


PURDUE
UNIVERSITY


CS54100: Database Systems

Building a DBMS
10 February 2012
Prof. Chris Clifton

*Slides adapted from those developed by
Stanford University Prof. Hector Garcia-Molina*




Indiana
Center for
Database
Systems



Isn't Implementing a Database System Simple?

Relations \Rightarrow Statements \Rightarrow Results

Spring 2012 Chris Clifton - CS54100 2




Introducing the

MEGATRON 3000

Database Management System

- The latest from Megatron Labs
- Incorporates latest relational technology
- UNIX compatible

Spring 2012 Chris Clifton - CS54100 3



Megatron 3000 Implementation Details

! First sign non-disclosure agreement **!**

Spring 2012 Chris Clifton - CS54100 4



Megatron 3000 Implementation Details

- Relations stored in files (ASCII)
e.g., relation R is in /usr/db/R

```
Smith # 123 # CS
Jones # 522 # EE
⋮
```

Spring 2012

Chris Clifton - CS54100

5



Megatron 3000 Implementation Details

- Directory file (ASCII) in /usr/db/directory

```
R1 # A # INT # B # STR ...
R2 # C # STR # A # INT ...
⋮
```

Spring 2012

Chris Clifton - CS54100

6



Megatron 3000 Sample Sessions

```
% MEGATRON3000
  Welcome to MEGATRON 3000!
&
:
& quit
%
```

Spring 2012

Chris Clifton - CS54100

7



Megatron 3000 Sample Sessions

```
& select *
  from R #

  Relation R
  A      B      C
  SMITH   123    CS

&
```

Spring 2012

Chris Clifton - CS54100

8



Megatron 3000 Sample Sessions

```
& select A,B  
  from R,S  
  where R.A = S.A and S.C > 100 #
```

<u>A</u>	<u>B</u>
123	CAR
522	CAT

```
&
```

Spring 2012

Chris Clifton - CS54100

9



Megatron 3000 Sample Sessions

```
& select *  
  from R | LPR #  
&
```

Result sent to LPR (printer).

Spring 2012

Chris Clifton - CS54100

10



Megatron 3000 Sample Sessions

```
& select *  
  from R  
  where R.A < 100 | T #  
&
```

New relation T created.

Spring 2012

Chris Clifton - CS54100

11



Megatron 3000

- To execute "**select * from R where condition**":
 - (1) Read dictionary to get R attributes
 - (2) Read R file, for each line:
 - (a) Check condition
 - (b) If OK, display

Spring 2012

Chris Clifton - CS54100

12



Megatron 3000

- To execute “**select * from R where condition | T**”:
 - (1) Process select as before
 - (2) Write results to new file T
 - (3) Append new line to dictionary

Spring 2012

Chris Clifton - CS54100

13



Megatron 3000

- To execute “**select A,B from R,S where condition**”:
 - (1) Read dictionary to get R,S attributes
 - (2) Read R file, for each line:
 - (a) Read S file, for each line:
 - (i) Create join tuple
 - (ii) Check condition
 - (iii) Display if OK

Spring 2012

Chris Clifton - CS54100

14



What's wrong with the Megatron 3000 DBMS?

- Tuple layout on disk
 - Change string from 'Cat' to 'Cats' and we have to rewrite file
 - ASCII storage is expensive
 - Deletions are expensive
- Search expensive; no indexes
 - Cannot find tuple with given key quickly
 - Always have to read full relation

Spring 2012

Chris Clifton - CS54100

15



What's wrong with the Megatron 3000 DBMS?

- No buffer manager
 - Need caching
- Brute force query processing
 - ```
select *
 from R,S
 where R.A = S.A and S.B > 1000
```
  - Do select first?
  - More efficient join?

Spring 2012

Chris Clifton - CS54100

16





## What's wrong with the Megatron 3000 DBMS?

- No concurrency control
- No reliability
  - Can lose data
  - Can leave operations half done
- No security
  - File system insecure
  - File system security is coarse

Spring 2012

Chris Clifton - CS54100

17



## What's wrong with the Megatron 3000 DBMS?

- No application program interface (API)
  - How can a payroll program get at the data?
- No GUI
- Cannot interact with other DBMSs.
- Poor dictionary facilities
  - How do we know what is in the database?
- Lousy salesman!!

Spring 2012

Chris Clifton - CS54100

18



## What do we need to know?

---

- **File & System Structure**  
Records in blocks, dictionary, buffer management,...
- **Indexing & Hashing**  
B-Trees, hashing,...
- **Query Processing**  
Query costs, join strategies,...
- **Crash Recovery**  
Failures, stable storage,...

Spring 2012

Chris Clifton - CS54100

19



## What do we need to know?

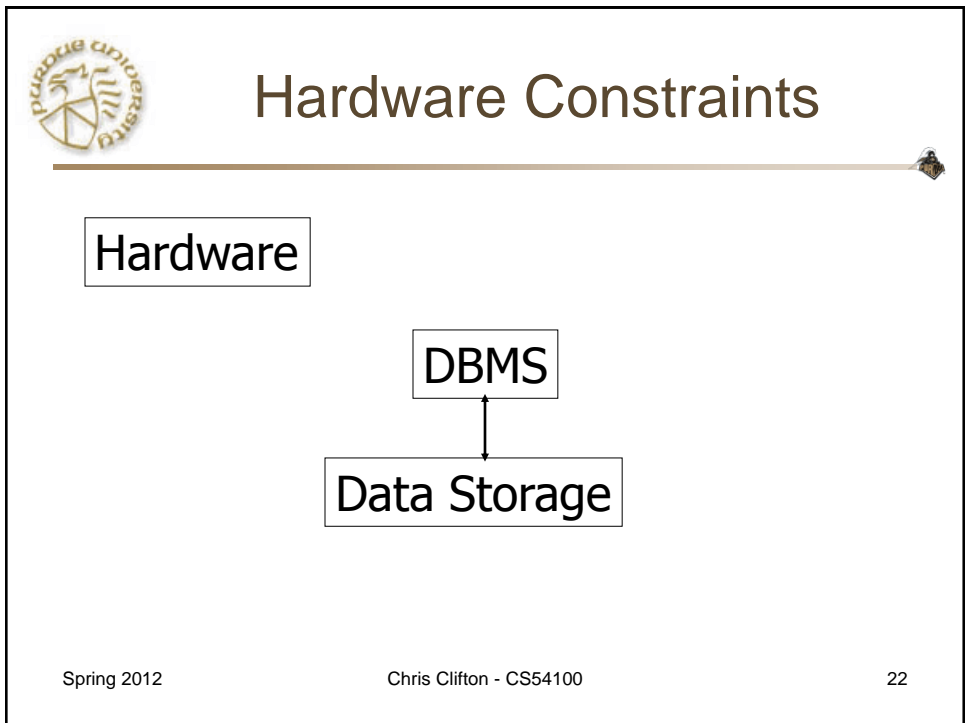
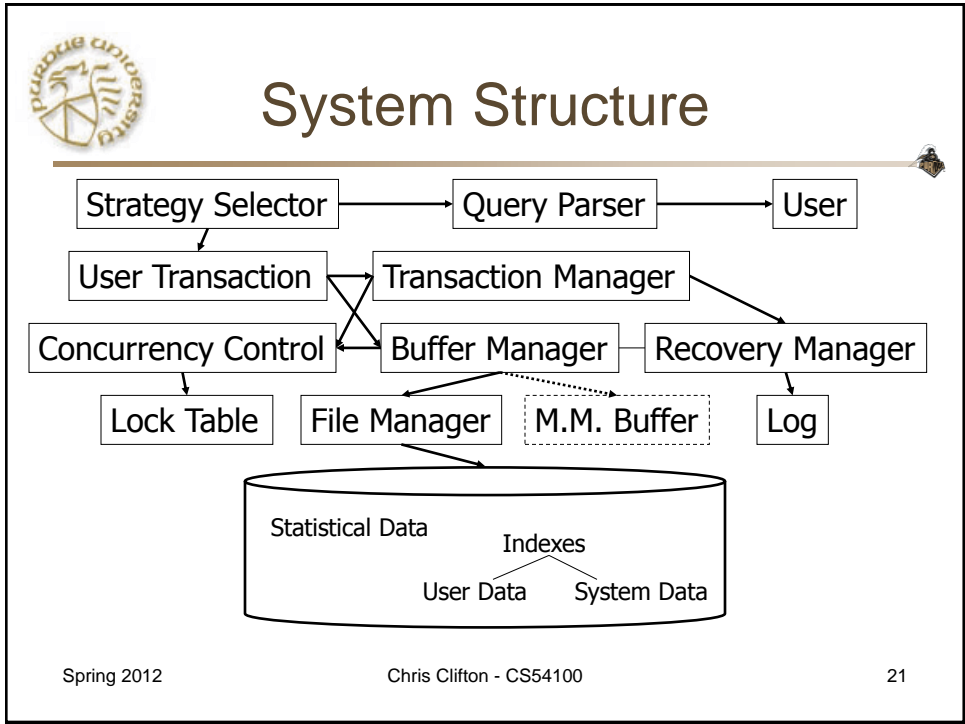
---

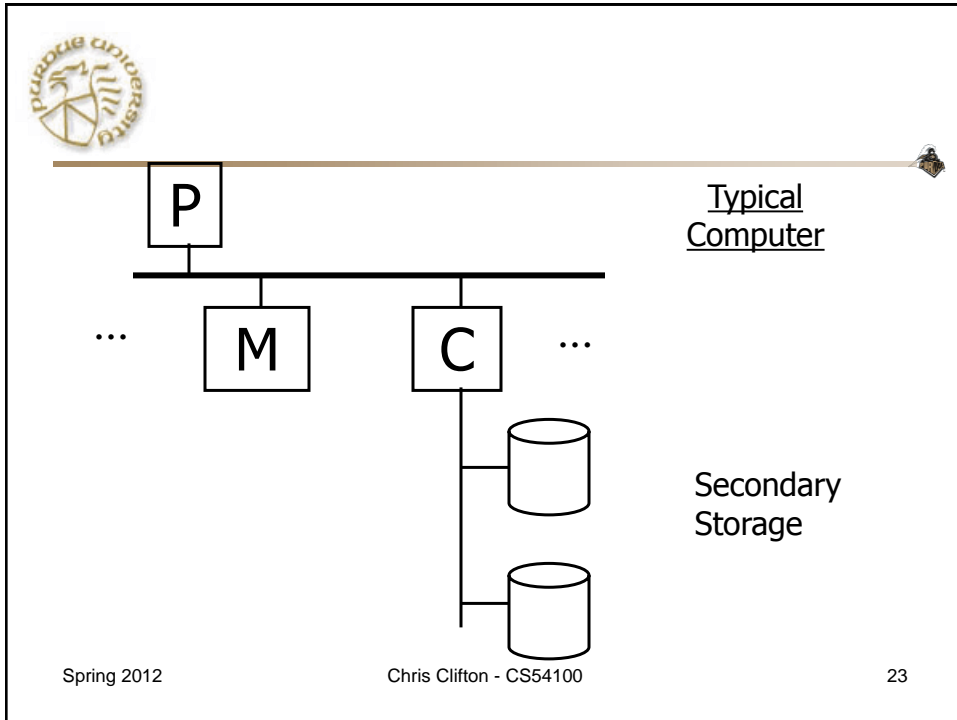
- **Concurrency Control**  
Correctness, locks,...
- **Transaction Processing**  
Logs, deadlocks,...
- **Security & Integrity**  
Authorization, encryption,...
- **Distributed Databases**  
Interoperation, distributed recovery,...

Spring 2012

Chris Clifton - CS54100

20





The diagram illustrates a typical computer architecture. At the top, a horizontal line represents the system bus. A box labeled 'P' (Processor) is connected to this bus. Below the bus, a vertical line connects to a horizontal line representing the memory bus. From this horizontal line, three vertical lines lead to boxes labeled 'M' (Memory) and 'C' (Cache), with ellipses on either side indicating other components. Below the 'C' box, a vertical line connects to two cylindrical icons representing secondary storage (disks).

Processor  
Fast, slow, reduced instruction set,  
with cache, pipelined...  
Speed: 100 → 500 → 1000 MIPS

Memory  
Fast, slow, non-volatile, read-only,...  
Access time:  $10^{-6}$  →  $10^{-9}$  sec.  
1  $\mu$ s → 1 ns

Spring 2012

Chris Clifton - CS54100

24

## Secondary storage

Many flavors:

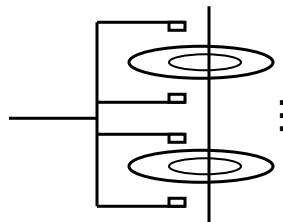
- Disk: Floppy (hard, soft)  
Removable Packs  
Winchester  
Ram disks  
Optical, CD-ROM...  
Arrays
- Tape Reel, cartridge  
Robots

Spring 2012

Chris Clifton - CS54100

25

## Focus on: "Typical Disk"

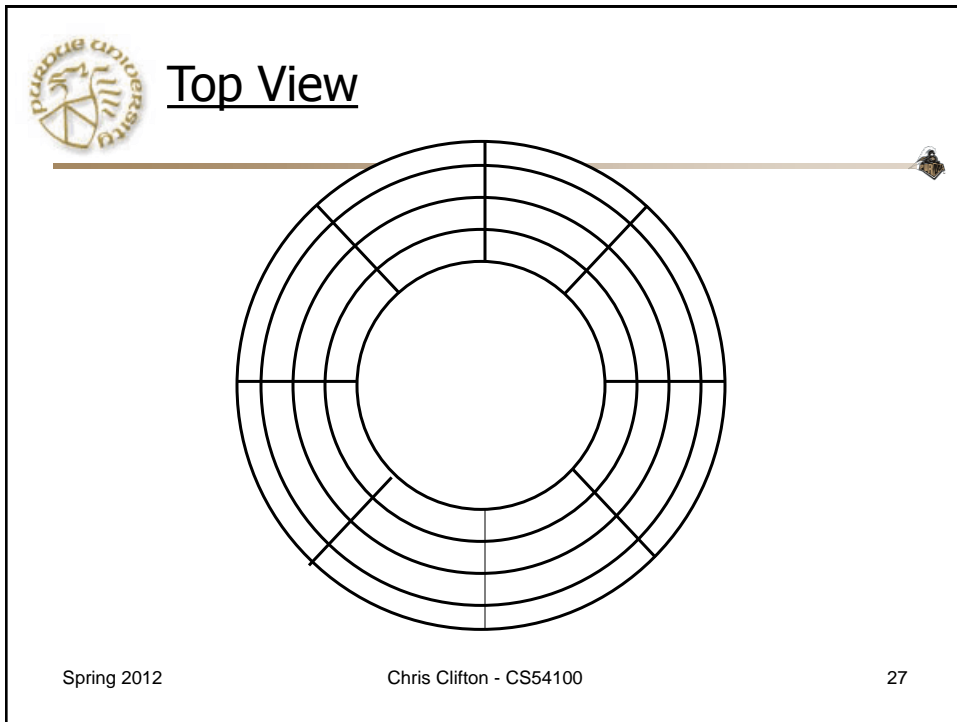


Terms: Platter, Head, Actuator  
Cylinder, Track  
Sector (physical),  
Block (logical), Gap

Spring 2012

Chris Clifton - CS54100


26



“Typical” Numbers

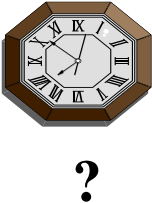
|              |                     |
|--------------|---------------------|
| Diameter:    | 1 inch → 15 inches  |
| Cylinders:   | 100 → 2000          |
| Surfaces:    | 1 (CDs) →           |
| (Tracks/cyl) | 2 (floppies) → 30   |
| Sector Size: | 512B → 50K          |
| Capacity:    | 360 KB (old floppy) |
|              | → TB                |

Spring 2012 Chris Clifton - CS54100 28




---

## Disk Access Time

I want block X →  → block x in memory

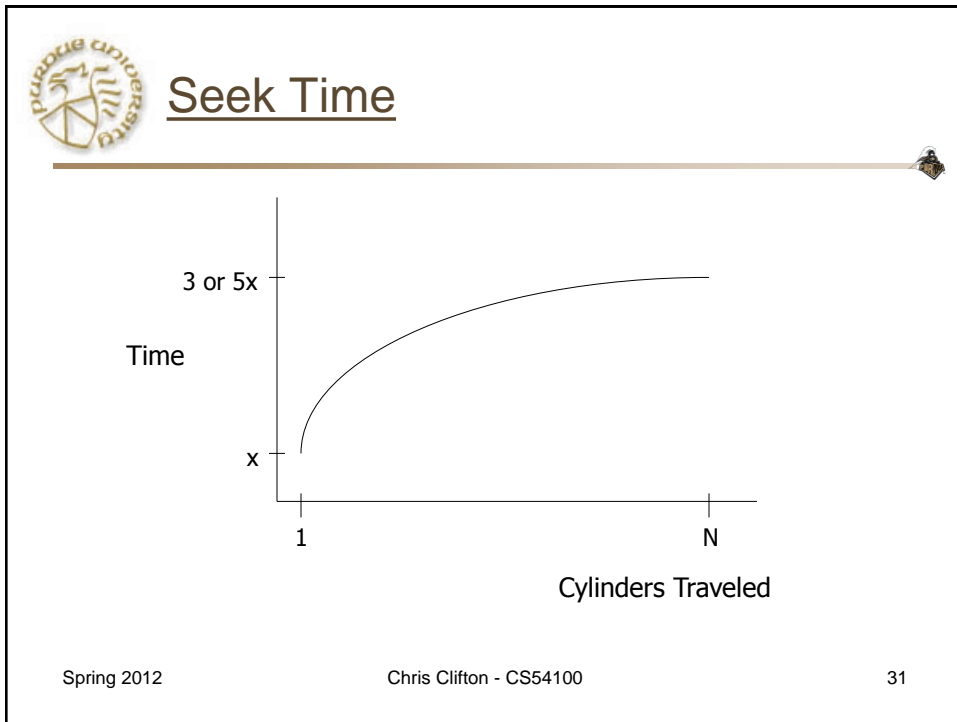
Spring 2012 Chris Clifton - CS54100 29



---

Time = Seek Time +  
Rotational Delay +  
Transfer Time +  
Other

Spring 2012 Chris Clifton - CS54100 30




**Average Random Seek Time**

$$S = \frac{\sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N \text{SEEKTIME}(i \rightarrow j)}{N(N-1)}$$

“Typical” S: 5 ms  $\rightarrow$  10 ms

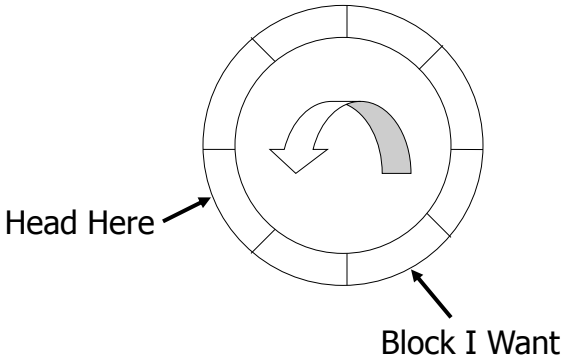
Spring 2012 Chris Clifton - CS54100 32





## Rotational Delay

---




Head Here

Block I Want

Spring 2012

Chris Clifton - CS54100

33



## Average Rotational Delay

---

$R = 1/2$  revolution

"typical"  $R = 2$  ms (15000 RPM)

Spring 2012

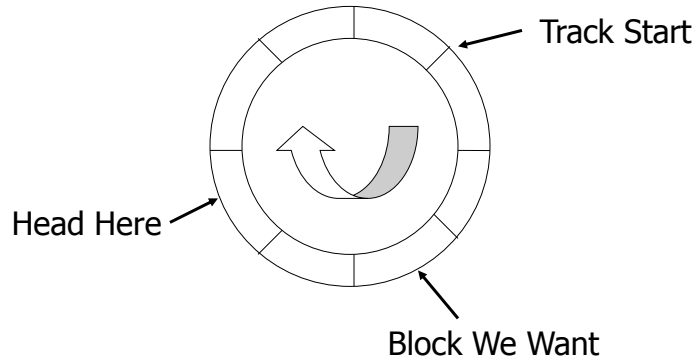
Chris Clifton - CS54100

34



## Complication

- May have to wait for start of track before we can read desired block



Spring 2012

Chris Clifton - CS54100

35



## Transfer Rate: $t$

- "typical"  $t$ : 1 → 3 MB/second
- transfer time: block size

$$t$$

Spring 2012

Chris Clifton - CS54100

36



## Other Delays

- CPU time to issue I/O
- Contention for controller
- Contention for bus, memory

“Typical” Value: 0

Spring 2012

Chris Clifton - CS54100

37



- So far: Random Block Access
- What about: Reading “Next” block?

Spring 2012

Chris Clifton - CS54100

38



## If we do things right

(e.g., Double Buffer, Stagger Blocks...)

$$\text{Time to get block} = \frac{\text{Block Size}}{t} + \text{Negligible}$$



- skip gap
- switch track
- once in a while,  
next cylinder

Spring 2012

Chris Clifton - CS54100

39



## Rule of Thumb

Random I/O: Expensive  
Sequential I/O: Much less

- Ex: 1 KB Block
  - » Random I/O: ~ 10 ms.
  - » Sequential I/O: ~ 1 ms.

Spring 2012

Chris Clifton - CS54100

40



## Cost for Writing similar to Reading

---

.... unless we want to verify!  
need to add (full) rotation + Block size  
\_\_\_\_\_ t

Spring 2012

Chris Clifton - CS54100

41



## To Modify a Block?

---

### To Modify Block:

- (a) Read Block
- (b) Modify in Memory
- (c) Write Block
- [(d) Verify?]

Spring 2012

Chris Clifton - CS54100

42



## Block Address:

- Physical Device
- Cylinder #
- Surface #
- Sector

Spring 2012

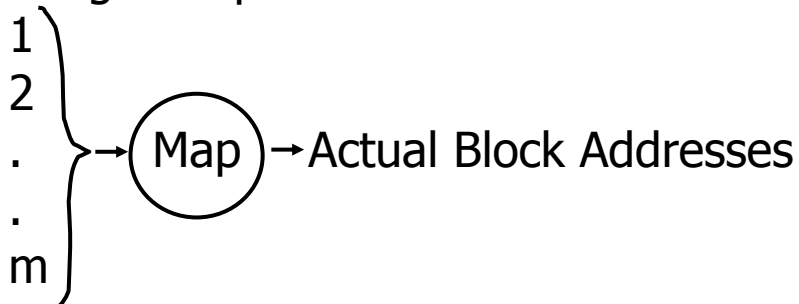
Chris Clifton - CS54100

43



## Complication: Bad Blocks

- Messy to handle
- May map via software to integer sequence



Spring 2012

Chris Clifton - CS54100

44