

---

# EECS 252 Graduate Computer Architecture

## Lec 18 – Storage

David Patterson

Electrical Engineering and Computer Sciences  
University of California, Berkeley

<http://www.eecs.berkeley.edu/~patterson>  
<http://vlsci.cs.berkeley.edu/ee252-s08>

---

## Review

- **Virtual Machine Revival**
  - Overcome security flaws of modern O/Ses
  - Processor performance no longer highest priority
  - Manage Software, Manage Hardware
- “... VMMs give OS developers another opportunity to develop functionality no longer practical in today’s complex and ossified operating systems, where innovation moves at geologic pace.”  
[Rosenblum and Garfinkel, 2006]
- **Virtualization challenges for processor, virtual memory, I/O**
  - Paravirtualization, ISA upgrades to cope with those difficulties
- **Xen as example VMM using paravirtualization**
  - 2006 performance on non-I/O bound, I/O intensive apps: 80% of native Linux without driver VM, 84% with driver VM
- **Opteron memory hierarchy still critical to performance**

4/12/2006

CS252 s08 Storage

2

---

## Case for Storage

- **Shift in focus from computation to communication and storage of information**
  - E.g., Cray Research/Thinking Machines vs. Google/Yahoo
  - “The Computing Revolution” (1980s to 1990s)
  - “The Information Age” (1990 to today)
- **Storage emphasizes reliability and scalability as well as cost-performance**
- **What is “Software king” that determines which HW actually features used?**
  - Operating System for storage
  - Compiler for processor
- **Also has own performance theory—queuing theory—balances throughput vs. response time**

4/12/2006

CS252 s08 Storage

3

---

## Outline

- **Magnetic Disks**
- **RAID**
- **Administrivia**
- **Advanced Dependability/Reliability/Availability**
- **I/O Benchmarks, Performance and Dependability**
- **Intro to Queuing Theory (if we have time)**
- **Conclusion**

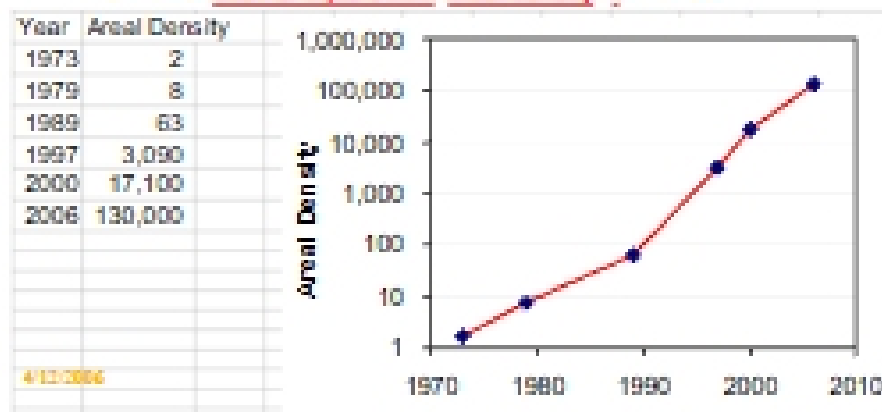
4/12/2006

CS252 s08 Storage

4

## Disk Figure of Merit: Areal Density

- Bits recorded along a track
  - Metric is **Bits/Per/Inch (BPI)**
- Number of tracks per surface
  - Metric is **Tracks/Per/Inch (TPI)**
- Disk Designs Brag about **bit density per unit area**
  - Metric is **Bits/Per/Square/Inch: Areal Density = BPI x TPI**



4/12/2006

1970 1980 1990 2000 2010

## Historical Perspective

- 1956 IBM Ramac — early 1970s Winchester
  - Developed for mainframe computers, proprietary interfaces
  - Steady shrink in form factor: 27 in. to 14 in.
- Form factor and capacity drives market more than performance
- 1970s developments
  - 5.25 inch floppy disk formfactor (microcode into mainframe)
  - Emergence of industry standard disk interfaces
- Early 1980s: PCs and first generation workstations
- Mid 1980s: Client/server computing
  - Centralized storage on file server
    - o accelerated disk downsizing: 8 inch to 5.25
  - Mass market disk drives become a reality
    - o industry standards: SCSI, IDE, IDE
    - o 5.25 inch to 3.5 inch drives for PCs, end of proprietary interfaces
- 1990s: Laptops => 2.5 inch drives
- 2000s: What new devices leading to new drives?

4/12/2006

CS252 w/6 Storage

6

## Future Disk Size and Performance

- Continued advance in capacity (60%/yr) and bandwidth (40%/yr)
- Slow improvement in seek, rotation (8%/yr)
- Time to read whole disk

Year	Sequentially	Randomly (1 sector/seek)
1990	4 minutes	6 hours
2000	12 minutes	1 week(!)
2006	56 minutes	3 weeks (SCSI)
2006	171 minutes	7 weeks (SATA)

4/12/2006

CS252 w/6 Storage

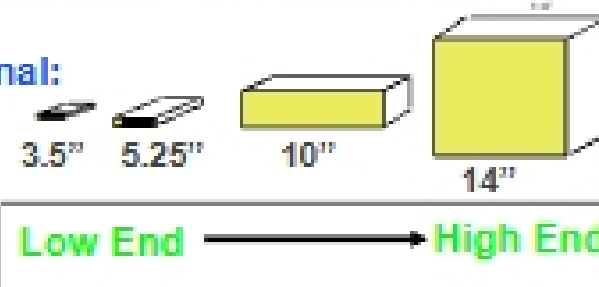
7

## Use Arrays of Small Disks?

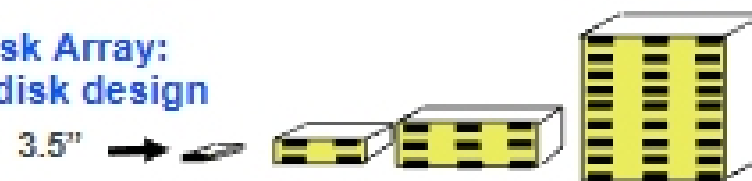
### •Katz and Patterson asked in 1987:

- Can smaller disks be used to close gap in performance between disks and CPUs?

Conventional:  
4 disk designs



Disk Array:  
1 disk design

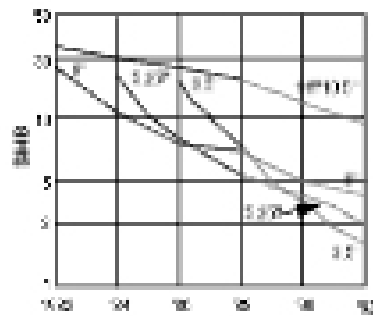


4/12/2006

CS252 w/6 Storage

8

## Advantages of Small Formfactor Disk Drives

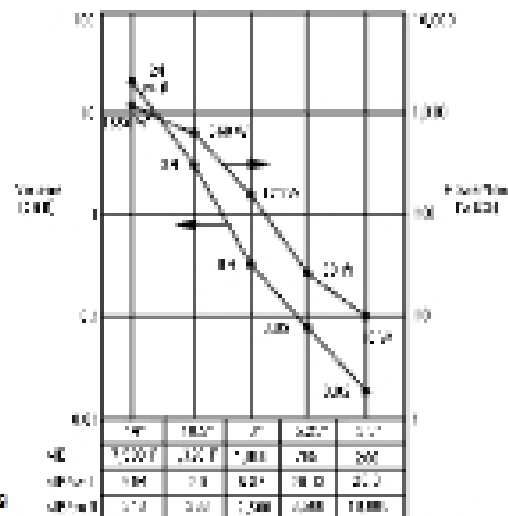


Low cost/MB  
High MB/volume  
High MB/watt  
Low cost/Actuator

Cost and Environmental Efficiencies

4/12/2006

CS252 s06 Storage



Year	1982	1987	1992	1997	2002
MB/KB	10000	10000	10000	10000	10000
MB/W	100	100	100	100	100
MB/KW	100	100	100	100	100

9

## Replace Small Number of Large Disks with Large Number of Small Disks! (1988 Disks)

	IBM 3390K	IBM 3.5" 0061	x70
Capacity	20 GBytes	320 MBytes	23 GBytes
Volume	97 cu. ft.	0.1 cu. ft.	11 cu. ft. <b>9X</b>
Power	3 KW	11 W	1 KW <b>3X</b>
Data Rate	15 MB/s	1.5 MB/s	120 MB/s <b>8X</b>
I/O Rate	600 I/Os/s	55 I/Os/s	3900 I/Os/s <b>6X</b>
MTTF	250 KHrs	50 KHrs	??? Hrs
Cost	\$250K	\$2K	\$150K

Disk Arrays have potential for large data and I/O rates, high MB per cu. ft., high MB per KW, but what about reliability?

10

## Array Reliability

- Reliability of N disks = Reliability of 1 Disk ÷ N

50,000 Hours ÷ 70 disks = 700 hours

Disk system MTTF: Drops from 8 years to 1 month!

- Arrays (without redundancy) too unreliable to be useful!

Hot spares support reconstruction in parallel with access: very high media availability can be achieved

4/12/2006

CS252 s06 Storage

11

## Redundant Arrays of (Inexpensive) Disks

- Files are "striped" across multiple disks
- Redundancy yields high data availability
  - **Availability**: service still provided to user, even if some components failed
- Disks will still fail
- Contents reconstructed from data redundantly stored in the array
  - ⇒ Capacity penalty to store redundant info
  - ⇒ Bandwidth penalty to update redundant info

4/12/2006

CS252 s06 Storage

12