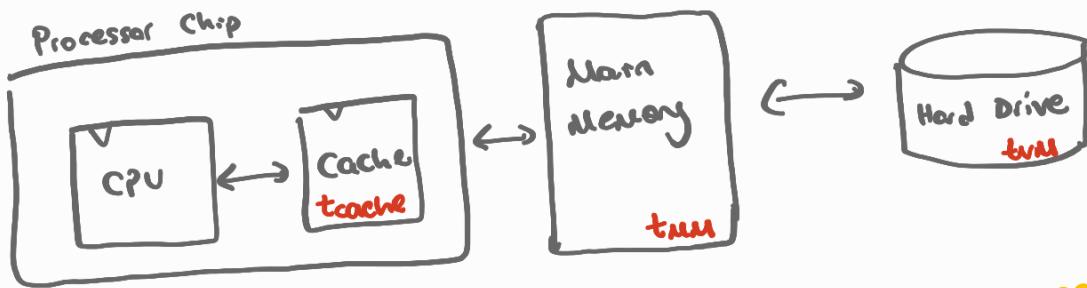


# MEMORY SYSTEMS & ORGANIZATION

## From Ece Güran Schmidt's Notes

## Overview

**Overview** In reality, memory is not as fast as the processor.



# locality

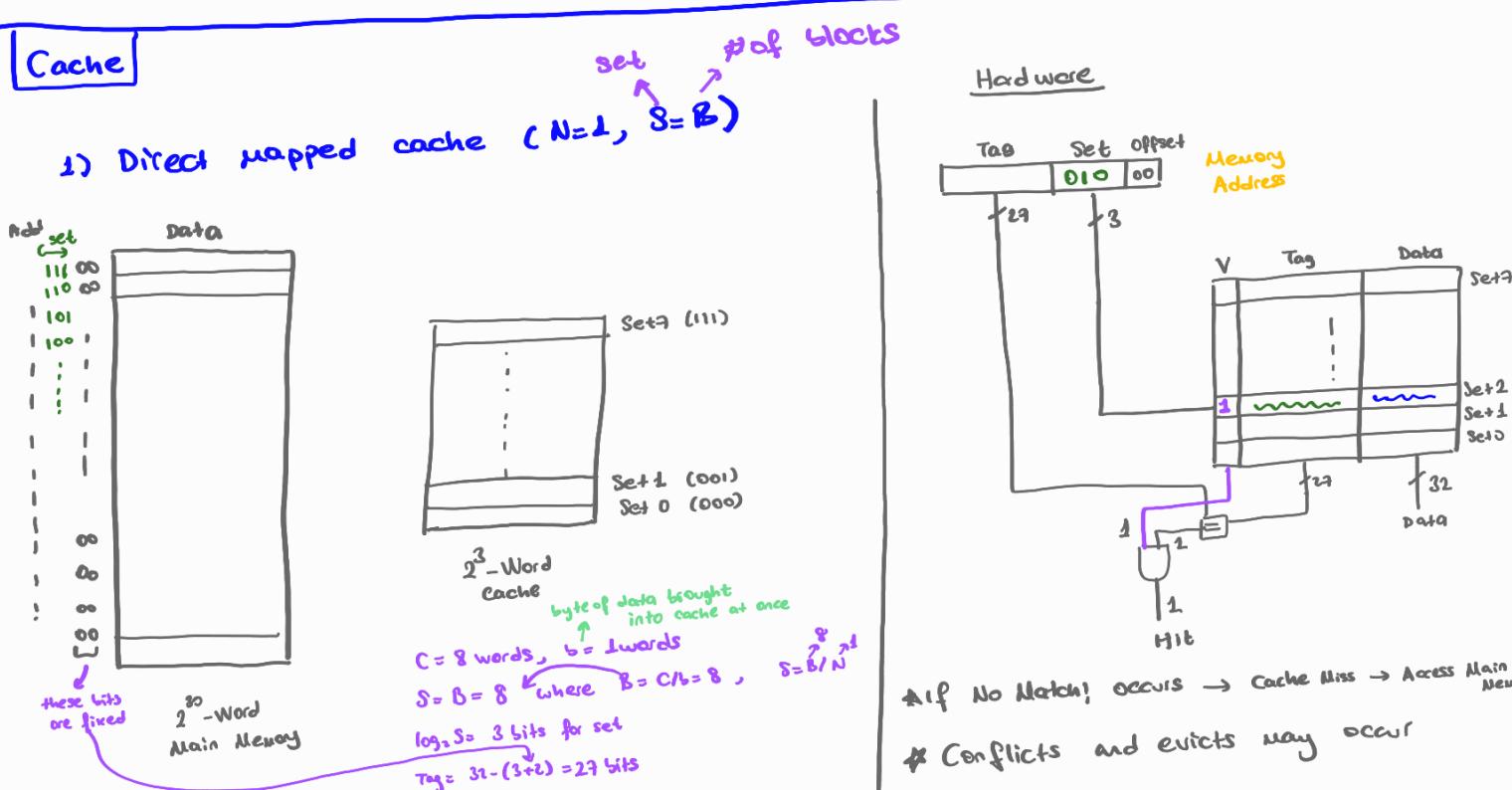
- Temporal → keep recently accessed data in higher levels of memory hierarchy
- Spatial → when access data, bring nearby data into higher levels of memory hierarchy too

$$\text{Miss Rate} = \frac{\# \text{ of misses}}{\# \text{ of total memory accesses}}$$

$$\text{Hit Rate} = 1 - \text{Miss Rate}$$

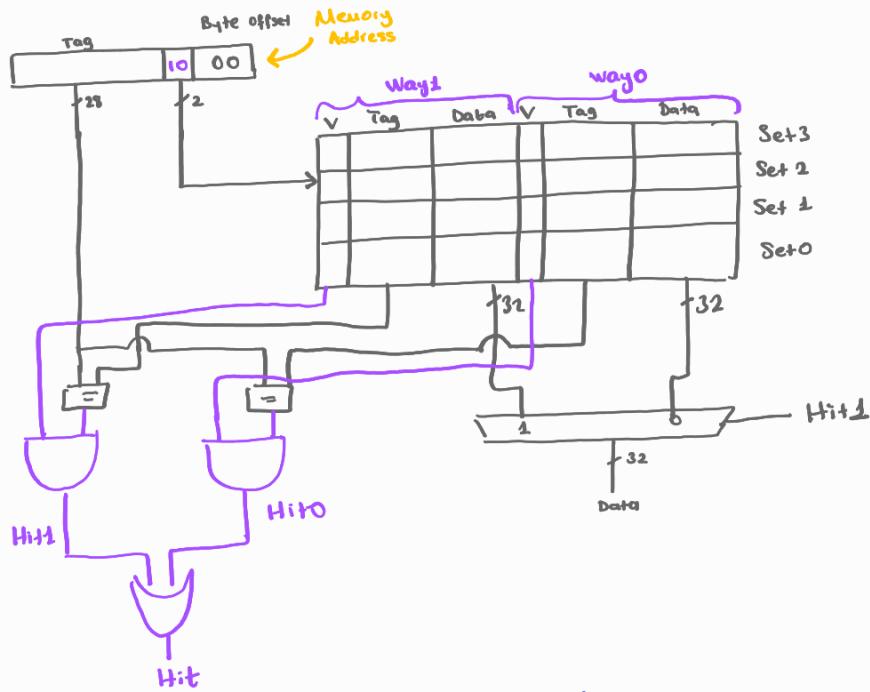
# Memory Performance

$$\text{Average memory access time (AMAT)} = t_{\text{cache}} + MR_{\text{cache}} (t_{\text{MM}} + MR_{\text{MM}} t_{\text{MM}})$$



## 2) Associative caches $N > 1$

2.1) N-way set associative cache:  $1 < N < B$ ,  $S = B/N$



$$\# \text{ of blocks} = \frac{C}{S}$$

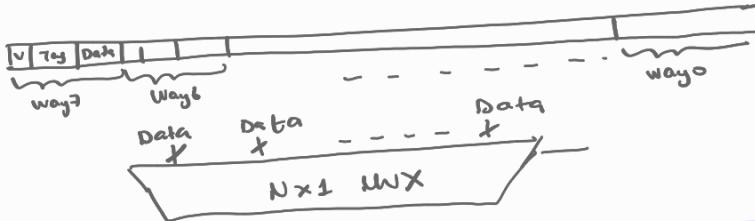
$$S = \frac{B}{N}$$

$$C = 8, S = 4, N = 2$$

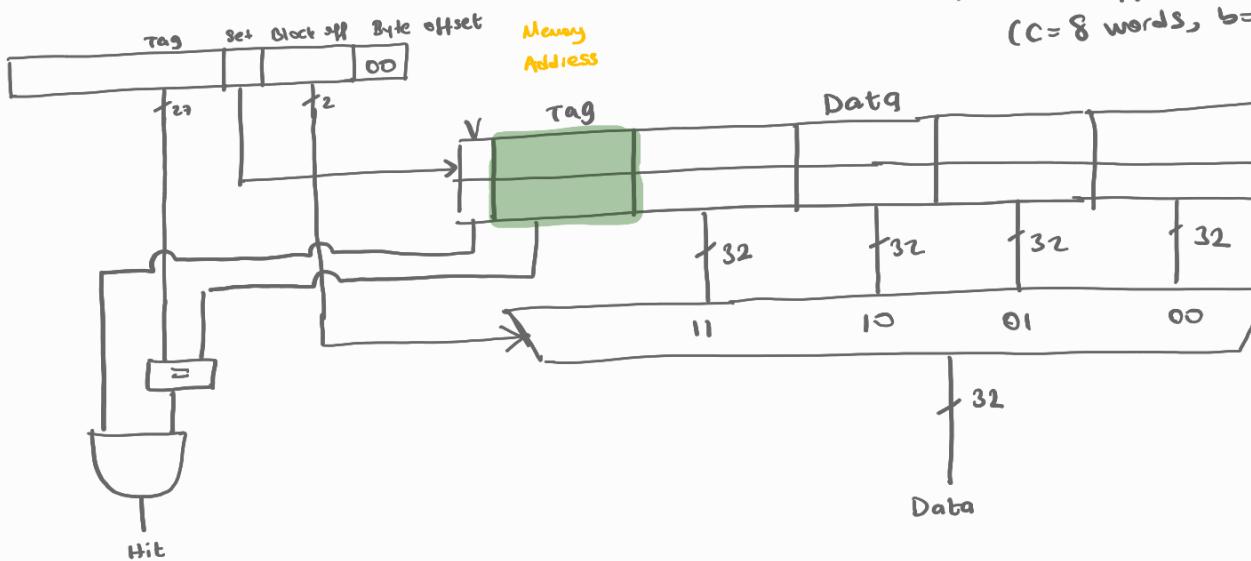
capacity      ↓      # of sets

degree of associativity ( $N$ )

2.2) Fully associative cache ( $N=B, S=1$ )



Increase block size (more than one word/block)



Direct Mapped Cache  
( $C = 8$  words,  $b = 4$  words,  $B = \frac{8}{4} = 2$ )

## Selecting the Data to Replace

Set associative cache,  $N=2$ , LRU:

Each time one of the way is used, U is adjusted to indicate the other way.

\* Use additional U bit for it.

## Example

$s = 16 \text{ bit}$   
 $W = 2 \text{ Byte (words)}$   
 $C = 64 \text{ Bytes}$   
 $b = 2 \text{ words}$   
 $A = 16 \text{ bits long}$   
 $N = 2$

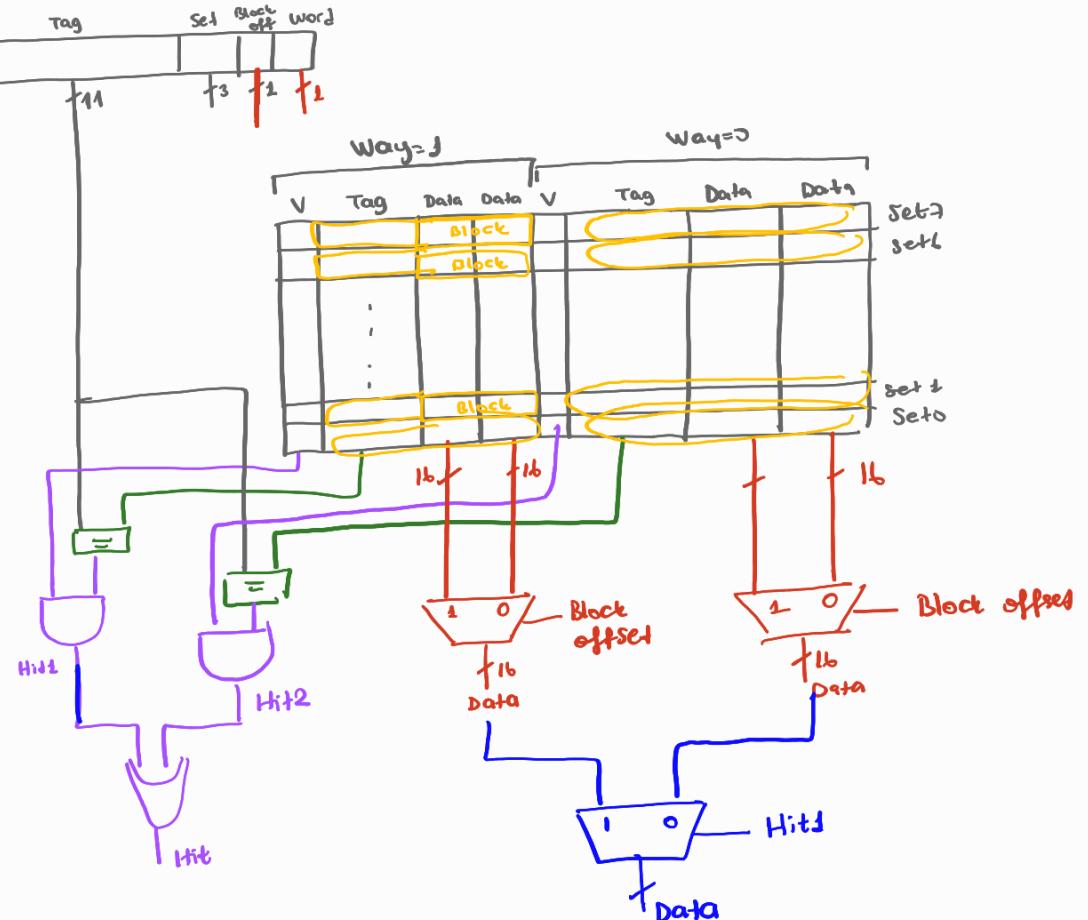
$$\frac{64 \text{ Byte}}{2 \text{ Byte}} = 2^5 \text{ word Cache}$$

$$B = \frac{C}{b} = \frac{32 \text{ word}}{2 \text{ word}} = 2^5$$

$$S = \frac{B}{N} = \frac{16}{2} = 8 \rightarrow \log_2 8 = 3 \text{ bit}$$

$$3+2=5$$

$$16-5=11 \text{ (Tag bits)}$$



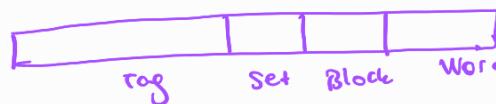
## Review

→ Capacity = # words in the cache

→ Cache block = cache line = fixed size unit of data transferred in one shot

$$\rightarrow \text{Cap} = C = \underbrace{S \times N}_{B} \times b \times W$$

$\underbrace{A \text{ bits}}_{\text{Memory}} = \# \text{tag bits} + \# \text{bits for set} + \# \text{bits for block} + \text{word offset}$



$$C = S \times N \times b \times W$$

$$64 = S \times 2 \times 2 \times 2 \Rightarrow S = 8 \Rightarrow \boxed{\log_2 8 = 3}$$

$$\# \text{tag bits} = 16 - (\underbrace{\log_2 8}_{\text{Set}} + \underbrace{\log_2 2}_{\text{word}} + \underbrace{\log_2 2}_{\text{block}}) = 11 \text{ bits}$$

memory address

$$\# \text{of blocks} = \frac{C}{b} = \frac{64 \text{ bytes}}{2 \frac{\text{words}}{\text{Byte}}} = 16 \text{ blocks}$$

$$\text{Each block} = 16 \times \left( \underbrace{11 \text{ bits}}_{\text{tag}} + \underbrace{32 \text{ bits}}_{\text{data}} \right) = 16 \times 43 //$$

(Tag + Data)

## Example]

Given  $C=16$  words

Construct the cache and map the following addresses ; 7C, 9C in :

a- Direct Mapped ,  $b=1$  ( $N=1$ )

b- 2 way ,  $b=1$  ( $N=2$ )

c- Direct Mapped ,  $b=2$  ( $N=1$ )

Then Map 78, 74, 70, 20, 1C, 98, ...

a)

$$C = \underbrace{S \times N}_{16} \times \underbrace{b}_{1} \times \underbrace{w}_{1 \text{ word}} \rightarrow S=16 \rightarrow 4 \text{ bits to encode}$$

(word=4 byte)

$$7C = 0111100 \xrightarrow{\text{set}} \text{set } 15$$

$$9C = 10011100 \rightarrow \text{set } 7$$

b) 2 way  $b=1$

$$C = \underbrace{S \times N}_{8} \times \underbrace{b}_{2} \times \underbrace{w}_{1} \rightarrow S=8 \rightarrow \log_2 8 = 3 \text{ bits}$$

$$7C = 0111 \quad 1100 \rightarrow \text{set } 7$$

$$9C = 1001 \quad 1100 \rightarrow \text{set } 7$$

c)  $C = \underbrace{S \times N}_{8} \times \underbrace{b}_{1} \times \underbrace{w}_{2} \rightarrow S=8 \rightarrow 3 \text{ bits for set}$

1 bit for block

$$7C = 0111 \quad 1100 \xrightarrow{\text{block}} \text{set } 7$$

$$9C = 1001 \quad 1100 \xrightarrow{\text{byte offset}} \text{set } 3$$

## Cache Misses

- Compulsory Miss: The first request to a cache block
- Capacity Miss: The cache is too small to hold all concurrently used data
- Conflict Miss:

## Cache Writes

Write-back Cache Writes

Dirty bit (D)

Write-through Caches

no dirty bit

Handling Write Misses:

→ Write Allocate

→ No-Write Allocate

No-write allocate

STR R1, #100 → Not in cache, write miss

STR R2, #100 → Still not in cache, write miss

LDR R3, #200 → Not in cache, read miss, loaded to cache

STR R1, #200 → In cache, Hit

STR R2, #100 → Not in cache, write miss

Write-Allocate

STR R1, #100 → Not in cache, write miss, written in cache

STR R2, #100 → In cache, hit

LDR R3, #200 → Not in cache, read miss, loaded in cache

STR R1, #200 → In cache, hit

STR R2, #100 → In cache, hit