

CMU SCS

**Carnegie Mellon Univ.**  
**Dept. of Computer Science**  
**15-415/615 - DB Applications**

C. Faloutsos & A. Pavlo  
 Lecture #8 (R&G ch9)  
*Storing Data: Disks and Files*

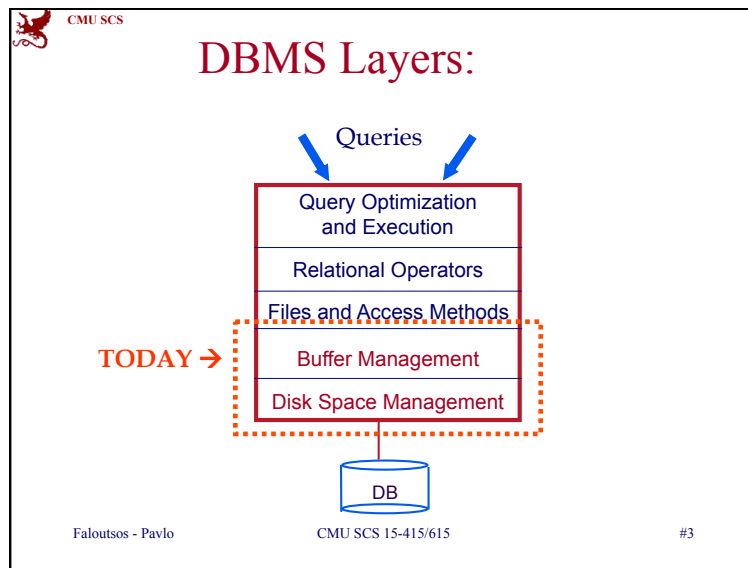
Faloutsos - Pavlo CMU SCS 15-415/615 #1

CMU SCS

## Overview

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

Faloutsos - Pavlo CMU SCS 15-415/615 #2



CMU SCS

## Leverage OS for disk/file management?

- Layers of abstraction are good ... but:

Faloutsos - Pavlo CMU SCS 15-415/615 #4

CMU SCS

## Leverage OS for disk/file management?

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

Faloutsos - Pavlo CMU SCS 15-415/615 #5

CMU SCS

## Leverage OS for disk/file management?

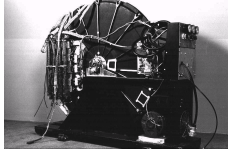

- DBMS wants/needs to do things “its own way”
  - **Specialized prefetching**
  - **Control over buffer replacement policy**
    - LRU not always best (sometimes worst!!)
  - **Control over thread/process scheduling**
    - “Convoy problem”
      - Arises when OS scheduling conflicts with DBMS locking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

Faloutsos - Pavlo CMU SCS 15-415/615 #6

CMU SCS

## Disks and Files

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

Faloutsos - Pavlo CMU SCS 15-415/615 #7

CMU SCS

## Disks and Files

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

Faloutsos - Pavlo CMU SCS 15-415/615 #8

CMU SCS

## Why Not Store It All in Main Memory?

Faloutsos - Pavlo CMU SCS 15-415/615 #9

CMU SCS

## Why Not Store It All in Main Memory?


- *Costs too much.*
  - disk: ~\$0.1/Gb; memory: ~\$10/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 - Pavlo CMU SCS 15-415/615 #10

CMU SCS

## The Storage Hierarchy

Smaller, Faster



Bigger, Slower

Faloutsos - Pavlo CMU SCS 15-415/615 #11

CMU SCS

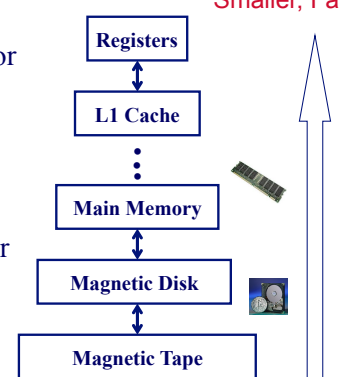
## The Storage Hierarchy

–Main memory (RAM) for currently used data.

–Disk for the main database (secondary storage).

–Tapes for archiving older versions of the data (tertiary storage).

Smaller, Faster



Bigger, Slower

Faloutsos - Pavlo CMU SCS 15-415/615 #12

### The Storage Hierarchy

Smaller, Faster

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

Registers  
 L1 Cache  
 Main Memory  
 SSD  
 Magnetic Disk  
 Magnetic Tape

Bigger, Slower

Faloutsos - Pavlo CMU SCS 15-415/615 #13

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

|                   |               |  |         |
|-------------------|---------------|--|---------|
| 10**9 tape        | Andromeda     |  | 2,000yr |
| 10**6 disk        | Pluto         |  | 2yr     |
| 100 Memory        | Pittsburgh    |  | 1.5h    |
| 10 On board cache | This building |  | 10min   |
| 2 on chip cache   | This room     |  | 1min    |
| 1 registers       | In my head    |  |         |

Faloutsos - Pavlo CMU SCS 15-415/615 #14

### 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 - Pavlo CMU SCS 15-415/615 #15


### Anatomy of a Disk

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

Faloutsos - Pavlo CMU SCS 15-415/615 #16

CMU SCS

## Accessing a Disk Page




- Time to access (read/write) a disk block:
  - .
  - .
  - .

Faloutsos - Pavlo      CMU SCS 15-415/615      #17

CMU SCS

## Accessing a Disk Page

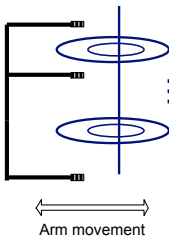


- Time to access (read/write) a disk block:
  - *seek time*: moving arms to position disk head on track
  - *rotational delay*: waiting for block to rotate under head
  - *transfer time*: actually moving data to/from disk surface

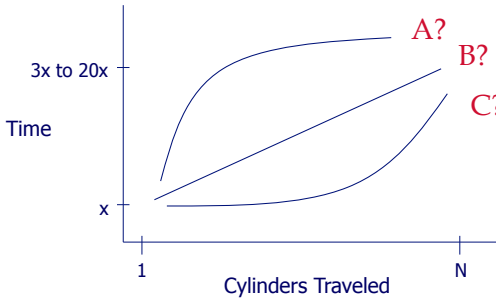
Faloutsos - Pavlo      CMU SCS 15-415/615      #18

CMU SCS

## Seek Time



Arm movement

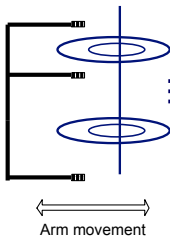


The graph shows Time on the y-axis and Cylinders Traveled on the x-axis. The x-axis has markers at 1 and N. The y-axis has markers at x and 3x to 20x. Three curves are shown, labeled A?, B?, and C? in red. Curve A? is the highest, B? is in the middle, and C? is the lowest. All curves start at (1, x) and increase as they approach N.

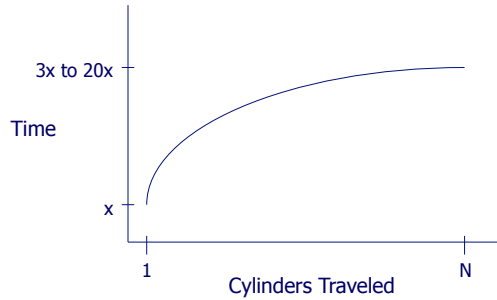
Faloutsos - Pavlo      CMU SCS 15-415/615      #19

CMU SCS

## Seek Time



Arm movement



The graph shows Time on the y-axis and Cylinders Traveled on the x-axis. The x-axis has markers at 1 and N. The y-axis has markers at x and 3x to 20x. A single curve is shown, starting at (1, x) and increasing as it approaches N.

Faloutsos - Pavlo      CMU SCS 15-415/615      #20

CMU SCS

## Rotational Delay

Head Here

Block I Want

Faloutsos - Pavlo      CMU SCS 15-415/615      #21

CMU SCS

## Accessing a Disk Page

- Relative times?
  - *seek time:*
  - *rotational delay:*
  - *transfer time:*

Faloutsos - Pavlo      CMU SCS 15-415/615      #22

CMU SCS

## Accessing a Disk Page

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

Seek

Rotate

transfer

Faloutsos - Pavlo      CMU SCS 15-415/615      #23

CMU SCS

## Seek time & rotational delay dominate

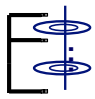
- Key to lower I/O cost: **reduce seek/rotation delays!**
- Also note: For shared disks, much time spent waiting in queue for access to arm/controller

Seek

Rotate

transfer

Faloutsos - Pavlo      CMU SCS 15-415/615      #24



## 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

Q1: Why not equal?

Q2: Why?

Faloutsos - Pavlo
CMU SCS 15-415/615
#25

## Rules of thumb...

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

write on blackboard

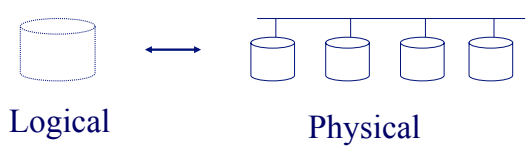
Faloutsos - Pavlo
CMU SCS 15-415/615
#26


## Overview

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

Faloutsos - Pavlo
CMU SCS 15-415/615
#27

## Disk Arrays: RAID





Prof. Garth Gibson

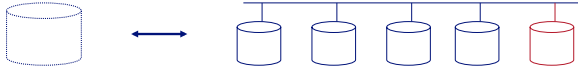
- Benefits:
  - Higher throughput (via data “striping”)
  - Longer MTTF

(Why?)

Faloutsos - Pavlo
CMU SCS 15-415/615
#28

CMU SCS

## Disk Arrays: RAID



Logical                      Physical

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

Faloutsos - Pavlo                      CMU SCS 15-415/615                      #29

CMU SCS

## Overview

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

Faloutsos - Pavlo                      CMU SCS 15-415/615                      #30

CMU SCS

## 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 - Pavlo                      CMU SCS 15-415/615                      #31

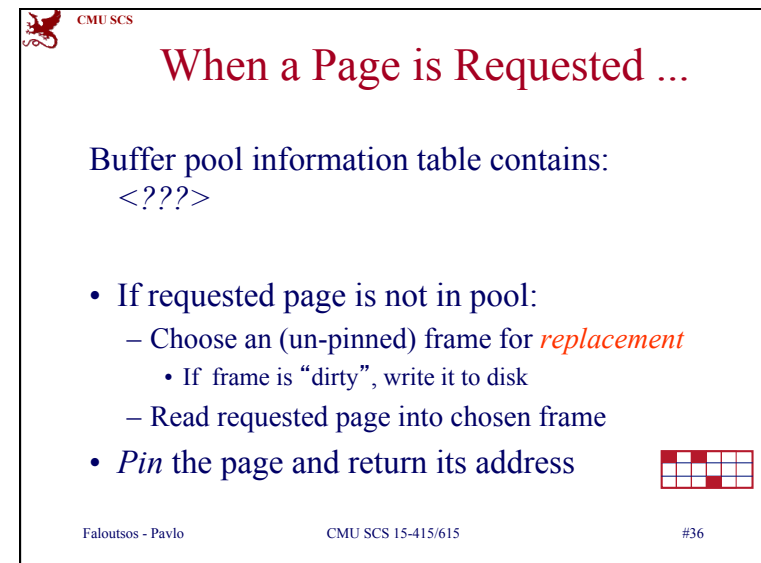
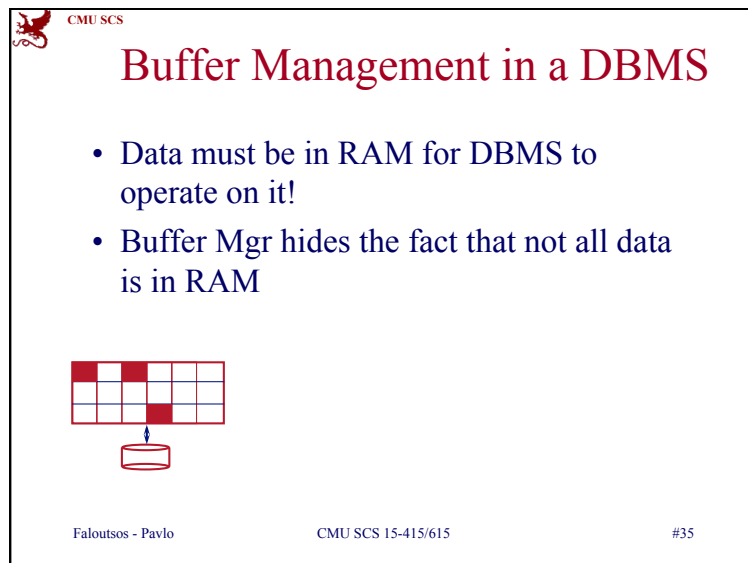
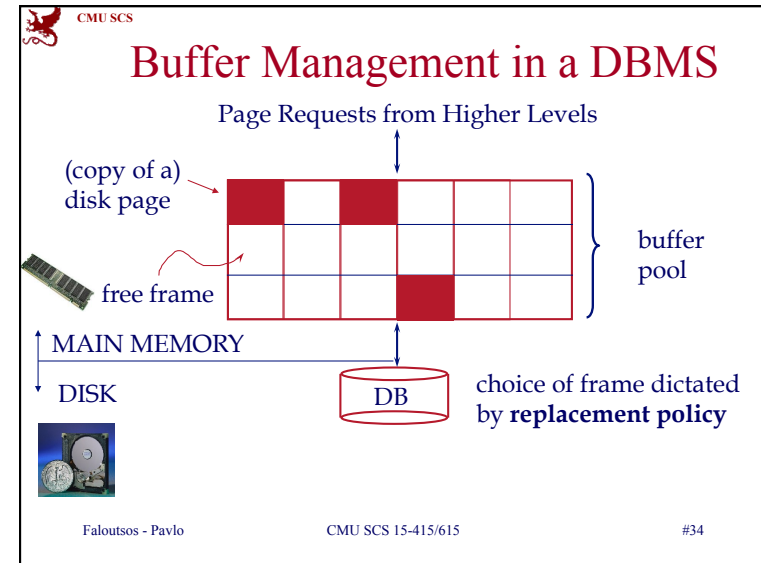
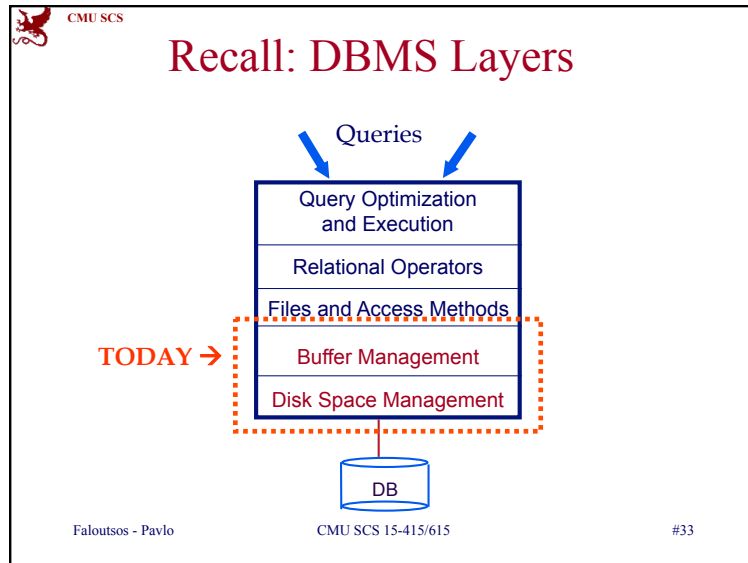
CMU SCS

## Overview

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

Faloutsos - Pavlo                      CMU SCS 15-415/615                      #32






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 - Pavlo      CMU SCS 15-415/615      #37

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 - Pavlo      CMU SCS 15-415/615      #38

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 - Pavlo      CMU SCS 15-415/615      #39

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 - Pavlo      CMU SCS 15-415/615      #40

CMU SCS

## Buffer Replacement Policy

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




Faloutsos - Pavlo CMU SCS 15-415/615 #41

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 - Pavlo CMU SCS 15-415/615 #42

CMU SCS

## 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 - Pavlo CMU SCS 15-415/615 #43

CMU SCS

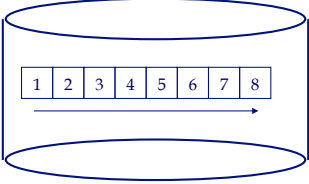
## Sequential Flooding – Illustration

LRU: 

| BUFFER POOL |     |     |     |
|-------------|-----|-----|-----|
| 102         | 116 | 242 | 105 |

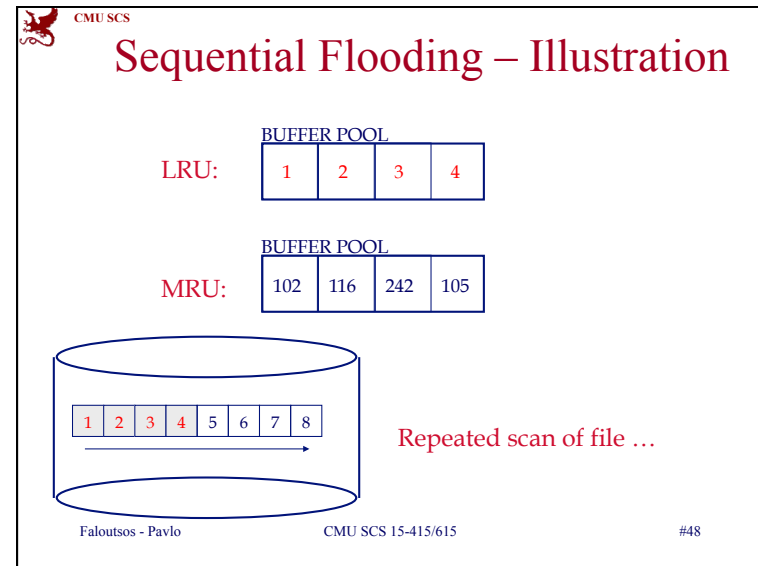
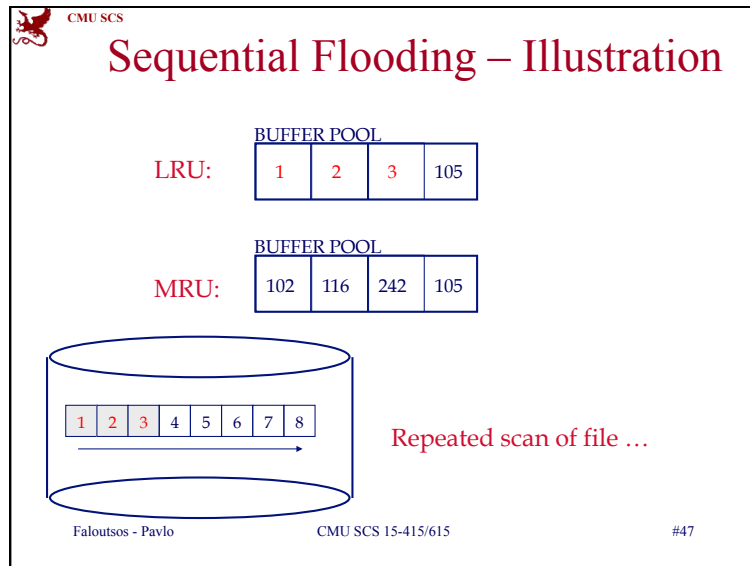
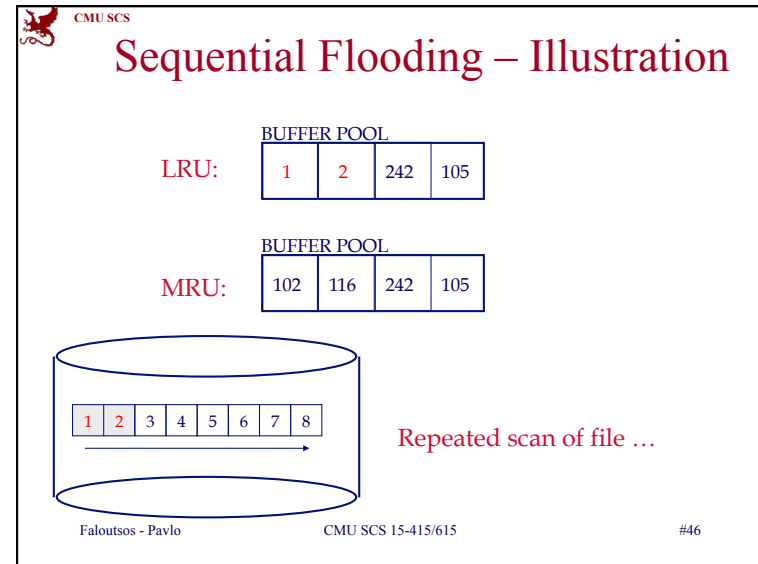
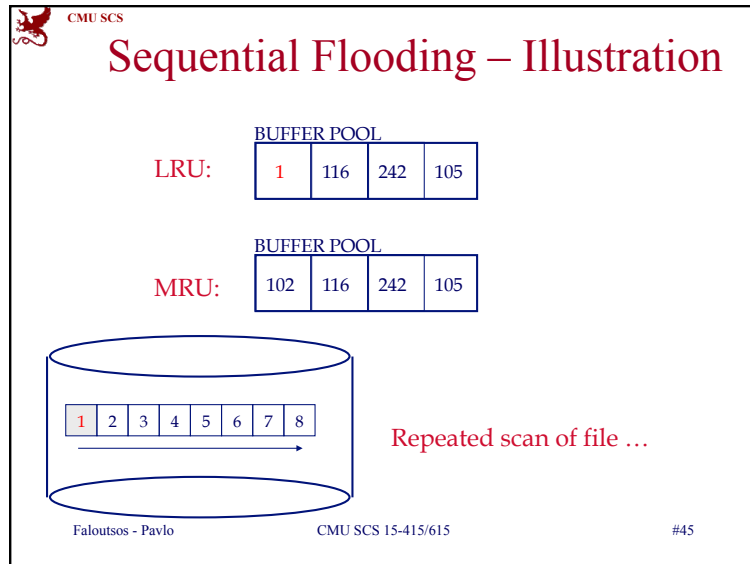
MRU: 

| BUFFER POOL |     |     |     |
|-------------|-----|-----|-----|
| 102         | 116 | 242 | 105 |



Repeated scan of file ...

Faloutsos - Pavlo CMU SCS 15-415/615 #44



CMU SCS

## How will MRU work?

Faloutsos - Pavlo      CMU SCS 15-415/615      #49

CMU SCS

## Sequential Flooding – Illustration

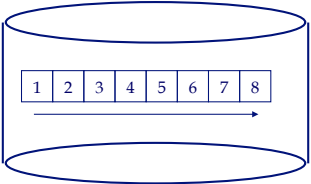
LRU: 

|   |   |   |   |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
|---|---|---|---|

 will not re-use these pages;

MRU: 

|     |     |     |     |
|-----|-----|-----|-----|
| 102 | 116 | 242 | 105 |
|-----|-----|-----|-----|



Faloutsos - Pavlo      CMU SCS 15-415/615      #50

CMU SCS

## Sequential Flooding – Illustration

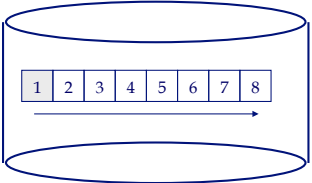
LRU: 

|   |   |   |   |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
|---|---|---|---|

 will not re-use these pages;

MRU: 

|   |     |     |     |
|---|-----|-----|-----|
| 1 | 116 | 242 | 105 |
|---|-----|-----|-----|



Faloutsos - Pavlo      CMU SCS 15-415/615      #51

CMU SCS

## Sequential Flooding – Illustration

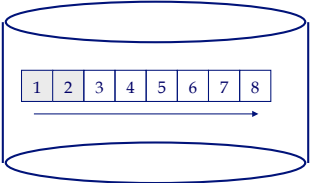
LRU: 

|   |   |   |   |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
|---|---|---|---|

 will not re-use these pages;

MRU: 

|   |     |     |     |
|---|-----|-----|-----|
| 2 | 116 | 242 | 105 |
|---|-----|-----|-----|



Faloutsos - Pavlo      CMU SCS 15-415/615      #52

CMU SCS

## Sequential Flooding – Illustration

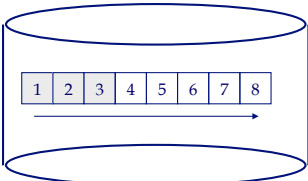
LRU: 

|   |   |   |   |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
|---|---|---|---|

 will not re-use these pages;

MRU: 

|   |     |     |     |
|---|-----|-----|-----|
| 3 | 116 | 242 | 105 |
|---|-----|-----|-----|



Repeated scan of file ...

Faloutsos - Pavlo CMU SCS 15-415/615 #53

CMU SCS

## Sequential Flooding – Illustration

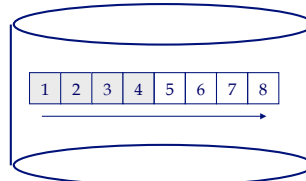
LRU: 

|   |   |   |   |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
|---|---|---|---|

 will not re-use these pages;

MRU: 

|   |     |     |     |
|---|-----|-----|-----|
| 4 | 116 | 242 | 105 |
|---|-----|-----|-----|



Repeated scan of file ...

Faloutsos - Pavlo CMU SCS 15-415/615 #54

CMU SCS

## Other policies?

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

Faloutsos - Pavlo CMU SCS 15-415/615 #55

CMU SCS

## “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 - Pavlo CMU SCS 15-415/615 #56

CMU SCS

## “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 - Pavlo CMU SCS 15-415/615 #57

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 - Pavlo CMU SCS 15-415/615 #58

CMU SCS

## “Clock” Replacement Policy

Faloutsos - Pavlo CMU SCS 15-415/615 #59

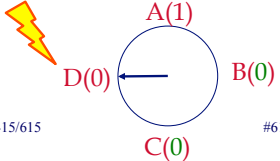
CMU SCS

## “Clock” Replacement Policy

Faloutsos - Pavlo CMU SCS 15-415/615 #60

CMU SCS

## “Clock” Replacement Policy



Faloutsos - Pavlo CMU SCS 15-415/615 #61

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 - Pavlo CMU SCS 15-415/615 #62

CMU SCS

## Overview

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

Faloutsos - Pavlo CMU SCS 15-415/615 #63

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 - Pavlo CMU SCS 15-415/615 #64



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 - Pavlo CMU SCS 15-415/615 #65

CMU SCS

## Files of records

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

Faloutsos - Pavlo CMU SCS 15-415/615 #66

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.

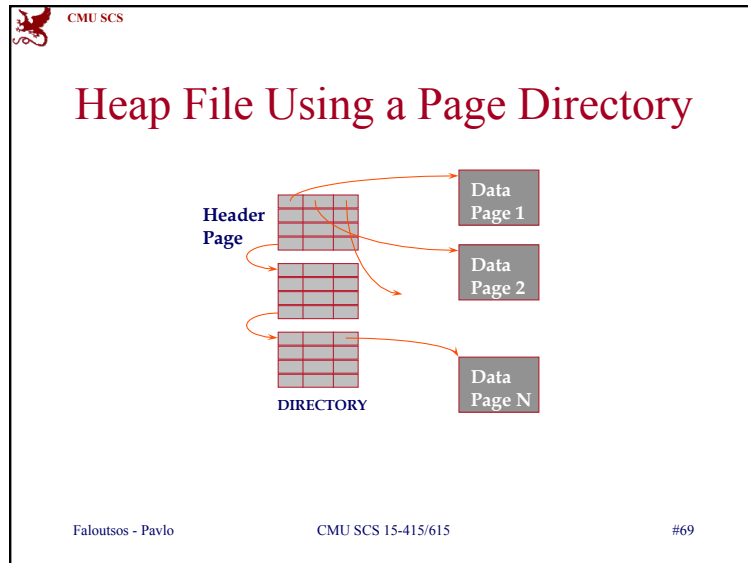
Faloutsos - Pavlo CMU SCS 15-415/615 #67

CMU SCS

## Heap File Using Lists

- Any problems?

Faloutsos - Pavlo CMU SCS 15-415/615 #68



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 - Pavlo

CMU SCS 15-415/615

#70

CMU SCS

## Overview

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

Faloutsos - Pavlo

CMU SCS 15-415/615

#71

CMU SCS

## Page Formats

- fixed length records
- variable length records

Faloutsos - Pavlo

CMU SCS 15-415/615

#72

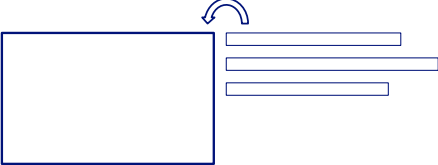
CMU SCS

## Problem definition

Q: How would you store records on a page/  
file, such that

1. you can point to them
2. you can insert/delete records with few disk  
accesses

4kb page



The diagram shows a square representing a 4kb page. To its right, three horizontal bars of varying lengths represent records. A curved arrow points from the top of the page to the top of the first record, indicating a pointer.

Faloutsos - Pavlo CMU SCS 15-415/615 #73

CMU SCS


## 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 - Pavlo CMU SCS 15-415/615 #74

CMU SCS

## 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:

⇒ Q2: Why not its byte offset in the file?

A2:



Faloutsos - Pavlo CMU SCS 15-415/615 #75

CMU SCS

## 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:

Q2: Why not its byte offset in the file?

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



Faloutsos - Pavlo CMU SCS 15-415/615 #76

CMU SCS

## 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.

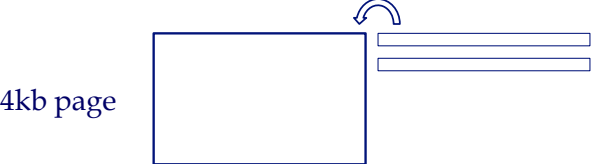


Faloutsos - Pavlo      CMU SCS 15-415/615      #77

CMU SCS

## Fixed length records

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



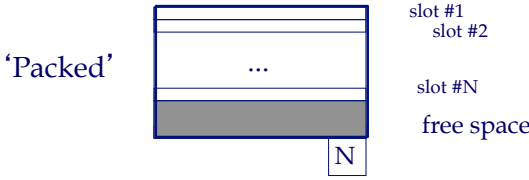
4kb page

Faloutsos - Pavlo      CMU SCS 15-415/615      #78

CMU SCS

## Fixed length records

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



'Packed'

slot #1  
slot #2  
...  
slot #N  
free space

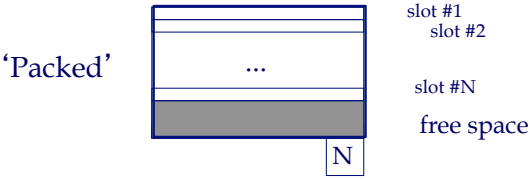
N  
number of full slots

Faloutsos - Pavlo      CMU SCS 15-415/615      #79

CMU SCS

## Fixed length records

- OK – how about insertion?



'Packed'

slot #1  
slot #2  
...  
slot #N  
free space

N  
number of full slots

Faloutsos - Pavlo      CMU SCS 15-415/615      #80

CMU SCS

## Fixed length records

- OK – how about insertion?

'Packed'

slot #1  
slot #2  
...  
slot #N  
free space  
N+1  
number of full slots

Faloutsos - Pavlo      CMU SCS 15-415/615      #81

CMU SCS

## Fixed length records

- How about deletion?

'Packed'

slot #1  
slot #2  
xxxxxxx  
...  
slot #N  
free space  
N  
number of full slots

Faloutsos - Pavlo      CMU SCS 15-415/615      #82

CMU SCS

## Fixed length records

- How about deletion?
- Bad - we have too much to reorganize/update

B-tree

slot #1  
slot #2  
xxxxxxx  
...  
slot #N  
free space  
N  
number of full slots

Faloutsos - Pavlo      CMU SCS 15-415/615      #83

CMU SCS

## Fixed length records

- What would you do?

Faloutsos - Pavlo      CMU SCS 15-415/615      #84

CMU SCS

## Fixed length records

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

free slots

slot #1  
slot #2  
...  
slot #N

page header

Faloutsos - Pavlo      CMU SCS 15-415/615      #85

CMU SCS

## Fixed length records

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

free slots

slot #1  
slot #2  
...  
slot #N

page header

Faloutsos - Pavlo      CMU SCS 15-415/615      #86

CMU SCS

## Variable length records

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

occupied records

page header

Faloutsos - Pavlo      CMU SCS 15-415/615      #87

CMU SCS

## Variable length records

**SLOTTED PAGE**

- Q: How would you store them on a page/file?
  - pack them
  - keep ptrs to them

occupied records

page header

slot directory

other info (# slots etc)

Faloutsos - Pavlo      CMU SCS 15-415/615      #88

CMU SCS **SLOTTED PAGE**

## Variable length records

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

occupied records

- pack them
- keep ptrs to them
- mark start of free space

The diagram shows a page layout. At the top, there are several colored bars representing records of varying lengths. Below these is a grey area representing free space. At the bottom, there is a 'page header' and a 'slot directory' consisting of a row of vertical bars. A double-headed arrow indicates the width of the slot directory. Below the slot directory, it says 'other info (# slots etc)'. A blue arrow points from the slot directory to the records above.

Faloutsos - Pavlo CMU SCS 15-415/615 #89

CMU SCS **SLOTTED PAGE**

## Variable length records

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

occupied records

- how many disk accesses to insert a record?
- to delete one?

This diagram is similar to the previous one but includes blue arrows showing the path from the slot directory to the records, illustrating the disk access process for inserting or deleting records.

Faloutsos - Pavlo CMU SCS 15-415/615 #90

CMU SCS **SLOTTED PAGE**

## Variable length records

- SLOTTED PAGE organization - popular.

occupied records

The diagram shows a simplified view of a slotted page with records at the top and a page header at the bottom. Blue arrows point from the records to the page header.

Faloutsos - Pavlo CMU SCS 15-415/615 #91

CMU SCS

## Overview

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

Faloutsos - Pavlo CMU SCS 15-415/615 #92

CMU SCS

## Formats of records

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

Faloutsos - Pavlo CMU SCS 15-415/615 #93

CMU SCS

## 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 - Pavlo CMU SCS 15-415/615 #94

CMU SCS

## Formats of records

- Fixed length records: straightforward - store info in catalog
- Variable length records: encode the length of each field
  - ?
  - ?

Faloutsos - Pavlo CMU SCS 15-415/615 #95

CMU SCS

## 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 - Pavlo CMU SCS 15-415/615 #96



CMU SCS

## Variable Length records

- Two alternative formats (# fields is fixed):

Fields Delimited by Special Symbols

Array of Field Offsets

Pros and cons?

Faloutsos - Pavlo      CMU SCS 15-415/615      #97

CMU SCS

## Variable Length records

- Two alternative formats (# fields is fixed):

Fields Delimited by Special Symbols

Array of Field Offsets

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

Faloutsos - Pavlo      CMU SCS 15-415/615      #98

CMU SCS

## 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 - Pavlo      CMU SCS 15-415/615      #99