# Procesory i Architektura Systemów Komputerowych

## Pamięć podręczna (cache)

**IET**
**Katedra Elektroniki**
**Kraków 2015**
**dr inż. Roman Rumian**

**temporal locality** The principle stating that if a data location is referenced then it will tend to be referenced again soon.

**spatial locality** The locality principle stating that if a data location is referenced, data locations with nearby addresses will tend to be referenced soon.

**memory hierarchy** A structure that uses multiple levels of memories; as the distance from the processor increases, the size of the memories and the access time both increase.

**block (or line)** Th e minimum unit of information that can be either present or not present in a cache.
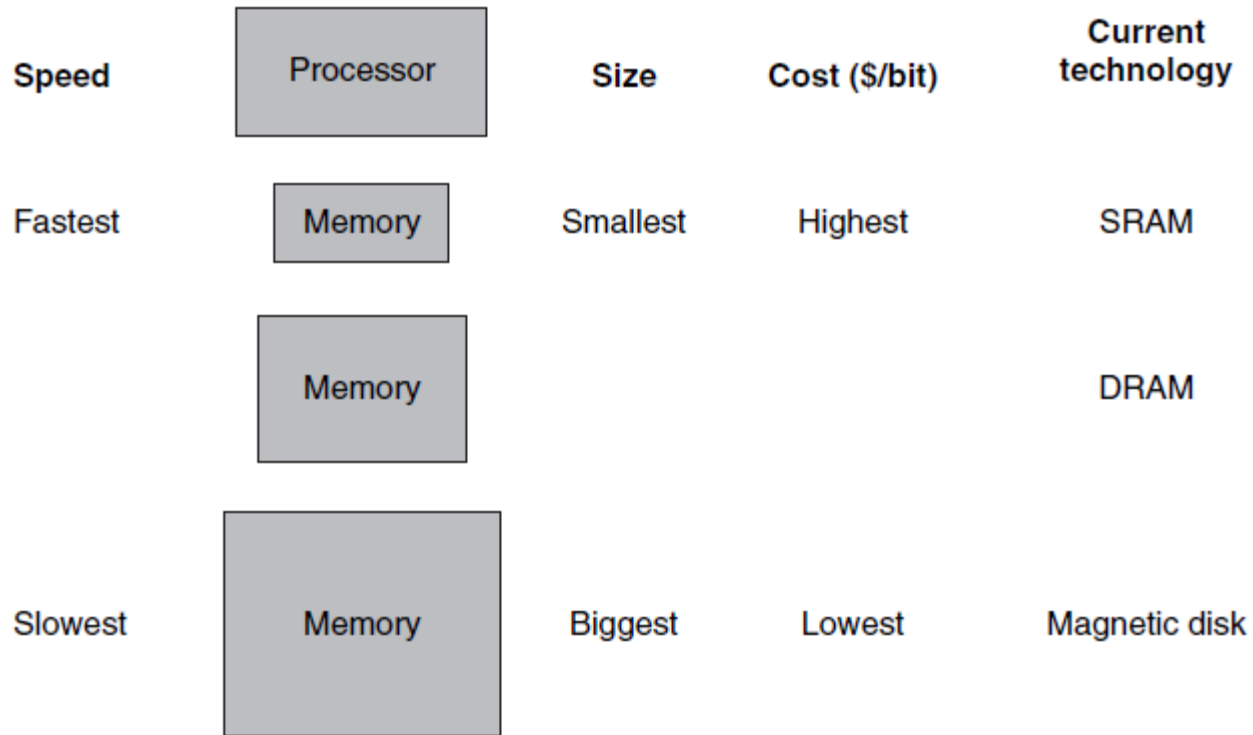
**hit rate** The fraction of memory accesses found in a level of the memory hierarchy.

**miss rate** The fraction of memory accesses not found in a level of the memory hierarchy.

**hit time** The time required to access a level of the memory hierarchy, including the time needed to determine whether the access is a hit or a miss.
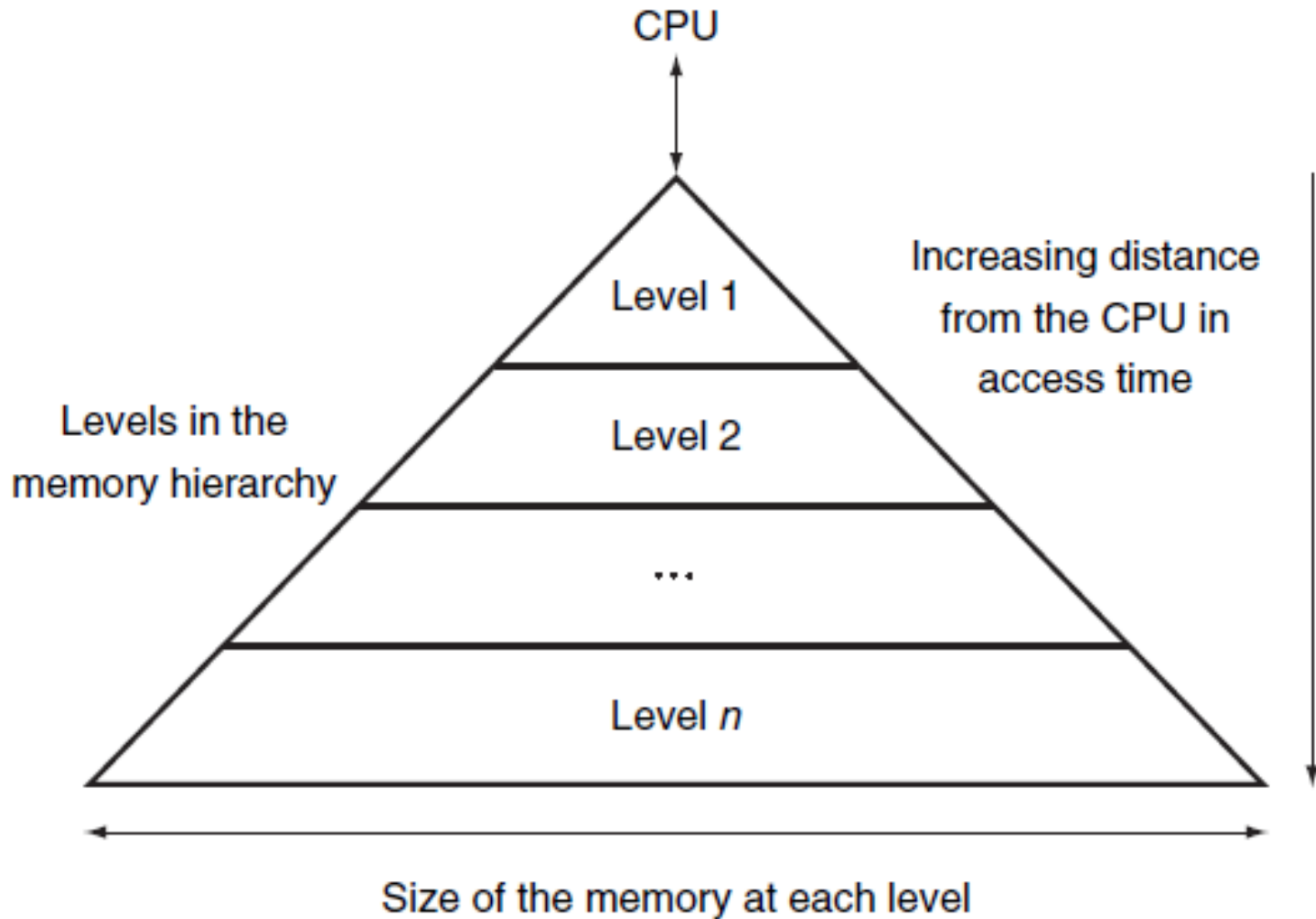
**miss penalty** The time required to fetch a block into a level of the memory hierarchy from the lower level, including the time to access the block, transmit it from one level to the other, insert it in the level that experienced the miss, and then pass the block to the requestor.
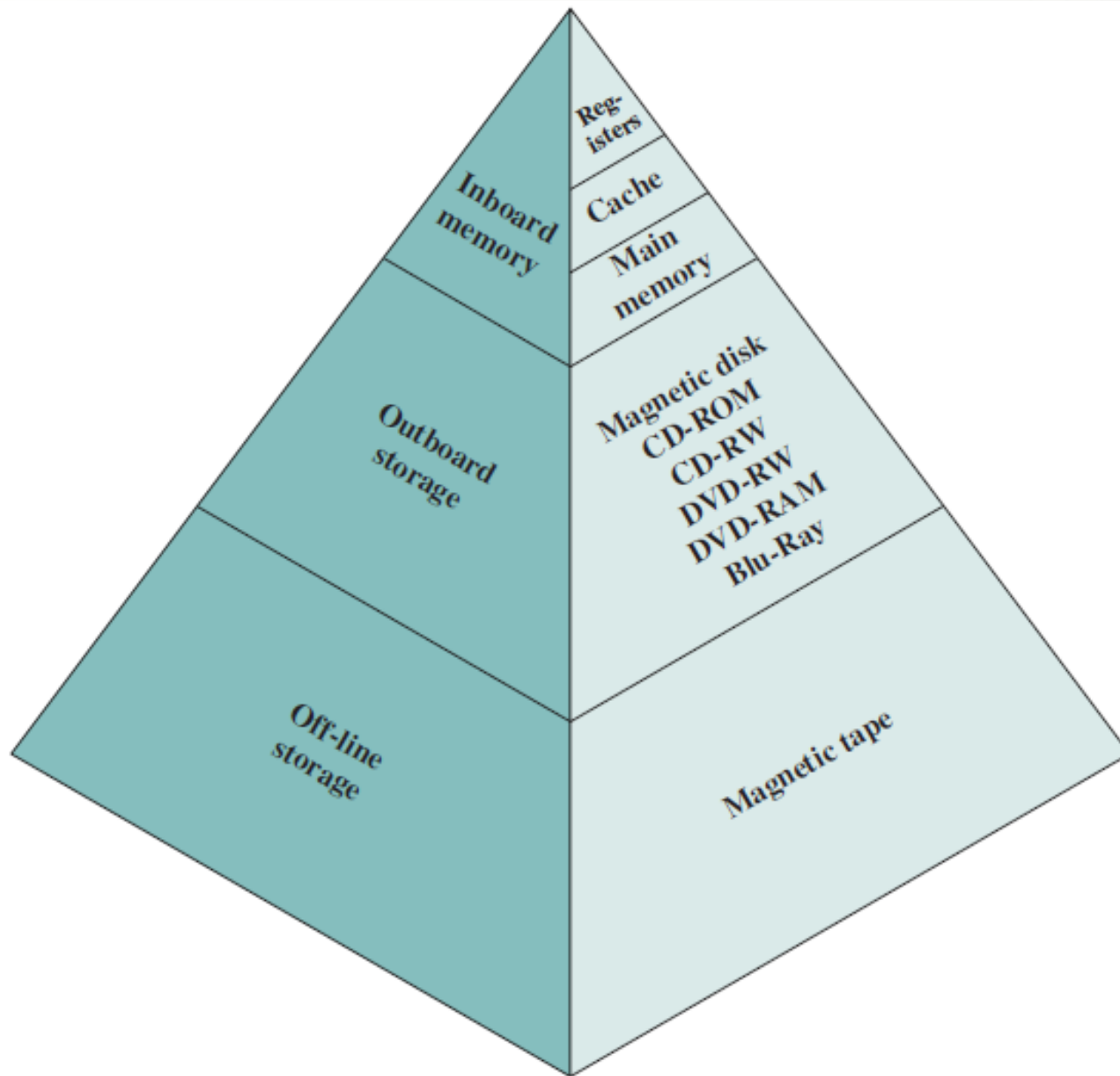
# The basic structure of a memory hierarchy

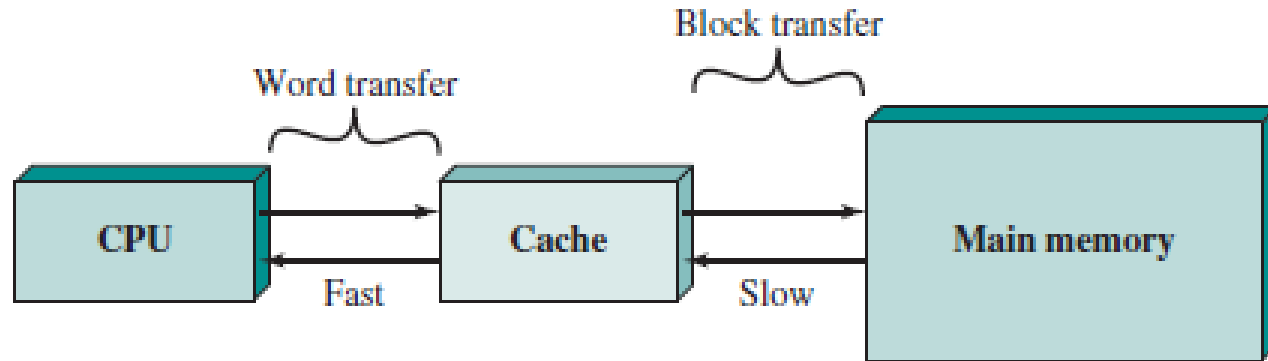| Speed | | Size | Cost ($/bit) | Current technology |
|---|---|---|---|---|
| | Processor | | | |
| Fastest | Memory | Smallest | Highest | SRAM |
| | Memory | | | DRAM |
| Slowest | Memory | Biggest | Lowest | Magnetic disk |

By implementing the memory system as a hierarchy, the user has the illusion of a memory that is as large as the largest level of the hierarchy, but can be accessed as if it were all built from the fastest memory. Flash memory has replaced disks in many personal mobile devices, and may lead to a new level in the storage hierarchy for desktop and server computers

# The memory hierarchy

CPU

Level 1

Level 2

...

Level n

Increasing distance
from the CPU in
access time

Levels in the
memory hierarchy

Size of the memory at each level

The memory hierarchy pyramid:

Inboard memory:
- Registers
- Cache
- Main memory

Outboard storage:
- Magnetic disk
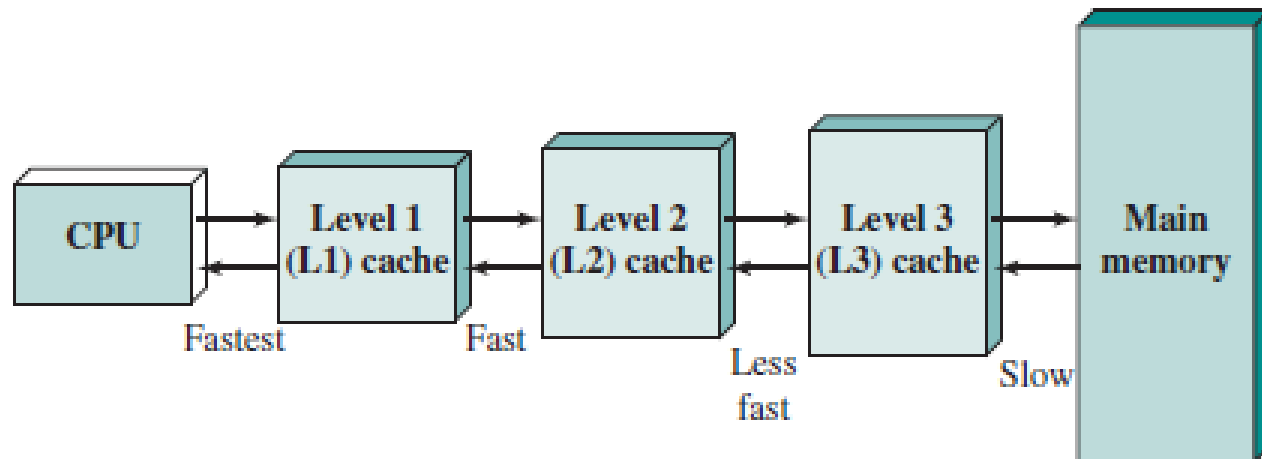- CD-ROM
- CD-RW
- DVD-RW
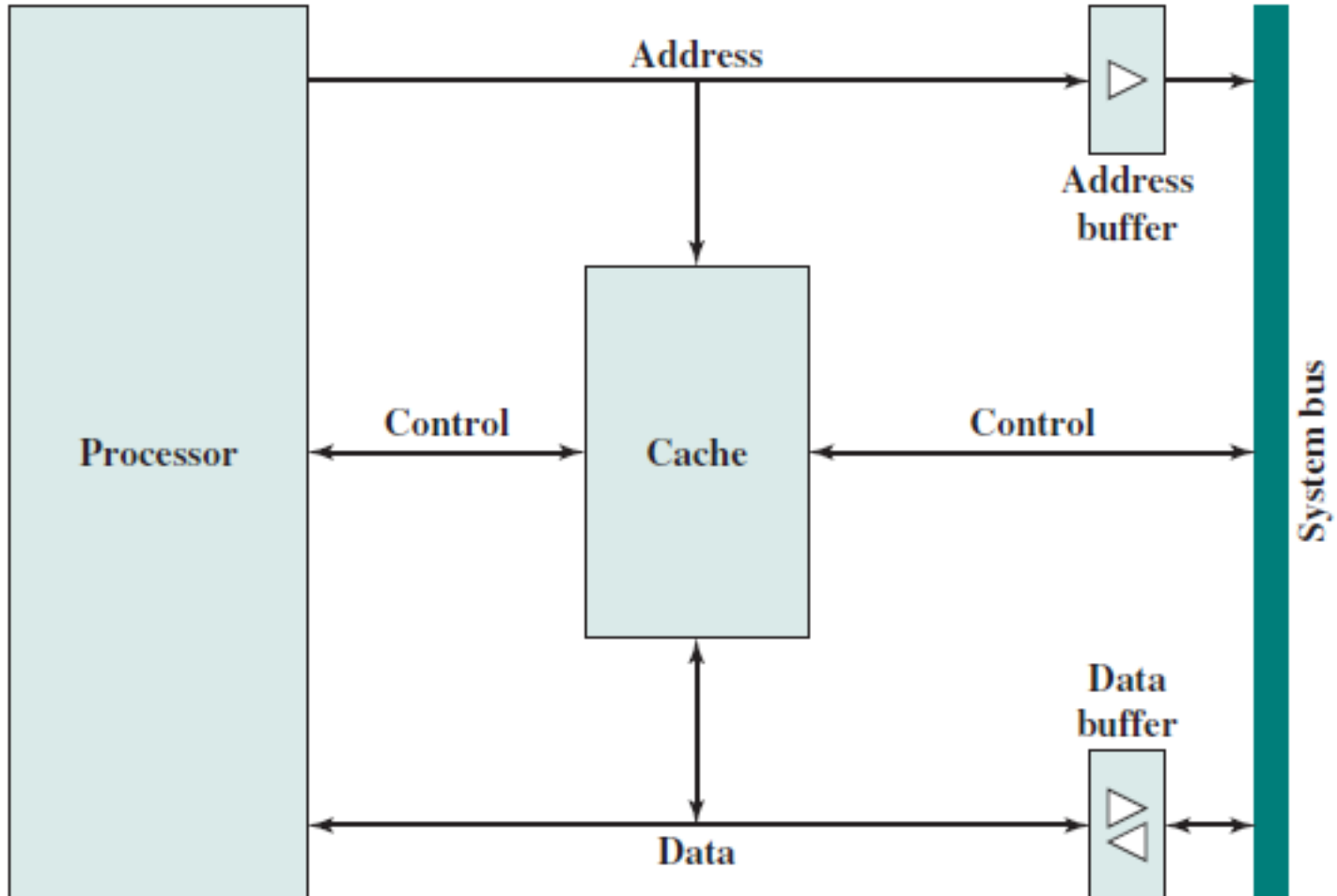- DVD-RAM
- Blu-Ray

Off-line storage:
- Magnetic tape

# Cache and Main Memory



(a) Single cache
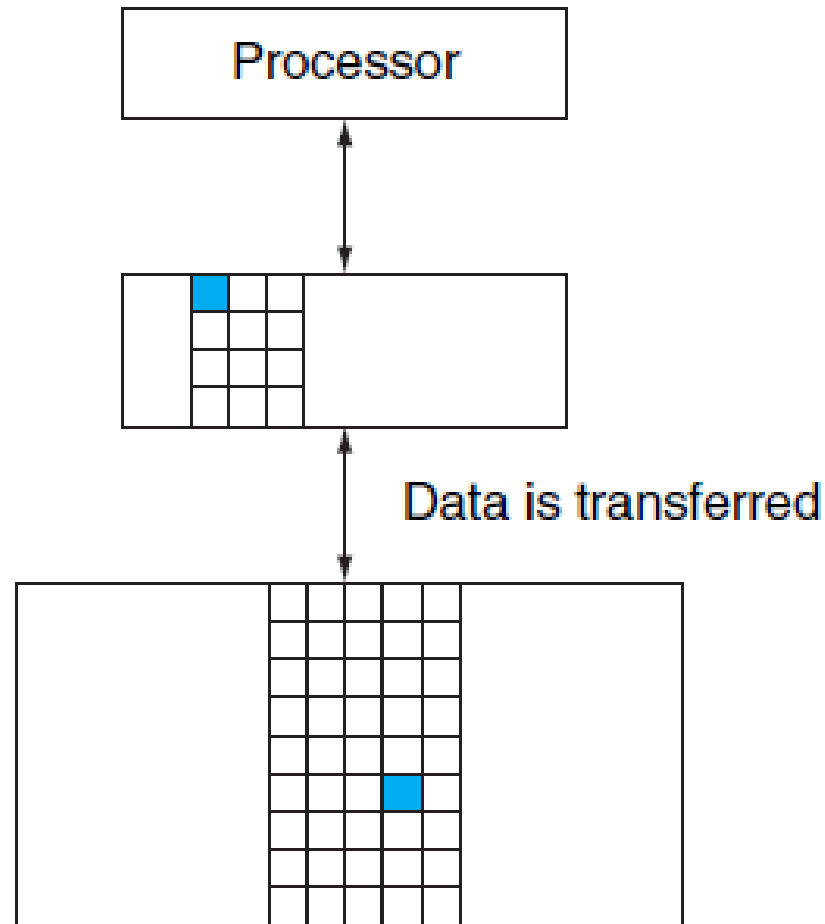
(b) Three-level cache organization

**Every pair of levels in the memory hierarchy can be thought of as having an upper and lower level.** Within each level, the unit of information that is present or not is called a *block* or a *line*. Usually we transfer an entire block when we copy something between levels
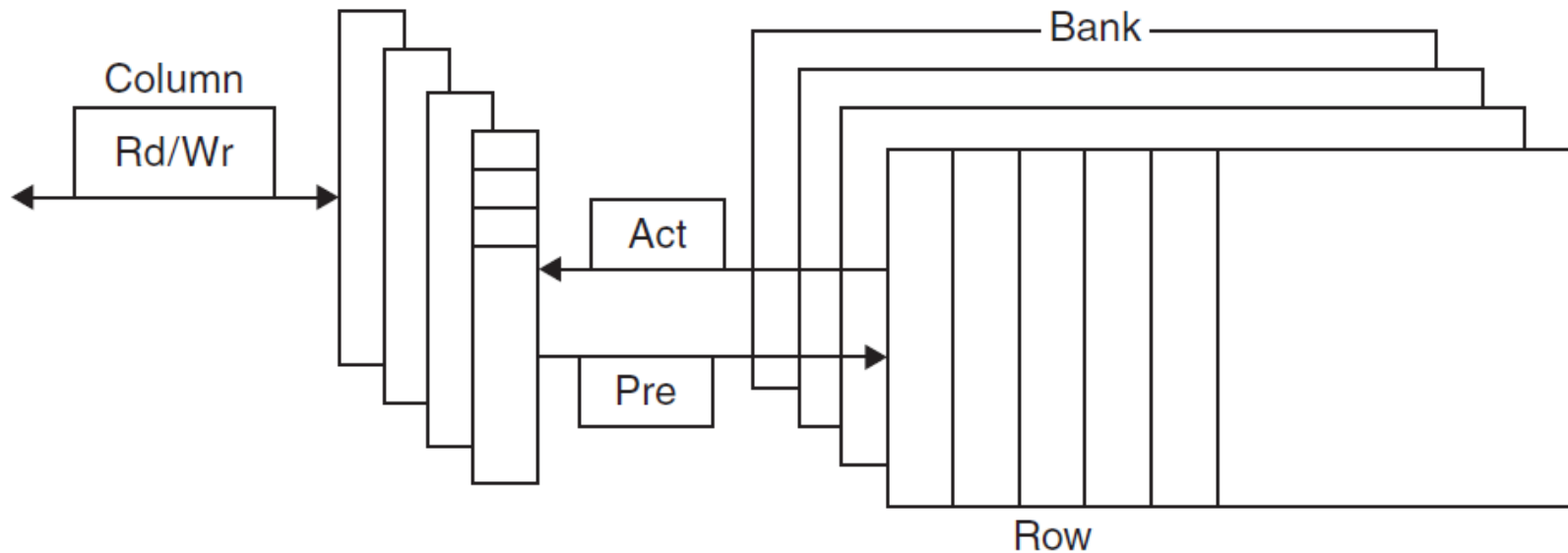


Processor

Data is transferred

# Memory Technologies

| Memory technology | Typical access time | $ per GiB in 2012 |
|---|---|---|
| SRAM semiconductor memory | 0.5–2.5 ns | $500–$1000 |
| DRAM semiconductor memory | 50–70 ns | $10–$20 |
| Flash semiconductor memory | 5,000–50,000 ns | $0.75–$1.00 |
| Magnetic disk | 5,000,000–20,000,000 ns | $0.05–$0.10 |

| Year Introduced | Chip size | $ per GIB | Total access time to a new row/column | Average column access time to existing row |
|---|---|---|---|---|
| 1980 | 64 Kibibit | $1,500,000 | 250 ns | 150 ns |
| 1983 | 256 Kibibit | $500,000 | 185 ns | 100 ns |
| 1985 | 1 Mebibit | $200,000 | 135 ns | 40 ns |
| 1989 | 4 Mebibit | $50,000 | 110 ns | 40 ns |
| 1992 | 16 Mebibit | $15,000 | 90 ns | 30 ns |
| 1996 | 64 Mebibit | $10,000 | 60 ns | 12 ns |
| 1998 | 128 Mebibit | $4,000 | 60 ns | 10 ns |
| 2000 | 256 Mebibit | $1,000 | 55 ns | 7 ns |
| 2004 | 512 Mebibit | $250 | 50 ns | 5 ns |
| 2007 | 1 Gibibit | $50 | 45 ns | 1.25 ns |
| 2010 | 2 Gibibit | $30 | 40 ns | 1 ns |
| 2012 | 4 Gibibit | $1 | 35 ns | 0.8 ns |

**DRAM size increased by multiples of four approximately once every three years until 1996, and thereafter considerably slower**
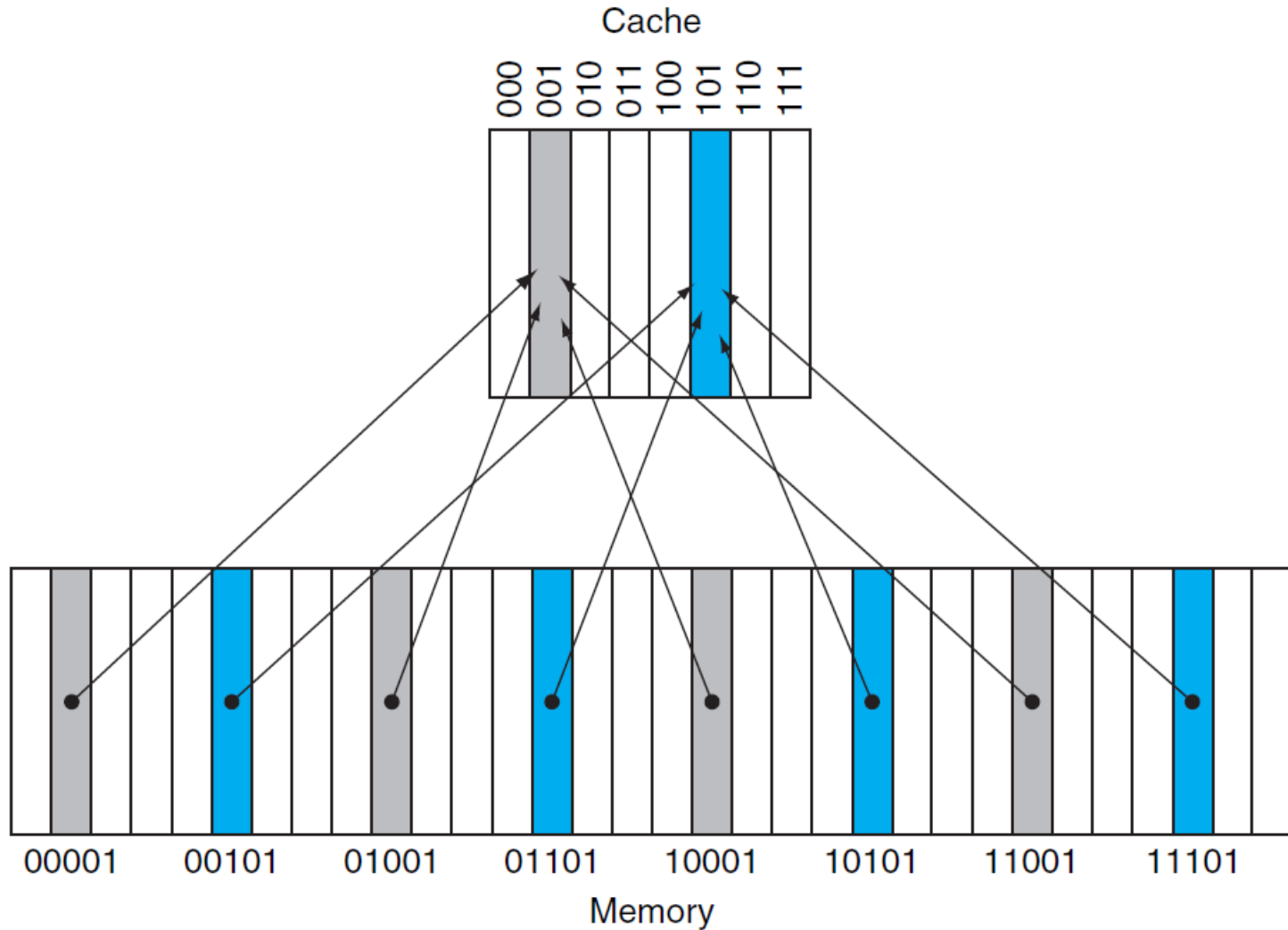
# Internal organization of a DRAM.



Modern DRAMs are organized in banks, typically four for DDR3. Each bank consists of a series of rows. Sending a PRE (precharge) command opens or closes a bank. A row address is sent with an Act (activate), which causes the row to transfer to a buffer. When the row is in the buffer, it can be transferred by successive column addresses at whatever the width of the DRAM is (typically 4, 8, or 16 bits in DDR3) or by specifying a block transfer and the starting address. Each command, as well as block transfers, is synchronized with a clock.
A DDR4-3200 DRAM can do 3200 million transfers per second, which means it has a 1600 MHz clock.
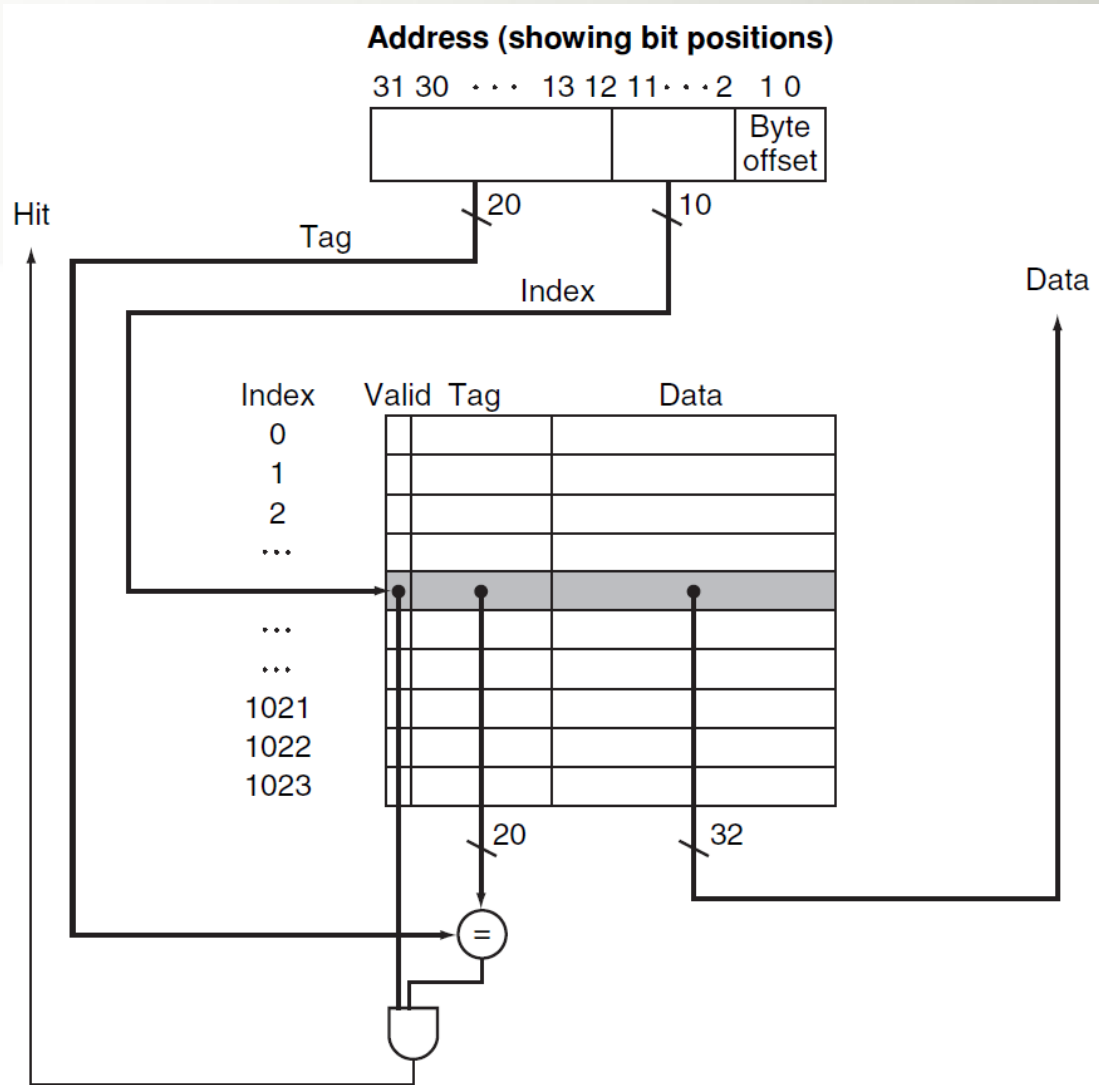
# direct-mapped cache

A cache structure in which each memory location is mapped to exactly one location in the cache.

**tag** A field in a table used for a memory hierarchy that contains the address information required to identify whether the associated block in the hierarchy corresponds to a requested word.

**valid bit** A field in the tables of a memory hierarchy that indicates that the associated block in the hierarchy contains valid data.

Caching is perhaps the most important example of the big idea of **prediction**. It relies on the principle of locality to try to find the desired data in the higher levels of the memory hierarchy, and provides mechanisms to ensure that when the prediction is wrong it finds and uses the proper data from the lower levels of the memory hierarchy. The hit rates of the cache prediction on modern computers are often higher than 95%.
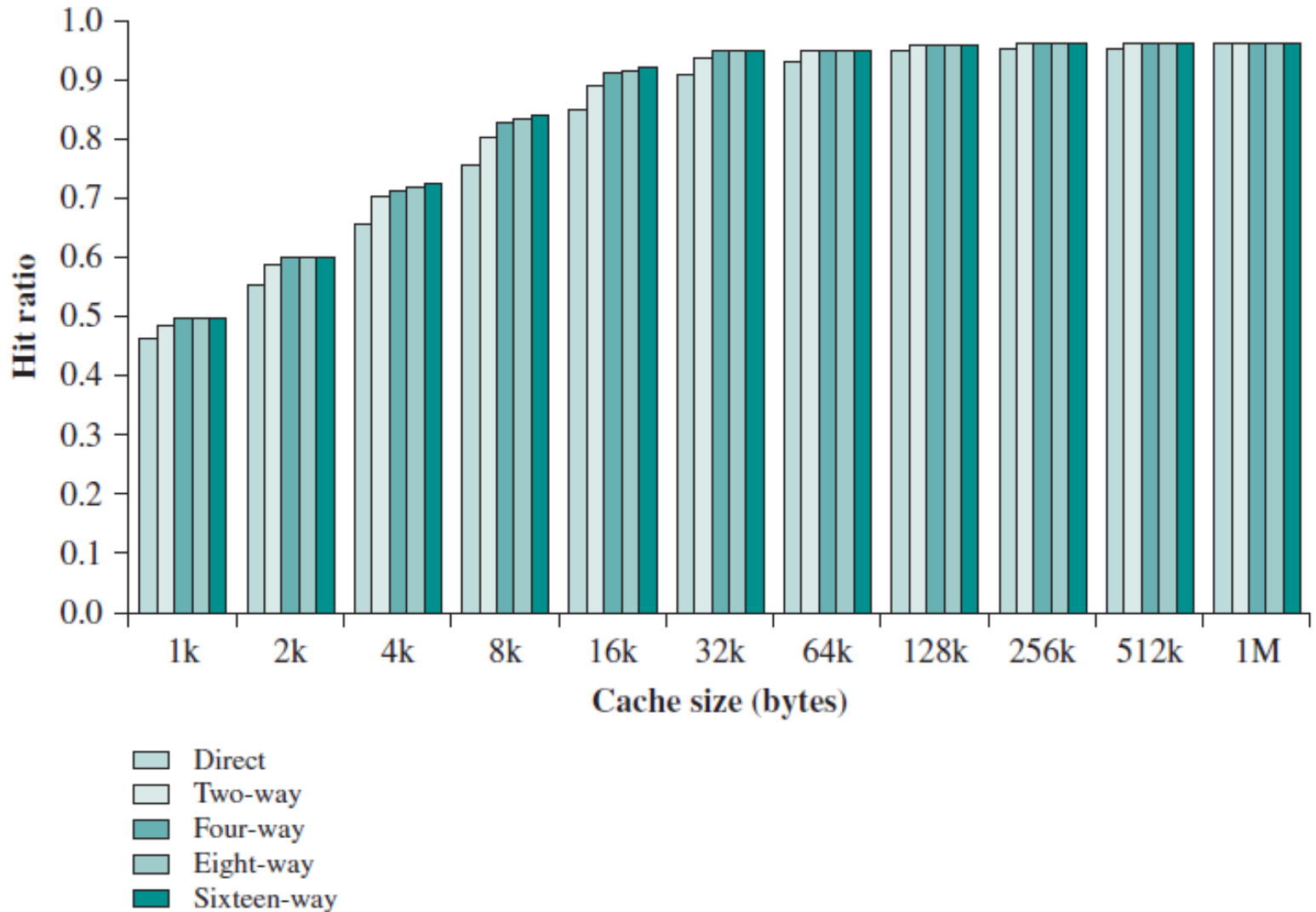
**Address (showing bit positions)**

31 30 · · · 13 12 11 · · · 2   1 0

Byte offset

Hit

Tag

20  10

Index

Data

Index  Valid  Tag  Data
0
1
2
...
...
...
1021
1022
1023

20  32

=

**For this cache, the lower portion of the address is used to select a cache entry consisting of a data word and a tag.** This cache holds 1024 words or 4 KiB. We assume 32-bit addresses in this chapter. The tag from the cache is compared against the upper portion of the address to determine whether the entry in the cache corresponds to the requested address. Because the cache has $2_{10}$ (or 1024) words and a block size of one word, 10 bits are used to index the cache, leaving $32 - 10 - 2 = 20$ bits to be compared against the tag. If the tag and upper 20 bits of the address are equal and the valid bit is on, then the request hits in the cache, and the word is supplied to the processor. Otherwise, a miss occurs.
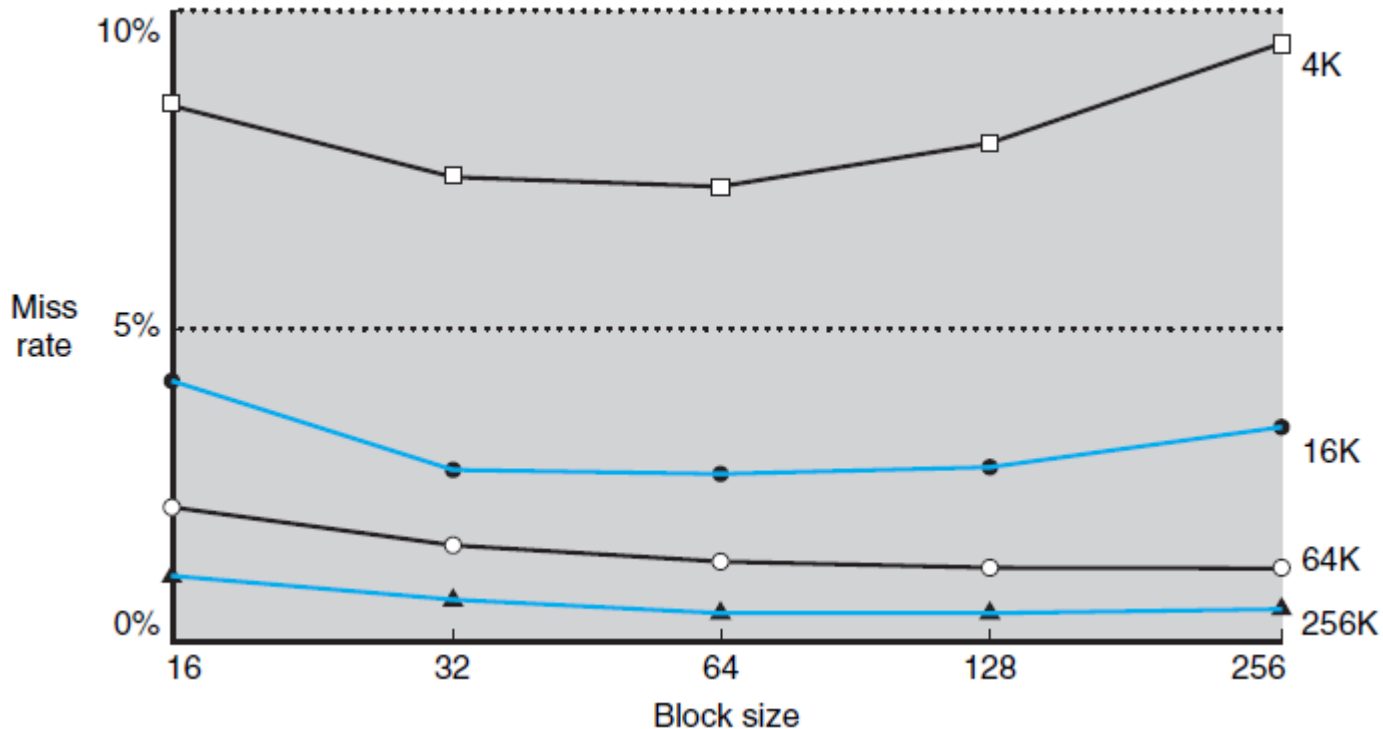
For the associative and set-associative techniques, a replacement algorithm is needed. To achieve high speed, such an algorithm must be implemented in hardware. A number of algorithms have been tried. We mention four of the most common. Probably the most effective is least recently used (LRU): Replace that block in the set that has been in the cache longest with no reference to it. For two-way set associative, this is easily implemented. Each line includes a USE bit. When a line is referenced, its USE bit is set to 1 and the USE bit of the other line in that set is set to 0. When a block is to be read into the set, the line whose USE bit is 0 is used. Because we are assuming that more recently used memory locations are more likely to be referenced, LRU should give the best hit ratio. LRU is also relatively easy to implement for a fully associative cache. The cache mechanism maintains a separate list of indexes to all the lines in the cache. When a line is referenced, it moves to the front of the list. For replacement, the line at the back of the list is used. Because of its simplicity of implementation, LRU is the most popular replacement algorithm.
Another possibility is first-in-first-out (FIFO): Replace that block in the set
that has been in the cache longest. FIFO is easily implemented as a round-robin or circular buffer technique. Still another possibility is least frequently used (LFU):
Replace that block in the set that has experienced the fewest references. LFU could be implemented by associating a counter with each line. A technique not based on usage (i.e., not LRU, LFU, FIFO, or some variant) is to pick a line at random from among the candidate lines. Simulation studies have shown that random replacement provides only slightly inferior performance to an algorithm based on usage

# Varying Associativity over Cache Size

# Miss rate versus block size.



Note that the miss rate actually goes up if the block size is too large relative to the cache size. Each line represents a cache of different size. (This figure is independent of associativity, discussed soon.) Unfortunately, SPEC CPU2000 traces would take too long if block size were included, so this data is based on SPEC92.

**write-through**
A scheme in which writes always update both the cache and the next lower level of the memory hierarchy, ensuring that data is always consistent between the two.

**write-back**
A scheme that handles writes by updating values only to the block in the cache, then writing the modified block to the lower level of the hierarchy when the block is replaced.

**write buffer**
A queue that holds data while the data is waiting to be written to memory.

**fully associative cache**
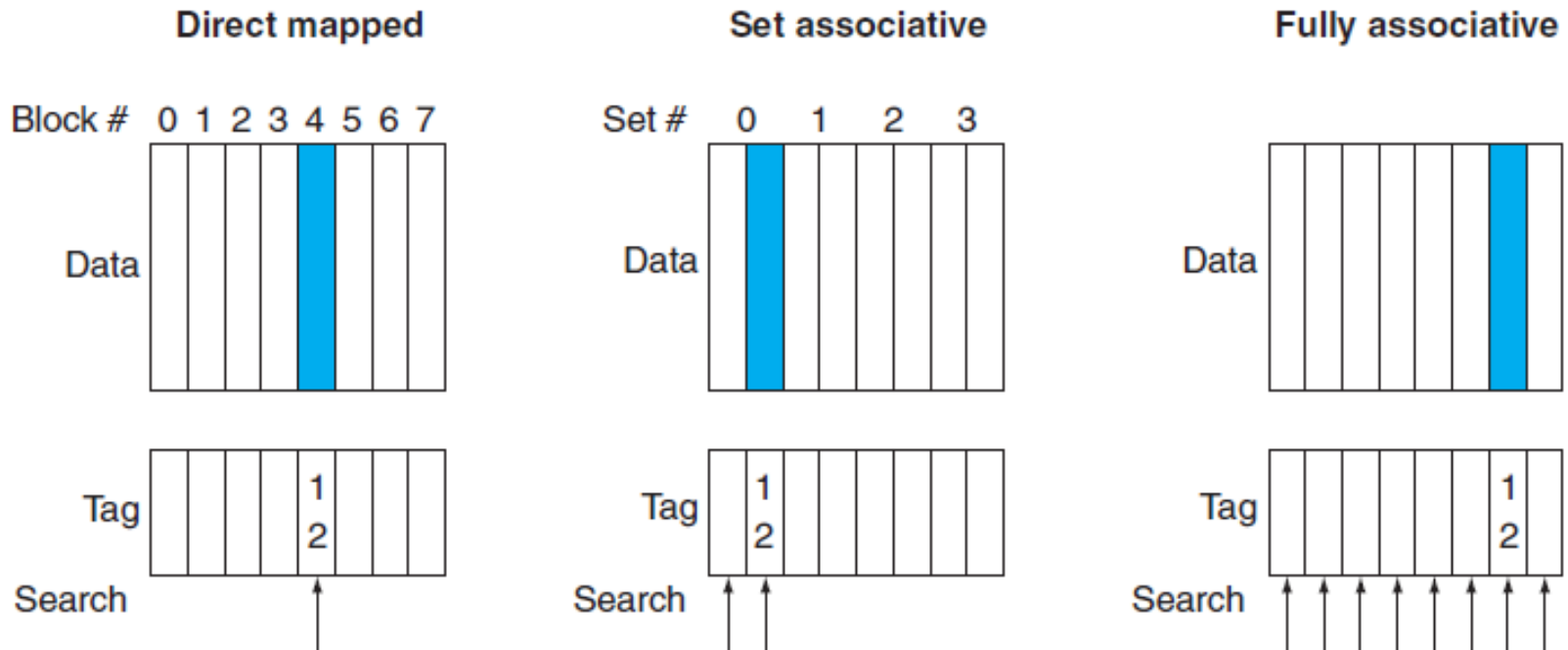A cache structure in which a block can be placed in any location in the cache.

**set-associative cache**
A cache that has a fixed number of locations (at least two) where each block can be placed.

**least recently used (LRU)**
A replacement scheme in which the block replaced is the one that has been unused for the longest time.

# Cache organization



**Direct mapped**

Block # 0 1 2 3 4 5 6 7

Data

Tag 1 2

Search

**Set associative**

Set # 0 1 2 3

Data

Tag 1 2

Search

**Fully associative**

Data

Tag 1 2

Search

**The location of a memory block whose address is 12 in a cache with eight blocks varies for direct-mapped, set-associative, and fully associative placement.** In directmapped placement, there is only one cache block where memory block 12 can be found, and that block is given by (12 modulo 8) 4. In a two-way set-associative cache, there would be four sets, and memory block 12 must be in set (12 mod 4) 0; the memory block could be in either element of the set. In a fully associative placement, the memory block for block address 12 can appear in any of the eight cache blocks.

**One-way set associative (direct mapped)**

| Block | Tag | Data |
|-------|-----|------|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |

**Two-way set associative**

| Set | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

**Four-way set associative**

| Set | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|-----|------|-----|------|
| 0 | | | | | | | | |
| 1 | | | | | | | | |

**Eight-way set associative (fully associative)**

| Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|
| | | | | | | | | | | | | | | | |

**An eight-block cache configured as direct mapped, two-way set associative, four-way set associative, and fully associative.** The total size of the cache in blocks is equal to the number of sets times the associativity. Thus, for a fixed cache size, increasing the associativity decreases the number of sets while increasing the number of elements per set. With eight blocks, an eight-way set associative cache is the same as a fully associative cache.
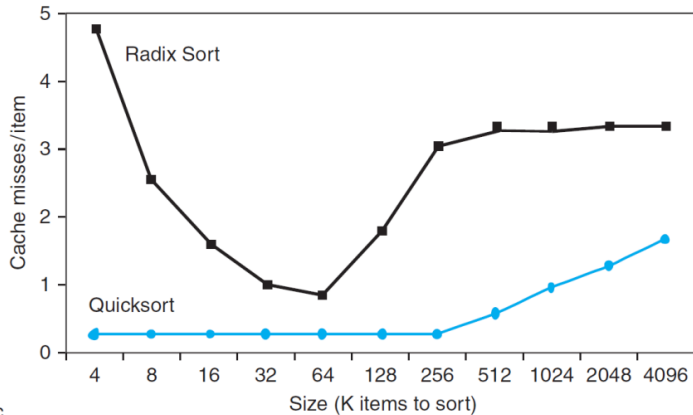
**The implementation of a four-way set-associative cache requires four comparators and a 4-to-1 multiplexor.** The comparators determine which element of the selected set (if any) matches the tag. The output of the comparators is used to select the data from one of the four blocks of the indexed set, using a multiplexor with a decoded select signal. In some implementations, the Output enable signals on the data portions of the cache RAMs can be used to select the entry in the set that drives the output. The Output enable signal comes from the comparators, causing the element that matches to drive the data outputs. This organization eliminates the need for the multiplexor.

**Comparing Quicksort and Radix Sort by (a) instructions executed per item sorted, (b) time per item sorted, and (c) cache misses per item sorted.** This data is from a paper by LaMarca and Ladner [1996]. Due to such results, new versions of Radix Sort have been invented that take memory hierarchy into account, to regain its algorithmic advantages. The basic idea of cache optimizations is to use all the data in a block repeatedly before it is replaced on a miss.

**The key quantitative design parameters that characterize the major elements of memory hierarchy in a computer.** These are typical values for these levels as of 2012. Although the range of values is wide, this is partially because many of the values that have shifted over time are related; for example, as caches become larger to overcome larger miss penalties, block sizes also grow. While not shown, server microprocessors today also have L3 caches, which can be 2 to 8 MiB and contain many more blocks than L2 caches. L3 caches lower the L2 miss penalty to 30 to 40 clock cycles.
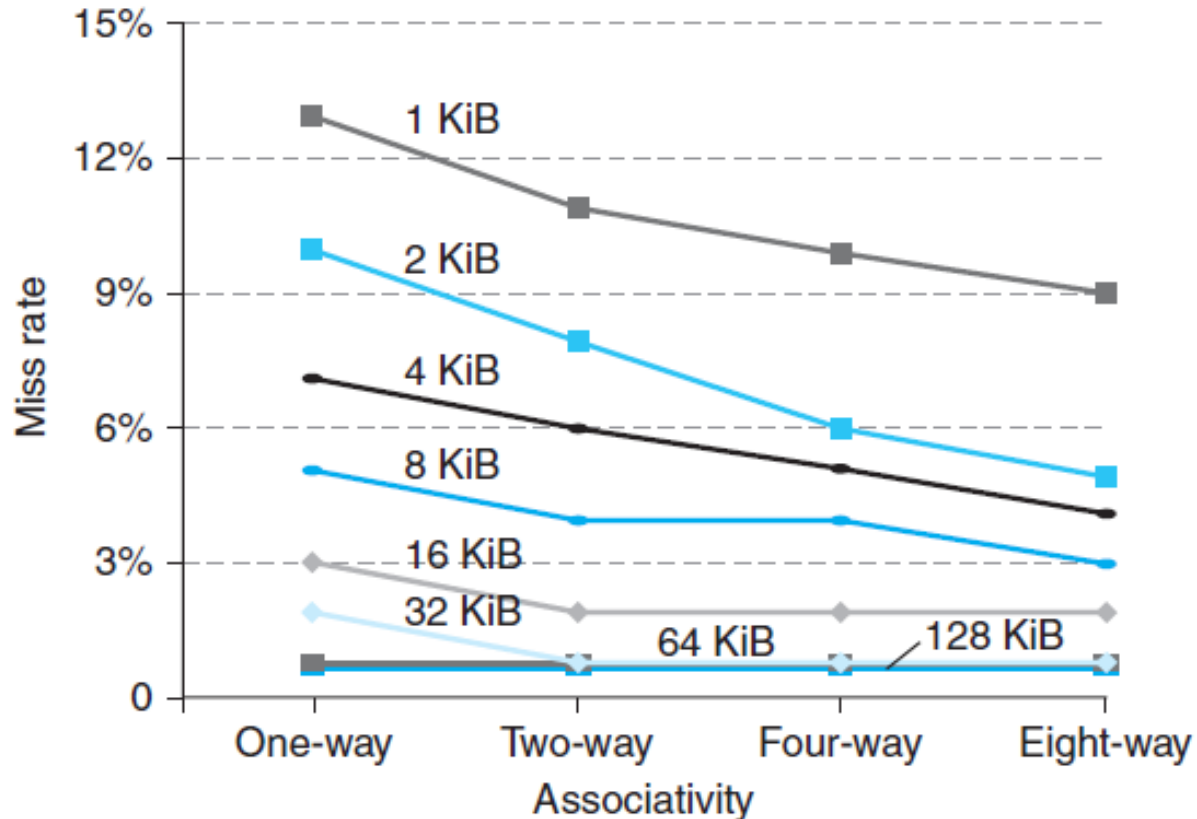
| Feature | Typical values for L1 caches | Typical values for L2 caches | Typical values for paged memory | Typical values for a TLB |
|---|---|---|---|---|
| Total size in blocks | 250–2000 | 2,500–25,000 | 16,000–250,000 | 40–1024 |
| Total size in kilobytes | 16–64 | 125–2000 | 1,000,000–1,000,000,000 | 0.25–16 |
| Block size in bytes | 16–64 | 64–128 | 4000–64,000 | 4–32 |
| Miss penalty in clocks | 10–25 | 100–1000 | 10,000,000–100,000,000 | 10–1000 |
| Miss rates (global for L2) | 2%–5% | 0.1%–2% | 0.00001%–0.0001% | 0.01%–2% |

| Scheme name | Number of sets | Blocks per set |
|---|---|---|
| Direct mapped | Number of blocks in cache | 1 |
| Set associative | $\dfrac{\text{Number of blocks in the cache}}{\text{Associativity}}$ | Associativity (typically 2–16) |
| Fully associative | 1 | Number of blocks in the cache |

The advantage of increasing the degree of associativity is that it usually decreases the miss rate. The improvement in miss rate comes from reducing misses that compete for the same location.

**The data cache miss rates for each of eight cache sizes improve as the associativity increases.**
While the benefit of going from one-way (direct mapped) to two-way set associative is significant, the benefits of further associativity are smaller (e.g., 1%–10% improvement going from two-way to four-way versus 20%–30% improvement going from one-way to two-way). There is even less improvement in going from four-way to eight-way set associative, which, in turn, comes very close to the miss rates of a fully associative cache. Smaller caches obtain a significantly larger absolute benefit from associativity because the base miss rate of a small cache is larger. Figure 5.16 explains how this data was collected.

| Processor | Type | Year of Introduction | L1 Cache[a] | L2 Cache | L3 Cache |
|---|---|---|---|---|---|
| IBM 360/85 | Mainframe | 1968 | 16–32 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 | 128–256 kB | — | — |
| Intel 80486 | PC | 1989 | 8 kB | — | — |
| Pentium | PC | 1993 | 8 kB/8 kB | 256–512 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 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 |
| 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 | 1 MB | — |
| IBM POWER6 | PC/server | 2007 | 64 kB/64 kB | 4 MB | 32 MB |
| IBM z10 | Mainframe | 2008 | 64 kB/128 kB | 3 MB | 24–48 MB |
| Intel Core i7 EE 990 | Workstation/ server | 2011 | $6 \times 32$ kB/ 32 kB | 1.5 MB | 12 MB |
| IBM zEnterprise 196 | Mainframe/ server | 2011 | $24 \times 64$ kB/ 128 kB | $24 \times 1.5$ MB | 24 MB L3 192 MB L4 |

| Characteristic | ARM Cortex-A8 | Intel Core i7 |
|---|---|---|
| Virtual address | 32 bits | 48 bits |
| Physical address | 32 bits | 44 bits |
| Page size | Variable: 4, 16, 64 KiB, 1, 16 MiB | Variable: 4 KiB, 2/4 MiB |
| TLB organization | 1 TLB for instructions and 1 TLB for data<br><br>Both TLBs are fully associative, with 32 entries, round robin replacement<br><br>TLB misses handled in hardware | 1 TLB for instructions and 1 TLB for data per core<br><br>Both L1 TLBs are four-way set associative, LRU replacement<br><br>L1 I-TLB has 128 entries for small pages, 7 per thread for large pages<br><br>L1 D-TLB has 64 entries for small pages, 32 for large pages<br><br>The L2 TLB is four-way set associative, LRU replacement<br><br>The L2 TLB has 512 entries<br><br>TLB misses handled in hardware |

**Address translation and TLB hardware for the ARM Cortex-A8 and Intel Core i7 920.**
Both processors provide support for large pages, which are used for things like the operating system or mapping a frame buffer. The large-page scheme avoids using a large number of entries to map a single object that is always present.

# Caches in the ARM Cortex-A8 and Intel Core i7 920.

| Characteristic | ARM Cortex-A8 | Intel Nehalem |
|---|---|---|
| L1 cache organization | Split instruction and data caches | Split instruction and data caches |
| L1 cache size | 32 KiB each for instructions/data | 32 KiB each for instructions/data per core |
| L1 cache associativity | 4-way (I), 4-way (D) set associative | 4-way (I), 8-way (D) set associative |
| L1 replacement | Random | Approximated LRU |
| L1 block size | 64 bytes | 64 bytes |
| L1 write policy | Write-back, Write-allocate(?) | Write-back, No-write-allocate |
| L1 hit time (load-use) | 1 clock cycle | 4 clock cycles, pipelined |
| L2 cache organization | Unified (instruction and data) | Unified (instruction and data) per core |
| L2 cache size | 128 KiB to 1 MiB | 256 KiB (0.25 MiB) |
| L2 cache associativity | 8-way set associative | 8-way set associative |
| L2 replacement | Random(?) | Approximated LRU |
| L2 block size | 64 bytes | 64 bytes |
| L2 write policy | Write-back, Write-allocate (?) | Write-back, Write-allocate |
| L2 hit time | 11 clock cycles | 10 clock cycles |
| L3 cache organization | – | Unified (instruction and data) |
| L3 cache size | – | 8 MiB, shared |
| L3 cache associativity | – | 16-way set associative |
| L3 replacement | – | Approximated LRU |
| L3 block size | – | 64 bytes |
| L3 write policy | – | Write-back, Write-allocate |
| L3 hit time | – | 35 clock cycles |