

Computer Architecture License: <https://creativecommons.org/licenses/by-nc-nd/4.0/>

**SRAM (Static RAM):**

- It is used to implement the **cache memory**.
- SRAM uses flip-flops (or latches) to store bits (see Digital Circuits lecture notes).
- It requires 6 transistors/bit (more expensive, fewer bits per chip). (-)
- SRAM requires no refresh. (+)
- Access time = cycle time. (+)

**DRAM vs. SRAM:**

- Both are volatile. Power must be continuously supplied.
- DRAM is more dense (smaller cells, more cells per unit area) and less expensive than a corresponding SRAM.
- DRAM requires the supporting refresh circuitry. The latency of DRAM is high.
- DRAM is cheap, SRAM is fast.

http://akademilitu.edu.tr/en/buzluca  
http://www.buzluca.info/

2013 - 2020 Feza BUZLUCA 6.7

Computer Architecture

**Associative Memory / Content Addressable Memory (CAM):**

- Random access memory (SRAM) + circuitry for search
- It consists of SRAM to store data and digital circuits for parallel search.
- It is used in high speed searching applications.

Cache memory is one of the primary uses for associative memory.

- A Content Addressable Memory (CAM) is an SRAM-based memory, which can be accessed in parallel to search for a given search word, providing as result the address of the matching data.
- A word is retrieved based on a portion of its contents rather than its address.
- The user supplies a data word (Argument A) (not the address), and the CAM searches its entire memory simultaneously (not sequentially) to see if that word is stored anywhere in it.
- The user also supplies a key value (K) to determine the portion of data to search for.
- If the search data (or the required portion) is found in a row of the memory, then the match bit for that row is set, and the data stored in this row is output.

http://akademilitu.edu.tr/en/buzluca  
http://www.buzluca.info/

2013 - 2020 Feza BUZLUCA 6.8

Computer Architecture

**Associative Memory / Content Addressable Memory (CAM) (cont'd):**

On power-up, the memory contains random values.

An additional flag (valid bit V) is necessary to indicate whether or not a memory line has been loaded with valid data.

On power-up, the hardware sets valid bits of all lines to "invalid" (V=0).

When data is written to line i, the corresponding valid bit is set (V=1).

Example:  
A: 10111111  
K: 11110000  
Memory: 11001111 M=0  
          10110110 M=1  
Out: 10111010

One bit for each row (m bits)

http://akademilitu.edu.tr/en/buzluca  
http://www.buzluca.info/

2013 - 2020 Feza BUZLUCA 6.9

Computer Architecture

**Internal Structure:**

$A_i$  and  $K_j$  values are broadcasted to all cells in parallel.

1. word  $x_{11} \dots x_{1n}$   
i. word  $x_{i1} \dots x_{in}$   
m. word  $x_{m1} \dots x_{mn}$

AND Gates

$M_i = \prod x_{ij}$

$V_i$

The structure of a cell ( $c_{ij}$ ):

Input

Write E

Read

D Latch

Match Logic

$x_{ij}$

$F_{ij}$

$K_j$

$A_i$

$x_{ij} = K_j + A_i \odot F_{ij}$   
 $x_{ij} = K_j + A_i \odot F_{ij}$   
 $x_{ij} = K_j + A_i \odot F_{ij} + A_i \odot F_{ij}$

$M_i = \prod_{j=1}^n x_{ij}$

http://akademilitu.edu.tr/en/buzluca  
http://www.buzluca.info/

2013 - 2020 Feza BUZLUCA 6.10

Computer Architecture **6.5 Cache Memory , The internal memory system**

**6.5.1 Cache Memory Principles**

The goal is to store copies of frequently accessed data in the fast cache memory to reduce the average access time.

**Hit:** The requested word is found in cache memory.

**Miss:** A miss occurs if the requested data is not found in cache memory.

Example:  
For read operations only:  
Cache memory access time: 20 ns.  
Main memory access time: 100 ns.  
Hit ratio:  $H=0.9$  (90% hit)  
Average memory access time:  
 $t_a = 0.9 \cdot 20 + 0.1 \cdot 100 = 28$  ns.

In systems with asynchronous bus operations like the MC68000, the DTACK signal is sent to the CPU at different times based on there being a hit or not. Thus, when data is retrieved from cache memory, the bus cycle is completed faster.

http://akademilitu.edu.tr/en/buzluca  
http://www.buzluca.info/

2013 - 2020 Feza BUZLUCA 6.11

Computer Architecture **6.5.1 Cache Memory Principles (cont'd)**

**Block transfer:**

Upon a cache miss (the requested element is not found in cache), a **block** of elements that contains the requested element is brought (copied) from main memory to cache memory.

**Reason for block** (instead of a single word) **transfer: spatial locality.**

The next element to be requested will most likely be located near the currently requested element.

The disadvantage of transferring a block instead of a single element is that it takes longer.

To reduce the block transfer time between main and cache memories, the **memory interleaving technique** is used.

Data are stored in different memory modules; that is, consecutive memory addresses are located in successive memory modules, so they can be accessed at the same time (similar to RAID, section 7.2).

consecutive bytes

M7 M6 M5 M4 M3 M2 M1 M0

B7 B6 B5 B4 B3 B2 B1 B0

Interleaved main memory

One block of cache memory

http://akademilitu.edu.tr/en/buzluca  
http://www.buzluca.info/

2013 - 2020 Feza BUZLUCA 6.12



Computer Architecture License: <https://creativecommons.org/licenses/by-nc-nd/4.0/>

**Example:** Array in a cache memory system (cont'd)

Main Memory: 256K x words Cache Memory: 2K x words Block size: 16 words

**Case B)** A program in this system accesses an array with the starting address \$0000A and a size of 10 words. The array is not in cache memory. Assume that, the least recently used frames are Frames #0 and #2.

Starting address of the array: 00 0000 0000 0000 1010 (\$0000A)  
End address of the array: 00 0000 0000 0001 0011 (\$00013)

Each block is processed separately. Different blocks have different tag values. In this example, two block replacements are necessary.

Although the size of the array is smaller than the block (frame) size, it occupies two frames in cache memory.

<http://akademik.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.19

Computer Architecture

**Example:** Array in a cache memory system (cont'd)

- When the CPU accesses the first word of the array (\$0000A), a miss occurs.
- The cache management system transfers Block 0 in its entirety (16 words) to Frame 0.
- Next 5 accesses to the array generate hits.
- When the CPU accesses the 7th word of the array (\$00010), a miss occurs because this word is in another block and its address has a new tag value.
- The cache management system transfers Block 1 with its entirety (16 words) to Frame 2 (because of LRU).
- Next 3 accesses to the array generate hits.
- In total: 2 misses, 8 hits

If the starting address of the array was selected properly (for example \$00000), the array would occupy only one block and one frame in the cache memory as in case A.

Therefore, placement policies for arrays in main memory affect the performance of cache systems.

<http://akademik.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.20

Computer Architecture

## 2. Direct Mapping

An incoming main memory block is always placed into a specific, fixed cache frame location.

It is not necessary to search for the location of a block in the cache because it is predetermined and fixed.

Therefore, associative memory is not necessary.

As the size of main memory is greater than that of the cache, several blocks of main memory map to the same cache frame.

It is necessary to determine which main memory block is currently residing in a frame.

The cache memory control unit divides the address from the CPU into three fields:

a bits

Tag	Cache Frame number	Word number
a-(f+w) bits	f bits	w bits

It indicates which of the blocks that can be placed in this frame is currently in cache.

This field determines the number of the cache frame that will hold this data. The blocks that have this filed in common (same) try to reside in the same frame. Only one of them can reside in the cache any given moment.

<http://akademik.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.21

Computer Architecture

a bit

Tag	Cache Frame num.	Word num.
a-(f+w) bit	f bit	w bit

V Tag: a-(f+w) bit

The frame is determined directly.

$i = B \bmod F$

i : Frame number  
B : Block number  
F : Number of frames  
 $F = 2^f$

Main memory  $2^b$  blocks

One block  $2^w$  words

Block 0  
Block 1  
:  
Block  $2^f$   
:  
Block  $2^b - 1$

Block 0  
Frame 0  
Frame 1  
:  
Frame  $2^f - 1$

$F = 2^f$  frames

A main memory block can only reside in the cache frame with the number that is determined by the "Cache Frame number" field of the address.

Even if there are unused (empty) frames in the cache, two blocks with the same "Cache Frame Number" field cannot be in the cache at the same time.

During replacement, a decision algorithm is not necessary.

Since the frame of the incoming block is fixed, this frame will be replaced.

**Associative memory is not necessary** because there is no need to search in the cache.

<http://akademik.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.22

Computer Architecture

**Example (Direct Mapping):**

Main Memory: 256K x word Address: a = 18 bits  
Block size: 16 words w = 4 bits Main memory contains  $2^{14}$  blocks. b = 14  
Cache Memory: 2K x word Cache memory contains  $2^7 = 128$  frames. f = 7  
data capacity

Main memory address:

Tag	Frame number	Word number
7 bits	7 bits	4 bits

In this system, the data in the following two addresses try to reside in the same cache frame.

Tag	Frame num.	Word num.
0000000	0000000	XXXX
0000001	0000000	XXXX

The "Frame number" fields of both addresses are the same: 0000000. They will be placed in Frame 0. At a specific point in time, only one of these data can be in cache memory.

To determine which data is currently in Frame 0 of cache memory, the tag value of the address is compared to the tag value stored in cache memory.

<http://akademik.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.23

Computer Architecture

**Example (Direct Mapping):**

Main address from CPU: 0010001 0000001 1000

Tag Frame Word

Hit/Miss ← Compare

Cache Data memory is addressed directly.

Cache Tag memory is addressed directly.

If the tag in the address is equal to the tag in the cache, then the data of the referenced address is in cache memory.

Tag Memory

1	1111111
1	0010001
0	0000000

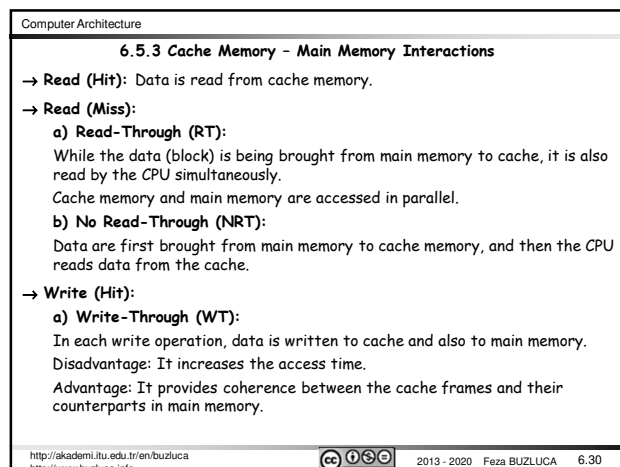
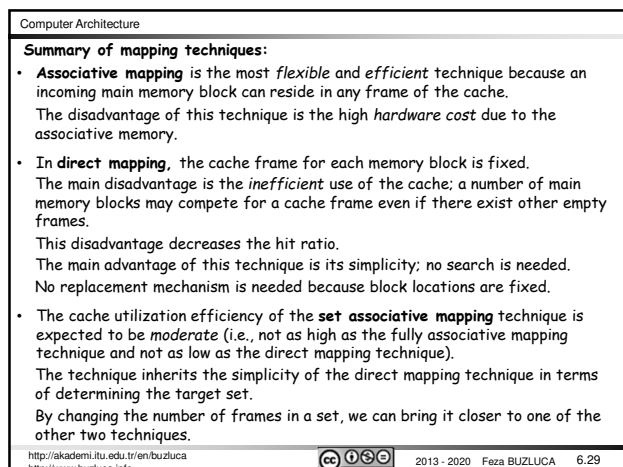
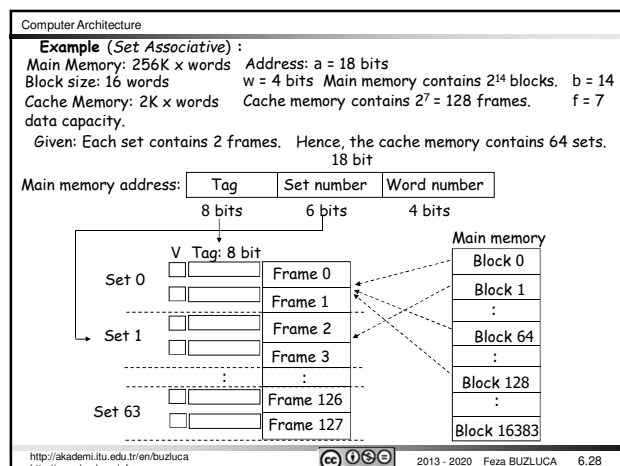
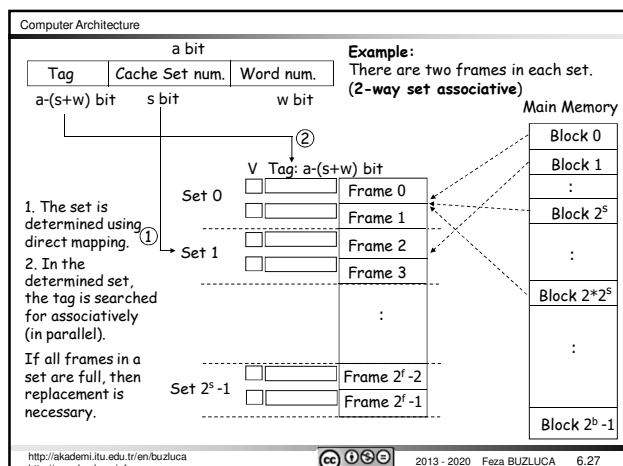
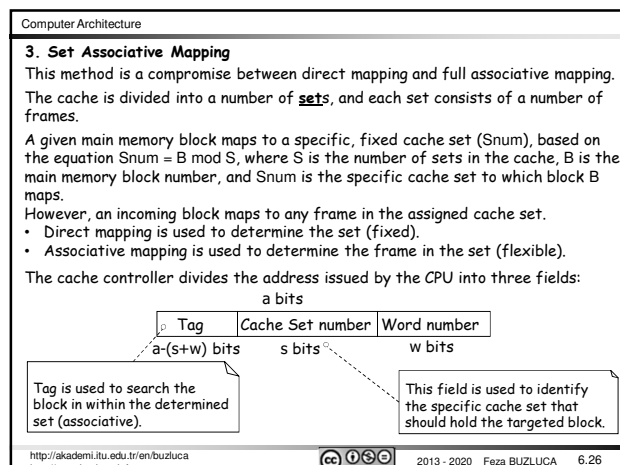
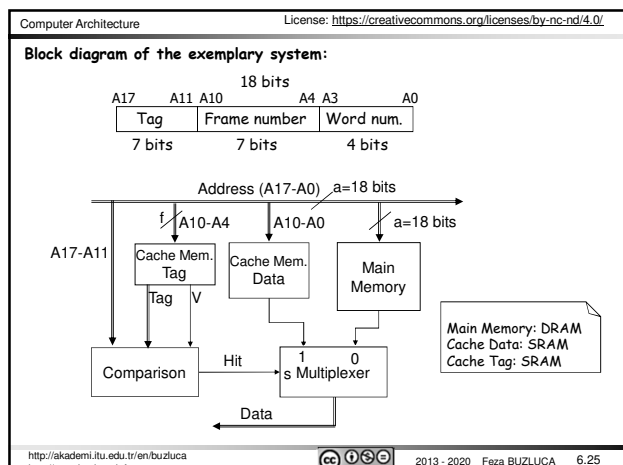
Data Memory

Frame 0
Frame 1
:
Frame 127

**Cache Memory**

<http://akademik.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.24



Computer Architecture License: <https://creativecommons.org/licenses/by-nc-nd/4.0/>

→ **Write (Hit)** (cont'd):

b) **Write-Back (WB):**

Writes are done only to the cache.  
A block is written back to main memory only when a replacement is needed.  
There are two types of write-back policies: *Simple write-back* and *flagged write-back*.

- **Simple Write-Back (SWB):**  
The replaced frame is always written back to main memory.  
It is not checked whether the frame was changed or not.
- **Flagged Write-Back (FWB):**  
Every cache frame is assigned a bit, called the *dirty bit*, to indicate that at least one write operation has been made to the block while residing in cache.  
At replacement time, the dirty bit is checked: if it is set, then the block is written back to main memory; otherwise, it is simply overwritten by the incoming block.  
The dirty bit is stored in the tag memory of the cache.

<http://akademil.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.31

Computer Architecture

→ **Write (Miss):**

a) **Write Allocate (WA):**  
The main memory block is updated and brought to cache.

a) **No Write Allocate (NWA):**  
The missed main memory block is updated in main memory and not brought to cache.  
If an attempt is made to read this block later, a miss will occur, and data will be brought to cache.

The write-through (WT) policy can be used together with write-allocate (WA) or no-write-allocate (NWA) methods. WTWA, WTNWA

In write-back (WB) policy, to maintain coherence between cache and main memory at the beginning, the write-allocate (WA) method is used (WBWA).

**Information held in Tag memory:**  
In addition to Valid (V) and tag bits, depending on the method used, the following data must be also kept in tag memory:

- If LRU is used *aging counters*,
- If flagged Write-Back (FWB) method is used "dirty" bit (D).

A single line of a tag memory: 

V	D	Counter	Tag
---	---	---------	-----

<http://akademil.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

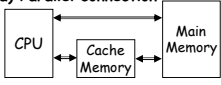
2013 - 2020 Feza BUZLUCA 6.32

Computer Architecture

### 6.5.4 CPU - Cache Memory - Main Memory Interconnection

CPU - cache memory - main memory interconnection can be implemented in two ways:

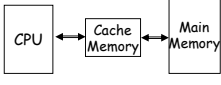
a) **Parallel connection**



- **Read-Through (RT):** While a block is brought from main memory to cache it is also read by the CPU simultaneously (in parallel).
- **Load Through (LT):** A data is written to cache and to main memory simultaneously.

The parallel connection is suitable for the *Write-Through (WT)* method.  
It can also be used with the *Write-Back (WB)* method.

b) **Serial connection**



- **No Read-Through (NRT):** The block is first brought from main memory to cache, and then the CPU reads the data from cache.
- **No Load-Through (NLT):** First, the data in cache is updated, then the frame is written to main memory.

This serial connection is suitable for the *Write-Back (WB)* method.

<http://akademil.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.33

Computer Architecture

### 6.5.5 Access Time:

$t_a$ : Average Memory Access Time  
W: Write ratio (number of write accesses / total number of all accesses)  
h: Hit ratio  
 $t_{cache}$ : Cache memory access time  
 $t_{main}$ : Main memory access time  
 $t_{trans}$ : Time to transfer a block between main memory and cache  
 $W_d$ : The probability that a block in a cache is updated

WT, RT/LT (Write-through, Parallel read/write)			WB, WA, NRT/NLT (Write-back, Serial read/write)	
Probability	NWA	WA	SWB	FWB
<b>Read Hit</b> (1-w)h	$t_{cache}$	$t_{cache}$	$t_{cache}$	$t_{cache}$
<b>Read Miss</b> (1-w)(1-h)	$t_{trans}$	$t_{trans}$	$2t_{trans} + t_{cache}$	$W_d(2t_{trans} + t_{cache}) + (1-W_d)(t_{trans} + t_{cache})$
<b>Write Hit</b> wh	$t_{main}$	$t_{main}$	$t_{cache}$	$t_{cache}$
<b>Write Miss</b> w(1-h)	$t_{main}$	$t_{main} + t_{trans}$	$2t_{trans} + t_{cache}$	$W_d(2t_{trans} + t_{cache}) + (1-W_d)(t_{trans} + t_{cache})$

<http://akademil.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.34

Computer Architecture

### Access Time Calculation:

• **Write Through with Write Allocate, Read/Load Through (WTWA, RT/LT):**  
Read Hit + Read Miss + Write Hit + Write Miss  
 $t_a = (1-w)h t_{cache} + (1-w)(1-h)t_{trans} + w h t_{main} + w(1-h)(t_{main} + t_{trans})$   
 $t_a = (1-w)h t_{cache} + (1-h)t_{trans} + W t_{main}$

• **Write Through with No Write Allocate, Read/Load Through (WTNWA, RT/LT)**  
 $t_a = (1-w)h t_{cache} + (1-w)(1-h)t_{trans} + w h t_{main} + w(1-h)t_{main}$   
 $t_a = (1-w)h t_{cache} + (1-w)(1-h)t_{trans} + W t_{main}$

• **Simple Write Back with Write Allocate, No Read Through (SWBWA, NRT/NLT)**  
Read Hit + Read Miss + Write Hit + Write Miss  
 $t_a = (1-w)h t_{cache} + (1-w)(1-h)(2t_{trans} + t_{cache}) + w h t_{cache} + w(1-h)(2t_{trans} + t_{cache})$   
 $t_a = t_{cache} + (1-h)2t_{trans}$   
One  $t_{trans}$  is needed to transfer a frame from cache to main memory and the second one is needed to bring the new block from main memory to cache.

• **Flagged Write Back, Write Allocate, No Read Through (FWBWA, NRT/NLT):**  
 $t_a = t_{cache} + (1-h)t_{trans} + W_d(1-h)t_{trans}$   
 $t_a = t_{cache} + (1-h)(1+W_d)t_{trans}$

<http://akademil.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.35

Computer Architecture

### Exemplary processors with cache memories:

- **Intel386™:** Cache memory is outside of the CPU chip. SRAM memory.
- **Intel486™ (1989)**  
8-KByte on-chip (L1)
- **Intel® Pentium® (1993)**  
L1 on-chip: 8 KB instruction, 8 KB data cache (Harvard architecture)
- **Intel P6 Family: (1995-1999)**
  - **Intel Pentium Pro:**  
L1 on-chip: 8 KB instruction, 8 KB data cache (Harvard architecture)  
First L2 cache memory in the CPU chip.  
L2 on-chip: 256 KB. Different interconnections between L1, L2 and the CPU.
  - **Intel Pentium II:**  
L1 on-chip: 16 KB instruction, 16 KB data cache (Harvard architecture)  
L2 on-chip: 256 KB, 512 KB, 1 MB
- **Intel® Pentium® M (2003)**  
L1 on-chip: 32 KB instruction, 32 KB data cache  
L2 on-chip: up to 2 MByte
- **Intel® Core™ i9-9900 (2019)**  
Multicore: 8 cores. Private caches (L1: 512KiB) and shared caches (L2: 2 MiB)  
L3: 16 MiB smartcache: All cores share this cache.

<http://akademil.itu.edu.tr/en/buzluca>  
<http://www.buzluca.info>

2013 - 2020 Feza BUZLUCA 6.36