

Introduction to Database Systems CSE 414

Lecture 17: Basics of Query Optimization and Query Cost Estimation

CSE 414 - Autumn 2018

1

Announcements

- Midterm will be released by end of day today
- Need to start one HW6 step NOW:
 - <https://aws.amazon.com/education/awseducate/apply/>
 - Need to make an AWS account, can use existing Amazon account
 - Click on application button under Students and fill out form with your @uw.edu email
 - Will then be sent email for verification, must click to verify your email address

CSE 414 - Autumn 2018

2

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?

3

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

4

Cost of Reading Data From Disk

CSE 414 - Autumn 2018

5

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

CSE 414 - Autumn 2018

6

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

One_year(did, month) e.g. (1, Jan), (2, Jan)...(365, Dec)

Selectivity Factors for Conditions

How many tuples would this select:

```
SELECT *
FROM One_year
WHERE did = 32
```

1 tuple (out of 365)

One_year(did, month) e.g. (1, Jan), (2, Jan)...(365, Dec)

Selectivity Factors for Conditions

How many tuples would this select:

```
SELECT *
FROM One_year
WHERE month = Jan
```

31 tuples (out of 365)

This is roughly 1/12 of the tuples, because 12 distinct values equally distributed.

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 $f = 1/V(R, A)$

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

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

- $Cond1 \wedge Cond2 \wedge Cond3 \wedge \dots$
 - Selectivity $f = f1 * f2 * f3 * \dots$ (assumes independence)

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: $B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$ 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
Why: we know we can scan a full block to get the desired range
 - 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
Why: we know we can scan a full block to get the desired range
 - If index is unclustered: $T(R) * 1/V(R,a) = 5,000$ I/Os

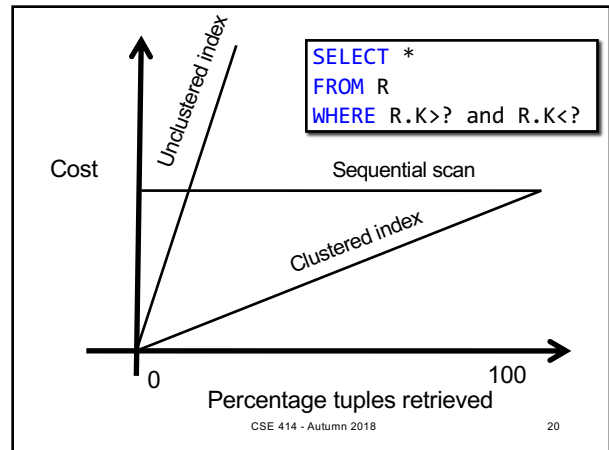
Index Based Selection

- Example: $B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$ $\text{cost of } \sigma_{a=v}(R) = ?$
- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:
 - If index is clustered: $B(R) * 1/V(R,a) = 100$ I/Os
 Why: we know we can scan a full block to get the desired range
 - 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 !

CSE 414 - Autumn 2018

19



Cost of Executing Operators (Focus on Joins)

CSE 414 - Autumn 2018

21

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

22

Join Algorithms

- Hash join
- Nested loop join

CSE 414 - Autumn 2018

23

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

24

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

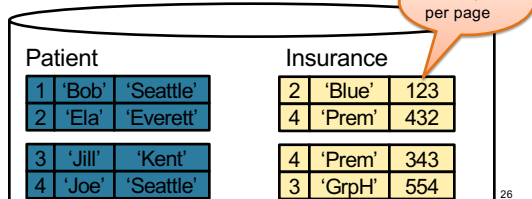
25

Hash Join Example

Patient(pid, name, address)

Insurance(pid, provider, policy_nb)

Patient \bowtie Insurance



26

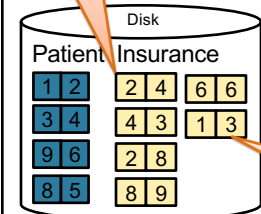
Hash Join Example

Patient \bowtie Insurance

Some large-enough #

Memory M = 21 pages

Showing pid only



This is one page with two tuples

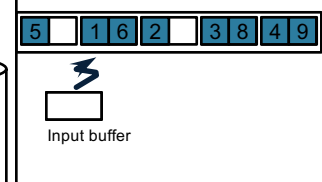
27

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



Input buffer

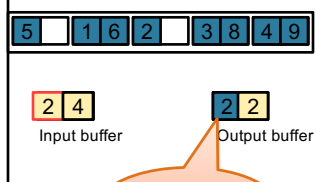
28

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



Write to disk or pass to next operator

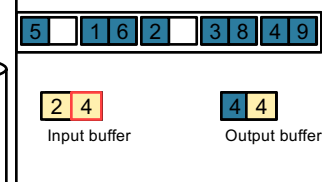
29

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



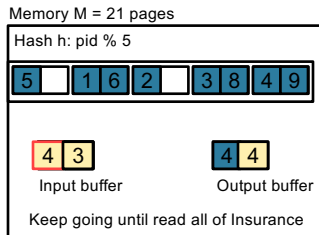
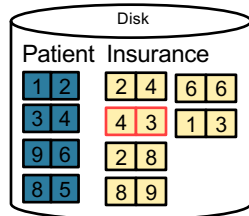
Input buffer

Output buffer

30

Hash Join Example

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



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

31

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

32

Nested Loop Joins

- Tuple-based nested loop $R \bowtie S$
- R is the outer relation, S is the inner relation

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

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

What is the Cost?

CSE 414 - Autumn 2018

33

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)
```

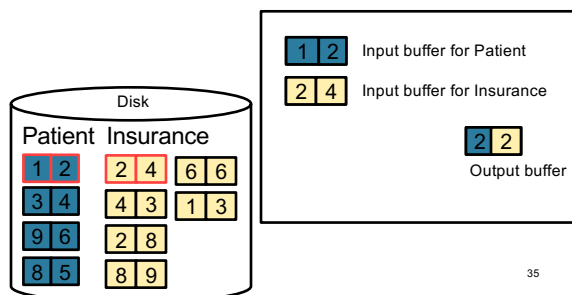
- Cost: $B(R) + B(R)B(S)$

What is the Cost?

CSE 414 - Autumn 2018

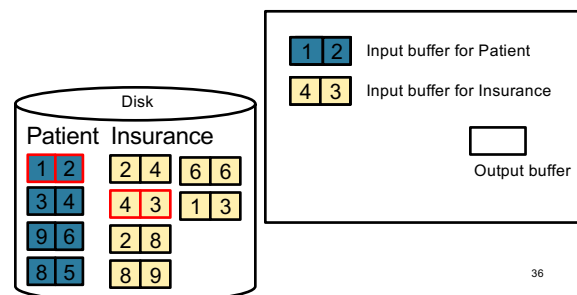
34

Page-at-a-time Refinement



35

Page-at-a-time Refinement



36

Page-at-a-time Refinement

Disk

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

1 2	Input buffer for Patient	
2 8	Input buffer for Insurance	
2 2	Output buffer	

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

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

37

INDEX JOINS

R	a	b	c
1	7	4	
...	
98	3	2	

S	c	d	e
3	43	7	
...	
9	24	9	

INDEX JOINS

44

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 2018 45

Index Nested Loop Join

If index on S is clustered:
 $B(R) + T(R) * (B(S) * 1/V(S,a))$

Still have to scan in R

Why is the multiplier term $T(R)$?

What does $1/V(S,a)$ represent?

$T(R)$ must be used because we cannot assume that a whole block of R ($B(R)$) will have the same attribute to join on, and thus use the same index access on S for.

$1/V(S,a)$ represents the nature of the B+ Tree index. We are only scanning as much as we need. Note that the performance of the index join will decrease as V decreases.

46

Index Nested Loop Join

If index on S is unclustered:
 $B(R) + T(R) * (T(S) * 1/V(S,a))$

Why did this change from $B(R)$ to $T(R)$?

Remember that tuples are stored on contiguous blocks. In a clustered index from before we know we can scan a single chunk of the disk to get the entire desired range. In an unclustered index we no longer can assume contiguous access. Thus we estimate that every tuple needs its own I/O operation.

47

GENERATING QUERY PLANS (REVIEW)

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

CSE 414 - Spring 2018 48

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

CSE 414 - Spring 2018 49

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"

CSE 414 - Spring 2018 50

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

Review: Physical Query Plan 1

(On the fly) TT_sname

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

(Nested loop)

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'
```

51

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

Review: Physical Query Plan 2

(On the fly) TT_sname

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

(Hash join)

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'
```

52

Query Optimization: Overview

- Compute cost of each operator
 - This depends on:
 - Table statistics (# of tuples etc)
 - Algorithm used
- Cost of a physical plan = sum(each operator cost)
- Cost each plan and choose the one with lowest cost

CSE 414 - Spring 2018 54

Cost of Query Plans

CSE 414 - Autumn 2018 55

