

PS2 out today. Lab 2 out today. Lab 1 due today - how was it?

Project Teams Due Wednesday

*Those of you who don't have groups -- send us email, or hand in a sheet with just your name on it saying you are looking for partners.*

Team Meetings In 2 Weeks

Show slide -- where are we

Last time:

Iterator model for executing queries, and intro to heap files and B+Trees

Today:

cost estimation basics  
access methods

What's the "cost" of a particular plan?

CPU cost (# of instructions)	- 1 ghz == 1 billions instrs / sec, 1 nsec / instr
Spinning Disk	
I/O cost (# of pages read, # of seeks)	- 100 MB / sec = 10 nsec / byte
Random I/O disk = page read + seek	- 10 msec / io operation = 100 io ops / sec
Flash disk	
	- 300-500 MB/sec seq rd/write
	- 10K + random iops / sec
	(100 usec / io operation)
	4K pages ~= 40 MB/sec

Random I/O can be a real killer (10 million instrs/seek) . When does a disk need to seek?

## Which do you think dominates in most database systems?

(Depends. Not always disk I/O. Typically vendors try to configure machines so they are 'balanced'. Can vary from query to query. Because of sequential access methods, buffer pool, we can often end up CPU bound.)

For example, fastest TPC-H benchmark result (data warehousing benchmark), on 10 TB of data, uses 1296 74 GB disks, which is 100 TB of storage. Add'l storage is partly for indices, but partly just because they needed add'l disk arms. 72 processors, 144 cores -- ~10 disks / processor!

But, if we do a bad job, random I/O can kill us!

100 tuples/page	select * from
10 pages RAM	emp, dept., kids
10 KB/page	where e.sal > 10k
	emp.dno = dept.dno
10 ms seek time	e.eid = kids.eid
100 MB/sec I/O	

ldeptl = 100 records = 1 page = 10 KB

lemp = 10K = 100 pages = 1 MB

lkidsl = 30K = 300 pages = 3 MB

⋈ (NL Join)

1000 / k (30000)

⋈(NL Join)

100 l \ 1000  
d σ<sub>sal>10k</sub>  
|  
e

1st Nested loops join -- 100,000 predicate ops; 2nd nested loops join -- 30,000,000 predicate ops

## Let's look at # disk I/Os assuming LRU and no indices

### Join 1

if d is outer:

- 1 scan of d

- 100 sec scans of e

- (100 x 100 pg. reads) -- cache doesn't benefit since e doesn't fit

- 1 scan of e: 1 seek + read in 1MB

  - $10 \text{ ms} + 1 \text{ MB} / 100 \text{ MB/sec} = 20 \text{ msec}$

- $20 \text{ ms} \times 100 \text{ depts} = 2 \text{ sec}$

- 10 msec seek to start of d and read into memory

2.1 secs

if d is inner

- read page of e -- 10 msec

- read all of d into RAM -- 10 msec

- seek back to e -- 10 sec

- scan rest of e -- 10 msec, joining with d in memory

Because d fits into memory, total cost is just 40 msec

NL # predicates evaluated doesn't depend on inner vs outer, but I/O cost does (due to caching effect!)

### Join 2

k inner:

- 1000 scans of 300 pages

- $3 / 100 = 30 \text{ msec} + 10 \text{ msec seek} = 40 \times 1000 = 40 \text{ sec}$

- if plan is pipelined, k must be inner

### So what will be cached?

That's the job of the buffer pool.

**Buffer pool** is a cache for memory access. Typically caches pages of files / indices.

Convenient "bottleneck" through which references to underlying pages go. When page is in buffer pool, don't need to read from disk. Updates can also be cached.

What is optimal buffer pool caching strategy? Always LRU?

No.

What if some relation doesn't fit into memory? (MRU preferred)

2 pages RAM, 3 pages of a relation R -- a, b c, accessed sequentially in a loop

	Access			
RAM	1	2	3	4
1	a	a	c	c
2		b	b	a

==> always evicting page we are about to access! Note that MRU would hit on 2/3

What if I know I will not be accessing a relation again in a query?

Are some records higher value than others?

Because data is large, we may want to use "access methods" that allow us to get the data we want with a minimum of I/Os. Especially true if just trying to look up one record (e.g., a particular employee) in a large database.

These access methods are dictionary structures we are familiar with (e.g., hash tables, trees), that are specially built to be "external" -- that is, disk resident. These are indices.

Hence, database system uses indices when it can.

Why do I say "when it can"?

(Some indices are good for different things -- e.g., hash tables don't support range queries)

Why not create indices on every attribute?

Update costs! (Indices are big, and so are typically stored on disk.)

Lets look at different types of access methods:

types of access methods

- heapfiles

- (sorted files)

- index files

  - leaves as data (primary index)

  - leaves as rids (secondary index)

  - clustered vs. unclustered

types of indexes

- hash files

- b-trees

- r-trees

- ...

tradeoffs between different storage representations

- scan vs. search vs. insert vs. delete

N - number of records

P - pages in index/file

R - pages in range

(units are # I/ Os)	Heap File	Hash File	B+Tree
<b>Search</b>	$P/2$	$O(1)$	$O(\log_b N)$
<b>Scan</b>	$P$	-	$O(\log_b N) + R$
<b>Delete</b>	$P/2$	$O(1)$	$O(\log_b N)$

## heap files

search cost  $\sim$  scan cost  $\sim$  delete cost

- linked list
- directory
- array of objs

\*\*\*\* STUDY BREAK \*\*\*\*

## hash files

search cost  $\sim$  1

insert cost  $\sim$  1

delete cost  $\sim$  1

range scan is equal to sequential scan

map(key)  $\rightarrow$  {rid} -- could also store map(key)  $\rightarrow$  {record}

suppose we have a hash-table that for employees stored on disk hashed on the 'name' attribute.

$h(\text{name}) \rightarrow \{1..k\}$

$h(x) = x \bmod k$ ; //not very uniform but works

(consider performance when number of distinct values is small, or all numbers even)

We have  $k$  buckets, and one page per bucket. Store to pages by hashing, appending the record to that page.  
 $k$  chosen to fit in memory.

Then, if a query looks for 'mike',  
we can find him in a single I/O (we know exactly where to go to find him.)

(Show slide)

How many buckets do we need? (hard to tell!)

Too many? -- Wasteful

Too few? -- Long chains that we have to follow

Extensible hashing:

Gave example before -- in principle, we can overflow a page, in which case we need an overflow chain. These overflow chains can get long and slow down the performance of our algorithm.

So, instead, we split hash buckets when they get full. We do this with a family of hash functions, where we switch to the next level when we get too many in the current level.

define  $h(v) = \{0, 1, \dots, b\}$

$h_k(v) = h(v) \bmod 2^k$

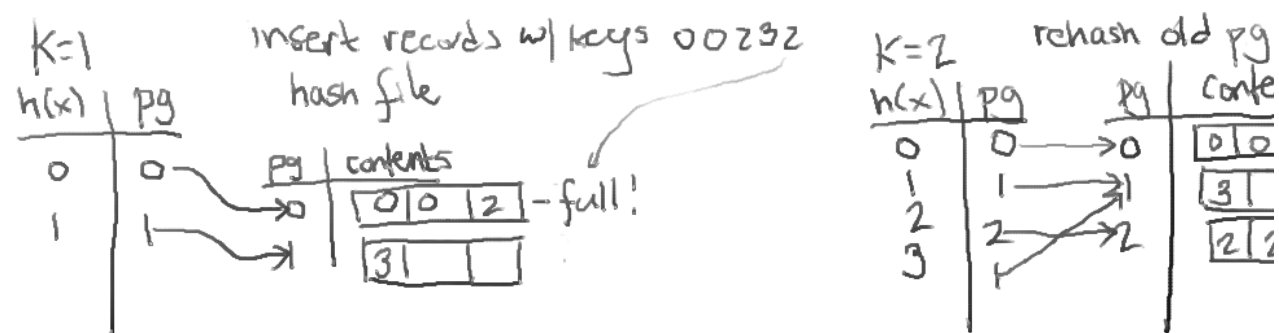
$h_1(v) = \{0, 1\}$

$h_2(v) = \{0, 1, 2, 3\}$

...

Maintain a current hash function, its "levels", and a directory that tells you where each bucket is on disk

example:

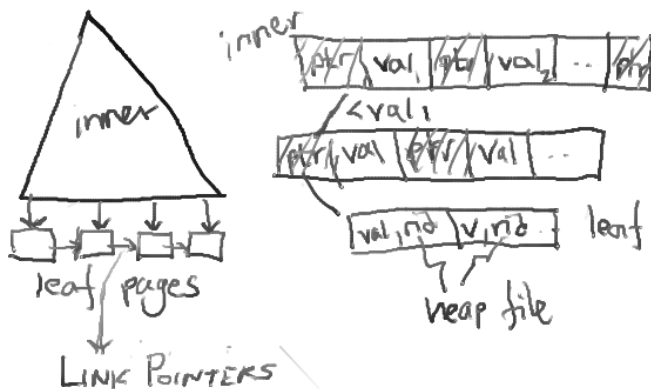


(Show slide)

(Lecture ends here -- B+Trees next time)

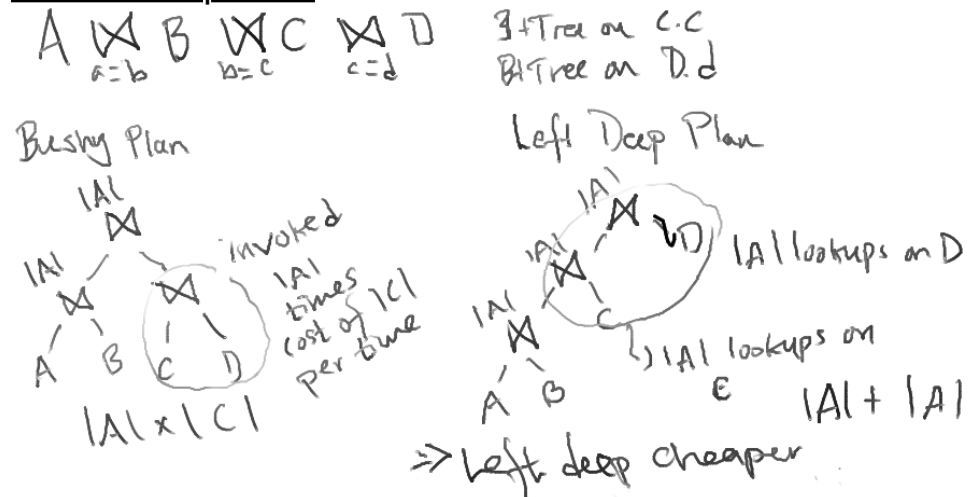
### B+-Trees:

Show example tree:



Not going to discuss details of how insertion/splitting/rebalancing work.

### Discussion points:



- What is point of link pointers?
- Why not store data in intermediate nodes
- Page occupancy (why does this matter)
- Key size?



- Balanced (why important?)
- How many levels in practice  
 disk page = 4096 bytes; values 4 bytes; ptrs 4 bytes = 512 ptrs /  
 node -->  $512^4 = 68$  billion
- $\log(n)$  for lookup/insert/delete; how many I/Os in practice?  
 (1, maybe 2, since top levels will fit in RAM,  
 E.g., top 2 levels are just  $513 * 4096$  bytes == 2MB)
- Scan is random, unless index is clustered
- Node format  
 "Fill factor": percentage of each page that is used.  
 Every node except root node is at least 50% full
- Why would we not want this to be 100%?  
 Ideally, 67% full (where from)?

Can be clustered or unclustered  
 clustered means that data is physically stored in order of index  
 often, leaves of index contain actual data records

Lots of other indexes, e.g., R+Trees.