

only if a miss occurs in both caches.

With a two-level cache, central memory has to be accessed

on the same chip as the CPU.

which is often located off-chip, whereas the L1 cache is located

in the secondary cache.

Most high-performance microprocessors include an L2 cache

the secondary cache is referred to as the L2 (level 2) cache.

The primary cache is referred to as the L1 (level 1) cache and

response to a miss in the primary cache.

One way in which this penalty can be reduced is to provide

another cache, the secondary cache, which is accessed in

one way for a cache miss is the extra time that it takes to

obtain the requested item from central memory.

The penalty for a cache miss is the extra time that it takes to

Multi-level Caches

- The primary cache is referred to as the L1 (level 1) cache and
- The secondary cache is referred to as the L2 (level 2) cache.
- One way in which this penalty can be reduced is to provide another cache, the secondary cache, which is accessed in one way for a cache miss is the extra time that it takes to obtain the requested item from central memory.
- The penalty for a cache miss is the extra time that it takes to obtain the requested item from central memory.

Example:

A computer system employs a write-back cache with a 70% hit ratio for writes. The cache operates in look-side mode and has a 90% read hit ratio. Reads account for 80% of all memory references and writes account for 20%. If the main memory cycle time is 200ns and the cache access time is 20ns, what would be the average access time for all references (reads as well as writes)?

$$\text{The average write time} = 0.7 * 20\text{ns} + 0.3 * 200\text{ns} = 74\text{ns}$$
$$\text{The average access time for reads} = 0.9 * 20\text{ns} + 0.1 * 200\text{ns} = 38\text{ns}.$$
$$\text{Hence the overall average access time for combined reads and writes is}$$
$$0.8 * 38\text{ns} + 0.2 * 74\text{ns} = 45.2\text{ns}$$

The average write time = $0.7 * 20\text{ns} + 0.3 * 200\text{ns} = 74\text{ns}$

Example:

used for pipelineing processes – the instruction processor and the execution unit –

- Split caches eliminate contention for cache between
- Between instruction and data
- Higher hit rate for unified cache as it balances
- Same cache for instruction and data
- Unified cache
- D-cache (data) – mostly random access
- Cache (instruction) – mostly accessed sequentially
- Separate caches for instructions and data
- Split cache
- Unified cache

Cache organization

- Conflict Misses: If the cache mapping is such that multiple blocks are mapped to the same cache entry
- Capacity Misses: If the cache cannot contain all the blocks needed during the execution of a program, capacity misses will occur due to blocks being discarded and later retrieved.
- Cold Misses: The very first access to a block will result in a miss because the block is not brought into cache until it is referenced.
- Compulsory Misses: These are misses that are caused by the cache being empty initially.

Sources of Cache Misses

- Conflict Misses
- Cold Misses
- Capacity Misses
- Compulsory Misses

Sources of Cache Misses

$$\begin{aligned} \Rightarrow h &= 1190/1200 \\ \Rightarrow 1-h &= 10/1200 \\ \Rightarrow 0.1 * 100 &= (1-h) * 1200 \\ \Rightarrow 0.1T_c &= (1-h) * T_m \\ \Rightarrow 1.1T_c &= T_c + (1-h) * T_m \\ \Rightarrow T^a &= T_c + (1-h) * T_m \end{aligned}$$

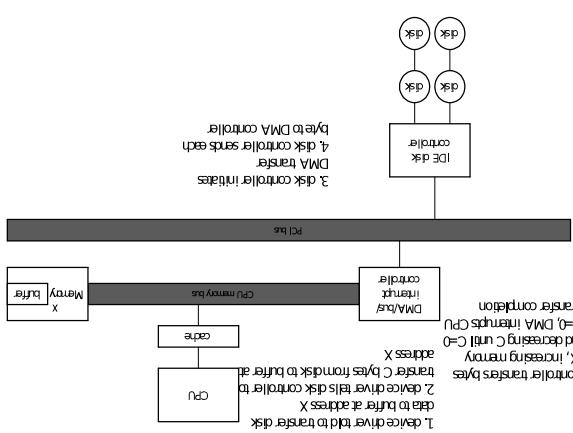
Look-through cache?

than the cache access time, what is the hit ratio H in 1200ns, if the effective access time is 10% greater

• Consider a memory system with $T_c = 100\text{ns}$ and $T_m =$

Problem

Introduction of Direct Memory Access (DMA)



DMA Progress (cont.)

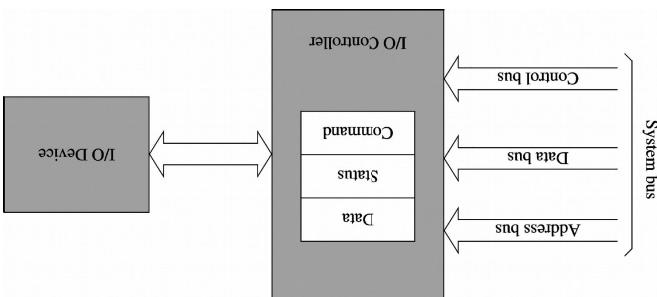
- The DMA controller proceeds to operate the memory bus directly without CPU help.
 - Handshaking exists between DMA controller and device controller.
 - When the entire transfer is finished, the DMA controller will interrupt the CPU.

DMA Progress

- To initiate a DMA transfer, the host writes a DMA command into memory:
 - A pointer to the source of a transfer
 - A count of the number of bytes to be transferred
 - The CPU writes the address of the DMA command block to the DMA controller.

Why is DMA?

- It is wasteful to feed data into a controller register 1 bytes at a time. (PIO)
 - The DMA unit is word.
 - In the high loading environment, a system with DMA has better improvement.



Introduction (cont'd)

- I/O devices serve two main purposes
 - * To communicate with outside world
 - * To store data
- I/O controller acts as an interface between the systems bus and I/O device
 - * Relieves the processor of low-level details
 - * Takes care of electrical interface
 - * Stores care of low-level details
- I/O controllers have three types of registers
 - * Data
 - * Command
 - * Status

Introduction

- Extreme interface
 - * Parallel interface
 - * Serial transmission
 - * USB
 - * I/O data transfer
 - * Keyboard
 - * Motivation
 - * USB architecture
 - * IEEE 1394
 - * Error detection and correction
 - * Parity encoding
 - * Transactions
 - * Advantages
 - * IEEE 1394
 - * Configuration
 - * CRC
 - * Bus arbitration
 - * Transaction
 - * Error correction
 - * Transceiving
 - * Advantages
 - * IEEE 1394

Outline

Input/Output Organization

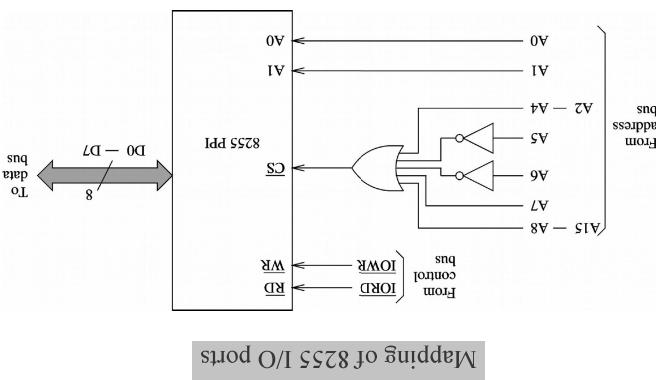
Chapter 19

- DMA-request and DMA-acknowledge
 - When the device controller receives the DMA-acknowledge signal, it transfers the word of data to memory - and remove the DMA-request signal.

Handshaking (cont.)

- DMA-request and DMA-acknowledge
 - When a word of data is available, the device controller places a signal on the DMA-request wire.
 - The signal causes the DMA controller to seize the memory bus.
 - To place the desired address on the memory-address wire.
 - To place a signal on the DMA-acknowledge wire

Handshaking



An Example I/O Device (cont'd)

- * Keyboard
 - ↳ Key's scan code depends on its position on the keyboard
 - ↳ Scan code is like a sequence number of the key
 - ↳ Supports key identity as a scan code
 - ↳ Key depressions and releases
- * Keyboard controller scans and reports
- * Interface through an 8-bit parallel I/O port
 - ↳ No relation to the ASCII value of the key
- * Originally supported by 8255 programmable peripheral interface chip (PPI)

An Example I/O Device

- * I/O address mapping
 - ↳ Memory-mapped I/O
 - ↳ Reading and writing are similar to memory read/write
 - ↳ Uses same memory read and write signals
 - ↳ Most processors use this I/O mapping
- * Isolated I/O
 - ↳ Separate I/O read and write signals are needed
 - ↳ Pentium supports isolated I/O
 - ↳ 64 KB address space
 - ↳ Key depressions and releases
- * Ports
 - ↳ Can be any combination of 8-, 16- and 32-bit I/O
 - ↳ Also supports memory-mapped I/O

Accessing I/O Devices

Port address	Port address
PA (input port)	60H
PB (output port)	61H
PC (input port)	62H
Command register	63H

* These ports are mapped as follows

- » Port A (PA)
- » Port B (PB)
- » Port C (PC)

• 8255 PPI has three 8-bit registers

An Example I/O Device (cont'd)

- * Register/I/O instructions
 - ↳ DX gives the port address
 - ↳ DX prefix to access first 256 ports
 - ↳ Both take no operands--as in string instructions
 - ↳ ins and outs
- * Block I/O instructions
 - ↳ DX gives the port address
 - ↳ ins accumulator, port8 : direct format
 - ↳ ins accumulator, DX : indirect format
- * Accessing I/O ports in Pentium
 - » Three types
 - ↳ Direct memory access (DMA)
 - ↳ Programmed I/O
 - ↳ Interrupt-driven I/O
 - » Two basic ways
 - ↳ This access depends on I/O mapping
 - ↳ Access to various registers (data, status,...)
 - » A protocol to communicate (to send data, ...)

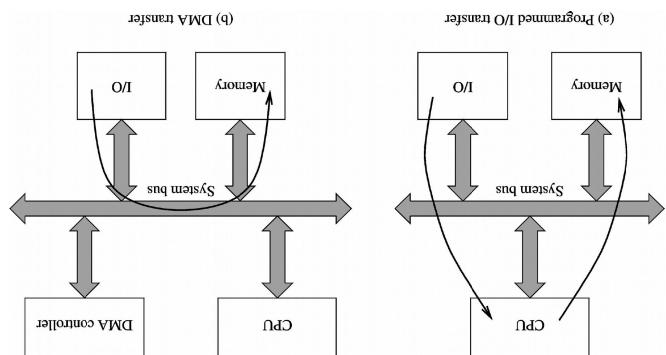
Accessing I/O Devices (cont'd)

- * To communicate with an I/O device, we need
 - ↳ Memory-mapped I/O
 - ↳ Isolated I/O
 - ↳ DMA
 - ↳ Programmed I/O
 - ↳ Interrupt-driven I/O
- * Access to various registers (data, status,...)
 - ↳ This access depends on I/O mapping
 - ↳ Two basic ways
 - ↳ Access to various registers (data, status,...)
- * A protocol to communicate (to send data, ...)

Introduction (cont'd)

- DMA is implemented using a DMA controller
 - * DMA controller
 - » Acts as slave to processor
 - » Receives instructions from processor
 - » Example: Reading from an I/O device
 - Processor gives details to the DMA controller
 - » I/O device number
 - » Main memory buffer address
 - » Number of bytes to transfer
 - » Direction of transfer (memory \rightarrow I/O device, or vice versa)

I/O Data Transfer (cont'd)



I/O Data Transfer (cont'd)

- DMA
 - * DMA devices
 - Disk read or write
 - » Waiting for it to be released
 - » Waiting for a key to be pressed, - In our example
 - » Processor wastes time polling
 - » Frees the processor of the data transfer responsibility
- Direct memory access (DMA)
 - * Problems with programmed I/O
 - » Waiting for PATA bit to go low
 - » Reading a key from the keyboard involves
 - Indicating that a key is pressed
 - » Translating the key scan code to the ASCII value
 - » Waiting until the key is released
 - » Programmed I/O uses this process to read input from the keyboard

I/O Data Transfer (cont'd)

- Example
 - * Done by busy-waiting
 - » This process is called *polling*
 - * Reading a key from the keyboard involves
 - Indicating that a key is pressed
 - » Waiting for PATA bit to go low
 - » Reading a key scan code to the ASCII value
 - » Waiting until the key is released
 - » Programmed I/O uses this process to read input from the keyboard

I/O Data Transfer (cont'd)

- Mapping I/O ports is similar to mapping memory
 - * Partial mapping
 - » It can be done either by
 - DMA
 - Programmed I/O
 - An end-notification phase
 - » Interrupt
 - Three basic techniques
 - * Programmed I/O
 - * DMA
 - * Interrupt-driven I/O (discussed in Chapter 20)

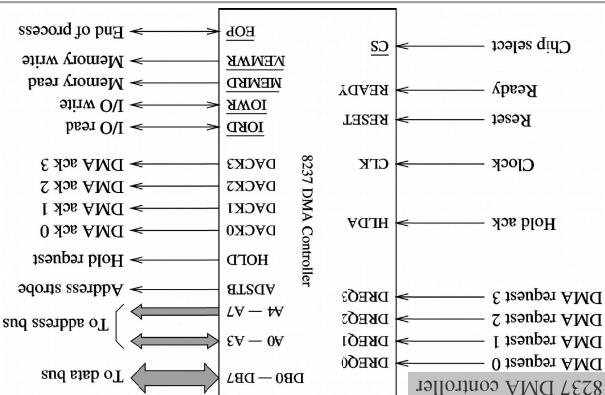
I/O Data Transfer

- Keyboard scan code is available from port 60H
 - * Key status is available from PA7
 - » PA7 = 0 – key depressed
 - » PA7 = 1 – key released
 - * 7-bit scan code is available from PA0 – PA6
 - » See our discussion in Chapter 16
- Full mapping
 - * Partial mapping
 - » It can be done either by
 - DMA
 - Programmed I/O

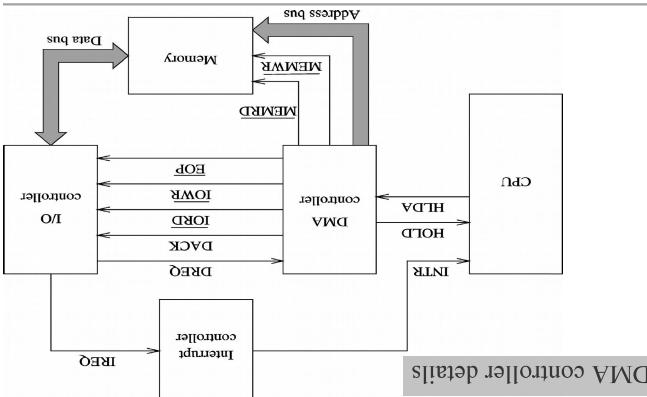
An Example I/O Device (cont'd)

- * Mode register
 - Each channel can be programmed to
 - Read or write
 - Automation or autodecrement the address
 - Automatize the channel
 - * Reduced register
 - For software-initiated DMA
 - * Mask register
 - Mask register
 - » Used to disable a specific channel
 - * Status register
 - Used for memory-to-memory transfers
 - * Temporary register
 - Used for memory-to-memory transfers

I/O Data Transfer (cont'd)



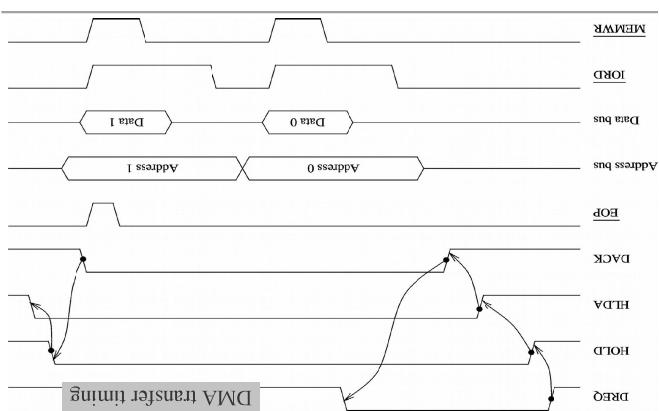
I/O Data Transfer (cont'd)



I/O Data Transfer (cont'd)

- **8237** supports four DMA channels
 - * Mode register
 - » Each channel can be programmed to
 - Read or write
 - Automatic or auto-decrement
 - Autoinitialize the channel
 - * Request register
 - » For software-initiated DMA
 - » Used to disable a specific channel
 - * Mask register
 - » Used to disable a specific channel
 - * Status register
 - » Used for memory-to-memory transfers
 - * Temporary register
 - » Used for priority-to-memory transfers
 - It has the following internal registers
 - * Current address register
 - » Holds address for the current DMA transfer
 - » One 16-bit register for each channel
 - * Current word register
 - » Holds address for the current DMA transfer
 - » Keeps the byte count
 - » Generates terminal count (TC) signal when the count goes from zero to FFFF
 - * Command register
 - » Used to program 8237 (type of priority, ...)

I/O Data Transfer (cont'd)



I/O Data Transfer (cont'd)

- * Steps in a DMA operation
 - * Processor initiates the DMA controller
 - * Gives device number, memory buffer pointer, ...
 - Called *channel initialization*
 - » Once initialized, it is ready for data transfer
 - * When ready, I/O device informs the DMA controller
 - DMA controller starts the data transfer process
 - Observes bus by going through bus arbitration
 - Places memory address and releases the bus
 - Completes transfer and appropriate control signals
 - Releases memory address and repeats the bus
 - Updates memory address and count value
 - If more to read, loops back to repeat the process
 - * Notify the processor when done
 - » Typically uses an interrupt

I/O Data Transfer (cont'd)

- Parity encoding
 - Add a parity bit such that the total number of 1s is odd (odd-parity)
 - Adds rudimentary error detection capability
 - Simplest mechanism
- Error detection and correction
 - How is error bit position computed?
 - Suppose P_1, P_2, \dots, P_8 are sum of weights of check bits in error
 - Error bit position = sum of weights of check bits in error
 - Suppose P_1, P_2, \dots, P_8 are sum of weights of check bits in error
 - Write the numbers $1, 2, 3, \dots, 8$ in binary
 - 0: 0000
 - 1: 0001
 - 2: 0010
 - 3: 0011
 - 4: 0100
 - 5: 0101
 - 6: 0110
 - 7: 0111
 - 8: 1000
 - Check those bit positions for which the second rightmost column has 1 (i.e., with weight $2^1 = 2$)
 - 4: 0100
 - 5: 0101
 - 6: 0110
 - 7: 0111
 - 8: 1000
 -
 - What is the logic?
 - Suppose P_1, P_2, \dots, P_8 are sum of weights of check bits in error
 - PI checks those bit positions for which the second rightmost column has 1 (i.e., with weight $2^1 = 2$)
 - 4: 0100
 - 5: 0101
 - 6: 0110
 - 7: 0111
 - 8: 1000
 - PI checks those bit positions for which the third rightmost column has 1 (i.e., with weight $2^2 = 4$)
 - 2: 0010
 - 3: 0011
 - 4: 0100
 - 5: 0101
 - 6: 0110
 - 7: 0111
 - 8: 1000
 - PI checks those bit positions for which the fourth rightmost column has 1 (i.e., with weight $2^3 = 8$)
 - 0: 0000
 - 1: 0001
 - 2: 0010
 - 3: 0011
 - 4: 0100
 - 5: 0101
 - 6: 0110
 - 7: 0111
 - 8: 1000

Error Detection and Correction (cont'd)

- I/O Data Transfer (cont'd)
 - 8237 supports four types of data transfer
 - Single cycle transfer
 - Only single transfer takes place
 - Useful for slow devices
 - Block transfer mode
 - Transfers data until TC is generated or external EQD signal is received
 - Similar to the block transfer mode
 - Demand transfer mode
 - In addition to TC and EQP, transfer can be terminated by deactivating DRQ signal
 - Cascade mode
 - Useful to expand the number of channels beyond four

- Parity encoding
 - Add a parity bit such that the total number of 1s is even (even-parity)
 - Simple to implement
 - Advantage:
 - Disadvantage
 - Cannot detect even number of bit errors
 - Can be used to detect single-bit errors
- Error detection and correction
 - Constructing codewords to correct single-bit errors
 - Count bit positions from left to right starting from 1
 - Parity bits are called **check bits**
 - Example for 8-bit data ($d = 8$)
 - Parity bits are in positions that are a power of 2
 - Codeword is 12 bits long
 - We need 4 check bits
 - Example for 8-bit data ($d = 8$)
 - Parity bits are in positions that are a power of 2
 - Codeword is 12 bits long

Error Detection and Correction

- I/O Data Transfer (cont'd)
 - 8237 supports four types of data transfer
 - Single cycle transfer
 - Only single transfer takes place
 - Useful for slow devices
 - Block transfer mode
 - Transfers data until TC is generated or external EQD signal is received
 - Similar to the block transfer mode
 - Demand transfer mode
 - In addition to TC and EQP, transfer can be terminated by deactivating DRQ signal
 - Cascade mode
 - Useful to expand the number of channels beyond four

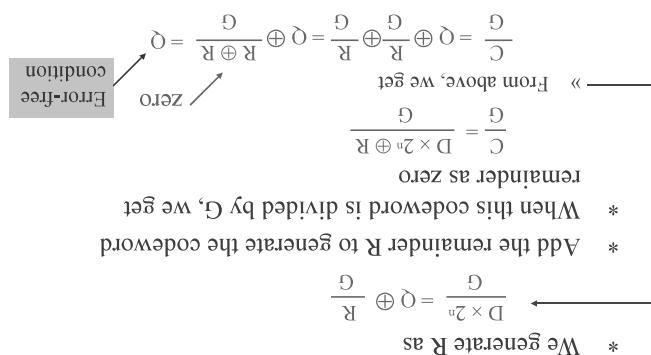
- Error correction
 - To correct single-bit errors in d data bits
 - Add p parity bits
 - Codeword $C = d + p$ bits
 - Depends on d
 - Hamming distance between codewords
 - Number of bit positions in which the two codewords differ
 - Hamming distance between codewords
 - Each check bit is responsible for checking certain bits
 - PI: checks bits 1, 3, 5, 7, 9, and 11
 - » Uses even parity
 - P2: checks bits 2, 3, 6, 7, 10, and 11
 - » PI: checks bits 8, 9, 10, 11, and 12
 - » P4: checks bits 4, 5, 6, 7, and 12
 - » P2: checks bits 1, 2, 3, 4, ..., 10, and 11
 - » P1: checks bits 1, 1, 0, 0, 1, 0, 1, 0
 - » Example
 - Uses even parity

Error Detection and Correction (cont'd)

| P1 | P2 | D7 | P4 | D6 | D5 | D4 | P8 | D3 | D2 | D1 | D0 |
|----|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |

- Correcting codewords to correct single-bit errors
 - Count bit positions from left to right starting from 1
 - Parity bits are called **check bits**
 - Example for 8-bit data ($d = 8$)
 - Parity bits are in positions that are a power of 2
 - Codeword is 12 bits long
 - We need 4 check bits
 - Example for 8-bit data ($d = 8$)
 - Parity bits are in positions that are a power of 2
 - Codeword is 12 bits long

Error Detection and Correction (cont'd)



Error Detection and Correction (cont'd)

- CRC uses a polynomial of degree n
- Example:
 - USB token polynomial: 110000000000000101
 - USB data polynomial: 100101
 - Polynomial identifies the 1 bit positions $X^5 + X^2 + 1$
 - USB polynomial for token packets $X^{16} + X^{15} + X^2 + 1$
 - USB polynomial for data packets $X^5 + X^2 + 1$
- » Such polynomials are called **polynomial generators**
- USB token polynomial: 100101
- USB data polynomial: 110000000000000101

Error Detection and Correction (cont'd)

- SECDED
- SECDED detection and correction scheme gives single-error detection and correction capability
- Previous example have 8 data bits and 5 check bits for SECDED capability
- » Example
 - This bit is added as the leftmost bit P_0 that is not used by the error correction scheme
 - Add an additional parity bit
 - To incorporate double error detection
- » Often used in high-performance systems
- Single-error correction and double error detection

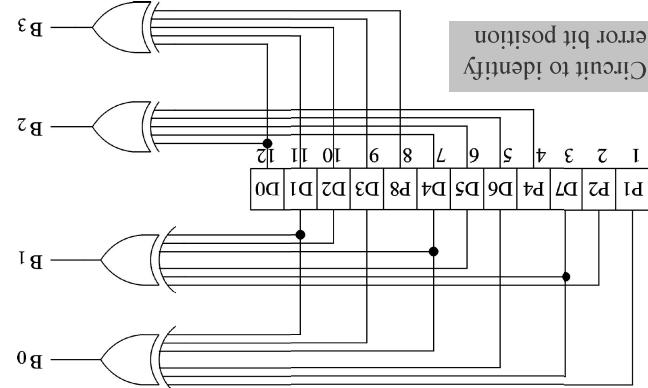
Error Detection and Correction (cont'd)

- A bit of theory behind CRC
- $D = d$ -bit data
- $C = (d+n)$ -bit codeword
- Goal: To regenerate C such that $C = d$ -bit remainder (*i.e.*, CRC code)
- $G = n$ -bit remainder polynomial (generator)
- $R = n$ -bit remainder (*i.e.*, CRC code)
- $D \times 2^n = Q \oplus R$
- $R = D \oplus QG$
- Since we appended n bits to the right
- » Remainder of $(C/G) = 0$
- G = degree n polynomial generator
- $D \times 2^n$ = remainder after division
- \oplus = XOR operation

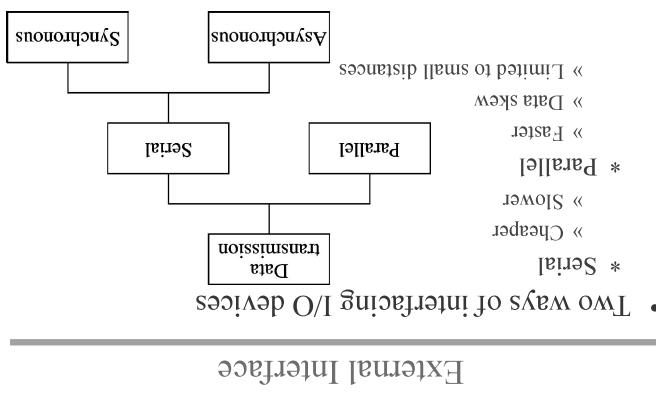
Error Detection and Correction (cont'd)

- Basic idea: If
 - Mostly 16 or 32 bits depending on the block size
 - Uses fixed number of bits
 - Widely used to detect burst errors
 - Complicated for a block of data
 - Cyclical redundancy check
- » Based on integer division
- $D \oplus R = Q$
- $D = Q + R$
- G

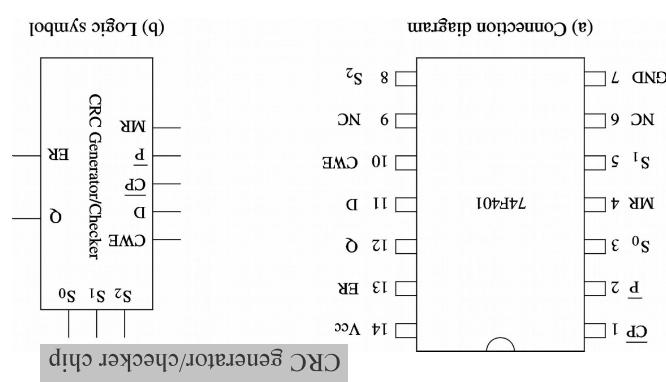
Error Detection and Correction (cont'd)



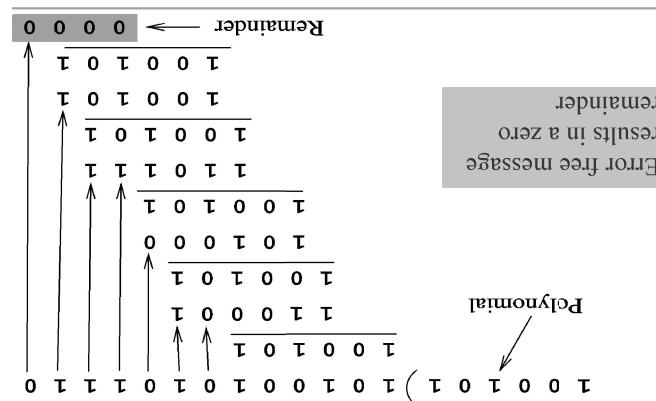
Error Detection and Correction (cont'd)



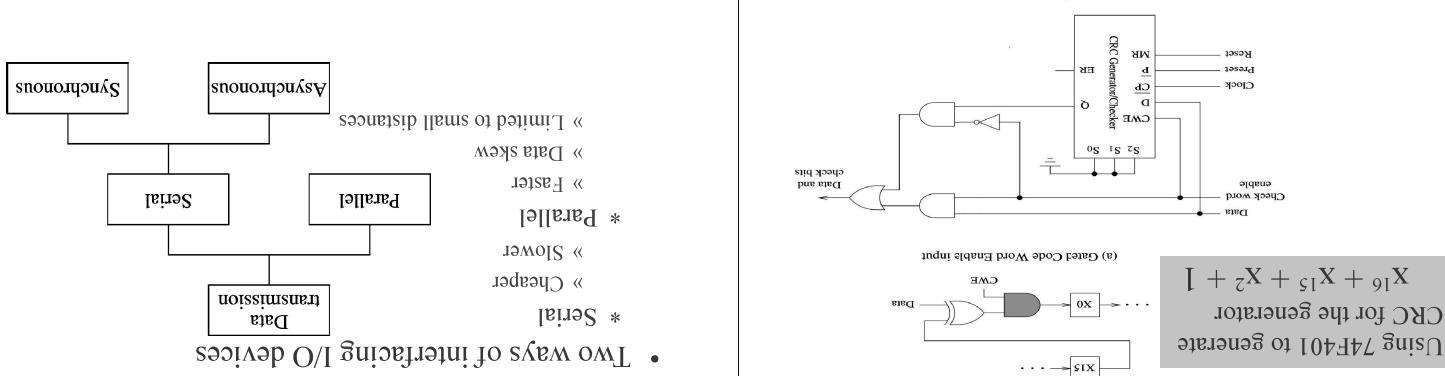
External Interface



Error Detection and Correction (cont'd)



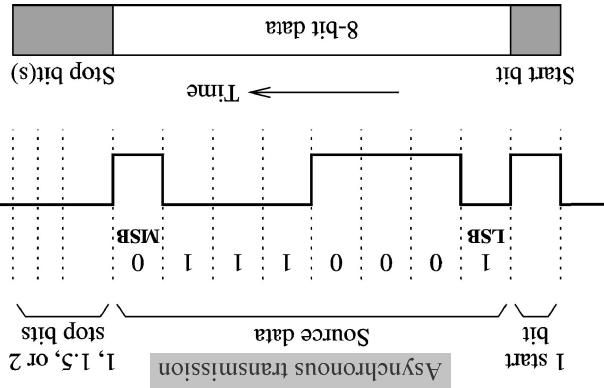
Error Detection and Correction (cont'd)



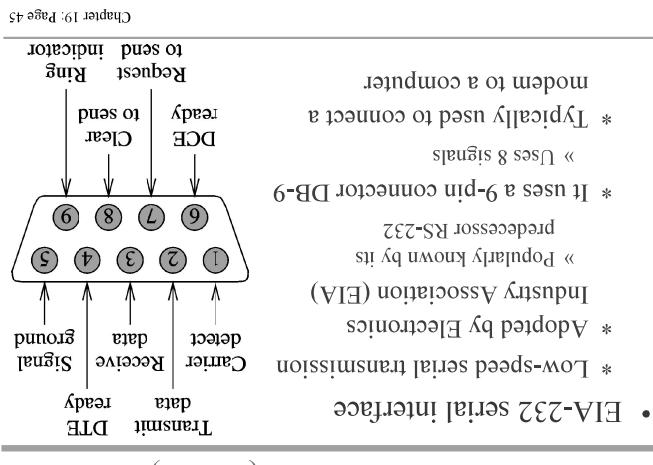
Error Detection and Correction (cont'd)

- Serial transmission
 - Start and stop bits
 - Each byte is encoded for transmission
 - No need for sender and receiver synchronization
 - Sender and receiver must synchronize
 - Done in hardware using phase locked loops (PLLs)
 - Block of data can be sent
 - » Major efficient
 - » Less overhead than asynchronous transmission
 - » Expensive
- Asynchronous
 - 1 start bit
 - Source data 1, 1.5, or 2 bits
 - Stop bits 1, 1.5, or 2 bits
 - Asynchronous transmission
- Transmitter protocol uses three phases
 - Connection setup
 - Data transmission
 - Connection termination
 - » Done by handshaking using requests-to-send (RTS) and clear-to-send (CTS) signals
 - » Done by deactivating RTS

External Interface (cont'd)



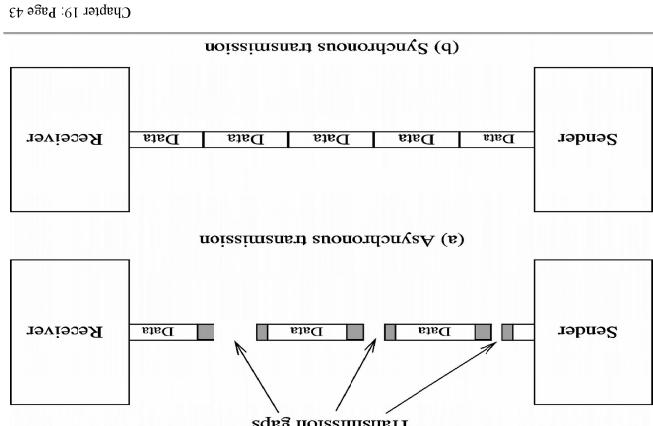
External Interface (cont'd)



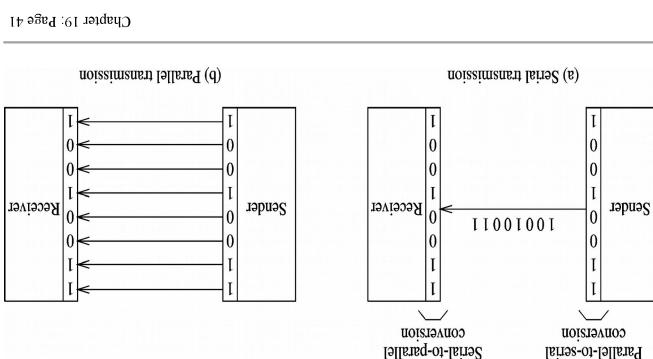
External Interface (cont'd)

- Serial transmission
 - Start and stop bits
 - Each byte is encoded for transmission
 - No need for sender and receiver synchronization
 - Sender and receiver must synchronize
 - Done in hardware using phase locked loops (PLLs)
 - Block of data can be sent
 - » Major efficient
 - » Less overhead than asynchronous transmission
 - » Expensive
- Synchronous
 - 1 start bit
 - Source data 1, 1.5, or 2 bits
 - Stop bits 1, 1.5, or 2 bits
 - Asynchronous transmission

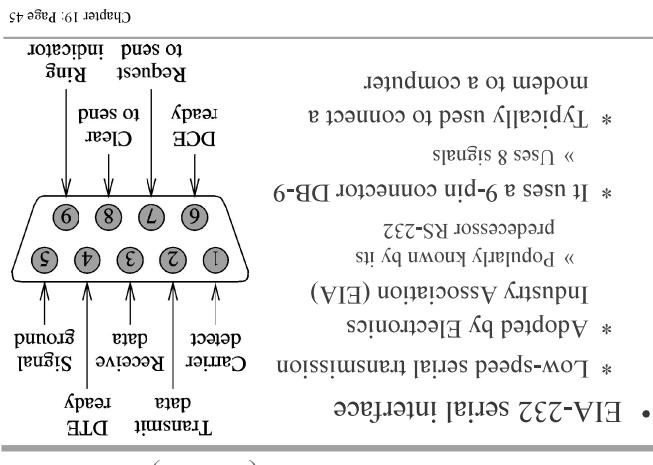
External Interface (cont'd)



External Interface (cont'd)



External Interface (cont'd)



External Interface (cont'd)

External Interface (cont'd)

External Interface (cont'd)

- Parallel printer interface

- » Data transfer uses simple handshaking
 - Latched by strobe (pin 1)
 - » ACK each byte, computer waits for ACK
 - Uses acknowledge (CK) signal

Chapter 19: Page 47

External Interface (cont'd)

• SCSI

- * Pronounced "scuzzy"
 - * Small Computer System Interface
 - » Comes in two bus widths
 - » Supports both internal and external connection
 - * SCSI
 - Known as *narrow SCSI*
 - Uses a 50-pin connector
 - Device id can range from 0 to 7
 - 8 bits
 - * Wide SCSI
 - Known as *wide SCSI*
 - Uses a 68-pin connector
 - Device id can range from 0 to 15
 - 16 bits

Table 19.4 Types of SCSI

| SCSI | Transfer rate
MB/s | Bus width
(bits) | SCSI type |
|--------------------------|-----------------------|---------------------|--------------------------|
| Fast SCSI | 10 | 8 | Ultra 1 SCSI |
| Ultra SCSI | 20 | 8 | Ultra 2 SCSI |
| Ultra 2 SCSI | 40 | 8 | Ultra 3 SCSI |
| Wide Ultra SCSI | 40 | 16 | Ultra 4 SCSI |
| Wide Ultra 2 SCSI | 80 | 16 | Ultra 4 (Ultra 320) SCSI |
| Ultra 3 (Ultra 160) SCSI | 160 | 16 | Ultra 4 (Ultra 320) SCSI |
| Ultra 4 (Ultra 320) SCSI | 320 | 16 | Ultra 4 (Ultra 320) SCSI |

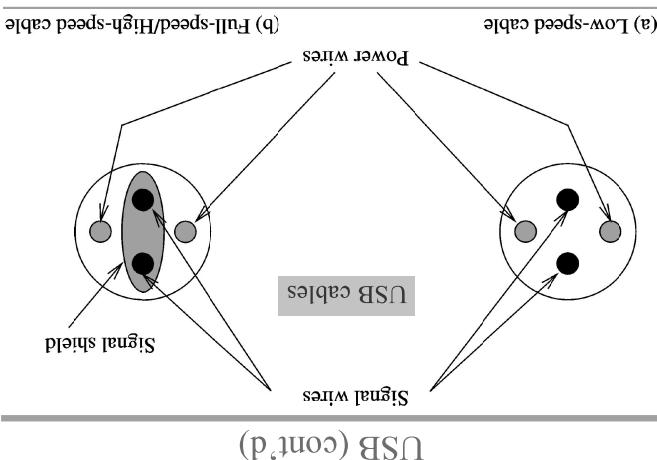
External Interface (cont'd)

| Description | SCSI signal | Pin | Pin | Signame | Description |
|---------------------|-------------|-----|-----|---------|-------------------------|
| Twisted pair ground | GND | 1 | 26 | D0 | Data 0 |
| Twisted pair ground | GND | 2 | 27 | D1 | Data 1 |
| Twisted pair ground | GND | 3 | 28 | D2 | Data 2 |
| Twisted pair ground | GND | 4 | 29 | D3 | Data 3 |
| Twisted pair ground | GND | 5 | 30 | D4 | Data 4 |
| Twisted pair ground | GND | 6 | 31 | D5 | Data 5 |
| Twisted pair ground | GND | 7 | 32 | D6 | Data 6 |
| Twisted pair ground | GND | 8 | 33 | D7 | Data 7 |
| Twisted pair ground | GND | 9 | 34 | DP | Data parity bit |
| Ground | GND | 10 | 35 | GND | Ground |
| Ground | GND | 11 | 36 | GND | Ground |
| Twisted pair ground | GND | 12 | 37 | GND | Reserived |
| No connection | | 13 | 38 | Termpwr | Termination power (+5V) |

12

External Interface (cont'd)

External Interface (cont'd)



USB (cont'd)

- * Allow hot attachment of devices
 - » Don't have to open the box to install and configure jumpers
- * Avoid installation and configuration problems
 - » USB does not require memory or address space
- * Avoid I/O address space and IRQ problems
 - » Standard interfaces support only one device
- * Avoid non-shareable interfaces
 - PS/2, serial, parallel, monitor, microphone, keyboard, ...
 - » Eliminates multitude of interfaces
- * Avoid device-specific interfaces
 - » Eliminates multitude of interfaces

Motivation for USB

- Minimizes adverse impact of cable propagation delay
 - Improves throughput
- ACK
 - » A number of bytes (e.g., 8) can be sent without waiting for REQ-ACK signals to be used for each byte
 - » Data are transferred synchronously
- * On a synchronous SCSI
- * SCSI uses asynchronous mode for all bus negotiations
 - » Uses handshaking using REQ and ACK signals for each byte of data

External Interface (cont'd)

- * SCSI uses client-server model
 - » Initiator and target for client and server
 - » Initiator issues commands to targets to perform a task
 - » Targets are typically SCSI host adapters
 - » Targets receive the command and perform the task
 - » Targets are SCSI devices like disk drives
- * SCSI transfer proceeds in phases
 - » Command
 - » Message in
 - » IN and OUT from initiator point
 - » Data in
 - » Message out
 - » Data out
 - » Status
- * Universal Serial Bus
 - » Originally developed in 1995 by a consortium including Compaq, HP, Intel, Lucent, Microsoft, and Philips
 - » USB 1.1 supports
 - » Low-speed devices (1.5 Mbps)
 - » Full-speed devices (12 Mbps)
 - » High-speed devices (480 Mbps over USB 2.0)
 - » USB 2.0 supports
 - » Standard interface (1.5 Mbps)
 - » Low-speed devices (12 Mbps)
 - » Full-speed devices (480 Mbps)
- * Additional advantages of USB
 - » Simple devices can be bus-powered
 - » Examples: mouse, keyboards, floppy disk drives, wireless LANs, ...
 - » Possibility because USB allows data to flow in both directions
 - » Expandable through hubs
 - » Power conservation
 - » Power suspended state if there is no activity for 3 ms
 - » Error detection and recovery
 - » Control peripherals
 - » Power distribution

USB (cont'd)

- * SCSI uses client-server model
 - » Initiator and target for client and server
 - » Initiator issues commands to targets to perform a task
 - » Targets are typically SCSI host adapters
 - » Targets receive the command and perform the task
 - » Targets are SCSI devices like disk drives
- * SCSI transfer proceeds in phases
 - » Command
 - » Message in
 - » IN and OUT from initiator point
 - » Data in
 - » Message out
 - » Data out
 - » Status
- * USB
 - » Originally developed in 1995 by a consortium including Compaq, HP, Intel, Lucent, Microsoft, and Philips
 - » USB 1.1 supports
 - » Low-speed devices (1.5 Mbps)
 - » Full-speed devices (12 Mbps)
 - » High-speed devices (480 Mbps over USB 2.0)
 - » USB 2.0 supports
 - » Standard interface (1.5 Mbps)
 - » Low-speed devices (12 Mbps)
 - » Full-speed devices (480 Mbps)
- * Additional advantages of USB
 - » Simple devices can be bus-powered
 - » Examples: mouse, keyboards, floppy disk drives, wireless LANs, ...
 - » Possibility because USB allows data to flow in both directions
 - » Expandable through hubs
 - » Power conservation
 - » Power suspended state if there is no activity for 3 ms
 - » Error detection and recovery
 - » Control peripherals
 - » Power distribution

External Interface (cont'd)

- * Bulk transfer
 - Recovery is by means of retries
 - » Error detection and recovery are used
 - Bulk transfers are deferred until load decreases bandwidth
 - » If the other three types of transfers take 100% of the bandwidth priority bandwidth allocation
 - Example: sending data to a printer requires no specific data transfer rate
 - » For devices with no specific transfer rate requirements

USB (cont'd)

- * Isochronous transfer
 - Polling interval can range from 1 ms to 25 ms
 - » Uses polling
 - * Interrupt transfer
 - » Four types of transfer
 - * Synchronous transfer
 - » Used in real-time applications that require constant data transfer rate
 - Example: Reading audio from CD-ROM
 - » These transfers are scheduled regularly
 - » Do not use error detection and recovery

USB (cont'd)

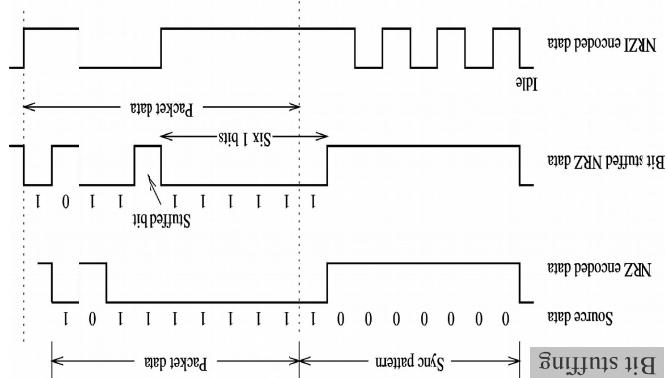
- * Still a problem
 - » Long strings of 1s do not cause signal change
 - » To solve this problem
 - » A zero is inserted after every six consecutive 1s
 - » Uses **bit stuffing**
- * Two desirable properties
 - » Non-Return-to-Zero-Inverted
 - » It is called **differential encoding**
 - » Signal transitions, not levels, need to be detected
 - » Long string of zeros causes signal changes
 - » Signal transitions, not levels, need to be detected

NRZI encoding

USB (cont'd)

- * Control transfer
 - » Used to configure and set up USB devices
 - » Three phases
 - Set up stage
 - Configuration stage
 - Data stage
 - » Control transfers that require data use this stage
 - Status stage
 - Options stage
 - Step phases
 - Used to request made to target device
 - » Connives type of request made to target device
 - Allocates a guaranteed bandwidth of 10%
 - » Checks the status of the operation
 - Error detection and recovery are used
 - Recovery is by means of retries

USB (cont'd)

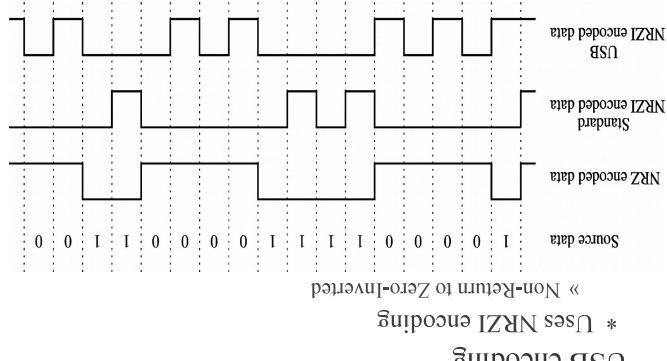


USB (cont'd)

- * Uses NRZI encoding
 - » Non-Return-to-Zero-Inverted
 - » It is called **differential encoding**
 - » A signal transition occurs if the next bit is zero
 - » Two desirable properties
 - » Signal transitions, not levels, need to be detected
 - » Long string of zeros causes signal changes
 - » Signal transitions, not levels, need to be detected

USB encoding

USB (cont'd)



USB (cont'd)

- * Most 4-port hubs are dual-powered
 - » Support 4 bus-powered USB hubs
 - » Support 4 high-powered devices
- * Self-powered
 - Support only low-powered devices
 - Number of ports is limited to four
- » Downstream ports can only supply 500 mA
 - » Must be connected to an upstream port that can supply 500 mA
 - » No extra power supply required
- * Bus-powered

USB (cont'd)

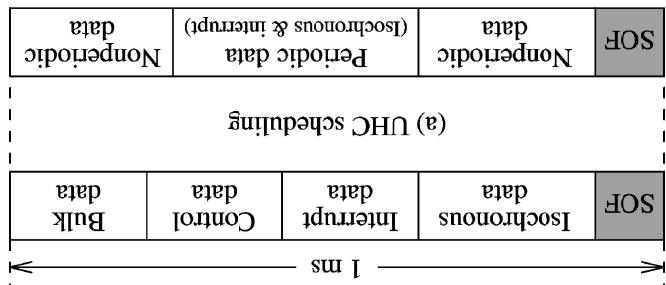
- * High-powered
 - » Between 100 mA and 500 mA
 - » Can be bus-powered
 - » Less than 100 mA
- * Low-power
 - » Between 100 mA and 500 mA
 - » Can be designed to have their own power
 - » Operate in three modes
 - » Unpowered (100 mA)
 - » Configured (500 mA)
 - » Suspended (about 2.5 mA)

USB (cont'd)

- * Left over bandwidth is allocated to non-periodic transfers
- * Next periodic transfers are scheduled
 - » Guarantees 90% bandwidth
- * Non-periodic transfers: control and bulk
 - » 10% bandwidth reserved
 - » Non-periodic transfers: control and bulk
- * Reserves space for non-periodic transfers first
- * Different from UHC scheduling
- * UHC scheduling

USB (cont'd)

(b) OHC scheduling



USB (cont'd)

- * Bulk transfers are scheduled only if there is bandwidth available
- * Bulk transfers are scheduled only if there is bandwidth
 - » Control transfers are guaranteed 10% of bandwidth
- * These transfers are followed by control and bulk transfers
 - » Can take up to 90% of bandwidth
 - » Periodic transfers: isochronous and interrupts
- * Schedules periodic transfers first
- * Root hub
- * Initiates transactions over USB
- * USB host controller
 - » Provides connection points
 - » Differentiates between the two
 - » Specified by National Semiconductor, Microsoft, Compaq
 - » Universal host controller (UHC)
 - » Open host controller (OHC)
 - » Defined by Intel
- * Two types of host controllers
 - » Provides connection points
 - » Initiates transactions over USB

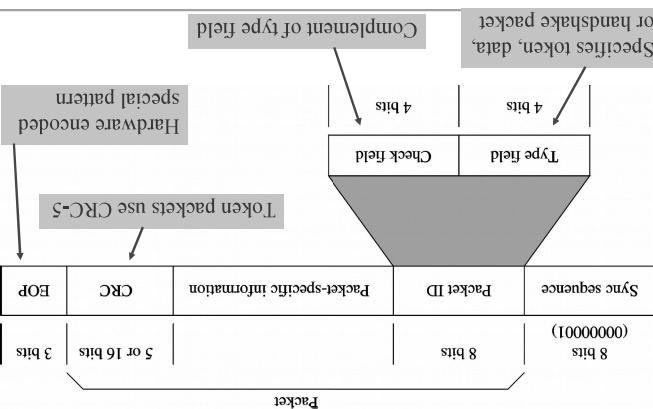
USB (cont'd)

- * How they schedule the four types of transfers
- * Difference between the two
- * Specified by National Semiconductor, Microsoft, Compaq
- * Universal host controller (UHC)
- * Open host controller (OHC)
- * Defined by Intel
- * Provides connection points
- * Initiates transactions over USB
- * USB host controller
 - » Provides connection points
 - » Initiates transactions over USB
- * Two types of host controllers
 - » Provides connection points
 - » Initiates transactions over USB

USB (cont'd)

- USB 2.0
 - * Supports 40X data rates
 - » Up to 480 Mbps
 - » 1/8 of USB 1.1
 - * USB 2.0 uses 125 μs frames
 - * USB 1.1 uses 1 ms frames
 - * Competitive with
 - » SCSI
 - » IEEE 1394 (FireWire)
 - * Widely available now

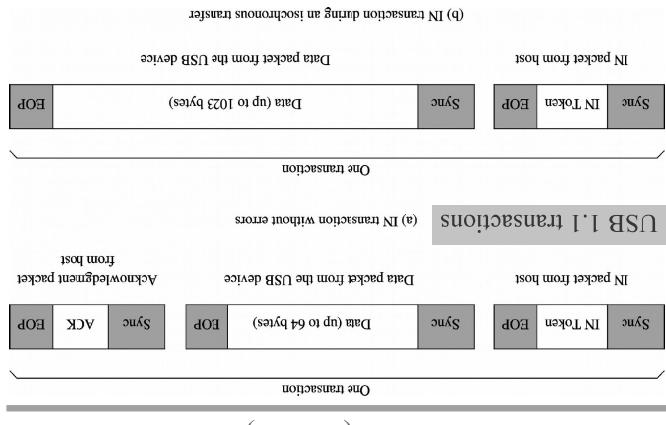
USB (cont'd)



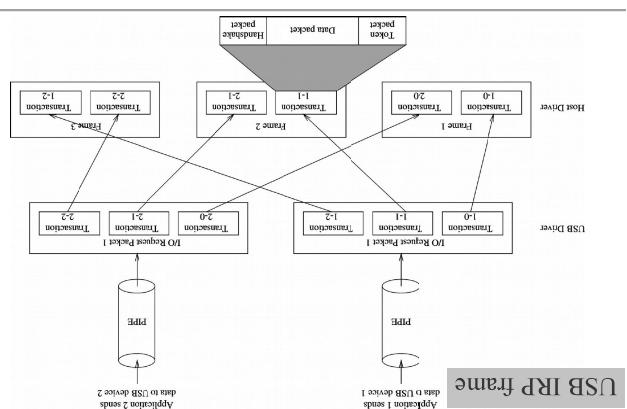
USB (cont'd)

- USB transactions
 - * Transfers are done in one or more transactions
 - » Each transaction consists of several packets
 - Data packet phase (optional)
 - Specifies transaction type and target device address
 - Token packet phase
 - » Tooken packet phase
 - » Handshake packet phase
 - Maximum of 1023 bytes are transferred
 - Provides feedback on whether data has been received without error
 - For guaranteed delivery
 - Except for isochronous transfers, others use error detection
 - » Handshake packet phase (optional)
 - Provides feedback on whether data has been received without error
 - * Transactions may have between 1 and 3 phases

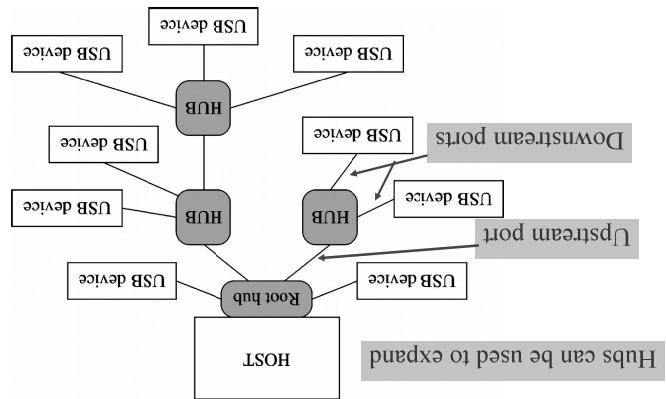
USB (cont'd)



USB (cont'd)



USB (cont'd)



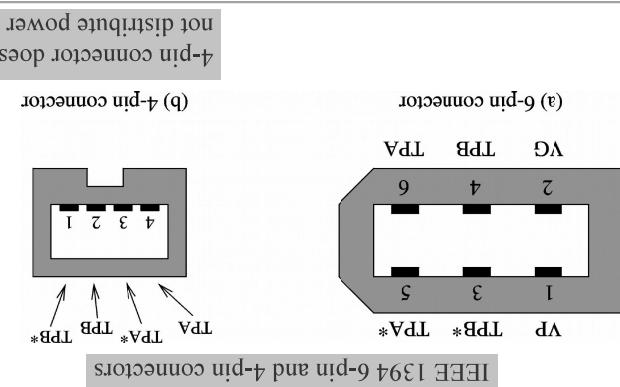
USB (cont'd)

- * Isochronous
 - » Cycle time: 125 ns
 - » Bandwidth allocation on a cycle-by-cycle basis
 - » Up to 80% of bandwidth allocated
 - » No acknowledgegment
 - » For real-time applications
- * Asynchronous
 - » Guaranteed bandwidth of 20%
 - » Uses an acknowledgegment to confirm delivery
 - Example: writing a file to a disk drive
 - » Applications that require correct delivery of data
- * Transfer types
 - » Isochronous
 - » Asynchronous

IEEE 1394 (cont'd)

- * Encoding
 - » NRZ encoded data
 - » Strobe signal is encoded
 - » Uses a simple NRZ encoding
 - » Changes the signal even if successive bits are the same

IEEE 1394 (cont'd)



IEEE 1394 (cont'd)

- * Expandable bus
 - » Devices can be connected in daisy-chain fashion
 - » Hubs can used to expand
 - » Power distribution
 - Like the USB, cables distribute power
 - Much higher power than USB
 - » Current can be up to 1.5 Amps
 - » Voltage between 8 and 33 V
 - » Error detection and recovery
 - As in USB, uses CRC
 - » Uses retransmission in case of error
 - » Long cables
 - Like the USB

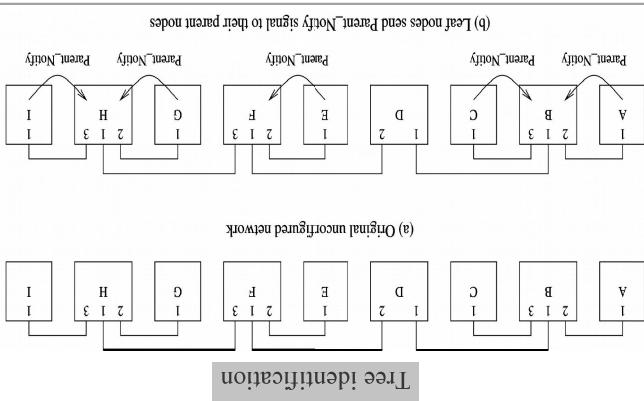
IEEE 1394 (cont'd)

- * Advantages
 - » High speed
 - » Supports three speeds
 - 100, 200, 400 Mbps
 - » Competes with USB 2.0
 - Plans to boost it to 3.2 Gbps
 - » Like USB
 - » No need to shut down power to attach devices
 - » Peer-to-peer support
 - » USB is processor-centric
 - » Supports peer-to-peer communication without involving the processor
 - » Hot attachment

IEEE 1394 (cont'd)

- * Speed peripherals
 - » Apple originally developed this standard for high-
 - » Known by a variety of names
 - FireWire
 - » Sony: iLINK
 - » First released in 1995 as IEEE 1394-1995
 - » A slightly revised version as 1394a
 - » Next revision 1394b
 - » Shares many of the features of USB

IEEE 1394



IEE 1394 (cont'd)

- * Does not require the host system
- * Consists of two main phases
 - » Tree identification
 - » Parent_notify and Child_Notify
- Uses two special signals
- Used to find the network topology
- Tree identification
- Does after the tree identification
- Assigns unique ids to nodes

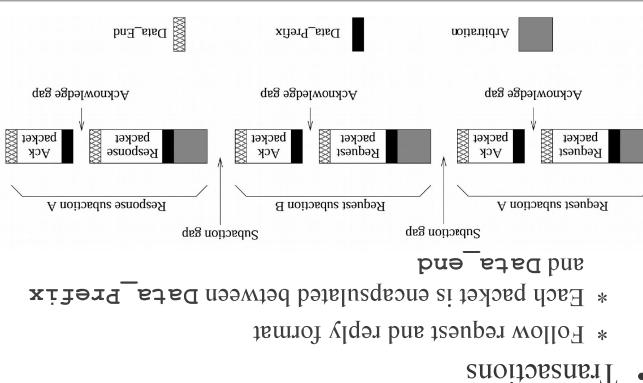
IEE 1394 (cont'd)

- * Bus arbitration
 - » Uses fairness interval
 - » Fairness-based allocation for asynchronous channels
 - » Bandwidth allocation to isochronous channels
 - » All nodes with pending asynchronous transaction are allowed bus ownership once
 - » Nodes with pending isochronous transactions go through arbitration during each cycle
 - » IRM is used for isochronous bandwidth allocations
- * Arbitration must reselect
 - » Needed because of peer-to-peer communication
 - » Arbitration must respecify
 - » Bandwidth allocation must respecify
 - » Fairness-based allocation for asynchronous channels
 - » All nodes with pending isochronous bandwidth allocations through each cycle
 - » Nodes with pending isochronous bandwidth allocations go through arbitration during each cycle

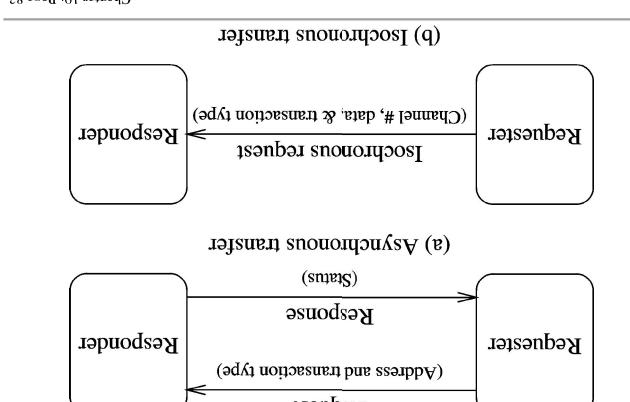
IEE 1394 (cont'd)

- * Main difference:
 - » No acknowledgement packets
 - » Similar to asynchronous transactions
- * Synchronous transactions
 - » Acknowledgment packets are transmitted on a separate channel called the 'Arbitration' channel.
 - » Data frames are transmitted on the 'Data' channel.
 - » Isochronous gaps are used to separate data frames.
 - » Data frames consist of a Data_Prefix, Data_End, and Arbitration.

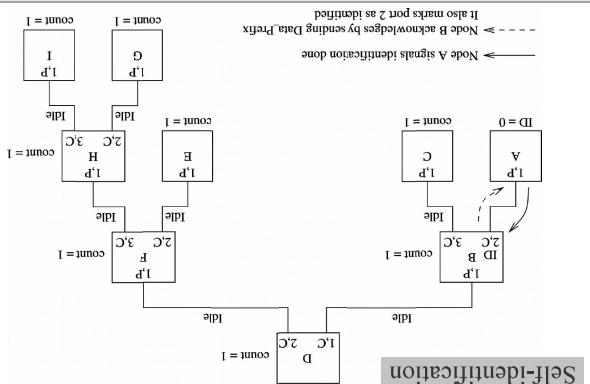
IEE 1394 (cont'd)



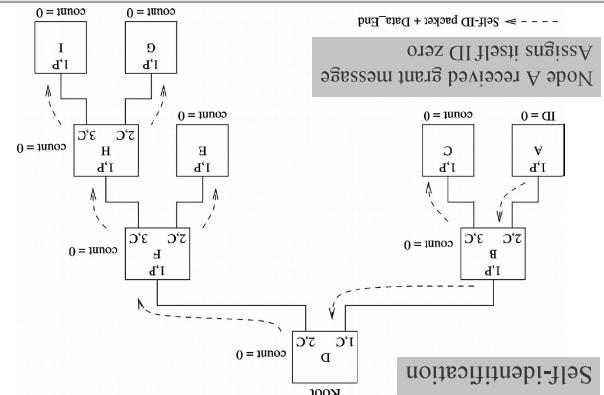
IEE 1394 (cont'd)



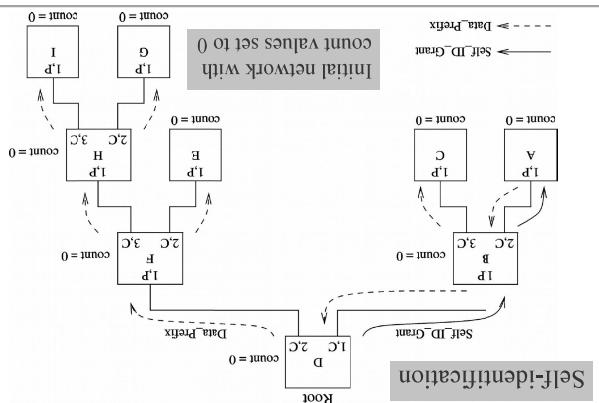
IEE 1394 (cont'd)



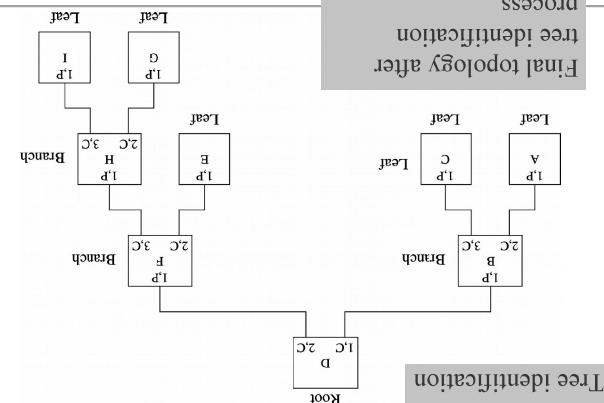
IEEE 1394 (cont'd)



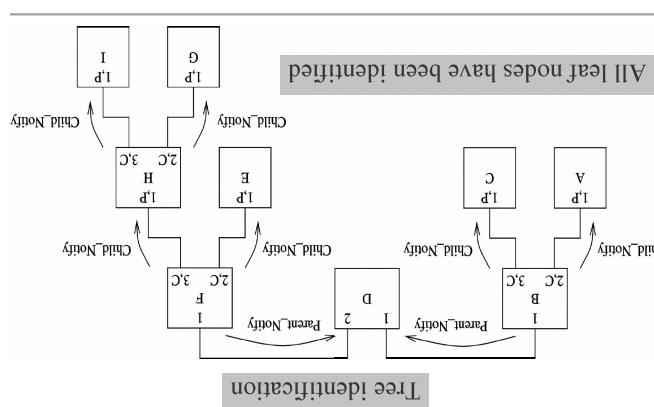
IEEE 1394 (cont'd)



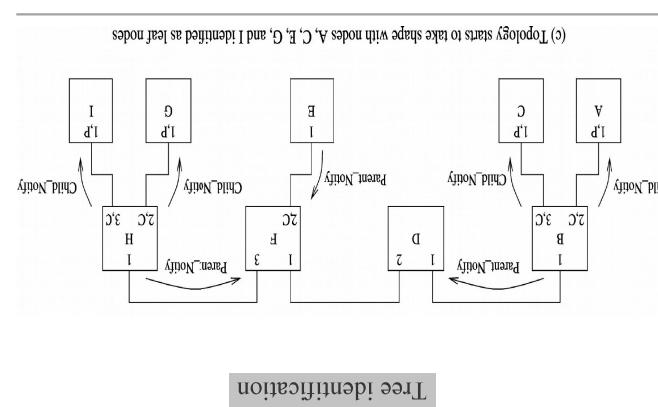
IEEE 1394 (cont'd)



IEEE 1394 (cont'd)

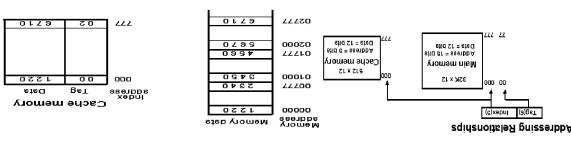


IEEE 1394 (cont'd)



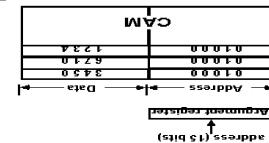
IEEE 1394 (cont'd)

On memory request, index field is used to access the cache. Tag field of CPU address is compared with the tag in the word read in the cache



- Each word in cache consists of data word and its associated tag
- Index bits are used to access the cache
- N-bit CPU memory address - k bits index field and (n-k) bits tag field

Direct Mapping



- CAM - content addressable memory
- If found data is read and sent to the CPU else main memory is accessed.
- Stores both the address and the content of the memory word the associative memory is searched for a matching address
- Any block location in cache can store any block in memory
- Fastest, most flexible but very expensive

Associative Mapping

- Bandwidth - time to retrieve the rest of this block
- Latency - time to retrieve the first word of the block
- Time required for the cache miss depends on both the latency and bandwidth

Parameters of Cache Memory

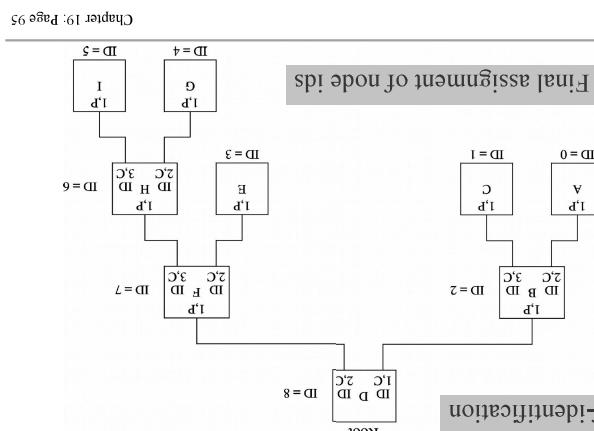
- Miss penalty
- Additional cycles required to serve the miss

- Hit ratio
- Cache Miss
- Cache Hit
- A referenced item is found in the cache by the processor
- A referenced item is not present in the cache
- Ratio of number of hits to total number of references => number of hits/(number of hits + number of Miss)

Parameters of Cache Memory

- SCSI is dominant in disk and storage device
- Parallel interface
 - * Its bandwidth could go up to 640 MB/s
- IEEE 1394
 - * Its bandwidth could go up to 640 MB/s
 - * Serial interface
 - * Supports peer-to-peer applications
 - * Dominant in video applications
- USB
 - * Useful in low-cost, host-to-peripheral applications
 - * Provides high-speed support

Buses Wars



IEEE 1394 (cont'd)

| | | | |
|---------------------|----|---|---|
| Main memory address | 14 | 8 | 4 |
| | | | |

be 22 bits, so the tag length is 11 bits and the word field length is 5 bits. Main memory consists of 64-Mbyte / 16 bytes = 2²² blocks. Therefore, the set plus tag lengths must contain 128 words. Therefore, 7 bits are needed to specify the word.

Similarly, the set number. Main memory consists of $4K = 2^2$ blocks. Therefore, the set plus tag lengths must be 12 bits and therefore the tag length is 8 bits. Each block contains 16 sets of 4 lines each. Therefore, 4 bits are needed to identify the set number. For the 64-Mbyte main memory, a 26-bit address is needed. Main memory consists of 256 sets of 2 lines each. Therefore 8 bits are needed to identify the set number. Thus the cache has 16 lines of 8 bytes each. There are a total of 8 kbytes / 16 bytes = 512 lines in the cache. Thus the cache

- A two-way set associative cache has lines of 16 bytes and a total size of 8k bytes. The 64-Mbyte main memory is byte addressable. Show the format of main memory addresses.

Problem 2

| | | | |
|---------------------|---|---|---|
| Main memory address | 8 | 4 | 7 |
| | | | |

The cache is divided into 16 sets of 4 lines each. Therefore, 4 bits are needed to identify the set number. Main memory consists of 4K blocks of 128 words each. Therefore, the set plus tag lengths must be 12 bits and therefore the tag length is 8 bits. Each block contains 16 words. Therefore, 7 bits are needed to specify the word.

- A set associative cache consists of 64 lines or slots, divided into four line sets. Main memory consists of 4K blocks of 128 words each. Show the format of main memory addresses.

Problems

The comparison logic is done by an associative search of the tags in the set similar to an associative memory search, thus the name "Set-Associative".

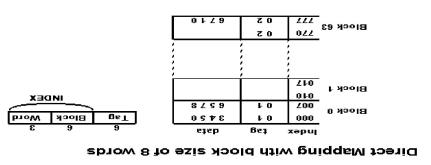
| Set Associative Mapping Cache with set size of two | | | | | | | | |
|--|-----|---------|-----|---------|-----|---------|-----|---------|
| Index | Ta0 | Data | Ta1 | Data | Ta0 | Data | Ta1 | |
| 000 | 0 1 | 3 4 5 0 | 0 2 | 3 6 7 0 | 0 0 | 2 3 4 0 | 0 2 | 6 7 1 0 |

- Each word under the same index address can store 2 or more words of memory under the same address.

Set-Associative Mapping

- Hit ratio increases as the set size increases but more complex comparisons logic is required when number of bits in words of cache increases
- When a miss occurs and set is full, one of tag-data items are replaced using block replacement policy

- It is time consuming but improves hit ratio because of main memory to cache.
- When a miss occurs, entire block is transferred from the sequential nature of programs.



Tags within the block are same.

EX: 512 words cache - 64 blocks of 8 words each - block field (6bits) and words field (3bits)

When memory is divided into blocks of words, index field - block field and word field

repeatedly.

Disadvantage: hit ratio drops if 2 or more words with same index but different tags are accessed

tag replacing the previous value

it is then stored in the cache together with the new

else miss - word is read from the memory

If match then there is hit

Block Replacement

- Least Recently Used: (LRU)
Replace that block in the set that has been in the cache longest with no reference to it.
- First Come First Out: (FIFO)
Replace that block in the set that has been in the cache longest.
- Least Frequently Used: (LFU)
Replace that block in the set that has experienced the fewest references.

Update Policies - Write Through

- Update main memory with every memory write operation
- Cache memory is updated in parallel if it contains the word at specific address.
- Advantage: main memory always contains the same data as the cache
- It is important during DMA transfers to ensure the data in main memory is valid
- Disadvantage: slow due to memory access time

Write Back

- Only cache is updated during write operation and marked by flag.
- When the word is removed from the cache, it is copied into main memory
- Memory is not up-to-date, i.e., the same item in cache and memory may have different value

Update Policies

- Write-Around
 - corresponds to items not currently in the cache (i.e. write without affecting the cache).
 - update the item in main memory and bring the block containing the updated item into the cache.
- Write-Allocate
 - misses the item could be updated in main memory only
 - update the item in main memory and bring the block containing the updated item into the cache.

Performance analysis

- Look through: The cache is checked first for a hit, and if a miss occurs then the access to main memory is started.
- Look aside: access to main memory in parallel with the cache lookup;

$$\begin{aligned} & \bullet \text{Miss Ratio } h = \frac{\text{number of references found in the cache}}{\text{total number of memory references}} \\ \\ & \bullet \text{TA} = h \cdot TC + (1-h) \cdot TM \\ & \quad \text{TC is the average cache access time} \\ & \quad \text{TM is the average memory access time} \\ & \quad (\text{Mean memory access time}) \\ & \bullet \text{Look through} \\ & \quad TA = TC + (1-h) \cdot TM \end{aligned}$$

• Miss Ratio $m = (1-h)$

Cache Organization

- Split cache for instruction and data
- Cache (data) – mostly random access
- Cache (instruction) – mostly accessed sequentially
- Separate caches for instructions and data
- Unified cache
- D-cache (data) – mostly random access
- L-cache (instruction) – mostly accessed sequentially
- Higher hit rate for unified cache as it balances between instruction and data
- Split caches eliminate contention for cache between the instruction processor and the execution unit – used for pipelineing processes

Multi-level Caches

- The primary cache is referred to as the L1 (level 1) cache and the secondary cache is called the L2 (level 2) cache.
- One way in which this penalty can be reduced is to provide another cache, the secondary cache, which is accessed in response to a miss in the primary cache.
- Most high-performance microprocessors include an L2 cache that is often located off-chip, whereas the L1 cache is located on the same chip as the CPU.
- With a two-level cache, central memory has to be accessed only if a miss occurs in both caches.
- The primary cache is detailed time that it takes to obtain the requested item from central memory.
- The penalty for a cache miss is the extra time it takes to detail the request from central memory.

Sources of Cache Misses

- Compulsory Misses: These are misses that are caused by the cache being empty initially.
- Cold Misses : The very first access to a block will result in a miss because the block is not brought into cache until it is referenced.
- Capacity Misses : If the cache cannot contain all the blocks needed during the execution of a program, capacity misses will occur due to blocks being discarded and later retrieved.
- Conflict Misses: If the cache mapping is such that multiple blocks are mapped to the same cache entry.

Sources of Cache Misses

- Consider a memory system with $T_c = 100\text{ns}$ and $T_m = 1200\text{ns}$. If the effective access time is 10% greater than the cache access time, what is the hit ratio H in look-through cache?
- $$T_{avg} = T_c + (1-h)*T_m$$
- $$T_{avg} = T_c + (1-h)*T_m$$
- $$T_{avg} = 1.1T_c = T_c + (1-h)*T_m$$
- $$0.1 = (1-h)$$
- $$h = 10/1200$$
- $$h = 1190/1200$$

Problem

Example assume that a computer system employs a cache with an access time of 20ns and a main memory with a cycle time of 200ns. Suppose that the hit ratio for reads is 50%.

$$T_{avg} = T_c + (1-h)*T_m$$

$$T_{avg} = 20\text{ns} + 0.10*200\text{ns} = 40\text{ns}$$

(a) what would be the average access time for reads if the cache is a "look-through" cache?

$$T_{avg} = T_c + (1-h)*T_m$$

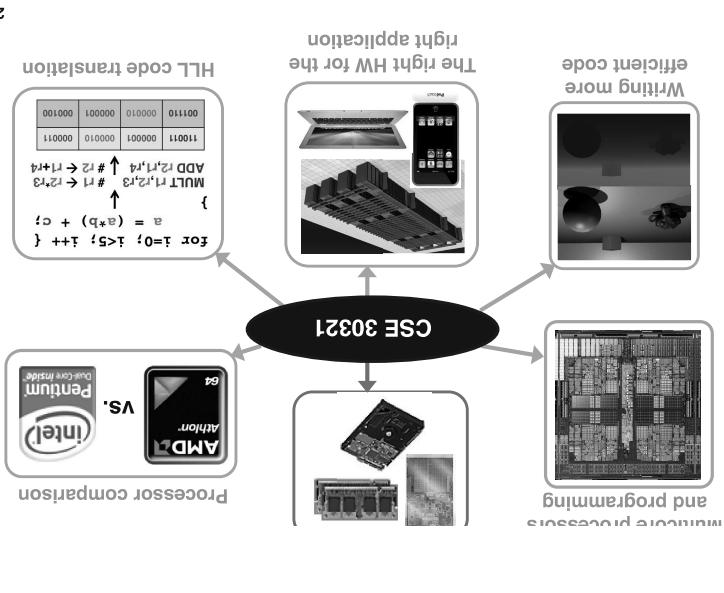
$$T_{avg} = 20\text{ns} + 0.10*200\text{ns} = 40\text{ns}$$

(b) what would be the average access time for reads if the cache is a "look-side" cache?

$$T_{avg} = h*T_c + (1-h)*T_m$$

$$T_{avg} = 0.5*20\text{ns} + 0.10*200\text{ns} = 38\text{ns}$$

Exercise: assume that a computer system employs a cache with an access time of 20ns and a main memory with a cycle time of 200ns. Suppose that the hit ratio for reads is 50%.



- J.L. Hennessy & D.A. Patterson, Computer architecture: A quantitative approach, Fourth Edition, Morgan Kaufmann, 2004.

References

- If you're a hardware designer OR a compiler writer, you need to be aware of how HLL is mapped to assembly instructions AND what HW is used to execute a sequence of assembly instructions.
- Otherwise, it is quite possible to always have built-in inefficiencies.

Why it's important...

- We must also account for the efficient pipelining of program results to ensure that the correct speedups offered from the introduction of pipelining.
- Additional HW is needed to ensure the correct control instructions (e.g. beq) to preserve performance gains and program correctness.

- Pipelining changes the timing as to when the result(s) of an instruction are produced.
- Additional HW is needed to ensure the correct sequencing of instructions while maintaining the correct pipeline.

- Additional HW is needed to ensure the correct sequencing of instructions while maintaining the correct pipeline.

Fundamental lesson(s)

- Example:**
- A computer system employs a write-back cache with a 70% hit ratio for writes. The cache operates in lookaside mode and has a 90% read hit ratio. Reads account for 20% of all memory references and writes account for 20%. If the main memory cycle time is 200ns and the cache access time is 20ns, what would be the average access time for all references (reads as well as writes)?
- The average access time for reads = $0.9 * 20\text{ns} + 0.1 * 200\text{ns} = 38\text{ns}$
- Hence the overall average access time for combined reads and writes is
- $$0.8 * 38\text{ns} + 0.2 * 74\text{ns} = 45.2\text{ns}$$

$$\text{The average write time} = 0.7 * 20\text{ns} + 0.3 * 200\text{ns} = 74\text{ns}$$

$$0.8 * 38\text{ns} + 0.2 * 74\text{ns} = 45.2\text{ns}$$

STRUCTURAL HAZARDS

- Pipeline hazards prevent next instruction from executing during designated clock cycle

Pipeline hazards

- There are 3 classes of hazards:
 - Structural Hazards:
 - HW cannot support all possible combinations of instructions
 - Arise from resource conflicts
 - Data Hazards:
 - Occur when given instruction depends on data from an instruction ahead of it in pipeline
 - Result from branch, other instructions that change flow of program (i.e. change PC)
 - Control Hazards:
 - Result from L1 cache miss

- If no stalls, speedup equal to # of pipeline stages in ideal case

$$\text{Speedup} = \frac{\text{CPI unpipelined}}{1 + \text{pipeline stall cycles per instruction}}$$

- Ignoring overhead and assuming stages are balanced:

Stalls and performance

- CPI pipelined =
 - Ideal CPI + Pipeline stall cycles per instruction
 - $1 + \text{Pipeline stall cycles per instruction}$
- Pipeline can be viewed to:
 - Decrease CPI or clock cycle time for instruction
 - Let's see what affect stalls have on CPI...

Stalls and performance

- If we say an instruction was "issued earlier than instruction x" and is further along in the pipeline

- If we say an instruction was "issued later than instruction x", we mean that it was issued after instruction x and is not as far along in the pipeline

- A note on terminology:

- Often, pipeline must be stalled
 - Stalling pipeline usually lets some instruction(s) in pipeline proceed, another/others wait for data, resource, etc.

How do we deal with hazards?



- ...and how they (generally) impact performance.

- Let's look at hazards...

On the board...

DATA HAZARDS

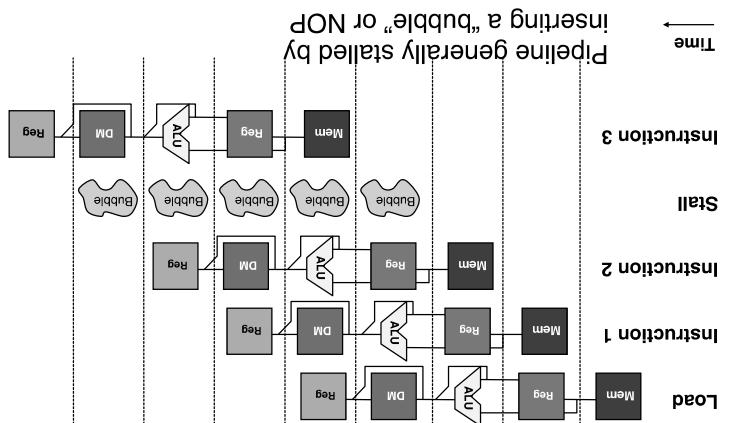
- CPI degrades quickly from our ideal “1” for even the simplest cases...
- (especially for the memory access example – i.e. the common case)
- Answer: Add more hardware.

What's the realistic solution?

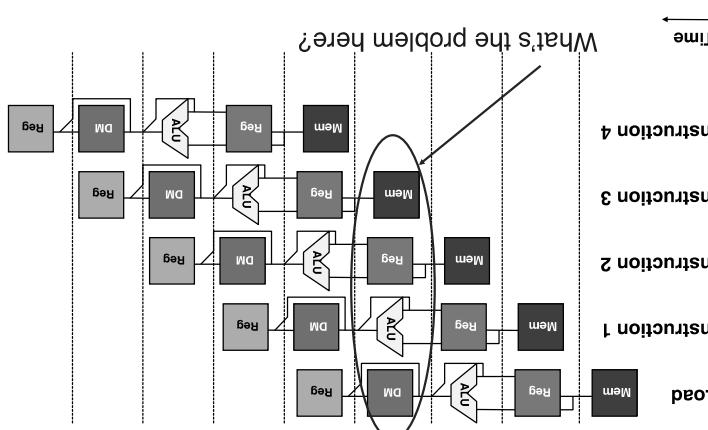
- Thus, no instruction completes on clock cycle 8
- LOAD instruction “steals” an instruction fetch cycle which will cause the pipeline to stall.



Or alternatively...



How is it resolved?

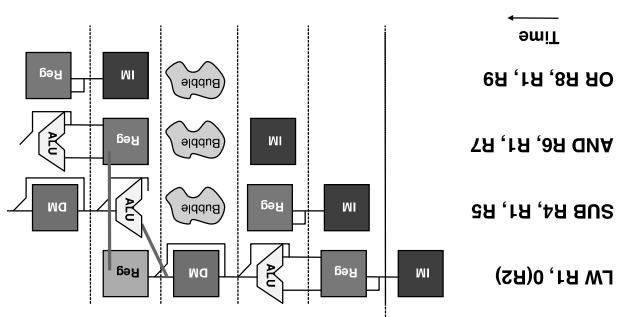


An example of a structural hazard

- Avoid structural hazards by duplicating resources
- e.g. an ALU to perform an arithmetic operation and an adder to increment PC
- If not all possible combinations of instructions can be executed, structural hazards can occur
- Pipelines stall result of hazards, CPI increased from the usual “1”

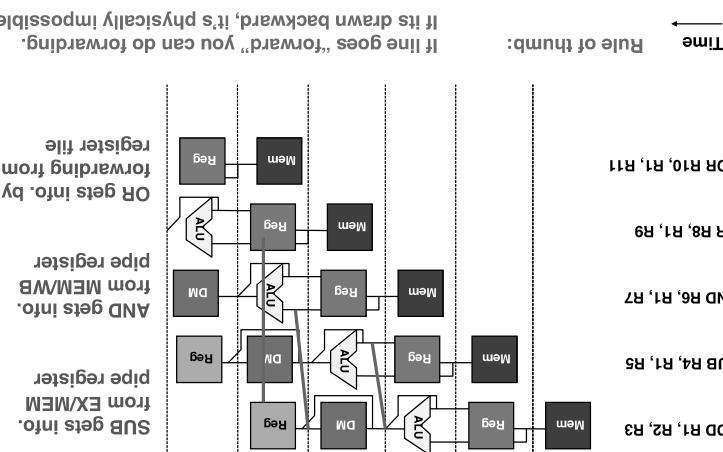
Structural hazards

insertion of bubble causes # of cycles to complete this sequence to grow by 1



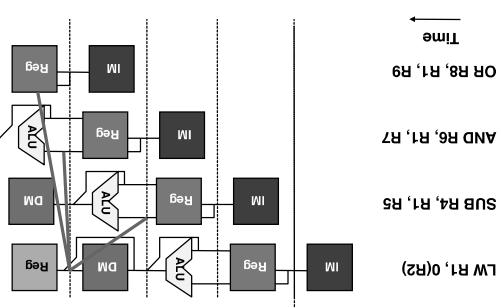
The solution pictorially

If line goes "forward", you can do forwarding.
If its drawn backward, it's physically impossible.



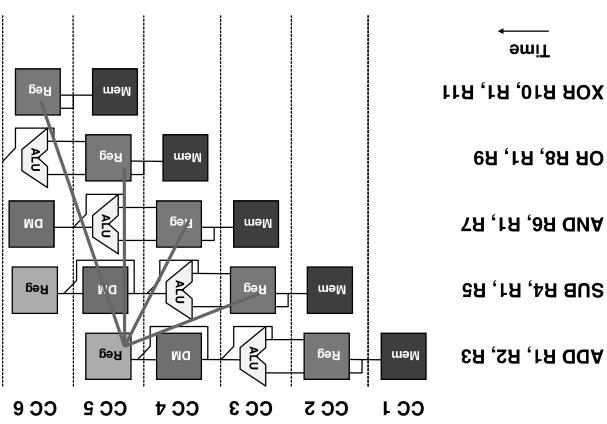
When can we forward?

Can't get data to subtract b/c result needed at beginning of CC #4.



Forwarding doesn't always work

b/c register not written until after those 3 read it.
ADD instruction causes a hazard in next 3 instructions



Illustrating a data hazard

Can't get data to subtract b/c result needed at beginning of CC #4.

- Consider this example:
- ADD R1, R2, R3
- SUB R4, R1, R5
- AND R6, R1, R7
- OR R8, R1, R9
- XOR R10, R1, R11
- Why do they exist??
- These exist because of pipelining

Instructions from un-pipeline machine
- Order differs from order seen by sequentially executing
- Pipelining changes when data operands are read, written
- Order of instructions on un-pipeline machine

• These exist because of pipelining

• Data hazards

- If hazard detected, control unit of pipeline must stall pipeline and prevent instructions in IF, ID from advancing at this stage, controls signals set whether forwarding is needed can also be determined.
- If data hazard, instruction stalled before its issued checked during ID phase of pipeline. For MIPS integer pipeline, all data hazards can be caused by data hazards.

What about control logic?

- In general though, more difficult than register hazards
 - In order, one at a time, read & write in same stage
 - In simple pipeline, memory hazards are easy
- | | | | | | |
|-----------------|---|---|----|---|----|
| Load R1, 0(SP) | F | D | EX | M | WB |
| Store R1, 0(SP) | 1 | 2 | 3 | 4 | 5 |
| | | | | | 6 |
- Seen register hazards, can also have memory hazards
 - RAW:
 - load R4, 0(SP)
 - store R1, 0(SP)

Memory Data Hazards

- With an in-order issue/in-order completion machine, we're not as concerned with WW, WAR
- There are actually 3 different kinds of data hazards:
 - Read After Write (RAW)
 - Write After Write (WAW)
 - Write After Read (WAR)

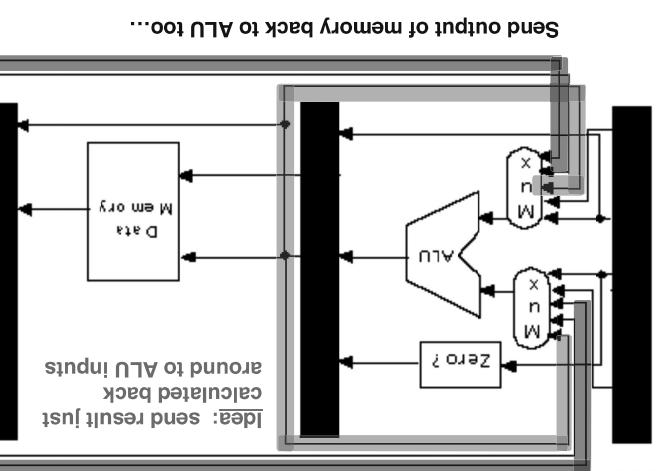
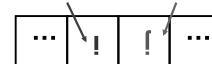
Data hazard specifics

- Thus, j would incorrectly receive old or incorrect value
- With RAW hazard, instruction j tries to read a source operand before instruction i writes it.
- Graphically/Example:

| | | | |
|-----|---|---|-----|
| ... | j | i | ... |
|-----|---|---|-----|

Instruction j is a read instruction issued before i. Instruction i is a write instruction issued after j. Thus, j would incorrectly receive old or incorrect value.
- Can use stalling or forwarding to resolve this hazard

j: ADD R1, R2, R3
i: SUB R4, R1, R6



HW Change for Forwarding

- i.e. compiler could not generate a LOAD instruction that is immediately followed by instruction that uses result of LOAD's destination register.

- Compiler should be able to help eliminate some stalls caused by data hazards.

Data hazards and the compiler

- What happens in the meantime?
- Effect?
- . How long will it take before the branch decision takes

```

    . 72: lw      $4, 50($7)
    . 52: add   $14, $2, $2
    . 48: or     $13, $6, $2
    . 44: and   $12, $2, $5
    . 40: beq   $1, $3, 28 # (28 leads to address 72)
  
```

- Example:
- . Also need to consider hazards involving branches:

- Data transfers

- Arithmetic/logic operations

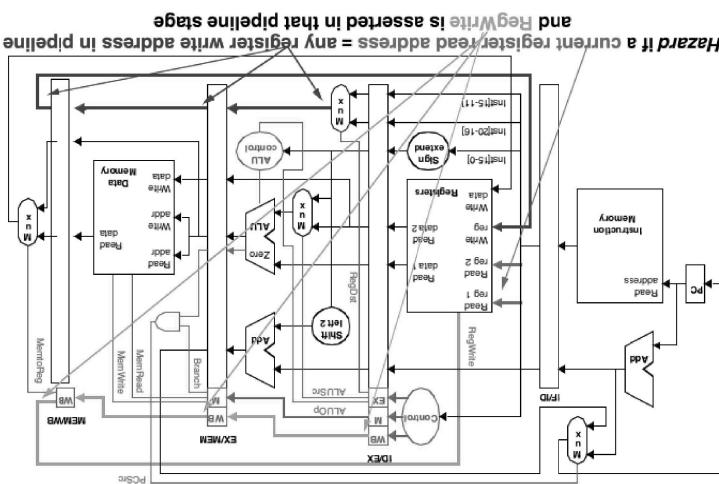
- . So far, we've limited discussion of hazards to:

Branch / Control Hazards



Examples 1-3

Examples...



Detecting Data Hazards

- . For example, can have RAW dependence with or without hazard
- depends on pipeline

- . pipelining depth
- . property of dynamic distance between instructions vs. "in-flight"
- . Potential only exists if instructions can be simultaneously issued
- implies potential for executing things in wrong order
- . hazard: property of program and processor organization

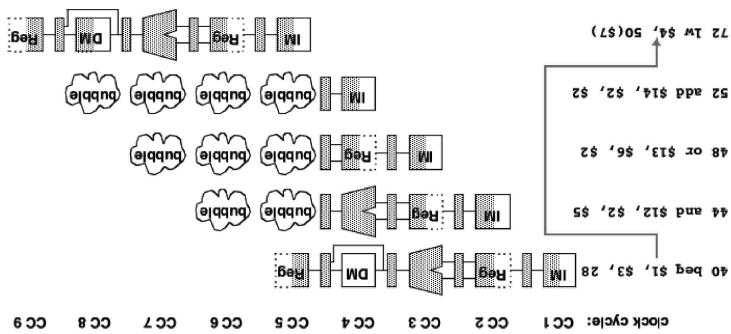
- (i.e., program)
- . dependence: fixed property of instruction stream

Hazards vs. Dependencies

| Situation | Example | Action | No Dependence | Dependence requiring stall | Dependence overcome by forwarding | Dependence overcome by forwarding | Accesses in order |
|---------------|---------------------------------|---|--|--|--|--|--|
| LW R1, 45(R2) | OR R9, R6, R7
SUB R8, R1, R7 | Comparators detect the use of R1 in the ADD and still the use of R1 in SUB and forward the result of LD to the ALU in time for SUB to begin with EX | LD R1, 45(R2)
ADD RS, R6, R7
No action is required because the read of R1 by OR occurs in the second half of the loaded data phase, while the write of the first half. | LD R1, 45(R2)
ADD RS, R6, R7
No action is required because the read of R1 by OR occurs in the second half of the loaded data phase, while the write of the first half. | LD R1, 45(R2)
ADD RS, R6, R7
No action is required because the read of R1 by OR occurs in the second half of the loaded data phase, while the write of the first half. | LD R1, 45(R2)
ADD RS, R6, R7
No action is required because the read of R1 by OR occurs in the second half of the loaded data phase, while the write of the first half. | LD R1, 45(R2)
ADD RS, R6, R7
No action is required because the read of R1 by OR occurs in the second half of the loaded data phase, while the write of the first half. |
| LW R1, 45(R2) | OR R9, R6, R7
SUB R8, R1, R7 | Comparators detect the use of R1 in the ADD and still the use of R1 in SUB and OR before the LD begins EX | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. |
| LW R1, 45(R2) | OR R9, R6, R7
SUB R8, R1, R7 | Comparators detect the use of R1 in the ADD and still the use of R1 in SUB and OR before the LD begins EX | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. | LD R1, 45(R2)
ADD RS, R6, R7
No hazard possible because no dependence exists on R1 in the immediate following three instructions. |

Some example situations

CONTROL HAZARDS



clock cycle: CC1 CC2 CC3 CC4 CC5 CC6 CC7 CC8 CC9

40 beq \$1, \$3, 28 IM Reg DM Reg

44 add \$12, \$55

48 or \$13, \$6, \$2

52 add \$14, \$2, \$2

72 jw \$4, 50(\$7)

- Bubbles injected into 3 stages during cycle 5

Impact of "predict not taken"

- Execution proceeds normally - no penalty

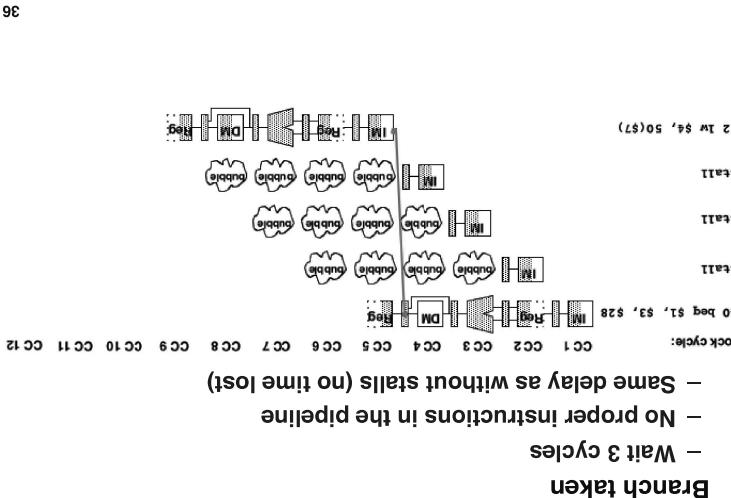
Impact of "predict not taken"

- If prediction is incorrect, just flush the pipeline
- Cuts overall time for branch processing in ½
- Always assume branch will NOT be taken
- One approach:
 - Need to flush improper instruction from pipeline
 - Else, if branch is taken...
 - If branch not taken...
 - On average, branches are taken ½ the time

Dealing w/branch hazards

- Still must wait 3 cycles
- Branch not taken
 - Time lost
 - Could have spent CS fetching, decoding next instructions
- On average, branches are taken ½ the time

Dealing w/branch hazards: always stall



clock cycle: CC1 CC2 CC3 CC4 CC5 CC6 CC7 CC8 CC9

40 beq \$1, \$3, 28 IM Reg DM Reg

- No proper instructions in the pipeline

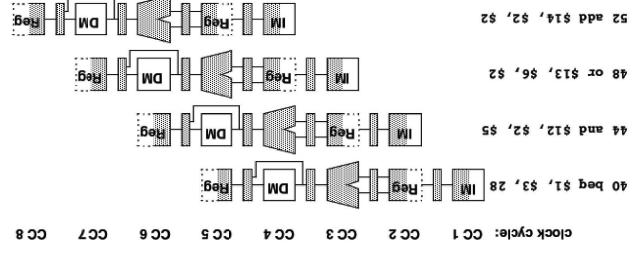
- Wait 3 cycles

- Same delay as without stalls (no time lost)

Dealing w/branch hazards: always stall

- If branch condition true, must skip 44, 48, 52
- But, these have already stalled down the pipeline
- They will complete unless we do something about it
- How do we deal with this?
 - We'll consider 2 possibilities

How branches impact pipelined instructions



clock cycle: CC1 CC2 CC3 CC4 CC5 CC6 CC7 CC8 CC9

44 and \$12, \$55

48 or \$13, \$6, \$2

52 add \$14, \$2, \$2

72 jw \$4, 50(\$7)

- Bubbles injected into 3 stages during cycle 5

Impact of "predict not taken"

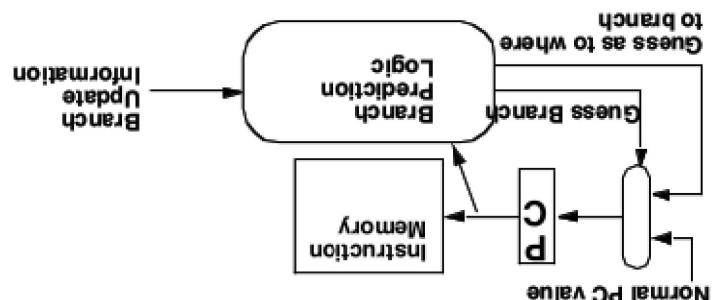
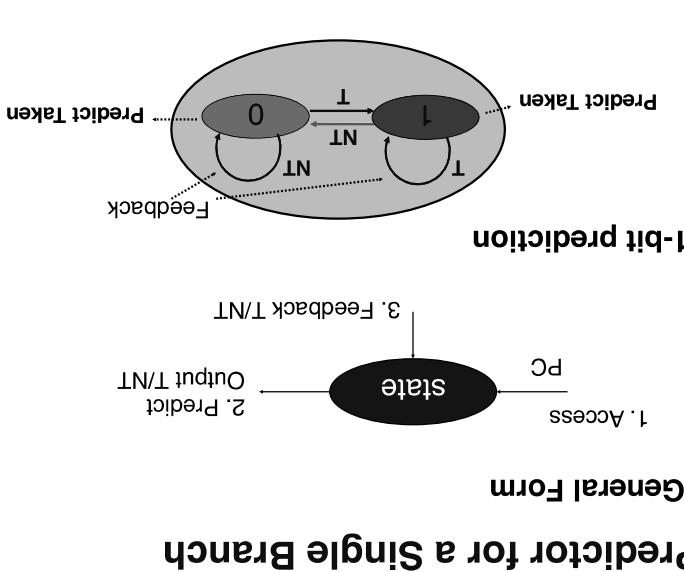
- Only 80% accuracy even if loop 90% of the time
 - When it predicts **exit** instead of looping
 - First time through loop on **next** time through code,
 as before
 }
 - End of loop case, when it exits instead of looping

```
a[i] = a[i] * 2.0;
for (...) {
    for (...) {
        ...
    }
}
```

- Consider a loop of 9 iterations before exit:

- in a loop, 1-bit BHT will cause 2 mispredictions

1 bit weakness

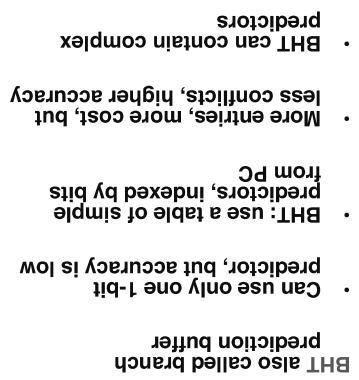


- Prior solutions are "ugly"
- Better (& more common): guess possible outcome
- "Predictor" to guess where / if instruction will branch
- "Recovery Mechanism":
 - (and to where)
 - "Recovery Mechanism":
 - i.e. a way to fix your mistake
- Alternative: flush instructions if branch taken
- Recovery: always guess branch never taken
- Predictor: always guess branch never taken
- To where it is taken
- Whether or not for any particular PC value a branch was taken next
- How to update with information from later stages



Branch Prediction

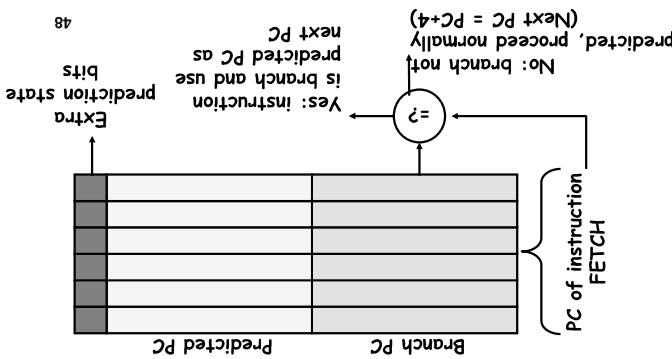
Branch Penalty Impact



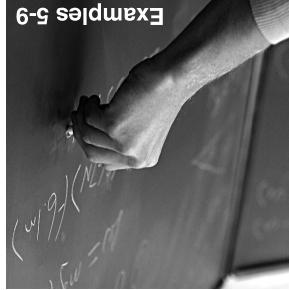
- How does instruction set design impact pipelining?
- Does increasing the depth of pipelining always increase performance?

Discussion

Branch target buffer

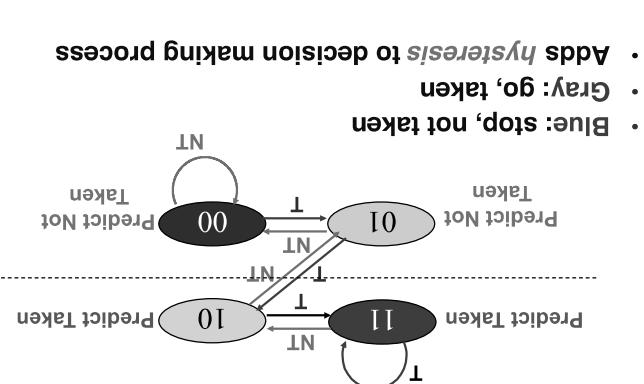


- Example: BTB combined with BHT
- Branch Target Buffer (BTB): Address of branch index to get prediction AND branch address (if taken)
 - Note: must check for branch match now, since can't use wrong branch address



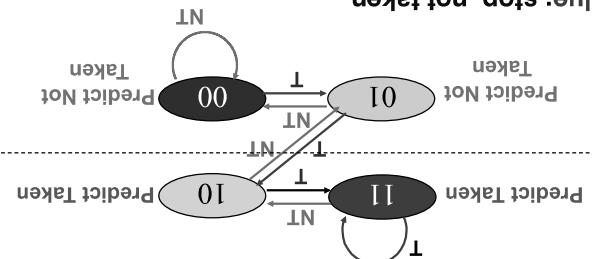
Examples...

- Solution: 2-bit scheme where change prediction only if get misprediction twice:

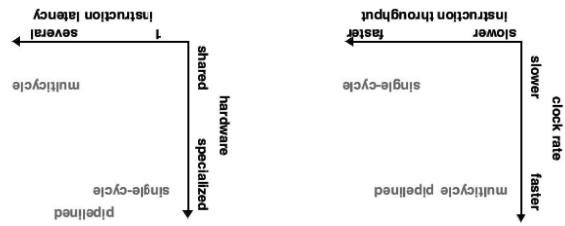


- If get misprediction twice:

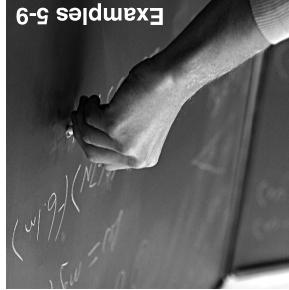
- Gray: go, taken
- Blue: stop, not taken
- Adds hysteresis to decision making process



- Throughput: instructions per clock cycle = $1/\text{CPI}$
- Pipeline has fast throughput and fast clock rate
 - Pipeline latency has fast throughput and fast clock rate
 - High latency causes problems in cycles
 - Increased time to resolve hazards



Comparative performance

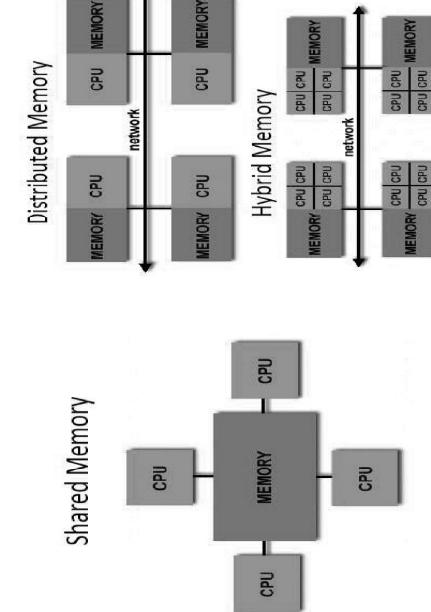


Instruction Cycle

- Parallel computers are those that emphasize the parallel processing between the operations in some way.
- Various Varieties
 - Shared Memory
 - Distributed Memory
 - Hybrid Memory
- Start
 - Calculate the address of the instruction to be executed
 - Fetch the instruction
 - Decode the instruction
 - Calculate the operand address
 - Fetch the operands
 - Execute the instructions
 - Store the results
 - If more instructions to be executed, go to step 2 else stop.

Flynn's Classification

- SISD: Single Instruction Single Data
 - Classical Von-Neumann architecture
- SIMD: Single Instruction Multiple data
- MISD: Multiple Instructions Single Data
- MIMD: Multiple Instructions Multiple Data
 - Most common and general parallel machine

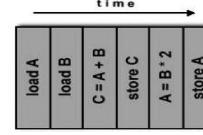


Flynn's Taxonomy

- This classification was first studied and proposed by Michael Flynn in 1972.
- Flynn did not consider the machine architecture for classification of parallel computers
- He introduced the concept of instruction and data streams for categorizing of computers.
- Instruction stream - a flow of instructions from main memory to the CPU
- Data stream - a flow of operands between processor and memory
- All the computers classified by Flynn are not parallel computers
- Let I_s and D_s are minimum number of streams flowing at any point in the execution

Single Instruction Single Data stream

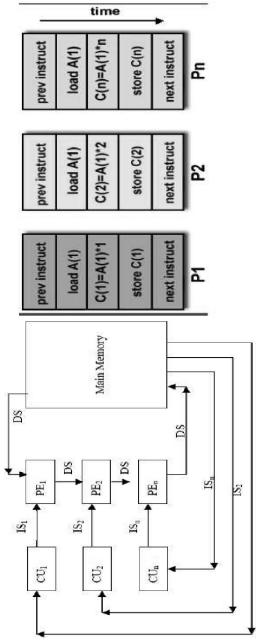
- A serial computer
- SISD machines are conventional serial computers that process only one stream of instructions and one stream of data.
- $I_s = D_s = 1$
- Examples
 - CDC 6600 which is unipipelined but has multiple functional units.
 - CDC 7600 which has a pipelined arithmetic unit.
 - Amithal 470/6 which has pipelined instruction processing.
 - Cray-1 which supports vector processing.



Single Instruction Multiple Data stream

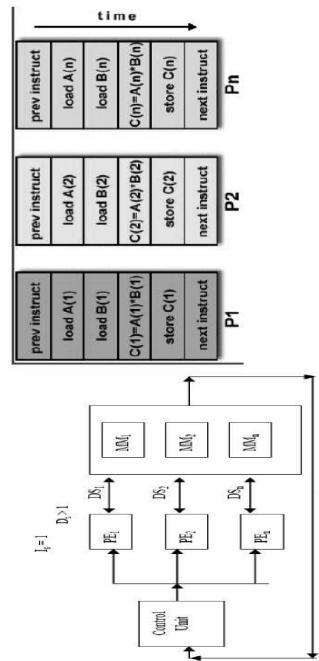
- Multiple processing elements work under the control of a single control unit.
- It has one instruction and multiple data stream
- Main memory can also be divided into modules for generating multiple data streams acting as a distributed memory
- Examples of SIMD organisation are ILLIAC-V, PEPE, BSP, STARAN, MPP, DAP and the Connection Machine (CM-1).
- A type of parallel computer
- Single instruction: All processing units execute the same instruction at any given clock cycle
- Multiple data: Each processing unit can operate on a different data element

MISD



Real time computers need to be fault tolerant where several processors execute the same data for producing the redundant data. This is also known as N-version programming. All these redundant data are compared as results which should be same; otherwise faulty unit is replaced. Thus MISD machines can be applied to fault tolerant real time computers.

SIMD



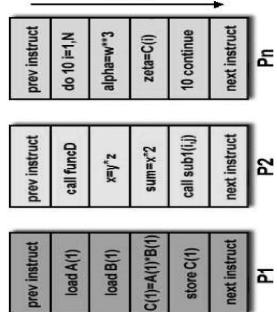
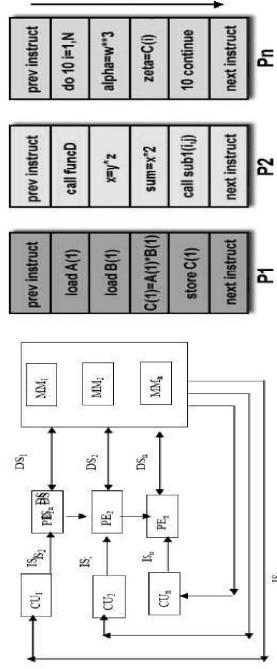
Multiple Instruction Multiple Data stream

- Multiple instruction streams operate on multiple data streams
- The processors work on their own data with their own instructions.
- Tasks executed by different processors can start or finish at different times.
- This classification actually recognizes the parallel computer.
- Examples include: C.mmp, Burroughs D825, Cray-2, S1, Cray X-MP, HEP, Pluribus, IBM 370/168 MP, Univac 1100/80, Tandem/16, IBM 3081/3084, C.m*, BBN Butterfly, Melko Computing Surface (CS-1), FPS/T40000, iPSC.
- MIMD organization is the most popular for a parallel computer.
- In the real sense, parallel computers execute the instructions in MIMD mode.
- $I_s > 1$, $D_s > 1$

Multiple Instruction Single Data stream

- Multiple processing elements are organised under the control of multiple control units.
- Each control unit is handling one instruction stream and processed through its corresponding processing element.
- A single data stream is fed into multiple processing units.
- Each processing unit operates on the data independently via independent instruction streams.
- But each processing element is processing only a single data stream at a time
- All processing elements are interacting with the common shared memory.

MIMD



What is RAID

- Redundant Array of Independent (Inexpensive) Disks
- A set of disk stations treated as one logical station
- Data are distributed over the stations
- Redundant capacity is used for parity allowing for data repair

Parity

- Way to do error checking and correction
- Add up all the bits that are 1
 - if even number, set parity bit to 0
 - if odd number, set parity bit to 1
- To actually implement this, do an exclusive OR of all the bits being considered
- Consider the following 2 bytes
 - byte parity
 - 10110011 1
 - 01101010 0
- If a single bit is bad, it is possible to correct it

RAID (Redundant Array of Independent Disks)

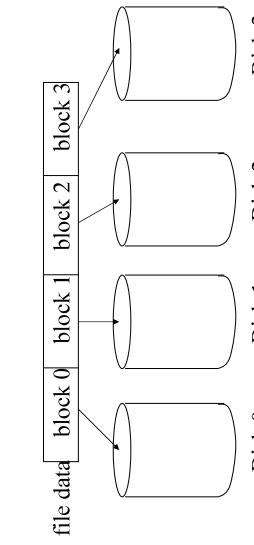
- Redundant array of inexpensive disks
- Multiple disk database design
- Set of physical disk drives viewed by the OS as a single logical drive
- Data are distributed across the physical drives of an array
- Improve access time and improve reliability
 - large storage capacity
 - redundant data
- 7 levels (6 levels in common use)
 - differing levels of redundancy, error checking, capacity, and cost

Mirroring

- Keep two copies of data on two separate disks
- Gives good error recovery
 - if some data is lost, get it from the other source
- Expensive
 - requires twice as many disks
- Write performance can be slow
 - have to write data to two different spots
- Read performance is enhanced
 - can read data from file in parallel

Striping

- Take file data and map it to different disks
- Allows for reading data in parallel



Levels of RAID

- 6 levels of RAID (0-5) have been accepted by the industry
- Other kinds have been proposed in literature
 - Level 2 and 4 are not commercially available, they are included for clarity

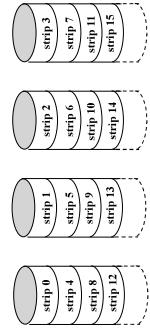
RAID 0

- All data (user and system) are distributed over the disks so that there is a reasonable chance for parallelism
- Disk is logically a set of strips (blocks, sectors,...). Strips are numbered and assigned consecutively to the disks (see picture.)

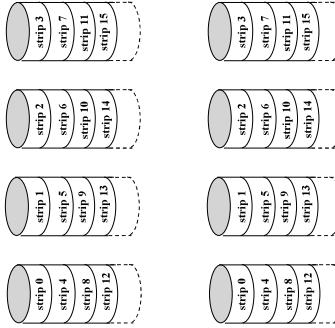
RAID 0

- Block level striping
- Without parity or mirroring
- No redundancy, No backup, No fault tolerance
- Improved performance, faster, as it is uses block level striping
- Maximum use of storage space as there is no backup
- Any drive failure destroys the array

Raid 0 (No redundancy)

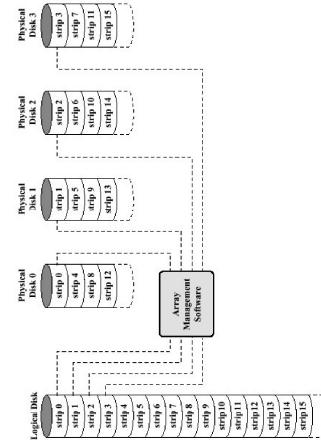


Raid 1 (mirrored)



Raid 1 (mirrored)

Data mapping Level 0



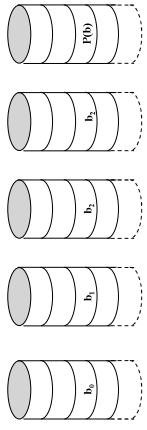
RAID 1

- RAID 1 does not use parity, it simply mirrors the data to obtain reliability
- Plus:
 - Reading request can be served by any of the two disks containing the requested data (minimum search time)
 - Writing request can be performed in parallel to the two disks: no “writting penalty”
 - Recovery from error is easy, just copy the data from the correct disk
- Minus:
 - Price for disks is doubled
 - Will only be used for system critical data that must be available at all times

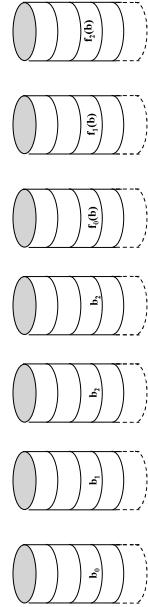
Figure 11.9 Data Mapping for a RAID Level 0 Array [MASS94]

RAID 3 (bit-interleaved parity)

- Mirroring
- Without parity
- Data is written identically to two drives (parallel write)
- Not slower, not faster (Write)
- Read can be faster by parallel access
- Fault tolerance
- 50% storage space can be used



Raid 2 (redundancy through Hamming code)



RAID 3

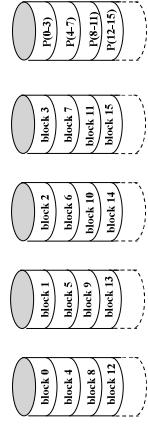
- Level 2 needs $\log_2(\text{number of disks})$ parity disks
 - Level 3 needs only one, for one parity bit
- In case one disk crashes, the data can still be reconstructed even on line ("reduced mode") and be written ($X1\text{-}4$ data, P parity):
$$P = X1+X2+X3+X4$$
$$X1=P+X2+X3+X4$$
- RAID 2-3 have high data transfer times, but perform only one I/O at the time so that response times in transaction oriented environments are not so good

$$P = X1+X2+X3+X4$$

$$X1=P+X2+X3+X4$$

RAID 4 (block-level parity)

- Small strips, one byte or one word
- Synchronized disks, each I/O operation is performed in a parallel way
 - Error correction code (Hamming code) allows for correction of a single bit error
- Controller can correct without additional delay
- Is still expensive



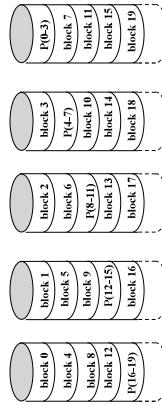
RAID 2

RAID 4

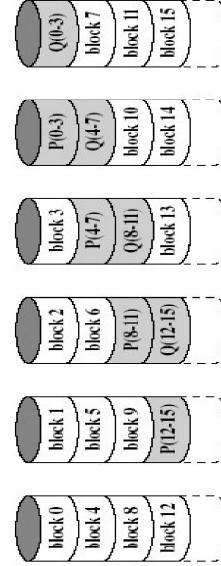
- Larger strips and one parity disk
- Blocks are kept on one disk, allowing for parallel access by multiple I/O requests
- Writing penalty: when a block is written, the parity disk must be adjusted (e.g. writing on X1):
$$\begin{aligned} P &= X4+X3+X2+X1 \\ &= X4+X3+X2+X1' + X1+X1 \end{aligned}$$
- Parity disk may be a bottleneck
- Good response times, less transfer rates

- Faster- Block level stripping
- Fault tolerance by Distributed parity (backup of data)
- Distribution of the parity strip to avoid the bottle neck.
- The array is not destroyed by the failure of single hard drive
 - Minimum three disks
 - Loose one hard disk space for parity
 - 75% space

RAID 5 (block-level distributed parity)



RAID 6



RAID 5

- Distribution of the parity strip to avoid the bottle neck.

- Two different parity calculations are carried out and stored in separate blocks on different disks.
 - Example: XOR and an independent data check algorithm => makes it possible to regenerate data even if two disks containing user data fail.
- No. of disks required = $N + 2$ (where N = number of disks required for data).
- Provides HIGH data availability.
- Incur substantial write penalty as each write affects two parity blocks.
- Three disks would have to fail within MTTR (mean time to repair) interval to cause data to be lost

RAID Level 6

- Faster-Block level striping
- Provides fault tolerance up to two failed drives Fault tolerance- double distributed parity (backup of backup of data)
- High availability systems
- Failure of hard drive, slow down the performance of the system
- Loose two hard disk space for parity
- Incurs substantial write penalty

Overview Raid 0-2

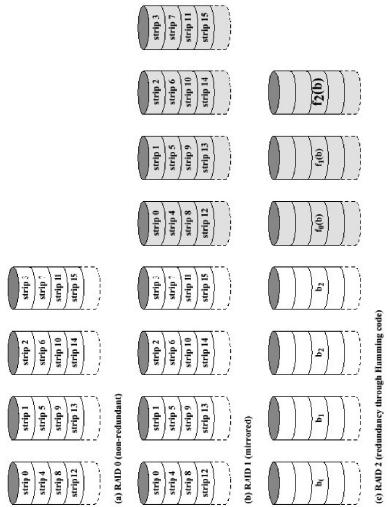


Figure 11.8 RAID Levels (page 1 of 2)

Overview Raid 3-5

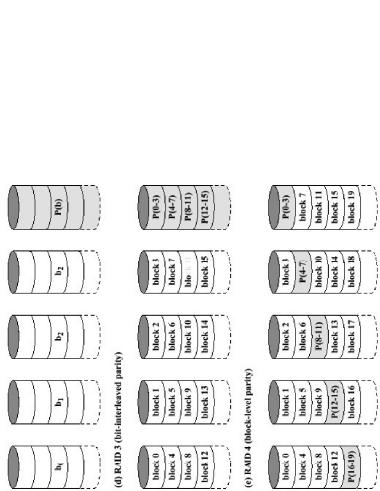


Figure 11.8 RAID Levels (page 2 of 2)