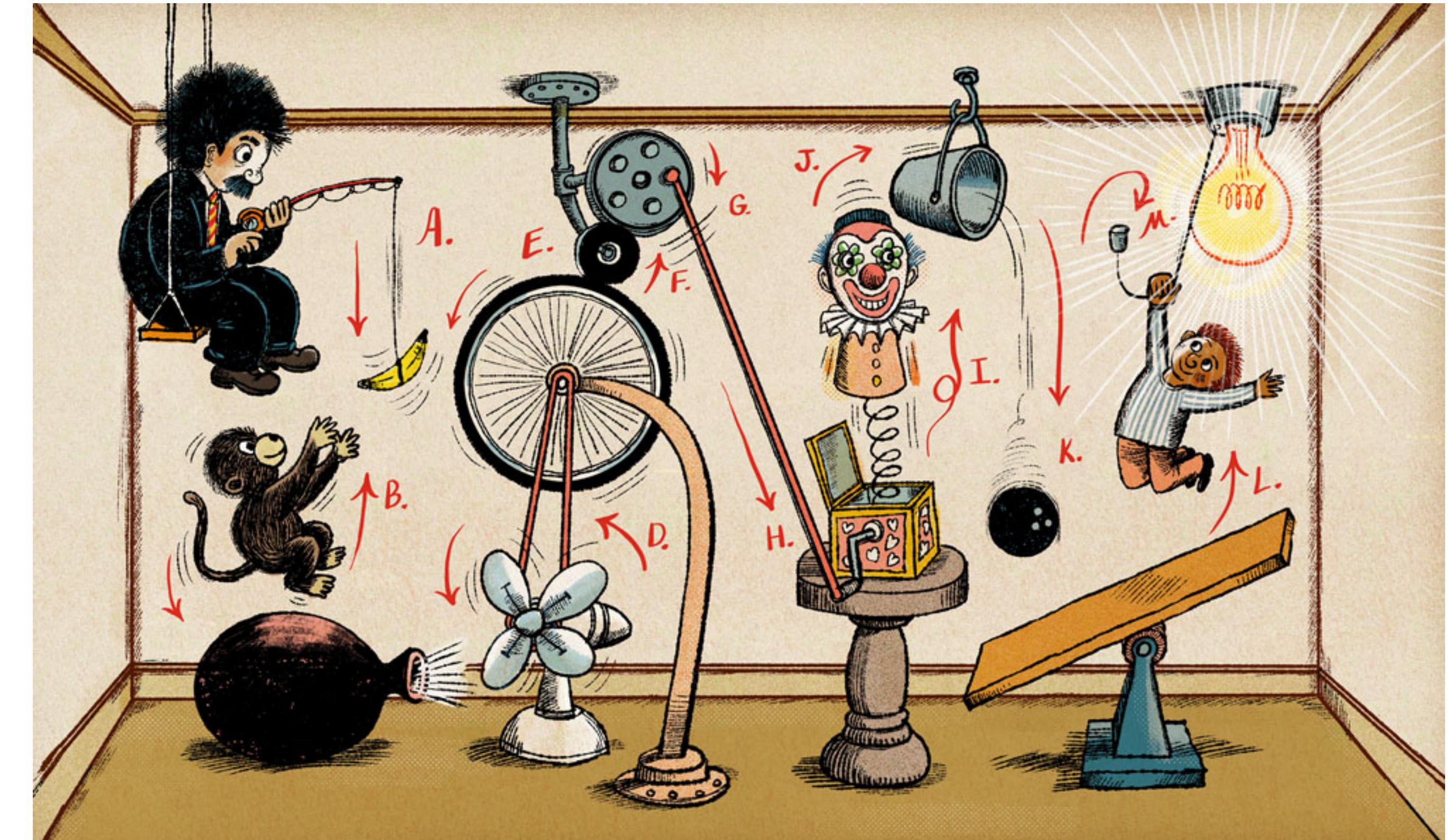


CS 61C: Great Ideas in Computer Architecture

Lecture 24: Dependability

Dependability: Reliability in Computer Hardware

- What could possibly go wrong?
- How do we mitigate the effects of faults and errors?



Types of Faults in Digital Designs

1. Design Bugs (function, timing, power draw)

- detected and corrected at design time through testing and verification (simulation, static checks)

2. Manufacturing Defects (violation of design rules, impurities in processing, statistical variations)

- post production testing for sorting
- spare on-chip resources for repair

3. Runtime Failures (physical effects and environmental conditions)

- “Hard Faults”: aging
- “Soft (transient) Faults”: electro-magnetic interference, cosmic particles

Intel Pentium FDIV Design Bug

A hardware bug affecting the floating point unit (FPU) of the early Intel Pentium processors.

The processor might return incorrect binary floating point results when dividing a number.
For example, the chip produced:

$$\frac{4,195,835}{3,145,727} = 1.333\textcolor{red}{739068902037589}$$

Instead of the correct:

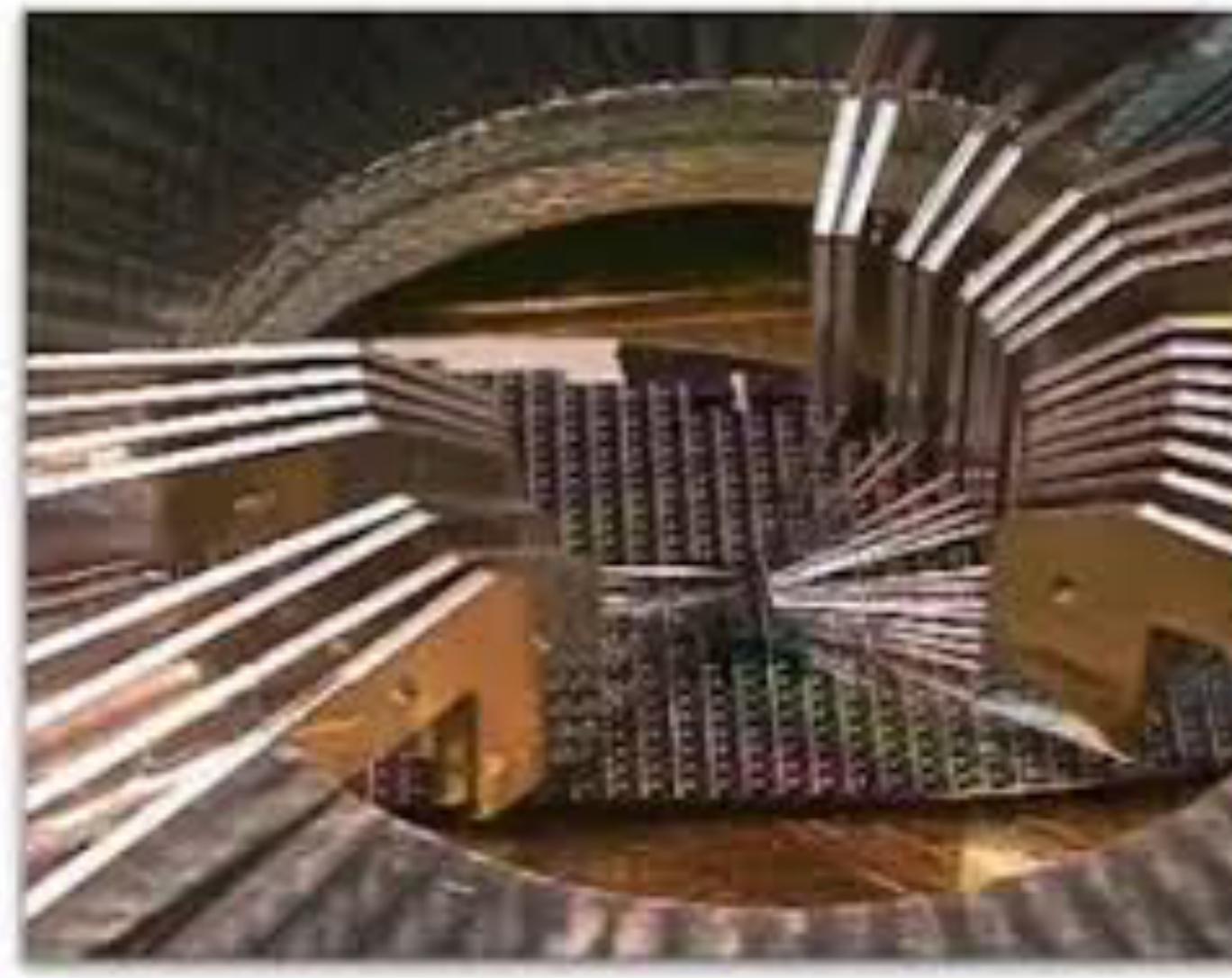
$$\frac{4,195,835}{3,145,727} = 1.333820449136241002$$


Intel attributed the error to missing entries in the lookup table used by the floating-point division circuitry.

In December 1994, Intel recalled the defective processors. In January 1995, Intel announced a pre-tax charge of \$475 million against earnings, the cost associated with replacement of the flawed processors.

Dealing with Manufacturing Faults in ICs

- Designers provide “test vectors”
 - Tools help with ATPG (Automatic Test Pattern Generation)
- Completed ICs are tested and “binned” for correct operation, and speed grade.



- Special on-chip circuits help speed the testing process
 - BIST (built in self test), Scan-chains

Today: Dealing with Runtime Failures (errors)

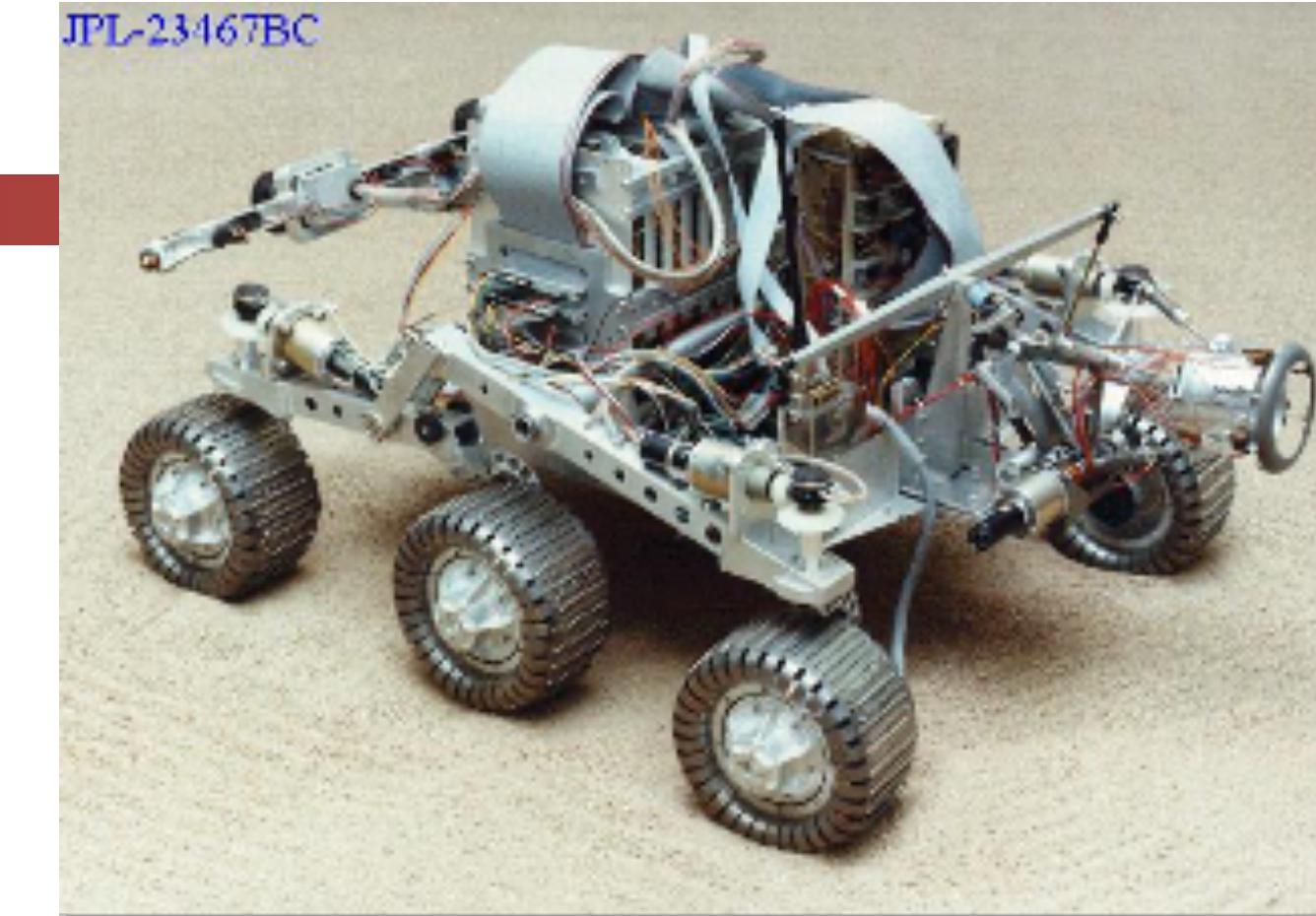
- Redundancy
- Measures of Dependability
- Codes for Error Detection/Correction
- Protecting Disk-Drives Against Errors
- And in Conclusion, ...

Great Idea #6: Dependability via Redundancy

- Applies to everything from data centers to memory
 - Redundant data centers so that can lose 1 datacenter but Internet service stays online
 - Redundant routes so can lose nodes but Internet doesn't fail
 - Or at least can recover quickly...
 - Redundant disks so that can lose 1 disk but not lose data (Redundant Arrays of Independent Disks (RAID))
 - Redundant memory bits so that can lose 1 bit but no data (Error Correcting Code (ECC) Memory)

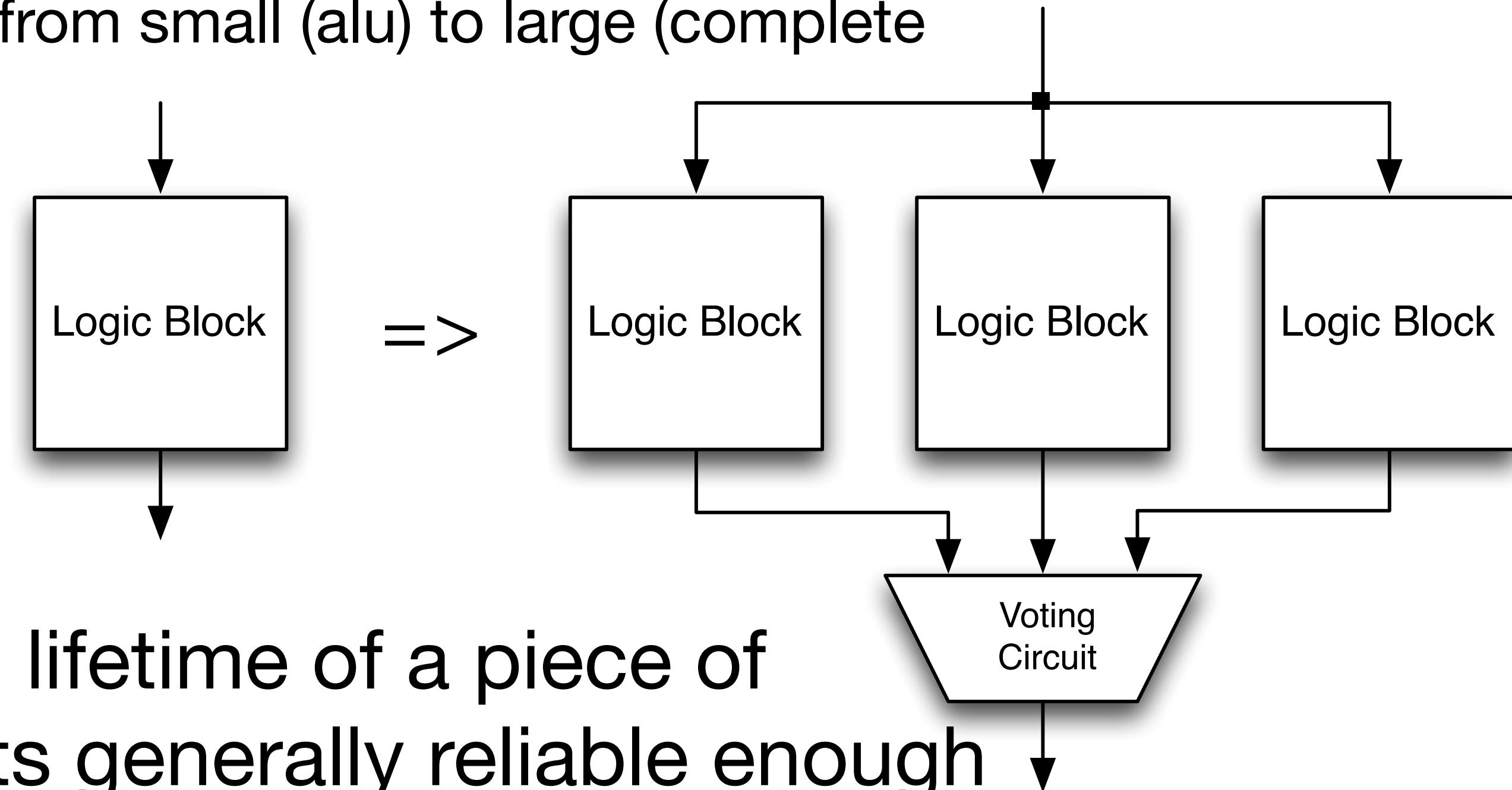


A Fault-Tolerant Design Methodology



Computer Science 61C Fall 2021

- Where hardware errors cannot be tolerated
 - Space-craft data systems, life-supporting medical devices, ...
- Triple Modular Redundancy
 - Relies on simple reliable voting circuit, assume $P(>1 \text{ error})$ is small
 - “Logic Block” could be all the way from small (alu) to large (complete processor)

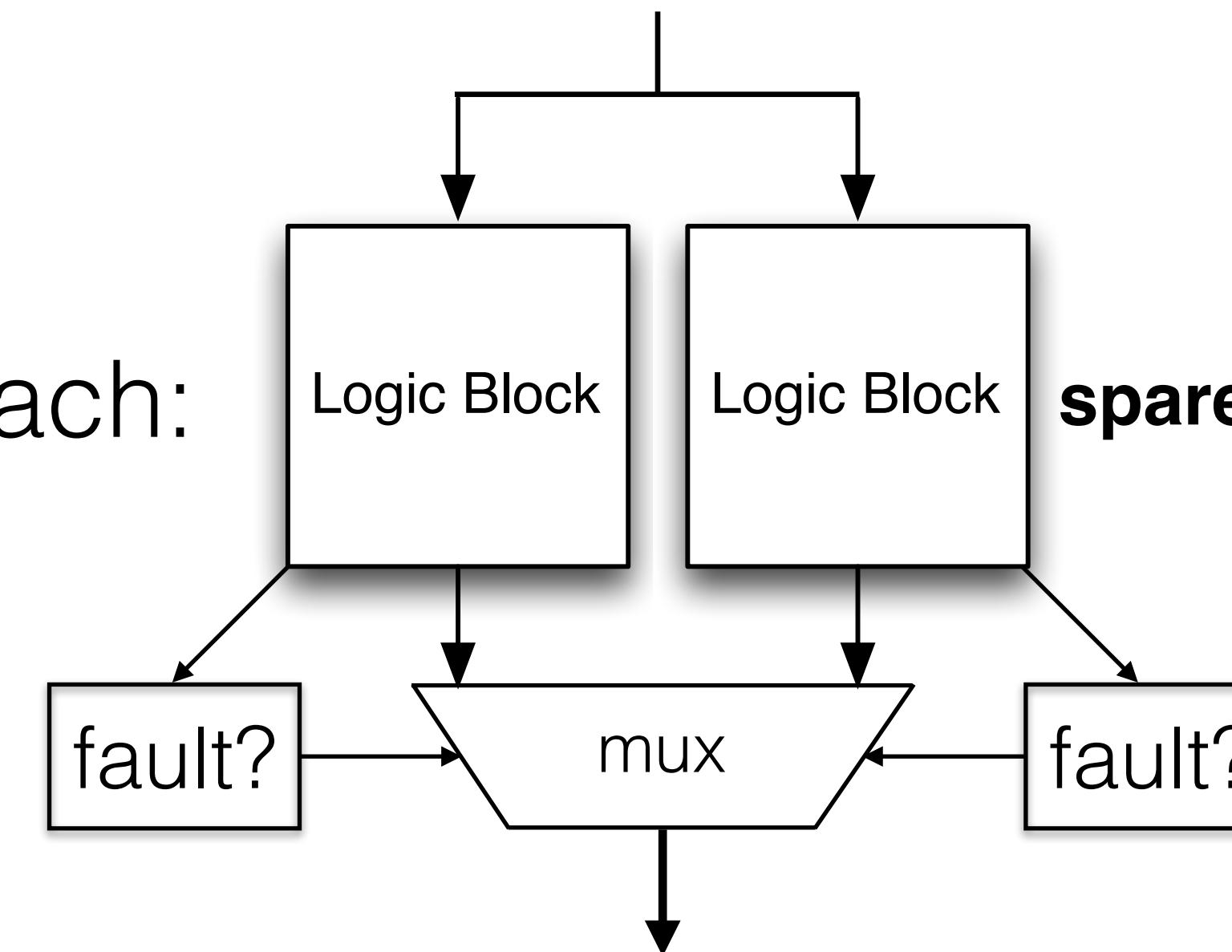


- On earth and for the normal lifetime of a piece of equipment, hardware circuits generally reliable enough
 - except for *large memory circuits* (we will see later how to protect)

Dependability Corollary: Fault Detection

- The ability to determine that ***something*** is wrong is often the key to effectively using redundancy:

Simple “sparing” approach:



- Often error detection is a necessary prerequisite to error correction

Dependability via Redundancy: Time vs. Space

- ***Spatial Redundancy*** – replicated data or extra information or hardware to handle hard and soft (transient) failures
- ***Temporal Redundancy*** – redundancy in time (retry) to handle soft (transient) failures
- "Insanity overcoming soft failures: repeatedly doing the same thing and expecting different results"

Dependability Measures

- *Reliability*: Mean Time To Failure (**MTTF**)
- *Service interruption*: Mean Time To Repair (**MTTR**)
- Mean time between failures (**MTBF**)
 - $\text{MTBF} = \text{MTTF} + \text{MTTR}$
- Availability = **MTTF / (MTTF + MTTR)**
- Improving Availability
 - Increase **MTTF**: More reliable hardware/software + Fault Tolerance
 - Reduce **MTTR**: improved tools and processes for diagnosis and repair

Availability Example

- Suppose your laptop dies on average once every 2 years
- Each time you take make an appointment at the repair shop and get it repaired, resulting in a total of 2 days downtime.
- What is the availability of your laptop?

$$MTTF = 2 * 365 \text{ days}$$

$$MTTR = 2 \text{ days}$$

$$\text{availability} = \frac{MTTF}{MTTF + MTTR} = \frac{730}{730 + 2} = .997 = 99.7\%$$

Availability Measures

- Availability = $MTTF / (MTTF + MTTR)$ as %
- Since we hope things are rarely down, shorthand is “number of 9s of availability per year”
- 1 nine: 90% => 36 days of repair/year
 - Airbears Reliability?
- 2 nines: 99% => 3.6 days of repair/year
- 3 nines: 99.9% => 526 minutes of repair/year
- 4 nines: 99.99% => 53 minutes of repair/year
- 5 nines: 99.999% => 5 minutes of repair/year
 - And serious \$\$\$ to do

Reliability Measures

- Another is average number of failures per year:

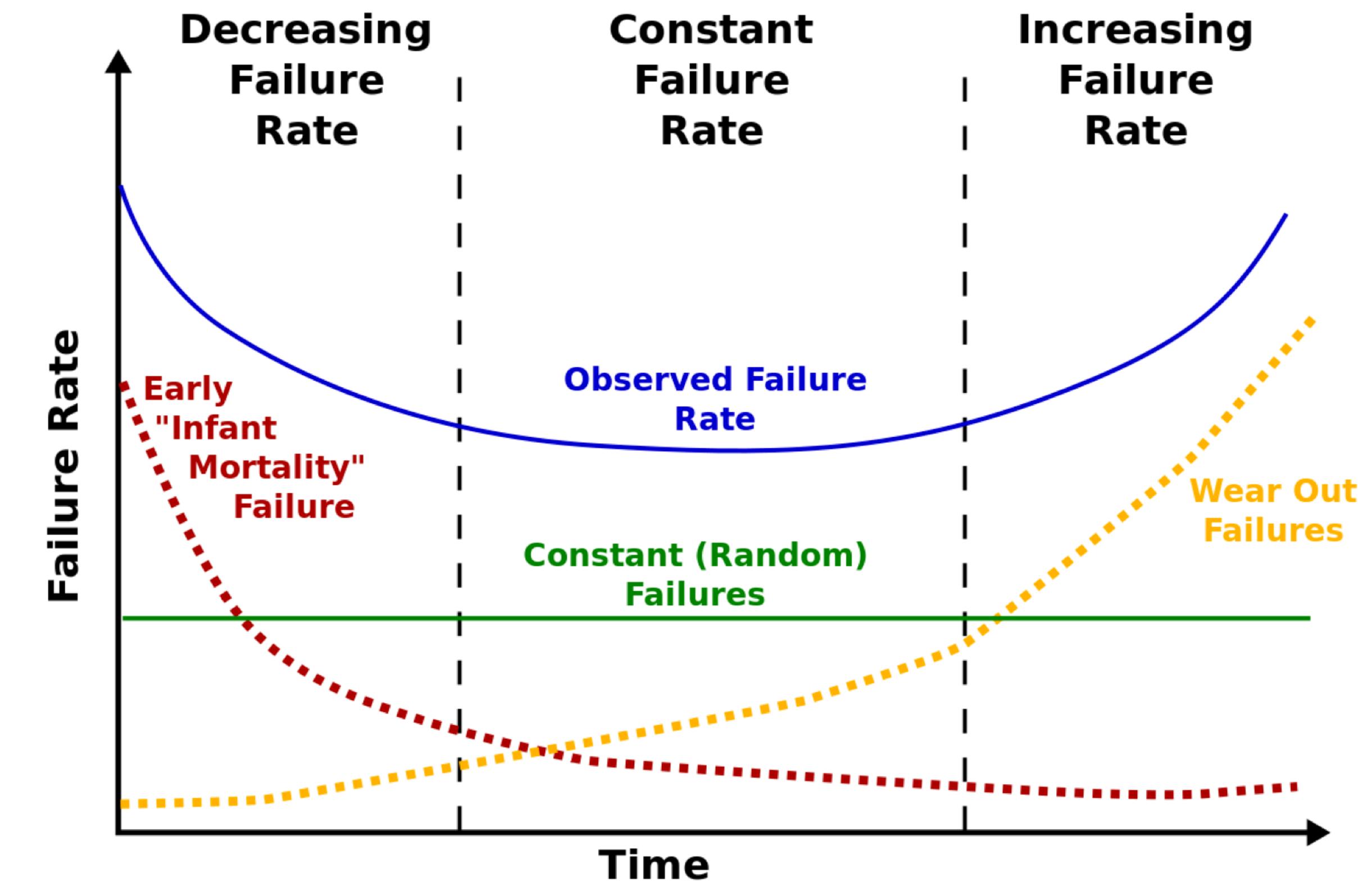
Annualized Failure Rate (AFR)

- E.g., 1000 disks with 100,000 hour MTTF
- $365 \text{ days/yr} * 24 \text{ hours} = 8760 \text{ hours/yr}$
- $(1000 \text{ disks} * 8760 \text{ hours/yr}) / 100,000 \text{ hours/failure} = 87.6 \text{ failed disks per year}$ on average
- $87.6/1000 = 8.76\%$ annual failure rate
- Google's 2007 study* found that actual AFRs for individual drives ranged from 1.7% for first year drives to over 8.6% for three-year old drives

*research.google.com/archive/disk_failures.pdf

The "Bathtub Curve"

- Often failures follow the "bathtub curve"
- Brand new devices, higher prob of failure
- Old devices, higher prob of failure
- In between, lower
- Some manufacturers will “burn in” products to eliminate early failures



Dependability Design Principle

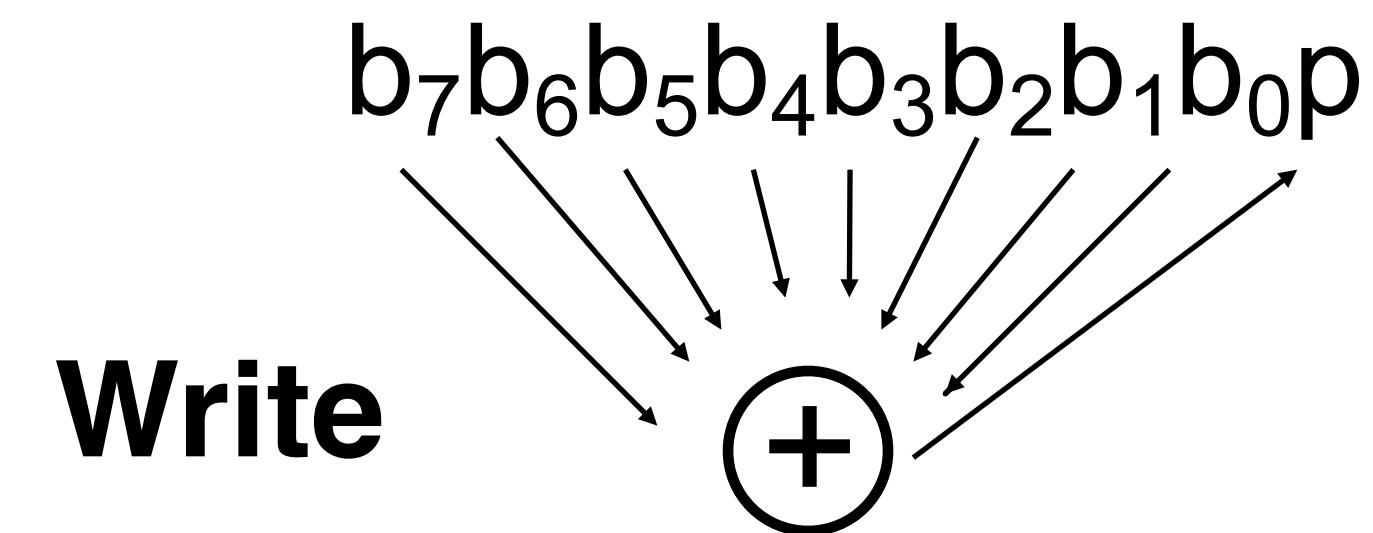
- Design Principle: No single points of failure
 - “Chain is only as strong as its weakest link”
- Dependability behaves like speedup of Amdahl’s Law
 - Doesn’t matter how dependable you make one portion of system
 - Dependability limited by part you do not improve

Error Detection/Correction Codes (EDC/ECC)

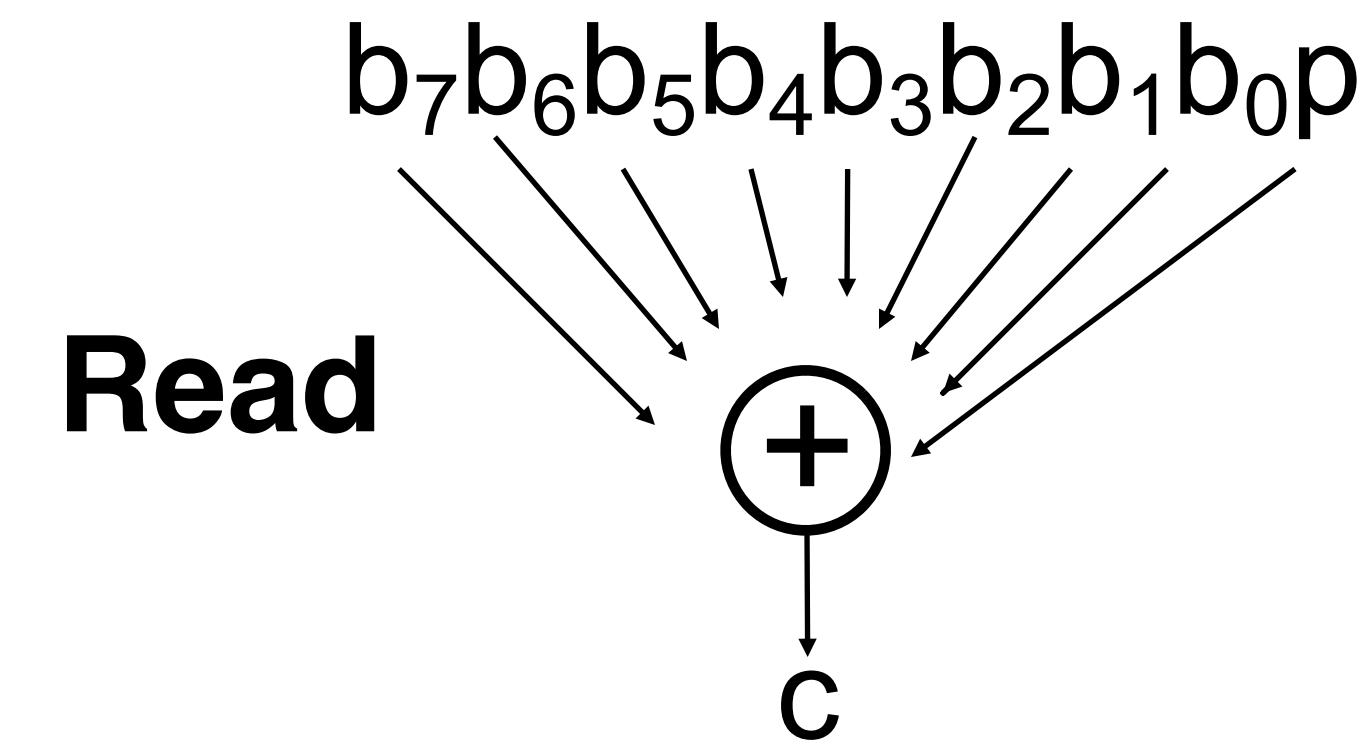
- Memory systems generate errors (accidentally flipped-bits)
 - DRAMs store very little charge per bit
 - “**Soft**” errors occur occasionally when cells are struck by alpha particles or other environmental upsets
 - “**Hard**” errors can occur when chips permanently fail
 - Problem gets worse as memories get denser and larger
- Memories protected against failures with EDC/ECC
- Extra bits are added to each data-word
 - Used to detect and/or correct faults in the memory system
 - Each data word value mapped to unique code word
 - A fault changes valid code word to invalid one, which can be detected

Simple Error Detection Coding: Parity Bit

- Each data value, before it is written to memory is “tagged” with an extra bit to force the stored word to have *even parity*:



- Each word, as it is read from memory is “checked” by finding its parity (including the parity bit).



- A non-zero parity check, c , indicates an error occurred:
 - two errors (on different bits) is not detected (nor any even number of errors)
 - odd numbers of errors are detected.

Parity Example

- $X = 0101\ 0101$
- 4 ones, *even* parity now
- Write to memory:
 $0101\ 0101\ 0$
to keep parity even
- $Y = 0110\ 0111$
- 5 ones, *odd* parity now
- Write to memory:
 $0110\ 0111\ 1$
to make parity even
- Read X from memory
 $0101\ 0101\ 0$
- 4 ones => even parity, so no error
- Read X from memory
 $1101\ 0101\ 0$
- 5 ones => odd parity, so error

What if error is parity bit?

Suppose Want to Correct Errors?

- Hamming came up with simple method using multiple parity bits to enable Error Detection and Correction
- Called “Hamming ECC”
 - Worked weekends on relay computer with unreliable card reader, frustrated with manual restarting
 - Got interested in error correction; published 1950
 - R. W. Hamming, “Error Detecting and Correcting Codes,” The Bell System Technical Journal, Vol. XXVI, No 2 (April 1950) pp 147-160.

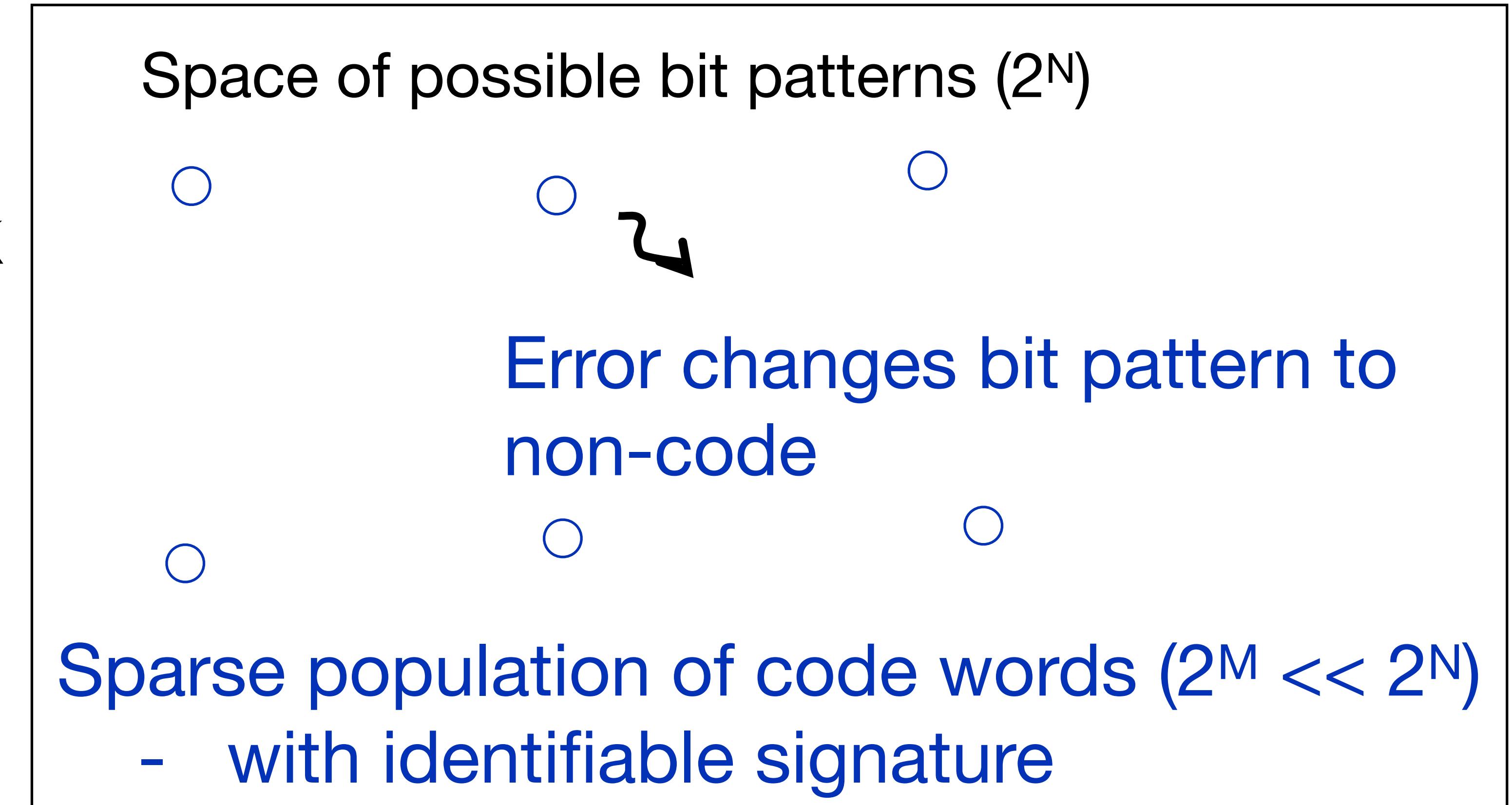


Richard Hamming
Turing Award Winner

Detecting/Correcting Code Concept

Use N-bits per symbol, and M-bits per valid code-word

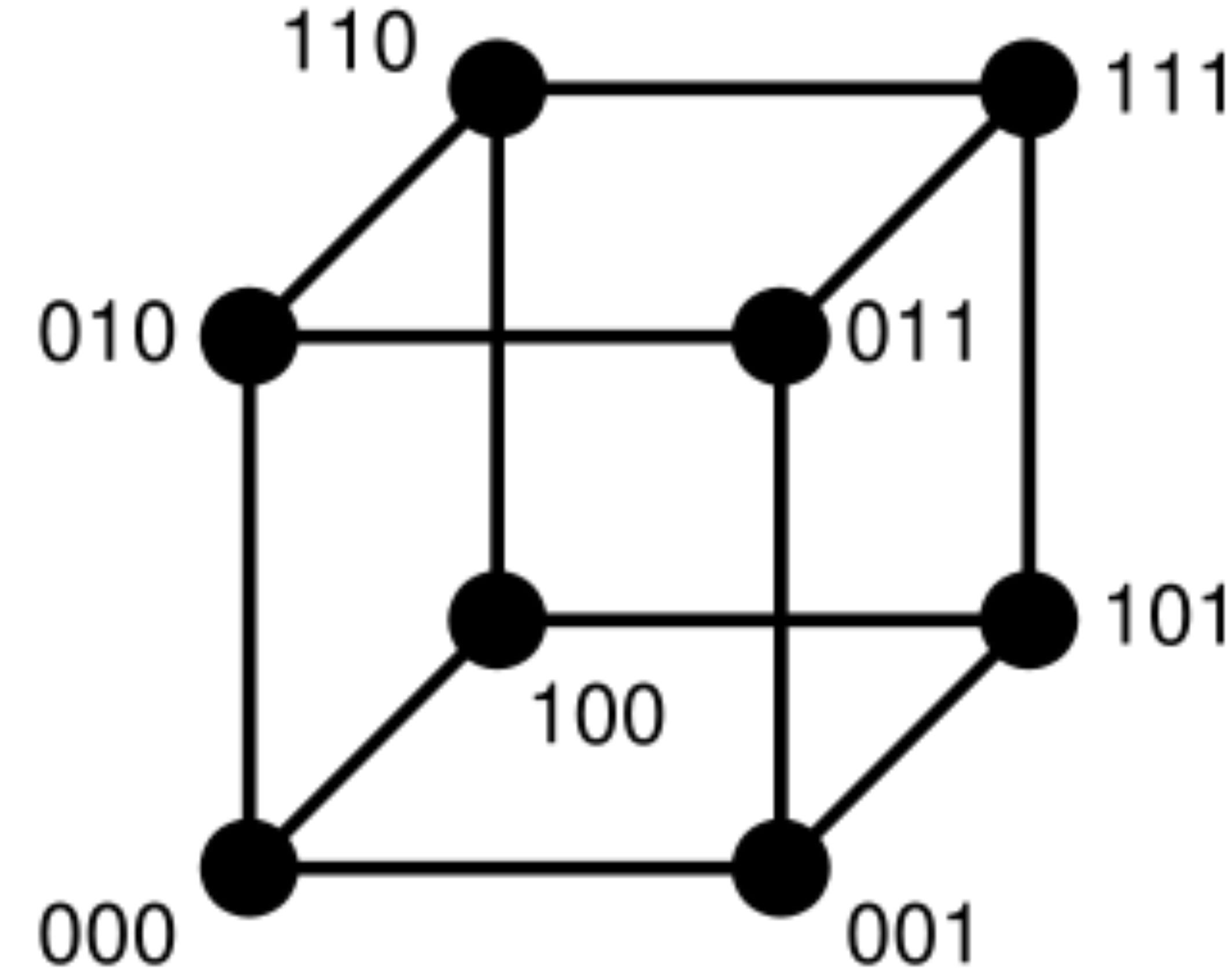
- **Detection:** bit pattern fails codeword check
- **Correction:** map to nearest valid code word



Block Code Principle

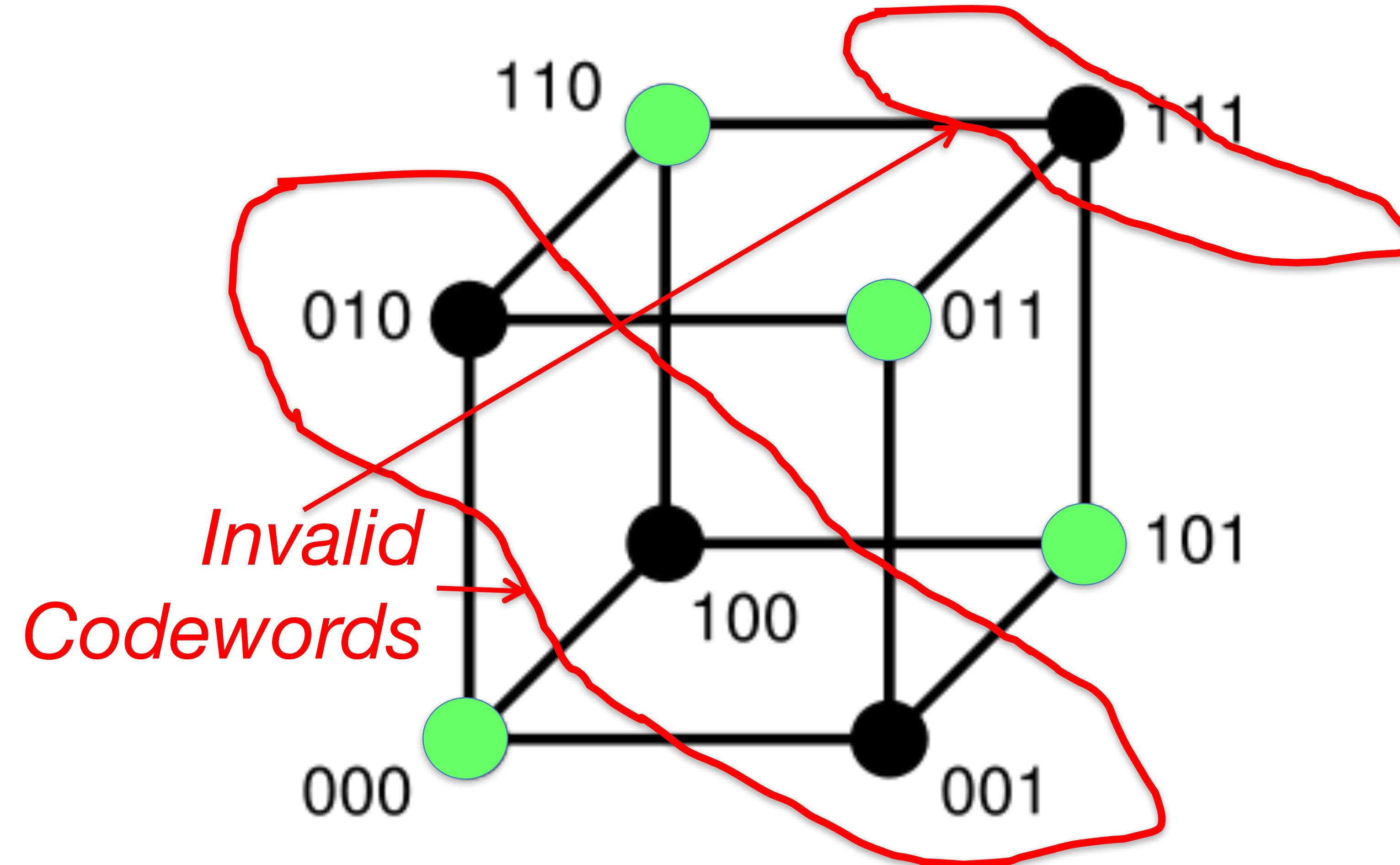
- Hamming distance = # of bits positions where words differ:
- $p = 0\underline{1}1011$, $q = 0\underline{0}1\underline{1}11$, Hamming distance $(p,q) = 2$
- $p = 011011$,
 $q = 110001$, Hamming distance $(p,q) = ?$
- Can think of adding extra bits to data to create a code word
- What if minimum distance between members of code is 2
and get a 1-bit error?

Hamming Distance 1: 8 code words



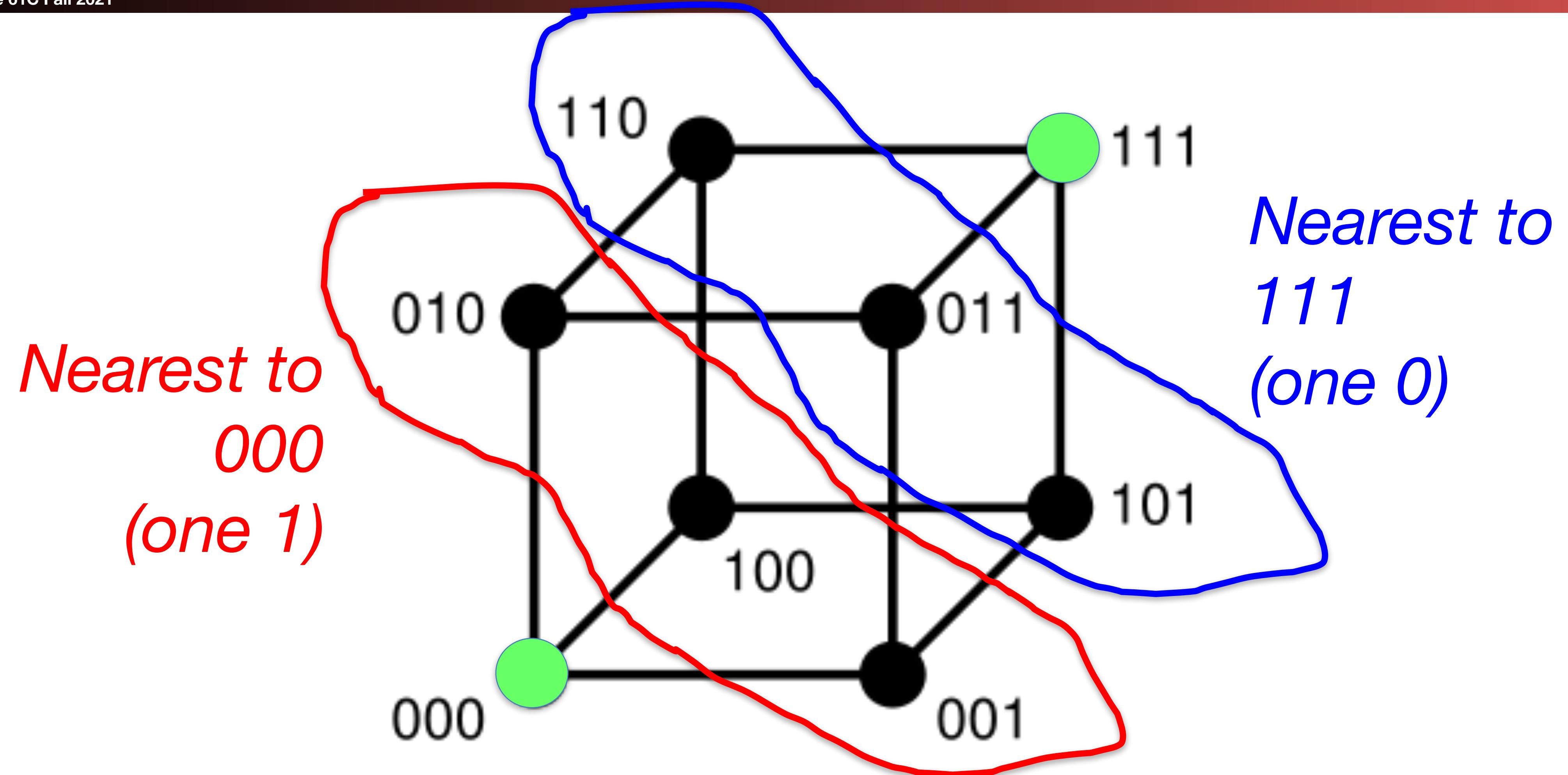
- Any single bit-error goes from valid codeword to another

Hamming Distance 2: Detect Single Bit-Errors



- No single bit-error goes to another valid codeword
- $\frac{1}{2}$ of symbols are valid codewords (in this case, simple parity)

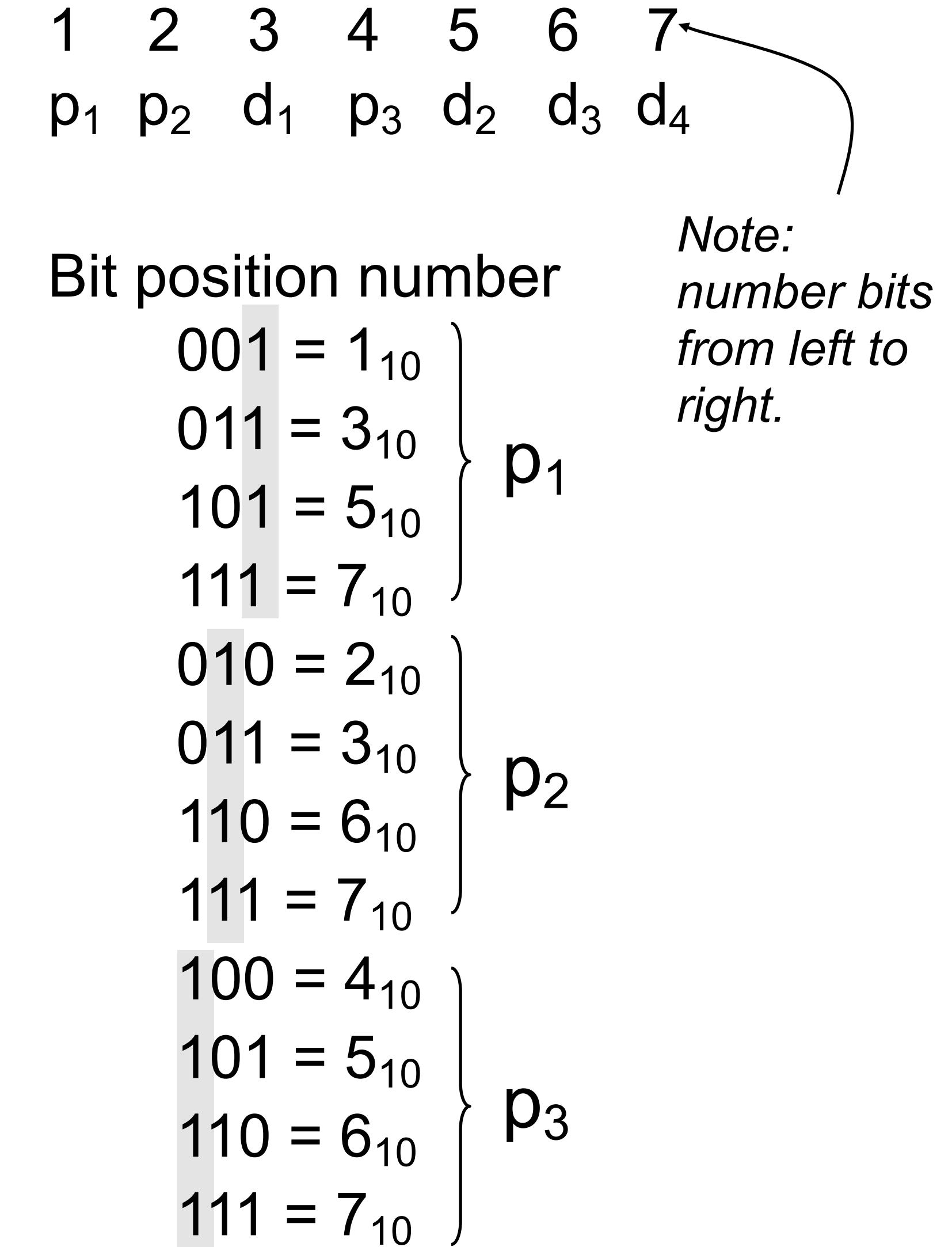
Hamming Distance 3: Correct Single Bit-Errors



- Any 1 bit-error can be corrected;
- 1/4 symbols are valid code-words
- Correct is voting in this case

Hamming Error Correcting Code

- Use more parity bits to pinpoint bit(s) in error, so they can be corrected.
- Example: Single error correction (SEC) on 4-bit data
 - use 3 parity bits, with 4-data bits results in 7-bit code word
 - 3 parity bits sufficient to identify any one of 7 code word bits
 - overlap the assignment of parity bits so that a single error in the 7-bit word can be corrected
- **Procedure:** group parity bits so they correspond to subsets of the 7 bits:
 - p_1 protects bits 1,3,5,7
 - p_2 protects bits 2,3,6,7
 - p_3 protects bits 4,5,6,7



Hamming Code Example

1	2	3	4	5	6	7
p ₁	p ₂	d ₁	p ₃	d ₂	d ₃	d ₄

- Note: parity bits occupy power-of-two bit positions in code-word.
- On writing to memory:
 - parity bits are assigned to force even parity over their respective groups.
- On reading from memory:
 - check bits (c_3, c_2, c_1) are generated by finding the parity of the group and its parity bit. If an error occurred in a group, the corresponding check bit will be 1, if no error the check bit will be 0.
 - check bits (c_3, c_2, c_1) form the position of the bit in error.

- Example: $c = c_3c_2c_1 = 101$
 - error in 4,5,6, or 7 (by $c_3=1$)
 - error in 1,3,5, or 7 (by $c_1=1$)
 - no error in 2, 3, 6, or 7 (by $c_2=0$)
- Therefore error must be in bit 5.
- *Note the check bits point to 5*
- By our clever positioning and assignment of parity bits, the check bits always address the position of the error!
- $c=000$ indicates no error

Hamming Error Correcting Code

- Overhead involved in single error correction code:
- let p be the total number of parity bits and d the number of data bits in a $p + d$ bit word.

If p error correction bits are to point to the error bit ($p + d$ cases) plus indicate that no error exists (1 case), we need:

$$2^p \geq p + d + 1,$$

$$\text{thus } p \geq \log(p + d + 1)$$

for large d , p approaches $\log(d)$

- Adding on extra parity bit covering the entire word can provide double error detection

1	2	3	4	5	6	7	8
p_1	p_2	d_1	p_3	d_2	d_3	d_4	p_4

- On reading, the C bits are computed (as usual) plus the parity over the entire word, P:

C=0 P=0, no error

C!=0 P=1, correctable single error

C!=0 P=0, a double error occurred

C=0 P=1, an error occurred in p_4 bit

Typical modern codes in DRAM memory systems:

64-bit data blocks (8 bytes) with 72-bit code words (9 bytes),
results in SEC, DED.



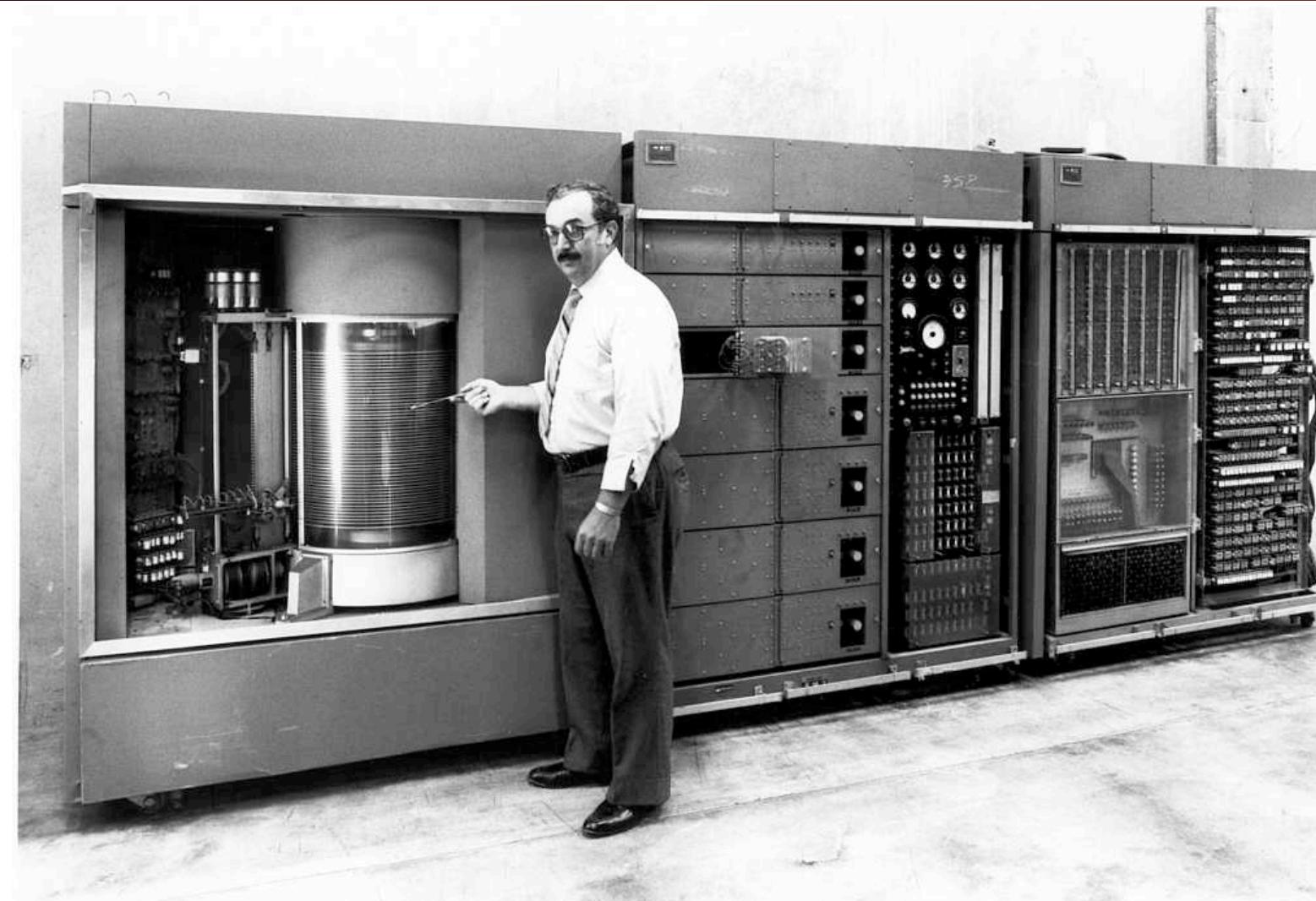
Other Codes

- **Many** different types of codes exist for applications in communications and data storage:
 - Cyclic Redundancy Code (CRC), for error detection Ethernet packets
 - “forward” error correction (FEC) codes in wireless telecomm and WiFi:
 - Low Density Parity Codes (LDPC), Verterbi/Turbo Codes, Polar Codes
 - Reed-Solomon Codes for communicating with Satellites billions of kilometers away, and for storing data on CDs (able to correct defects up to 1mm in size)
 - Cryptographic hash-functions, to guard against malicious tampering (e.g. SHA256)
 - An attacker is not be able to change, add, or remove any bits without changing the hash output
 - All use redundancy (extra bits) to represent or summarize the data

RAID: Redundancy for Disks

- Spinning disk is still a critical technology
 - Although worse latency than SSD...
- Disk has equal or greater bandwidth and an order of magnitude better storage density (bits/cm^3) and cost density ($\text{bits}/\$$)
- So when you need to store a petabyte or three...
 - You need to use disk, not SSDs
 - Oh, and SSDs can fail too

Evolution of the Disk Drive



IBM RAMAC 305, 1956

First commercial computer that used a moving-head hard disk drive (magnetic disk storage) for secondary storage.



up to 22.7 billion bytes
(gigabytes) of storage

IBM 3390K, 1989

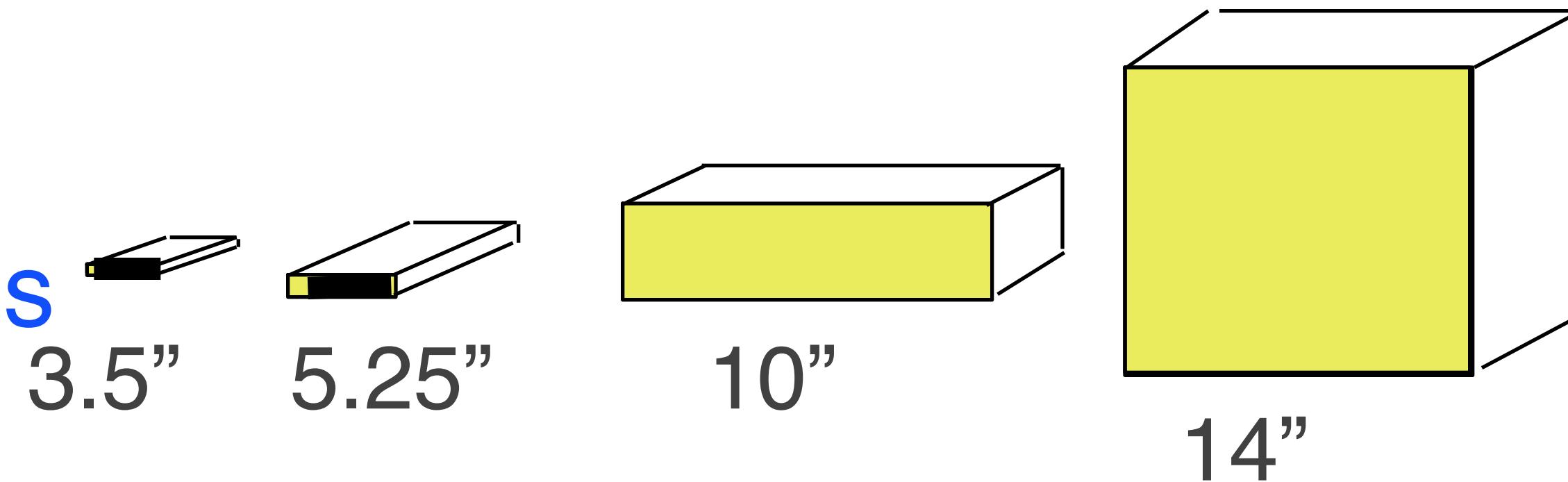


Apple SCSI, 1986

Arrays of Small Disks

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

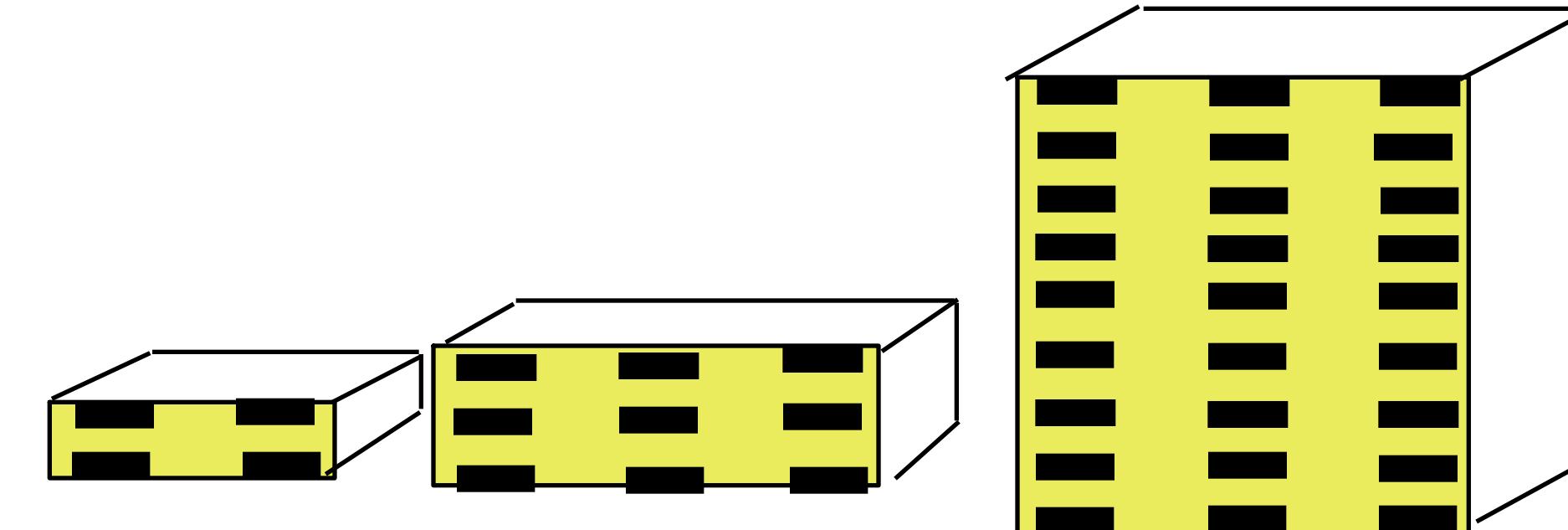
Conventional:
4 disk designs



Low End → High End

Disk Array:
1 disk design

3.5" →



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.	9X
Power	3 KW	11 W	1 KW	3X
Data Rate	15 MB/s	1.5 MB/s	105 MB/s	7X
I/O Rate	600 I/Os/s	55 I/Os/s	3900 IOs/s	6X
MTTF	250 K Hrs	50 K Hrs	??? 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?

But MTTF goes through the roof...

- If 1 disk has MTTF of 50k hours...
 - 70 disks will have a MTTF of ~700 hours!!!
 - This is assuming failures are independent...
- But fortunately we know when failures occur!
 - Disks use a lot of CRC coding, so we don't have corrupted data, just no data
- We can have both “Soft” and “Hard” failures
 - Soft failure just the read is incorrect/failed, the disk is still good
 - Hard failures kill the disk, necessitating replacement
 - Most RAID setups are “Hot swap”: Unplug the disk and put in a replacement while things are still going
 - Most modern RAID arrays also have “hot spares”: An already installed disk that is used automatically if another disk fails.

RAID: Redundant Arrays of (Inexpensive) Disks

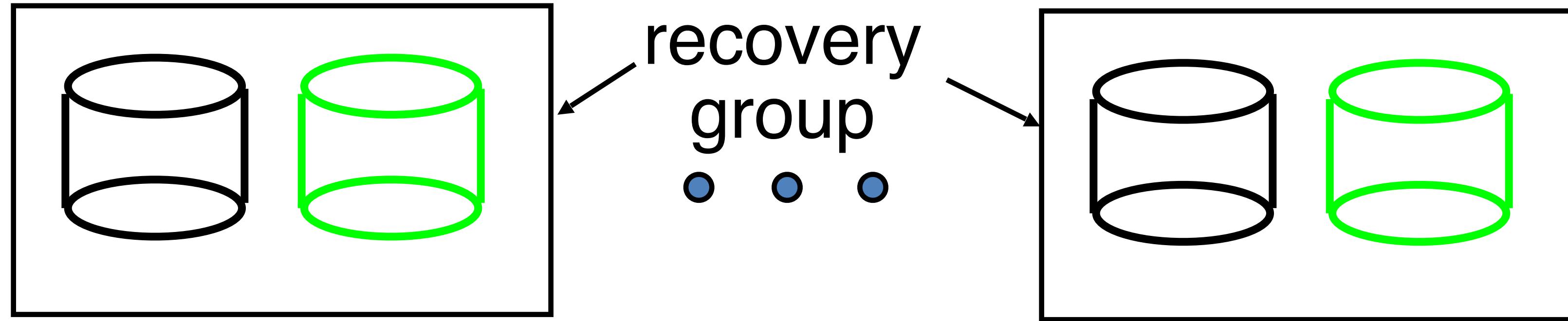
- Files are "striped" across multiple disks, ex:
- Redundancy yields high data availability
 - Availability: service still provided to user, even if some components failed
- Disks will still fail
- Contents reconstructed from data redundantly stored in the array
 - Capacity penalty to store redundant info
 - Bandwidth penalty to update redundant info on writes
- 6 Raid *Levels*, 0, 1, 5, 6 most common today

RAID 0: Striping

- "RAID 0" is not actually RAID (no redundancy)
 - It is simply spreading the data across multiple disks
 - So, e.g, for 4 disks, stripe-unit address 0 is on disk 0, address 1 is on disk 1, address 2 is on disk 2, address 4 on disk 0...
- Improves bandwidth linearly
 - With 4 disks you have 4x the disk bandwidth
- Doesn't really help latency
 - Still have the individual disks seek and rotation time
- Failures will happen...

You can set the stripe-unit size of an IBM® SAS Disk Array to 16 KB, 64 KB, 256 KB, or 512 KB.

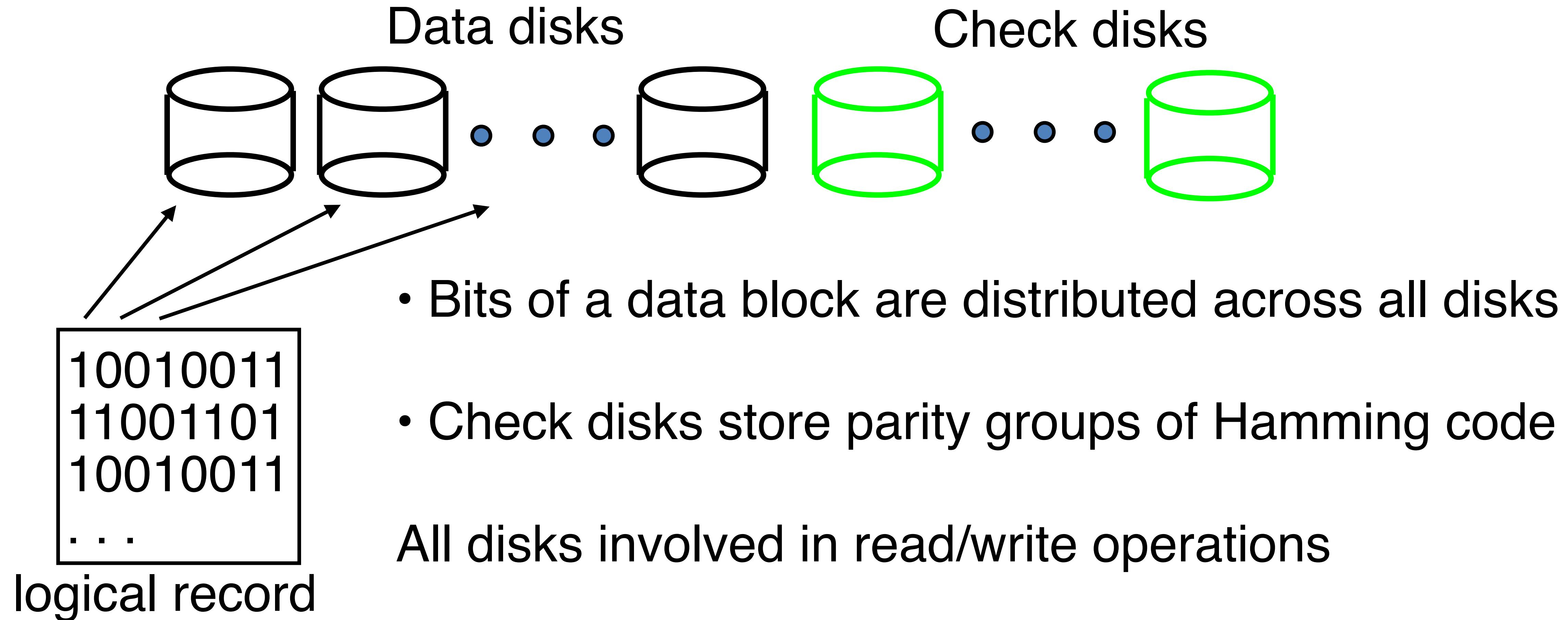
RAID 1: Disk Mirroring/Shadowing (online sparing)



- Each disk is fully duplicated onto its “mirror”
Very high availability can be achieved
- Writes go to disk and mirror - limited by single-disk speed
- Reads from original disk, unless failure

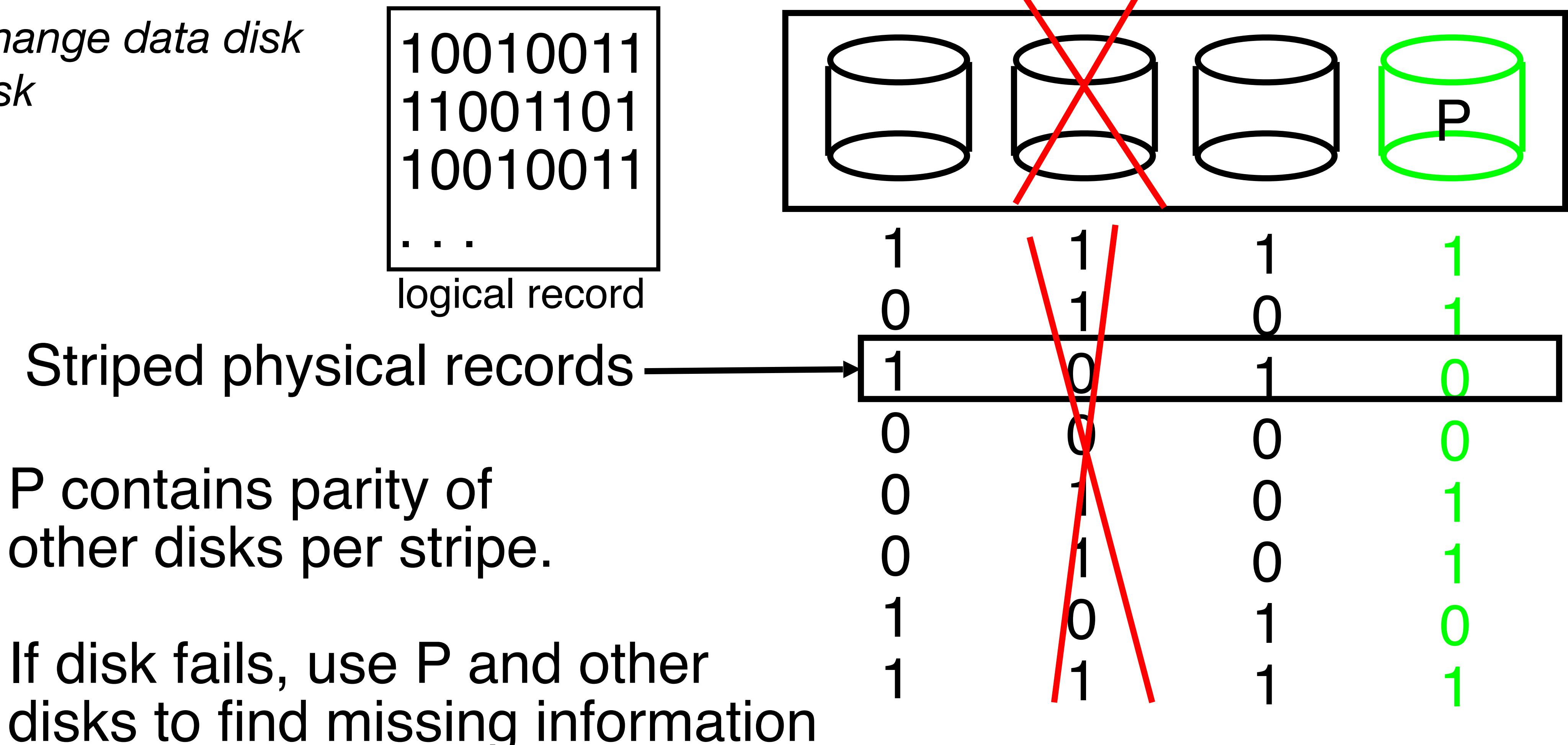
Most expensive solution: 100% (2x) capacity overhead

RAID 2: Hamming Code for Error Correction



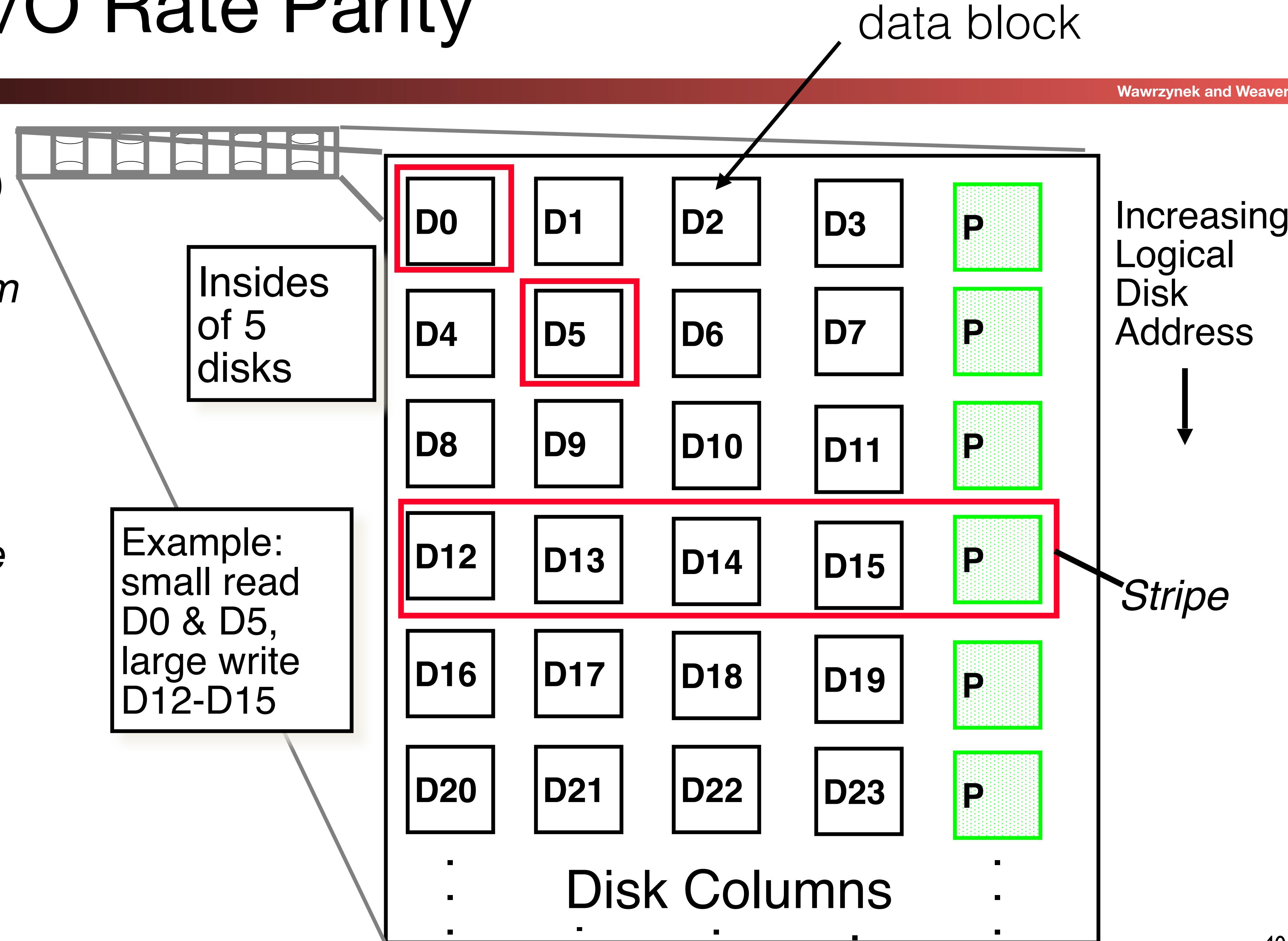
RAID 3: Single Parity Disk

- Disk drives themselves code data and detect failures
- Reconstruction of data can be done with single parity disk if we know which disk failed
- Writes change data disk and P disk



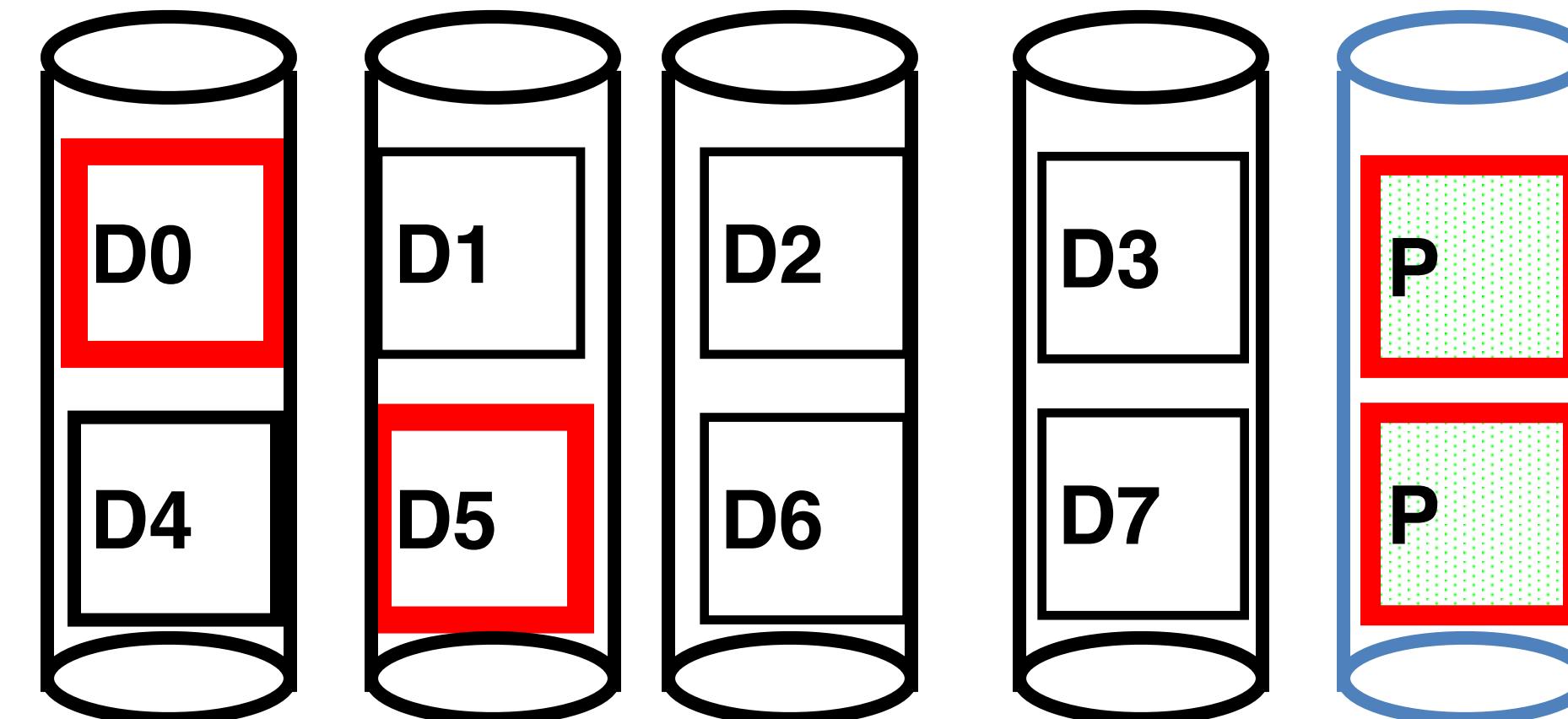
RAID 4: High I/O Rate Parity

- Interleave data at the sector level (data block) rather than bit level - permits more parallelism (independent small reads, parallel large reads)
- Reconstruction of data can be done with single parity disk
- Reading (w/o) fault involve only data disk
- Writing involves data disk + parity disk

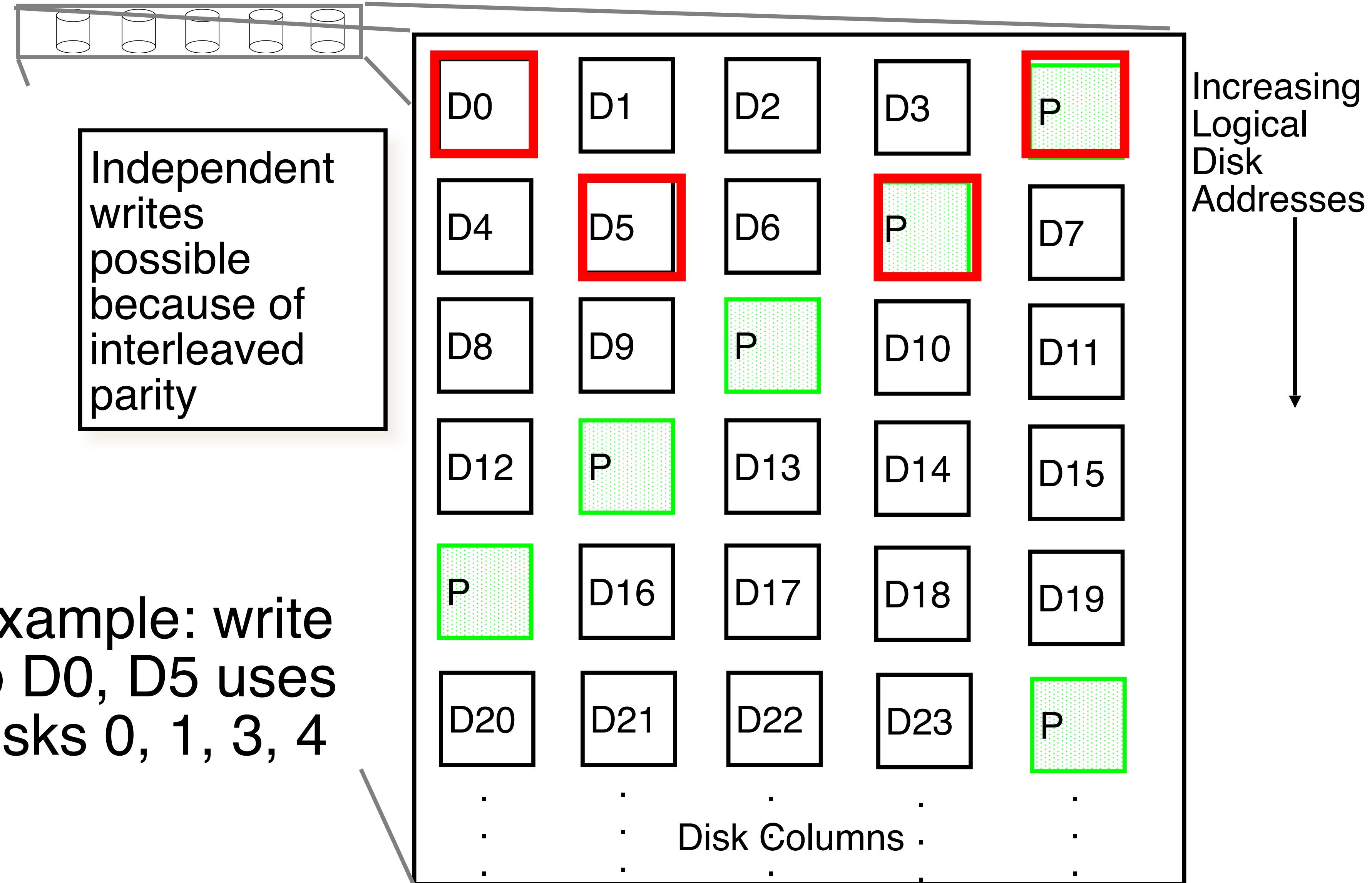


Inspiration for RAID 5

- RAID 4 works well for reads, but
- Parity Disk is the bottleneck for writes: Write to D0, D5 both also write to P disk



RAID 5: High I/O Rate Interleaved Parity



RAID 6

- RAID 5 is no longer the “gold standard”
- Can experience 1 disk failure and continue operation
 - RAID array is in a “degraded” state
- But disk failures are not actually independent!
 - When one disk has failed, there’s a decent chance another will fail soon
- RAID 6: Add another parity block per stripe
 - Now 2 blocks per stripe rather than 1
 - Sacrifice capacity for increased redundancy
 - Now the array can tolerate **2** disk failures and continue operating

Berkeley's Role in Definition of RAID (December 1987)

A Case for Redundant Arrays of Inexpensive Disks (RAID)

The screenshot shows a Google Scholar search results page. The search query 'Case for Raid' is entered in the search bar. The results section is titled 'Scholar' and indicates 'About 138,000 results (0.08 sec)'. The top result is a book titled '[book] A case for redundant arrays of inexpensive disks (RAID)' by DA Patterson, G Gibson, and RH Katz, published in 1988 at dl.acm.org. The abstract discusses the need for RAID to match the performance increase of CPUs and memories. It mentions that while the capacity of Single Large Expensive Disks (SLED) has grown rapidly, the performance improvement of SLED has been modest. The paper introduces five levels of RAID and compares them to an IBM 3380 and a Fujitsu Super Eagle. The page also includes filters for 'Articles', 'Legal documents', and 'Any time'.

Google Case for Raid

Scholar About 138,000 results (0.08 sec)

Articles Legal documents Any time

[book] [A case for redundant arrays of inexpensive disks \(RAID\)](#)
DA Patterson, G Gibson, RH Katz - 1988 - dl.acm.org

Abstract Increasing performance of CPUs and memories will be squandered if not matched by a similar performance increase in I/O. While the capacity of Single Large Expensive Disks (SLED) has grown rapidly, the performance improvement of SLED has been modest. ...

Cited by 2814 Related articles All 239 versions Cite More▼

Increasing performance of CPUs and memories will be squandered if not matched by a similar performance increase in I/O. While the capacity of Single Large Expensive Disk (SLED) has grown rapidly, the performance improvement of SLED has been modest. Redundant Arrays of Inexpensive Disks (RAID), based on the magnetic disk technology developed for personal computers, offers an attractive alternative to SLED, promising improvements of an order of magnitude in performance, reliability, power consumption, and scalability.

This paper introduces five levels of RAID, giving their relative cost/performance, and compares RAIDs to an IBM 3380 and a Fujitsu Super Eagle.

RAID Version 1

- RAID-I (1989)
 - Consisted of a Sun 4/280 workstation with 128 MB of DRAM, four dual-string SCSI controllers, 28 5.25-inch SCSI disks and specialized disk striping software



RAID Version 2

- 1990-1993
- Early Network Attached Storage (NAS) System running a Log Structured File System (LFS)
- Impact:
 - \$25 Billion/year in 2002
 - Over \$150 Billion in RAID device sold since 1990-2002
 - 200+ RAID companies (at the peak)
 - Software RAID a standard component of modern OSs



RAID Is Not Enough By Itself

- You don't just have one disk die...
 - You can have more die in a short period of time
 - Thank both the "bathtub curve" and common environmental conditions
- If you care about your data, RAID isn't sufficient
 - You need to also consider a separate backup solution
- A good practice in clusters/warehouse scale computers:
 - RAID-6 in each cluster node with auto-failover and a hot spare
 - Distributed filesystem on top
 - Replicates amongst the cluster nodes so that nodes can fail
 - And then distribute to a different WSC...

In Conclusion ...

- We have methods to mitigate faults in electronic systems:
 - Design bugs, Manufacturing defects, and Runtime Faults
- Dependability Measures let us quantify
- Dealing with Runtime Faults requires redundancy
 - either more hardware (cost) or more time (performance)
- Redundancy most commonly used in memory systems (DRAM, SRAM, Disks, SSD), also for communications