

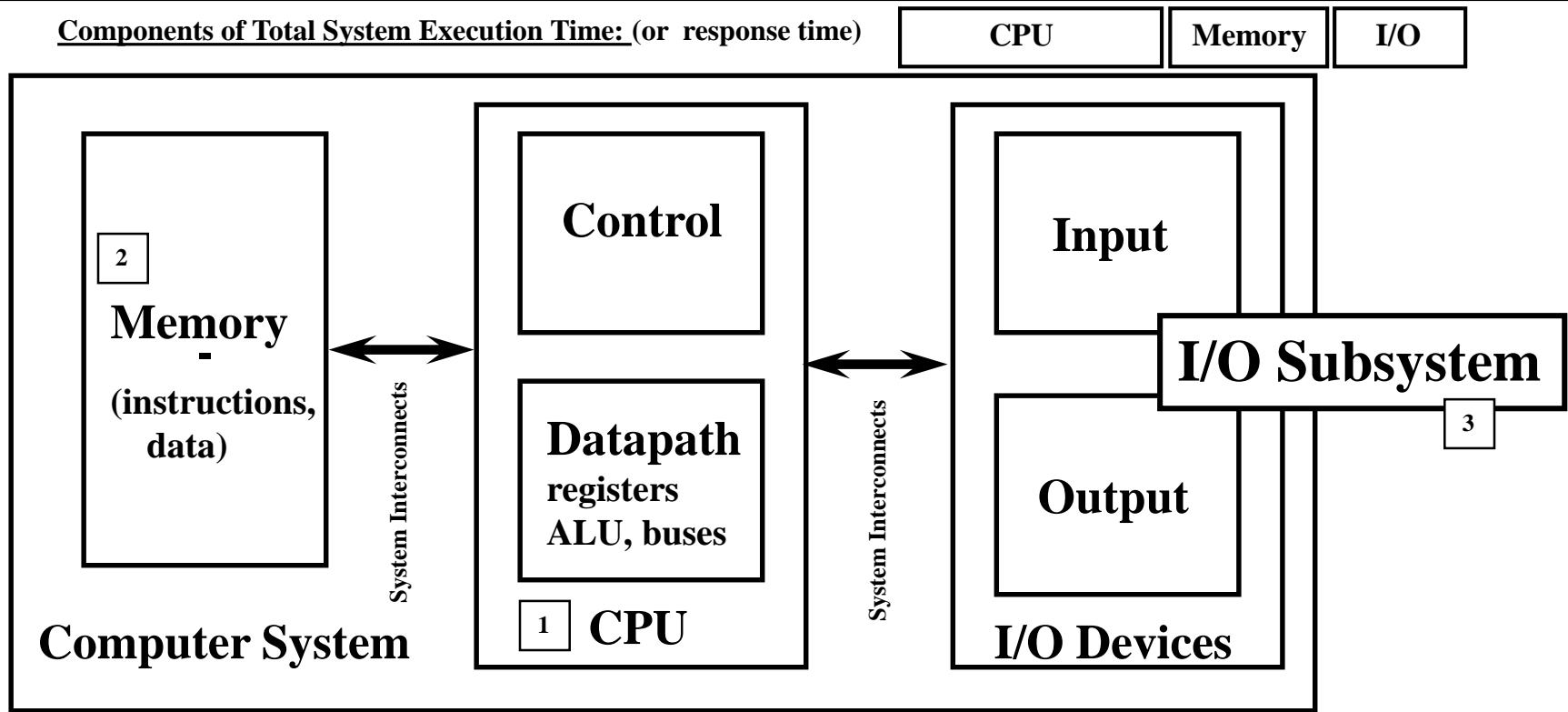
I/O Subsystem

The Von-Neumann Computer Model

- **Partitioning of the computing engine into components:**

- 1 – **Central Processing Unit (CPU):** Control Unit (instruction decode, sequencing of operations), Datapath (registers, arithmetic and logic unit, buses).
- 2 – **Memory:** Instruction (program) and operand (data) storage.
- 3 – **Input/Output (I/O):** Communication between the CPU/memory and the outside world.

System Architecture = System components and how the components are connected (system interconnects)



System performance depends on many aspects of the system
("limited by weakest link in the chain"): The system performance bottleneck

Input and Output (I/O) Subsystem

- The I/O subsystem provides the mechanism for communication between the CPU and the outside world (I/O devices).

Including users

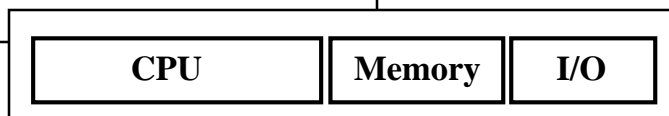
Including memory

- **Design factors:**

- I/O device characteristics (input, output, storage, etc.) /Performance.
- I/O Connection Structure (degree of separation from memory operations). →

Isolated I/O System Architecture
- I/O interface (the utilization of dedicated I/O and bus controllers).
- Types of buses/system interconnects (processor-memory vs. I/O buses/interconnects).
- I/O data transfer or synchronization method (programmed I/O, interrupt-driven, DMA).

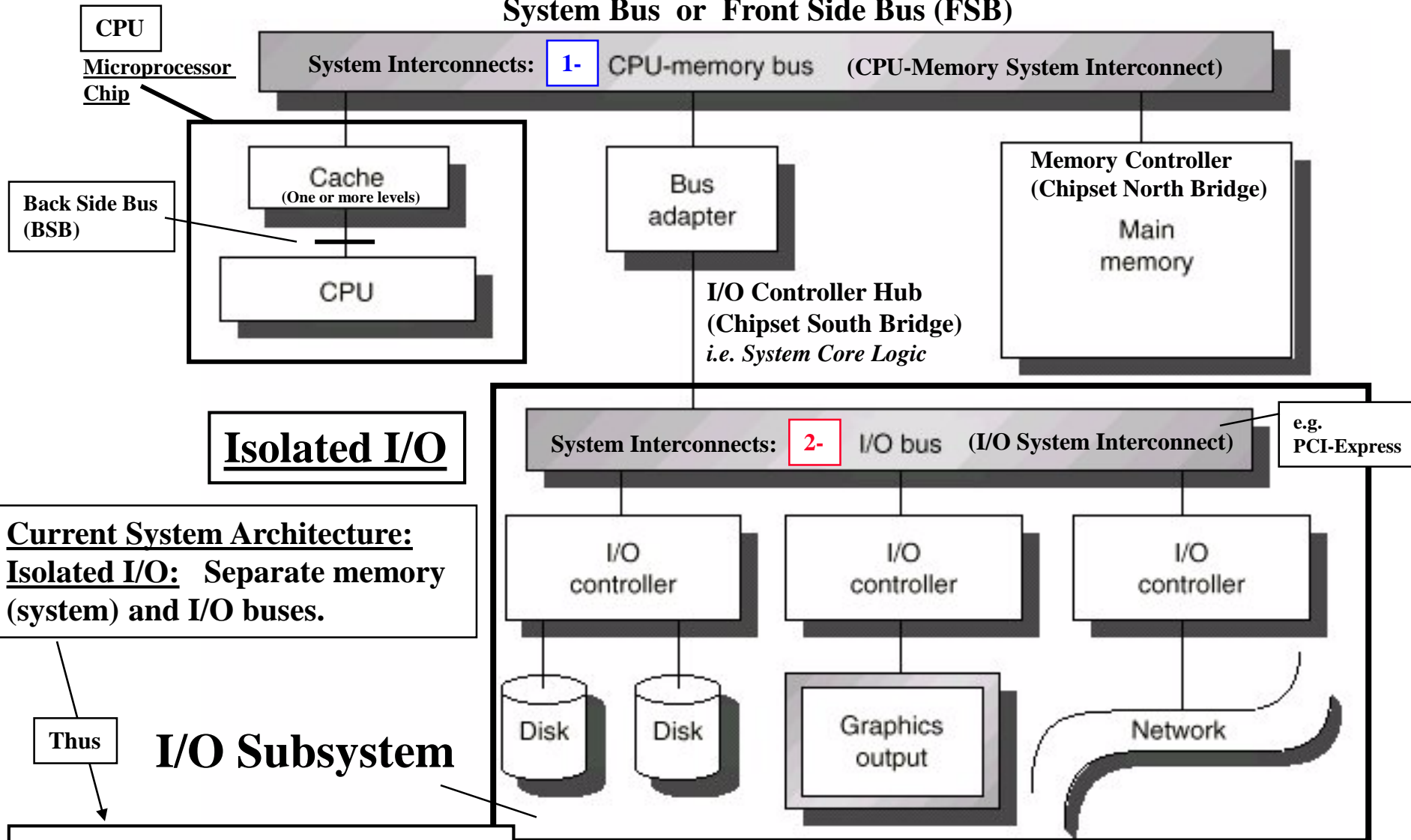
Components of Total System Execution Time:
(or response time)



Typical FSB-Based System Architecture

System Architecture = System Components + System Component Interconnects

System Bus or Front Side Bus (FSB)

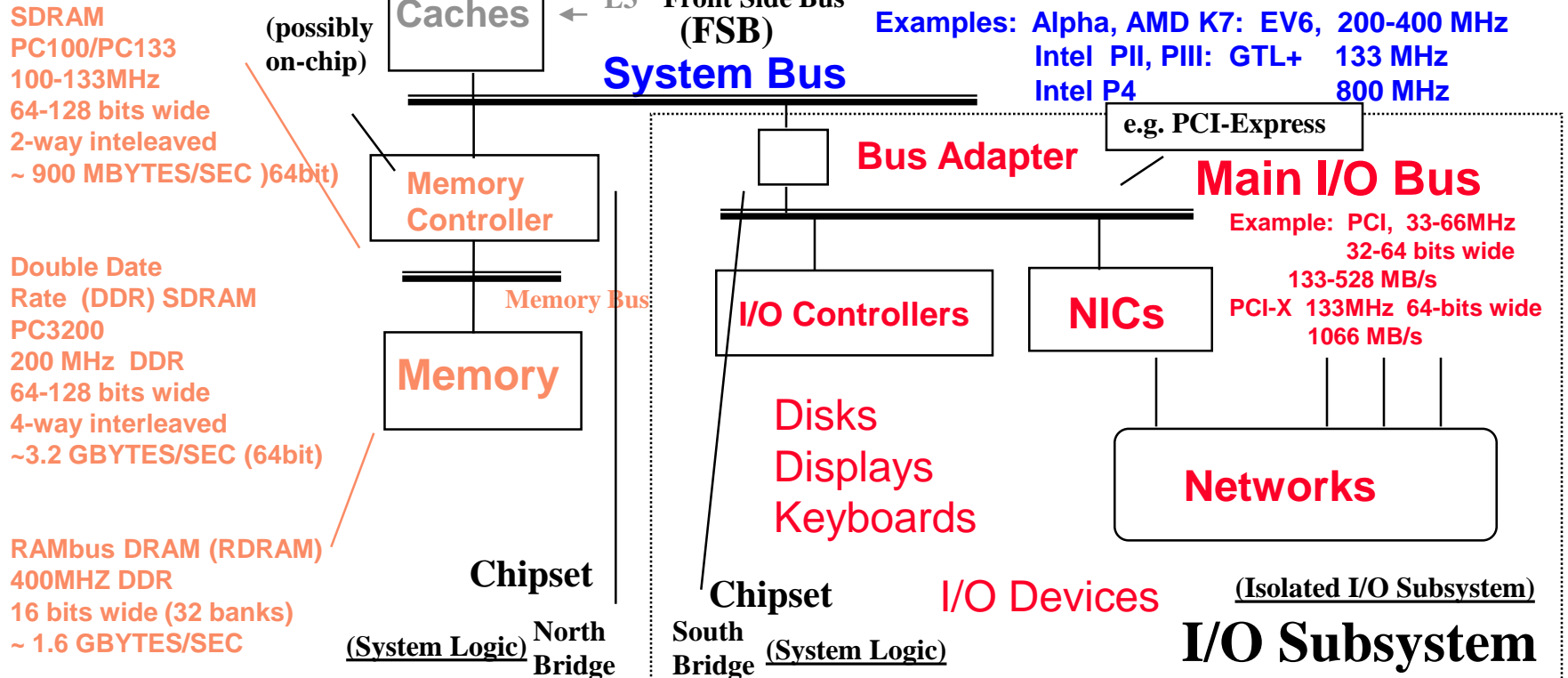


Two Types of System Interconnects/Buses:
1- CPU-Memory Bus or interconnect
2 – I/O Buses/interfaces

- CPU Core**
- 1 GHz - 3.8 GHz**
- 4-way Superscaler**
- RISC or RISC-core (x86):**
 - Deep Instruction Pipelines**
 - Dynamic scheduling**
 - Multiple FP, integer FUs**
 - Dynamic branch prediction**
 - Hardware speculation**

All Non-blocking caches

L1	16-128K	1-2 way set associative (on chip), separate or unified
L2	256K- 2M	4-32 way set associative (on chip) unified
L3	2-16M	8-32 way set associative (on or off chip) unified



Current System Architecture:
Isolated I/O: Separate memory (system) and I/O buses.

Thus

Two Types of System Interconnects/Buses:

- 1- CPU-Memory Bus or interconnect**
- 2 – I/O Buses/interfaces**

Important issue: Which component creates a system performance bottleneck?

Main Types of Buses/Interconnects in The System

1 Processor-Memory Bus/Interconnect: AKA System Bus, Front Side Bus, (FSB)

- Should offer very high-speed (bandwidth) and low latency.
- Matched to the memory system performance to maximize memory-processor bandwidth.
- Usually system design-specific (not an industry standard).
- Examples: Alpha EV6 (AMD K7), Peak bandwidth = $400 \text{ MHz} \times 8 = 3.2 \text{ GB/s}$
Intel GTL+ (P3), Peak bandwidth = $133 \text{ MHz} \times 8 = 1 \text{ GB/s}$
Intel P4, Peak bandwidth = $800 \text{ MHz} \times 8 = 6.4 \text{ GB/s}$
HyperTransport 2.0: 200Mhz-1.4GHz, Peak bandwidth up to 22.8 GB/s

QPI

Also Intel's QuickPath Interconnect (QPI)
used in Core i7 system architecture

(point-to-point system interconnect not a bus)

Dedicated
Links

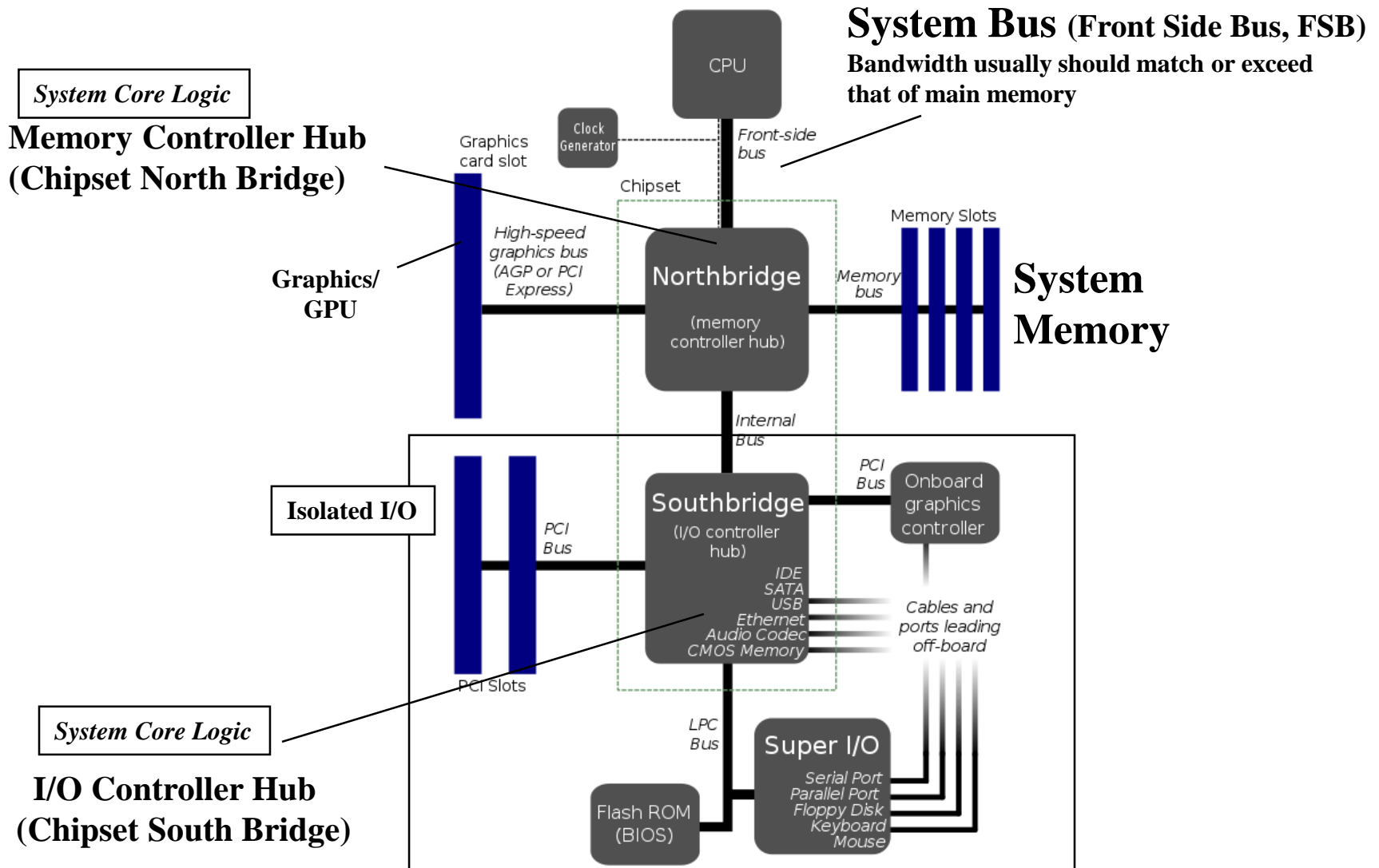
2 I/O buses/Interconnects: Sometimes called I/O *channels or interfaces*

- Follow bus/interface industry standards.
- Usually formed by I/O interface adapters to handle many types of connected I/O devices.
- Wide range in the data bandwidth and latency
- Not usually interfaced directly to memory instead connected to processor-memory bus via a bus adapter (system chipset south bridge).
- Examples: Main system I/O bus: PCI, PCI-X, PCI Express
Storage Interfaces: SATA, PATA, SCSI.

Isolated I/O System Architecture

System Architecture = System Components + System Component Interconnects

FSB-Based Single Processor Socket System Architecture



Intel Pentium 4 System Architecture

And Core 2

(Using The Intel 925 Chipset)

CPU

(Including cache)

Intel® Pentium® 4
Processor
Extreme Edition

System Bus (Front Side Bus, FSB)

Bandwidth usually should match or exceed
that of main memory

System Core Logic

Memory Controller Hub
(Chipset North Bridge)

Graphics/GPU

PCI Express®
x16 Graphics

8.0
GB/s

82925X
MCH

DDR2

8.5 GB/s

DDR2

System
Memory

Two 8-byte DDR2 Channels

Graphics I/O Bus (PCI Express)

Isolated I/O

2 GB/s

DMI

Intel® High
Definition Audio

4 PCI
Express® x1

500
MB/s

8 Hi-Speed
USB 2.0 Ports

60
MB/s

Misc.
I/O
Interfaces

ICH6RW

150
MB/s

133
MB/s

Storage I/O (Serial ATA)

4 Serial
ATA Ports

6
PCI

Main
I/O Bus
(PCI)

Intel® Matrix
Storage Technology

Intel® Wireless
Connect Technology

Misc.
I/O
Interfaces

BIOS Supports
HT Technology

Basic Input Output System (BIOS)

I/O Controller Hub
(Chipset South Bridge)

I/O Subsystem

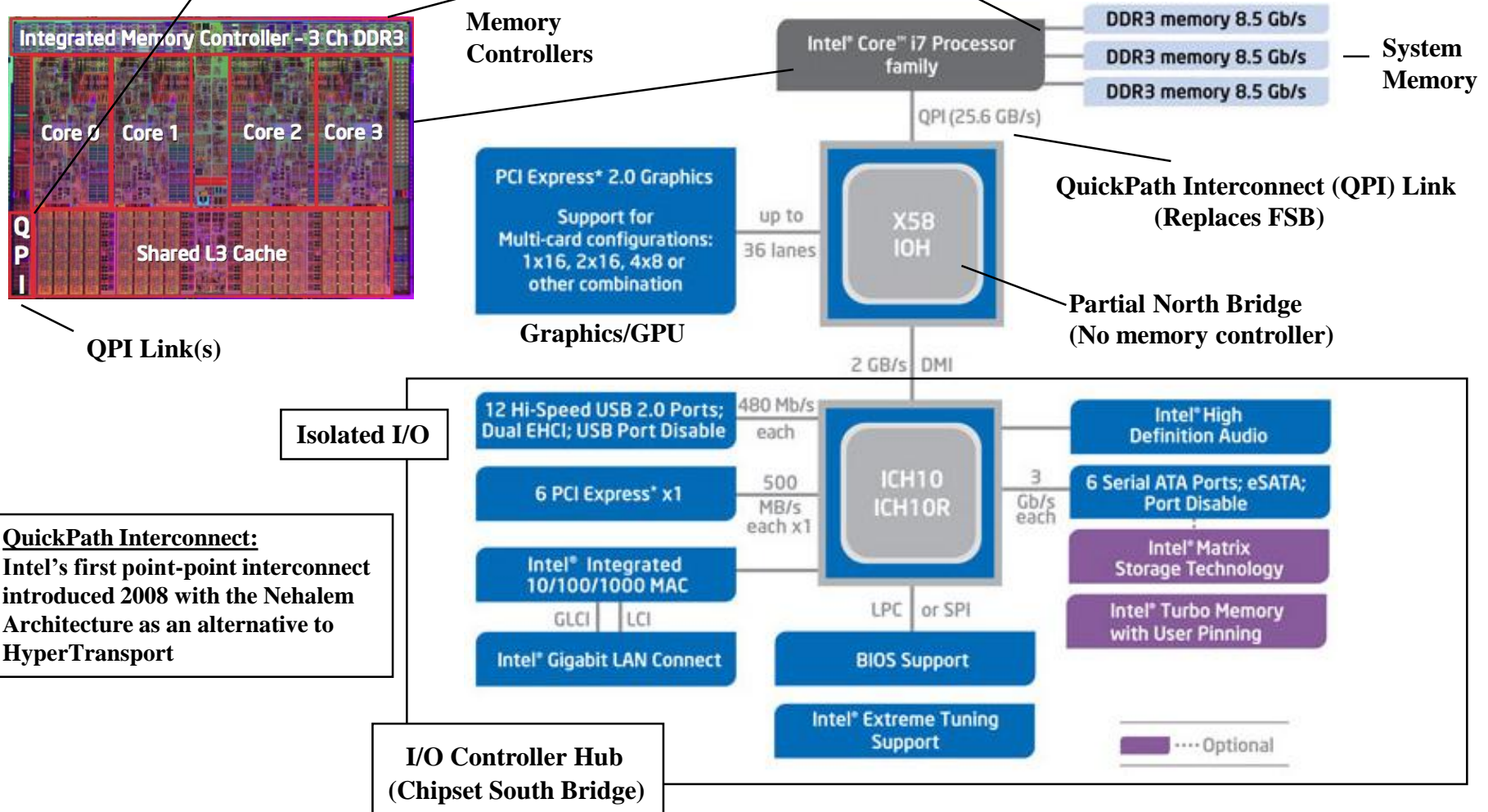
Current System Architecture:

Isolated I/O: Separate memory and I/O buses.

Intel Core i7 “Nehalem” System Architecture

Intel's QuickPath Interconnect (QPI) Point-to-point system interconnect used instead of Front Side Bus (FSB)

+ Memory controller integrated on processor chip (three DDR3 channels)



(e.g . FSB)

Bus Characteristics

<i>Option</i>	<i>High performance</i>	<i>Low cost/performance</i>
Bus width	Separate address & data lines	Multiplex address & data lines
Data width	Wider is faster (e.g., 64 bits)	Narrower is cheaper (e.g., 16 bits)
Transfer size	Multiple words has less bus overhead	Single-word transfer is simpler
Bus masters	Multiple (requires arbitration)	Single master (no arbitration)
Split	Yes, separate Request and Reply packets gets higher bandwidth (needs multiple masters)	No , continuous transaction? connection is cheaper and has lower latency
Clocking	Synchronous	Asynchronous

Example CPU-Memory System Buses (Front Side Buses, FSBs)

Bus	Summit	Challenge	XDBus	SP	P4
Originator	HP	SGI	Sun	IBM	Intel
Clock Rate (MHz)	60	48	66	111	800
Split transaction?	Yes	Yes	Yes	Yes	Yes
Address lines	48	40	??	??	??
Data lines	128	256	144	128	64
Clocks/transfer	4	5	4	??	??
Peak (MB/s)	960	1200	1056	1700	6400
Master	Multi	Multi	Multi	Multi	Multi
Arbitration	Central	Central	Central	Central	Central
Addressing	Physical	Physical	Physical	Physical	Physical
Length	13 inches	12 inches	17 inches	??	??

FSB Bandwidth matched with single 8-byte channel SDRAM

FSB Bandwidth matched with dual channel PC3200 DDR SDRAM

Main System I/O Bus Example: PCI, PCI-Express

	Specification	Bus Width (bits)	Bus Frequency (MHz)	Peak Bandwidth (MB/sec)
Legacy PCI	PCI 2.3	32	33.3	133
	PCI 2.3	64	33.3	266
	PCI 2.3	64	66.6	533
	PCI-X 1.0	64	133.3	1066
	PCI-X 2.0	64	266, 533	2100 , 4200
Not Implemented Yet				
Formerly Intel's 3GIO	PCI-Express	1-32	???	500-16,000

Addressing	Physical	PCI Bus Transaction Latency:
Master	Multi	PCI requires 9 cycles @ 33Mhz (272ns)
Arbitration	Central	PCI-X requires 10 cycles @ 133MHz (75ns)

PCI = Peripheral Component Interconnect

Storage IO Interfaces/Buses

EIDE/Parallel ATA (PATA)

Data Width	16 bits
Clock Rate	Upto 100MHz
Bus Masters	1
Max no. devices	2
Peak Bandwidth	200 MB/s
Target Application	Desktop

EIDE = Enhanced Integrated Drive Electronics
ATA = Advanced Technology Attachment
PATA = Parallel ATA
SATA = Serial ATA

SCSI

8 or 16 bits (wide)
10MHz (Fast)
20MHz (Ultra)
40MHz (Ultra2)
80MHz (Ultra3)
160MHz (Ultra4)
Multiple
7 (8-bit bus)
15 (16-bit bus)
320MB/s (Ultra4)
Servers_____

SCSI = Small Computer System Interface

I/O Data Transfer Methods

1

• Programmed I/O (PIO): Polling (For low-speed I/O)

- The I/O device puts its status information in a status register.
- The processor must periodically check the status register.
- The processor is totally in control and does all the work.
- Very wasteful of processor time.
- Used for low-speed I/O devices (mice, keyboards etc.)

Memory-mapped
register

2

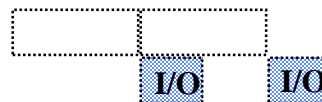
• Interrupt-Driven I/O (For medium-speed I/O):

- An interrupt line from the I/O device to the CPU is used to generate an I/O interrupt indicating that the I/O device needs CPU attention. (e.g data is ready)
- The interrupting device places its identity in an interrupt vector.
- Once an I/O interrupt is detected the current instruction is completed and an I/O interrupt handling routine (by OS) is executed to service the device.
- Used for moderate speed I/O (optical drives, storage, networks ..)
- Allows overlap of CPU processing time and I/O processing time

$$\text{Time(workload)} = \text{Time(CPU)} + \text{Time(I/O)} - \text{Time(Overlap)}$$



No overlap



Overlap of CPU processing
Time and I/O processing time

I/O data transfer methods:

3 Direct Memory Access (DMA) (For high-speed I/O):

- **Implemented with a specialized controller that transfers data between an I/O device and memory independent of the processor.**
- **The DMA controller becomes the bus master and directs reads and writes between itself and memory.**
- **Interrupts are still used only on completion of the transfer or when an error occurs.**
- **Even lower CPU overhead, used in high speed I/O (storage, network interfaces)**
- **Allows more overlap of CPU processing time and I/O processing time than interrupt-driven I/O.**
- **DMA transfer steps:**
 - 1 – The CPU sets up DMA by supplying device identity, operation, memory address of source and destination of data, the number of bytes to be transferred.**
 - 2 – The DMA controller starts the operation. When the data is available it transfers the data, including generating memory addresses for data to be transferred.**
 - 3 – Once the DMA transfer is complete, the controller interrupts the processor, which determines whether the entire operation is complete.**

I/O Interface/Controller

I/O Interface, I/O controller or I/O bus adapter:

- **Specific to each type of I/O device/interface standard.**
- **To the CPU, and I/O device, it consists of a set of control and data registers (usually memory-mapped) within the I/O address space.**
- **On the I/O device side, it forms a localized I/O bus which can be shared by several I/O devices**
 - (e.g IDE, SCSI, USB ...) Industry-standard interfaces

Why? – **Handles I/O details (originally done by CPU) such as:**

- Assembling bits into words,
- Low-level error detection and correction
- Accepting or providing words in word-sized I/O registers.
- Presents a uniform interface to the CPU regardless of I/O device.
- Handles DMA I/O data transfers.

**Low-level
I/O Processing
off-loaded
from CPU**

I/O Controller Architecture

Part of System Core Logic

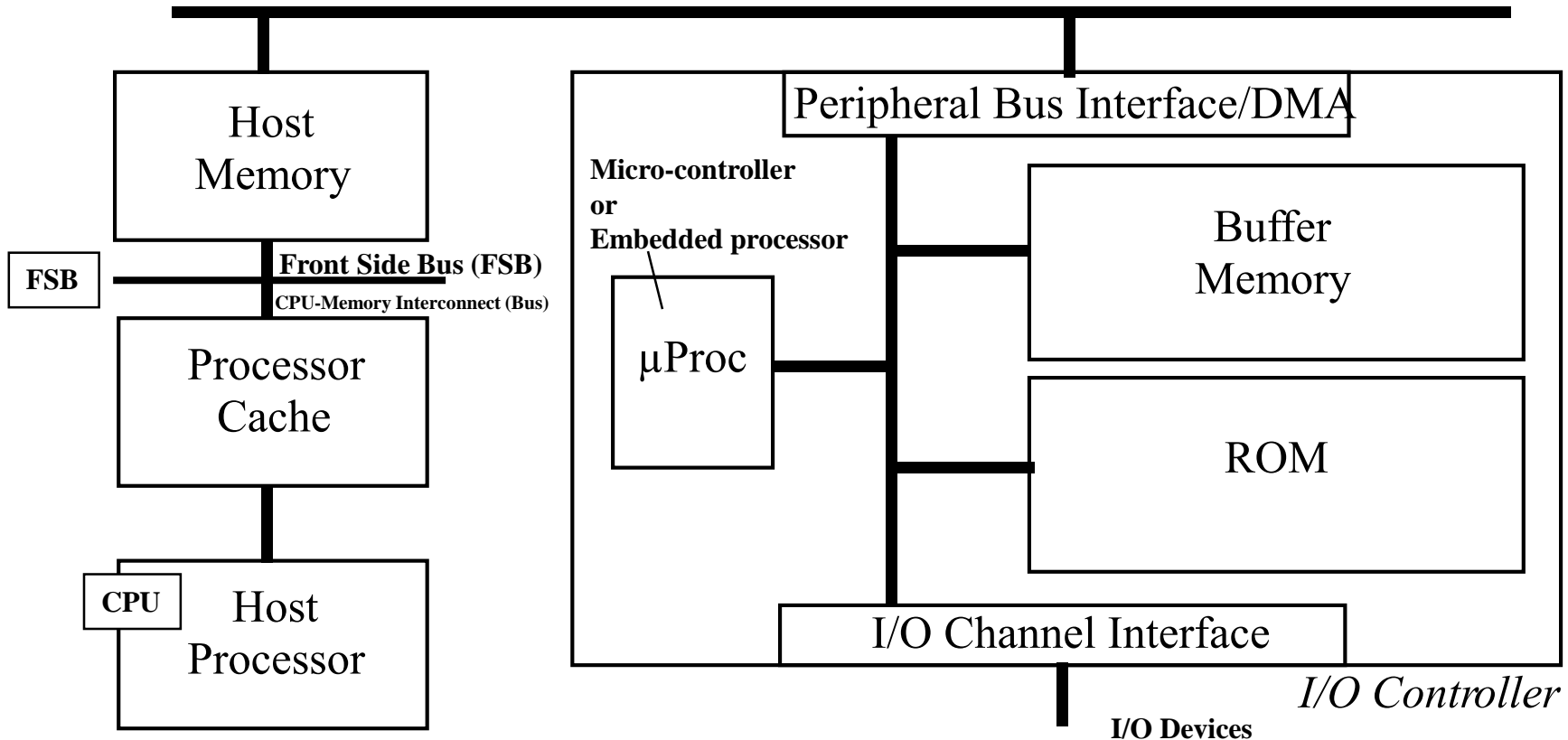
Chipset /
North Bridge

e.g. PCI-Express

Part of System Core Logic

Chipset /
South Bridge

Peripheral or Main I/O Bus (PCI, PCI-X, etc.)



I/O Controller

I/O Devices

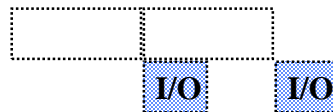
SCSI, IDE, USB,

Industry-standard interfaces

$$\text{Time(workload)} = \text{Time(CPU)} + \text{Time(I/O)} - \text{Time(Overlap)}$$



No overlap



Overlap of CPU processing
Time and I/O processing time

I/O: A System Performance Perspective

- CPU Performance: Improvement of ~ 60% per year.
- I/O Sub-System Performance: Limited by *mechanical* delays (disk I/O). Improvement less than 10% per year (IO rate per sec or MB per sec).

i.e storage devices (hard drives)

CPU

Memory

I/O

- From Amdahl's Law: overall system speed-up is limited by the slowest component:

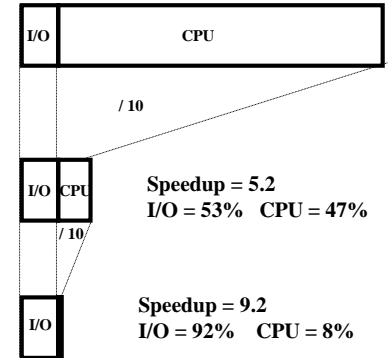
To Start With

→ If I/O is 10% of current processing time:

- Increasing CPU performance by 10 times
⇒ 5 times system performance increase (50% loss in performance)
- Increasing CPU performance by 100 times
⇒ ~ 10 times system performance (90% loss of performance)

Originally: CPU-bound

Originally: I/O = 10% CPU = 90%



After: I/O-bound

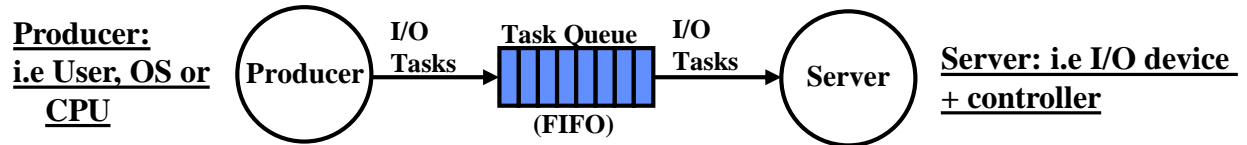
- The I/O system performance bottleneck diminishes the benefit of faster CPUs on overall system performance.

System performance depends on many aspects of the system
("limited by weakest link in the chain"): The system performance bottleneck

System & I/O Performance Metrics/Modeling

- **Diversity**: The variety of I/O devices that can be connected to the system.
- **Capacity**: The maximum number of I/O devices that can be connected to the system.

I/O Performance Modeling:



- **Producer/server Model of I/O**: The producer (CPU, human etc.) creates tasks to be performed and places them in a task buffer (queue); the server (I/O device or controller) takes tasks from the queue and performs them.

I/O (or Entire System) Performance Metrics:

- 1 • **I/O Throughput**: The maximum data rate that can be transferred to/from an I/O device or sub-system, or the maximum number of I/O tasks or transactions completed by I/O in a certain period of time
⇒ Maximized when task queue is never empty (server always busy).
- 2 • **I/O Latency or response time**: The time an I/O task takes from the time it is placed in the task buffer or queue until the server (I/O system) finishes the task. Includes I/O device service time and buffer waiting (or queuing time).
⇒ Minimized when task queue is always empty (no queuing time).

$$\text{Response Time} = \text{Service Time} + \text{Queuing Time}$$

System & I/O Performance Metrics: Throughput

- Throughput is a measure of speed—the rate at which the I/O or storage system delivers data.
- I/O Throughput is measured in two ways:

1

- I/O task rate:

- Measured in:

I/O Tasks/sec

- Accesses/second or,
 - Transactions Per Second (TPS) or,
 - I/O Operations Per Second (IOPS).

- I/O rate is generally used for applications where the size of each request is small, such as in transaction processing.

i.e server applications

2

- Data rate, measured in *bytes/second* or *megabytes/second* (*MB/s*, *GB/s* ...).

- Data rate is generally used for applications where the size of each request is large, such as in scientific and multimedia applications.

System & I/O Performance Metrics: Response time

- Response time measures how long a storage (or I/O) system takes to process an I/O request and access data. Or entire system

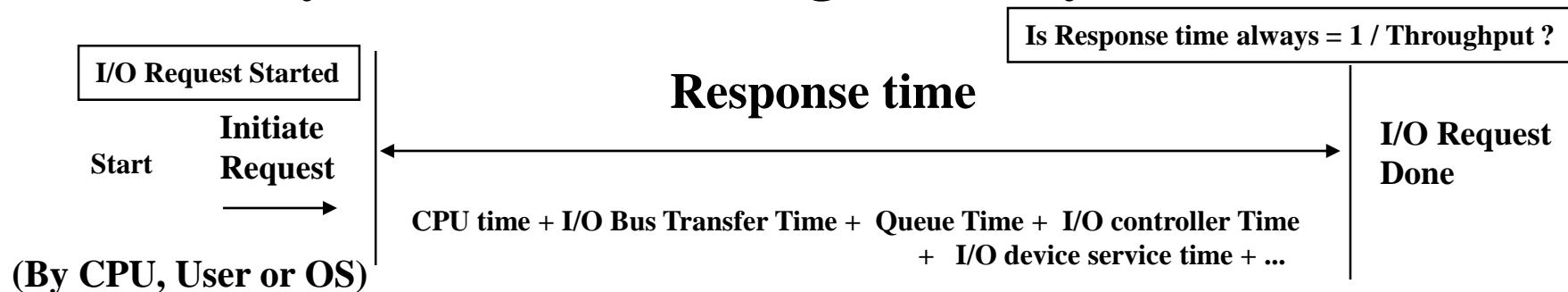
- I/O request latency or total processing time per I/O request.

- This time can be measured in several ways.

For example:

i.e. Time it takes the system to process an average task

- One could measure time from the user's perspective,
 - the operating system's perspective,
 - or the disk controller's perspective, depending on what you view as the storage or I/O system.

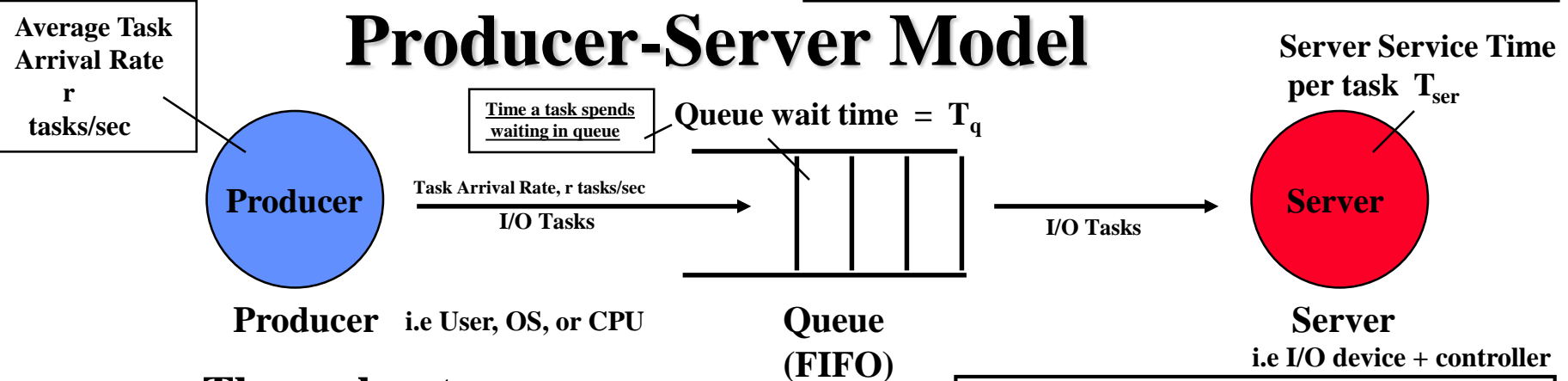


The utilization of DMA and I/O device queues and multiple I/O devices servicing a queue may make throughput $\gg 1 / \text{response time}$

I/O Modeling:

$$\text{Time}_{\text{system}} = \text{Time in System for a task} = \text{Response Time} = \text{Queuing Time} + \text{Service Time}$$

Producer-Server Model



- **Throughput:**

- The number of tasks completed by the server in unit time.
- In order to get the highest possible throughput:

Shown above: Single Queue + Single Server

Throughput
is maximized when:

- 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:

Response Time
is minimized when:

- The queue should be empty (no waiting time in queue).
- The server will be idle at times.

Shown above is a (Single Queue + Single Server) Producer-Server Model

Producer-Server Model

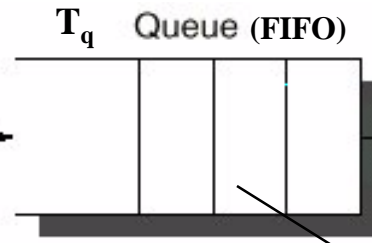
Single Queue +
Single Server

User or CPU



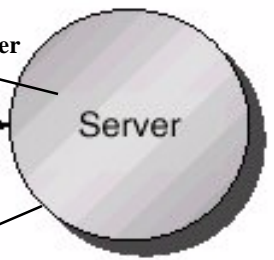
Task Arrival Rate, r

I/O
Tasks



T_{ser}

I/O
Tasks



I/O device +
controller

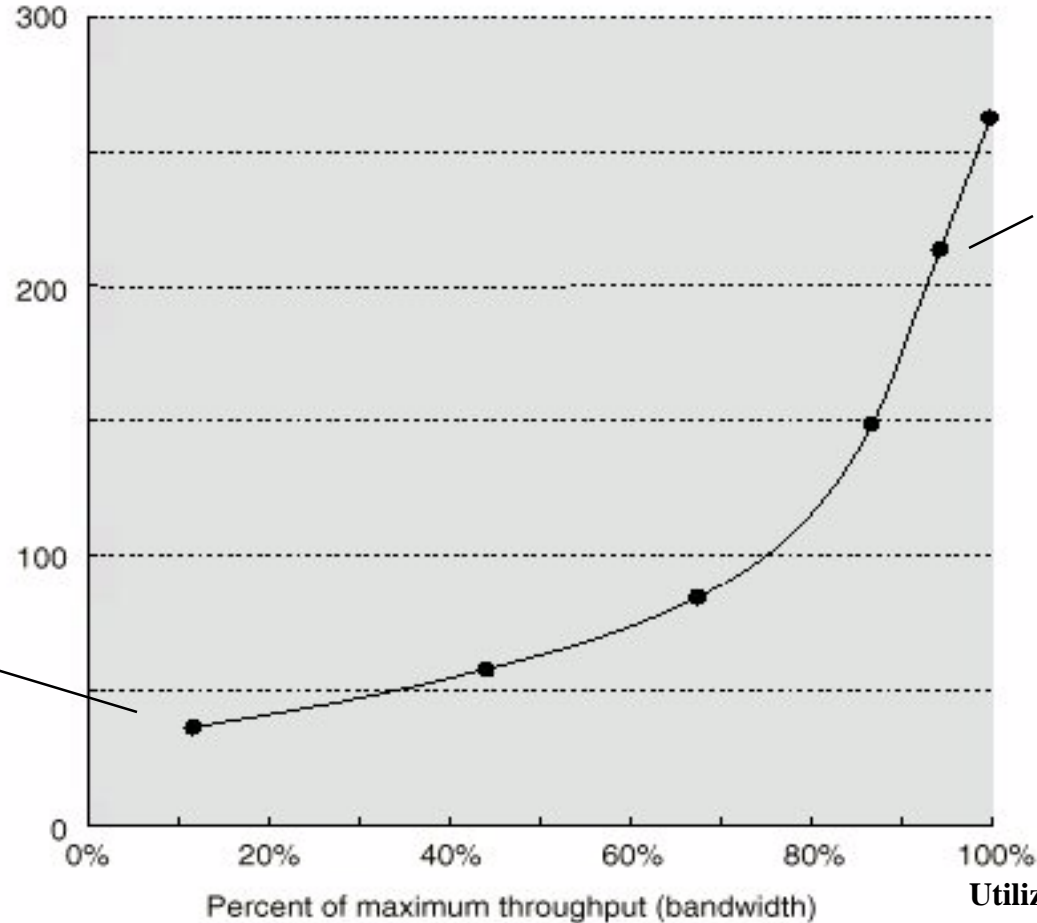
$$\text{Response Time} = \text{Time}_{\text{System}} = \text{Time}_{\text{Queue}} + \text{Time}_{\text{Server}} = T_q + T_{ser}$$

Throughput VS. Response Time

Response time
(latency)
in ms

Queue almost empty
most of the time
Less time in queue

Queue
full
most of the
time.
More time
in queue



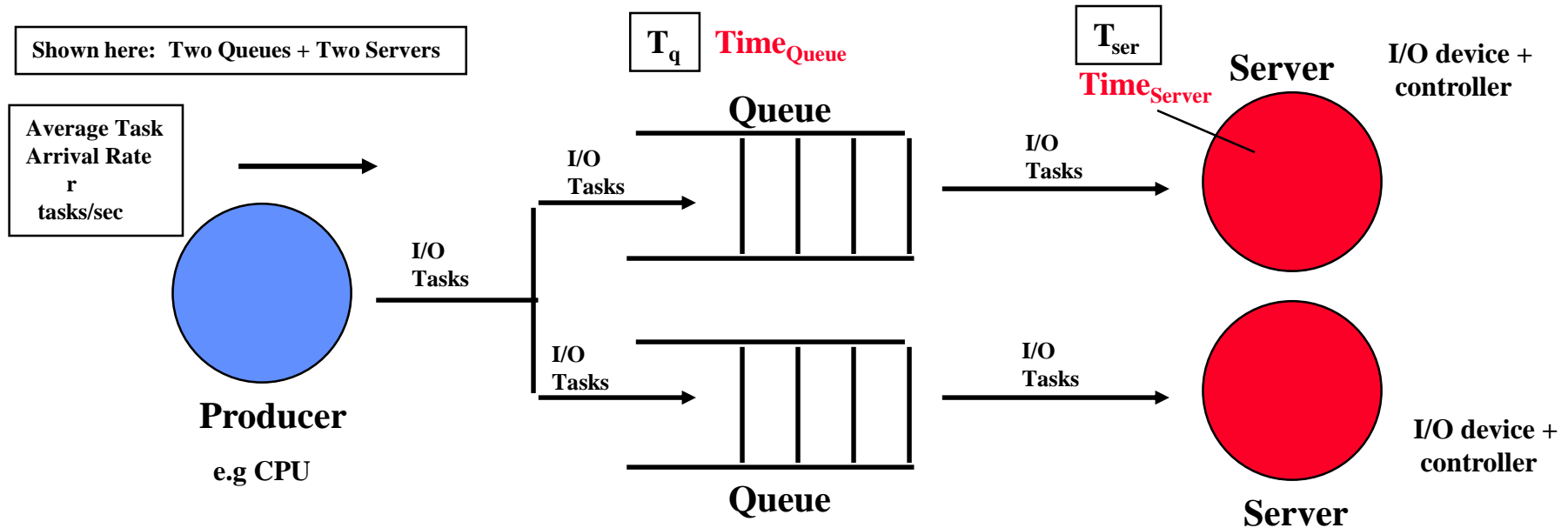
Shown here is a (Single Queue + Single Server)
Producer-Server Model

AKA Loading Factor

i.e Utilization = U ranges from 0 to 1 (0 % to 100%)

I/O Performance:

Throughput Enhancement



- In general throughput can be improved by:
 - Throwing more hardware at the problem.
 - Reduces load-related latency. Less queuing time
- Response time is much harder to reduce.
 - e.g. Faster I/O device (i.e server)

Ignoring CPU I/O processing time and other system delays

$$\text{Response Time} = \text{Time}_{\text{System}} = \text{Time}_{\text{Queue}} + \text{Time}_{\text{Server}} = T_q + T_{ser}$$

Magnetic Disks

Characteristics:

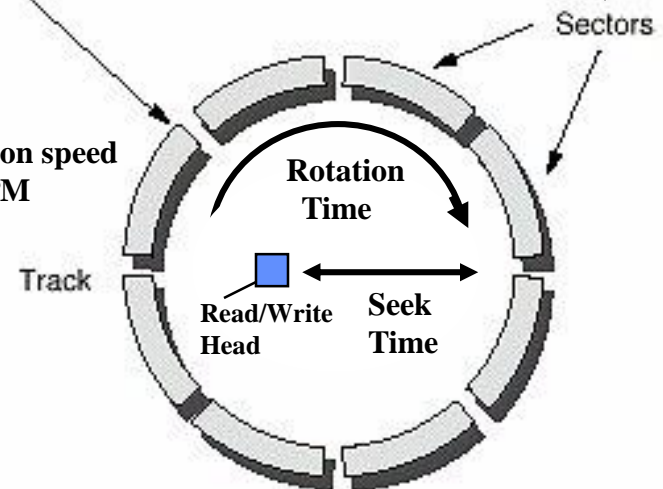
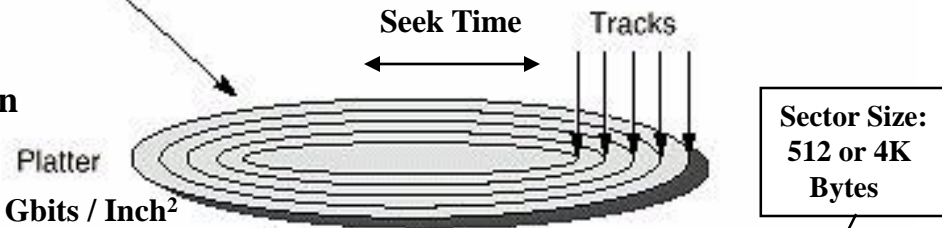
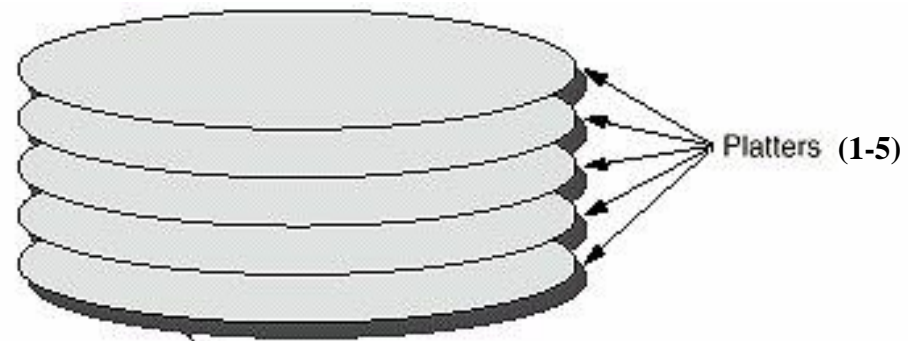
- **Diameter (form factor):** 1.8in - 3.5in
- **Rotational speed:** 5,400 RPM-15,000 RPM
- **Tracks per surface.**
- **Sectors per track:** Outer tracks contain more sectors.
- **Recording or Areal Density:** Tracks/in X Bits/in
- **Cost Per Megabyte.** Bits/ Inch²
- **Seek Time:** (2-12 ms) Current Areal Density ~ 640 Gbits / Inch²
- **The time needed to move the read/write head arm.**
- **Reported values: Minimum, Maximum, Average.**
- **Rotation Latency or Delay: (2-8 ms)**
- **The time for the requested sector to be under the read/write head. (~ time for half a rotation)**
- **Transfer time:** The time needed to transfer a sector of bits.
- **Type of controller/interface:** SCSI, EIDE (PATA, SATA)
- **Disk Controller delay or time.**
- **Average time to access a sector of data =**

average seek time + average rotational delay + transfer time

+ disk controller overhead

(ignoring queuing time)

Access time = average seek time + average rotational delay



Current Rotation speed
7200-15000 RPM

Basic Disk Performance Example

- Given the following Disk Parameters:

- Average seek time is 5 ms
- Disk spins at 10,000 RPM
- Transfer rate is 40 MB/sec

i.e. $T_{\text{queue}} = T_q = 0$

- Controller overhead is 0.1 ms

- Assume that the disk is idle, so no queuing delay exist.

i.e.
 T_{Ser}

- What is Average Disk read or write service time for a 500-byte (.5 KB) Sector?

Time for half a rotation

$$\begin{aligned}
 &\text{Ave. seek} + \text{ave. rot delay} + \text{transfer time} + \text{controller overhead} \\
 = & 5 \text{ ms} + 0.5 / (10000 \text{ RPM} / 60) + 0.5 \text{ KB} / 40 \text{ MB/s} + 0.1 \text{ ms} \\
 = & 5 + 3 + 0.13 + 0.1 = 8.23 \text{ ms}
 \end{aligned}$$

Access Time

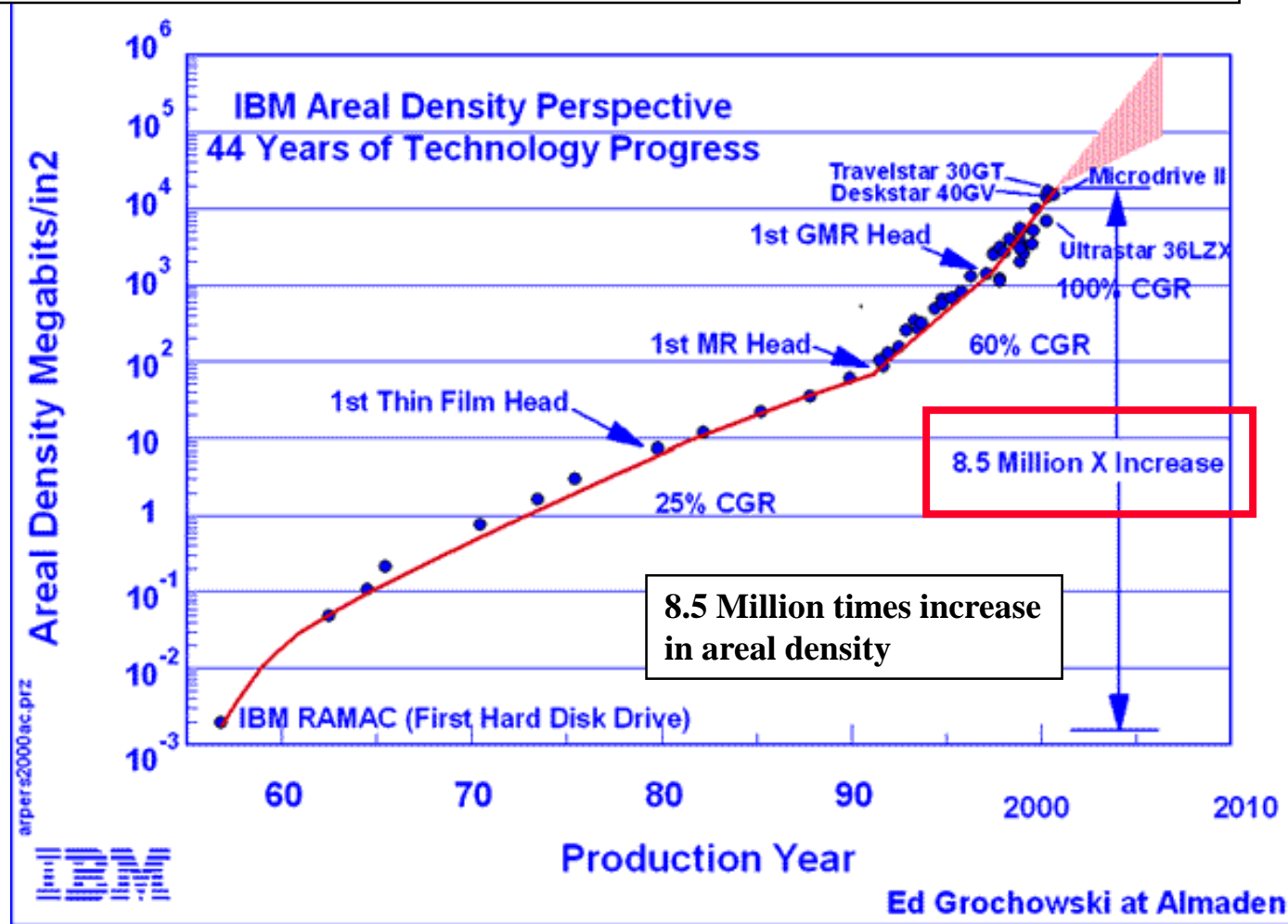
Actual time to process the disk request is greater and may include CPU I/O processing Time and queuing time

T_{service} (Disk Service Time for this request)

Or T_{ser}

Here: 1KBytes = 10^3 bytes, MByte = 10^6 bytes, 1 GByte = 10^9 bytes

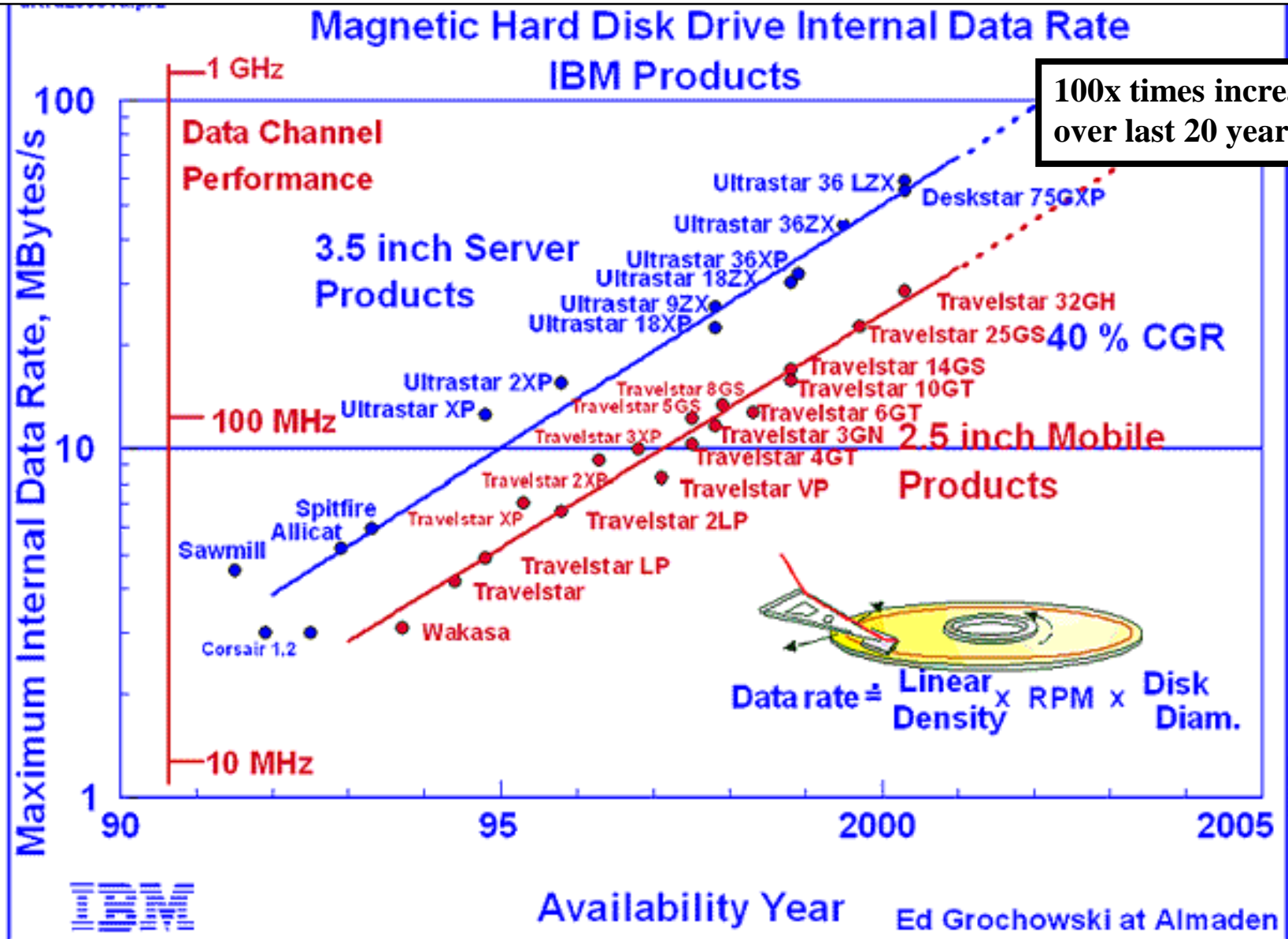
Historic Perspective of Hard Drive Characteristics Evolution: Areal Density



Drive areal density has increased by a factor of 8.5 million since the first disk drive, IBM's RAMAC, was introduced in 1957. Since 1991, the rate of increase in areal density has accelerated to 60% per year, and since 1997 this rate has further accelerated to an incredible 100% per year.

Current Areal Density ~ 640 Gbits / In²

Historic Perspective of Hard Drive Characteristics Evolution: Internal Data Transfer Rate



Internal data transfer rate increase is influenced by the increase in areal density

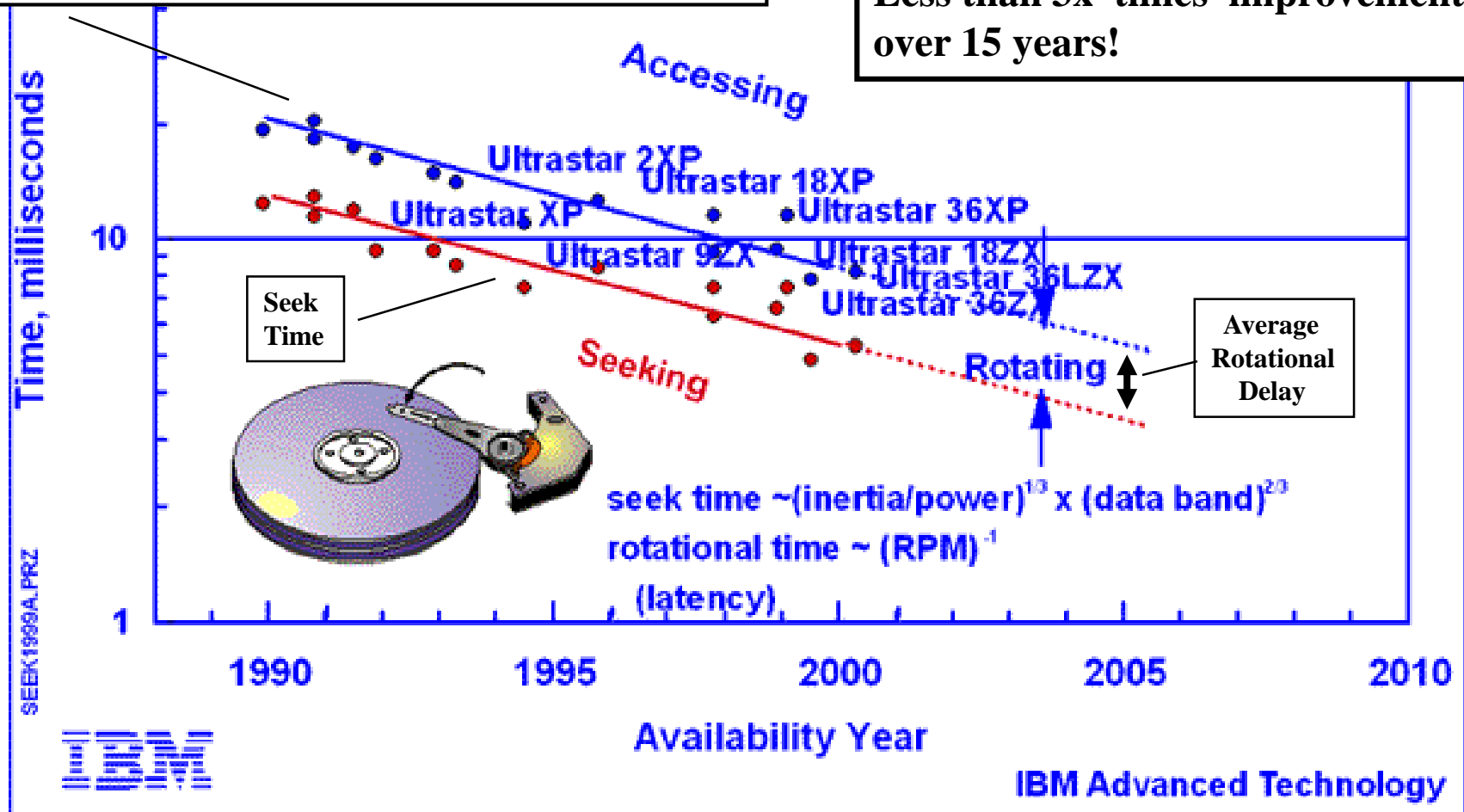
Historic Perspective of Hard Drive Characteristics Evolution: Access/Seek Time

100

IBM HDD Access/Seek Time-Performance Increase

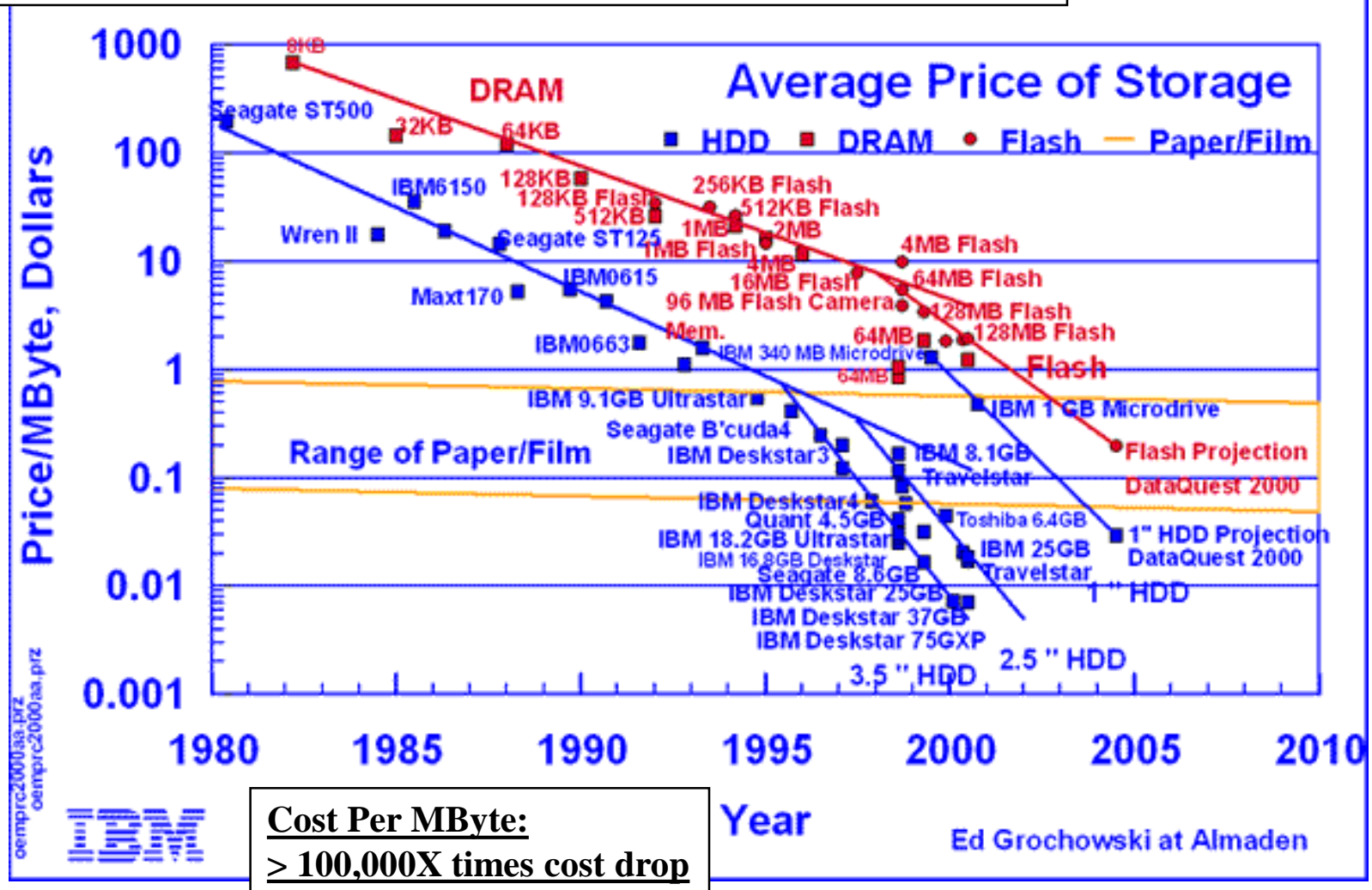
Access time = average seek time + average rotational delay

Less than 3x times improvement over 15 years!



Access/Seek Time is a big factor in service(response) time for small/random disk requests. Limited improvement due to mechanical rotation speed + seek delay

Historic Perspective of Hard Drive Characteristics Evolution: Cost



Cost Per MByte:

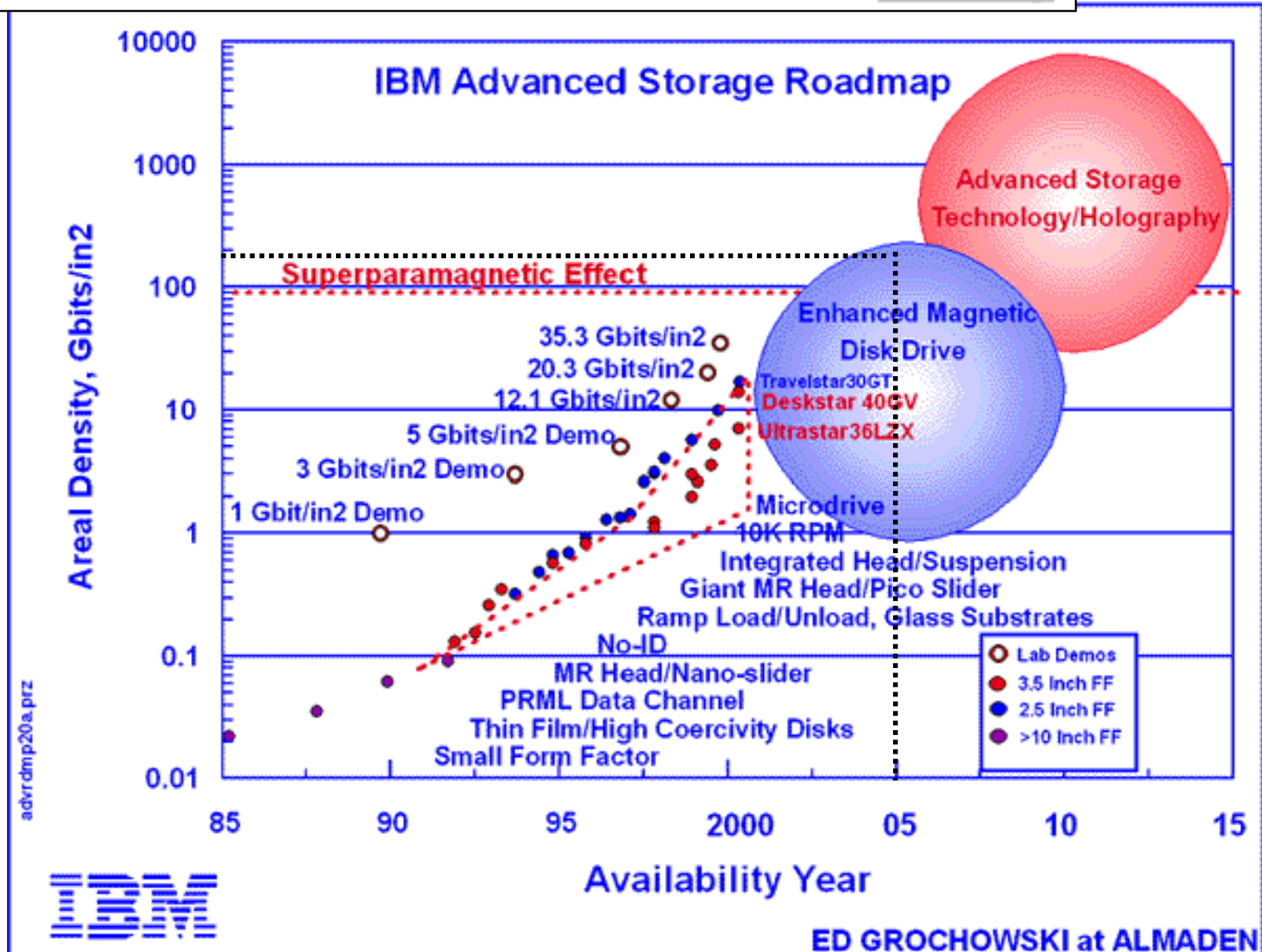
> 100,000X times cost drop

The price per megabyte of disk storage has been decreasing at about 40% per year based on improvements in **data density**,-- even faster than the price decline for flash memory chips. Recent trends in HDD price per megabyte show an even steeper reduction.

Actual Current Hard Disk Storage Cost (Fourth Quarter 2014):

~ 0.00005 dollars per MByte or about 20 GBytes /Dollar

Historic Perspective of Hard Drive Characteristics Evolution: Roadmap



Current Areal Density ~ 640 Gbits / In²

Factors Affecting System & I/O Performance

- I/O processing computational requirements:

- CPU computations available for I/O operations.
- Operating system I/O processing policies/routines.
- I/O Data Transfer/Processing Method used.
 - CPU cycles needed: Polling >> Interrupt Driven > DMA

CPU

- I/O Subsystem performance:

- Raw performance of I/O devices (i.e magnetic disk performance).
- I/O bus capabilities.
- I/O subsystem organization. i.e number of devices, array level ..
- Loading level (u) of I/O devices (queuing delay, response time).

Service Time, Tser, Throughput

I/O

- Memory subsystem performance:

- Available memory bandwidth for I/O operations (For DMA)

Tq

Memory

- Operating System Policies:

- File system vs. Raw I/O.
- File cache size and write Policy.
- File pre-fetching, etc.

OS

Components of Total System Execution Time:

CPU

Memory

I/O

System performance depends on many aspects of the system

(“limited by weakest link in the chain”): The system performance bottleneck