



Carnegie Mellon Univ.
Dept. of Computer Science
15-415 - Database Applications

Lecture #9
Storing Data: Disks and Files

#1



Overview

- Memory hierarchy
- RAID (briefly)
- Disk space management
- Buffer management
- Files of records
- Page Formats
- Record Formats

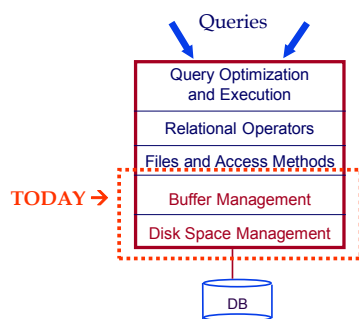
Faloutsos

CMU SCS 15-415

#2



DBMS Layers:



Faloutsos

CMU SCS 15-415

#3



Leverage OS for disk/file management?

- Layers of abstraction are good ... but:

Faloutsos

CMU SCS 15-415

#4



Leverage OS for disk/file management?

- Layers of abstraction are good ... but:
 - Unfortunately, OS often **gets in the way** of DBMS

Faloutsos

CMU SCS 15-415

#5



Leverage OS for disk/file management?

- DBMS wants/needs to do things “its own way”
 - Specialized prefetching
 - Control over buffer replacement policy
 - LRU not always best (sometimes worst!!)
 - Control over thread/process scheduling
 - “Convoy problem”
 - Arises when OS scheduling conflicts with DBMS locking
 - Control over flushing data to disk
 - WAL protocol requires flushing log entries to disk

Faloutsos

CMU SCS 15-415

#6



Disks and Files

- DBMS stores information on disks.
 - but: disks are (relatively) VERY slow!
- Major implications for DBMS design!



Faloutsos

CMU SCS 15-415

#7



Disks and Files

- Major implications for DBMS design:
 - **READ**: disk \rightarrow main memory (RAM).
 - **WRITE**: reverse
 - Both are high-cost operations, relative to in-memory operations, so must be planned carefully!

Faloutsos

CMU SCS 15-415

#8



Why Not Store It All in Main Memory?

- {Note to self: Get PowerPoint to insert reality here}

Faloutsos

CMU SCS 15-415

#9



Why Not Store It All in Main Memory?

- *Costs too much.*
 - disk: $\sim \$0.85/\text{Gb}$; memory: $\sim \$100/\text{Gb}$
 - High-end Databases today in the 10-100 TB range.
 - Approx 60% of the cost of a production system is in the disks.
- *Main memory is volatile.*
- *Note*: some specialized systems do store entire database in main memory.

Faloutsos

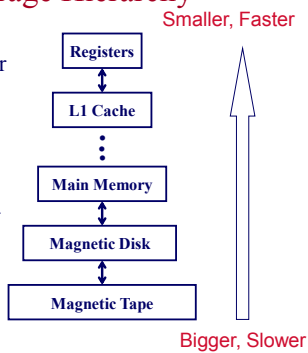
CMU SCS 15-415

#10



The Storage Hierarchy

- Main memory (RAM) for currently used data.
- Disk for the main database (secondary storage).
- Tapes for archiving older versions of the data (tertiary storage).



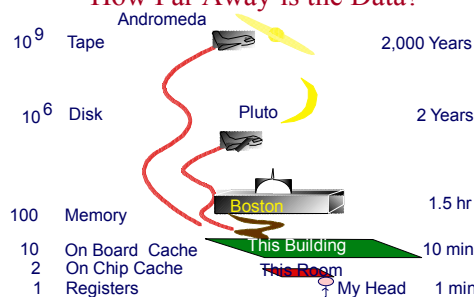
Faloutsos

CMU SCS 15-415

#11



Jim Gray's Storage Latency Analogy: How Far Away is the Data?



Faloutsos

CMU SCS 15-415

#12



Disks

- Secondary storage device of choice.
- Main advantage over tapes: *random access* vs. *sequential*.
- Data is stored and retrieved in units called *disk blocks* or *pages*.
- Unlike RAM, time to retrieve a disk page varies depending upon location on disk.
 - relative placement of pages on disk is important!

Faloutsos

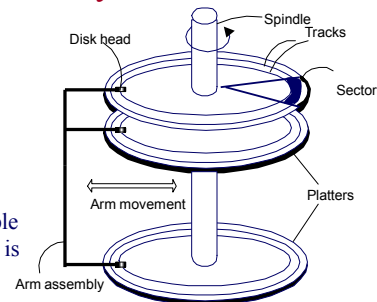
CMU SCS 15-415

#13



Anatomy of a Disk

- Sector
- Track
- Cylinder
- Platter
- Block size = multiple of sector size (which is fixed)



Faloutsos

CMU SCS 15-415

#14



Accessing a Disk Page

- Time to access (read/write) a disk block:
 - *seek time*: moving arms to position disk head on track
 - *rotational delay*: waiting for block to rotate under head
 - *transfer time*: actually moving data to/from disk surface

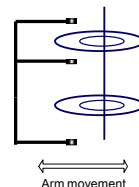
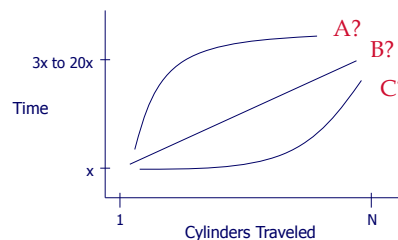
Faloutsos

CMU SCS 15-415

#15



Seek Time



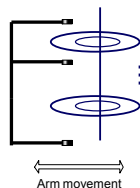
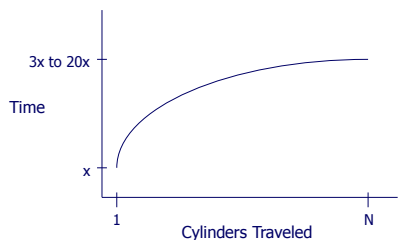
Faloutsos

CMU SCS 15-415

#16



Seek Time



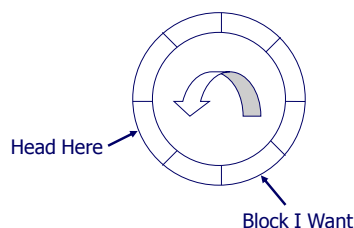
Faloutsos

CMU SCS 15-415

#17



Rotational Delay



Faloutsos

CMU SCS 15-415

#18



Accessing a Disk Page

- Relative times?
 - *seek time*: about 1 to 20msec
 - *rotational delay*: 0 to 10msec
 - *transfer time*: < 1msec per 4KB page



Faloutsos

CMU SCS 15-415

#19



Seek time & rotational delay dominate

- Key to lower I/O cost:
reduce seek/rotation delays!
- Also note: For shared disks, much time spent waiting in queue for access to arm/controller



Faloutsos

CMU SCS 15-415

#20



Arranging Pages on Disk

- “*Next*” block concept:
 - blocks on same track, followed by
 - blocks on same cylinder, followed by
 - blocks on adjacent cylinder
- Accessing ‘next’ block is cheap
- An important optimization: pre-fetching
 - See R&G page 323

Faloutsos

CMU SCS 15-415

#21



Rules of thumb...

- Memory access much faster than disk I/O
(~ 1000x)
- “Sequential” I/O faster than “random” I/O
(~ 10x)

Faloutsos

CMU SCS 15-415

#22



Overview

- Memory hierarchy
- RAID (briefly)
- Disk space management
- Buffer management
- Files of records
- Page Formats
- Record Formats

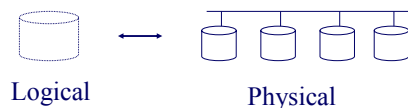
Faloutsos

CMU SCS 15-415

#23



Disk Arrays: RAID



- Benefits:
 - Higher throughput (via data “striping”)
 - Longer MTTF (via redundancy)

Faloutsos

CMU SCS 15-415

#24



Overview

- Memory hierarchy
- RAID (briefly)
- Disk space management
- Buffer management
- Files of records
- Page Formats
- Record Formats

Faloutsos

CMU SCS 15-415

#25



Disk Space Management

- Lowest layer of DBMS software manages space on disk
- Higher levels call upon this layer to:
 - allocate/de-allocate a page
 - read/write a page
- Best if requested pages are stored **sequentially** on disk! Higher levels don't need to know if/how this is done, nor how free space is managed.

Faloutsos

CMU SCS 15-415

#26



Overview

- Memory hierarchy
- RAID (briefly)
- Disk space management
- Buffer management
- Files of records
- Page Formats
- Record Formats

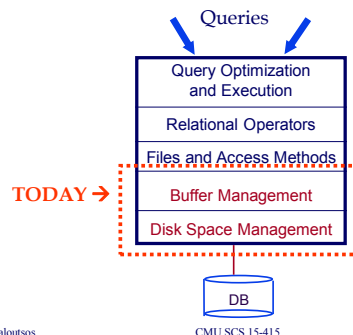
Faloutsos

CMU SCS 15-415

#27



Recall: DBMS Layers



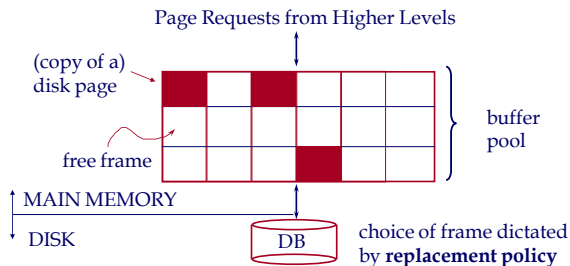
Faloutsos

CMU SCS 15-415

#28



Buffer Management in a DBMS



Faloutsos

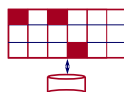
CMU SCS 15-415

#29



Buffer Management in a DBMS

- Data must be in RAM for DBMS to operate on it!
- Buffer Mgr hides the fact that not all data is in RAM



Faloutsos

CMU SCS 15-415

#30



CMU SCS

When a Page is Requested ...

Buffer pool information table contains:

$\langle \text{frame\#}, \text{pageid}, \text{pin_count}, \text{dirty-bit} \rangle$

- If requested page is not in pool:
 - Choose an (un-pinned) frame for *replacement*
 - If frame is “dirty”, write it to disk
 - Read requested page into chosen frame
- *Pin* the page and return its address

Faloutsos

CMU SCS 15-415

#31



CMU SCS

When a Page is Requested ...

- If requests can be predicted (e.g., sequential scans)
- then pages can be pre-fetched several pages at a time!

Faloutsos

CMU SCS 15-415

#32



CMU SCS

More on Buffer Management

- When done, requestor of page must
 - unpin it, and
 - indicate whether page has been modified: *dirty* bit
- Page in pool may be requested many times:
 - *pin count*
- if *pin count* = 0 (“unpinned”), page is candidate for replacement

Faloutsos

CMU SCS 15-415

#33



CMU SCS

More on Buffer Management

- CC & recovery may entail additional I/O when a frame is chosen for replacement. (*Write-Ahead Log* protocol; more later.)

Faloutsos

CMU SCS 15-415

#34



CMU SCS

Buffer Replacement Policy

- Frame is chosen for replacement by a *replacement policy*:
 - Least-recently-used (LRU), MRU, Clock, etc.
- Policy can have big impact on # of I/O's; depends on the *access pattern*.

Faloutsos

CMU SCS 15-415

#35



CMU SCS

LRU Replacement Policy

- *Least Recently Used (LRU)*
 - for each page in buffer pool, keep track of time last *unpinned*
 - replace the frame which has the oldest (earliest) time
 - very common policy: intuitive and simple
- Problems?

Faloutsos

CMU SCS 15-415

#36



LRU Replacement Policy

- Problem: *Sequential flooding*
 - LRU + repeated sequential scans.
 - # *buffer frames* < # *pages in file* means each page request causes an I/O. MRU much better in this situation (but not in all situations, of course).

Faloutsos

CMU SCS 15-415

#37



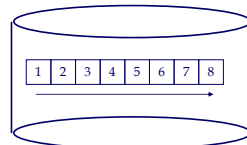
Sequential Flooding – Illustration

LRU:

BUFFER POOL			
102	116	242	105

MRU:

BUFFER POOL			
102	116	242	105



Repeated scan of file ...

Faloutsos

CMU SCS 15-415

#38



Sequential Flooding – Illustration

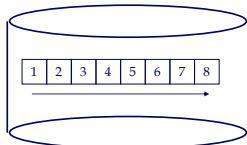
LRU:

BUFFER POOL			
1	2	3	4

will not re-use these pages;

MRU:

BUFFER POOL			
4	116	242	105



Repeated scan of file ...

Faloutsos

CMU SCS 15-415

#39



Other policies?

- LRU is often good - but needs timestamps and sorting on them
- something easier to maintain?

Faloutsos

CMU SCS 15-415

#40



“Clock” Replacement Policy

Main ideas:

- Approximation of LRU.
- Instead of maintaining & sorting time-stamps, find a ‘reasonably old’ frame to evict.
- How? by round-robin, and marking each frame - frames are evicted the second time they are visited.
- Specifically:

Faloutsos

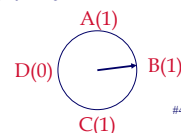
CMU SCS 15-415

#41



“Clock” Replacement Policy

- Arrange frames into a cycle, store one “reference bit” per frame
- When pin count goes to 0, reference bit set on (= ‘one life left’ - not ready for eviction yet)
- When replacement necessary, get the next frame that has reference-bit = 0



Faloutsos

CMU SCS 15-415

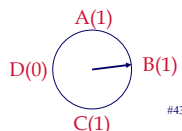
#42



CMU SCS

“Clock” Replacement Policy

```
do {
  if (pincount == 0 && ref bit is off)
    choose current page for replacement;
  else if (pincount == 0 && ref bit is on)
    turn off ref bit;
    advance current frame;
} until a page is chosen for replacement;
```



Faloutsos

CMU SCS 15-415

#43



CMU SCS

Summary

- Buffer manager brings pages into RAM.
- Very important for performance
 - Page stays in RAM until released by requestor.
 - Written to disk when frame chosen for replacement (which is sometime after requestor releases the page).
 - Choice of frame to replace based on *replacement policy*.
 - Good to *pre-fetch* several pages at a time.

Faloutsos

CMU SCS 15-415

#44



CMU SCS

Overview

- Memory hierarchy
- RAID (briefly)
- Disk space management
- Buffer management
- Files of records
- Page Formats
- Record Formats

Faloutsos

CMU SCS 15-415

#45



CMU SCS

Files

- FILE: A collection of pages, each containing a collection of records.
- Must support:
 - insert/delete/modify record
 - read a particular record (specified using *record id*)
 - scan all records (possibly with some conditions on the records to be retrieved)

Faloutsos

CMU SCS 15-415

#46



CMU SCS

Alternative File Organizations

Several alternatives (w/ trade-offs):

- Heap files: Suitable when typical access is a file scan retrieving all records.
 - Sorted Files:
 - Index File Organizations:
- } later

Faloutsos

CMU SCS 15-415

#47



CMU SCS

Files of records

- Heap of pages
 - as linked list or
 - directory of pages

Faloutsos

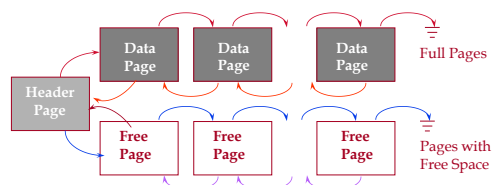
CMU SCS 15-415

#48



CMU SCS

Heap File Using Lists



- The header page id and Heap file name must be stored someplace.
- Each page contains 2 'pointers' plus data.

Faloutsos

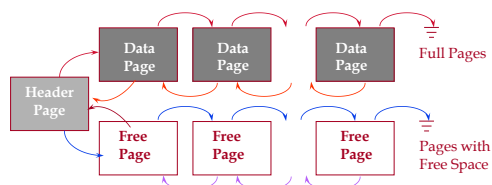
CMU SCS 15-415

#49



CMU SCS

Heap File Using Lists



- Any problems?

Faloutsos

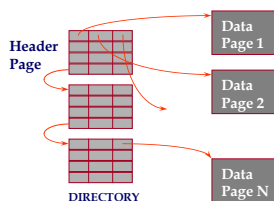
CMU SCS 15-415

#50



CMU SCS

Heap File Using a Page Directory



Faloutsos

CMU SCS 15-415

#51



CMU SCS

Heap File Using a Page Directory

- The entry for a page can include the number of free bytes on the page.
- The directory is a collection of pages; linked list implementation is just one alternative.

– *Much smaller than linked list of all HF pages!*

Faloutsos

CMU SCS 15-415

#52



CMU SCS

Overview

- Memory hierarchy
- RAID (briefly)
- Disk space management
- Buffer management
- Files of records
- Page Formats
- Record Formats

Faloutsos

CMU SCS 15-415

#53



CMU SCS

Page Formats

- fixed length records
- variable length records

Faloutsos

CMU SCS 15-415

#54



Page Formats

Important concept: *rid* == record id

Q0: why do we need it?

Q1: How to mark the location of a record?

Q2: Why not its byte offset in the file?

Faloutsos

CMU SCS 15-415

#55



Page Formats

Important concept: *rid* == record id

Q0: why do we need it?

A0: eg., for indexing

Q1: How to mark the location of a record?

A1: $rid = record\ id = page-id \ \& \ slot-id$

Q2: Why not its byte offset in the file?

A2: too much re-organization on ins/del.

Faloutsos

CMU SCS 15-415

#56



Fixed length records

- Q: How would you store them on a page/file?

Faloutsos

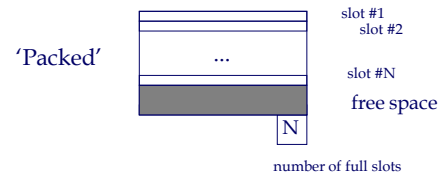
CMU SCS 15-415

#57



Fixed length records

- Q: How would you store them on a page/file?
- A1: How about:



Faloutsos

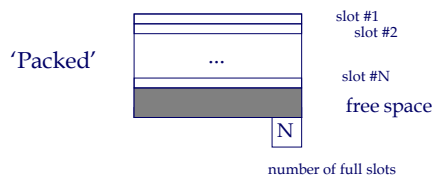
CMU SCS 15-415

#58



Fixed length records

- A1: How about: **BUT**: On insertion/deletion, we have too much to reorganize/update



Faloutsos

CMU SCS 15-415

#59



Fixed length records

- What would you do?

Faloutsos

CMU SCS 15-415

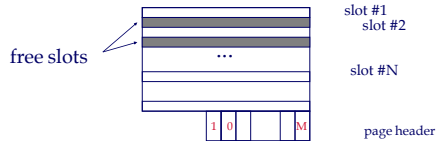
#60



CMU SCS

Fixed length records

- Q: How would you store them on a page/file?
- A2: Bitmaps



Faloutsos

CMU SCS 15-415

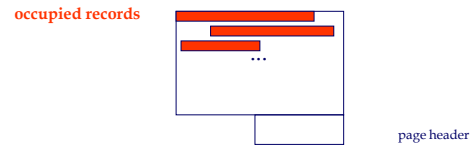
#61



CMU SCS

Variable length records

- Q: How would you store them on a page/file?



Faloutsos

CMU SCS 15-415

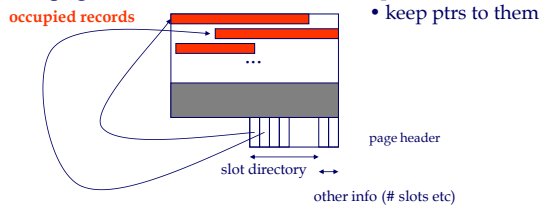
#62



CMU SCS

Variable length records

- Q: How would you store them on a page/file?



Faloutsos

CMU SCS 15-415

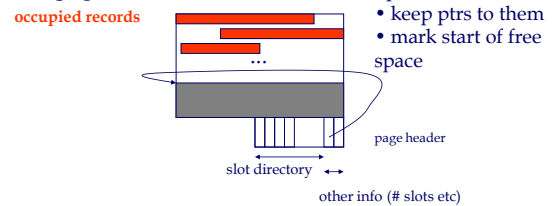
#63



CMU SCS

Variable length records

- Q: How would you store them on a page/file?



Faloutsos

CMU SCS 15-415

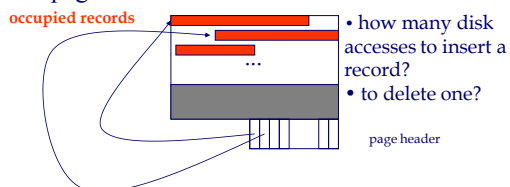
#64



CMU SCS

Variable length records

- Q: How would you store them on a page/file?



Faloutsos

CMU SCS 15-415

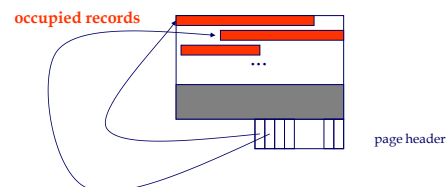
#65



CMU SCS

Variable length records

- SLOTTED PAGE organization - popular.



Faloutsos

CMU SCS 15-415

#66



Overview

- Memory hierarchy
- RAID (briefly)
- Disk space management
- Buffer management
- Files of records
- Page Formats
- Record Formats

Faloutsos

CMU SCS 15-415

#67



Formats of records

- Fixed length records
 - How would you store them?
- Variable length records

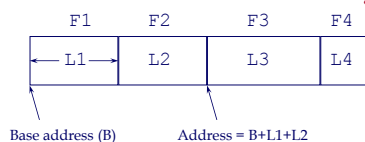
Faloutsos

CMU SCS 15-415

#68



Record Formats: Fixed Length



- Information about field types same for all records in a file; stored in *system catalogs*.
- Finding *i*'th field done via arithmetic.

Faloutsos

CMU SCS 15-415

#69



Formats of records

- Fixed length records: straightforward - store info in catalog
- Variable length records: encode the length of each field
 - store its length or
 - use a field delimiter

Faloutsos

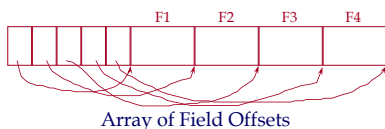
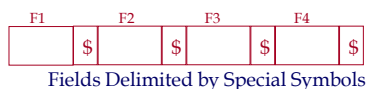
CMU SCS 15-415

#70



Variable Length records

- Two alternative formats (# fields is fixed):



Offset approach: usually superior (direct access to *i*-th field)

Faloutsos

CMU SCS 15-415

#71



Conclusions

- Memory hierarchy
- Disks: (>1000x slower) - thus
 - pack info in blocks
 - try to fetch nearby blocks (sequentially)
- Buffer management: very important
 - LRU, MRU, Clock, etc
- Record organization: Slotted page

Faloutsos

CMU SCS 15-415

#72