

Introduction to Data Management CSE 414

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**

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

Student

| 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

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

CSE 414 - Autumn 2018

7

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

CSE 414 - Autumn 2018

8

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



CSE 414 - Autumn 2018



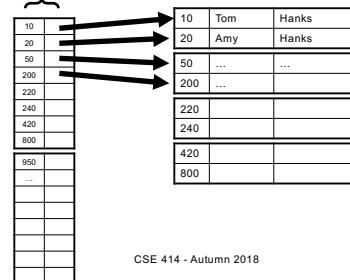
Example 1: Index on ID

Student

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

Index Student_ID on Student.ID

Data File Student



CSE 414 - Autumn 2018

10

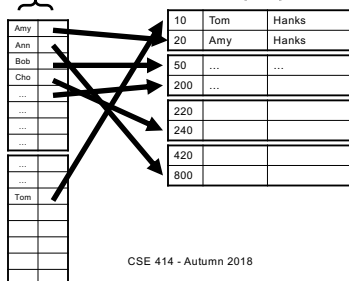
Example 2: Index on fName

Index Student_fName
on Student.fName

Data File Student

Student

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



CSE 414 - Autumn 2018

11

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

CSE 414 - Autumn 2018

12

Hash table example

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

Index Student_ID on Student.ID Data File Student

Index File (preferably in memory)

Data file (on disk)

CSE 414 - Autumn 2018 13

B+ Tree Index by Example

d = 2

Find the key 40

CSE 414 - Autumn 2018 14

Clustered vs Unclustered

CLUSTERED

UNCLUSTERED

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

CSE 414 - Autumn 2018 15

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

CSE 414 - Autumn 2018 16

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

CSE 414 - Autumn 2018 17

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

CSE 414 - Autumn 2018 18

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 = sequential scan entire file;
 - 1-2% decreases over time, because of increased density
- Solid state (SSD): **still, too expensive today**



CSE 414 - Autumn 2018

19

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

CSE 414 - Autumn 2018

20

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



Example

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

CSE 414 - Autumn 2018

21

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



Example

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

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

CSE 414 - Autumn 2018

22

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



Example

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

```
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

CSE 414 - Autumn 2018

23

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



Example

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

```
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
```

CSE 414 - Autumn 2018

24

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'

CSE 414 - Autumn 2018 25

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'

Index join

for x' in Student_ID where x'.ID = y.studentID
x = fetch the Student record pointed to by x'

CSE 414 - Autumn 2018 26

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'

Index join

for x' in Student_ID where x'.ID = y.studentID
x = fetch the Student record pointed to by x'

CSE 414 - Autumn 2018 27

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'

Index join

for x' in Student_ID where x'.ID = y.studentID
x = fetch the Student record pointed to by x'

CSE 414 - Autumn 2018 28

**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)
```

CSE 414 - Autumn 2018 29

**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)
```

CSE 414 - Autumn 2018 30

Getting Practical: Creating Indexes in SQL

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

yes

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

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

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

CSE 414 - Autumn 2018 31

Getting Practical: Creating Indexes in SQL

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

yes

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

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

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

select * from V where P=55

CSE 414 - Autumn 2018 32

Getting Practical: Creating Indexes in SQL

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

yes

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

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

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

CSE 414 - Autumn 2018 33

Getting Practical: Creating Indexes in SQL

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

yes

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

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

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

select * from V where P=55

yes

select * from V where M=77

CSE 414 - Autumn 2018 34

Getting Practical: Creating Indexes in SQL

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

yes

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

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

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

select * from V where M=77

no

CSE 414 - Autumn 2018 35

Getting Practical: Creating Indexes in SQL

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

yes

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

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

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

select * from V where P=55

yes

select * from V where M=77

no

Not supported in SQLite

CSE 414 - Autumn 2018 36

Which Indexes?

| Student | | |
|---------|-------|-------|
| ID | fName | IName |
| 10 | Tom | Hanks |
| 20 | Amy | Hanks |
| ... | | |

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

CSE 414 - Autumn 2018

37

Which Indexes?

| Student | | |
|---------|-------|-------|
| ID | fName | IName |
| 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

CSE 414 - Autumn 2018

38

Which Indexes?

| Student | | |
|---------|-------|-------|
| ID | fName | IName |
| 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

CSE 414 - Autumn 2018

39

Which Indexes?

| Student | | |
|---------|-------|-------|
| ID | fName | IName |
| 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



CSE 414 - Autumn 2018

40

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

CSE 414 - Autumn 2018

41

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=?
```

CSE 414 - Autumn 2018

42

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 ?

CSE 414 - Autumn 2018

43

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)

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 ?

CSE 414 - Autumn 2018

45

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)

CSE 414 - Autumn 2018

46

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 ?

CSE 414 - Autumn 2018

47

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)

CSE 414

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 ?

CSE 414 - Autumn 2018

49

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

CSE 414 - Autumn 2018

50

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?

CSE 414 - Autumn 2018

51

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

CSE 414 - Autumn 2018

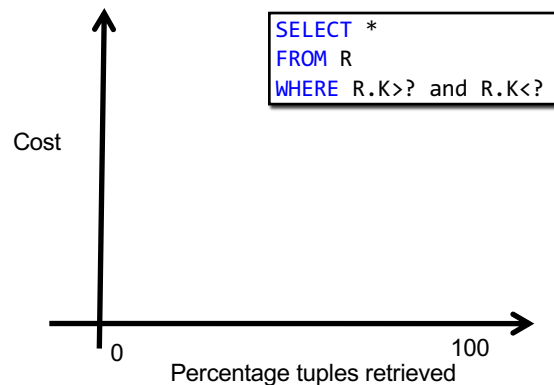
52

To Cluster or Not

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

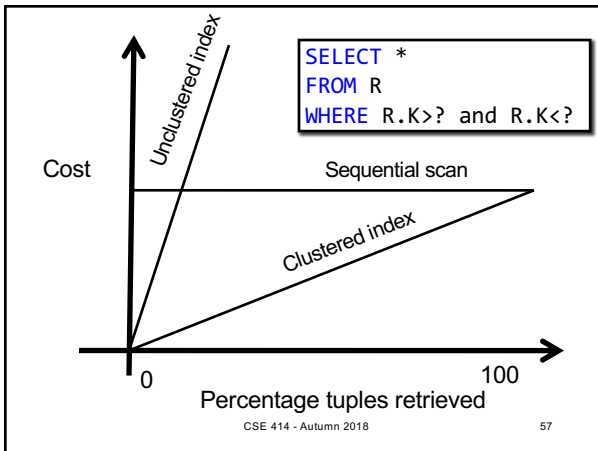
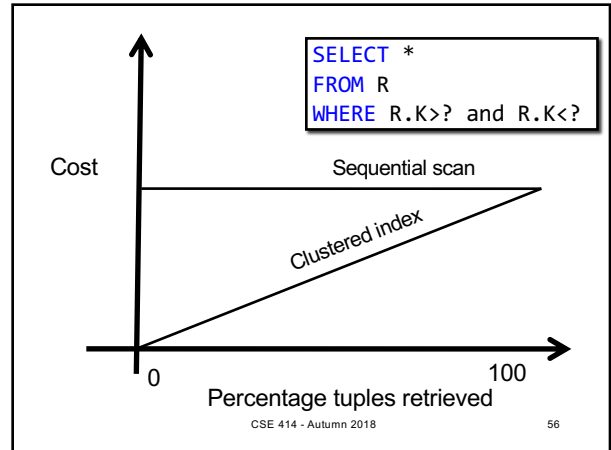
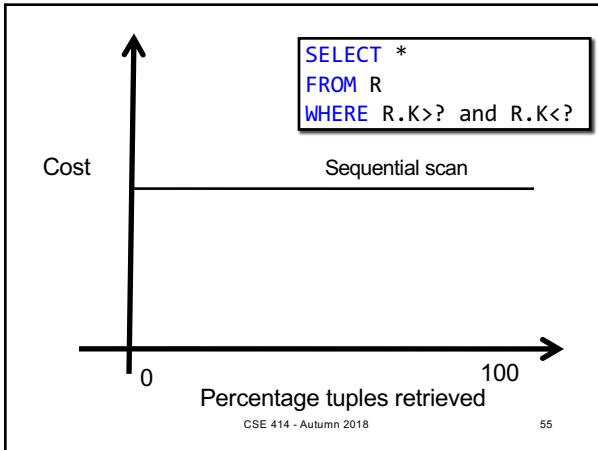
CSE 414 - Autumn 2018

53



CSE 414 - Autumn 2018

54



Introduction to Database Systems
CSE 344

Lecture 17:
Basics of Query Optimization and
Query Cost Estimation

CSE 414 - Autumn 2018 58

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

CSE 414 - Autumn 2018 59

Cost of Reading
Data From Disk

CSE 414 - Autumn 2018 60

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

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

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: $f * B(R)$
 - Unclustered index $f * T(R)$

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

CSE 414 - Autumn 2018

73

Join Algorithms

- Hash join
- Nested loop join
- Sort-merge join

CSE 414 - Autumn 2018

74

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?

CSE 414 - Autumn 2018

75

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

CSE 414 - Autumn 2018

76

Hash Join Example

Patient(pid, name, address)
Insurance(pid, provider, policy_nb)
Patient \bowtie Insurance

| Patient | | | Insurance | |
|---------|--------|-----------|-----------|------------|
| 1 | 'Bob' | 'Seattle' | 2 | 'Blue' 123 |
| 2 | 'Ela' | 'Everett' | 4 | 'Prem' 432 |
| 3 | 'Jill' | 'Kent' | 4 | 'Prem' 343 |
| 4 | 'Joe' | 'Seattle' | 3 | 'GrpH' 554 |

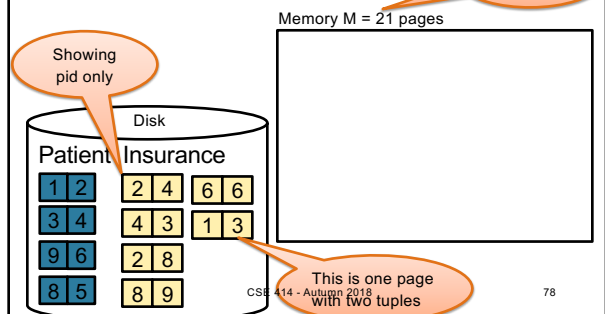
Two tuples per page

CSE 414 - Autumn 2018

77

Hash Join Example

Patient \bowtie Insurance



CSE 414 - Autumn 2018

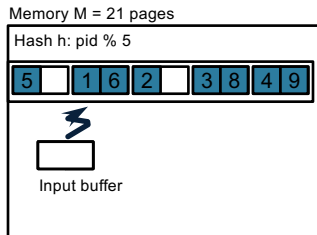
78

Hash Join Example

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

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |



CSE 414 - Autumn 2018

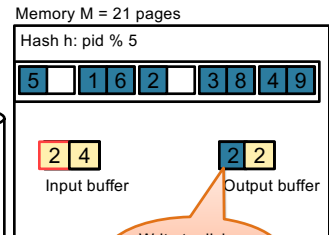
79

Hash Join Example

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

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |



Write to disk or pass to next operator

CSE 414 - Autumn 2018

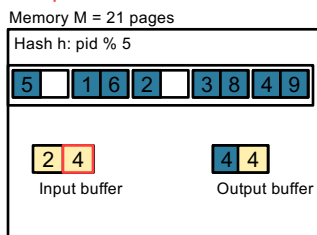
80

Hash Join Example

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

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |



CSE 414 - Autumn 2018

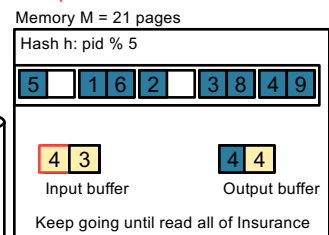
81

Hash Join Example

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

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |



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

CSE 414 - Autumn 2018

82

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?

CSE 414 - Autumn 2018

83

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?

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

CSE 414 - Autumn 2018

84

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?

CSE 414 - Autumn 2018 85

Page-at-a-time Refinement

CSE 414 - Autumn 2018 86

Page-at-a-time Refinement

CSE 414 - Autumn 2018 87

Page-at-a-time Refinement

Keep going until read all of Insurance
 Then repeat for next page of Patient... until end of Patient

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

CSE 414 - Autumn 2018 88

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?

CSE 414 - Autumn 2018 89

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

CSE 414 - Autumn 2018 90

Sort-Merge Join Example

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

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |

Memory M = 21 pages

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

CSE 414 - Autumn 201891

Sort-Merge Join Example

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

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |

Memory M = 21 pages

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

CSE 414 - Autumn 201892

Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |

Memory M = 21 pages

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

Output buffer

CSE 414 - Autumn 201893

Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance

Disk

| Patient | Insurance |
|---------|-----------|
| 1 2 | 2 4 6 6 |
| 3 4 | 4 3 1 3 |
| 9 6 | 2 8 |
| 8 5 | 8 9 |

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

CSE 414 - Autumn 201894

Index Nested Loop Join

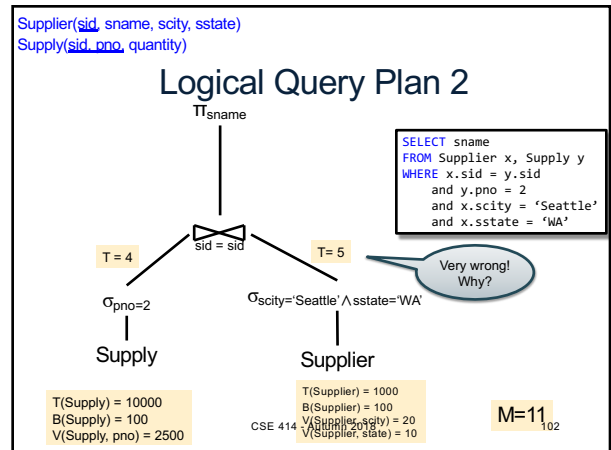
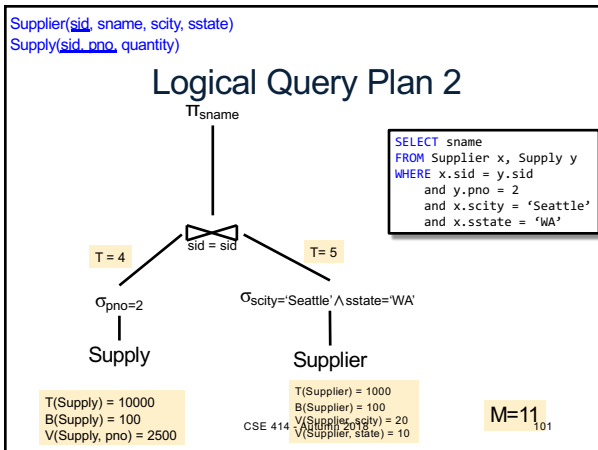
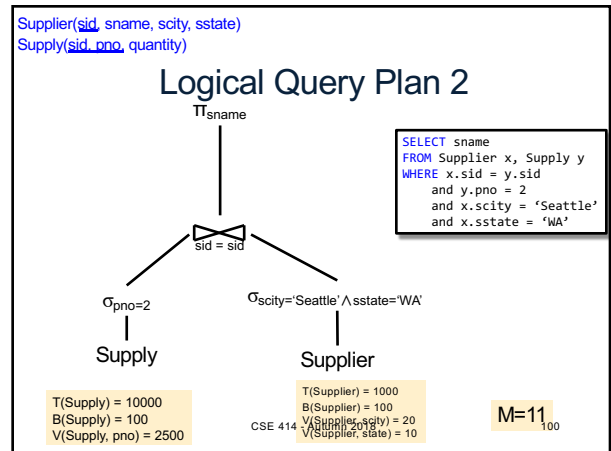
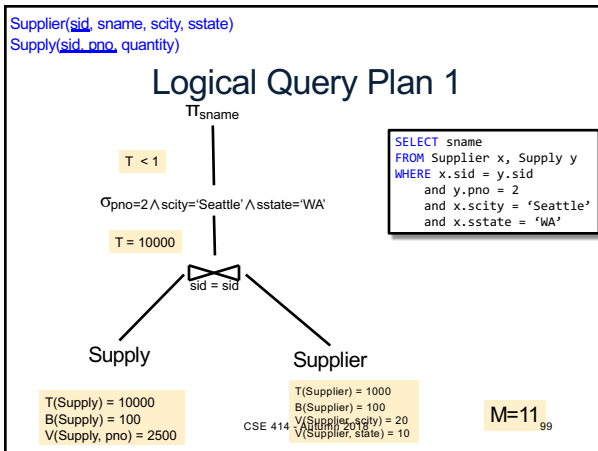
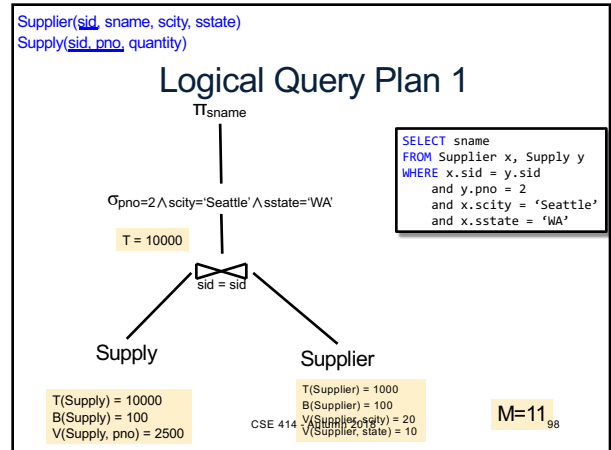
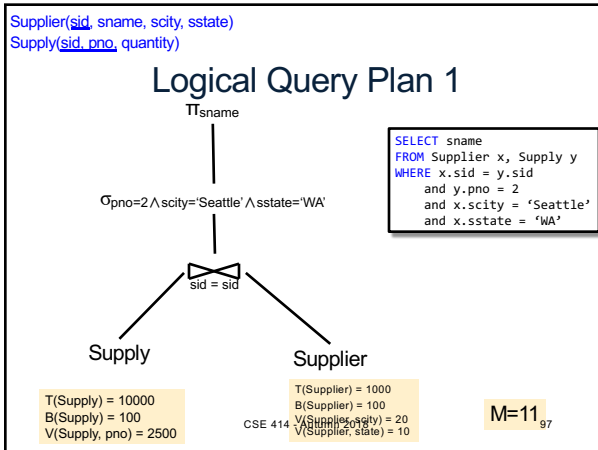
$R \bowtie S$

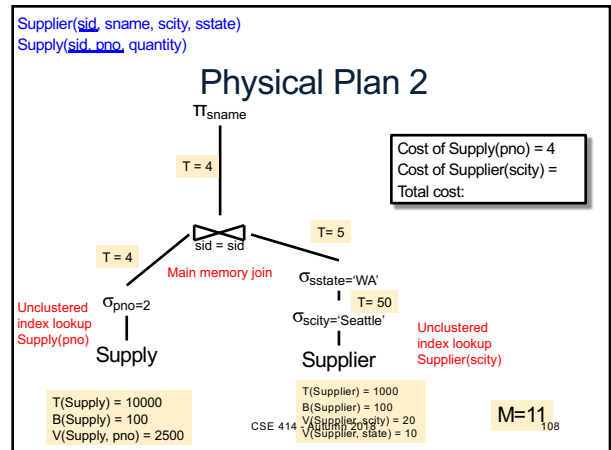
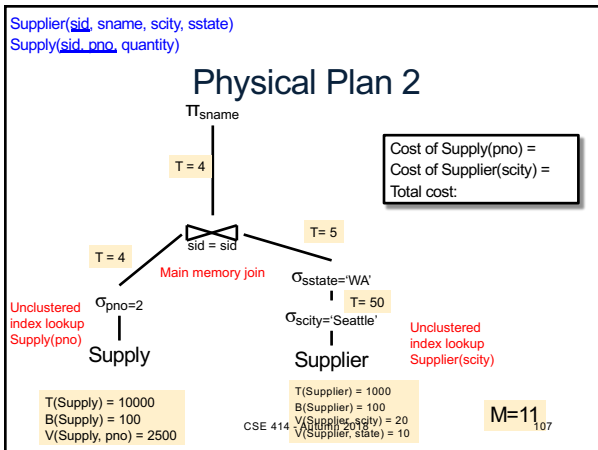
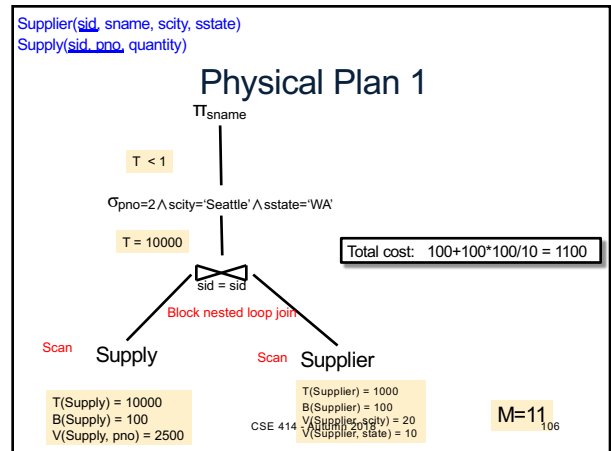
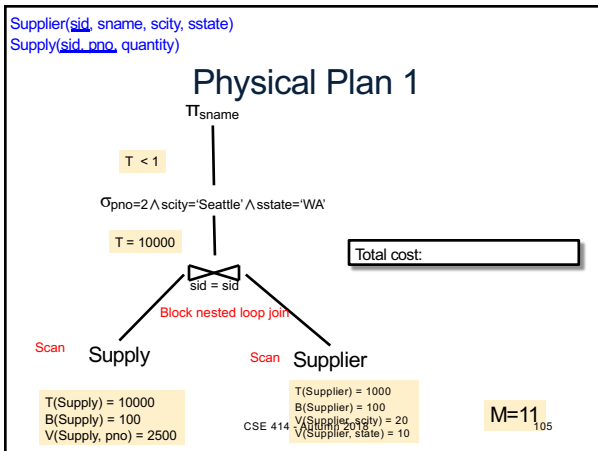
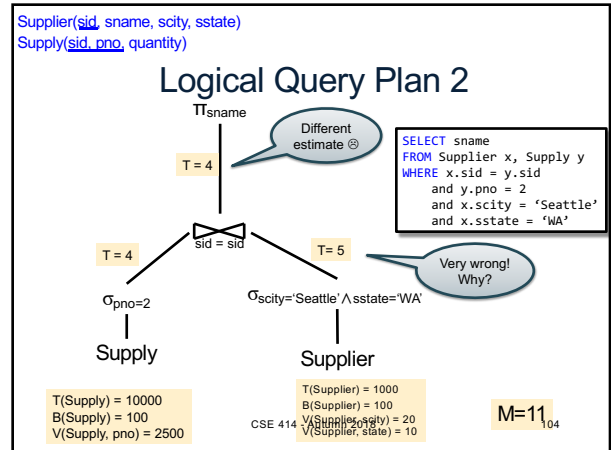
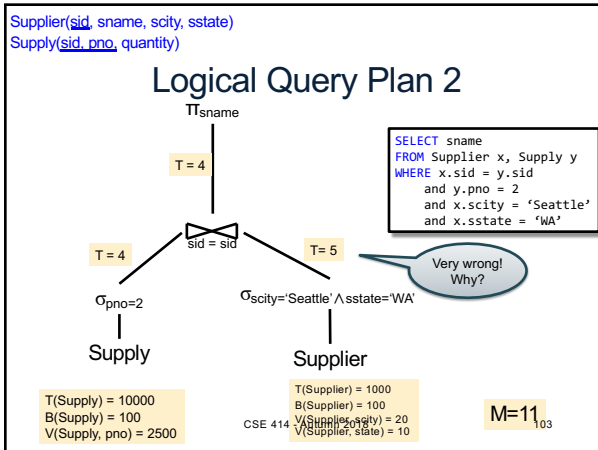
- Assume S has an index on the join attribute
- Iterate over R, for each tuple fetch corresponding tuple(s) from S
- **Cost:**
 - If index on S is clustered:
 $B(R) + T(R) * (B(S) * 1/V(S,a))$
 - If index on S is unclustered:
 $B(R) + T(R) * (T(S) * 1/V(S,a))$

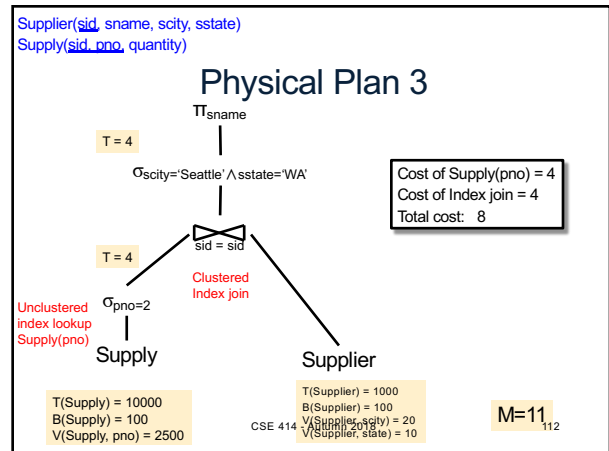
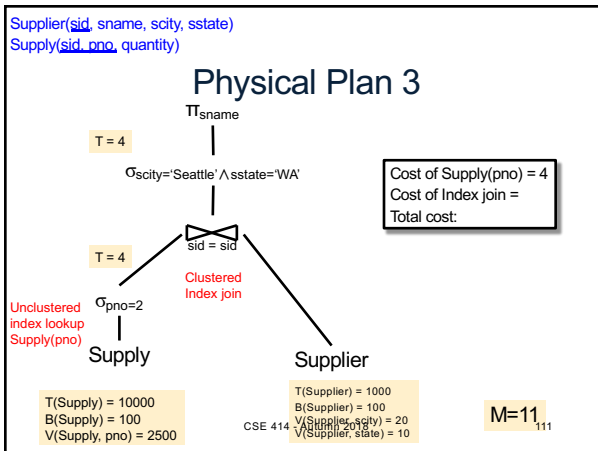
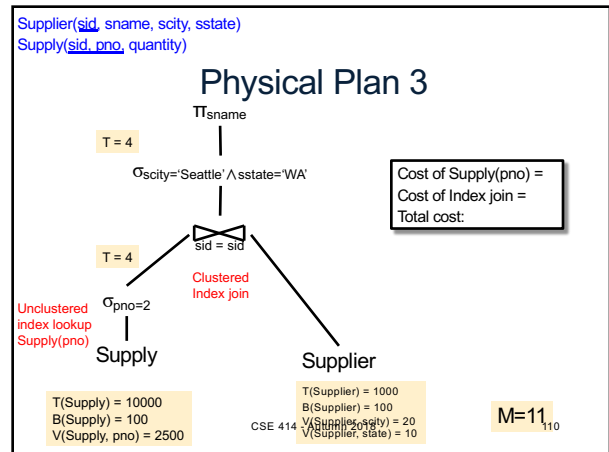
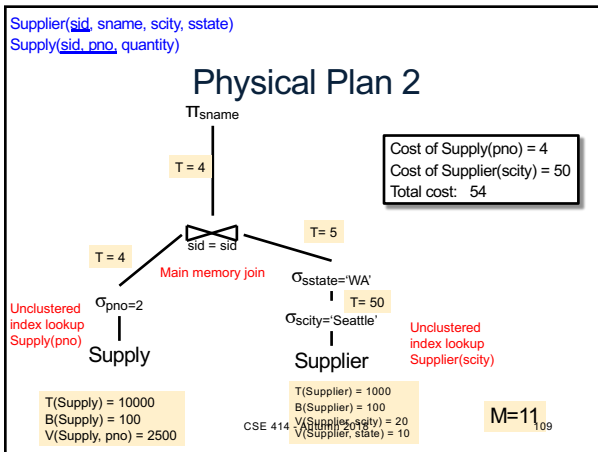
CSE 414 - Autumn 201895

Cost of Query Plans

CSE 414 - Autumn 201896







- ## 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
- CSE 414 - Autumn 2018 113