

Machine Called Computer

Part 7

Underlying Technologies and Evolution

References:

1. Computer Organization and Design & Computer Architecture, Hennessy and Patterson (slides are adapted from those by the authors)

Semiconductor Technology

학습요령:

- Scaling의 개념과 잇점 그리고 결과적인 추세 이해하실 것
 - 그리고 wafer, die, yield 용어를 이해
 - "technology" 라는 용어의 용법 이해
 - 숫자나 IC 제조공정은 암기대상이 아님

Why Transistor?

- ❑ Solid-state semiconductor device (반도체 장치)
 - Small, fast, reliable, energy-efficient, inexpensive
 - Integrated circuits (IC) 형태로 집적 가능

Image of cross-section of CMOS inverter (two transistors):

http://en.wikipedia.org/wiki/File:Cmos_impurity_profile.PNG

Semiconductor Technology

- ❑ Transistor: invented in Bell labs. in 1947
 - Took 10 years to commercialize
- ❑ IC (integrated circuits): invented in 1958
 - 5 years to commercialize
 - SSI, MSI, LSI, VLSI
- ❑ Major driving force behind computer performance evolution

CMOS NAND Gate

Image of CMOS NAND gate:

http://en.wikipedia.org/wiki/File:CMOS_NAND.svg

Image of CMOS NAND layout:

http://en.wikipedia.org/wiki/File:CMOS_NAND_Layout.svg

Image of CMOS transistor pair:

http://en.wikipedia.org/wiki/File:Cmos_impurity_profile.PNG

Technology Scaling (data from Wikipedia)

- ❑ Minimum feature size: $10\mu\text{m}$ (1970) to $0.022\mu\text{m}$ (2012)
 - Exponential decrease

Min. Feature Size	year
$10\mu\text{m}$	1970
$3\mu\text{m}$	1975
$1\mu\text{m}$	1985
350 nm	1995
130 nm	2002
45 nm	2008
22 nm	2012
7 nm	2018 (estimated)

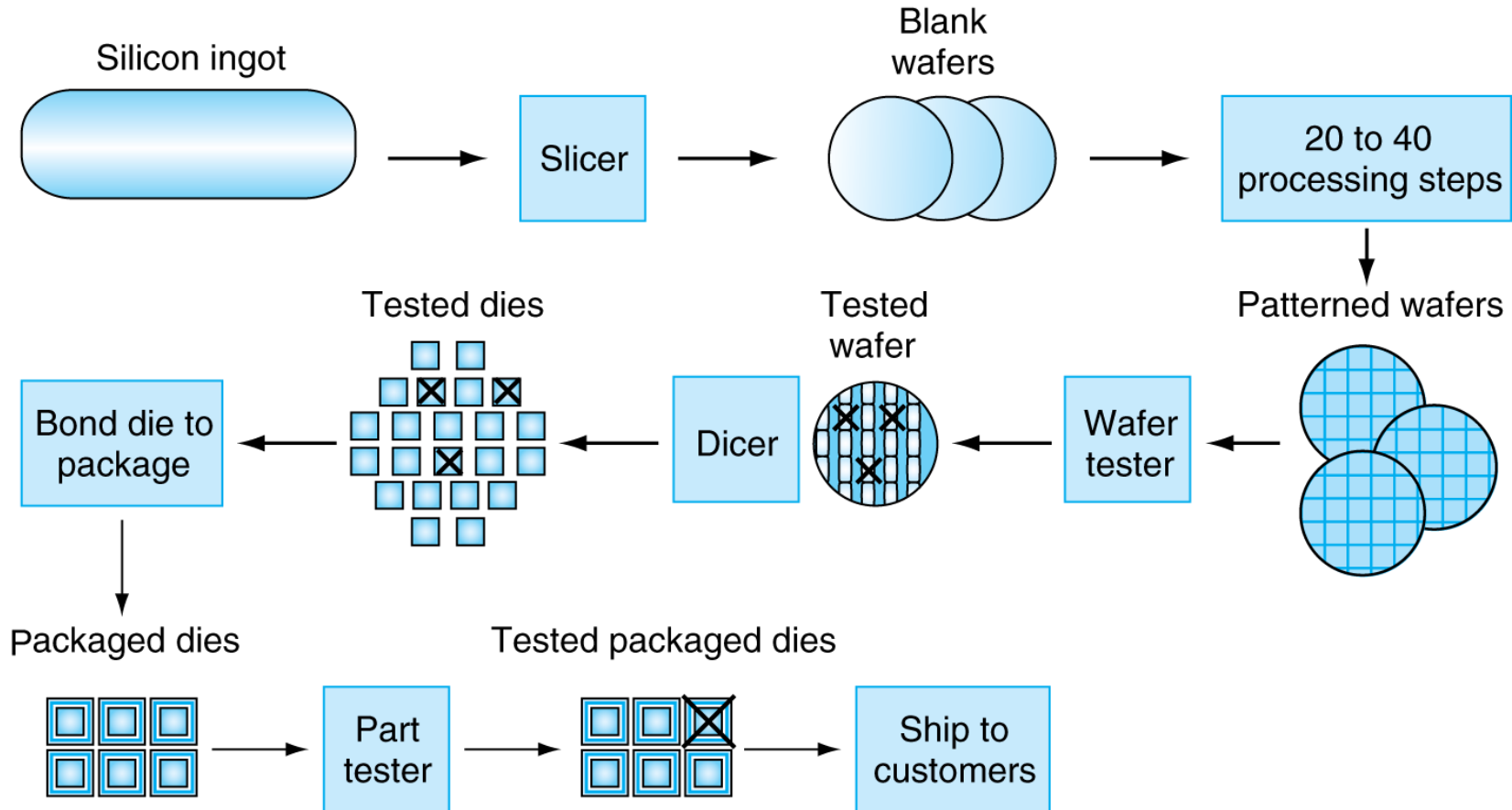
Technology Trends

- Smaller is faster

Year	Technology	Relative performance/cost
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuit (IC)	900
1995	Very large scale IC (VLSI)	2,400,000
2005	Ultra large scale IC	6,200,000,000

Manufacturing ICs

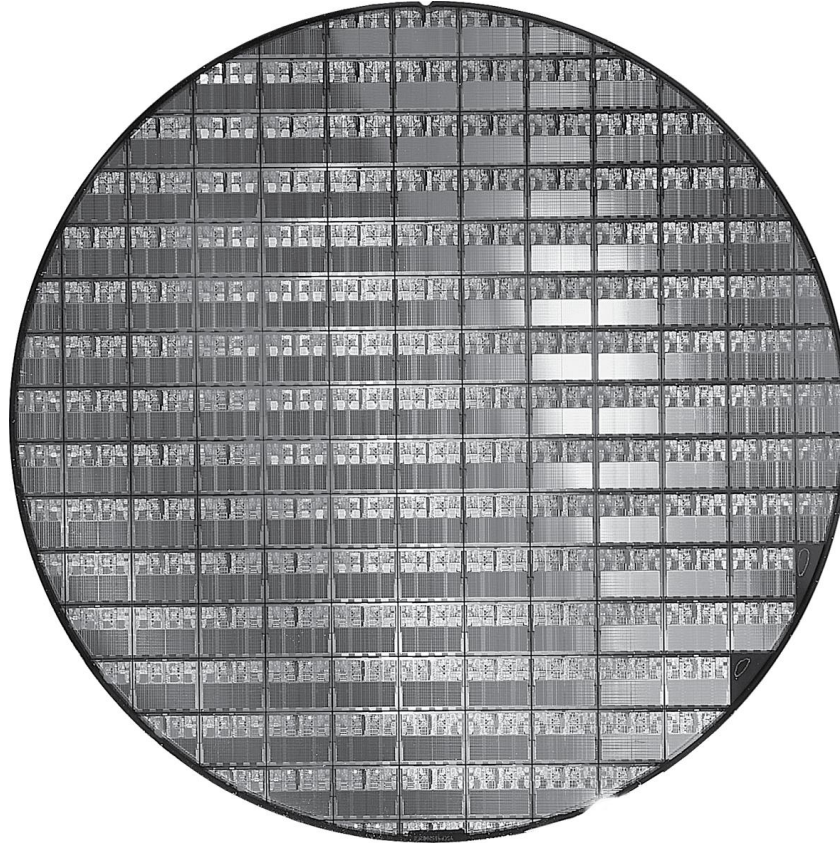
(Hennessy and Patterson slide, Computer Organization and Design, Morgan Kaufmann)



- **Yield**: proportion of working **dies** per **wafer**

AMD Opteron X2 Wafer

(Hennessy and Patterson slide, Computer Organization and Design, Morgan Kaufmann)

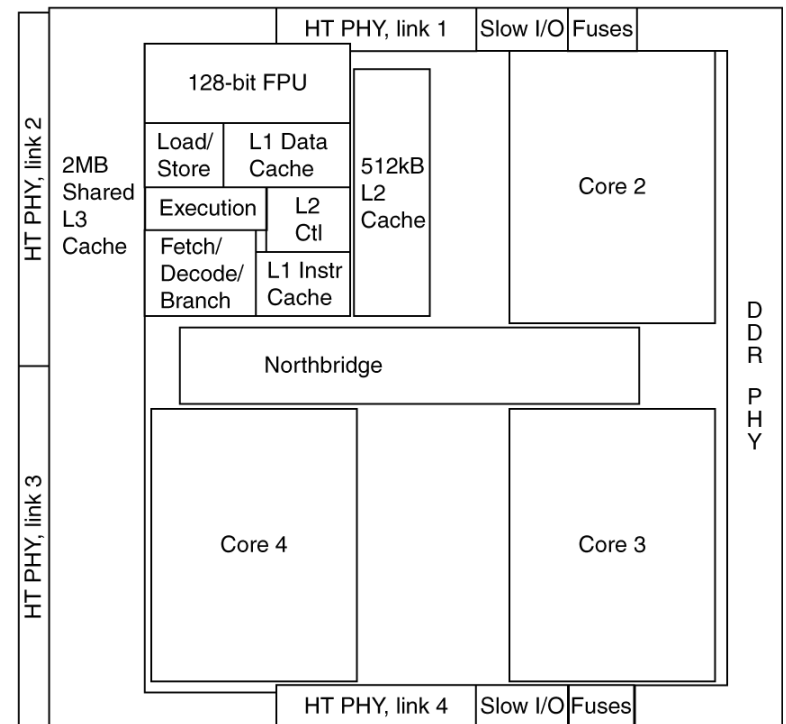
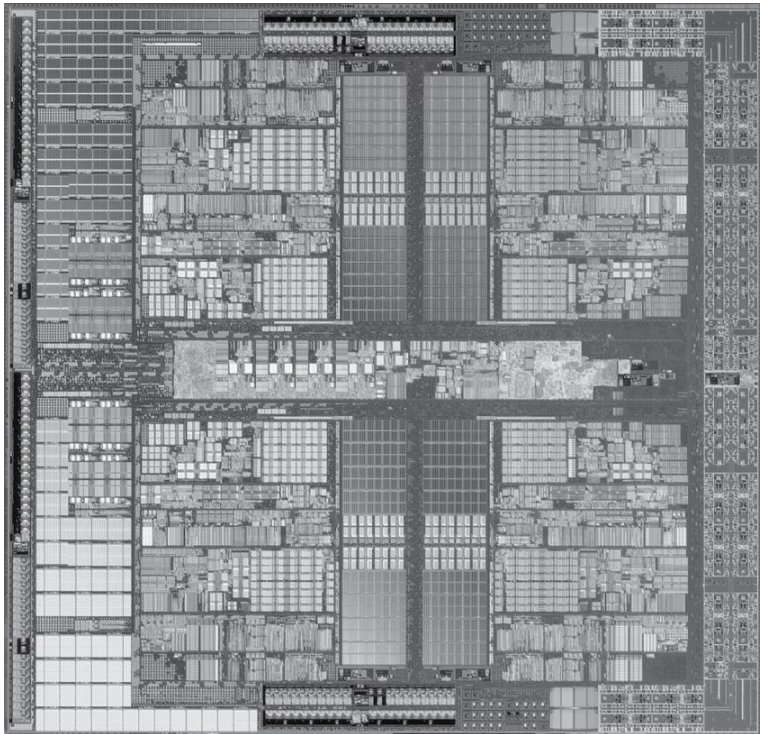


- ❑ X2: 300mm wafer, 117 chips, 90nm technology
- ❑ X4: 45nm technology

Inside the Processor

(Hennessy and Patterson slide, Computer Organization and Design, Morgan Kaufmann)

- AMD Barcelona: 4 processor cores



Semiconductor Technology

- ❑ What does it mean?
 - We have 32 nm technology
 - We have 300 mm wafer technology
 - What if we use larger (e.g., 450mm) wafer?
- ❑ Moore's law: exponential growth
 - Number of transistors per chip double every 18 (or 12 or 24) months
 - Cost of fabrication facility also increase exponentially over time

Intel and Processor Technology

학습요령:

- Microprocessor의 출현과 발전 흐름 이해
 - RISC 프로세서 용어 이해
- * 전체적인 흐름이 중요하고, 아주 세부적인 내용이나 숫자는 기억할 필요 없음

Semiconductor Technology

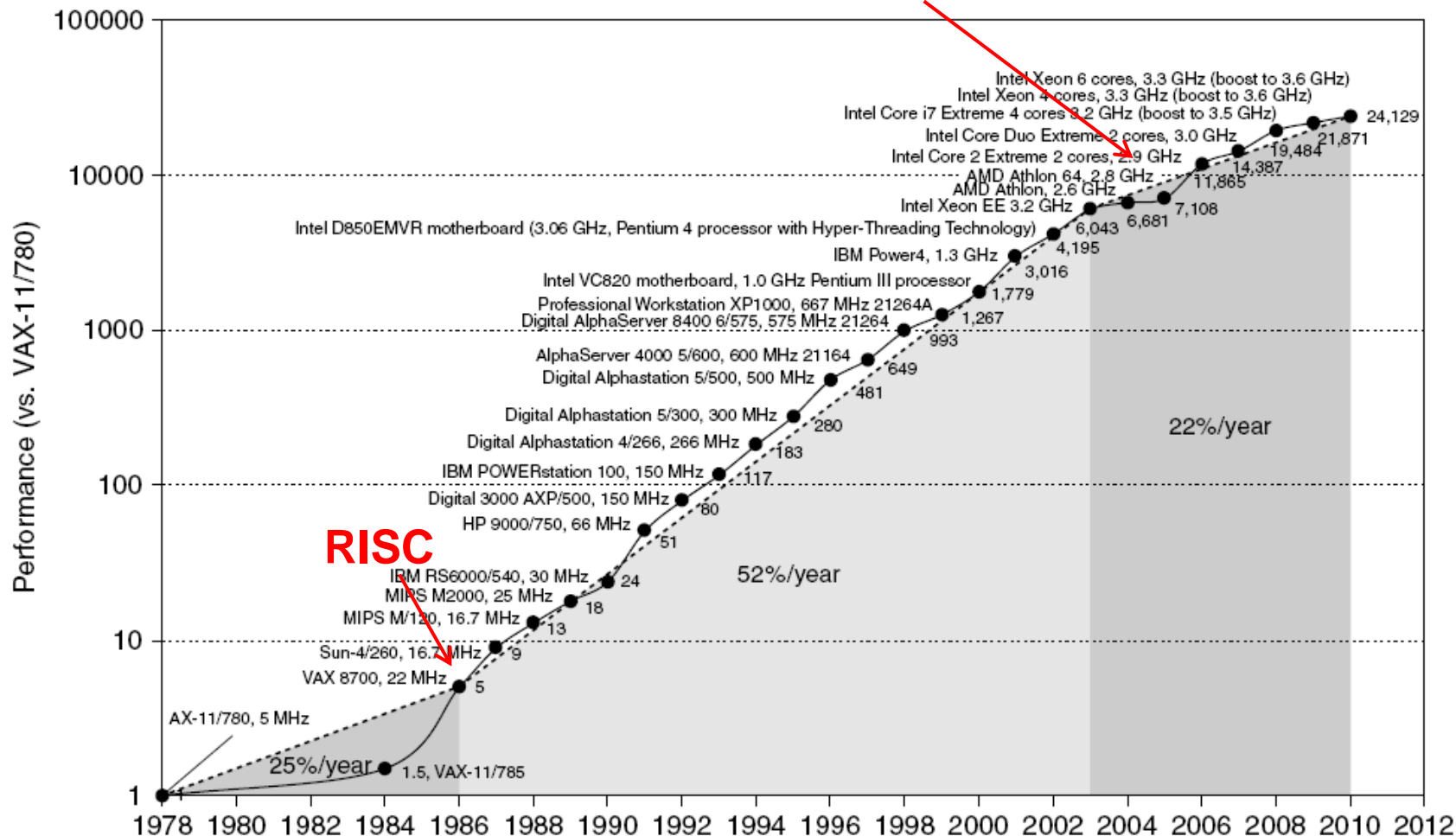
- ❑ Major driving force behind computer performance evolution
- ❑ Smaller transistor, increased die size
 - Processor perspective
 - Exponential growth in performance
- ❑ Around 2012
 - Highest transistor count in commercial CPU
 - 2.5B in Intel's 10-core Xeon Westmere-EX (32nm)
 - † FPGA: Xilinx 6.8B in Virtex-7 (28nm)
 - † Samsung 20nm 4Gb DRAM
 - † Samsung 10nm 64Gb flash memory

Transistor Count in Processor (Wikipedia)

Processor	#Transistors	Year	Process (μm)	Area (mm^2)	
Intel 4004	2,300	1971	10	12	First μp
Intel 8080	4,500	1974	6	20	
Intel 8088	29,000	1979	3	33	IBM PC, 16 bit
Intel 80286	134,000	1982	1.5	49	PC/AT
Intel 80386	275,000	1985	1.5	104	x86, IA-32
Intel 80486	1,180,000	1989	1	160	
Pentium	3,100,000	1993	0.8	294	
Pentium II	7,500,000	1997	0.35	195	
Pentium 4	42,000,000	2000	0.18		
Itanium 2	220,000,000	2003	0.13		IA-64, RISC
Core i7 (Quad)	731,000,000	2008	0.045	263	x86-64
10-core Xeon Westmere-EX	2,600,000,000	2011	0.032	512	x86

Single Processor Performance

Move to multi-processor



Trends in Technology

- ❑ Integrated circuit technology (e.g., processor)
 - Transistor density: 35%/year
 - Die size: 10-20%/year (why increase?)
 - Integration overall: 40-55%/year
- ❑ DRAM capacity: 25-40%/year
- ❑ Flash memory capacity: 50-60%/year
 - 15-20X cheaper/bit than DRAM
- ❑ Magnetic disk technology: 40%/year
 - 15-25X cheaper/bit than Flash
 - 300-500X cheaper/bit than DRAM

CPU in Mainframes in 1971

- ❑ Large printed circuit boards
 - Several of them for a CPU (IC technology immature)
- ❑ Design cycle: 5 years (HW + OS + Appl.)

Image of mainframe CPU in 1971:

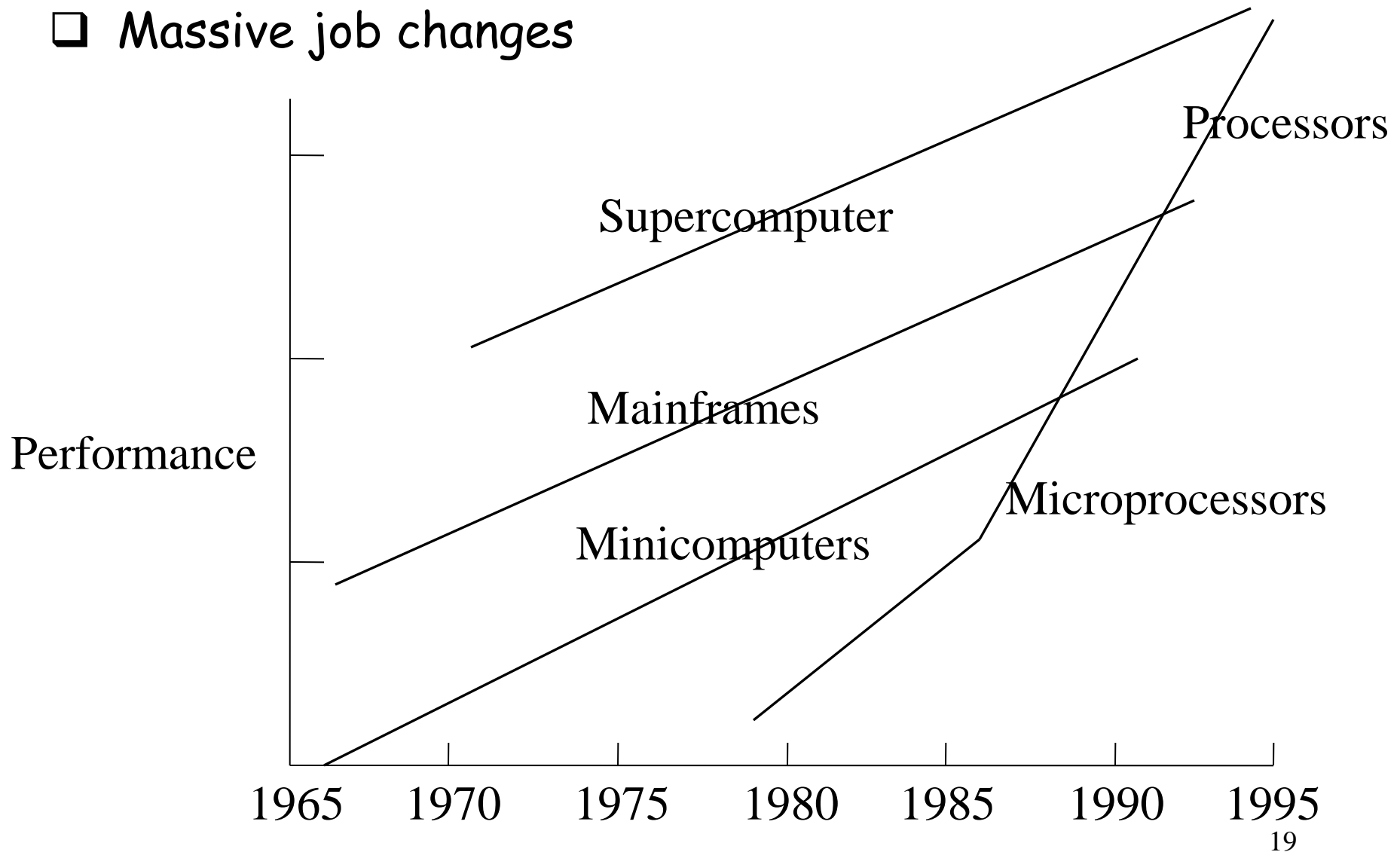
http://en.wikipedia.org/wiki/File:386DX40_MB_Jaguar_V.jpg

Microprocessor

- ❑ Breakthrough: 4004 microprocessor by Intel in 1971
 - Miniature version of minicomputer CPU
 - Designed for Busicom to build calculators
- ❑ Intel look for market
- ❑ Intel announce 8-bit 8008 in 1972, 8080 in 1974
- ❑ Altair in 1975: first personal computer with 8080
 - Along the style of best minicomputers
- ❑ IPM PC in 1981 use 16-bit 8088 and MS DOS
- ❑ Cloning: from IBM PC to PC (dominance by Intel, MS)
 - Intel microprocessors: single chip, short design cycle
 - Take full advantage of semiconductor technology¹⁸

(Approximate) Technology Trend

□ Massive job changes



Intel and Processors

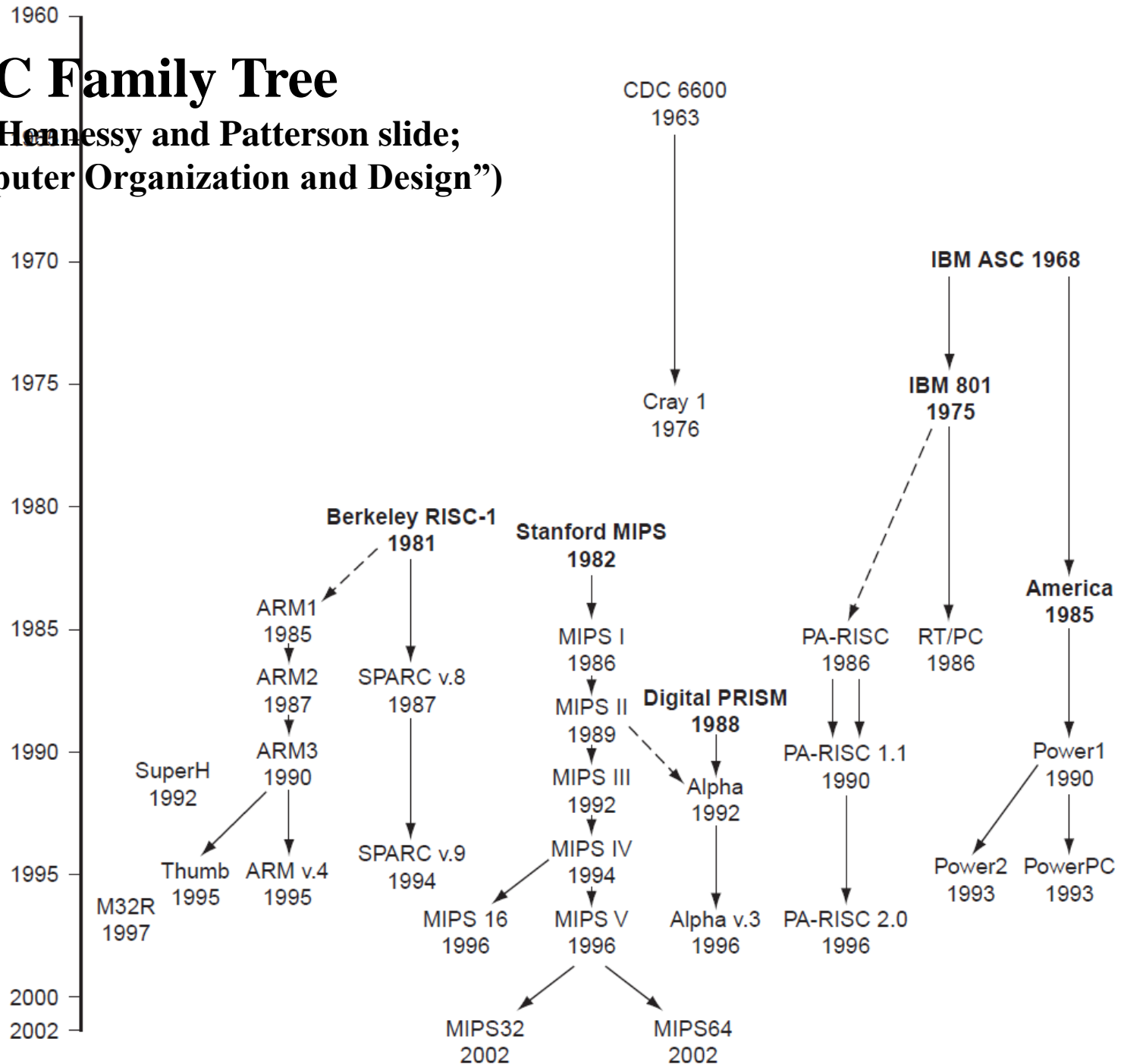
- ❑ Intel's microprocessors become powerful processors
- ❑ Computer manufacturers
 - Buy processors from processor vendors
 - Massive job changes in computer companies in 1980s
 - From hardware to software (5:5 -> 1:9)
 - Focus on systems, software, service
- † Small microprocessors still there for low-cost embedded systems (many many of them)

Instruction Set 변화 - RISC Processors

- ❑ Processor design in 1970s: what we call CISC
 - Constraint: memory expensive
- ❑ 1980s: renaissance of processor design (RISC style)
 - Semiconductor technology
 - Memory become cheaper
 - Open Unix operating system
 - High-level programming
- ❑ Emergence of powerful 32-bit RISC processors
 - PowerPC, PA-RISC, MIPS, SPARC, Alpha, ARM
 - † Exception is Intel Pentium

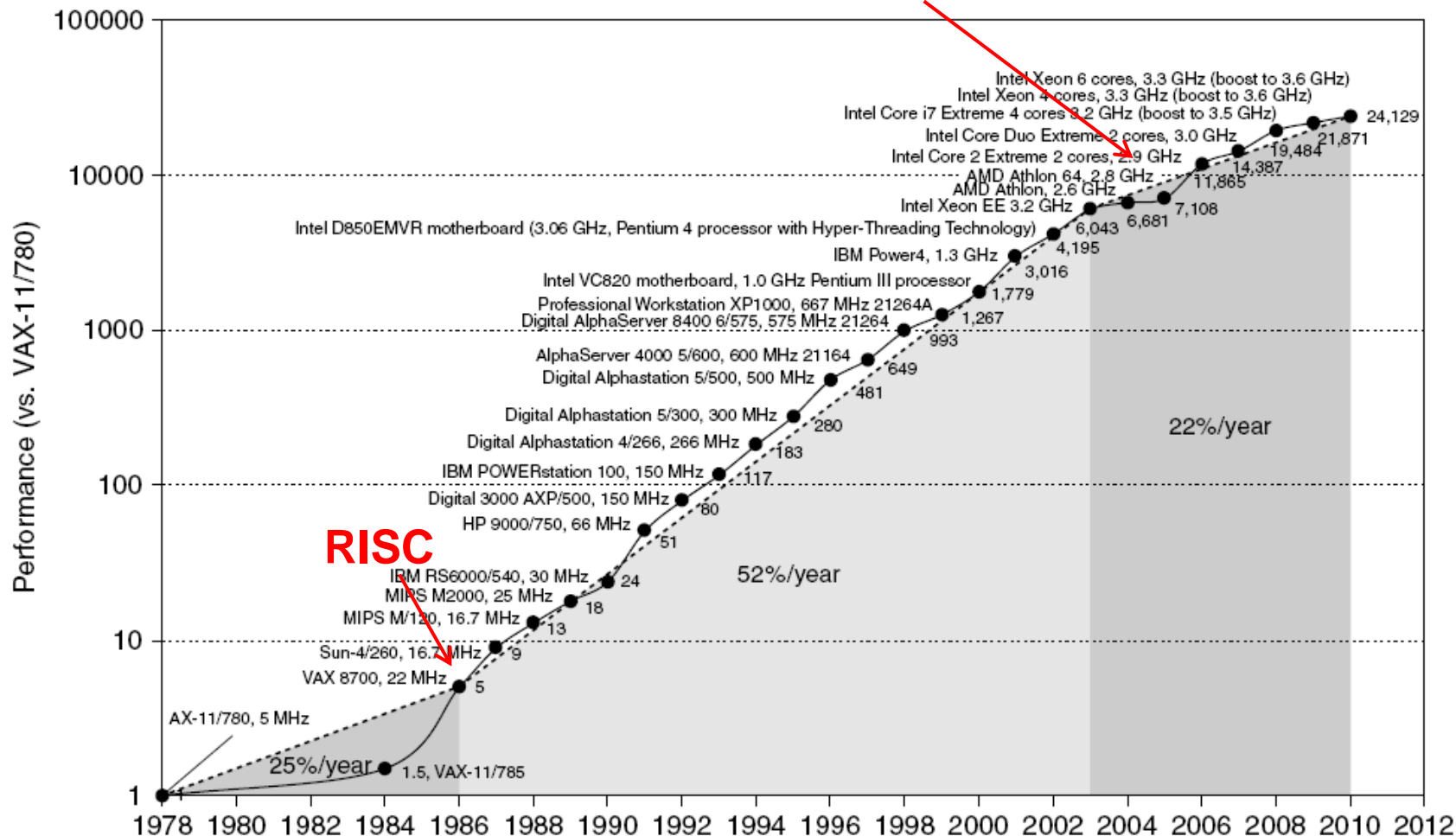
RISC Family Tree

(from Hennessy and Patterson slide;
“Computer Organization and Design”)



Single Processor Performance

Move to multi-processor

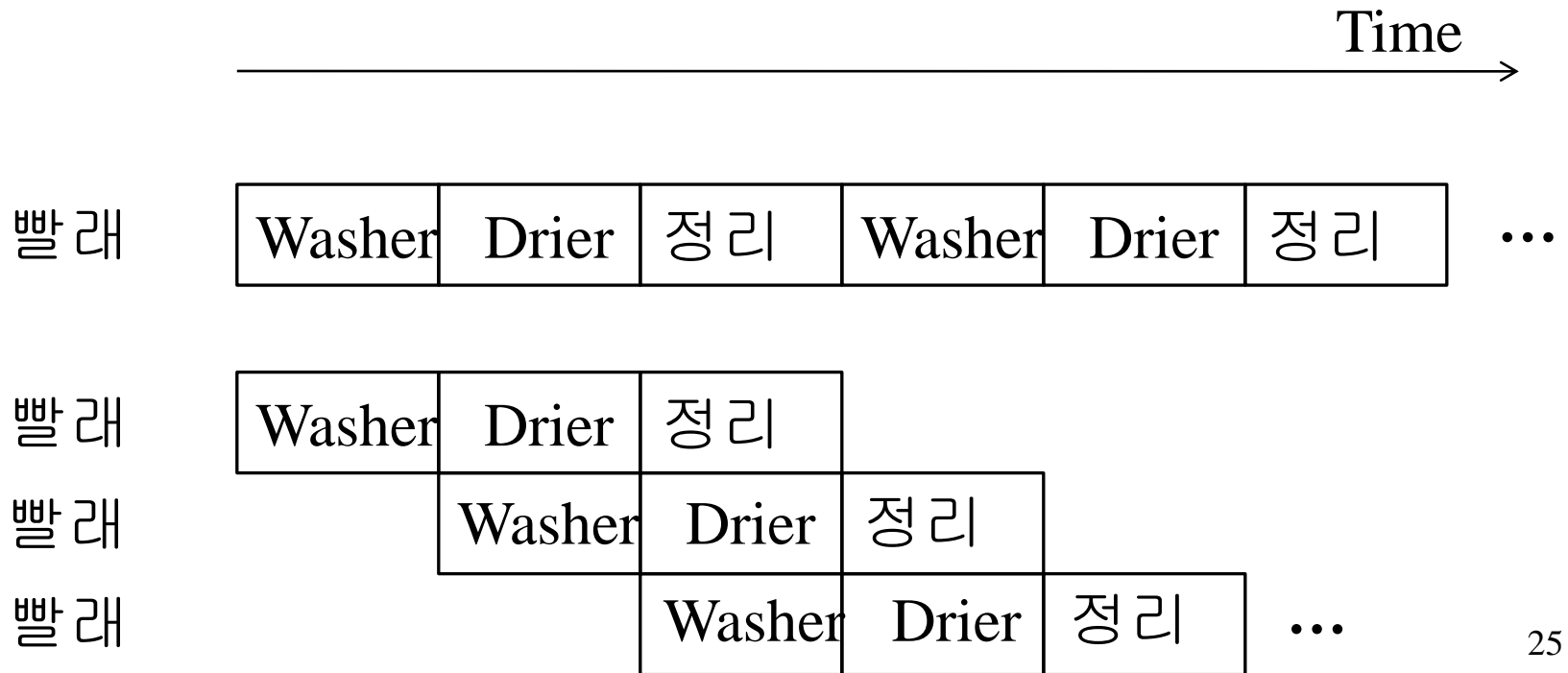


Key Speedup Techniques in CPU

- ❑ Pipelining
- ❑ Cache memory

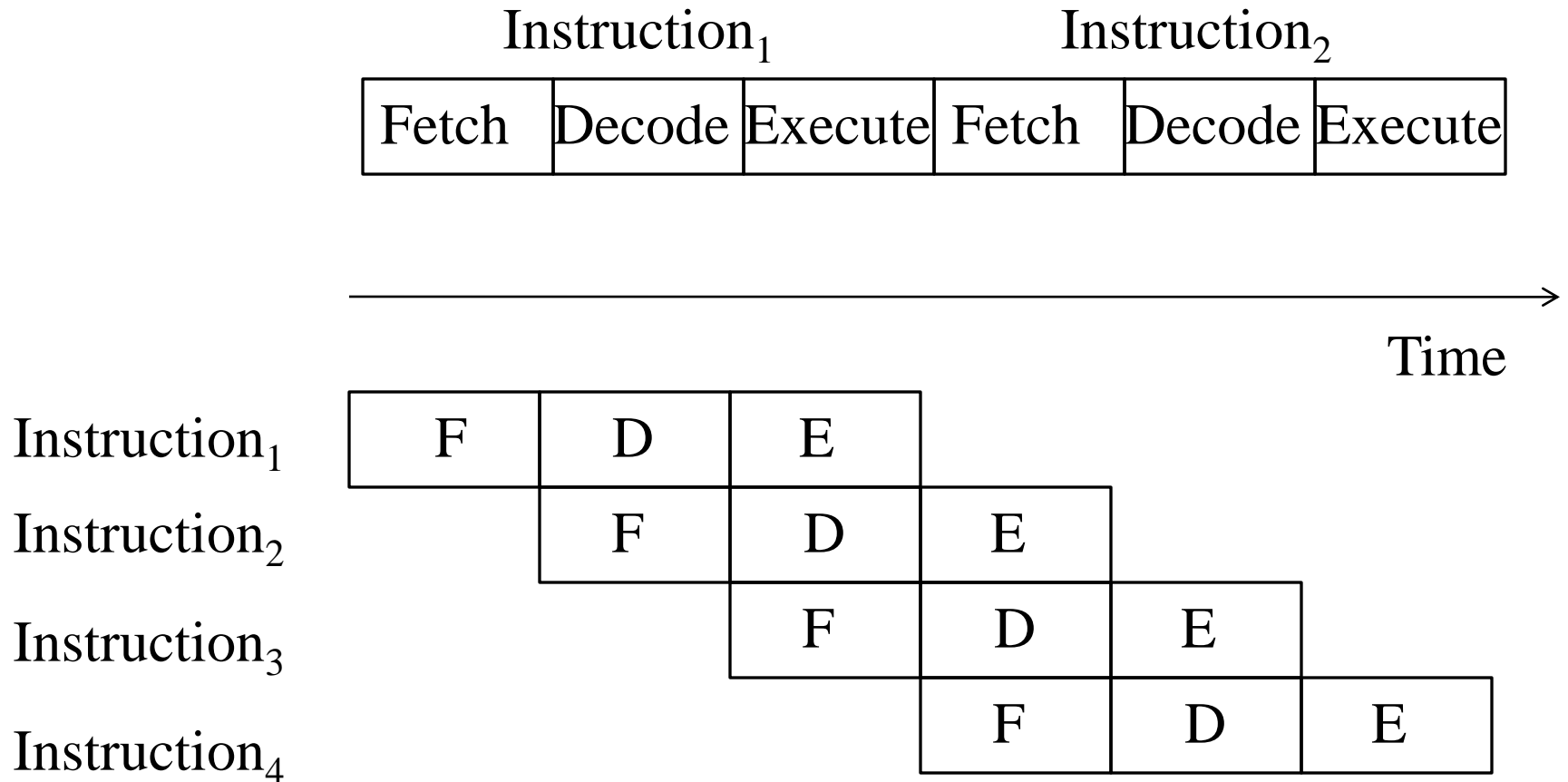
Pipelining - General Speedup Technique

- 3-stage pipeline (e.g., washer-dryer example)
 - Speedup?



Pipelining

- 3-stage pipeline for fetch-decode-execute

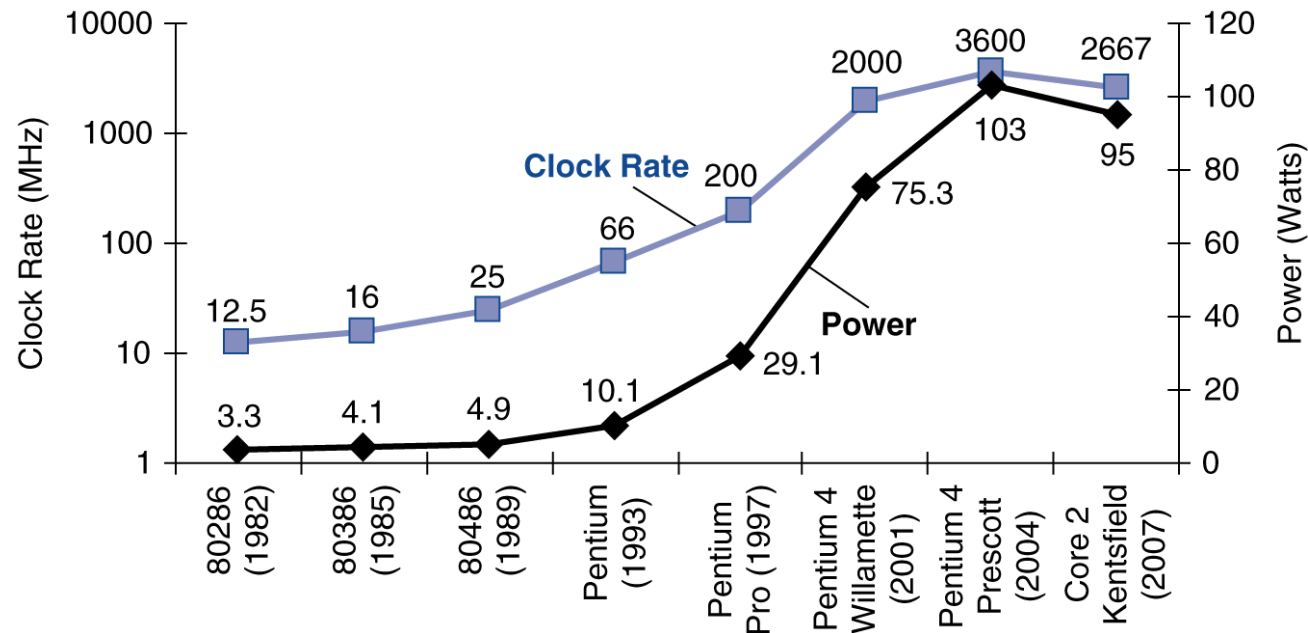


Advanced Pipelining

- ❑ Powerful server processors
 - Ideal speedup for 10-stage pipeline?
 - What if we build 4 pipelines per processor?
 - What is the ideal speedup?

Power Wall (skip)

- ❑ Intel 80386 consumed ~ 2 W
- ❑ 3.3 GHz Intel Core i7 consumes 130 W
- ❑ Heat must be dissipated from 1.5 x 1.5 cm chip
- ❑ This is the limit of what can be cooled by air



Memory Technology, Memory Systems

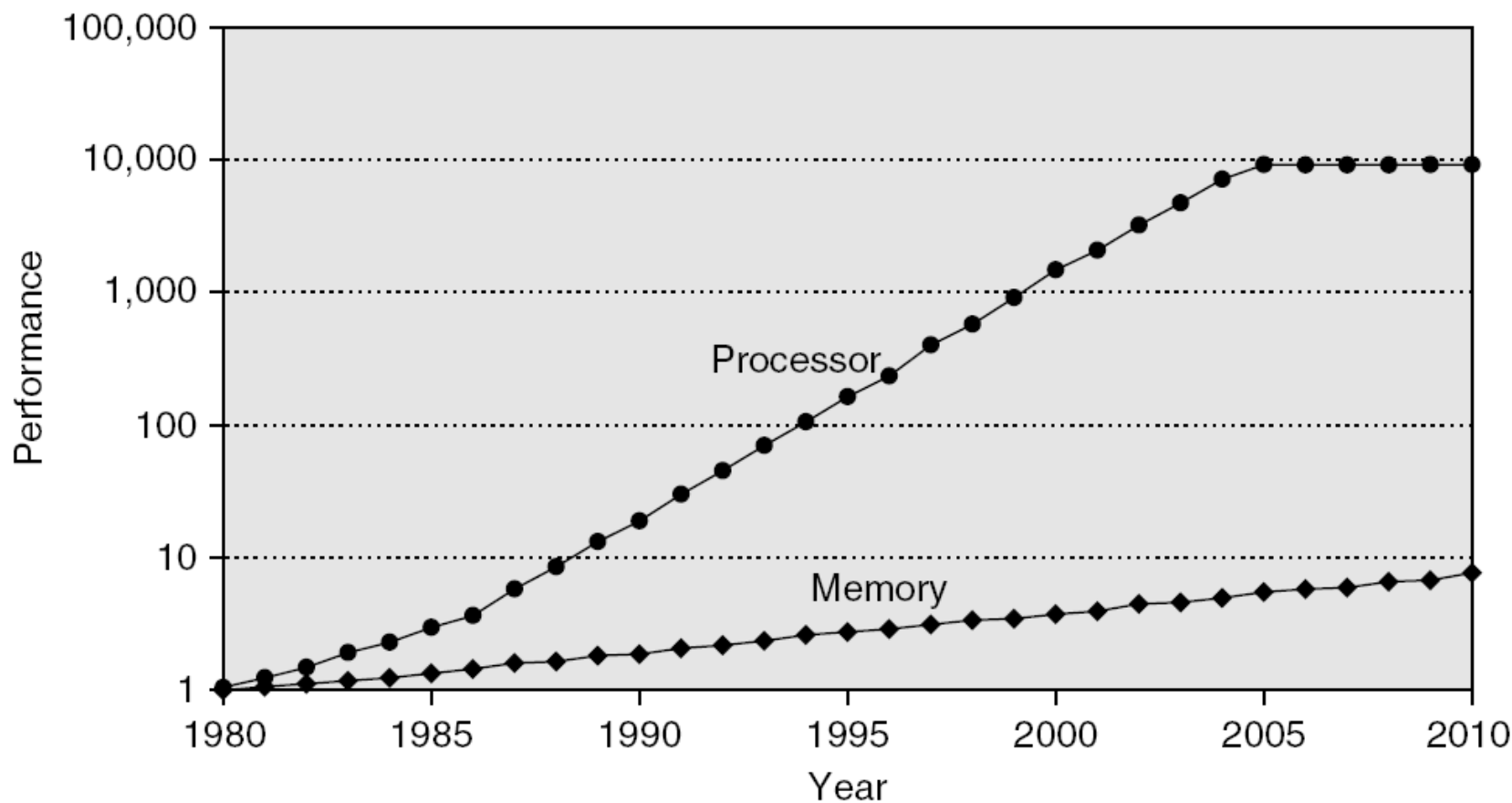
학습요령:

- Semiconductor memory 기술의 추세 이해
- Caching의 개념 이해

Semiconductor Technology

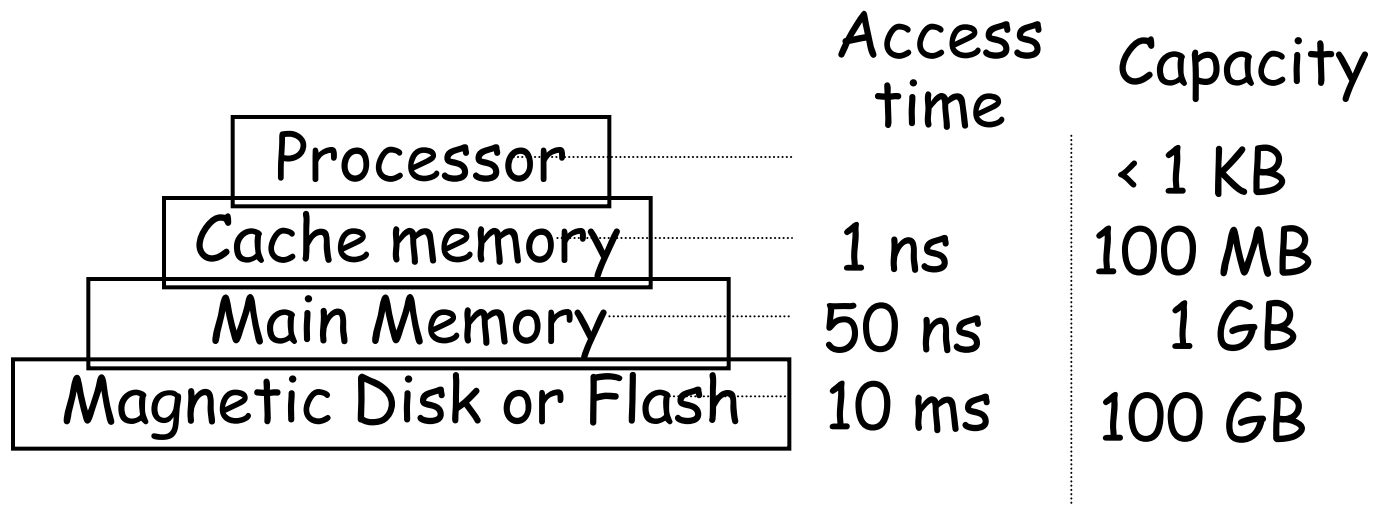
- ❑ Major driving force behind computer performance evolution
- ❑ Smaller transistor, increased die size
 - Processor perspective
 - Exponential growth in performance
 - Memory perspective
 - Exponential capacity growth
 - Speed not improve significantly

CPU-Memory Performance Gap

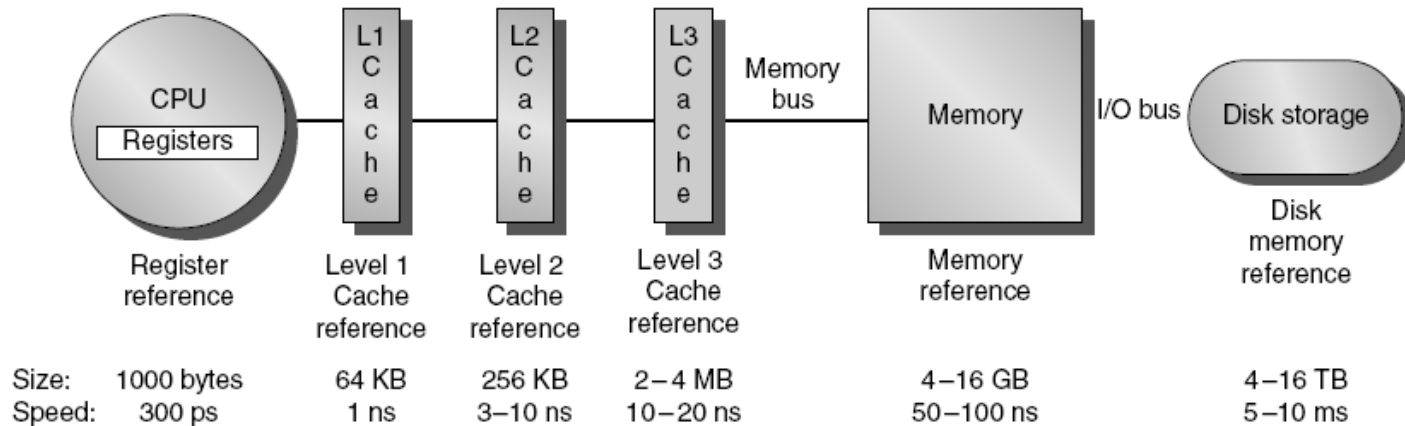


Memory Hierarchy

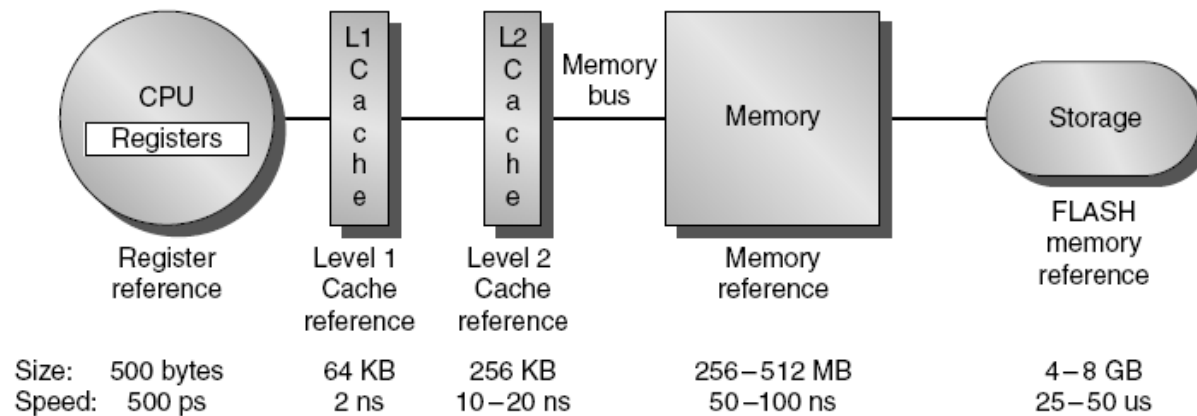
- ❑ Memory: performance bottleneck
- ❑ How to build (illusion of) “ideal memory”
 - Current technology: SRAM, DRAM, Disk (or flash)



Memory Hierarchy (skip)



(a) Memory hierarchy for server



(b) Memory hierarchy for a personal mobile device

Semiconductor Memory

❑ SRAM

- Flip-flop invented by Eccles and Jordan in 1918
- Cache memory, volatile

❑ DRAM invented in 1966, IBM

- Main memory, volatile

+ NAND Flash memory

- Non-volatile
- Compete with hard disk, especially in mobile market

Semiconductor Memory (skip)

Image of six-transistor SRAM cell:

[http://en.wikipedia.org/wiki/File:SRAM_Cell_\(6_Transistors\).svg](http://en.wikipedia.org/wiki/File:SRAM_Cell_(6_Transistors).svg)

Image of DRAM:

http://en.wikipedia.org/wiki/File:Square_array_of_mosfet_cells_read.png

Hard Disk (skip)

- ❑ Invented by IBM in 1953, first commercial use in 1956
- ❑ Secondary memory

Image of IBM hard disk in 1956:

http://en.wikipedia.org/wiki/File:IBM_350_RAMAC.jpg

Image of hard disk drive:

<http://en.wikipedia.org/wiki/File:HardDiskAnatomy.jpg>

Trends - Survival of Fittest (반복)

- ❑ Integrated circuit technology (e.g., processor)
 - Transistor density: 35%/year
 - Die size: 10-20%/year
 - Integration overall: 40-55%/year
- ❑ DRAM capacity: 25-40%/year
- ❑ Flash memory capacity: 50-60%/year
 - 15-20X cheaper/bit than DRAM
- ❑ Magnetic disk technology: 40%/year
 - 15-25X cheaper/bit than Flash
 - 300-500X cheaper/bit than DRAM

Semiconductor Flash Memory (skip)

- ❑ Invented around 1980, Toshiba
 - Toshiba announce NAND flash memory in 1989
 - Replace hard disc in mobile devices
 - † Expensive, but reliable, low-power
 - Memory cards, USB flash drives, solid-state drives

Image of NAND flash memory:

http://en.wikipedia.org/wiki/File:USB_flash_drive.JPG

Optical Disc (skip)

- ❑ Invented in 1958
 - First commercial use in 1972
 - Subsequent CD (1980), DVD(1995), CD-RW (1996)
 - Audio, video, computer archive storage

Image of the surface of compact disc:

http://en.wikipedia.org/wiki/File:CD_autolev_crop.jpg

Image of various optical storage media:

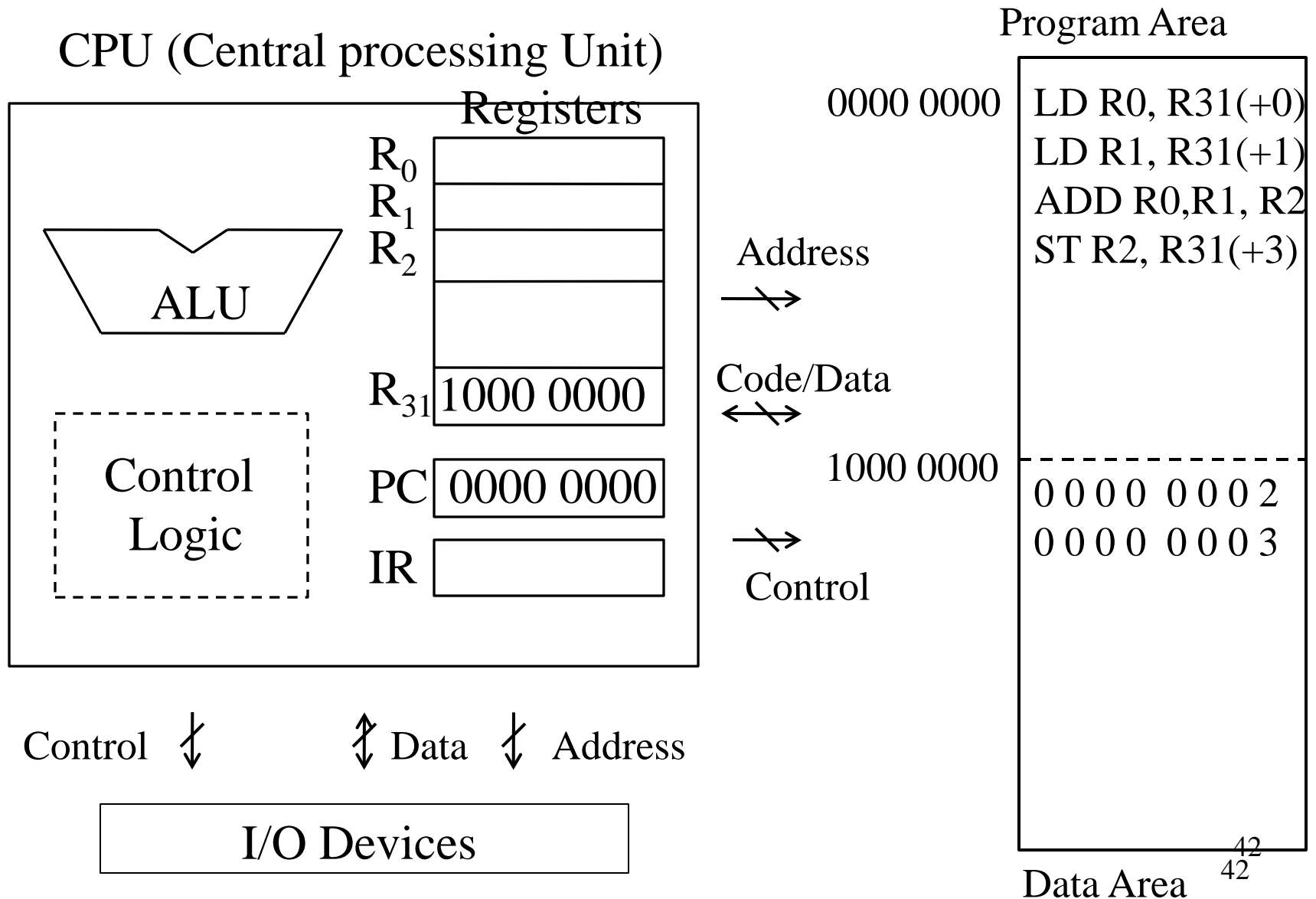
http://en.wikipedia.org/wiki/File:Comparison_CD_DVD_HDDVD_BD.svg

Question on Data Loss

- ❑ I don't want to lose the data in my PC
 - Backup in optical disk or in external hard disk?
 - How long would it last?
 - What is your solution?
- ❑ Financial companies in New York
 - Risk: war, earthquake, tsunami, ...
 - What is the state-of-the-art?
- ❑ Heard about company specialized in backup and archive?
 - What kind of facilities would they have?

More on Computer

Computer Hardware

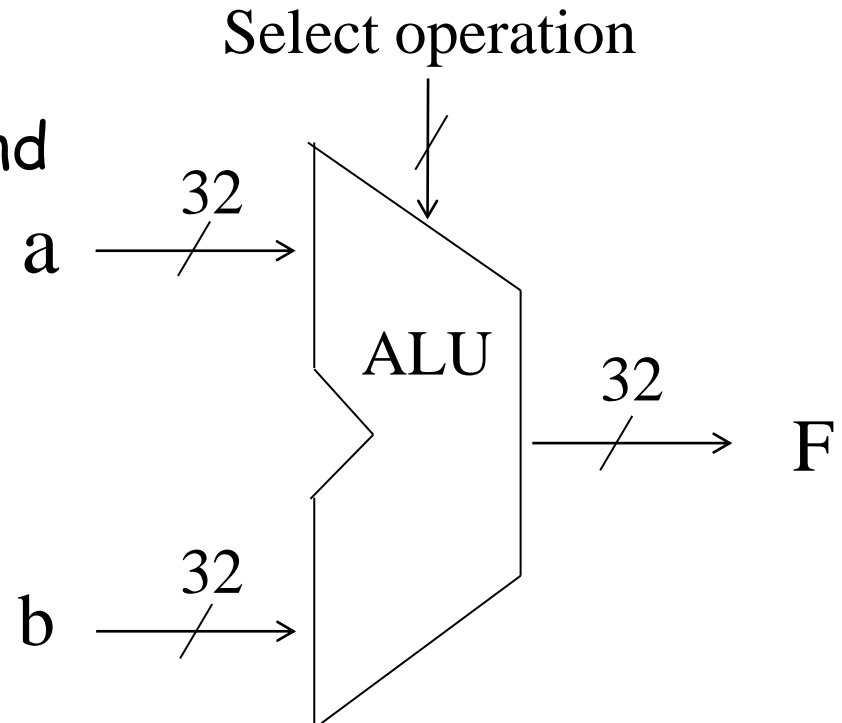


Quiz: X-Bit Computer

- ❑ What does 'x-bit' means?
 - Size of ALU input operand
 - Size of register
 - Width of processor internal bus
 - Width of processor-memory bus
 - Width of I/O bus
 - No. of data lines (pins) of processor
 - No. of address lines (or pins) of processor
 - Length of instruction
- ❑ Is 64-bit processor better than 32-bit processor?

X-Bit Computer

- ❑ What does 'x-bit' mean?
 - Size of ALU input operand



- ❑ Is 64-bit computer better than 32-bit computer?
 - Larger numbers
 - Speedup with parallel operation

Size of Address Space

❑ What else is important?

- Size of address space
 - What does that mean to programmers?

❑ Microprocessor history

Processor	data size	address size
• 8-bit	8	16
• 16-bit	16	16 (+a)
• 32-bit RISC	32	32
• 64-bit	64	?

Byte Addressing

- ❑ Viewed as a large, single-dimension array, with an address
- ❑ A memory address is an index into the array
- ❑ "Byte addressing" means that the index points to a byte of memory

0	8 bits of data
1	8 bits of data
2	8 bits of data
3	8 bits of data
4	8 bits of data
5	8 bits of data
6	8 bits of data

...

Byte Addressing

- ❑ Bytes are nice, but most data items use larger "words"
- ❑ For MIPS, a word is 32 bits or 4 bytes.

0	32 bits of data
4	32 bits of data
8	32 bits of data
12	32 bits of data
...	

Registers hold 32 bits of data

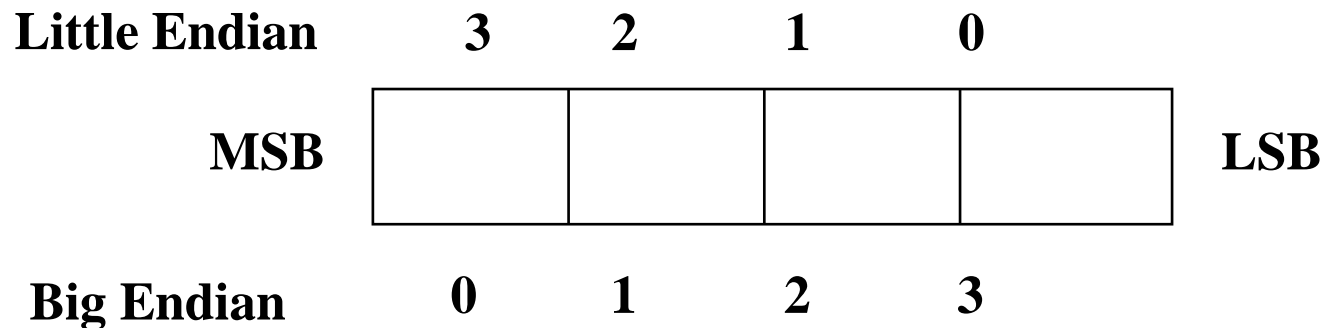
- 2^{32} bytes with byte addresses from 0 to $2^{32}-1$
- 2^{30} words with byte addresses 0, 4, 8, ... $2^{32}-4$
- ❑ Words are aligned
 - What are least 2 significant bits of a word address?

Alignment

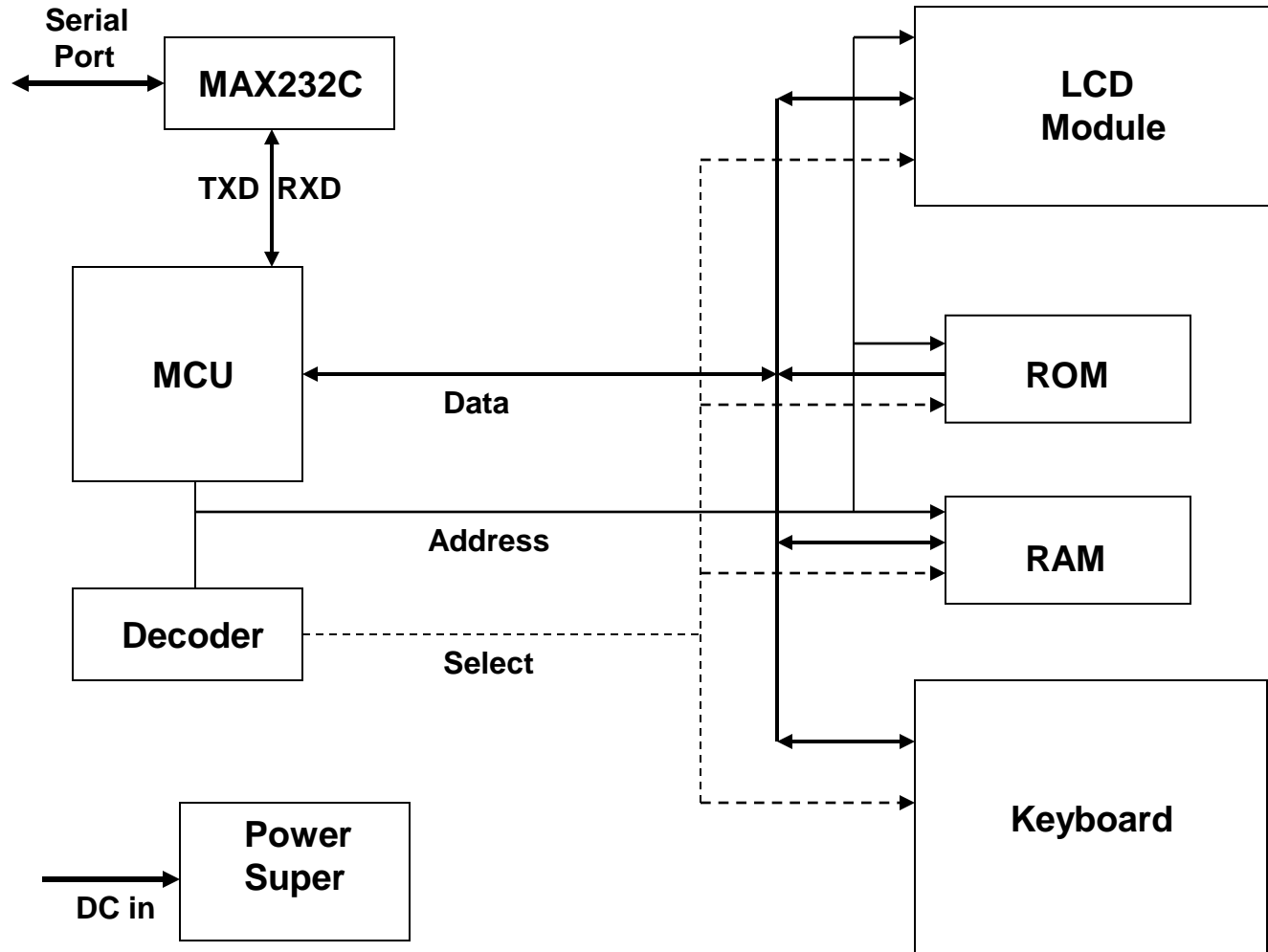
Width of object	0	1	2	3
1 byte (byte)	Aligned	Aligned	Aligned	Aligned
2 bytes (half word)	Aligned		Aligned	
2 bytes (half word)		Misaligned		Misaligned
4 bytes (word)	Aligned			
4 bytes (word)		Misaligned		
4 bytes (word)			Misaligned	
4 bytes (word)				

Little/Big Endian

- ❑ Around 1950s
 - A few mainframes in the world
 - "not invented here"
- ❑ Byte address: two conventions (refer to databook)



80196 Training Kit



Address Decoding

- ❑ To enable Chip Select using address lines
- ❑ Memory
 - Large number of address (to store program & data)
 - Each memory word has an address
- ❑ Peripheral (e.g., UART)
 - A few addresses
 - Control registers, data registers
- † Determine address map during hardware design

Address "Map" - 80196 example

FFFFH	External Memory or I/O Area	FF00H~FFFFH : Interrupt Double Vector	RAM
		FE00H~FEFFH : PTS Double Vector	
		8000H~FDFFH : User Program	
2080H	Special Purpose Memory Area	2080H~7FFFH : Monitor Program	ROM
207FH		2040H~205DH : PTS Vector	
		2030H~203FH : Upper Interrupt Vector	
		2018H : CCB	
		2010H~2013H : Special Interrupt Vector	
2000H		2000H~200FH : Lower Interrupt Vector	
1FFFFH	Port 4 Port 3	1FFE0H~1FFFFH : Address / Data BUS	On Chip
1FFFEH			
1FFDDH	External Memory or I/O Area	1F00H : User Select 3	
		1E00H : User Select 2	
		1D00H : User Select 1	
		1C00H : Key	
0200H		0C00H : Key IN	
		1A00H : LCD	
01FFFH	Register File	100H~1FFFH : Upper Register File	
		1AH~FFH : Register RAM	
		18H~19H : Stack Pointer	
0000H		00H~17H : SFR	

Processor Databook

❑ Processor data book - what do you expect to see?

† Must ask this kind of question whatever you do

- ISA
 - Instructions, addressing modes, encoding
- Physical interface
 - Pins, how to use them, timing
- Others
 - Environmental range, clock speed, power, voltage

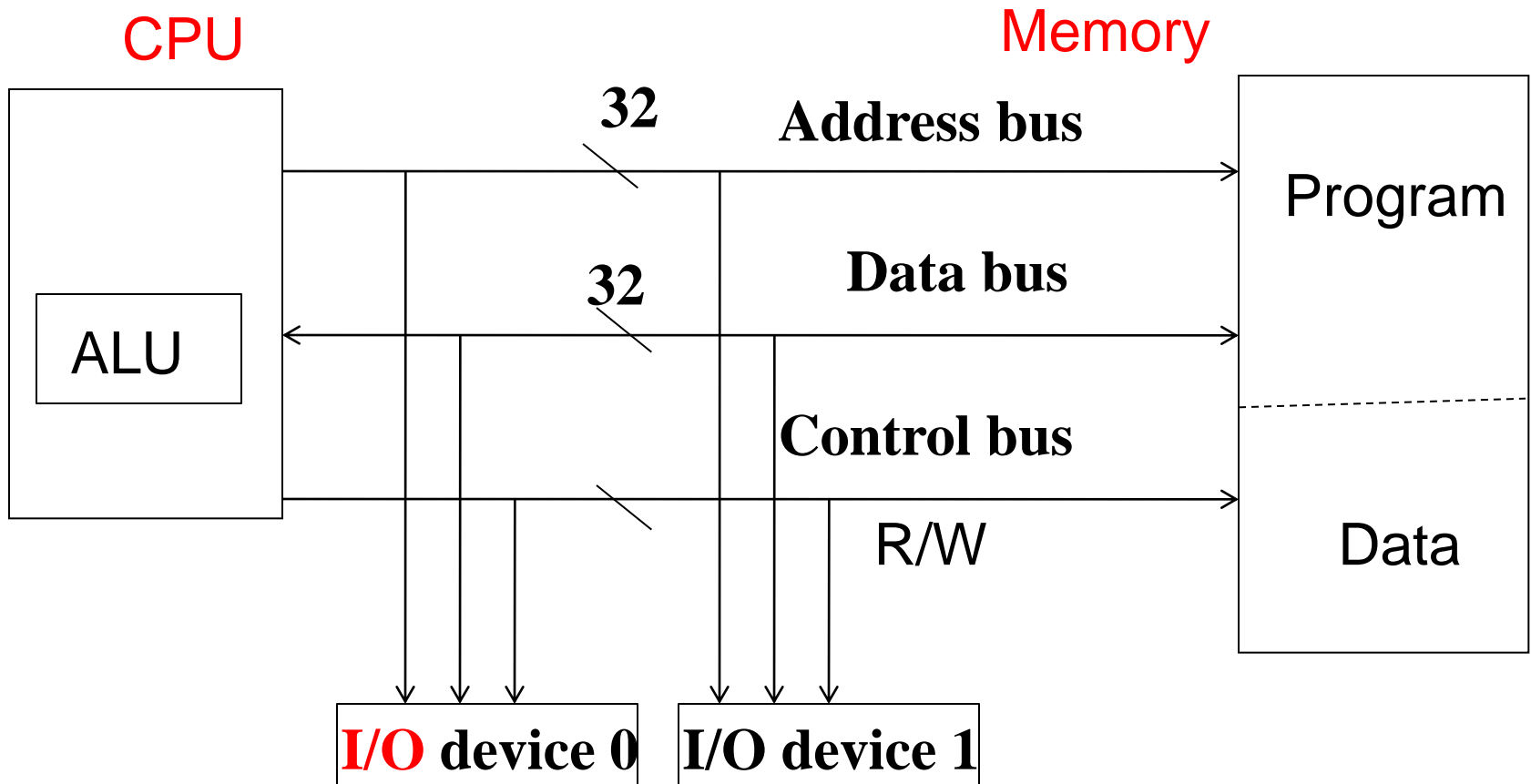
Microcontroller

- ❑ Microcontroller (versus microprocessor)
 - CPU core plus set of commonly-used peripherals
 - e.g., timer, memory, ADC/DAC, display controller
 - Ideal: single-chip solution
 - Faster, more reliable, less expensive
- ❑ What microcontroller data book additionally has
 - Functionality of each peripheral
 - Application examples
 - How to initialize/program each peripheral
 - Memory map for peripherals

Interconnection

- ❑ Data bus, address bus, control bus
- ❑ PCI, ISA, CAN, Ethernet, LAN
- ❑ What is bus?
 - Shared (broadcast) medium
 - Bus protocol, bus arbitration, bus controller
- ❑ Alternate topology
 - Mesh, tree, hypercube, complete connection
 - What about Internet?
- ❑ General interconnection issues
 - Unique address
 - Routing: how to deliver messages

32-bit Computer



- ❑ $4G = 2^{32}$ memory and I/O locations
- ❑ Given address, enable corresponding location

Interconnection Networks

(Hennessy and Patterson slide, Computer Organization and Design, Morgan Kaufmann)

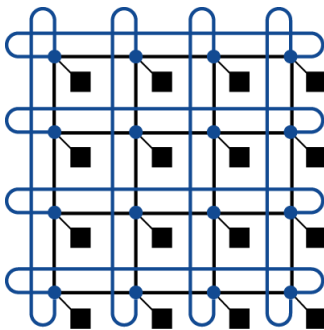
- Network topologies
 - Arrangements of processors, switches, and links



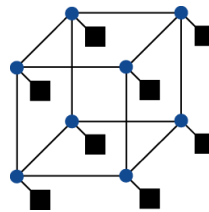
Bus



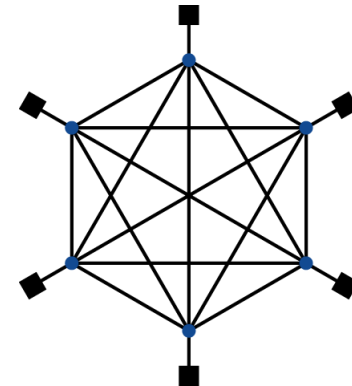
Ring



2D Mesh



N-cube ($N = 3$)

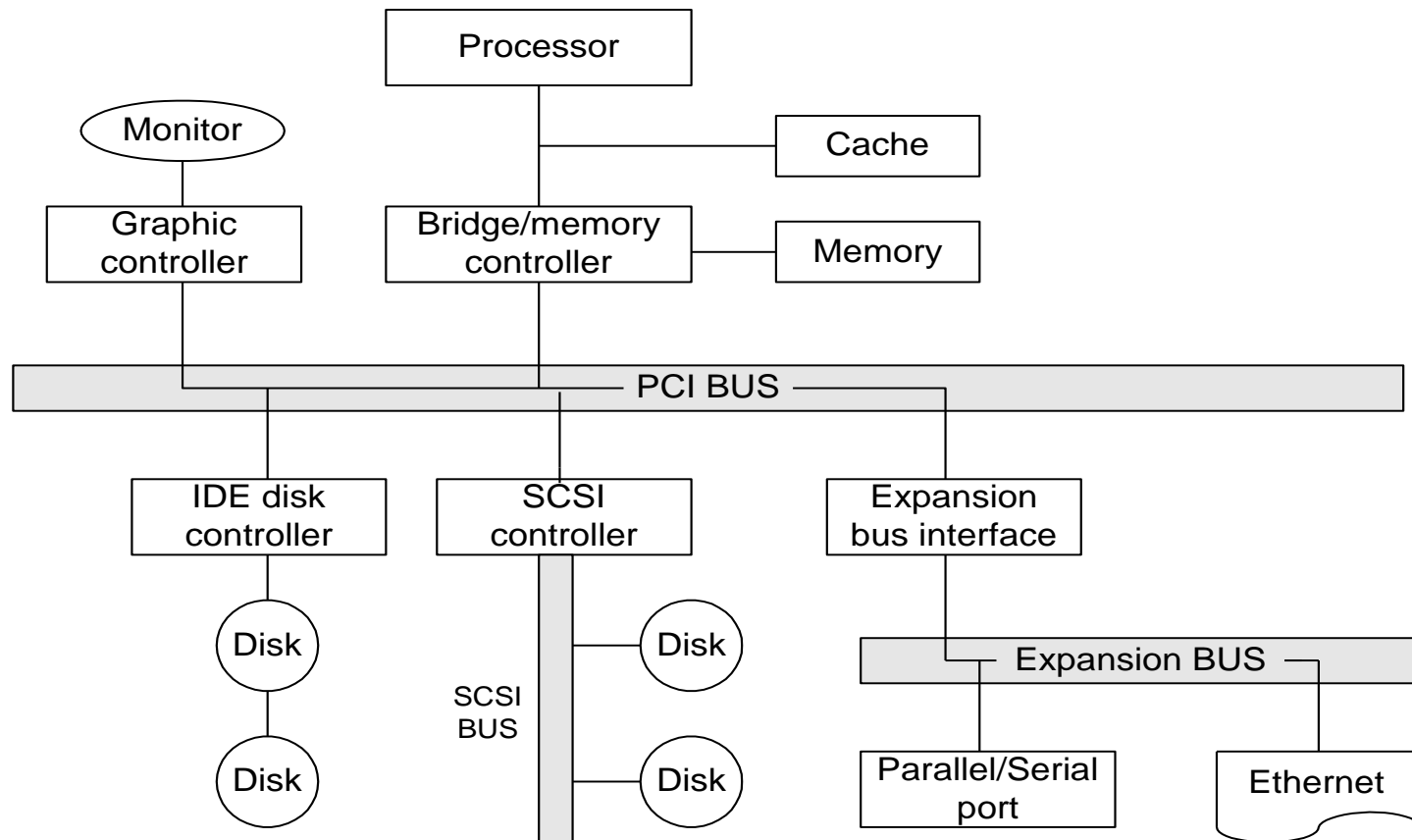


Fully connected

Interconnection

(Hennessy and Patterson slide, Computer Organization and Design, Morgan Kaufmann)

- ❑ Processor-memory bus: proprietary
- ❑ I/O bus: standard



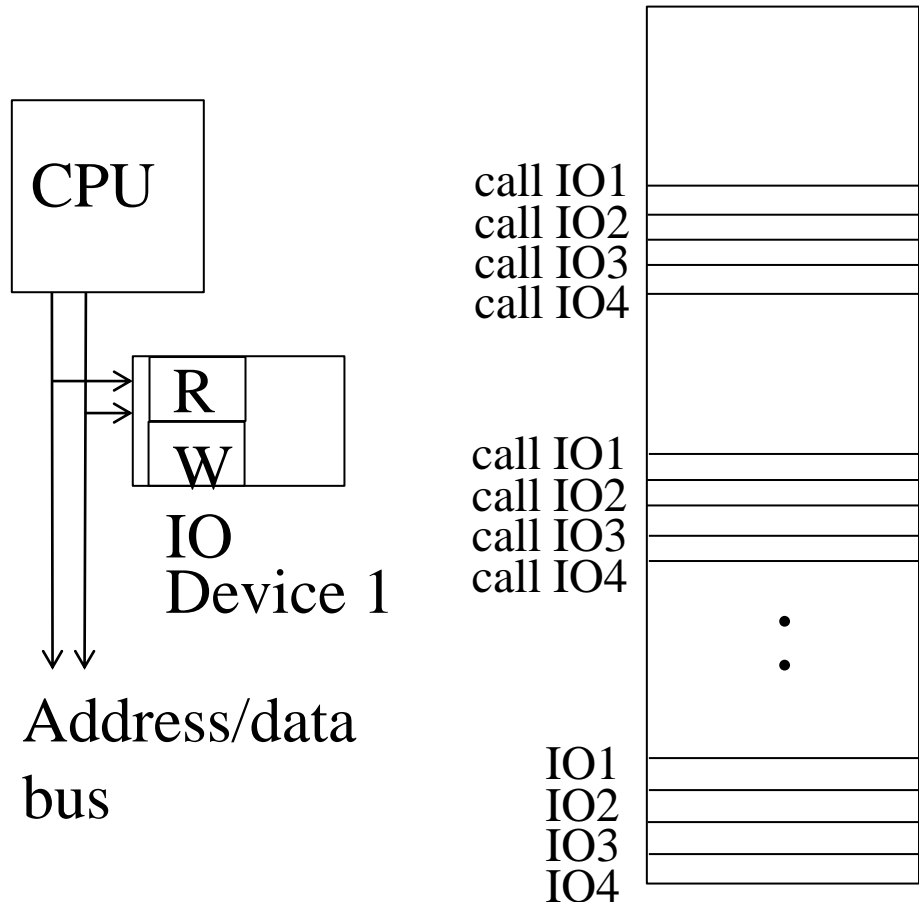
Interrupts (like Hardware Hotline)

- ❑ What is programmed I/O?
 - Periodic polling with existing hardware (address and data bus)
- ❑ Why programmed I/O not sufficient?
 - Can be burden to processor (in large systems)
 - Response time
- ❑ Interrupt
 - Extra hardware (hotline and machine support) to draw CPU attention

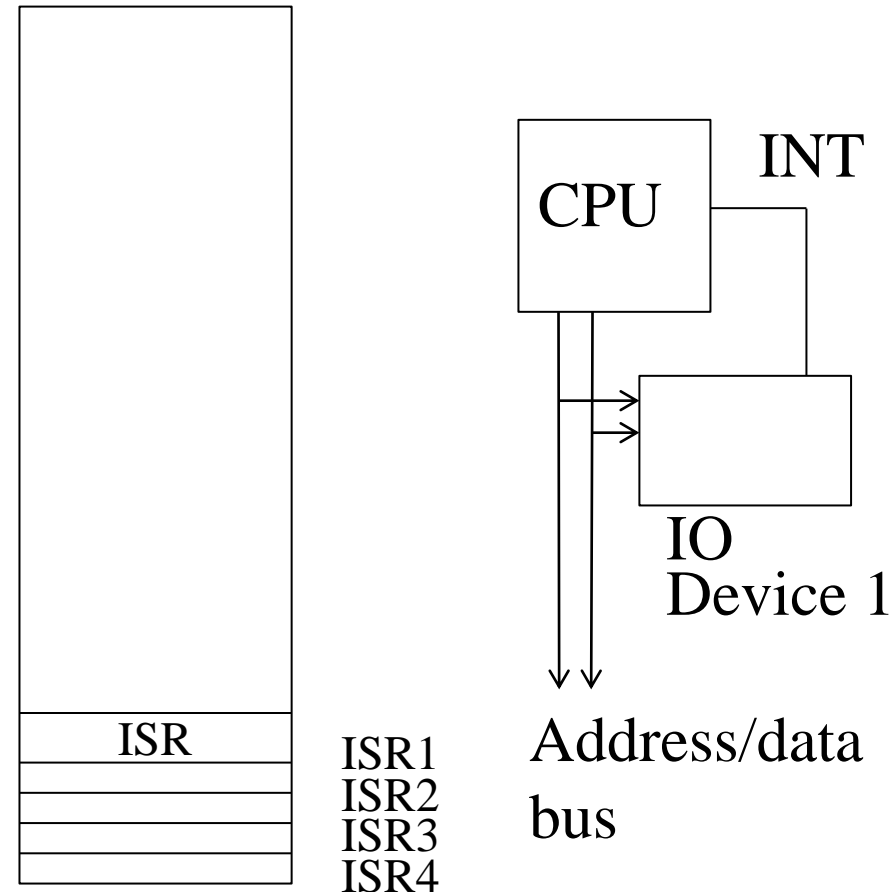
† I/O devices: many of them, slow, sporadic use

Programmed I/O vs. Interrupt

CPU: periodic IO check

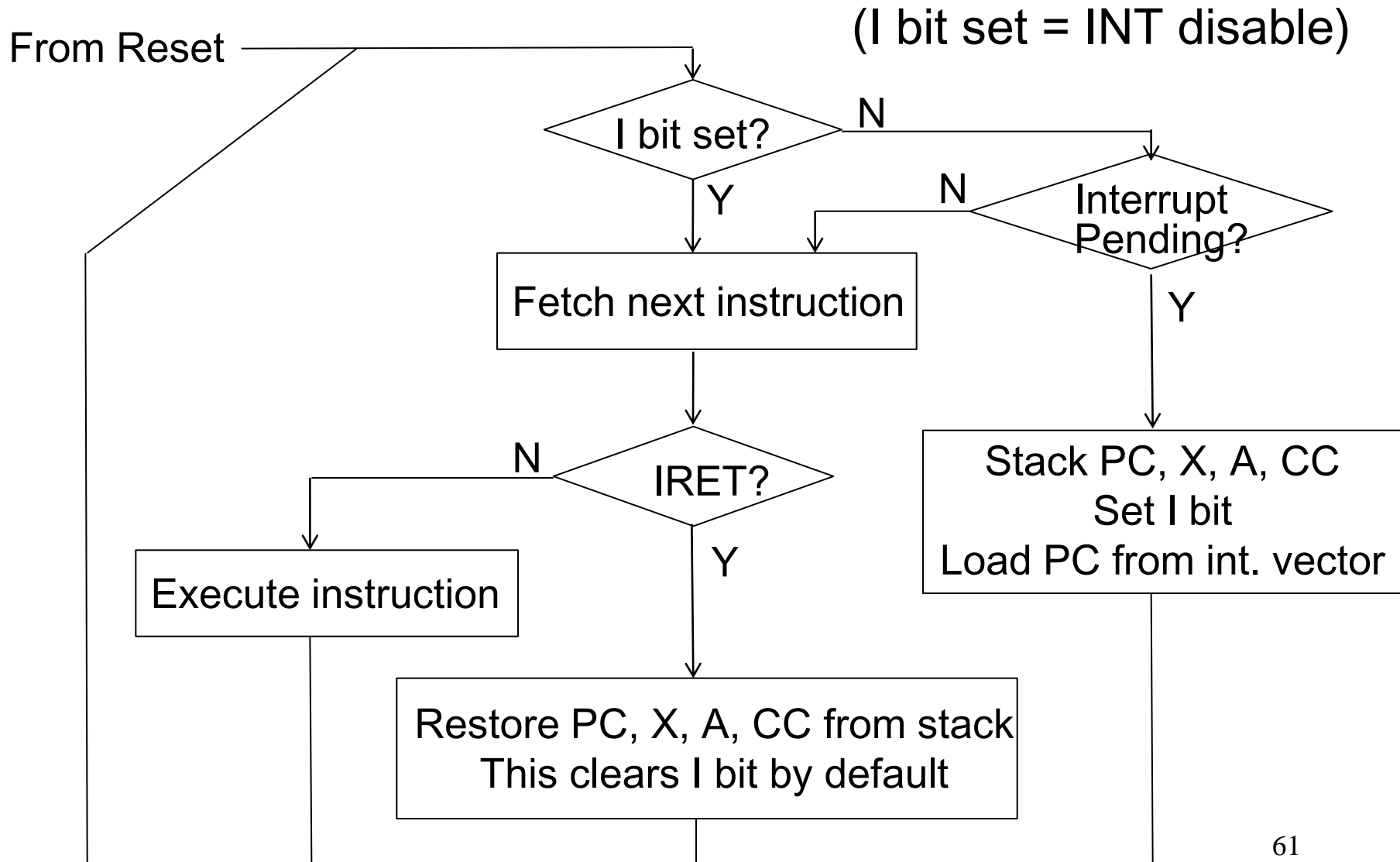


CPU: rely on INT hotline



- On accepting INT, jump to ⁶⁰ISR

Interrupt Processing (ST7 example)



Fetch-Execute and Interrupt

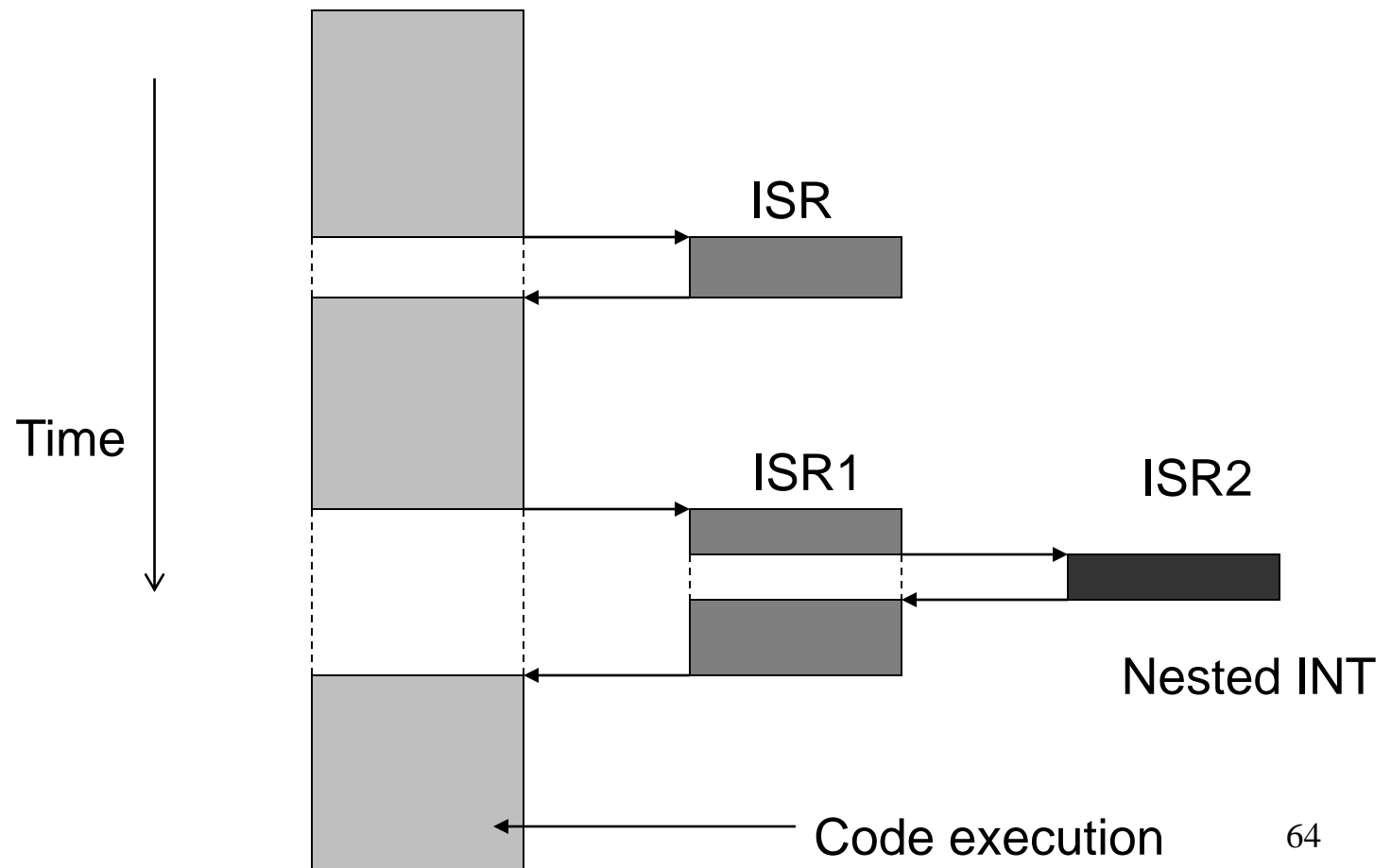
- ❑ Machine called computer
 - Fetch-decode-execute, adequate ISA, interrupt
 - † Special machine instructions: enable INT, disable INT
- ❑ Machine instruction
 - Atomic (all or nothing)
 - Interrupts checked after an instruction is finished
- ❑ Ticket reservation
 - Atomic
 - Locking, transaction

Related Terms

- ❑ Atomic operation (Mutex)
 - Timer interrupts
 - OS process scheduling
- ❑ Real-time systems (response time, deadline)
 - Hard
 - Soft
- ❑ RTOS (real-time OS)
 - Priority-based preemptive scheduling
 - † General-purpose OS: fairness

Time Diagram, Multiple INTs

- ❑ INT service routine (ISR)
- ❑ INT priority, INT vector



Summary

- ❑ Three-terminal digital switches (i.e., transistors)
- ❑ Semiconductor technology
- ❑ Intel and processor technology
- ❑ Memory technology and memory systems
- ❑ More on computer
 - x-bit computer, byte addressing, microcontroller, interconnection, interrupt