
Computer Architecture I/O Systems

Outline of Today's Lecture

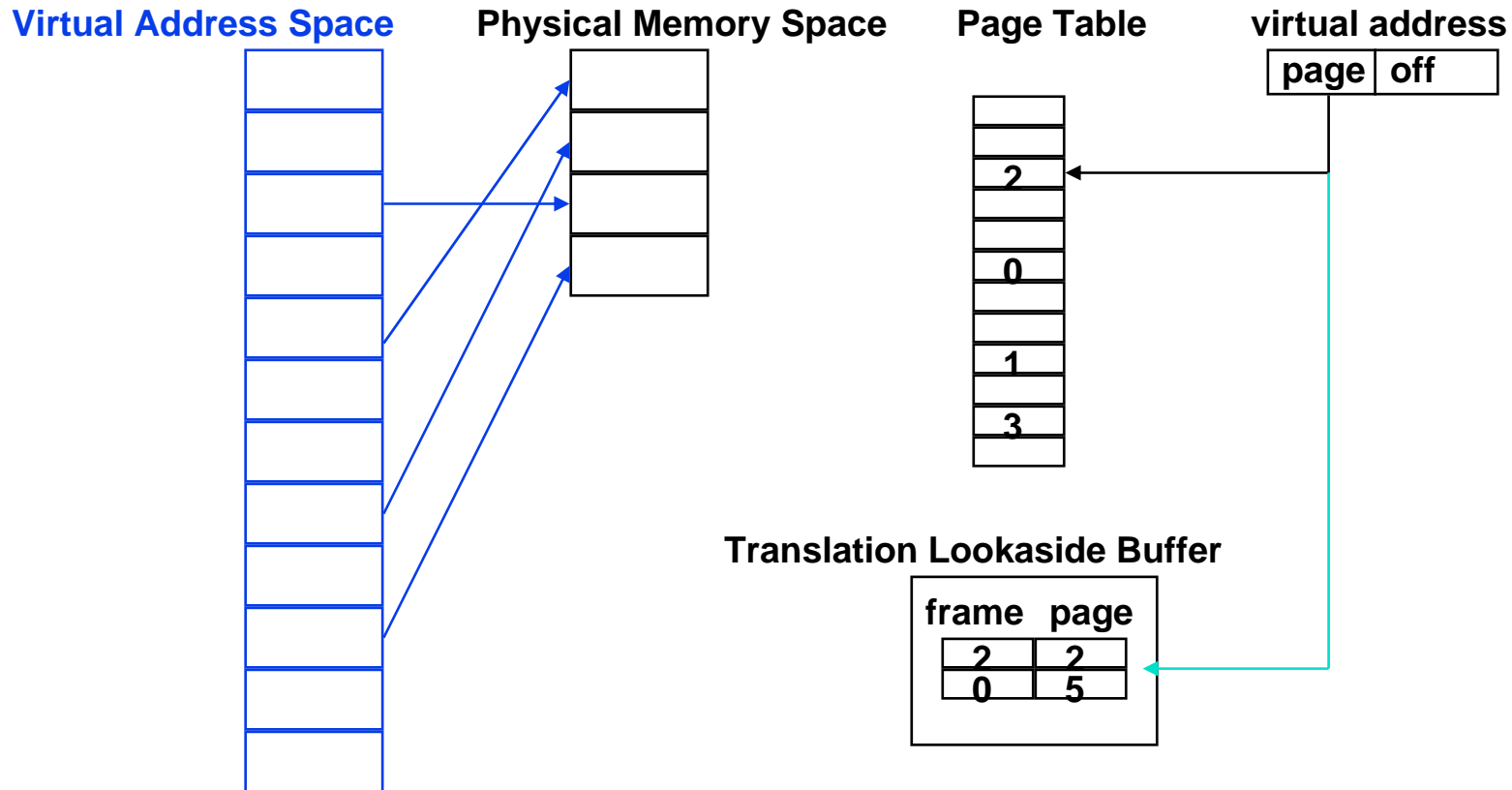
- **Recap Busses**
- **Finish VM: TLBs**
- **Questions and Administrative Matters**
- **I/O Performance Measures**
- **Types and Characteristics of I/O Devices**
- **Magnetic Disks**
- **Break**
- **DMA, Multimedia, OS**
- **Summary**

Recap: Busses

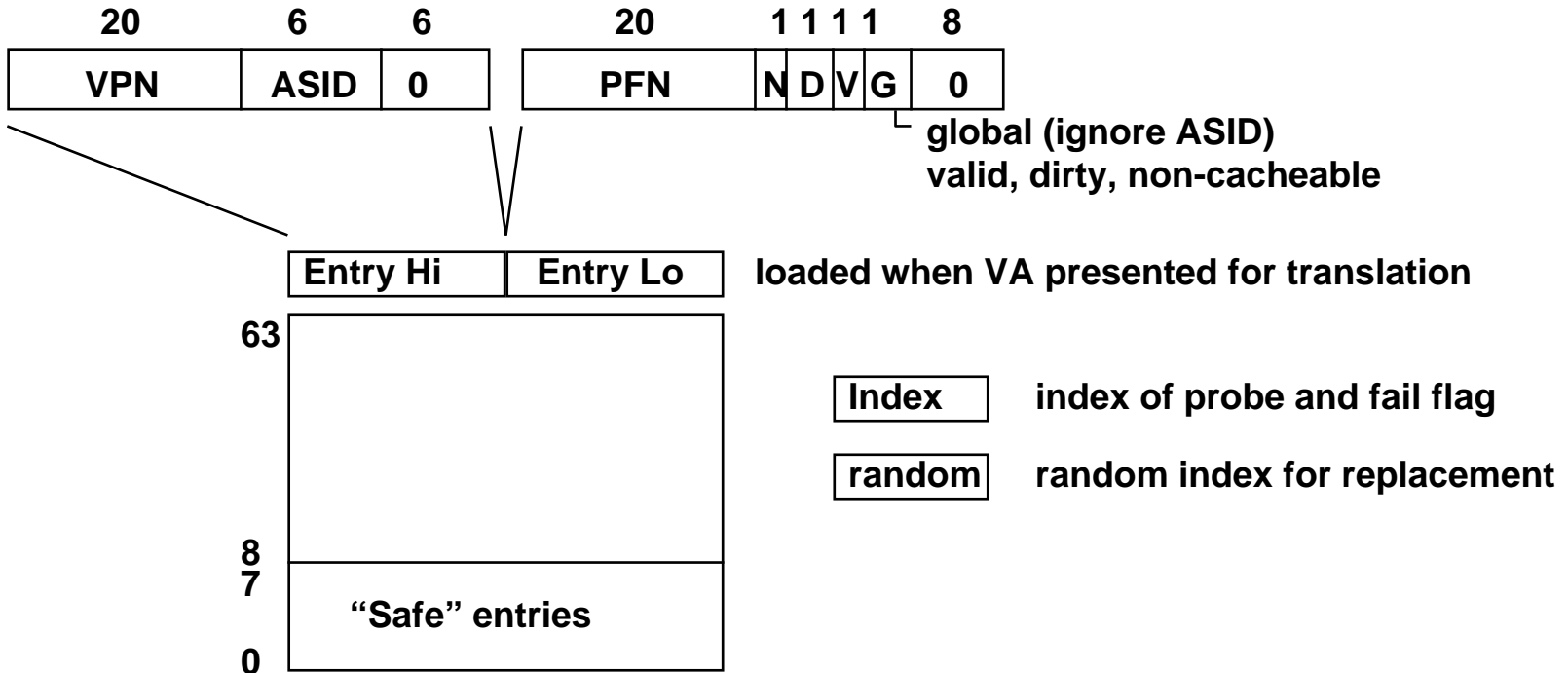
- **Fundamental tool for designing and building computer systems**
 - **divide the problem into independent components operating against a well defined interface**
 - **processor, memory, I/O**
 - **compose the components efficiently**
- **Shared collection of wires**
 - **command, address, data**
- **Communication path between multiple subsystems**
- **Inexpensive**
- **Limited bandwidth**
- **Layers of a bus specification**
 - **mechanical, electrical, signalling, timing, transactions**

Making address translation practical: TLB

- Virtual memory => memory acts like a cache for the disk
- Page table maps virtual page numbers to physical frames
- Translation Look-aside Buffer (TLB) is a cache of recent translations

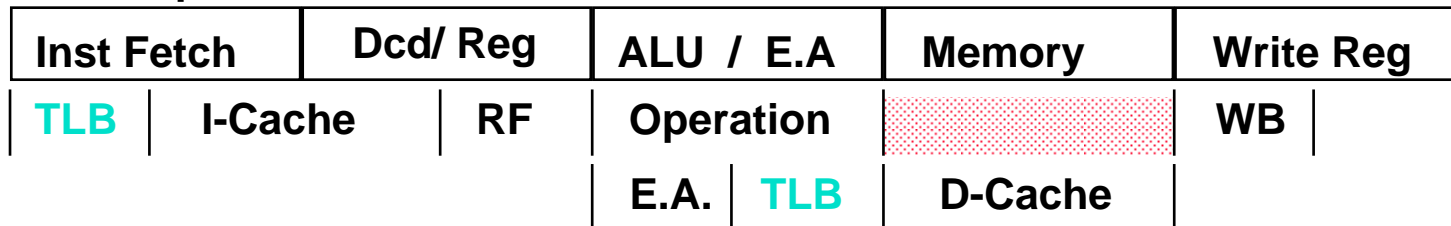


R3000 TLB & CPO (MMU)



Constraints on TLB organization

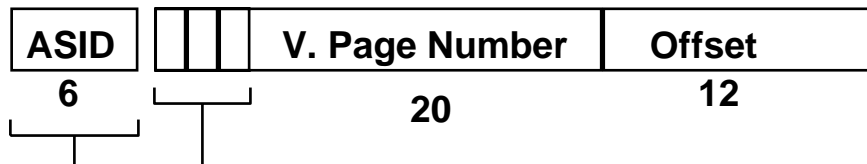
MIPS R3000 Pipeline



TLB

64 entry, on-chip, fully associative, software TLB fault handler

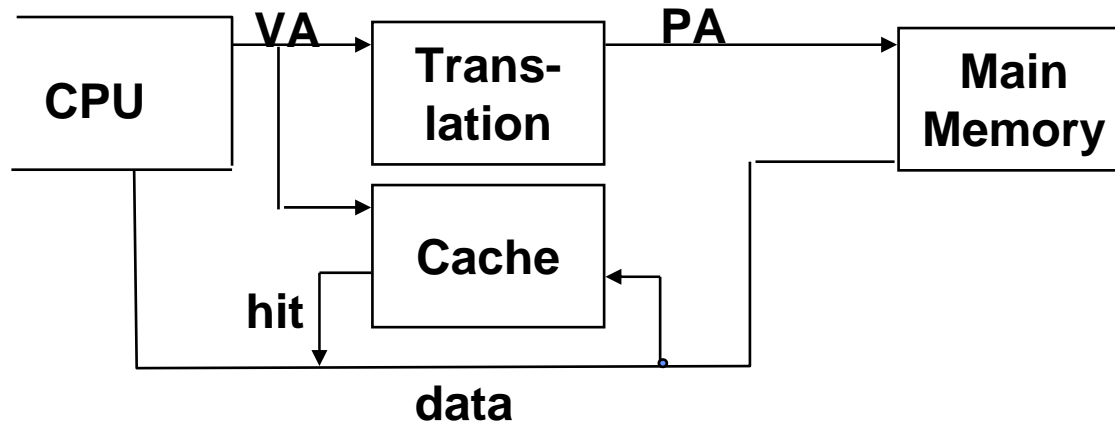
Virtual Address Space



- 0xx User segment (caching based on PT/TLB entry)
- 100 Kernel physical space, cached
- 101 Kernel physical space, uncached
- 11x Kernel virtual space

Allows context switching among
64 user processes without TLB flush

Virtually Addressed Cache



Only require address translation on cache miss!

synonym problem: two different virtual addresses map to same physical address => two different cache entries holding data for the same physical address!

nightmare for update: must update all cache entries with same physical address or memory becomes inconsistent

determining this requires significant hardware, essentially an associative lookup on the physical address tags to see if you have multiple hits.

(usually disallowed by fiat)

Optimal Page Size

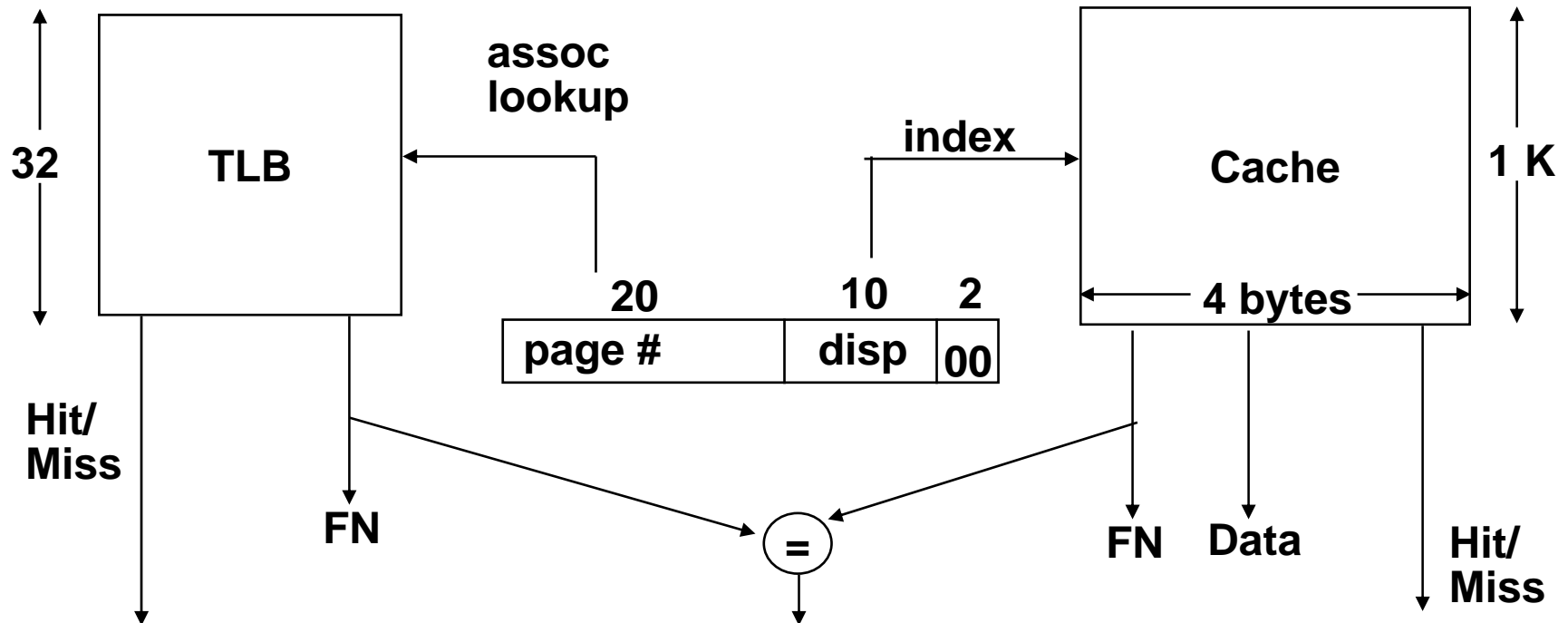
- **Mimimize wasted storage**
 - **small page minimizes internal fragmentation**
 - **small page increase size of page table**
- **Minimize transfer time**
 - **large pages (multiple disk sectors) amortize access cost**
 - **sometimes transfer unnecessary info**
 - **sometimes prefetch useful data**
 - **sometimes discards useless data early**

General trend toward larger pages because

- **big cheap RAM**
- **increasing mem / disk performance gap**
- **larger address spaces**

Overlapped TLB & Cache Access

- So far TLB access is serial with cache access
 - can we do it in parallel?
 - only if we are careful in the cache organization!



What if cache size is increased to 8KB?

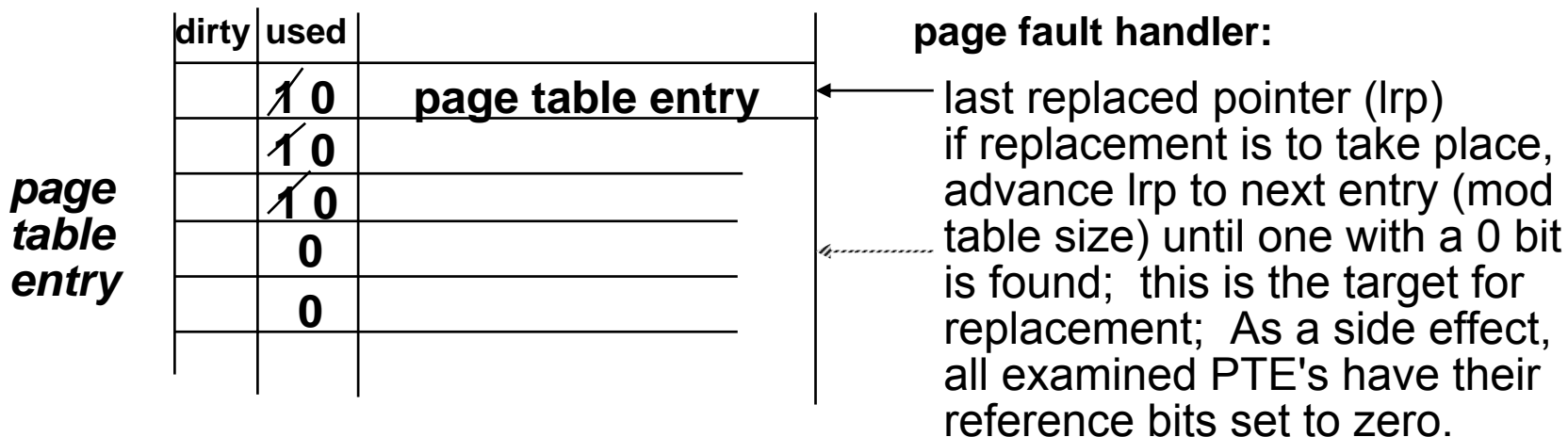
Page Fault: What happens when you miss?

- Not talking about TLB miss
 - TLB is HWs attempt to make page table lookup fast (on average)
- Page fault means that page is not resident in memory
- Hardware must detect situation
- Hardware cannot remedy the situation
- Therefore, hardware must trap to the operating system so that it can remedy the situation
 - pick a page to discard (possibly writing it to disk)
 - load the page in from disk
 - update the page table
 - resume to program so HW will retry and succeed!
- What is in the page fault handler?
 - see CS162
- What can HW do to help it do a good job?

Page Replacement: Not Recently Used (1-bit LRU, Clock)

Associated with each page is a reference flag such that
ref flag = 1 if the page has been referenced in recent past
= 0 otherwise

-- if replacement is necessary, choose any page frame such that its reference bit is 0. This is a page that has not been referenced in the recent past



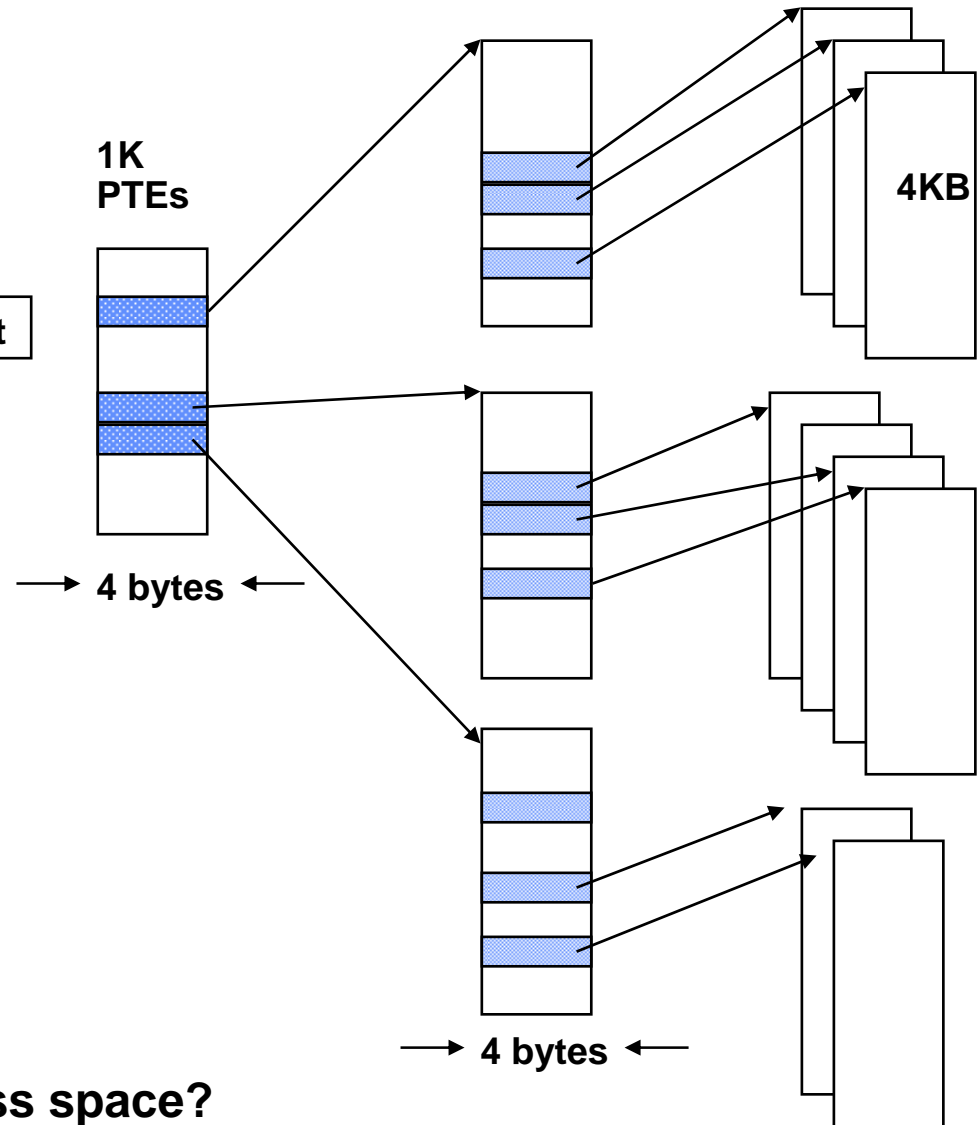
Or search for the a page that is both
not recently referenced AND not dirty.

**Architecture part: support dirty and used bits in the page table
=> may need to update PTE on any instruction fetch, load, store
How does TLB affect this design problem? Software TLB miss?**

Large Address Spaces

Two-level Page Tables

32-bit address:



- 2 GB virtual address space
- 4 MB of PTE2
 - paged, holes
- 4 KB of PTE1

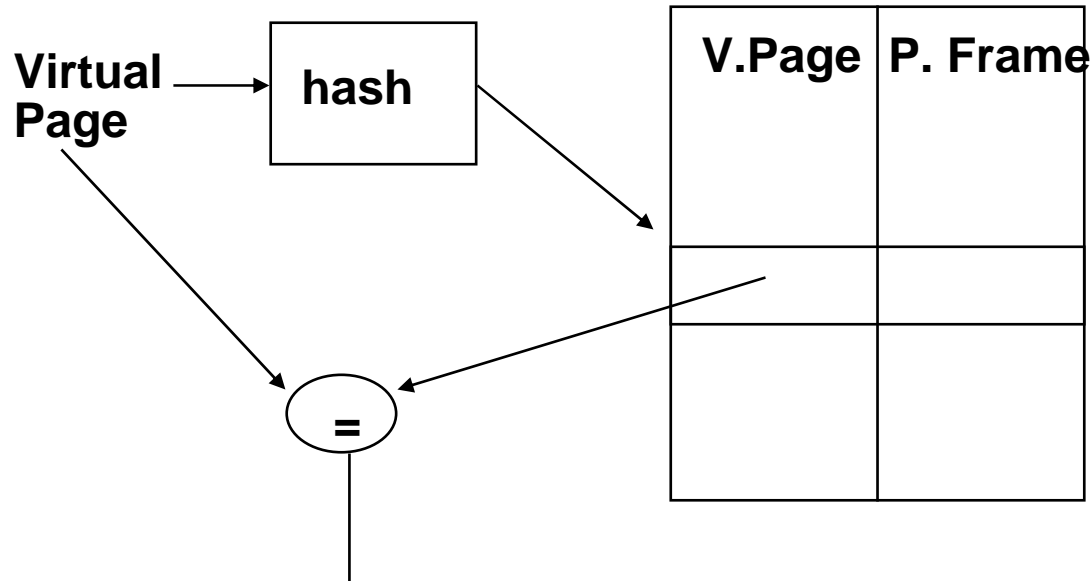
What about a 48-64 bit address space?

Inverted Page Tables

IBM System 38 (AS400) implements 64-bit addresses.

48 bits translated

start of object contains a 12-bit tag



=> TLBs or virtually addressed caches are critical

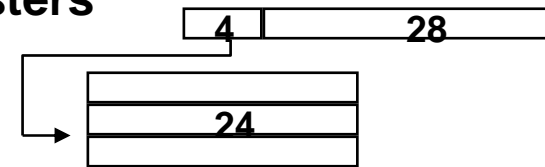
Survey

◦ R4000

- 32 bit virtual, 36 bit physical
- variable page size (4KB to 16 MB)
- 48 entries mapping page pairs (128 bit)

◦ MPC601 (32 bit implementation of 64 bit PowerPC arch)

- 52 bit virtual, 32 bit physical, 16 segment registers
- 4KB page, 256MB segment
- 4 entry instruction TLB
- 256 entry, 2-way TLB (and variable sized block xlate)
- overlapped lookup into 8-way 32KB L1 cache
- hardware table search through hashed page tables



◦ Alpha 21064

- arch is 64 bit virtual, implementation subset: 43, 47, 51, 55 bit
- 8, 16, 32, or 64KB pages (3 level page table)
- 12 entry ITLB, 32 entry DTLB
- 43 bit virtual, 28 bit physical octword address

Hardware / Software Boundary

- **What aspects of the Virtual -> Physical Translation is determined in hardware?**
 - **TLB Format**
 - **Type of Page Table**
 - **Page Table Entry Format**
 - **Disk Placement**
 - **Paging Policy**

Why virtual memory?

- **Generality**
 - ability to run programs larger than size of physical memory
- **Storage management**
 - allocation/deallocation of variable sized blocks is costly and leads to (external) fragmentation
- **Protection**
 - regions of the address space can be R/O, Ex, . . .
- **Flexibility**
 - portions of a program can be placed anywhere, without relocation
- **Storage efficiency**
 - retain only most important portions of the program in memory
- ◦ **Concurrent I/O**
 - execute other processes while loading/dumping page
- ◦ **Expandability**
 - can leave room in virtual address space for objects to grow.
- ◦ **Performance**

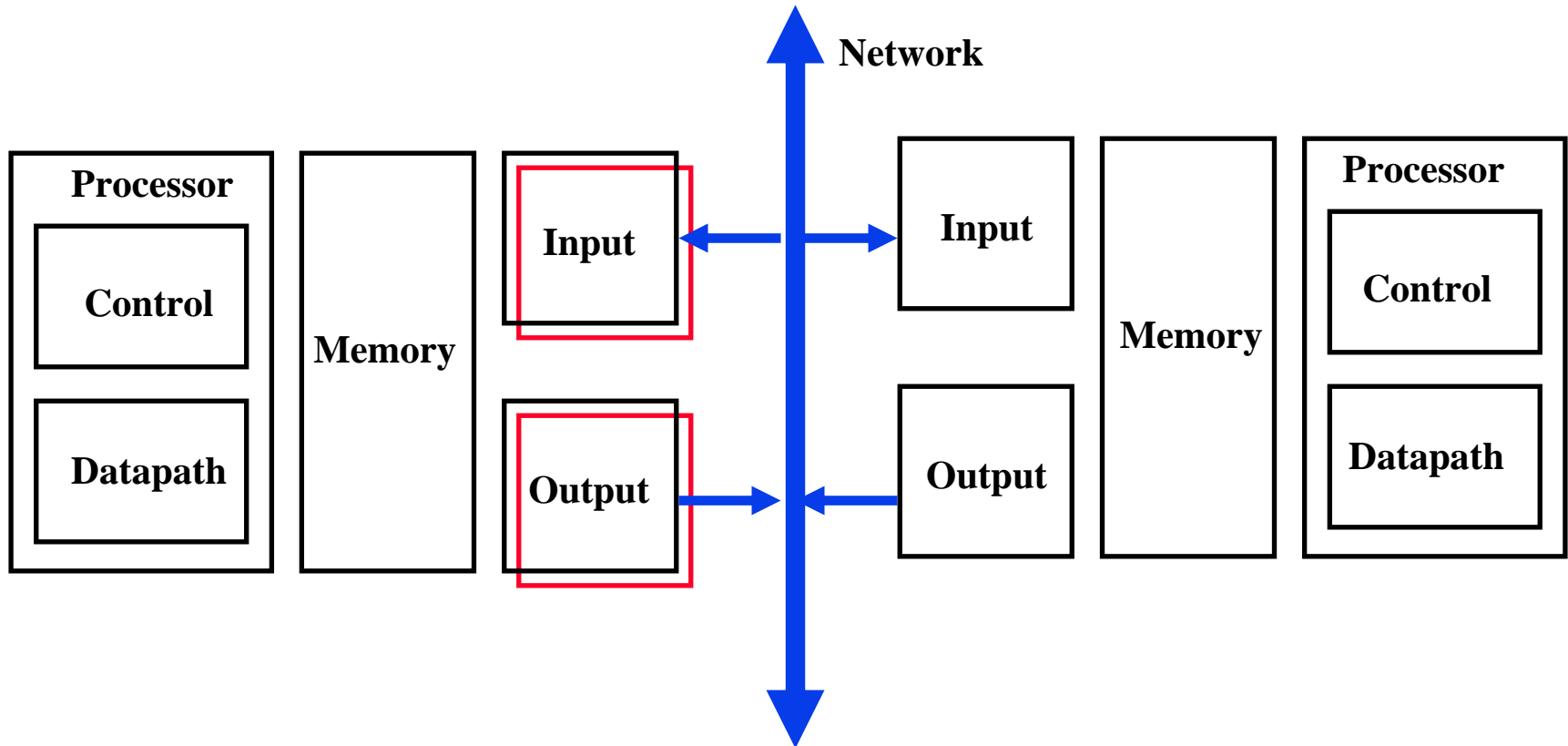
Observe: impact of multiprogramming, impact of higher level languages

Administrative Issues

- Read P&H Chapter 8
- Intel “field trip” Friday 11/21
 - bus at 7 am !

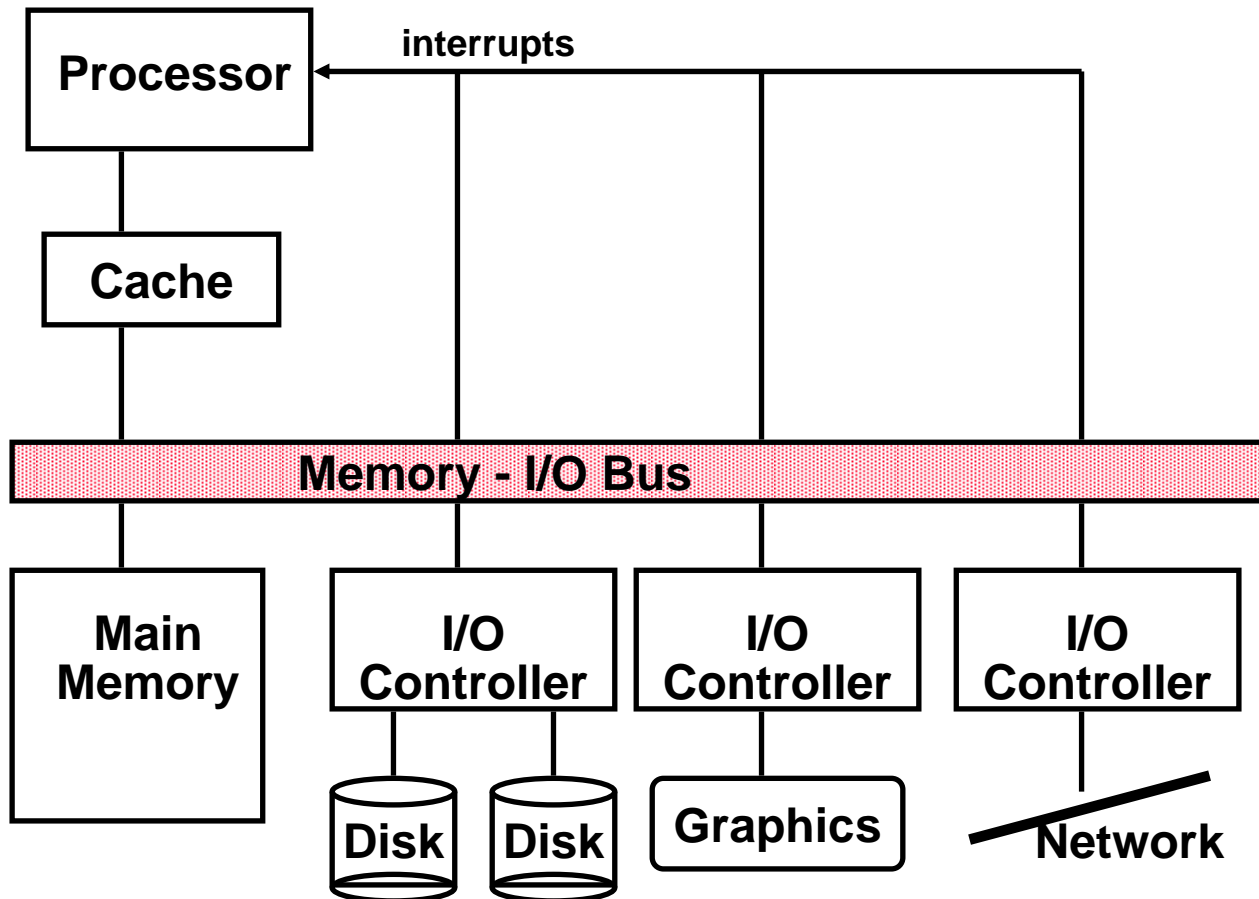
The Big Picture: Where are We Now?

- Today's Topic: I/O Systems



I/O System Design Issues

- Performance
- Expandability
- Resilience in the face of failure



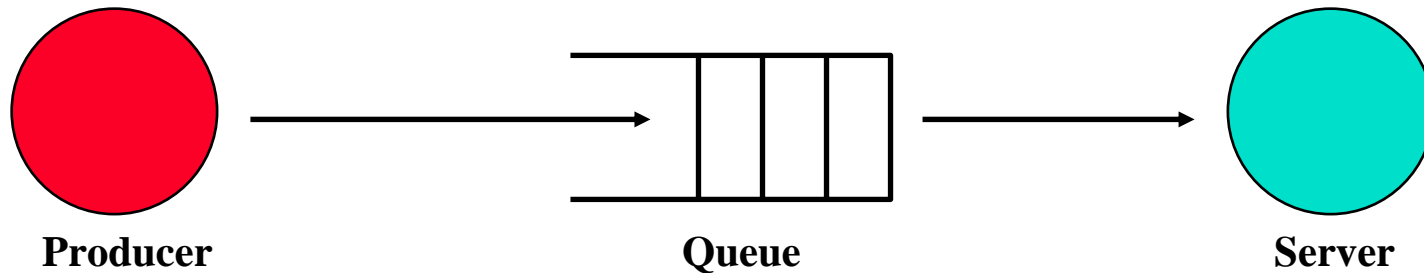
I/O Device Examples

Device	Behavior	Partner	Data Rate (KB/sec)
Keyboard	Input	Human	0.01
Mouse	Input	Human	0.02
Line Printer	Output	Human	1.00
Floppy disk	Storage	Machine	50.00
Laser Printer	Output	Human	100.00
Optical Disk	Storage	Machine	500.00
Magnetic Disk	Storage	Machine	5,000.00
Network-LAN	Input or Output	Machine	20 – 1,000.00
Graphics Display	Output	Human	30,000.00

I/O System Performance

- **I/O System performance depends on many aspects of the system (“limited by weakest link in the chain”):**
 - **The CPU**
 - **The memory system:**
 - **Internal and external caches**
 - **Main Memory**
 - **The underlying interconnection (buses)**
 - **The I/O controller**
 - **The I/O device**
 - **The speed of the I/O software (Operating System)**
 - **The efficiency of the software’s use of the I/O devices**
- **Two common performance metrics:**
 - **Throughput: I/O bandwidth**
 - **Response time: Latency**

Simple Producer-Server Model



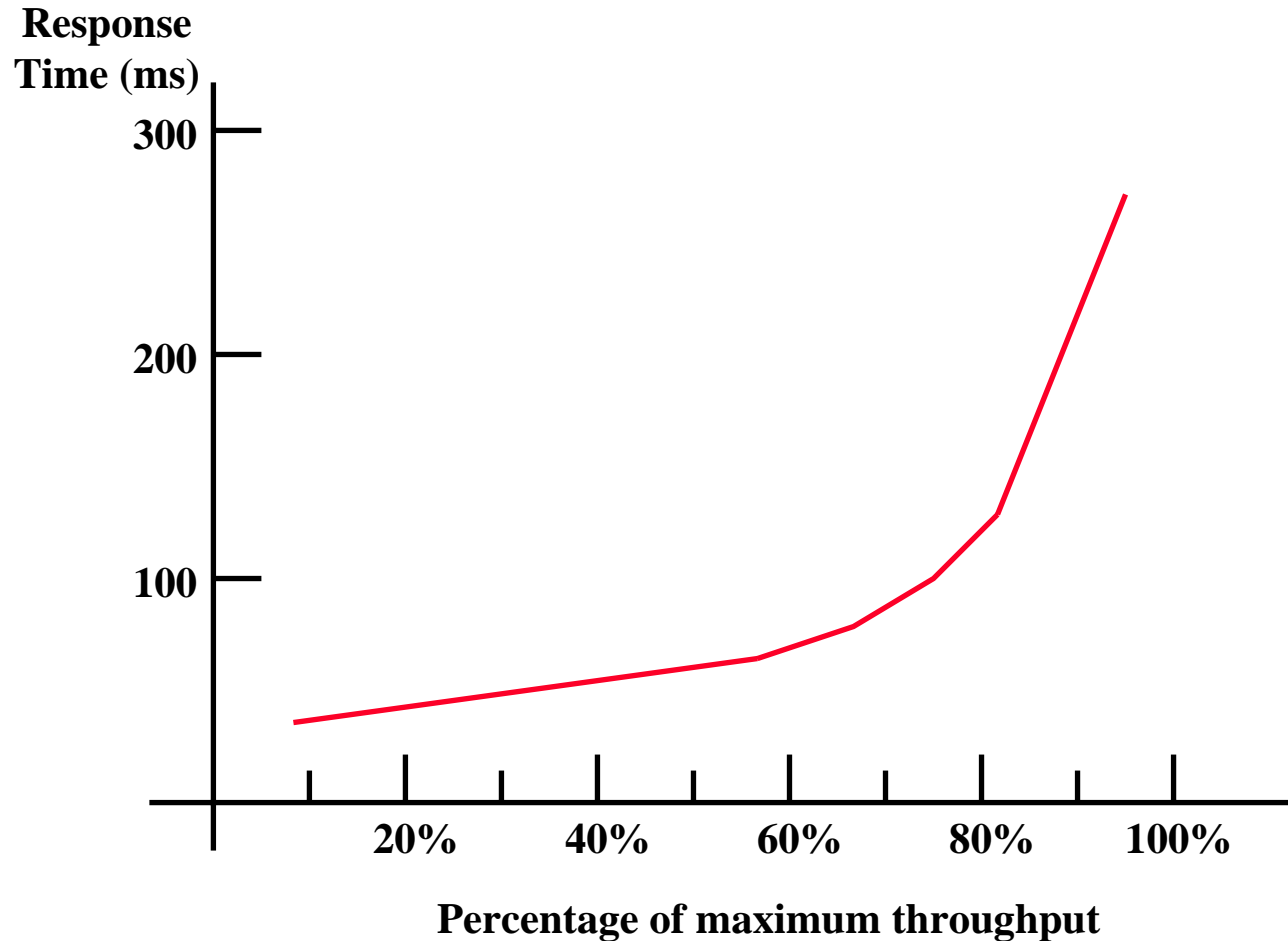
◦ Throughput:

- The number of tasks completed by the server in unit time
- In order to get the highest possible throughput:
 - The server should never be idle
 - The queue should never be empty

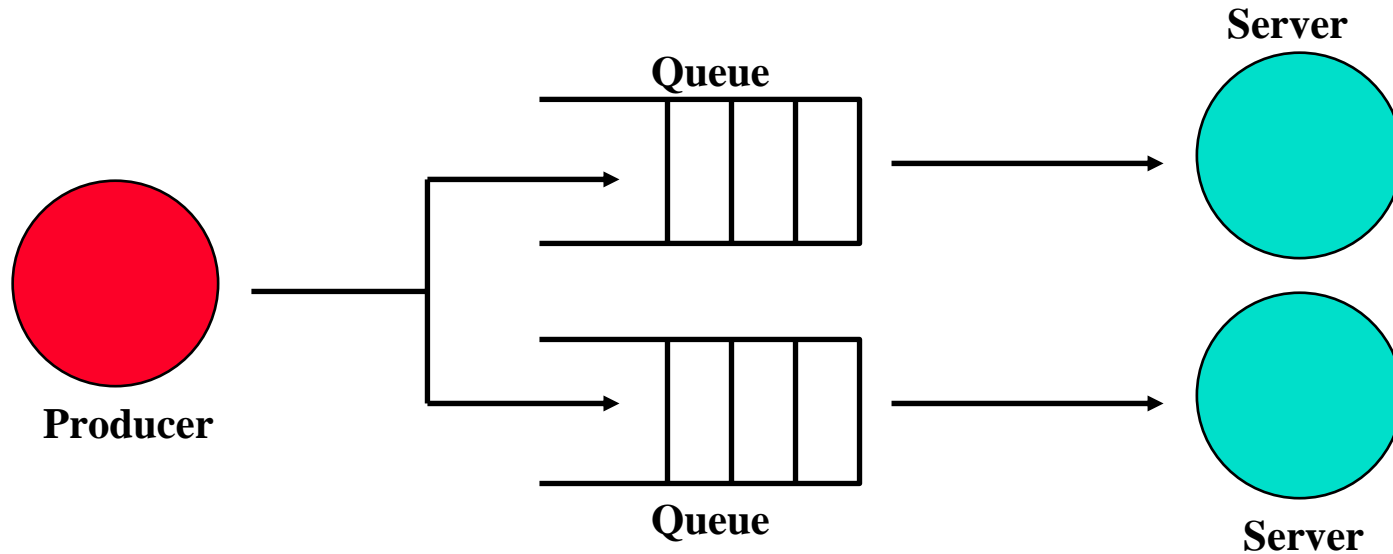
◦ Response time:

- Begins when a task is placed in the queue
- Ends when it is completed by the server
- In order to minimize the response time:
 - The queue should be empty
 - The server will be idle

Throughput versus Response Time



Throughput Enhancement



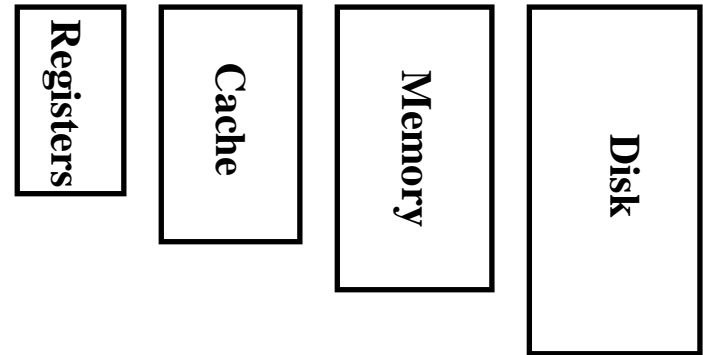
- In general throughput can be improved by:
 - Throwing more hardware at the problem
 - reduces load-related latency
- Response time is much harder to reduce:
 - Ultimately it is limited by the speed of light (but we're far from it)

I/O Benchmarks for Magnetic Disks

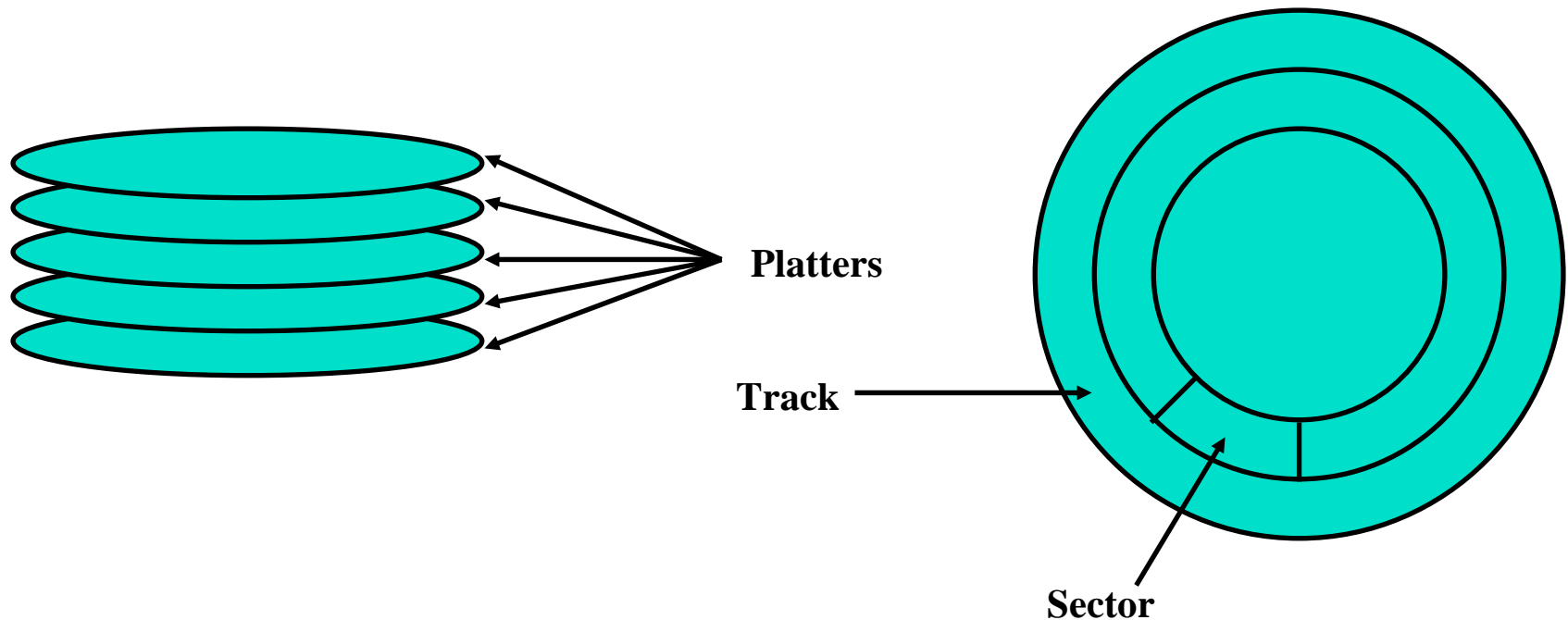
- Supercomputer application:
 - Large-scale scientific problems => large files
 - One large read and many small writes to snapshot computation
 - **Data Rate**: MB/second between memory and disk
- Transaction processing:
 - Examples: Airline reservations systems and bank ATMs
 - Small changes to large shared software
 - **I/O Rate**: No. disk accesses / second given upper limit for latency
- File system:
 - Measurements of UNIX file systems in an engineering environment:
 - 80% of accesses are to files less than 10 KB
 - 90% of all file accesses are to data with sequential addresses on the disk
 - 67% of the accesses are reads, 27% writes, 6% read-write
 - **I/O Rate & Latency**: No. disk accesses /second and response time

Magnetic Disk

- Purpose:
 - Long term, nonvolatile storage
 - Large, inexpensive, and slow
 - Lowest level in the memory hierarchy
- Two major types:
 - Floppy disk
 - Hard disk
- Both types of disks:
 - Rely on a rotating platter coated with a magnetic surface
 - Use a moveable read/write head to access the disk
- Advantages of hard disks over floppy disks:
 - Platters are more rigid (metal or glass) so they can be larger
 - Higher density because it can be controlled more precisely
 - Higher data rate because it spins faster
 - Can incorporate more than one platter



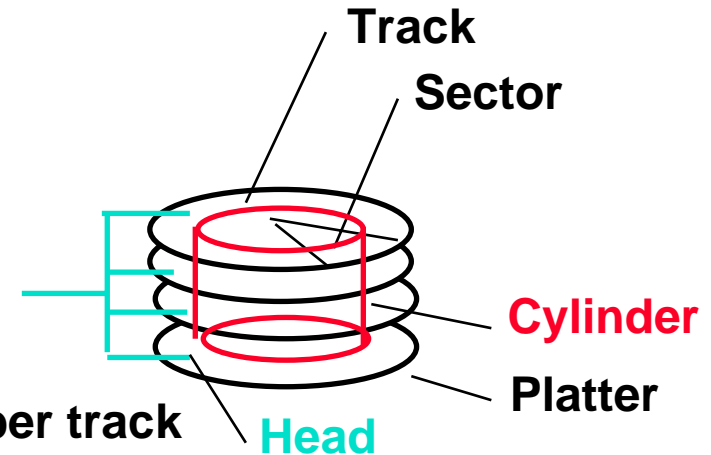
Organization of a Hard Magnetic Disk



- **Typical numbers (depending on the disk size):**
 - **500 to 2,000 tracks per surface**
 - **32 to 128 sectors per track**
 - **A sector is the smallest unit that can be read or written**
- **Traditionally all tracks have the same number of sectors:**
 - **Constant bit density: record more sectors on the outer tracks**
 - **Recently relaxed: constant bit size, speed varies with track location**

Magnetic Disk Characteristic

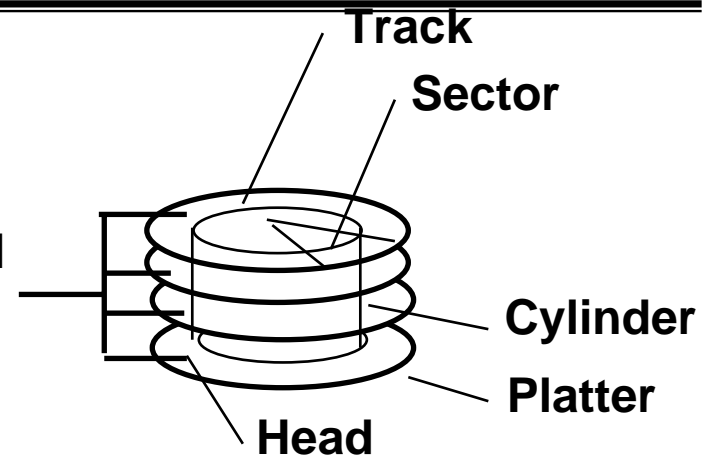
- **Cylinder:** all the tracks under the head at a given point on all surface
- **Read/write data is a three-stage process:**
 - **Seek time:** position the arm over the proper track
 - **Rotational latency:** wait for the desired sector to rotate under the read/write head
 - **Transfer time:** transfer a block of bits (sector) under the read-write head
- **Average seek time as reported by the industry:**
 - Typically in the range of 8 ms to 12 ms
 - $(\text{Sum of the time for all possible seek}) / (\text{total \# of possible seeks})$
- **Due to locality of disk reference, actual average seek time may:**
 - Only be 25% to 33% of the advertised number



Typical Numbers of a Magnetic Disk

◦ Rotational Latency:

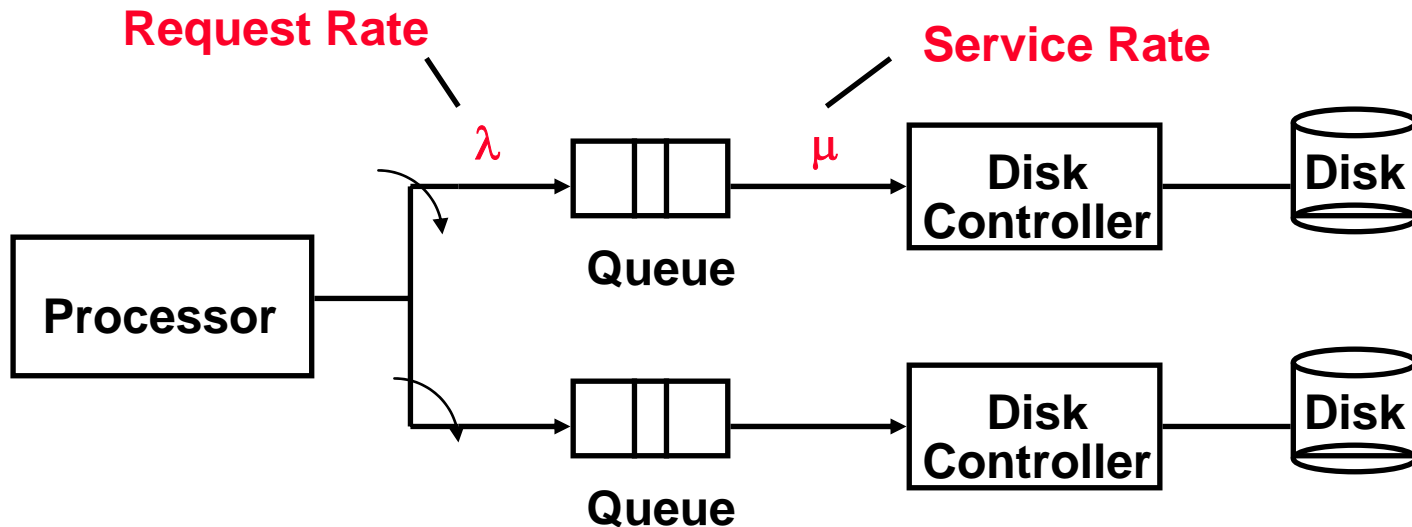
- Most disks rotate at 3,600 to 7200 RPM
- Approximately 16 ms to 8 ms per revolution, respectively
- An average latency to the desired information is halfway around the disk: 8 ms at 3600 RPM, 4 ms at 7200 RPM



◦ Transfer Time is a function of :

- Transfer size (usually a sector): 1 KB / sector
- Rotation speed: 3600 RPM to 7200 RPM
- Recording density: bits per inch on a track
- Diameter typical diameter ranges from 2.5 to 5.25 in
- Typical values: 2 to 12 MB per second

Disk I/O Performance



- **Disk Access Time = Seek time + Rotational Latency + Transfer time + Controller Time + Queueing Delay**
- **Estimating Queue Length:**
 - **Utilization = $U = \text{Request Rate} / \text{Service Rate}$**
 - **Mean Queue Length = $U / (1 - U)$**
 - **As Request Rate \rightarrow Service Rate**
 - **Mean Queue Length \rightarrow Infinity**

Example

- 512 byte sector, rotate at 5400 RPM, advertised seeks is 12 ms, transfer rate is 4 MB/sec, controller overhead is 1 ms, queue idle so no service time
- Disk Access Time = Seek time + Rotational Latency + Transfer time + Controller Time + Queueing Delay
- Disk Access Time = 12 ms + 0.5 / 5400 RPM + 0.5 KB / 4 MB/s + 1 ms + 0
- Disk Access Time = 12 ms + 0.5 / 90 RPS + 0.125 / 1024 s + 1 ms + 0
- Disk Access Time = 12 ms + 5.5 ms + 0.1 ms + 1 ms + 0 ms
- Disk Access Time = 18.6 ms
- If real seeks are 1/3 advertised seeks, then its 10.6 ms, with rotation delay at 50% of the time!

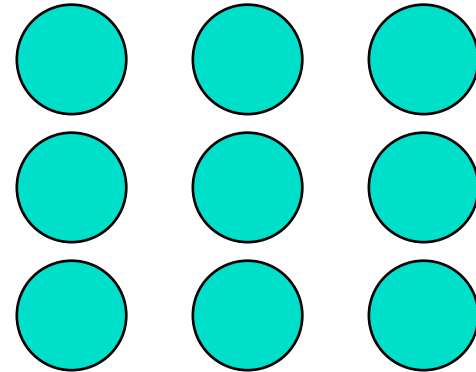
Magnetic Disk Examples

Characteristics	IBM 3090	IBM UltraStar	Integral 1820
Disk diameter (inches)	10.88	3.50	1.80
Formatted data capacity (MB)	22,700	4,300	21
MTTF (hours)	50,000	1,000,000	100,000
Number of arms/box	12	1	1
Rotation speed (RPM)	3,600	7,200	3,800
Transfer rate (MB/sec)	4.2	9-12	1.9
Power/box (watts)	2,900	13	2
MB/watt	8	102	10.5
Volume (cubic feet)	97	0.13	0.02
MB/cubic feet	234	33000	1050

Reliability and Availability

- **Two terms that are often confused:**
 - **Reliability: Is anything broken?**
 - **Availability: Is the system still available to the user?**
- **Availability can be improved by adding hardware:**
 - **Example: adding ECC on memory**
- **Reliability can only be improved by:**
 - **Bettering environmental conditions**
 - **Building more reliable components**
 - **Building with fewer components**
 - **Improve availability may come at the cost of lower reliability**

Disk Arrays



- **A new organization of disk storage:**
 - **Arrays of small and inexpensive disks**
 - **Increase potential throughput by having many disk drives:**
 - **Data is spread over multiple disk**
 - **Multiple accesses are made to several disks**
- **Reliability is lower than a single disk:**
 - **But availability can be improved by adding redundant disks (RAID):
Lost information can be reconstructed from redundant information**
 - **MTTR: mean time to repair is in the order of hours**
 - **MTTF: mean time to failure of disks is tens of years**

Optical Compact Disks

- **Disadvantage:**
 - It is primarily read-only media
- **Advantages of Optical Compact Disk:**
 - It is removable
 - It is inexpensive to manufacture
 - Have the potential to compete with new tape technologies for archival storage

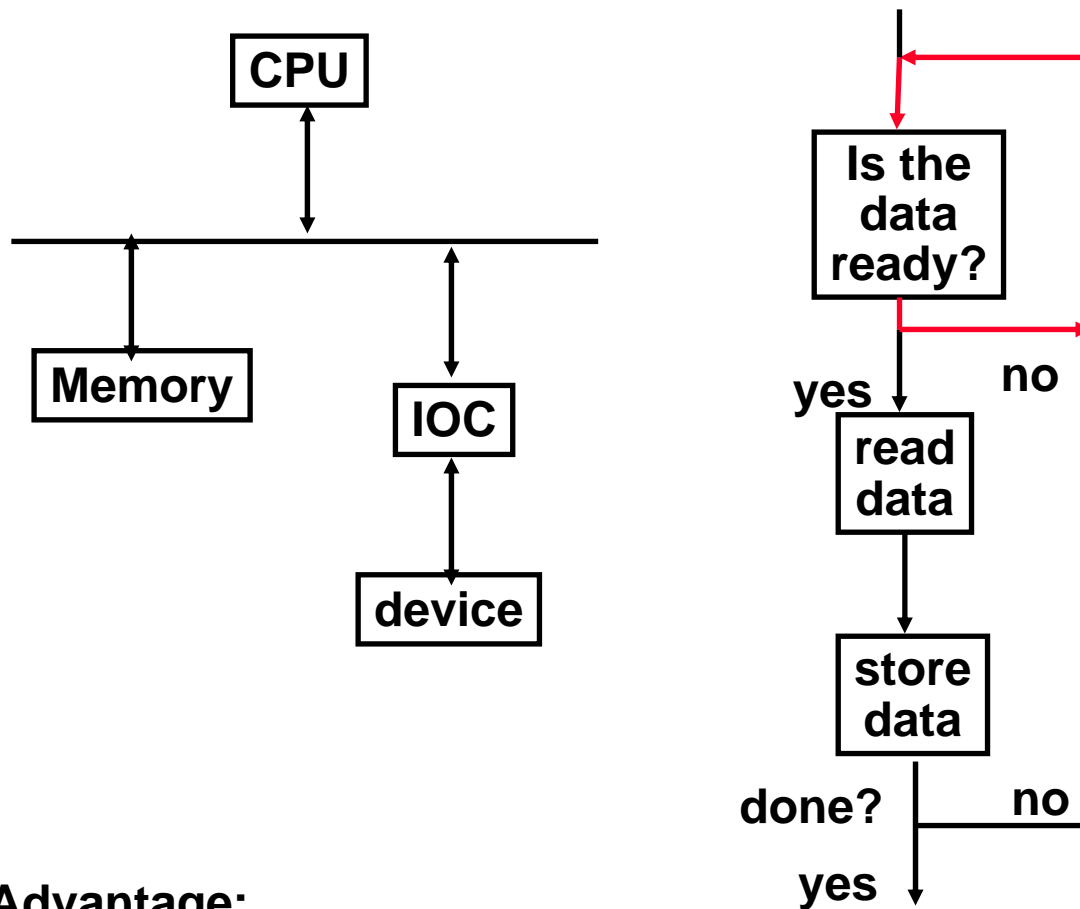
Giving Commands to I/O Devices

- **Two methods are used to address the device:**
 - **Special I/O instructions**
 - **Memory-mapped I/O**
- **Special I/O instructions specify:**
 - **Both the device number and the command word**
 - **Device number: the processor communicates this via a set of wires normally included as part of the I/O bus**
 - **Command word: this is usually send on the bus's data lines**
- **Memory-mapped I/O:**
 - **Portions of the address space are assigned to I/O device**
 - **Read and writes to those addresses are interpreted as commands to the I/O devices**
 - **User programs are prevented from issuing I/O operations directly:**
 - **The I/O address space is protected by the address translation**

I/O Device Notifying the OS

- **The OS needs to know when:**
 - **The I/O device has completed an operation**
 - **The I/O operation has encountered an error**
- **This can be accomplished in two different ways:**
 - **Polling:**
 - **The I/O device put information in a status register**
 - **The OS periodically check the status register**
 - **I/O Interrupt:**
 - **Whenever an I/O device needs attention from the processor, it interrupts the processor from what it is currently doing.**

Polling: Programmed I/O



**busy wait loop
not an efficient
way to use the CPU
unless the device
is very fast!**

**but checks for I/O
completion can be
dispersed among
computation
intensive code**

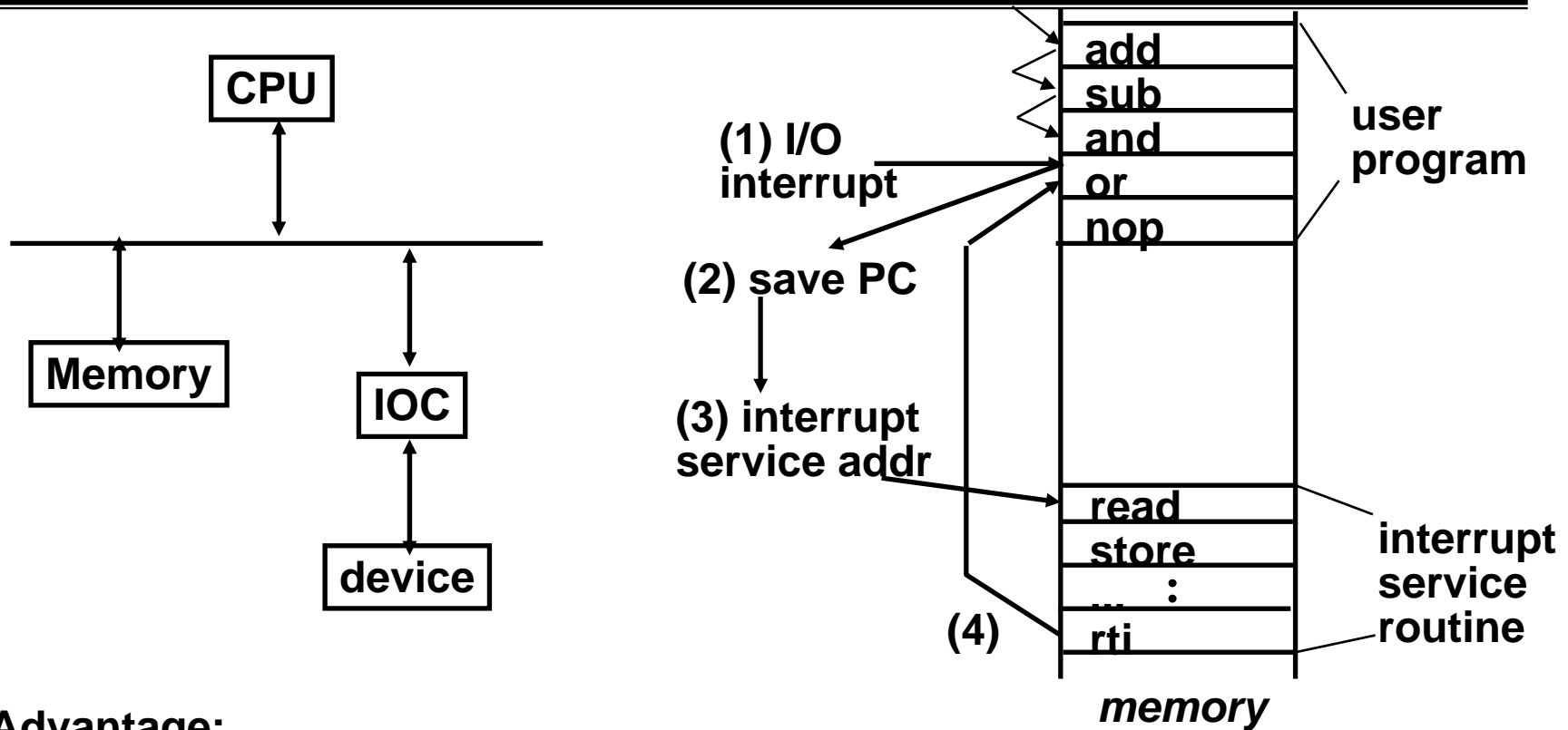
◦ **Advantage:**

- **Simple:** the processor is totally in control and does all the work

◦ **Disadvantage:**

- **Polling overhead can consume a lot of CPU time**

Interrupt Driven Data Transfer



◦ Advantage:

- User program progress is only halted during actual transfer

◦ Disadvantage, special hardware is needed to:

- Cause an interrupt (I/O device)
- Detect an interrupt (processor)
- Save the proper states to resume after the interrupt (processor)

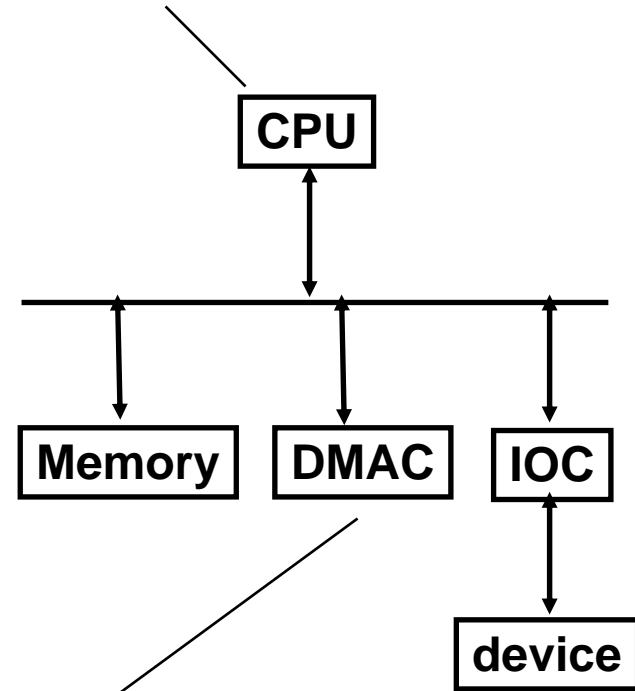
I/O Interrupt

- **An I/O interrupt is just like the exceptions except:**
 - **An I/O interrupt is asynchronous**
 - **Further information needs to be conveyed**
- **An I/O interrupt is asynchronous with respect to instruction execution:**
 - **I/O interrupt is not associated with any instruction**
 - **I/O interrupt does not prevent any instruction from completion**
 - **You can pick your own convenient point to take an interrupt**
- **I/O interrupt is more complicated than exception:**
 - **Needs to convey the identity of the device generating the interrupt**
 - **Interrupt requests can have different urgencies:**
 - **Interrupt request needs to be prioritized**

Delegating I/O Responsibility from the CPU: DMA

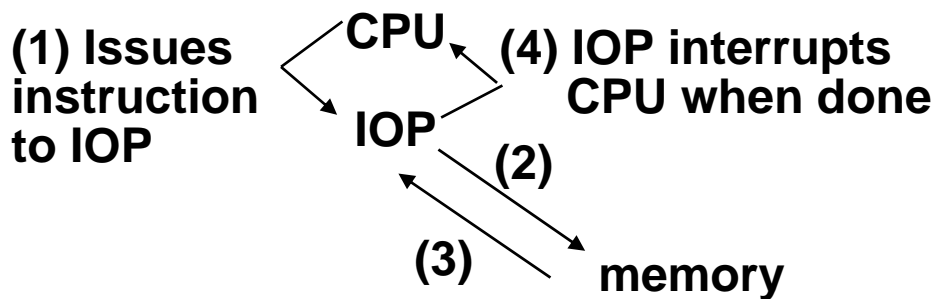
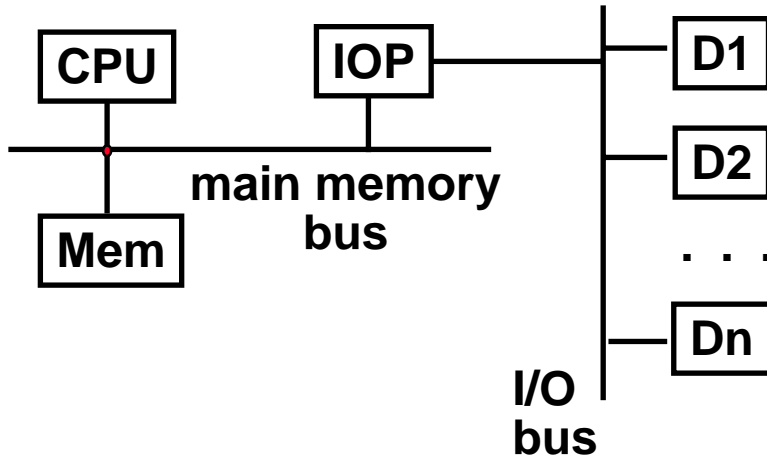
- **Direct Memory Access (DMA):**
 - External to the CPU
 - Act as a maser on the bus
 - Transfer blocks of data to or from memory without CPU intervention

CPU sends a starting address, direction, and length count to DMAC. Then issues "start".



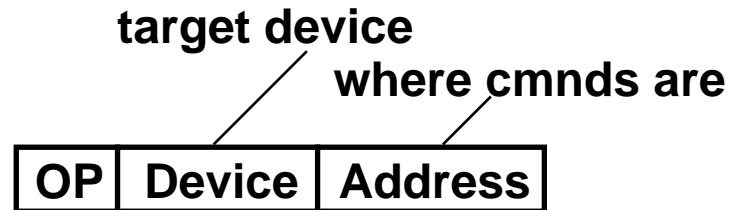
DMAC provides handshake signals for Peripheral Controller, and Memory Addresses and handshake signals for Memory.

Delegating I/O Responsibility from the CPU: IOP

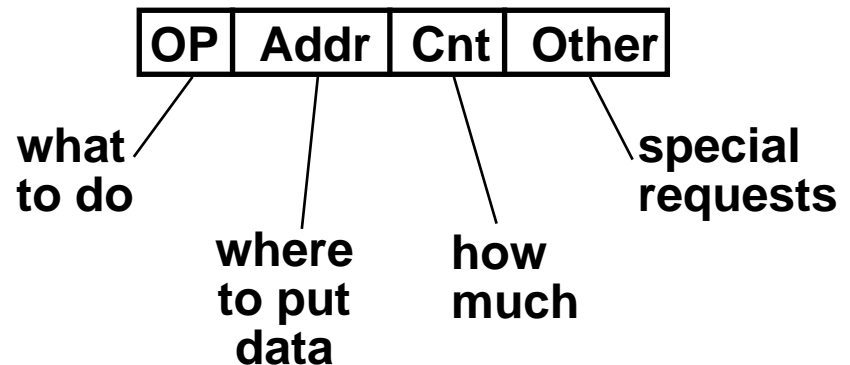


Device to/from memory transfers are controlled by the IOP directly.

IOP steals memory cycles.



IOP looks in memory for commands



Responsibilities of the Operating System

- **The operating system acts as the interface between:**
 - **The I/O hardware and the program that requests I/O**
- **Three characteristics of the I/O systems:**
 - **The I/O system is shared by multiple program using the processor**
 - **I/O systems often use interrupts (external generated exceptions) to communicate information about I/O operations.**
 - **Interrupts must be handled by the OS because they cause a transfer to supervisor mode**
 - **The low-level control of an I/O device is complex:**
 - **Managing a set of concurrent events**
 - **The requirements for correct device control are very detailed**

Operating System Requirements

- **Provide protection to shared I/O resources**
 - **Guarantees that a user's program can only access the portions of an I/O device to which the user has rights**
- **Provides abstraction for accessing devices:**
 - **Supply routines that handle low-level device operation**
- **Handles the interrupts generated by I/O devices**
- **Provide equitable access to the shared I/O resources**
 - **All user programs must have equal access to the I/O resources**
- **Schedule accesses in order to enhance system throughput**

OS and I/O Systems Communication Requirements

- **The Operating System must be able to prevent:**
 - **The user program from communicating with the I/O device directly**
- **If user programs could perform I/O directly:**
 - **Protection to the shared I/O resources could not be provided**
- **Three types of communication are required:**
 - **The OS must be able to give commands to the I/O devices**
 - **The I/O device must be able to notify the OS when the I/O device has completed an operation or has encountered an error**
 - **Data must be transferred between memory and an I/O device**

Multimedia Bandwidth Requirements

- **High Quality Video**

- **Digital Data = (30 frames / second) (640 x 480 pels) (24-bit color / pel) = 221 Mbps (75 MB/s)**

- **Reduced Quality Video**

- **Digital Data = (15 frames / second) (320 x 240 pels) (16-bit color / pel) = 18 Mbps (2.2 MB/s)**

- **High Quality Audio**

- **Digital Data = (44,100 audio samples / sec) (16-bit audio samples)**
- **(2 audio channels for stereo) = 1.4 Mbps**

- **Reduced Quality Audio**

- **Digital Data = (11,050 audio samples / sec) (8-bit audio samples) (1 audio channel for monaural) = 0.1 Mbps**

- **compression changes the whole story!**

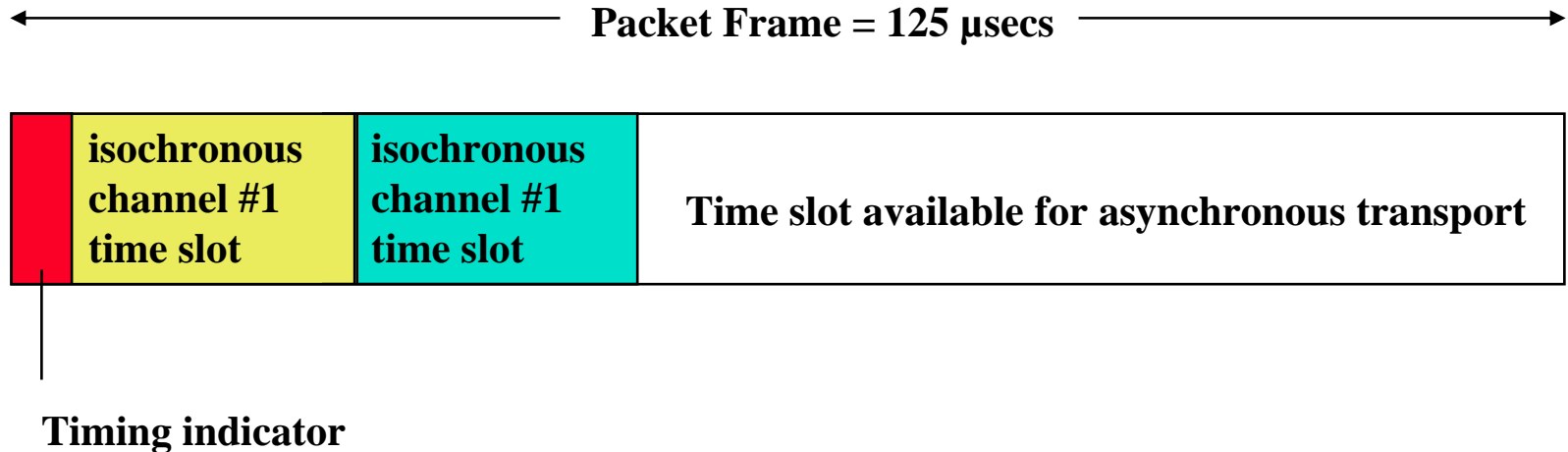
Multimedia and Latency

- **How sensitive is your eye / ear to variations in audio / video rate?**
- **How can you ensure constant rate of delivery?**
- **Jitter (latency) bounds vs constant bit rate transfer**
- **Synchronizing audio and video streams**
 - **you can tolerate 15-20 ms early to 30-40 ms late**

P1394 High-Speed Serial Bus (firewire)

- a digital interface – there is no need to convert digital data into analog and tolerate a loss of data integrity,
- physically small - the thin serial cable can replace larger and more expensive interfaces,
- easy to use - no need for terminators, device IDs, or elaborate setup,
- hot pluggable - users can add or remove 1394 devices with the bus active,
- inexpensive - priced for consumer products,
- scalable architecture - may mix 100, 200, and 400 Mbps devices on a bus,
- flexible topology - support of daisy chaining and branching for true peer-to-peer communication,
- fast - even multimedia data can be guaranteed its bandwidth for just-in-time delivery, and
- non-proprietary
- mixed asynchronous and isochronous traffic

Firewire Operations



- Fixed frame is divided into preallocated CBR slots + best effort asynchronous slot
- Each slot has packet containing “ID” command and data
- Example: digital video camera can expect to send one 64 byte packet every 125 μs
 - $80 * 1024 * 64 = 5\text{MB/s}$

Summary:

- **I/O performance is limited by weakest link in chain between OS and device**
- **Disk I/O Benchmarks: I/O rate vs. Data rate vs. latency**
- **Three Components of Disk Access Time:**
 - **Seek Time: advertised to be 8 to 12 ms. May be lower in real life.**
 - **Rotational Latency: 4.1 ms at 7200 RPM and 8.3 ms at 3600 RPM**
 - **Transfer Time: 2 to 12 MB per second**
- **I/O device notifying the operating system:**
 - **Polling: it can waste a lot of processor time**
 - **I/O interrupt: similar to exception except it is asynchronous**
- **Delegating I/O responsibility from the CPU: DMA, or even IOP**
- **wide range of devices**
 - **multimedia and high speed networking pose important challenges**