

Memory Hierarchy

Hamza Osman İLHAN

hoilhan@yildiz.edu.tr

YTU-CE D037

-
- Memory lies at the heart of the stored-program computer.
 - In this lecture, we focus on memory organization. A clear understanding of these ideas is essential for the analysis of system performance.

Memory Characteristics

Location

Processor
Internal (main)
External (secondary)

Capacity

Word size
Number of words

Unit of Transfer

Word
Block

Access Method

Sequential
Direct
Random
Associative

Performance

Access time
Cycle time
Transfer rate

Physical Type

Semiconductor
Magnetic
Optical
Magneto-Optical

Physical Characteristics

Volatile/nonvolatile
Erasable/nonerasable

Organization

Unit of Transfer

- Internal
 - Usually governed by data bus width
- External
 - Usually a block which is much larger than a word
- Addressable unit
 - Smallest location which can be uniquely addressed
 - Word internally
 - Cluster on M\$ disks

Access Methods (1)

- Sequential
 - Start at the beginning and read through in order
 - Access time depends on location of data and previous location
 - e.g. tape
- Direct
 - Individual blocks have unique address
 - Access is by jumping to vicinity plus sequential search
 - Access time depends on location and previous location
 - e.g. disk

Access Methods (2)

- Random
 - Individual addresses identify locations exactly
 - Access time is independent of location or previous access
 - e.g. RAM
- Associative
 - Data is located by a comparison with contents of a portion of the store
 - Access time is independent of location or previous access
 - e.g. cache

Performance

- Access time
 - Time between presenting the address and getting the valid data
- Memory Cycle time
 - Time may be required for the memory to “recover” before next access
 - Cycle time is access + recovery
- Transfer Rate
 - Rate at which data can be moved

Physical Types

- Semiconductor
 - RAM
- Magnetic
 - Disk & Tape
- Optical
 - CD & DVD
- Others
 - Bubble
 - Hologram

Physical Characteristics

- Decay
- Volatility
- Erasable
- Power consumption

Organisation

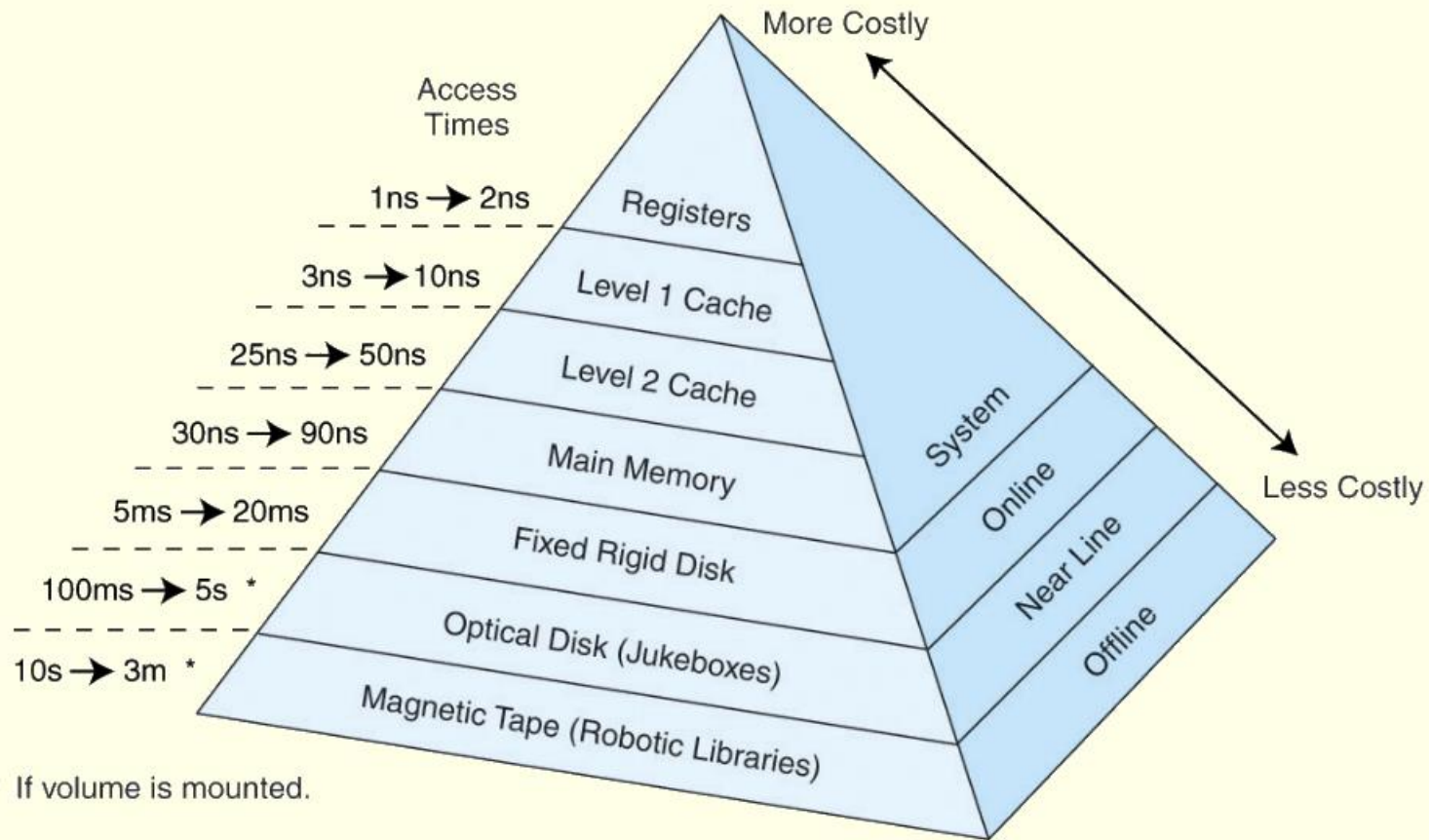
- Physical arrangement of bits into words
- Not always obvious
- e.g. interleaved

-
- Generally speaking, faster memory is more expensive than slower memory.
 - To provide the best performance at the lowest cost, memory is organized in a hierarchical fashion.
 - Small, fast storage elements are kept in the CPU, larger, slower main memory is accessed through the data bus.
 - Larger, (almost) permanent storage in the form of disk and tape drives is still further from the CPU.

Memory Hierarchy

- Registers
 - In CPU
- Internal or Main memory
 - May include one or more levels of cache
 - “RAM”
- External memory
 - Backing store

- This storage organization can be thought of as a pyramid:



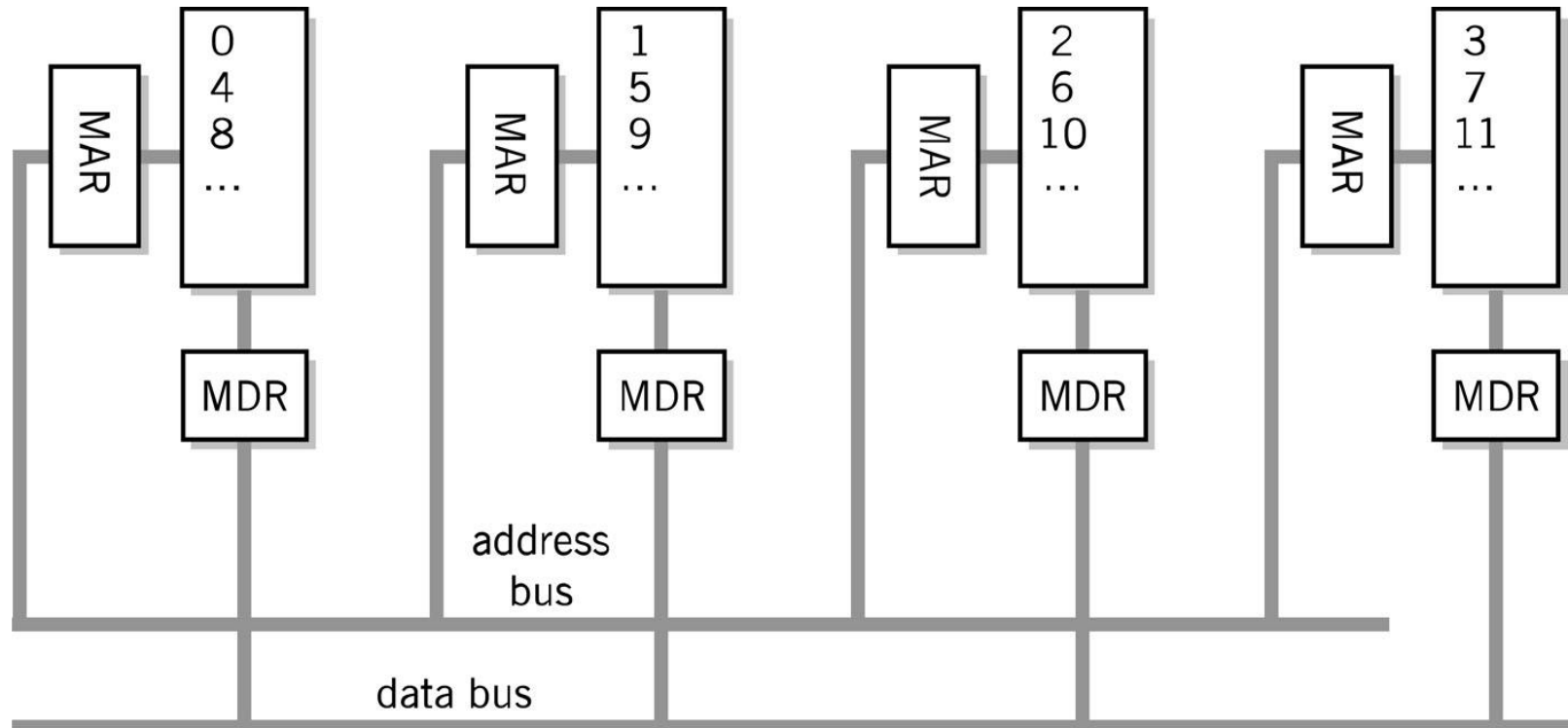
Hierarchy List

- Registers
- L1 Cache
- L2 Cache
- Main memory
- Disk cache
- Disk
- Optical
- Tape

Memory Enhancements

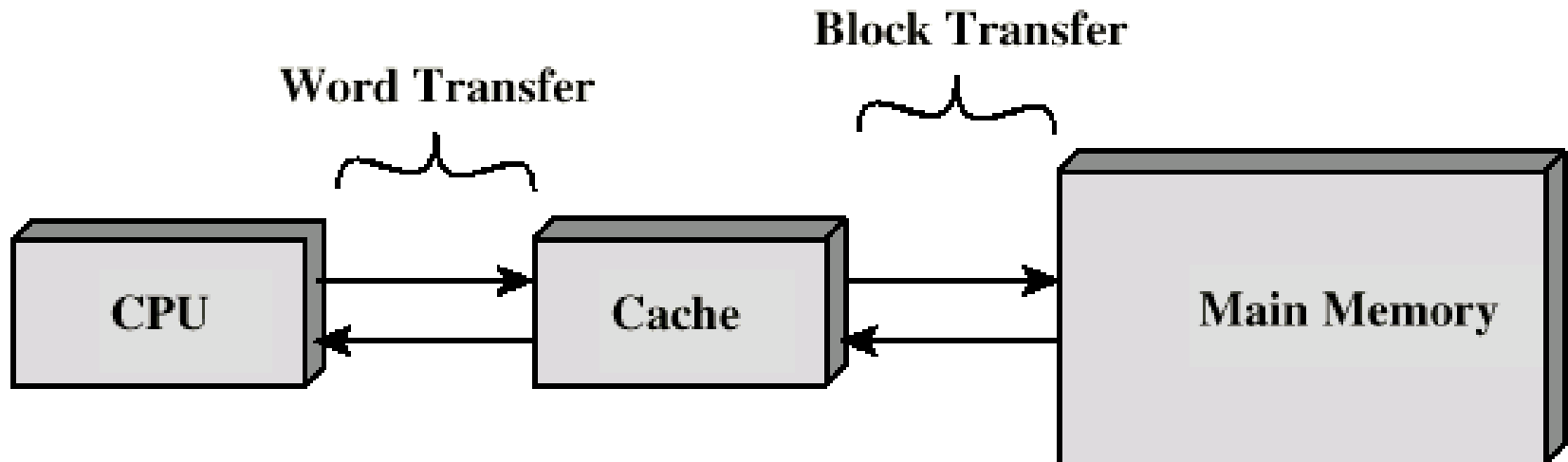
- Memory is slow compared to CPU processing speeds!
 - 2Ghz CPU = 1 cycle in $\frac{1}{2}$ of a billionth of a second
 - 70ns DRAM = 1 access in 70 millionth of a second
- Methods to improvement memory accesses
 - Wide Path Memory Access
 - Retrieve multiple bytes instead of 1 byte at a time
 - Memory Interleaving
 - Partition memory into subsections, each with its own address register and data register
 - Cache Memory

Memory Interleaving



Cache

- Small amount of fast memory
- Sits between normal main memory and CPU
- May be located on CPU chip or module



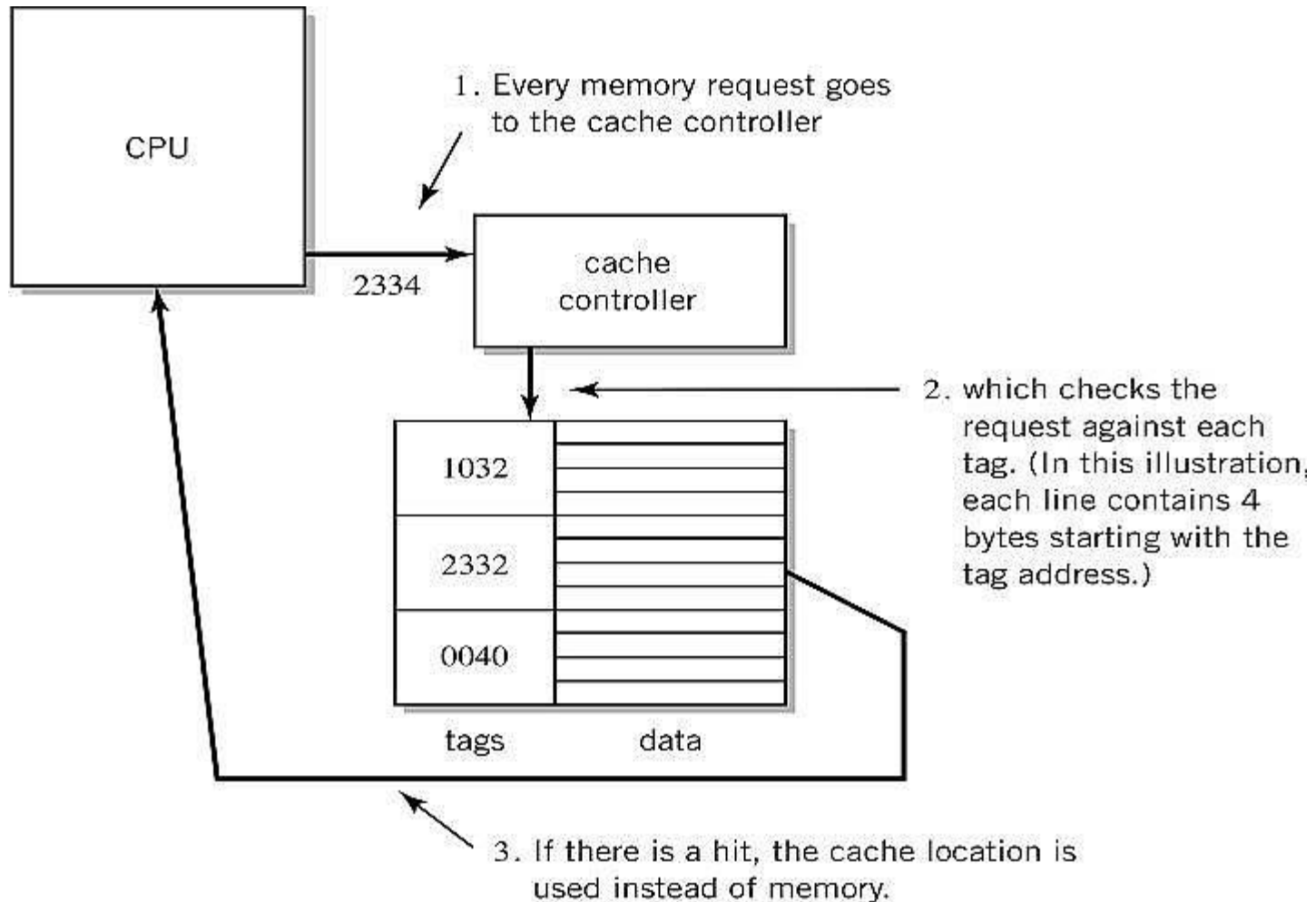
Why Cache?

- Even the fastest hard disk has an access time of about 10 milliseconds
- 2Ghz CPU waiting 10 milliseconds **wastes 20 million clock cycles!**

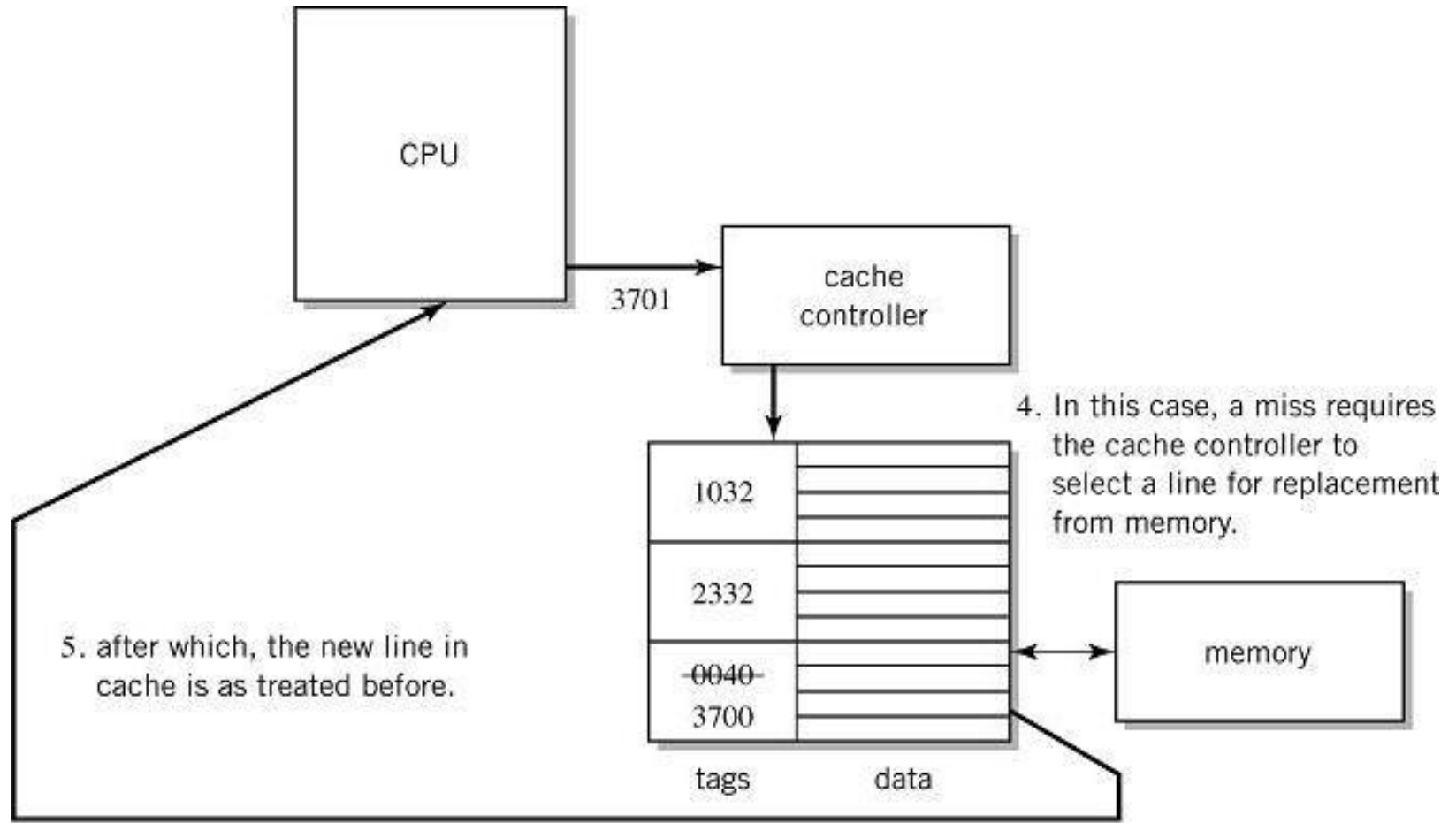
Cache Memory

- Blocks: 8 or 16 bytes
- Tags: location in main memory
- Cache controller
 - hardware that checks tags
- Cache Line
 - Unit of transfer between storage and cache memory
- Hit Ratio: ratio of hits out of total requests
- Synchronizing cache and memory
 - Write through
 - Write back

Step-by-Step Use of Cache



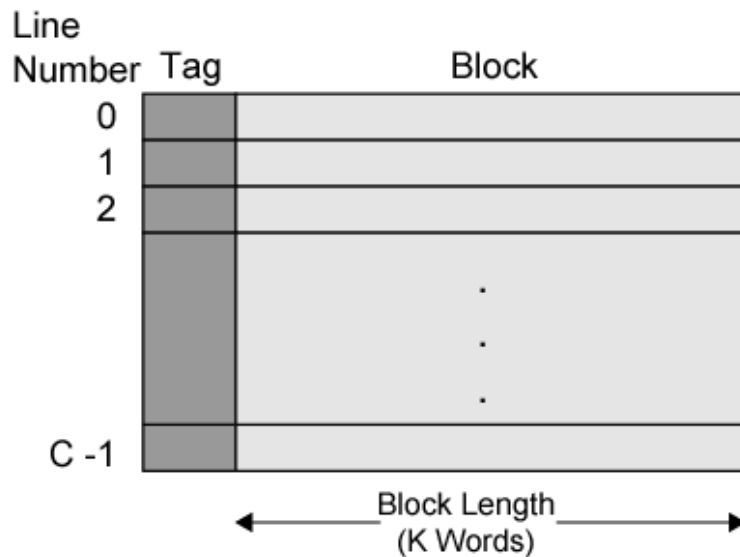
Step-by-Step Use of Cache



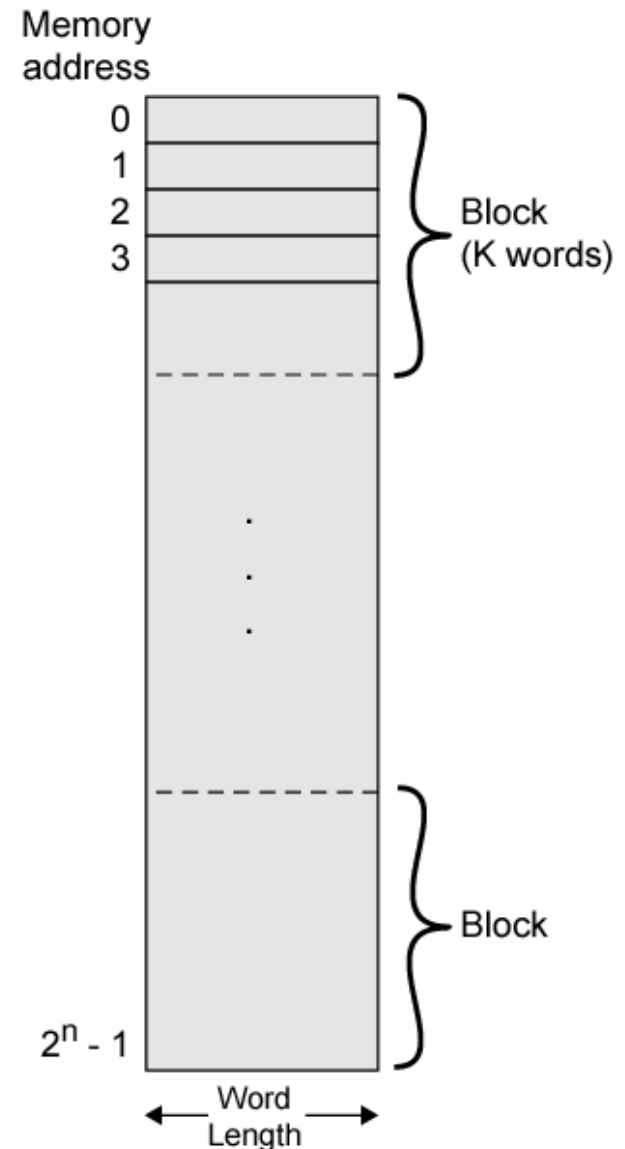
Cache vs. Virtual Memory

- Cache **speeds** up memory access
- Virtual memory increases amount of **perceived storage**
 - independence from the configuration and capacity of the memory system
 - low cost per bit

Cache/Main Memory Structure



(a) Cache



(b) Main memory

- Main memory size: upto 2^n words
- Each word has a unique n -bit address
- Fixed length blocks of K words each
- Number of blocks: $M=2^n/K$
- Cache consists of C lines
- Each line contains K words + tag
- $C \ll M$

Cache operation – overview

- CPU requests contents of memory location
- Check cache for this data
- If present, get from cache (fast)
- If not present, read required block from main memory to cache
- Then deliver from cache to CPU
- Cache includes tags to identify which block of main memory is in each cache slot

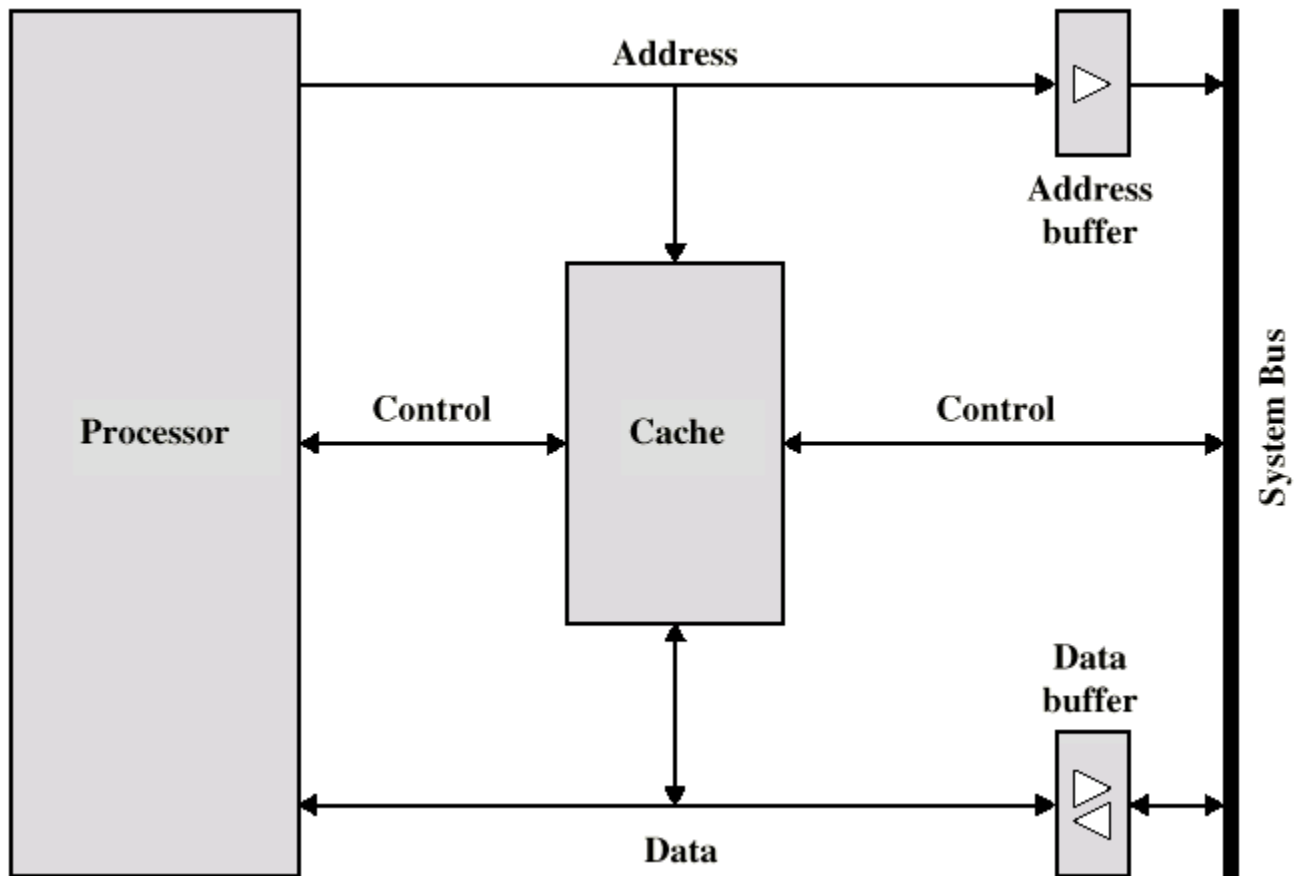
Cache Design

- Size
- Mapping Function
- Replacement Algorithm
- Write Policy
- Block Size
- Number of Caches

Size does matter

- Cost
 - More cache is expensive
- Speed
 - More cache is faster (up to a point)
 - Checking cache for data takes time

Typical Cache Organization



Comparison of Cache Sizes

Processor	Type	Year of Introduction	L1 cache ^a	L2 cache	L3 cache
IBM 360/85	Mainframe	1968	16 to 32 KB	—	—
PDP-11/70	Minicomputer	1975	1 KB	—	—
VAX 11/780	Minicomputer	1978	16 KB	—	—
IBM 3033	Mainframe	1978	64 KB	—	—
IBM 3090	Mainframe	1985	128 to 256 KB	—	—
Intel 80486	PC	1989	8 KB	—	—
Pentium	PC	1993	8 KB/8 KB	256 to 512 KB	—
PowerPC 601	PC	1993	32 KB	—	—
PowerPC 620	PC	1996	32 KB/32 KB	—	—
PowerPC G4	PC/server	1999	32 KB/32 KB	256 KB to 1 MB	2 MB
IBM S/390 G4	Mainframe	1997	32 KB	256 KB	2 MB
IBM S/390 G6	Mainframe	1999	256 KB	8 MB	—
Pentium 4	PC/server	2000	8 KB/8 KB	256 KB	—
IBM SP	High-end server/ supercomputer	2000	64 KB/32 KB	8 MB	—
CRAY MTA ^b	Supercomputer	2000	8 KB	2 MB	—
Itanium	PC/server	2001	16 KB/16 KB	96 KB	4 MB
SGI Origin 2001	High-end server	2001	32 KB/32 KB	4 MB	—
Itanium 2	PC/server	2002	32 KB	256 KB	6 MB
IBM POWER5	High-end server	2003	64 KB	1.9 MB	36 MB
CRAY XD-1	Supercomputer	2004	64 KB/64 KB	1MB	—

Mapping Function

- Because there are fewer lines than main memory blocks, an algorithm is needed for mapping main memory blocks into cache lines.
- Which main memory block currently occupies a cache line?
- Three techniques can be used:
 - Direct mapping
 - Associative mapping
 - Set associative mapping

Replacement Algorithms (1)

Direct mapping

- No choice
- Each block only maps to one line
- Replace that line

Replacement Algorithms (2)

Associative & Set Associative

- Hardware implemented algorithm (speed)
- Least Recently used (LRU)
 - e.g. in 2 way set associative
 - Which of the 2 block is lru?
- First in first out (FIFO)
 - replace block that has been in cache longest
- Least frequently used
 - replace block which has had fewest hits
- Random

Write Policy

- Must not overwrite a cache block unless main memory is up to date
- Multiple CPUs may have individual caches
- I/O may address main memory directly

Write through

- All writes go to main memory as well as cache
- Multiple CPUs can monitor main memory traffic to keep local (to CPU) cache up to date
- Lots of traffic
- Slows down writes

Write back

- Updates initially made in cache only
- Update bit for cache slot is set when update occurs
- If block is to be replaced, write to main memory only if update bit is set
- Other caches get out of sync
- I/O must access main memory through cache
- N.B. 15% of memory references are writes

Pentium 4 Cache

- 80386 – no on chip cache
- 80486 – 8k using 16 byte lines and four way set associative organization
- Pentium (all versions) – two on chip L1 caches
 - Data & instructions
- Pentium III – L3 cache added off chip
- Pentium 4
 - L1 caches
 - 8k bytes
 - 64 byte lines
 - four way set associative
 - L2 cache
 - Feeding both L1 caches
 - 256k
 - 128 byte lines
 - 8 way set associative
 - L3 cache on chip

PowerPC Cache Organization

- 601 – single 32kb 8 way set associative
- 603 – 16kb (2 x 8kb) two way set associative
- 604 – 32kb
- 620 – 64kb
- G3 & G4
 - 64kb L1 cache
 - 8 way set associative
 - 256k, 512k or 1M L2 cache
 - two way set associative
- G5
 - 32kB instruction cache
 - 64kB data cache