



15-440 Distributed Systems

Errors and Failures

Types of Errors



- **Hard errors:** The component is dead.
- **Soft errors:** A signal or bit is wrong, but it doesn't mean the component must be faulty
- Note: You can have recurring soft errors due to faulty, but not dead, hardware

Examples



- DRAM errors
 - Hard errors: Often caused by motherboard - faulty traces, bad solder, etc.
 - Soft errors: Often caused by cosmic radiation or alpha particles (from the chip material itself) hitting memory cell, changing value. (Remember that DRAM is just little capacitors to store charge... if you hit it with radiation, you can add charge to it.)

Some fun #s



- Both Microsoft and Google have recently started to identify DRAM errors as an increasing contributor to failures... Google in their datacenters, Microsoft on your desktops.
- We've known hard drives fail
 - Especially when students need to hand in HW/projects :)

E.g., See "DRAM Errors in the Wild: A Large-Scale Field Study"

Replacement Rates



HPC1		COM1		COM2	
Component	%	Component	%	Component	%
Hard drive	30.6	Power supply	34.8	Hard drive	49.1
Memory	28.5	Memory	20.1	Motherboard	23.4
Misc/Unk	14.4	Hard drive	18.1	Power supply	10.1
CPU	12.4	Case	11.4	RAID card	4.1
motherboard	4.9	Fan	8	Memory	3.4
Controller	2.9	CPU	2	SCSI cable	2.2
QSW	1.7	SCSI Board	0.6	Fan	2.2
Power supply	1.6	NIC Card	1.2	CPU	2.2
MLB	1	LV Pwr Board	0.6	CD-ROM	0.6
SCSI BP	0.3	CPU heatsink	0.6	Raid Controller	0.6

From “Disk failures in the real world: What does an MTTF of 1,000,000 hours mean to you?”

5

HPC1 is a five year log of hardware replacements collected from a 765 node high-performance computing cluster. Each of the 765 nodes is a 4-way SMP with 4 GB of memory and three to four 18GB 10K rpm SCSI drives.

COM1 is a log of hardware failures recorded by an internet service provider. data covers a population of 26,734 10K rpm SCSI disk drives

COM2 is a warranty service log of hardware failures recorded on behalf of an internet service provider aggregating events in multiple distributed sites. The size of the underlying system changed significantly during the measurement period, starting with 420 servers in 2004 and ending with 9,232 servers in 2006.

Measuring Availability



- Mean time to failure (MTTF)
- Mean time to repair (MTTR)
- $MTBF = MTTF + MTTR$

- $Availability = MTTF / (MTTF + MTTR)$
 - Suppose OS crashes once per month, takes 10min to reboot.
 - $MTTF = 720 \text{ hours} = 43,200 \text{ minutes}$
 $MTTR = 10 \text{ minutes}$
 - $Availability = 43200 / 43210 = 0.997$ (~“3 nines”)

Availability



Availability %	Downtime per year	Downtime per month*	Downtime per week
90% ("one nine")	36.5 days	72 hours	16.8 hours
95%	18.25 days	36 hours	8.4 hours
97%	10.96 days	21.6 hours	5.04 hours
98%	7.30 days	14.4 hours	3.36 hours
99% ("two nines")	3.65 days	7.20 hours	1.68 hours
99.50%	1.83 days	3.60 hours	50.4 minutes
99.80%	17.52 hours	86.23 minutes	20.16 minutes
99.9% ("three nines")	8.76 hours	43.8 minutes	10.1 minutes
99.95%	4.38 hours	21.56 minutes	5.04 minutes
99.99% ("four nines")	52.56 minutes	4.32 minutes	1.01 minutes
99.999% ("five nines")	5.26 minutes	25.9 seconds	6.05 seconds
99.9999% ("six nines")	31.5 seconds	2.59 seconds	0.605 seconds
99.99999% ("seven nines")	3.15 seconds	0.259 seconds	0.0605 seconds

Availability in practice



- Carrier airlines (2002 FAA fact book)
 - 41 accidents, 6.7M departures
 - 99.9993% availability
- 911 Phone service (1993 NRIC report)
 - 29 minutes per line per year
 - 99.994%
- Standard phone service (various sources)
 - 53+ minutes per line per year
 - 99.99+%
- End-to-end Internet Availability
 - 95% - 99.6%

Real Devices



PRODUCT OVERVIEW

Cheetah 15K.4

Mainstream enterprise disc drive

Simply the best price/
performance, lowest cost of
ownership disc drive ever



KEY FEATURES AND BENEFITS

- The Cheetah® 15K.4 is the highest performance drive ever offered by Seagate™, offering up to 10% better drive-to-disk (D2D) performance than any other 3.5" drive.
- The Cheetah 15K.4 offers performance advantages with the breakthrough benefits of parallel ATA (PATA) over SATA, such as the option to load data for less and optimize storage.
- Provides self-optimized background management functions to reduce media integrity, reduce drive off-line, reduce redundancy of redundant drives and improve data reliability.
- The Cheetah 15K.4 offers its maximum performance and features like with Cheetah 15K.7 and Cheetah™ 15K.8 to ensure greater future flexibility and reduced drive-to-disk.

KEY SPECIFICATIONS

- 15K, 7200 RPM, 300MB/s
- 3.5" form factor and 3.5" form factor with 3.5" form factor
- 20 to 30 MB/sec sustained transfer rate
- 1.8 Million Bytes (MB) data cache (MP3)
- Serial ATA (SATA) 3.0G, SATA 3.0G and 2.0G SATA 3.0G interface
- 5 year warranty

For more information on why 15K is the industry's best price-performance disc drive for use in high-end storage applications, visit <http://seagate.com/15K>

Real Devices – the small print



- The Cheetah 15K.4 is the highest-performance drive ever offered by Seagate, delivering maximum IOPS with fewer drives to yield lower TCO.
- The Cheetah 15K.4 price-per-performance value united with the breakthrough benefits of serial attached SCSI (SAS) make it the optimal 3.5-inch drive for rock solid enterprise storage.
- Proactive, self-initiated background management functions improve media integrity, increase drive efficiency, reduce incidence of integration failures and improve field reliability.
- The Cheetah 15K.4 shares its electronics architecture and firmware base with Cheetah 10K.7 and Savvio™ to ensure greater factory consistency and reduced time to market.

KEY SPECIFICATIONS

- 146-, 73- and 36-Gbyte capacities
- 3.3-msec average read and 3.8-msec average write seek times
- Up to 98-Mbytes/sec sustained transfer rate
- 1.4 million hours full duty cycle MTBF
- Serial Attached SCSI (SAS), Ultra320 SCSI and 2 Gbits/sec Fibre Channel interfaces
- 5-year warranty

For more information on why 15K is the industry's best price/performance disc drive for use in mainstream storage applications, visit <http://specials.seagate.com/15k>

10

Modern SSD 1.5M MTBF === 170 years

720hrs/month === 2 month = 1.5 k hrs === 2000 months = 1.5m hrs = 2000/12 == 170yrs

Real Devices – the small print

170 years....??!

Last Will and Testament
One hard drive for my great-grandkids to use

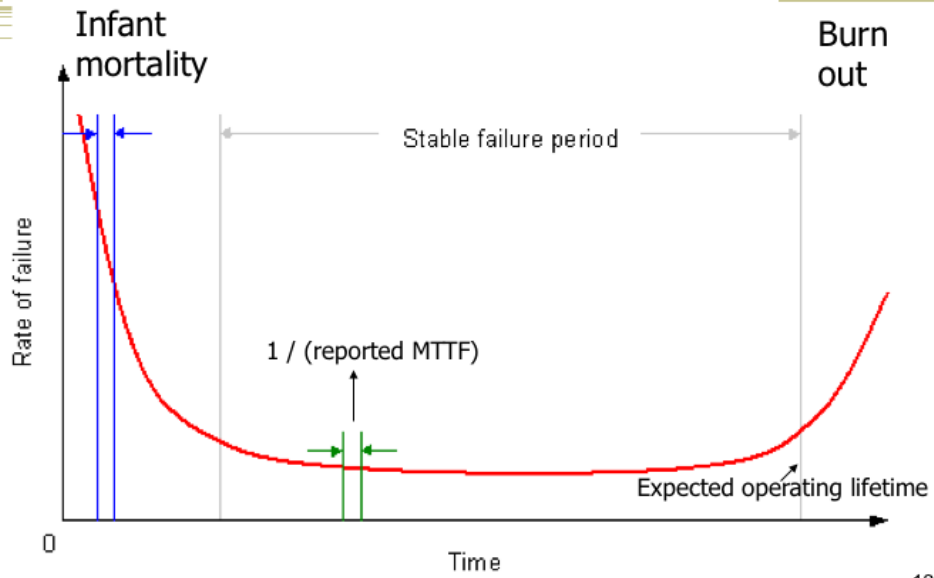
- The Cheetah 15K.4 is the highest-performance drive ever offered by Seagate, delivering maximum IOPS with fewer drives to yield lower TCO.
- The Cheetah 15K.4 price-per-performance value united with the breakthrough benefits of SAS (SAS) make it the optimal 3.5-inch drive for rock solid performance. Advanced background management functions improve media integrity, reduce migration failures and improve overall drive reliability. The Cheetah 15K.4 SAS drive is built on a proven architecture and firmware base with a long history of reliability, consistency and reduced downtime.
- Ultra-low latency random write seek times
- Up to 12 Gbits/sec Fibre Channel interfaces
- Best price/performance disc drive for SAS applications. Visit <http://specials.seagate.com/15k>

11

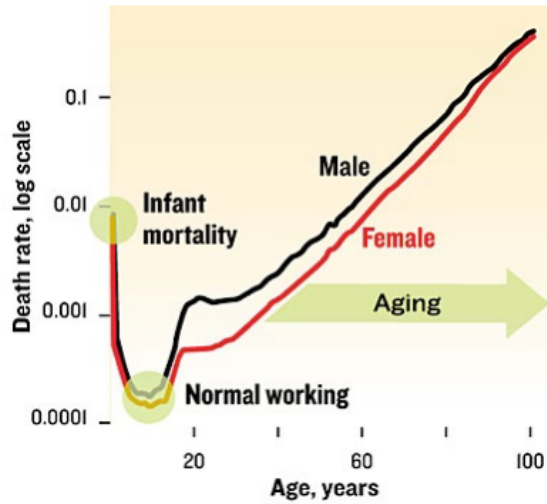
Modern SSD 1.5M MTBF === 170 years

720hrs/month === 2 month = 1.5 k hrs === 2000 months = 1.5m hrs = 2000/12 == 170yrs

Disk failure conditional probability distribution - Bathtub curve



Other (more morbid) Bathtub Curves



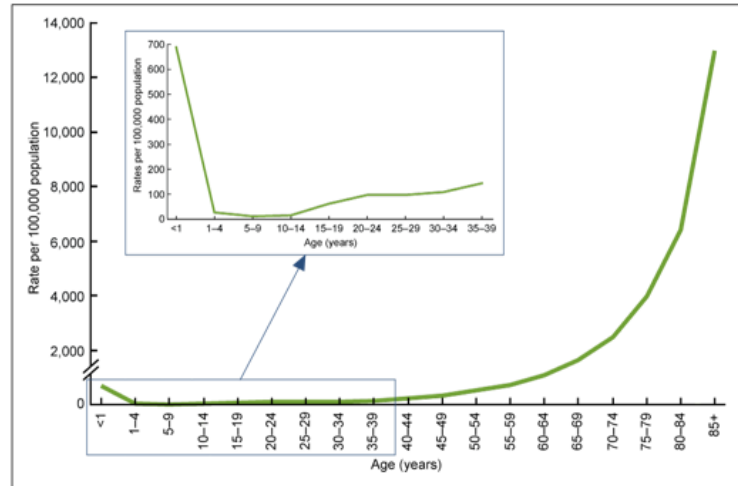
Human
Mortality
Rates
(US, 1999)

From: L. Gavrilov & N. Gavrilova, "Why We Fall Apart," IEEE Spectrum, Sep. 2004.
Data from <http://www.mortality.org>

Other Bathtub Curves



Figure 2. Age-specific death rates: United States, preliminary 2007



SOURCE: National Vital Statistics System, Mortality.

From: L. Gavrilov & N. Gavrilova, "Why We Fall Apart," *IEEE Spectrum*, Sep. 2004.
Data from <http://www.mortality.org>

So, back to disks...



- How can disks fail?
 - Whole disk failure (power supply, electronics, motor, etc.)
 - Sector errors - soft or hard
 - Read or write to the wrong place (e.g., disk is bumped during operation)
 - Can fail to read or write if head is too high, coating on disk bad, etc.
 - Disk head can hit the disk and scratch it.

Coping with failures...



- A failure
 - Let's say one bit in your DRAM fails.
- Propagates
 - Assume it flips a bit in a memory address the kernel is writing to. That causes a big memory error elsewhere, or a kernel panic.
 - This program is running one of a dozen storage servers for your distributed filesystem.
 - A client can't read from the DFS, so it hangs.
 - A professor can't check out a copy of your 15-440 assignment, so he gives you an F.

Recovery Techniques



- We've already seen some: e.g., retransmissions in TCP and in your RPC system
- Modularity can help in failure isolation: preventing an error in one component from spreading.
 - Analogy: The firewall in your car keeps an engine fire from affecting passengers
- Today: Redundancy and Retries
 - Next lecture: Specific techniques used in file systems, disks
 - This time: Understand how to quantify reliability
 - Understand basic techniques of replication and fault masking

What are our options?

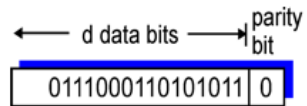


1. Silently return the wrong answer.
2. **Detect failure.**
3. Correct / mask the failure

Parity Checking



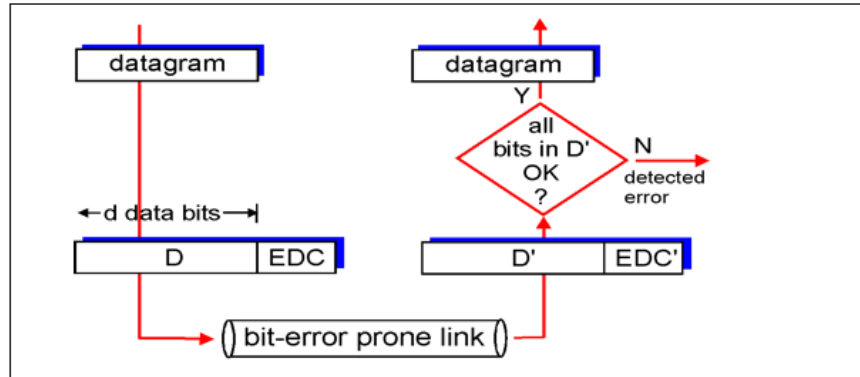
Single Bit Parity: Detect single bit errors



Block Error Detection



- EDC= Error Detection and Correction bits (redundancy)
- D = Data protected by error checking, may include header fields
- Error detection not 100% reliable!
 - Protocol may miss some errors, but rarely
 - Larger EDC field yields better detection and correction



Error Detection - Checksum



- Used by TCP, UDP, IP, etc..
- Ones complement sum of all words/shorts/bytes in packet
- Simple to implement
- Relatively weak detection
 - Easily tricked by typical error patterns – e.g. bit flips

Example: Internet Checksum



- Goal: detect “errors” (e.g., flipped bits) in transmitted segment

Sender

- Treat segment contents as sequence of 16-bit integers
- Checksum: addition (1's complement sum) of segment contents
- Sender puts checksum value into checksum field in header

Receiver

- Compute checksum of received segment
- Check if computed checksum equals checksum field value:
 - NO - error detected
 - YES - no error detected. But maybe errors nonetheless?

22

16 bit checksum = 2^{16} values

This means that 1 out 2^{16} packets have the same checksum

Error Detection – Cyclic Redundancy Check (CRC)

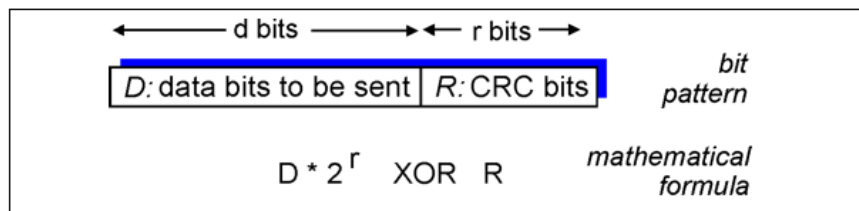


- Polynomial code
 - Treat packet bits as coefficients of n-bit polynomial
 - Choose r+1 bit generator polynomial (well known – chosen in advance)
 - Add r bits to packet such that message is divisible by generator polynomial
- Better loss detection properties than checksums
 - Cyclic codes have favorable properties in that they are well suited for detecting burst errors
 - Therefore, used on networks/hard drives

Error Detection – CRC



- View data bits, D , as a binary number
- Choose $r+1$ bit pattern (generator), G
- Goal: choose r CRC bits, R , such that
 - $\langle D, R \rangle$ exactly divisible by G (modulo 2)
 - Receiver knows G , divides $\langle D, R \rangle$ by G . If non-zero remainder: error detected!
 - Can detect all burst errors less than $r+1$ bits
- Widely used in practice



CRC Example



Want:

$$D \cdot 2^r \text{ XOR } R = nG$$

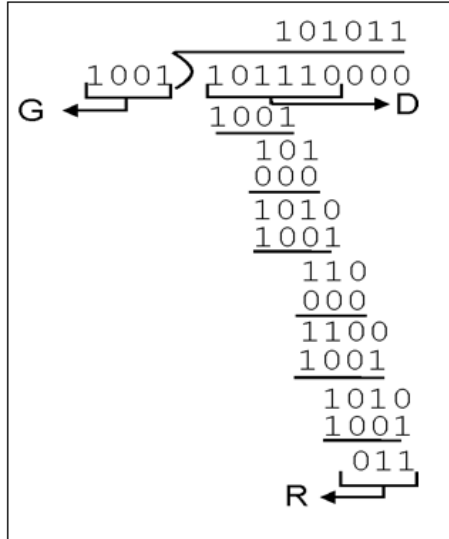
equivalently:

$$D \cdot 2^r = nG \text{ XOR } R$$

equivalently:

if we divide $D \cdot 2^r$ by G ,
want remainder R

$$R = \text{remainder} \left[\frac{D \cdot 2^r}{G} \right]$$



Break... and valuable feedback



10. You can't cheat on this exam because no one knows the answers.
9. This class was like milk, it was good for 2 weeks.
8. I would have been better off using the tuition money to heat my apartment last winter.
7. Emotional scarring may fade away, but this big fat F on my transcript won't.
6. Profs are evil computer science teaching robots who crush students for pleasure.
5. I learned there are 137 tiles on the ceiling.
4. Not only is the book a better teacher, it also has a better personality.
3. Teaches well, invites questions - then tells bad jokes for 20 minutes.
2. Bring a pillow.
1. Your pillow will need a pillow.

(really from RateMyProfessor.com)

Real Feedback



- Lecture
 - Speed, breaks, more examples
- Slides
 - More/Less detail, before class, with answers
- Links between projects/HW/lecture
 - E.g., DFS or mutual exclusion
 - Too much networking
- Projects are hard but rewarding
- Love Go
- More recent systems
- More recitations/office hours

What are our options?



1. Silently return the wrong answer.
2. Detect failure.
3. Correct / mask the failure

Error Recovery



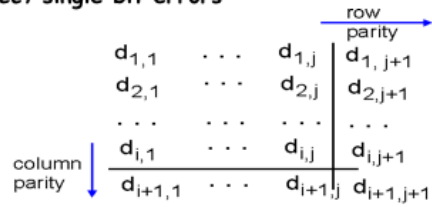
- Two forms of error recovery
 - Redundancy
 - Error Correcting Codes (ECC)
 - Replication/Voting
 - Retry
- ECC
 - Keep encoded redundant data to help repair losses
 - Forward Error Correction (FEC) – send bits in advance
 - Reduces latency of recovery at the cost of bandwidth

Error Recovery – Error Correcting Codes (ECC)



Two Dimensional Bit Parity:

Detect *and correct* single bit errors



10101 1	10101 1	
11110 0	1 01100 0	→ parity error
01110 1	01110 1	
00101 0	00101 0	
<i>no errors</i>	parity error	
	<i>correctable single bit error</i>	

Replication/Voting



- If you take this to the extreme
[r1] [r2] [r3]
- Send requests to all three versions of the software: Triple modular redundancy
 - Compare the answers, take the majority
 - Assumes no error detection
- In practice - used mostly in space applications; some extreme high availability apps (stocks & banking? maybe. But usually there are cheaper alternatives if you don't need real-time)
 - Stuff we cover later: surviving malicious failures through voting (byzantine fault tolerance)

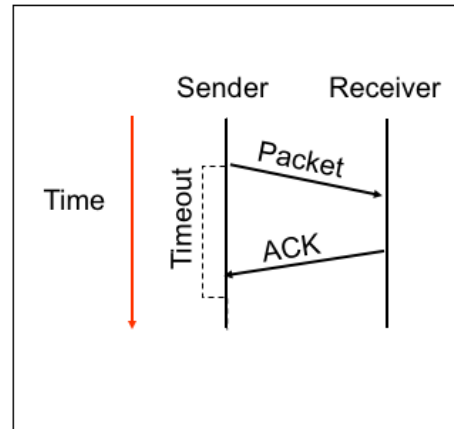
31

31

Retry – Network Example



- Sometimes errors are transient
- Need to have error detection mechanism
 - E.g., timeout, parity, checksum
 - No need for majority vote



One key question



- How correlated are failures?
- Can you assume independence?
 - If the failure probability of a computer in a rack is p ,
 - What is $p(\text{computer 2 failing}) \mid \text{computer 1 failed}$?
 - Maybe it's p ... or maybe they're both plugged into the same UPS...
- Why is this important?
 - Correlation reduces value of redundancy

Fault Tolerant Design



- Quantify probability of failure of each component
- Quantify the costs of the failure
- Quantify the costs of implementing fault tolerance

- This is all probabilities...

Summary



- Definition of MTTF/MTBF/MTTR: Understanding availability in systems.
- Failure detection and fault masking techniques
- Engineering tradeoff: Cost of failures vs. cost of failure masking.
 - At what level of system to mask failures?
 - Leading into replication as a general strategy for fault tolerance
- Thought to leave you with:
 - What if you have to survive the failure of entire computers? Of a rack? Of a datacenter?

Replacement Rates



HPC1		COM1		COM2	
Component	%	Component	%	Component	%
Hard drive	30.6	Power supply	34.8	Hard drive	49.1
Memory	28.5	Memory	20.1	Motherboard	23.4
Misc/Unk	14.4	Hard drive	18.1	Power supply	10.1
CPU	12.4	Case	11.4	RAID card	4.1
motherboard	4.9	Fan	8	Memory	3.4
Controller	2.9	CPU	2	SCSI cable	2.2
QSW	1.7	SCSI Board	0.6	Fan	2.2
Power supply	1.6	NIC Card	1.2	CPU	2.2
MLB	1	LV Pwr Board	0.6	CD-ROM	0.6
SCSI BP	0.3	CPU heatsink	0.6	Raid Controller	0.6

Back to Disks...

What are our options?



1. Silently return the wrong answer.
2. Detect failure.
 - Every sector has a header with a checksum. Every read fetches both, computes the checksum on the data, and compares it to the version in the header. Returns error if mismatch.
3. Correct / mask the failure
 - Re-read if the firmware signals error (may help if transient error, may not)
 - Use an error correcting code (what kinds of errors do they help?)
 - Bit flips? Yes. Block damaged? No
 - Have the data stored in multiple places (RAID)

Fail-fast disk



```
failfast_get (data, sn) {  
    get (s, sn);  
    if (checksum(s.data) = s.cksum) {  
        data ← s.data;  
        return OK;  
    } else {  
        return BAD;  
    }  
}
```

Careful disk



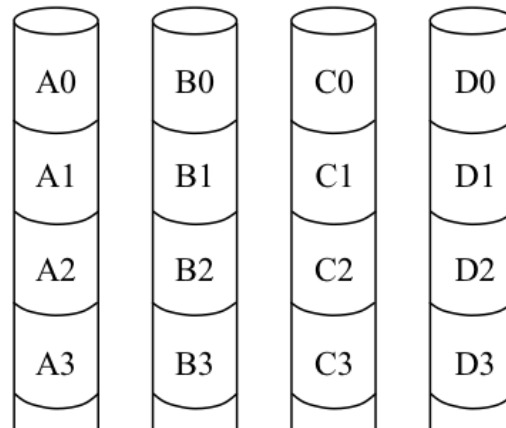
```
careful_get (data, sn) {  
    r ← 0;  
    while (r < 10) {  
        r ← failfast_get (data, sn);  
        if (r = OK) return OK;  
        r++;  
    }  
    return BAD;  
}
```

Use multiple disks?



- Capacity
 - More disks allows us to store more data
- Performance
 - Access multiple disks in parallel
 - Each disk can be working on independent read or write
 - Overlap seek and rotational positioning time for all
- Reliability
 - Recover from disk (or single sector) failures
 - Will need to store multiple copies of data to recover
- So, what is the simplest arrangement?

Just a bunch of disks (JBOD)



- Yes, it's a goofy name
 - industry really does sell "JBOD enclosures"

Disk Subsystem Load Balancing

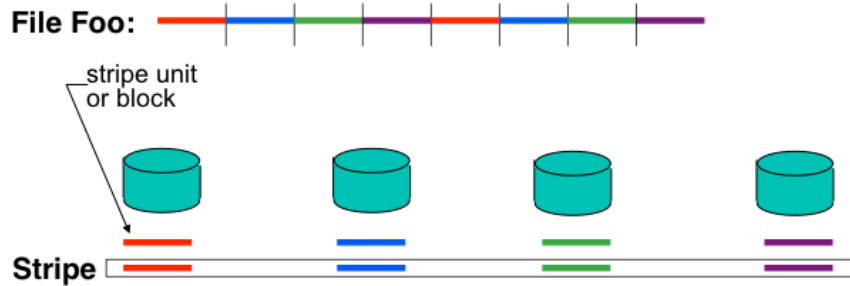


- I/O requests are almost never evenly distributed
 - Some data is requested more than other data
 - Depends on the apps, usage, time, ...
- What is the right data-to-disk assignment policy?
 - Common approach: Fixed data placement
 - Your data is on disk X, period!
 - For good reasons too: you bought it or you're paying more...
 - Fancy: Dynamic data placement
 - If some of your files are accessed a lot, the admin(or even system) may separate the "hot" files across multiple disks
 - In this scenario, entire files systems (or even files) are manually moved by the system admin to specific disks
 - Alternative: Disk striping
 - Stripe all of the data across all of the disks

Disk Striping



- Interleave data across multiple disks
 - Large file streaming can enjoy parallel transfers
 - High throughput requests can enjoy thorough load balancing
 - If blocks of hot files equally likely on all disks (really?)



Disk striping details



- How disk striping works
 - Break up total space into fixed-size stripe units
 - Distribute the stripe units among disks in round-robin
 - Compute location of block #B as follows
 - $\text{disk\#} = B \% N$ ($\% = \text{modulo}$, $N = \text{\#ofdisks}$)
 - $\text{LBN\#} = B / N$ (computes the LBN on given disk)

Now, What If A Disk Fails?



- In a JBOD (independent disk) system
 - one or more file systems lost
- In a striped system
 - a part of each file system lost
- Backups can help, but
 - backing up takes time and effort
 - backup doesn't help recover data lost during that day
 - Any data loss is a big deal to a bank or stock exchange

Tolerating and masking disk failures



- If a disk fails, it's data is gone
 - may be recoverable, but may not be
- To keep operating in face of failure
 - must have some kind of data redundancy
- Common forms of data redundancy
 - replication
 - erasure-correcting codes
 - error-correcting codes