

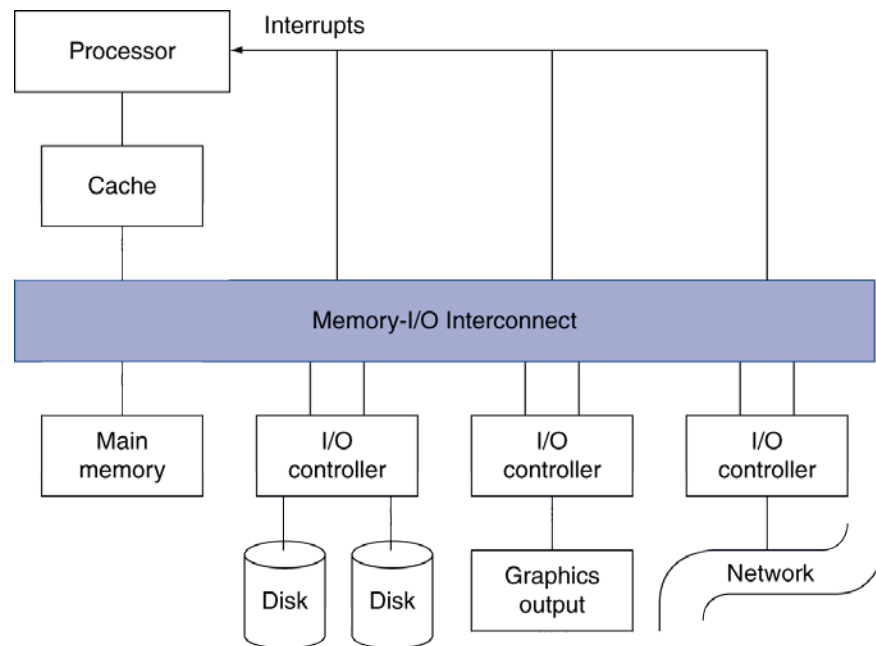


Chapter 6

Storage and Other I/O
Topics

Introduction

- I/O devices can be characterized by
 - Behaviour: input, output, storage
 - Partner: human or machine
 - Data rate: bytes/sec, transfers/sec
- I/O bus connections



I/O System Characteristics

- Performance measures
 - Latency (response time)
 - Throughput (bandwidth)
 - Desktops & embedded systems
 - Mainly interested in response time & diversity of devices
 - Servers
 - Mainly interested in throughput & expandability of devices
- Dependability is important
 - Particularly for storage devices



Diversity of I/O devices

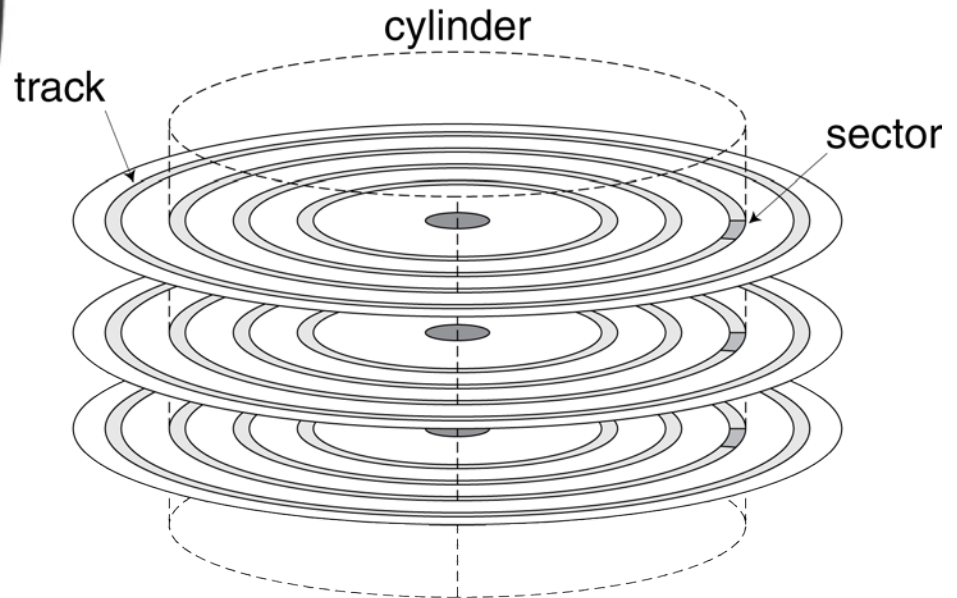
Device	Behavior	Partner	Data rate (Mbit/sec)
Keyboard	Input	Human	0.0001
Mouse	Input	Human	0.0038
Voice input	Input	Human	0.2640
Sound input	Input	Machine	3.0000
Scanner	Input	Human	3.2000
Voice output	Output	Human	0.2640
Sound output	Output	Human	8.0000
Laser printer	Output	Human	3.2000
Graphics display	Output	Human	800.0000–8000.0000
Cable modem	Input or output	Machine	0.1280–6.0000
Network/LAN	Input or output	Machine	100.0000–10000.0000
Network/wireless LAN	Input or output	Machine	11.0000–54.0000
Optical disk	Storage	Machine	80.0000–220.0000
Magnetic tape	Storage	Machine	5.0000–120.0000
Flash memory	Storage	Machine	32.0000–200.0000
Magnetic disk	Storage	Machine	800.0000–3000.0000

Embedded systems have even more diversity of devices:
sensors, actuators, etc.



Disk Storage

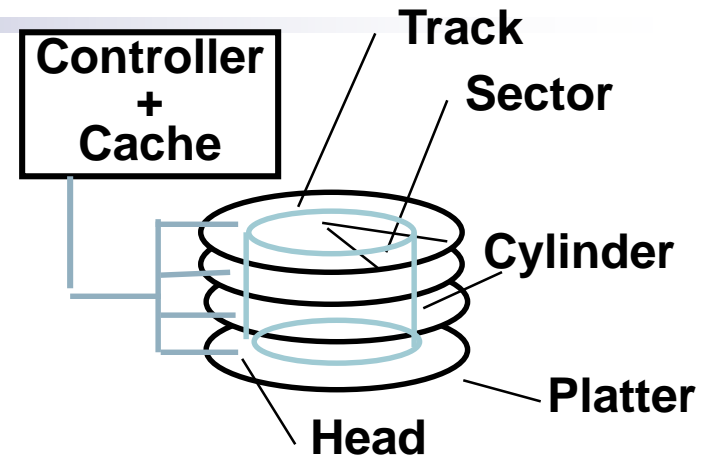
- Nonvolatile, rotating magnetic storage
 - Vary in #platters, rotational speed, density



Magnetic Disk Characteristics

■ Disk read/write components

1. **Seek time:** position the head over the proper track (3 to 13 ms avg)
 - due to locality of disk references the actual average seek time may be only 25% to 33% of the advertised number
2. **Rotational latency:** wait for the desired sector to rotate under the head ($\frac{1}{2}$ of $1/\text{RPM}$ converted to ms)
 - $0.5/5400\text{RPM} = 5.6\text{ms}$ to $0.5/15000\text{RPM} = 2.0\text{ms}$
3. **Transfer time:** transfer a block of bits (one or more sectors) under the head to the disk controller's cache (70 to 125 MB/s are typical disk transfer rates in 2008)
 - the disk controller's "cache" takes advantage of spatial locality in disk accesses
 - cache transfer rates are much faster (e.g., 375 MB/s)
4. **Controller time:** the overhead the disk controller imposes in performing a disk I/O access (typically $< .2$ ms)



Disk Sectors and Access

- Each sector records
 - Sector ID
 - Data (512 bytes, 4096 bytes proposed)
 - Error correcting code (ECC)
 - Used to hide defects and recording errors
 - Synchronization fields and gaps
- Access to a sector involves
 - Queuing delay if other accesses are pending
 - Seek: move the heads
 - Rotational latency
 - Data transfer
 - Controller overhead



Disk Access Example

- Given
 - 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk
- Average read time
 - 4ms seek time
 - + $\frac{1}{2} / (15,000/60) = 2\text{ms}$ rotational latency
 - + $512 / 100\text{MB/s} = 0.005\text{ms}$ transfer time
 - + 0.2ms controller delay
 - = 6.2ms
- If actual average seek time is 1ms
 - Average read time = 3.2ms

Disk Performance Issues

- Manufacturers quote average seek time
 - Based on all possible seeks
 - Locality and OS scheduling lead to smaller actual average seek times
- Smart disk controller allocate physical sectors on disk
 - Present logical sector interface to host
 - ATA/IDE, SATA, SCSI/SAS (serial SCSI)
 - Disk drives include caches
 - Prefetch sectors in anticipation of access
 - Avoid seek and rotational delay



Disk drive characteristics

Characteristics	Seagate ST33000655SS	Seagate ST31000340NS	Seagate ST973451SS	Seagate ST9160821AS
Disk diameter (inches)	3.50	3.50	2.50	2.50
Formatted data capacity (GB)	147	1000	73	160
Number of disk surfaces (heads)	2	4	2	2
Rotation speed (RPM)	15,000	7200	15,000	5400
Internal disk cache size (MB)	16	32	16	8
External interface, bandwidth (MB/sec)	SAS, 375	SATA, 375	SAS, 375	SATA, 150
Sustained transfer rate (MB/sec)	73–125	105	79–112	44
Minimum seek (read/write) (ms)	0.2/0.4	0.8/1.0	0.2/0.4	1.5/2.0
Average seek read/write (ms)	3.5/4.0	8.5/9.5	2.9/3.3	12.5/13.0
Mean time to failure (MTTF) (hours)	1,400,000 @ 25°C	1,200,000 @ 25°C	1,600,000 @ 25°C	—
Annual failure rate (AFR) (percent)	0.62%	0.73%	0.55%	—
Contact start-stop cycles	—	50,000	—	>600,000
Warranty (years)	5	5	5	5
Nonrecoverable read errors per bits read	<1 sector per 10 ¹⁶	<1 sector per 10 ¹⁵	<1 sector per 10 ¹⁶	<1 sector per 10 ¹⁴
Temperature, shock (operating)	5°–55°C, 60 G	5°–55°C, 63 G	5°–55°C, 60 G	0°–60°C, 350 G
Size: dimensions (in.), weight (pounds)	1.0" × 4.0" × 5.8", 1.5 lbs	1.0" × 4.0" × 5.8", 1.4 lbs	0.6" × 2.8" × 3.9", 0.5 lbs	0.4" × 2.8" × 3.9", 0.2 lbs
Power: operating/idle/standby (watts)	15/11/—	11/8/1	8/5.8/—	1.9/0.6/0.2
GB/cu. in., GB/watt	6 GB/cu.in., 10 GB/W	43 GB/cu.in., 91 GB/W	11 GB/cu.in., 9 GB/W	37 GB/cu.in., 84 GB/W
Price in 2008, \$/GB	~ \$250, ~ \$1.70/GB	~ \$275, ~ \$0.30/GB	~ \$350, ~ \$5.00/GB	~ \$100, ~ \$0.60/GB

Laptops

(others for servers)



Disk Latency & Bandwidth Milestones

	CDC Wren	SG ST41	SG ST15	SG ST39	SG ST37
RSpeed (RPM)	3600	5400	7200	10000	15000
Year	1983	1990	1994	1998	2003
Capacity (Gbytes)	0.03	1.4	4.3	9.1	73.4
Diameter (inches)	5.25	5.25	3.5	3.0	2.5
Interface	ST-412	SCSI	SCSI	SCSI	SCSI
Bandwidth (MB/s)	0.6	4	9	24	86
Latency (msec)	48.3	17.1	12.7	8.8	5.7

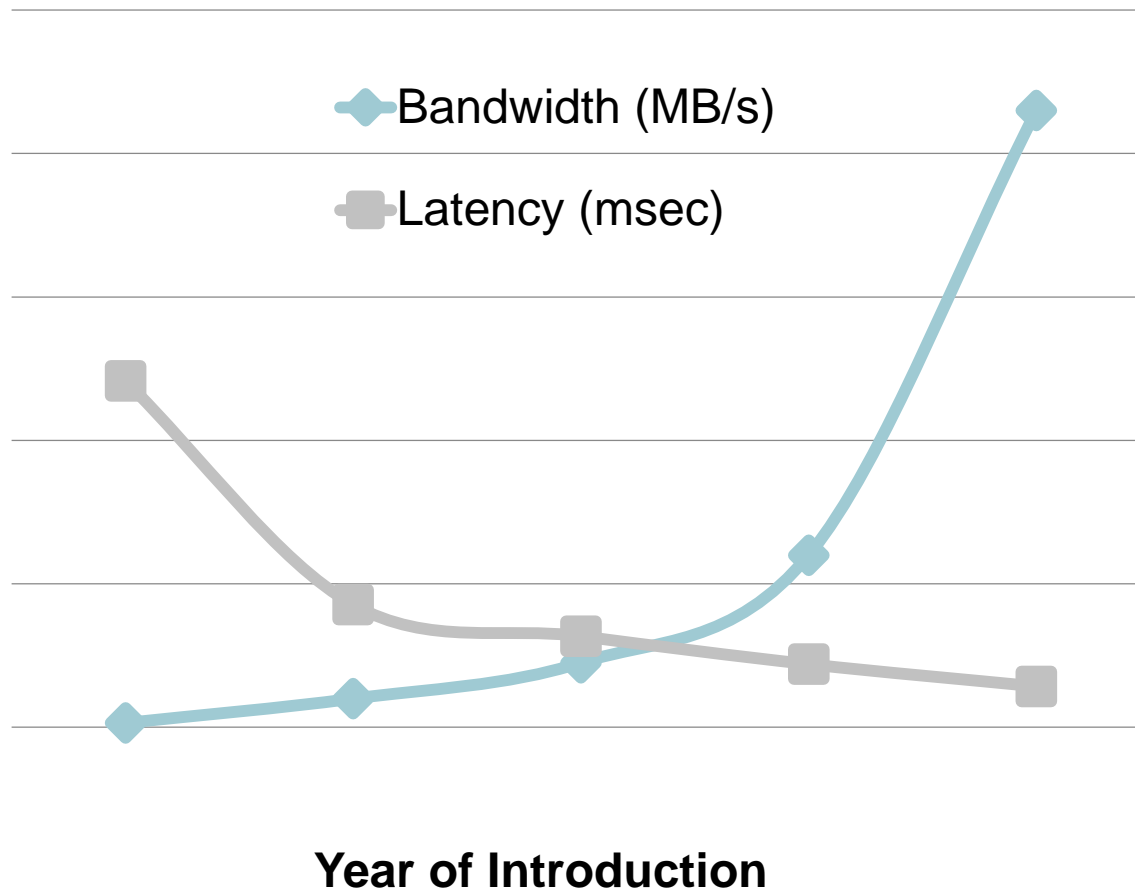
Patterson, CACM Vol 47, #10, 2004

- Disk **latency** is one average seek time plus the rotational latency.
- Disk **bandwidth** is the peak transfer time of formatted data from the media (not from the cache).



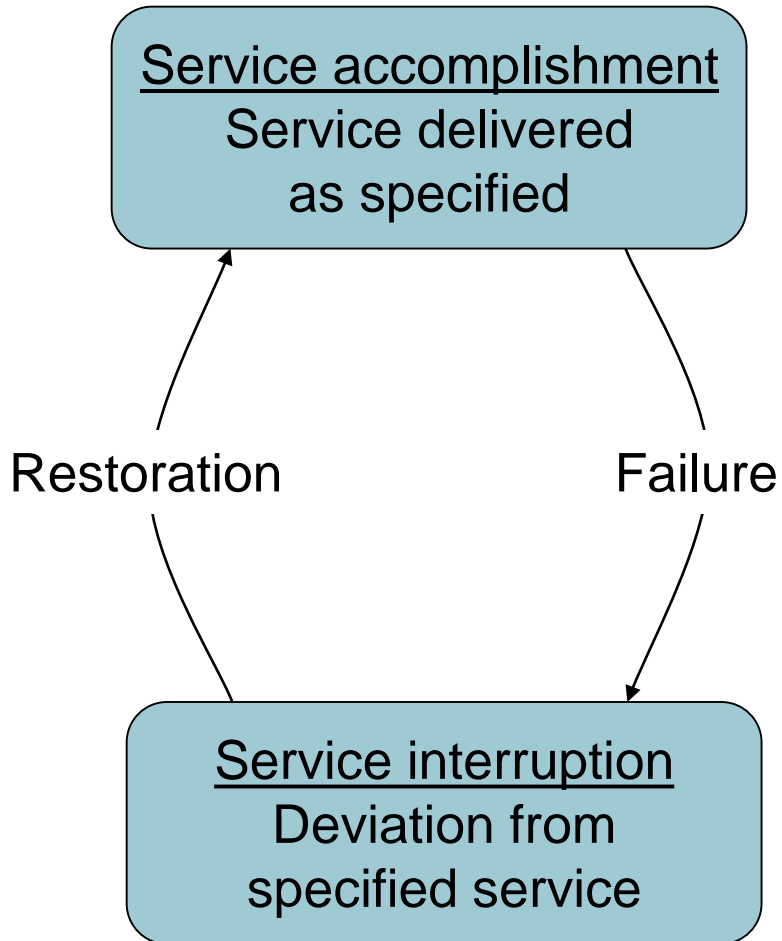
Latency & Bandwidth Improvements

- In the time that the disk **bandwidth doubles** the **latency** improves by a factor of only **1.2 to 1.4**



Mary Jane Irwin,
PSU, 2008

Dependability



- Fault: failure of a component
 - May or may not lead to system failure

Dependability Measures

- Reliability: mean time to failure (MTTF)
- Service interruption: mean time to repair (MTTR)
- Mean time between failures
 - $MTBF = MTTF + MTTR$
- $Availability = MTTF / (MTTF + MTTR)$
- Improving Availability
 - Increase MTTF: fault avoidance, fault tolerance, fault forecasting
 - Reduce MTTR: improved tools and processes for diagnosis and repair

RAID

- Redundant Array of Inexpensive (Independent) Disks
 - Use multiple smaller disks (c.f. one large disk)
 - Parallelism improves performance
 - Plus extra disk(s) for redundant data storage
- Provides fault tolerant storage system
 - Especially if failed disks can be “hot swapped”
- RAID 0
 - No redundancy (“AID”?)
 - Just stripe data over multiple disks
 - But it does improve performance

RAID 1 & 2

- RAID 1: Mirroring
 - N + N disks, replicate data
 - Write data to both data disk and mirror disk
 - On disk failure, read from mirror
- RAID 2: Error correcting code (ECC)
 - N + E disks (e.g., 10 + 4)
 - Split data at bit level across N disks
 - Generate E-bit ECC
 - Too complex, not used in practice

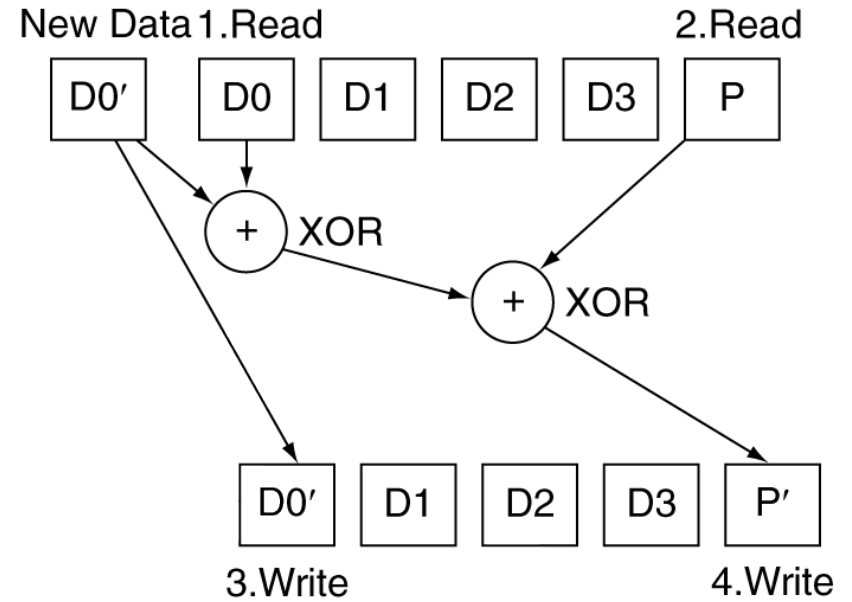
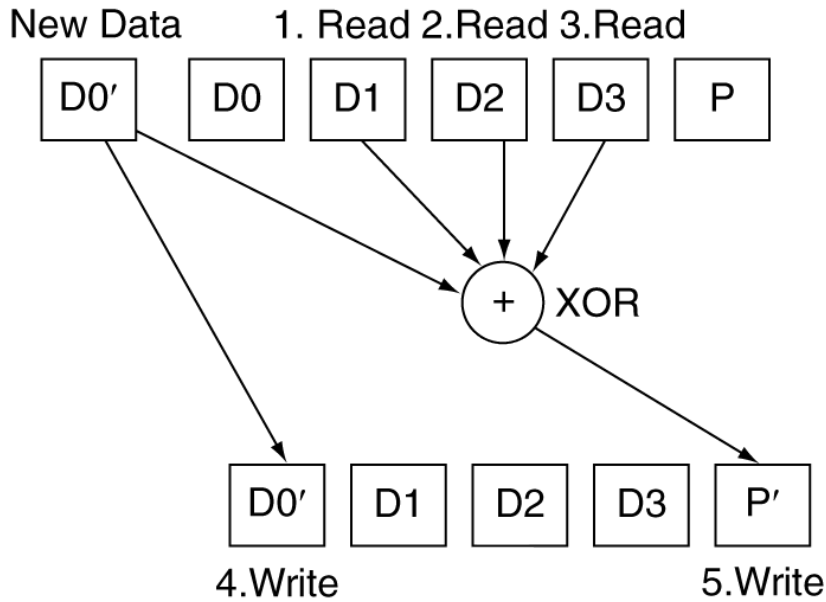
RAID 3: Bit-Interleaved Parity

- N + 1 disks
 - Data striped across N disks at byte level
 - Redundant disk stores parity
 - Read access
 - Read all disks
 - Write access
 - Generate new parity and update all disks
 - On failure
 - Use parity to reconstruct missing data
- Not widely used

RAID 4: Block-Interleaved Parity

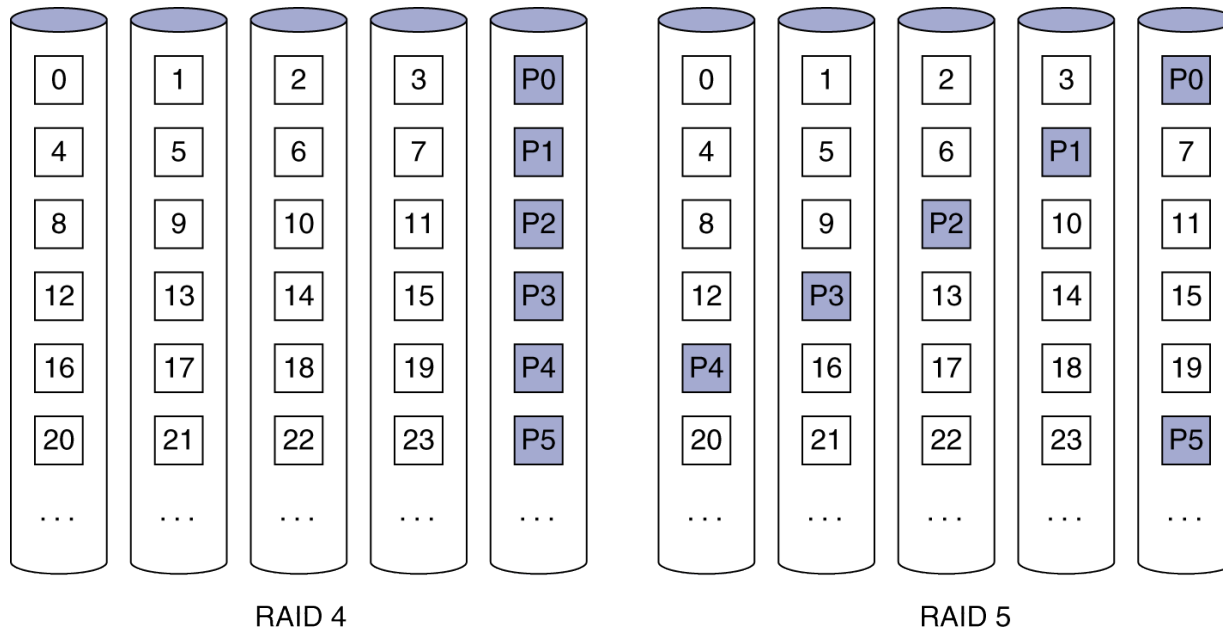
- N + 1 disks
 - Data striped across N disks at block level
 - Redundant disk stores parity for a group of blocks
 - Read access
 - Read only the disk holding the required block
 - Write access
 - Just read disk containing modified block, and parity disk
 - Calculate new parity, update data disk and parity disk
 - On failure
 - Use parity to reconstruct missing data
- Not widely used

RAID 3 vs RAID 4



RAID 5: Distributed Parity

- N + 1 disks
 - Like RAID 4, but parity blocks distributed across disks
 - Avoids parity disk being a bottleneck
- Widely used



RAID 6: P + Q Redundancy

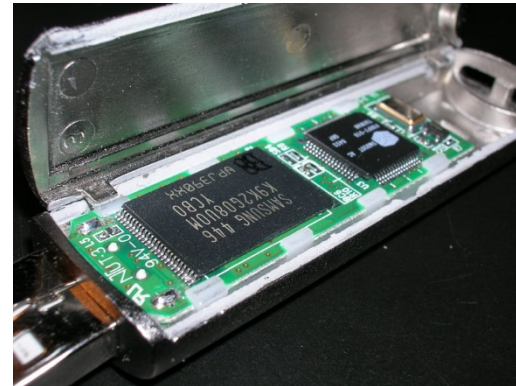
- N + 2 disks
 - Like RAID 5, but two lots of parity
 - Greater fault tolerance through more redundancy
- Multiple RAID
 - More advanced systems give similar fault tolerance with better performance

RAID Summary

- RAID can improve performance and availability
 - High availability requires hot swapping
- Assumes independent disk failures
 - Too bad if the building burns down!
- See “Hard Disk Performance, Quality and Reliability”
 - <http://www.pcguide.com/ref/hdd/perf/index.htm>

Flash Storage

- Nonvolatile semiconductor storage
 - 100× – 1000× faster than disk
 - Smaller, lower power, more robust
 - But more \$/GB (between disk and DRAM)



Flash Types

- NOR flash: bit cell like a NOR gate
 - Random read/write access
 - Used for instruction memory in embedded systems
- NAND flash: bit cell like a NAND gate
 - Denser (bits/area), but block-at-a-time access
 - Cheaper per GB
 - Used for USB keys, media storage, ...
- Flash bits wears out after 1000's of accesses
 - Not suitable for direct RAM or disk replacement
 - Wear leveling: remap data to less used blocks

Flash storage characteristics

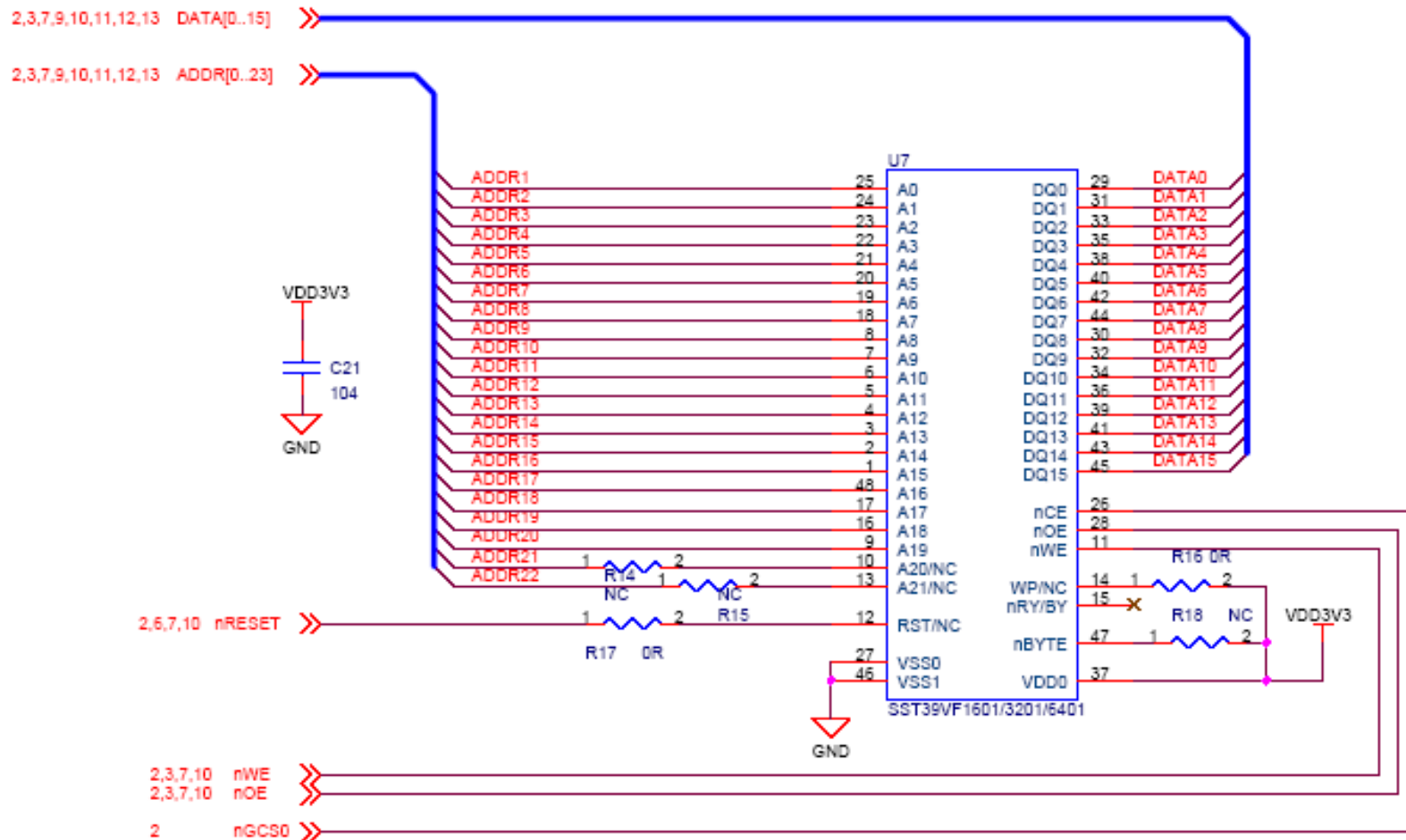
Characteristics	Kingston SecureDigital (SD) SD4/8 GB	Transend Type I CompactFlash TS16GCF133	RiDATA Solid State Disk 2.5 inch SATA
Formatted data capacity (GB)	8	16	32
Bytes per sector	512	512	512
Data transfer rate (read/write MB/sec)	4	20/18	68/50
Power operating/standby (W)	0.66/0.15	0.66/0.15	2.1/—
Size: height × width × depth (inches)	0.94 × 1.26 × 0.08	1.43 × 1.68 × 0.13	0.35 × 2.75 × 4.00
Weight in grams (454 grams/pound)	2.5	11.4	52
Mean time between failures (hours)	> 1,000,000	> 1,000,000	> 4,000,000
GB/cu. in., GB/watt	84 GB/cu.in., 12 GB/W	51 GB/cu.in., 24 GB/W	8 GB/cu.in., 16 GB/W
Best price (2008)	~ \$30	~ \$70	~ \$300

NOR vs NAND flash memory

Characteristics	NOR Flash Memory	NAND Flash Memory
Typical use	BIOS memory	USB key
Minimum access size (bytes)	512 bytes	2048 bytes
Read time (microseconds)	0.08	25
Write time (microseconds)	10.00	1500 to erase + 250
Read bandwidth (MBytes/second)	10	40
Write bandwidth (MBytes/second)	0.4	8
Wearout (writes per cell)	100,000	10,000 to 100,000
Best price/GB (2008)	\$65	\$4

SST39VF1601- 1M x 16 Flash

(on uCdragon board)



SST39VF1601 characteristics

- Organized as 1M x 16
 - 2K word sectors, 32K word blocks
- Performance:
 - Read access time = 70ns or 90ns
 - Word program time = 7us
 - Sector/block erase time = 18ms
 - Chip erase time = 40ms
- Check status of write/erase operation via read
 - DQ7 = complement of written value until write complete
 - DQ7=0 during erase, DQ7=1 when erase done



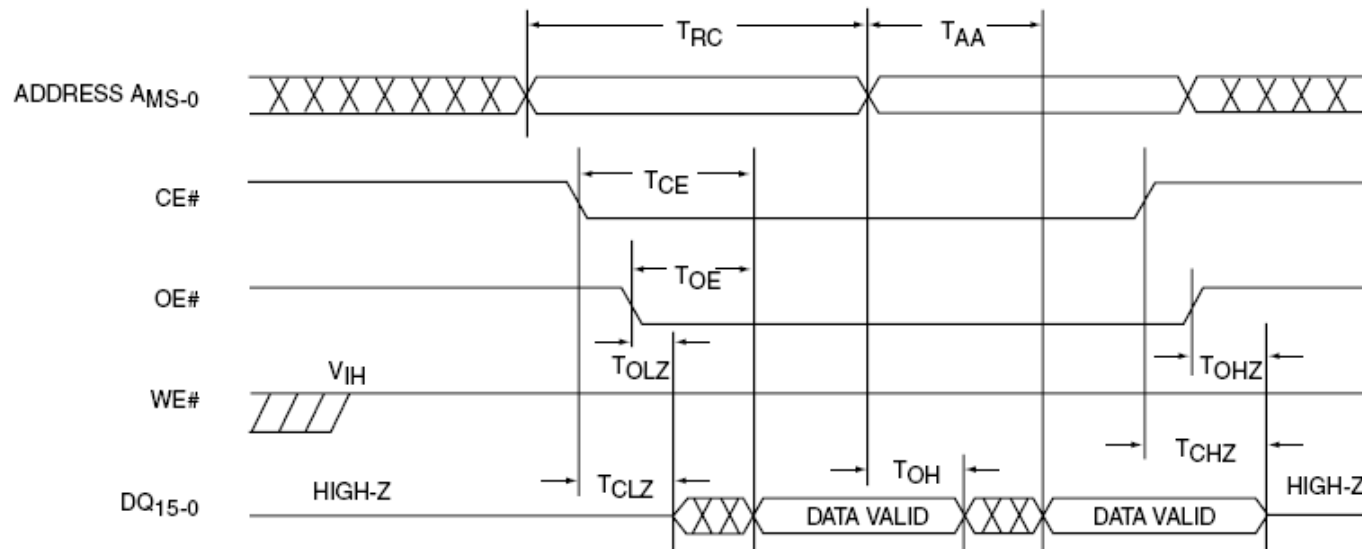
SST39VF1601 command sequences

(assert WE# and CE# to write commands)

Command Sequence	1st Bus Write Cycle		2nd Bus Write Cycle		3rd Bus Write Cycle		4th Bus Write Cycle		5th Bus Write Cycle		6th Bus Write Cycle	
	Addr ¹	Data ²	Addr ¹	Data ²	Addr ¹	Data ²	Addr ¹	Data ²	Addr ¹	Data ²	Addr ¹	Data ²
Word-Program	5555H	AAH	2AAAH	55H	5555H	A0H	WA ³	Data				
Sector-Erase	5555H	AAH	2AAAH	55H	5555H	80H	5555H	AAH	2AAAH	55H	SA _X ⁴	30H
Block-Erase	5555H	AAH	2AAAH	55H	5555H	80H	5555H	AAH	2AAAH	55H	BA _X ⁴	50H
Chip-Erase	5555H	AAH	2AAAH	55H	5555H	80H	5555H	AAH	2AAAH	55H	5555H	10H
Erase-Suspend	XXXXH	B0H										
Erase-Resume	XXXXH	30H										
Query Sec ID ⁵	5555H	AAH	2AAAH	55H	5555H	88H						
User Security ID Word-Program	5555H	AAH	2AAAH	55H	5555H	A5H	WA ⁶	Data				
User Security ID Program Lock-Out	5555H	AAH	2AAAH	55H	5555H	85H	XXH ⁶	0000H				
Software ID Entry ^{7,8}	5555H	AAH	2AAAH	55H	5555H	90H						
CFI Query Entry	5555H	AAH	2AAAH	55H	5555H	98H						
Software ID Exit ^{9,10} /CFI Exit/Sec ID Exit	5555H	AAH	2AAAH	55H	5555H	F0H						
Software ID Exit ^{9,10} /CFI Exit/Sec ID Exit	XXH	F0H										



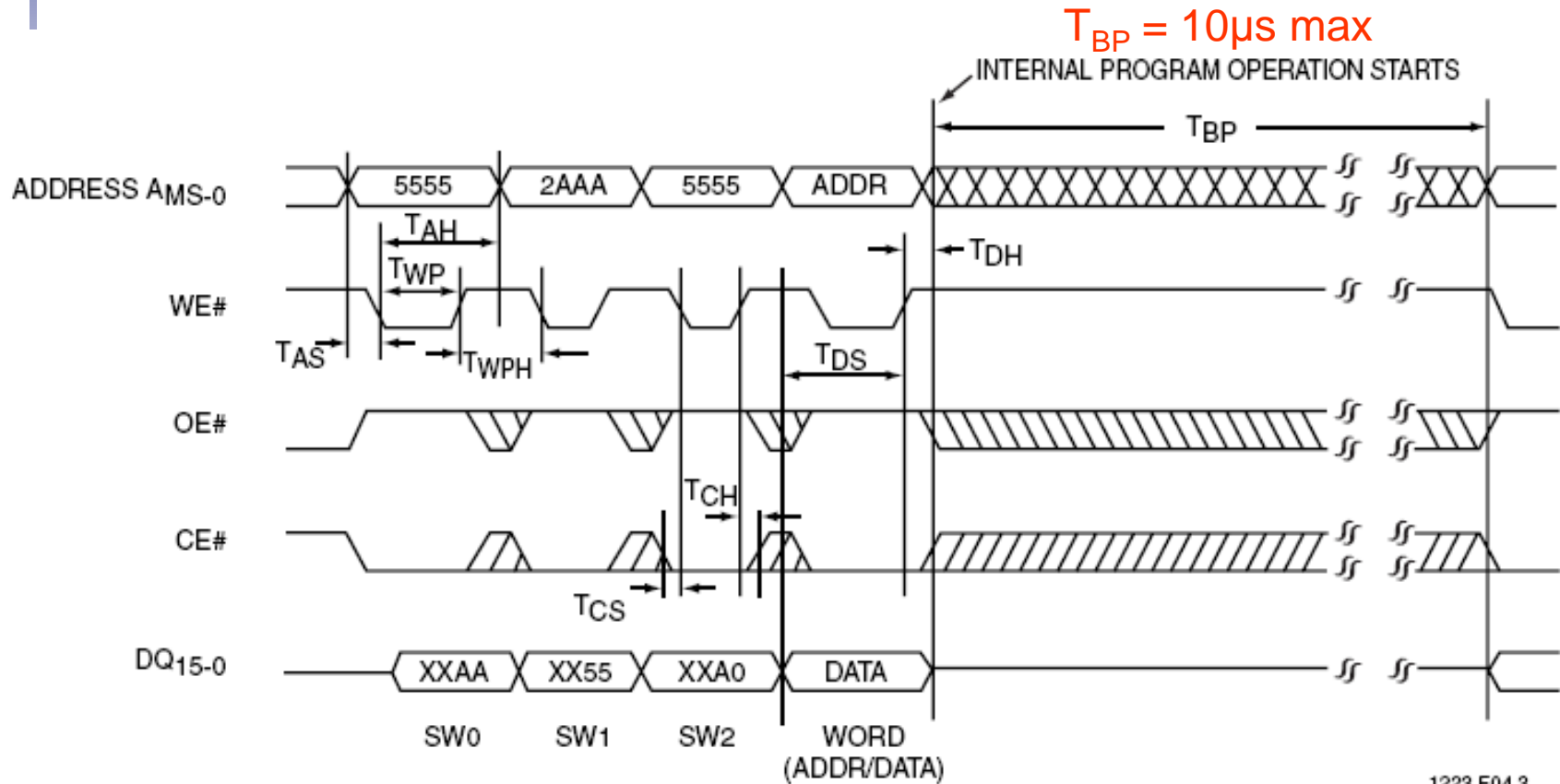
SST39VF1601 read cycle timing



Symbol	Parameter	SST39VFxx01/xx02-70		SST39VFxx01/xx02-90		Units
		Min	Max	Min	Max	
T_{RC}	Read Cycle Time	70		90		ns
T_{CE}	Chip Enable Access Time		70		90	ns
T_{AA}	Address Access Time		70		90	ns
T_{OE}	Output Enable Access Time		35		45	ns
T_{CLZ}^1	CE# Low to Active Output	0		0		ns
T_{OLZ}^1	OE# Low to Active Output	0		0		ns
T_{CHZ}^1	CE# High to High-Z Output		20		30	ns
T_{OHZ}^1	OE# High to High-Z Output		20		30	ns
T_{OH}^1	Output Hold from Address Change	0		0		ns
T_{RP}^1	RST# Pulse Width	500		500		ns
T_{RHR}^1	RST# High before Read	50		50		ns
$T_{RY}^{1,2}$	RST# Pin Low to Read Mode		20		20	μ s



SST39VF1601 word write



1223 F04.3

