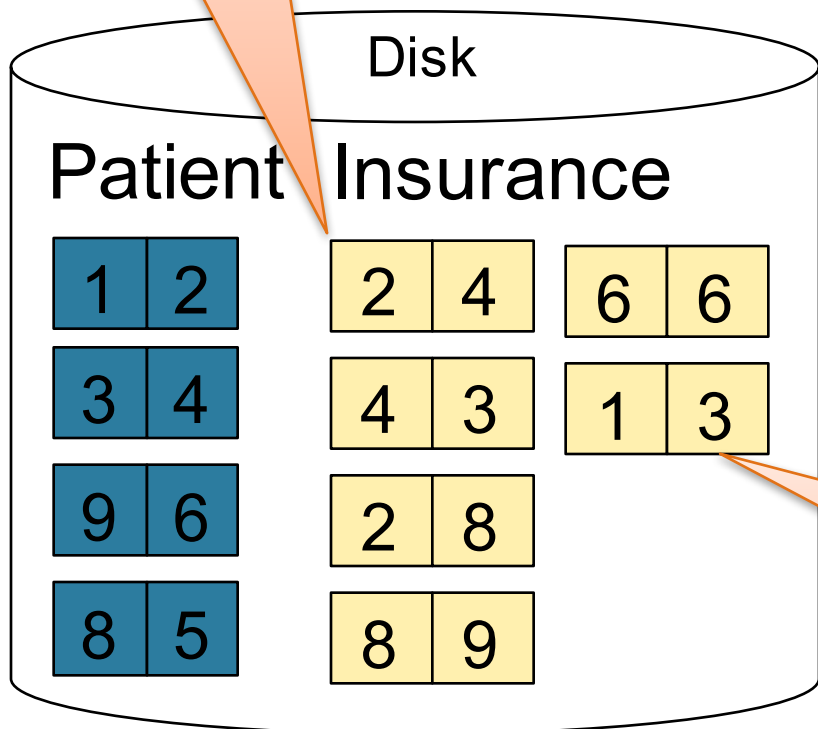
Hash Join Example

Patient \bowtie Insurance

Some large-enough #

Memory M = 21 pages

Showing pid only

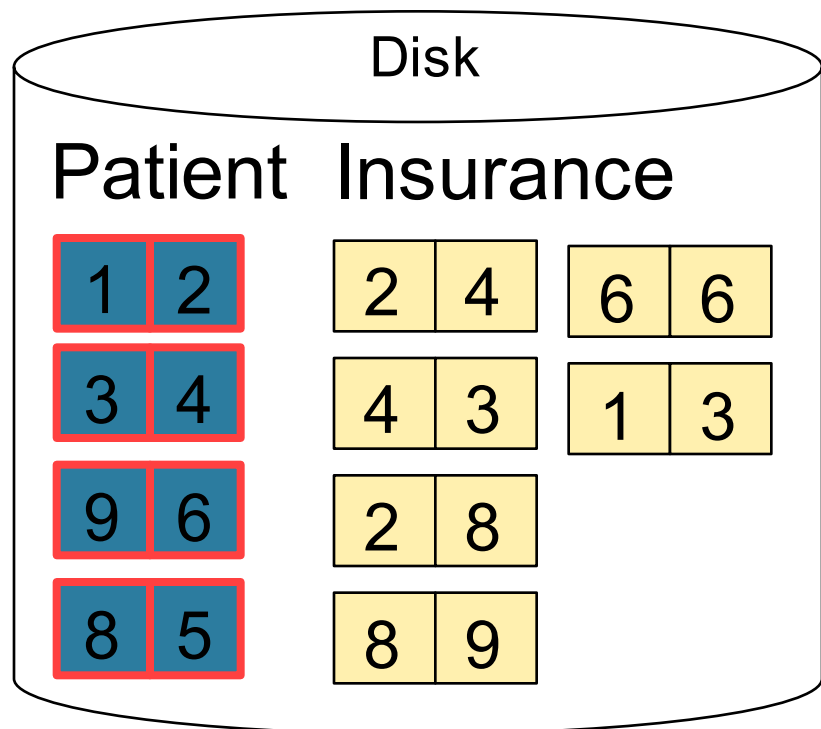
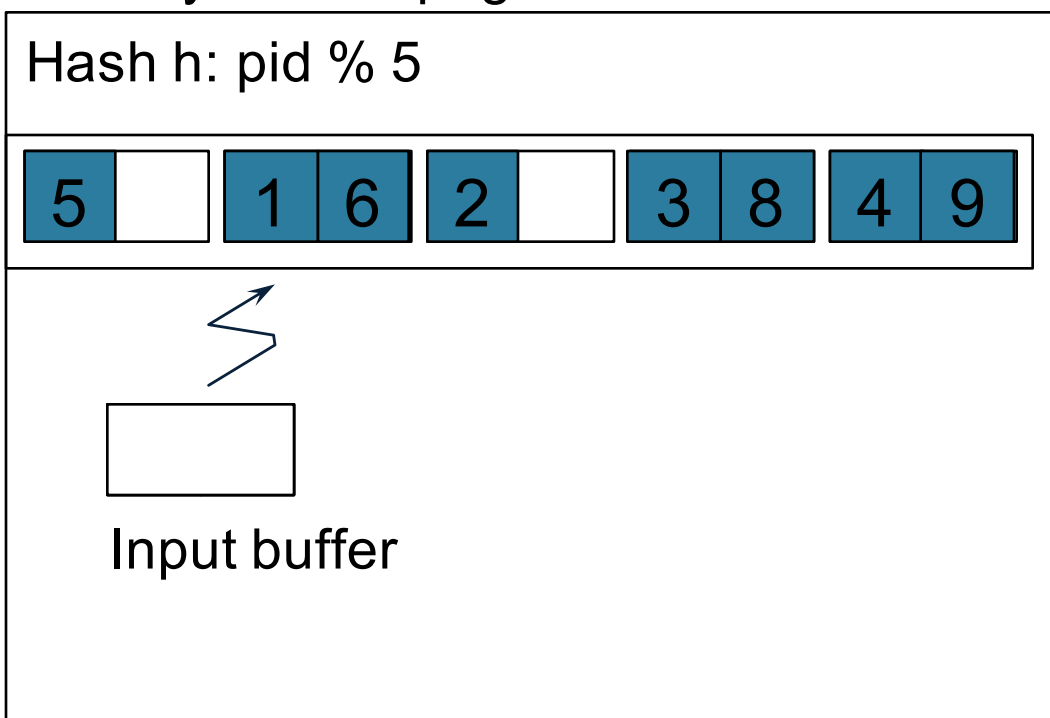


This is one page with two tuples

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

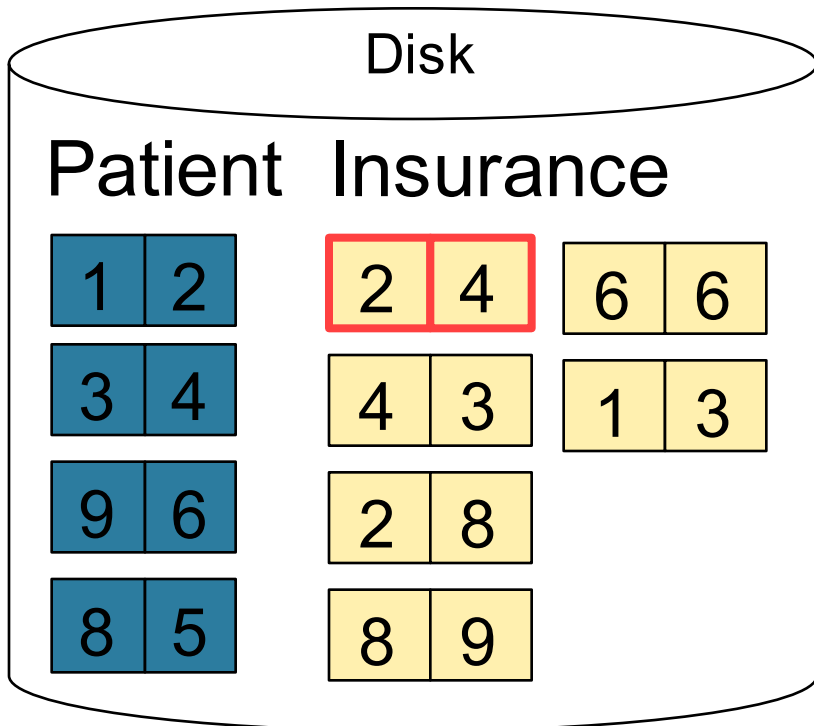
2	4
---	---

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

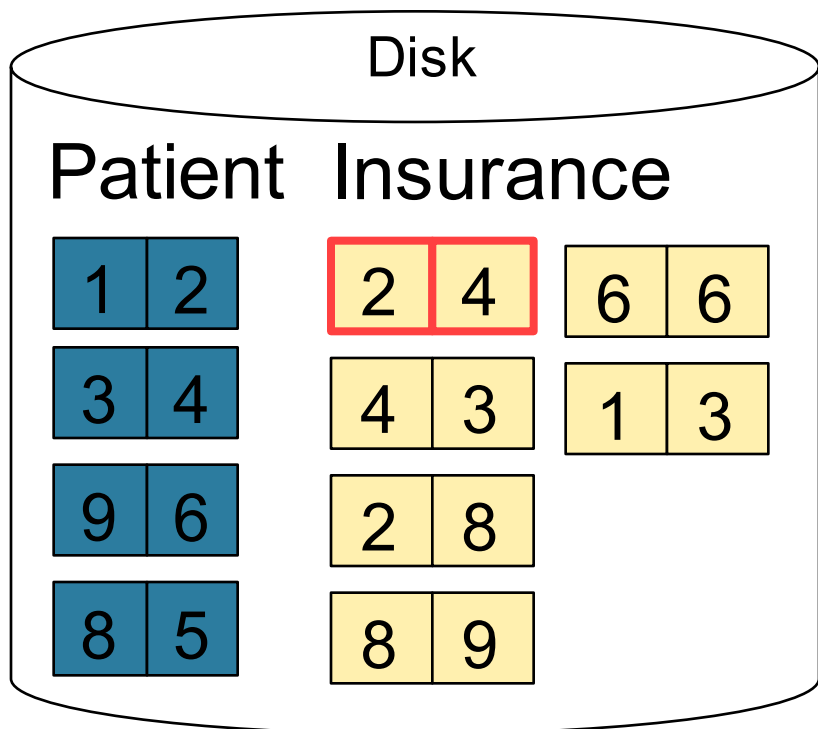
5		1	6	2		3	8	4	9
---	--	---	---	---	--	---	---	---	---

2	4
---	---

Input buffer

4	4
---	---

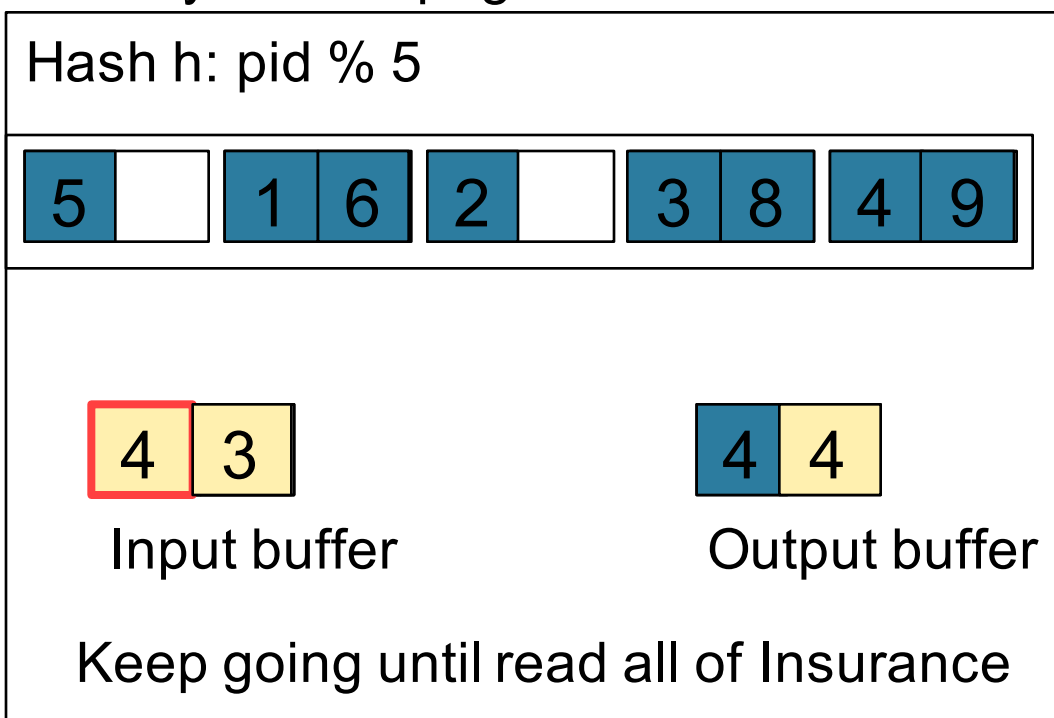
Output buffer



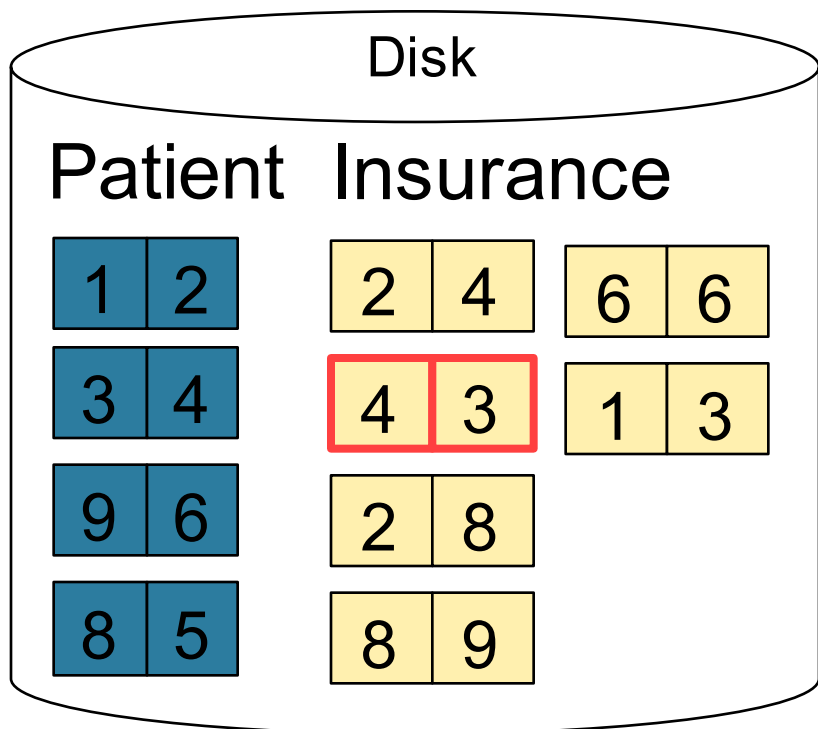
Hash Join Example

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

Memory M = 21 pages



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

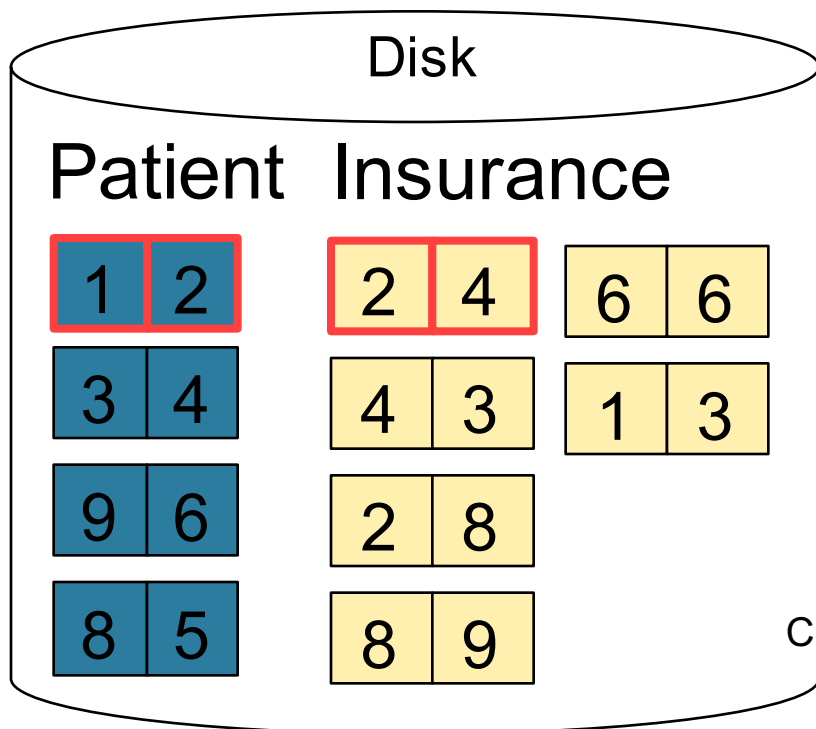
Page-at-a-time Refinement

```
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)
```

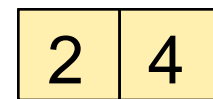
- Cost: $B(R) + B(R)B(S)$

What is the **Cost**?

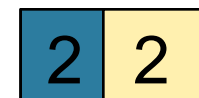
Page-at-a-time Refinement



Input buffer for Patient

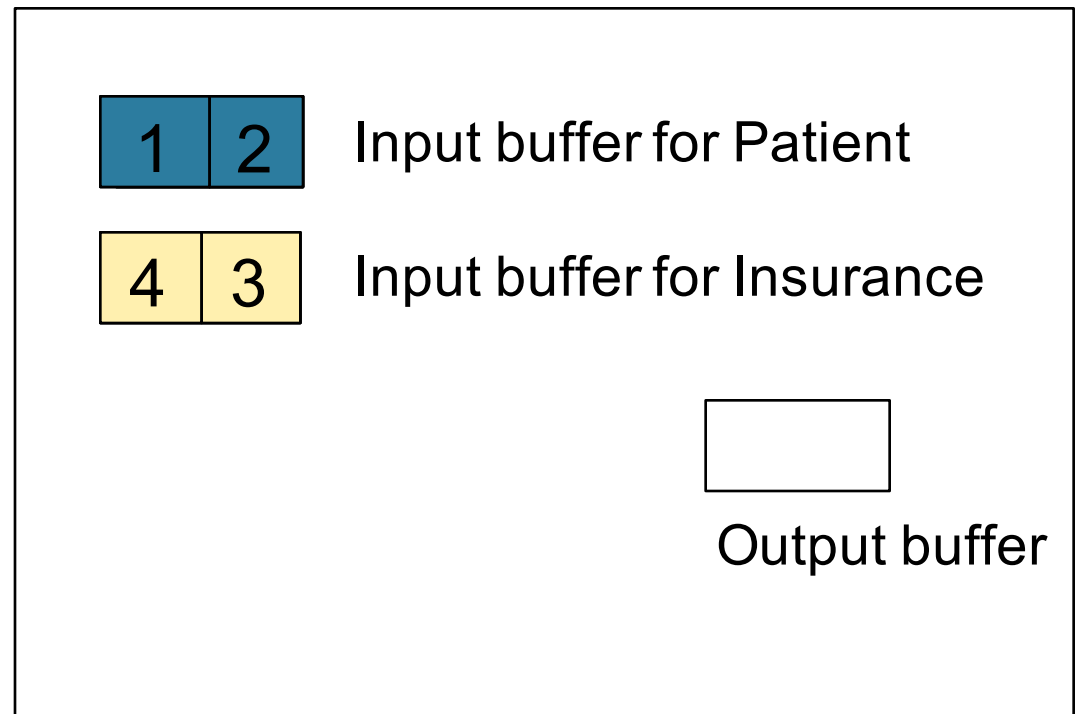
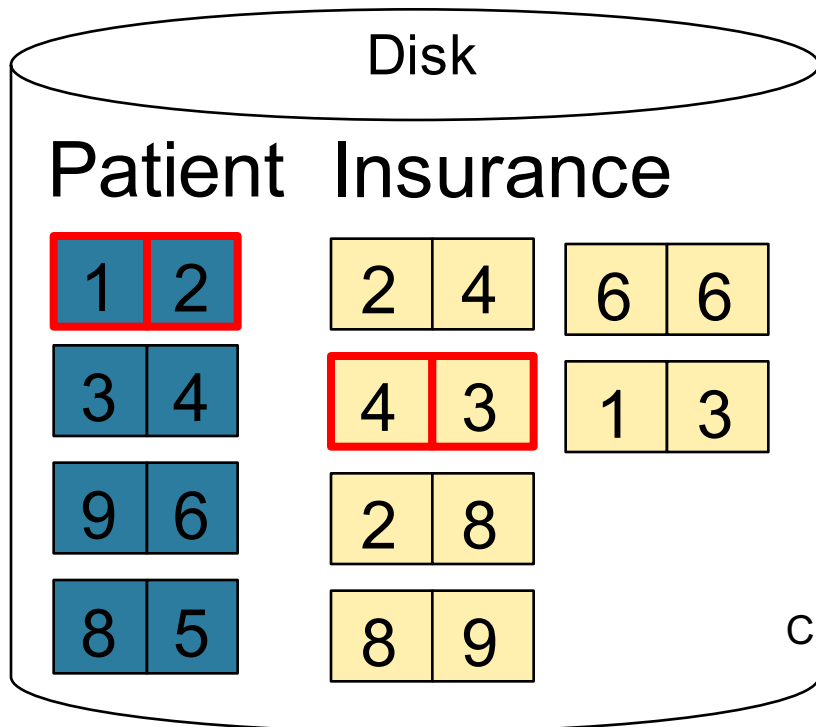


Input buffer for Insurance

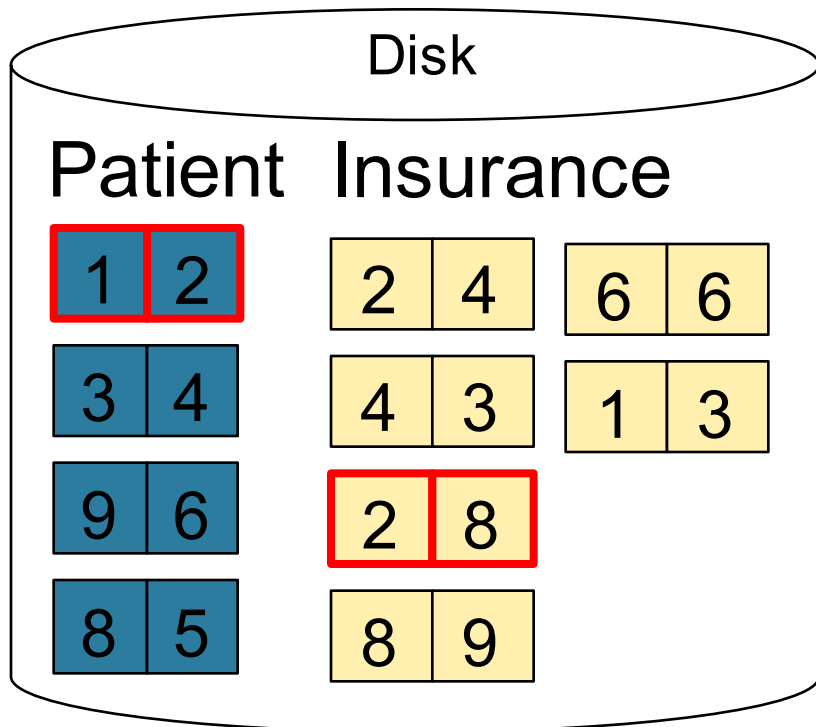


Output buffer

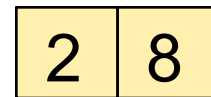
Page-at-a-time Refinement



Page-at-a-time Refinement

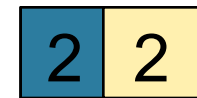


Input buffer for Patient



Input buffer for Insurance

Keep going until read all of Insurance



Output buffer

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

$$\text{Cost: } B(R) + B(R)B(S)$$

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**?