

Cache Memory (H.1)

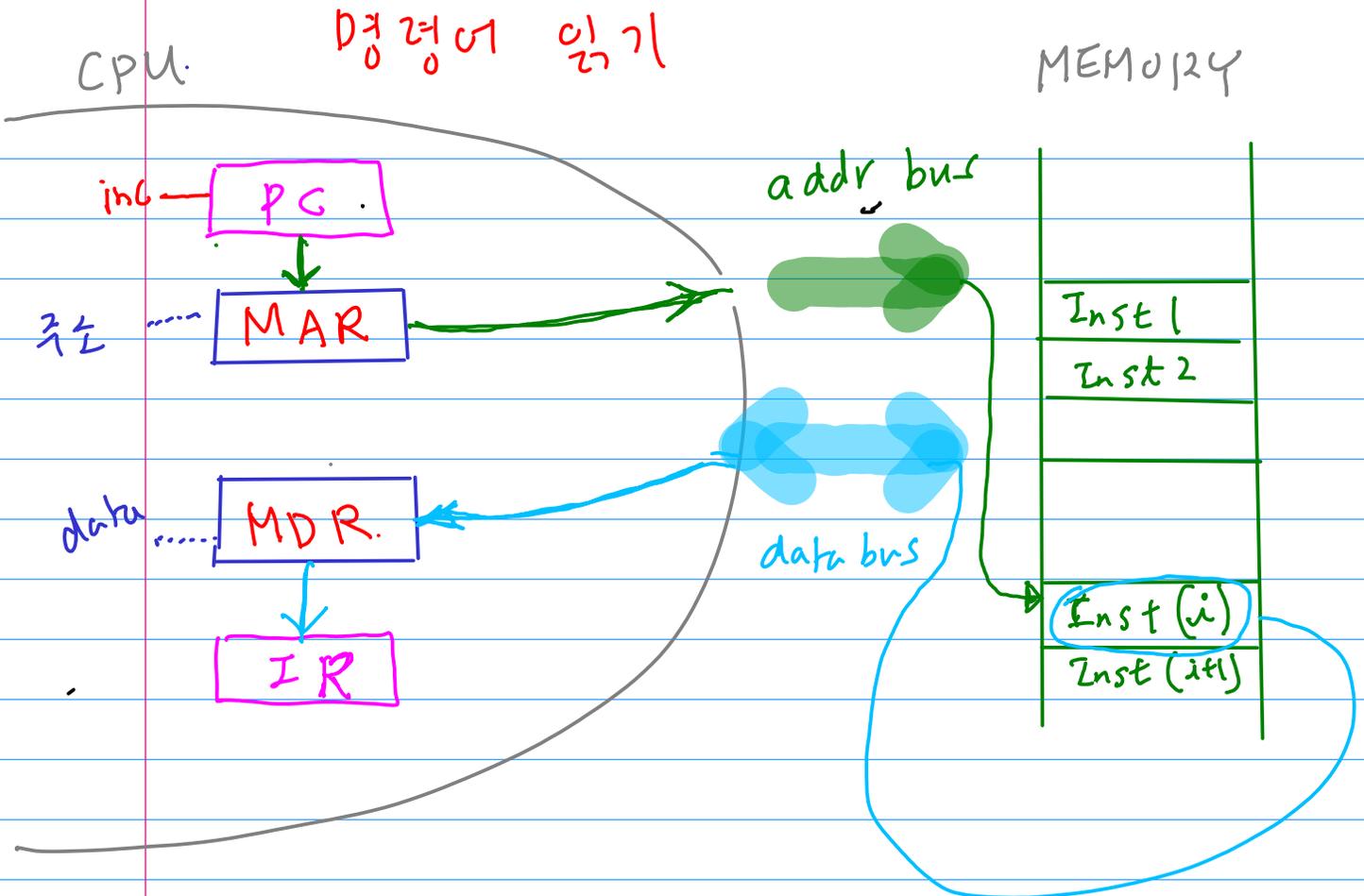
2015.06.11

www.wikiversity.org

The necessities in Computer Organization

Copyright (c) 2015 Young W. Lim.

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled "GNU Free Documentation License".



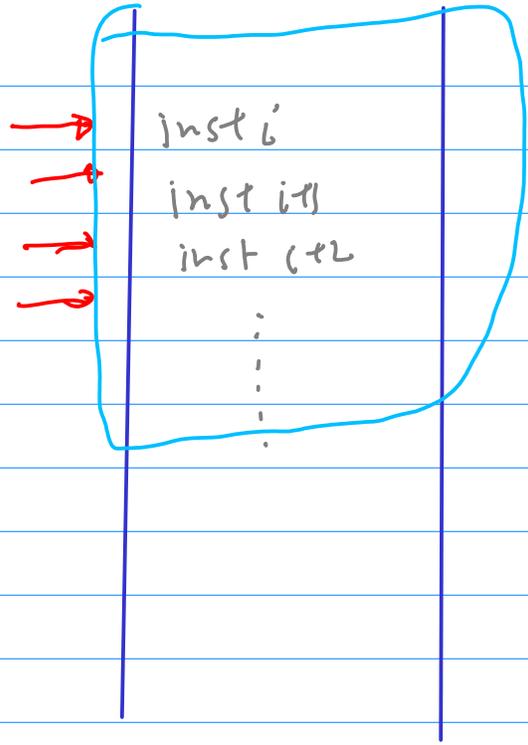
PC가 가리키는 주소에는
다음에 이어질 명령어가 있다

메모리에 있는 명령어는
순차적으로 수행하는 경우가 많다.

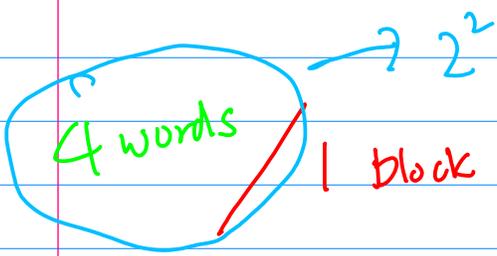


Locality of
reference

locality of reference

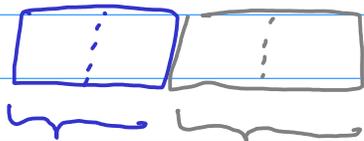


Fixed size block Main memory partition

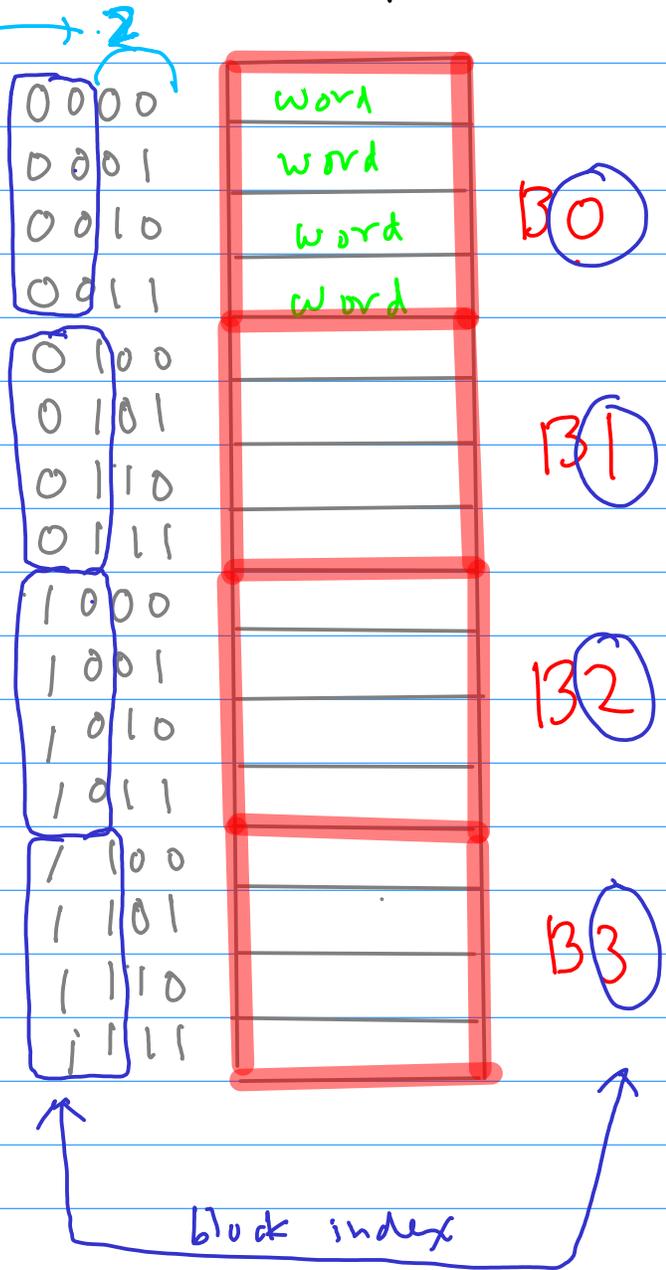


$\frac{2}{16}$ 16 words

$\frac{2}{4}$ 4 blocks

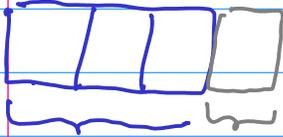


block index block 당
 2^2 words
 2-bit
 건너뛰기

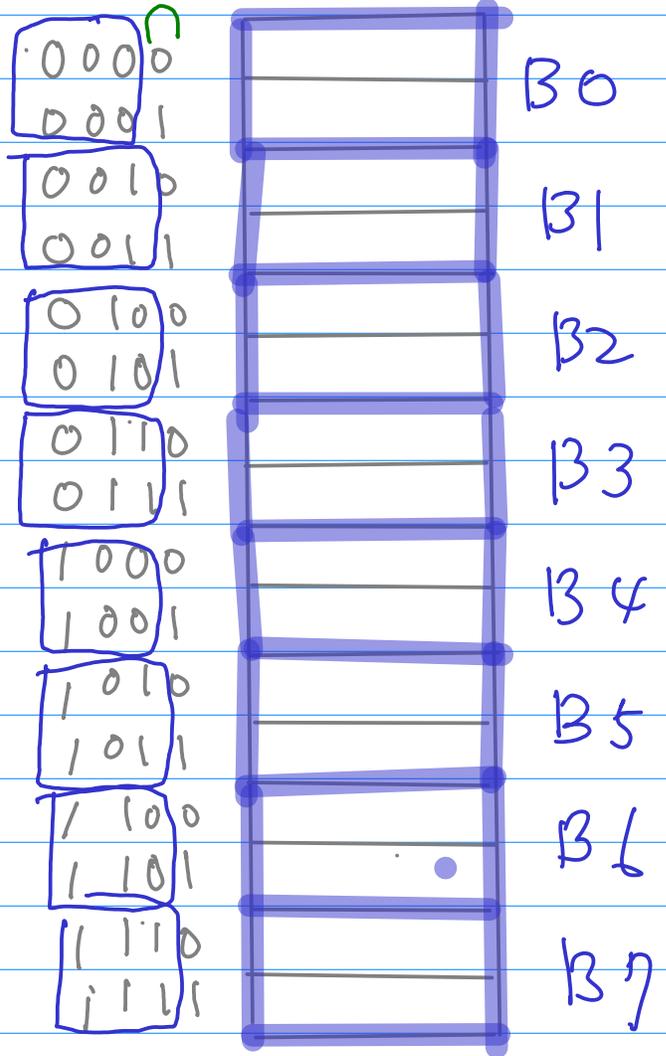


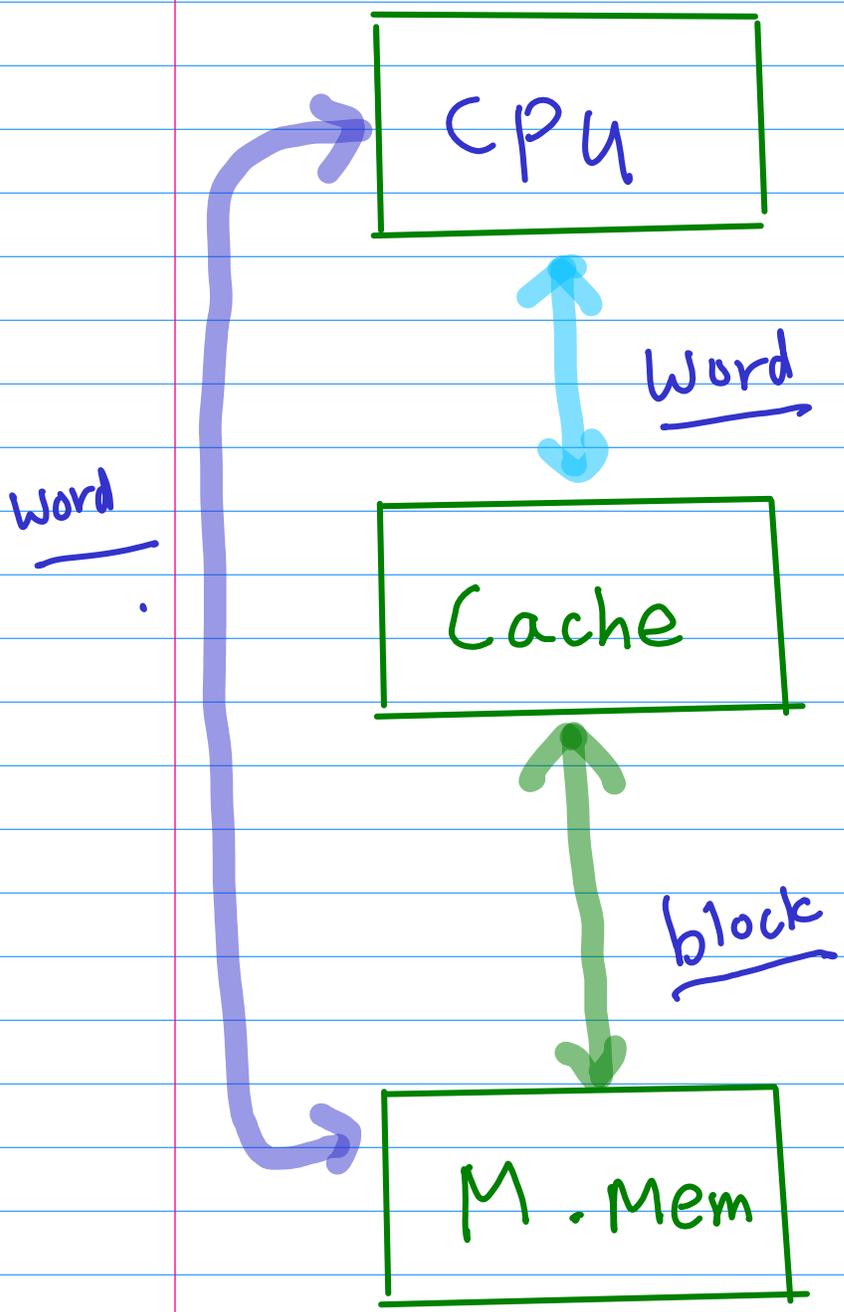
2^1 words/block.

2 words / 1 block



block index
1 block당
 2^1 words
1-bit을
걸어볼까나



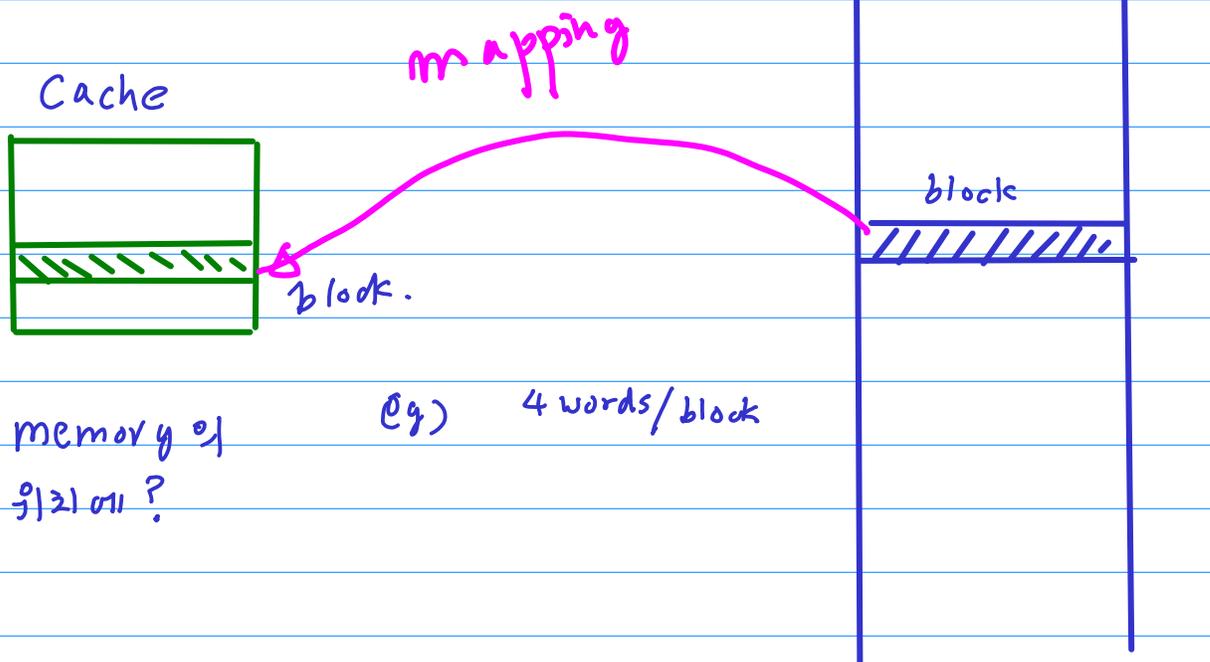


Cache Mapping Type

- ① Fully Associative Mapping
- ② Direct Mapping
- ③ k-way Set-Associative Mapping

- ① 완전 연관 사상
- ② 직접 사상
- ③ 세트 연관 사상

Mapping



Cache memory의
어느 위치에?

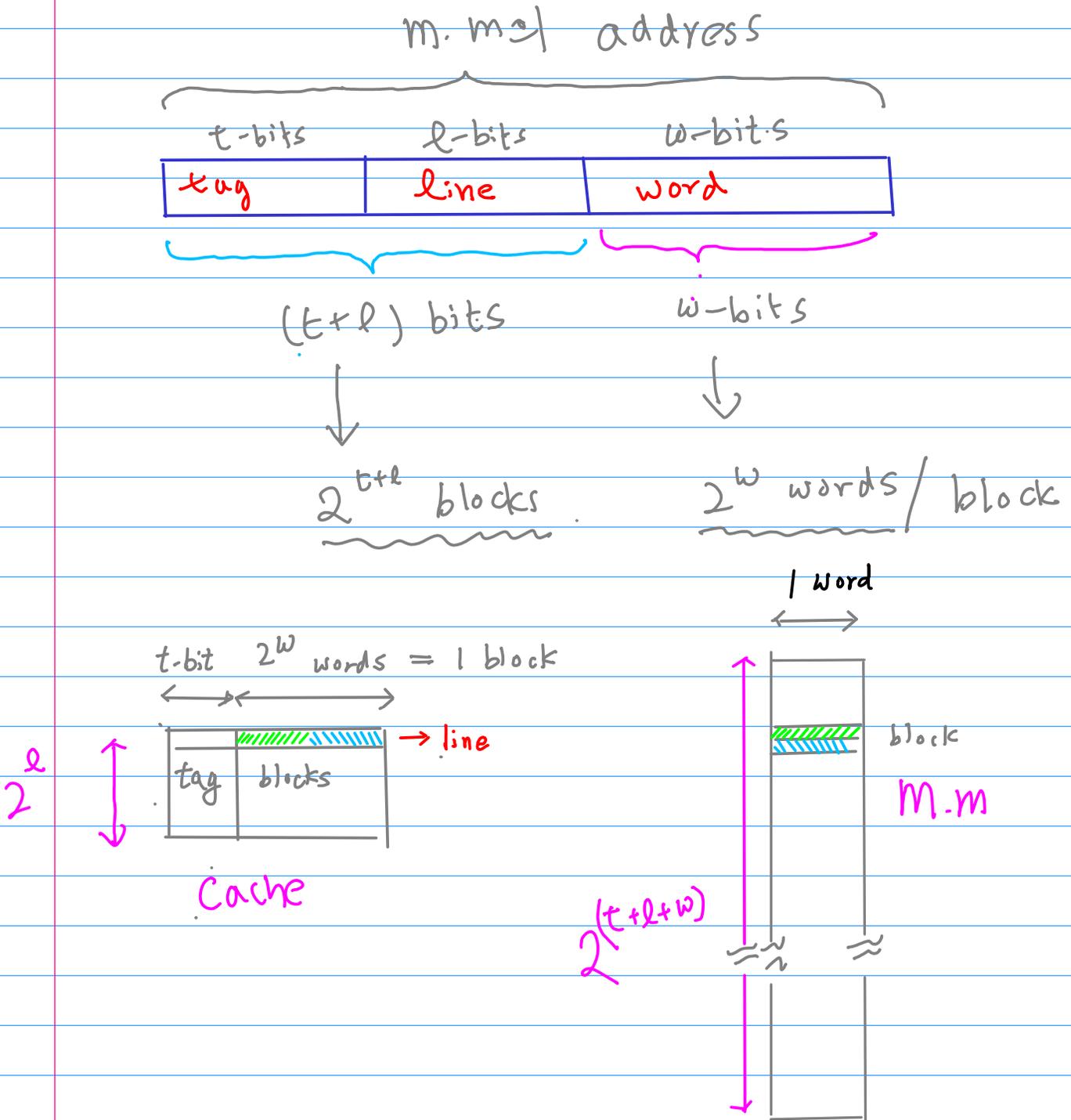
(eg) 4 words/block

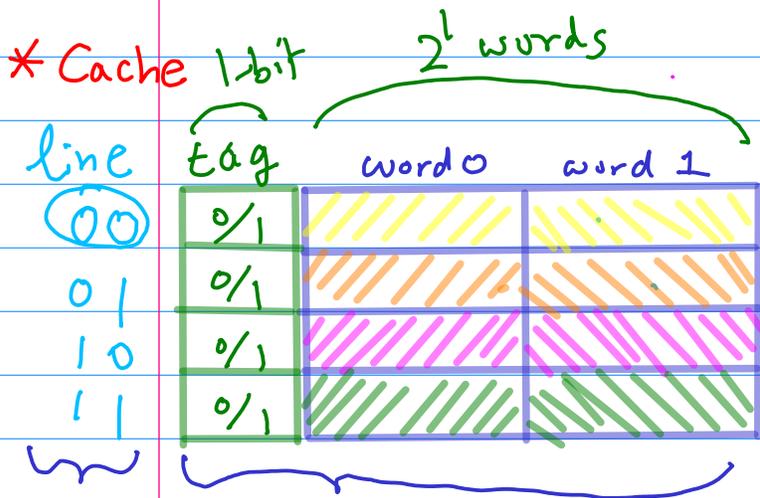
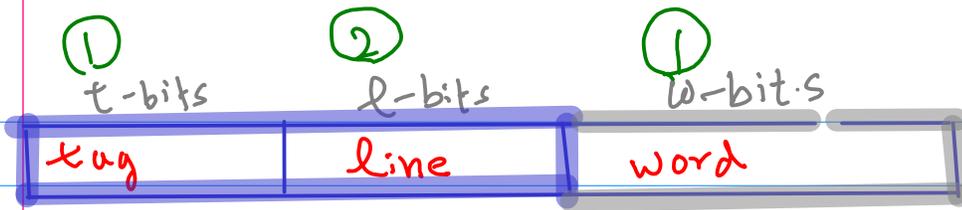
Direct Mapped Cache

m.m의 block들이

저장된 어느 **한** cache **line**으로 mapping

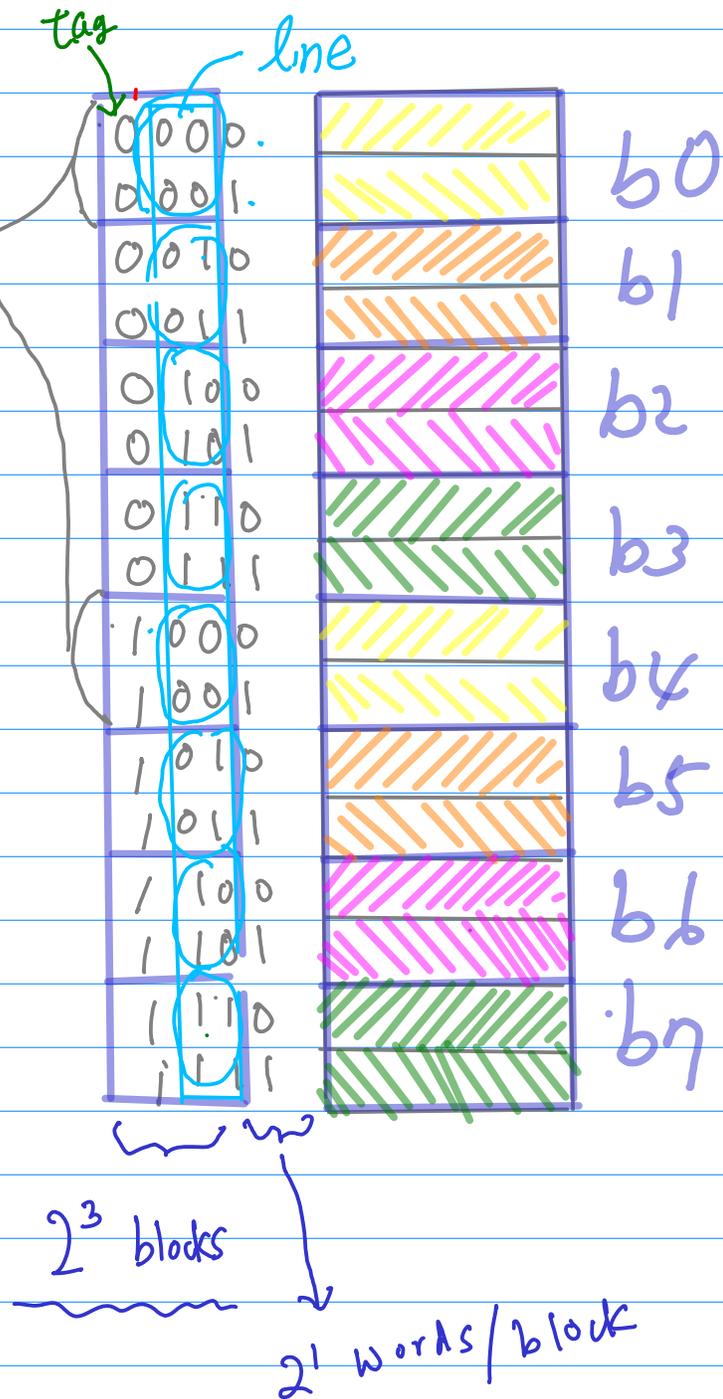
각 block들은 저장된 한개의 특정 line이 있다



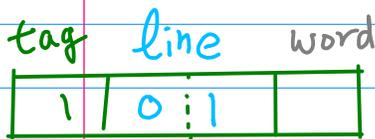
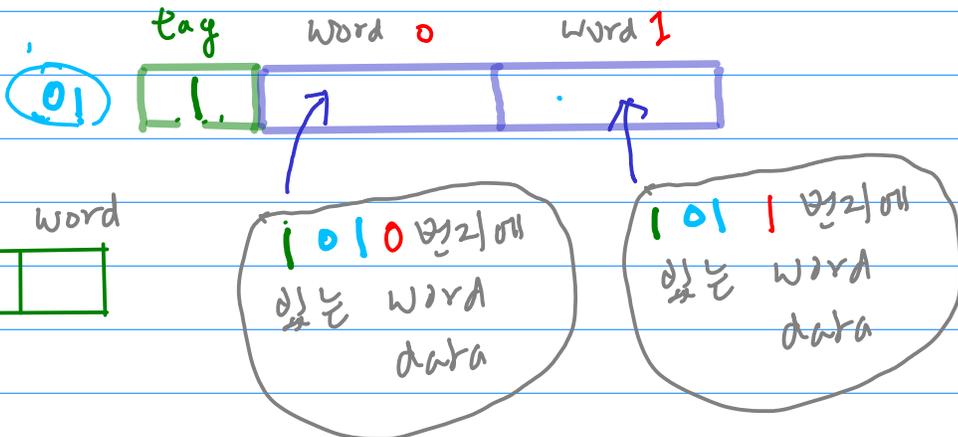
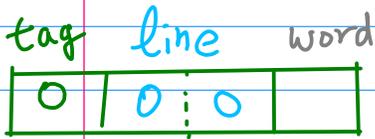
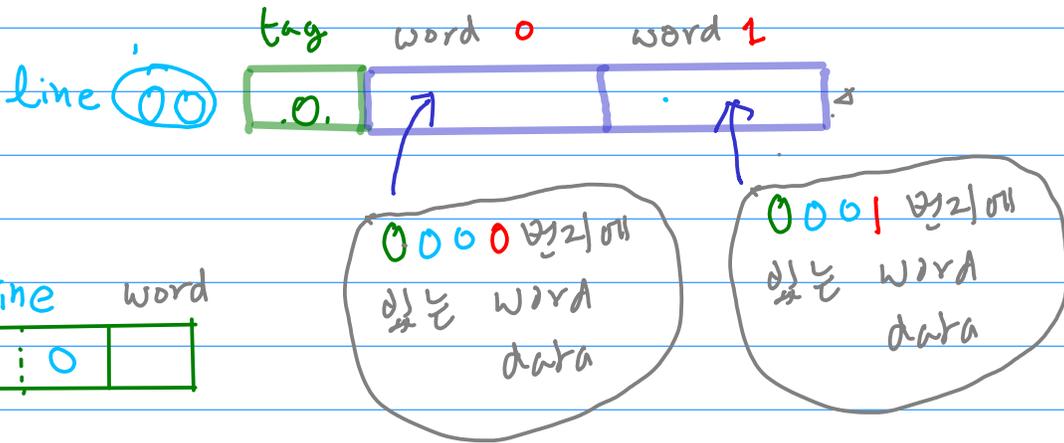
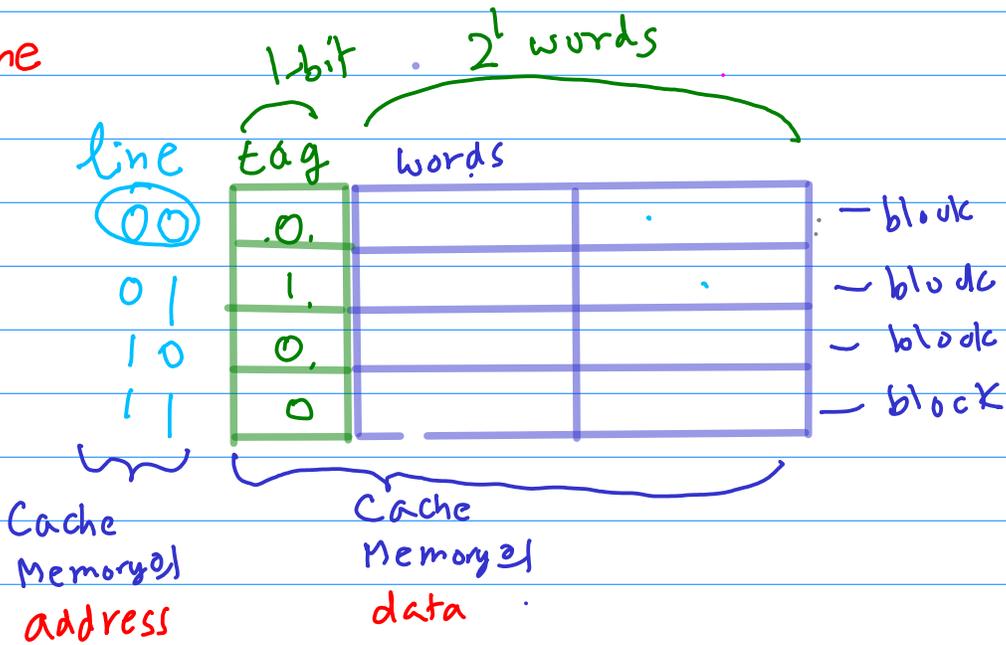


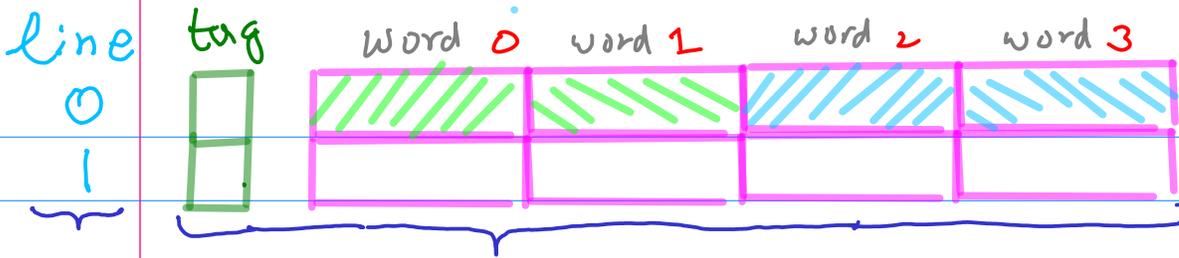
Cache Memory address

Cache Memory data



* Cache

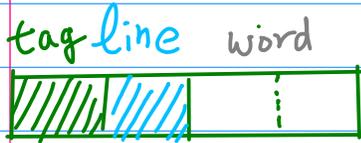




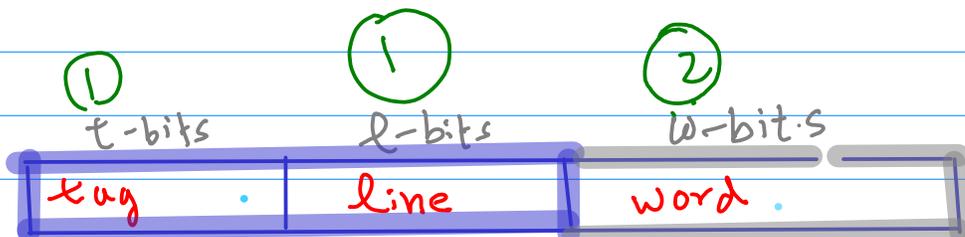
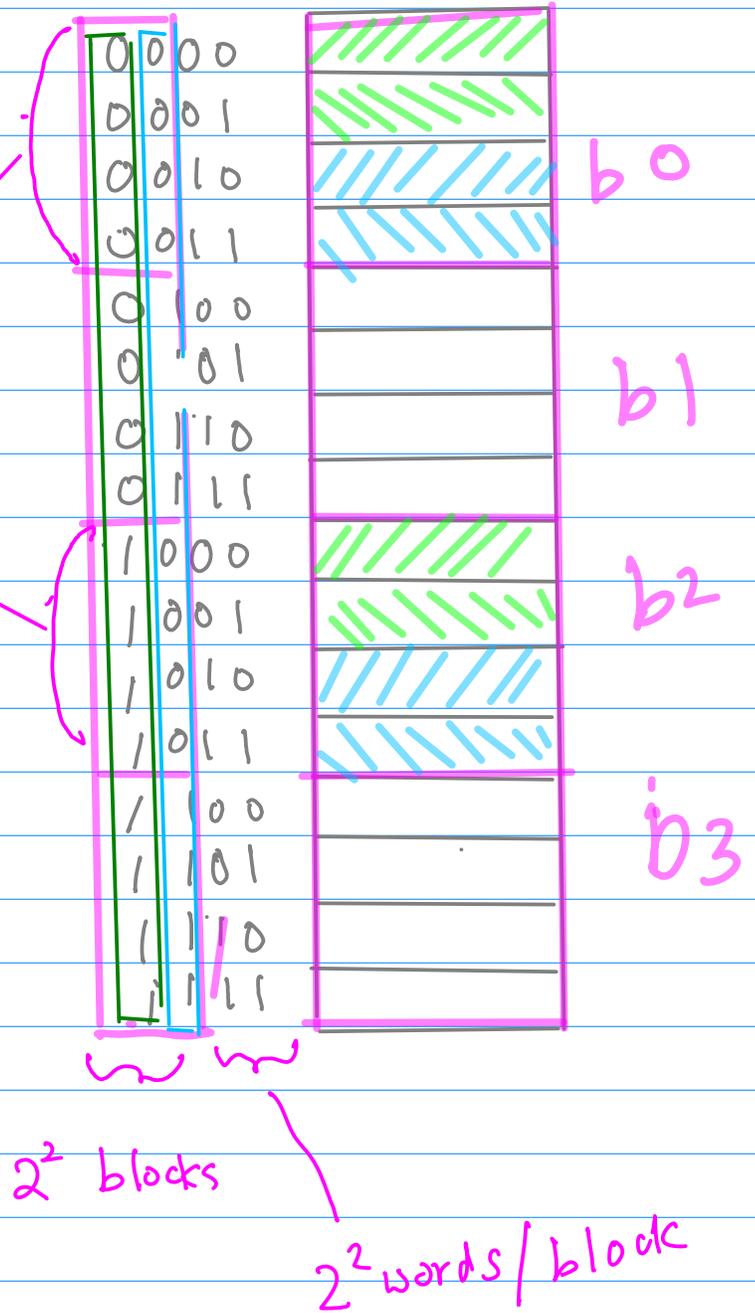
Cache Memory의 address

Cache Memory의 data

같은 line에 mapping
그러나 1개만 load 가능

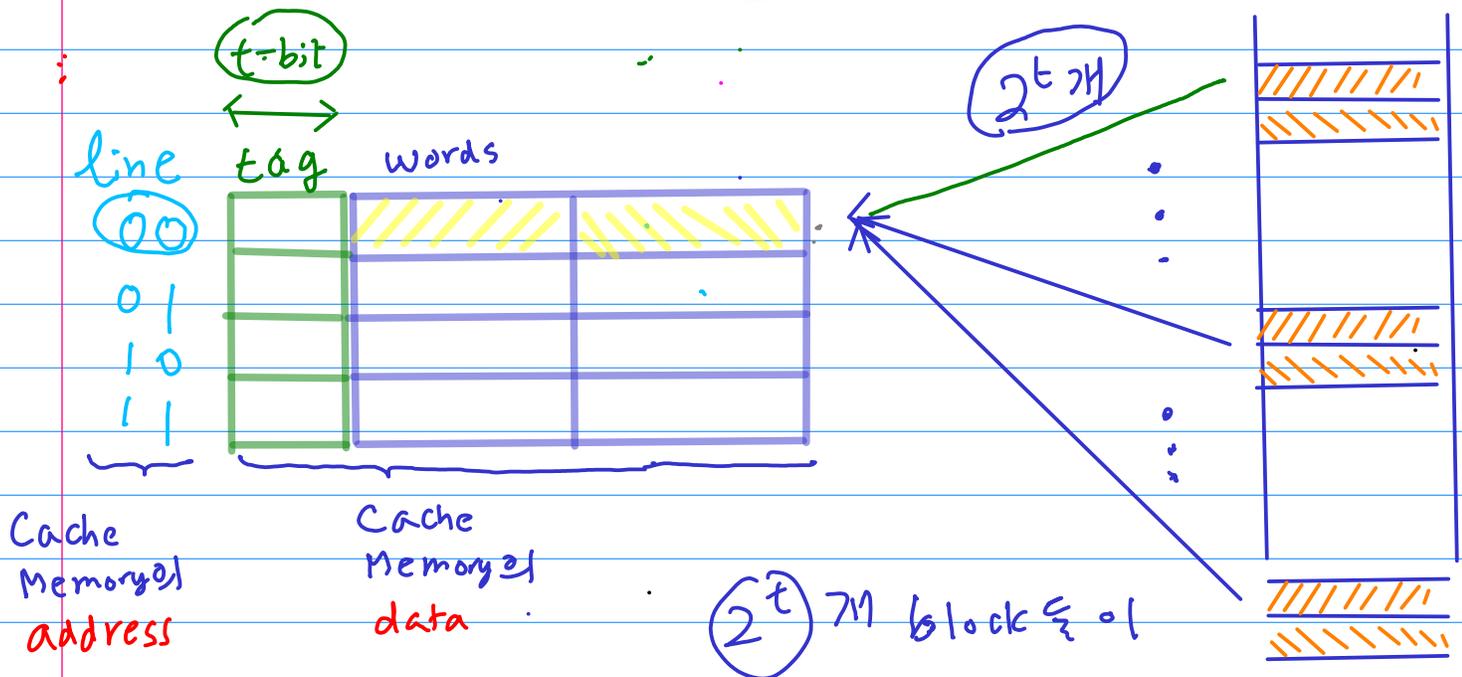


어떤 block이 load 되느냐를
나타내기 위하여
메모리 주소의 tag 필드를
cache memory에
data의 앞부분으로 저장

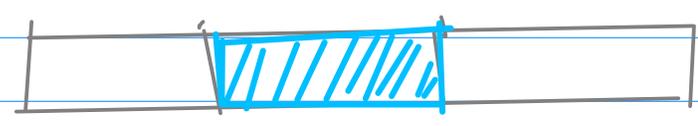
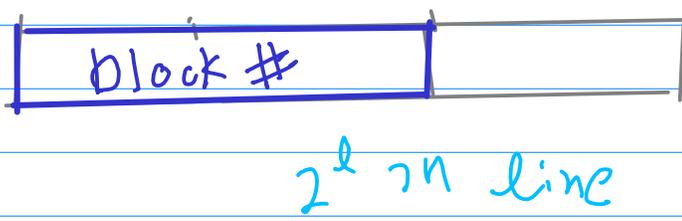
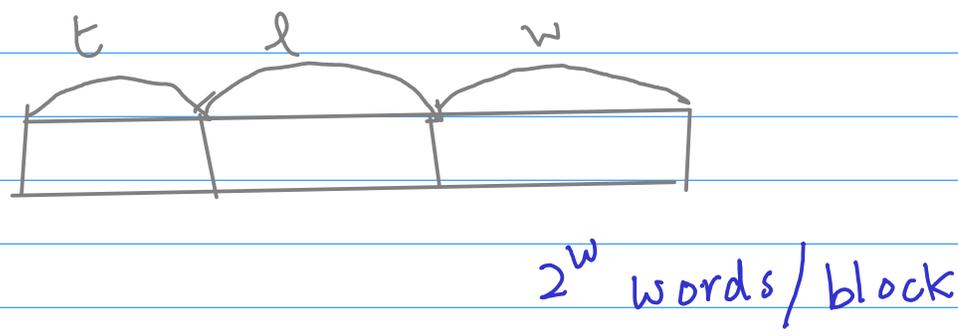


Direct Mapped Cache

각 cache line은 2^t block이 의하여 공유된다

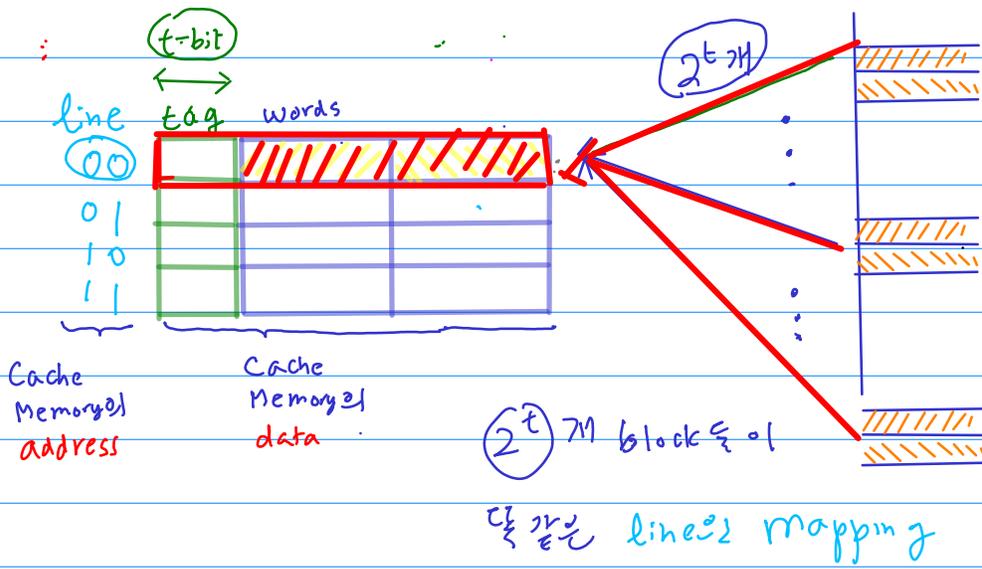


다들 같은 line에 mapping

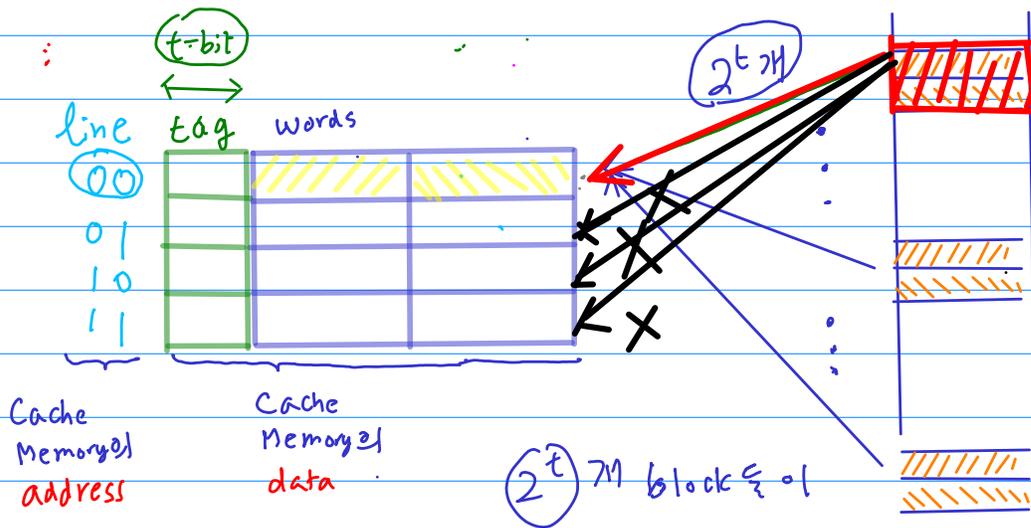


* 즉 어떤 line은 2^t 개 block이 올 수 있음.

* 1개의 block은 정해진 1개의 line에만 load 가능



* 1개의 line에는 서로 다른 2^t 개 block들이 올 수 있음.



* 1개의 block은 정해진 1개의 line에만 load 가능

이런 block이 지정된 line에 load 되었을 때

2 line에 load될 수 있는 다른 block 들라 구별하기 위해서

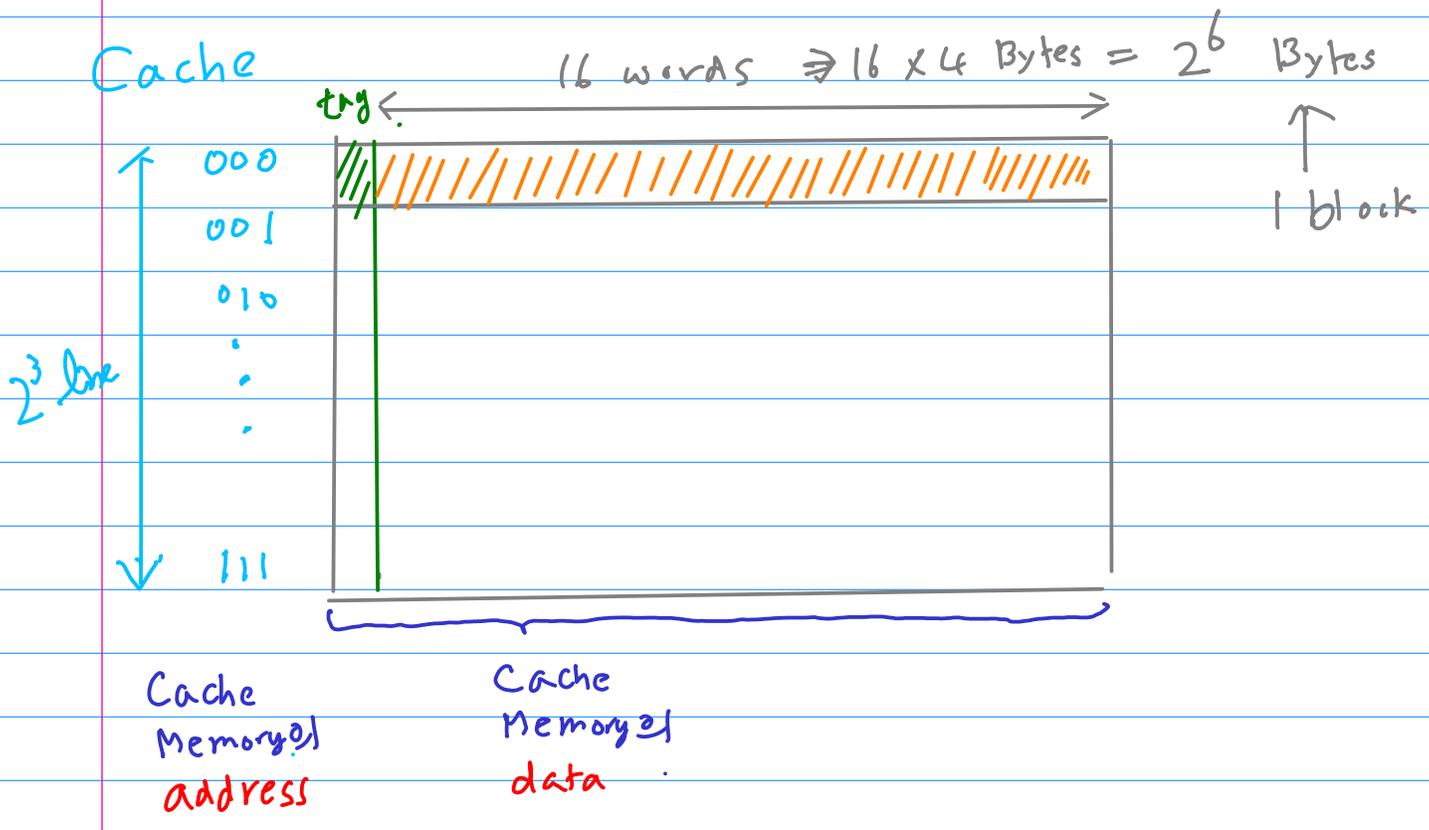
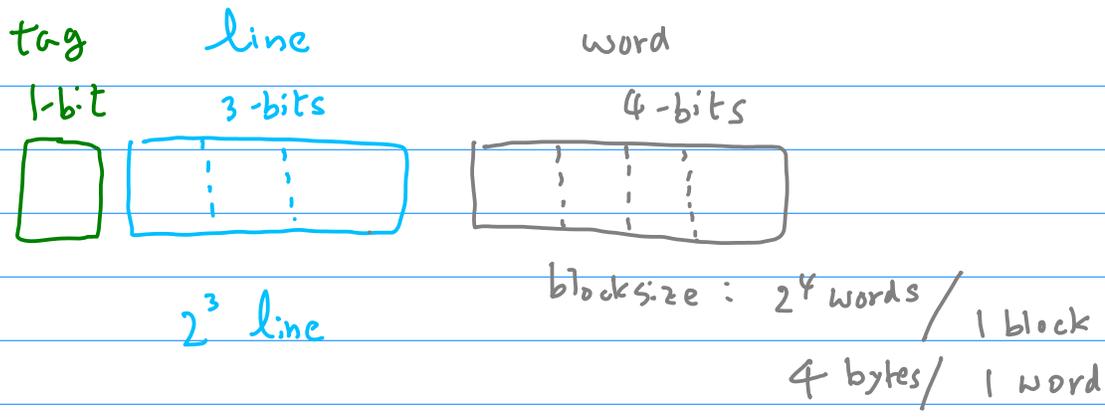
t-bit tag를 cache memory의 data로 저장한다.

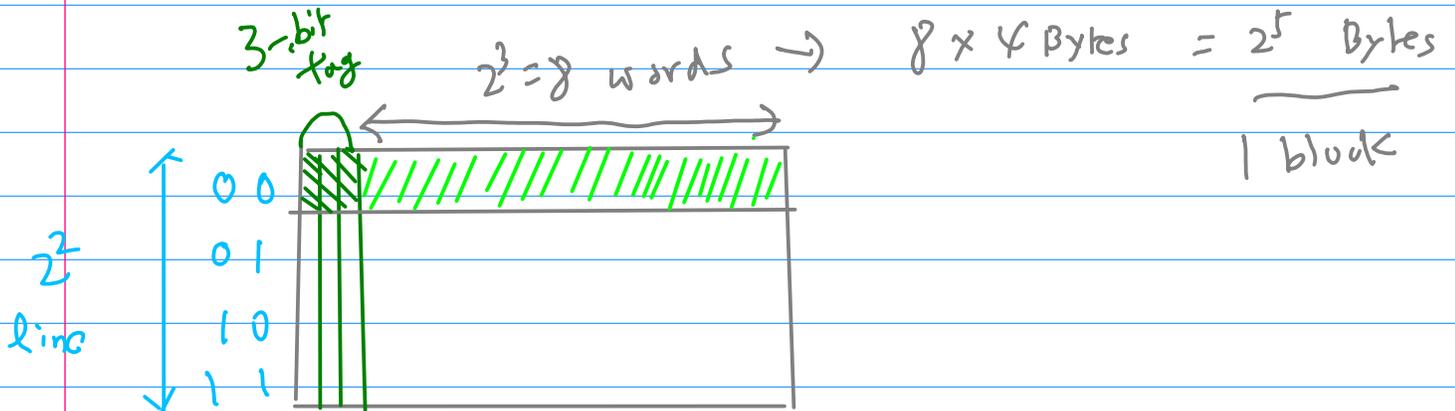
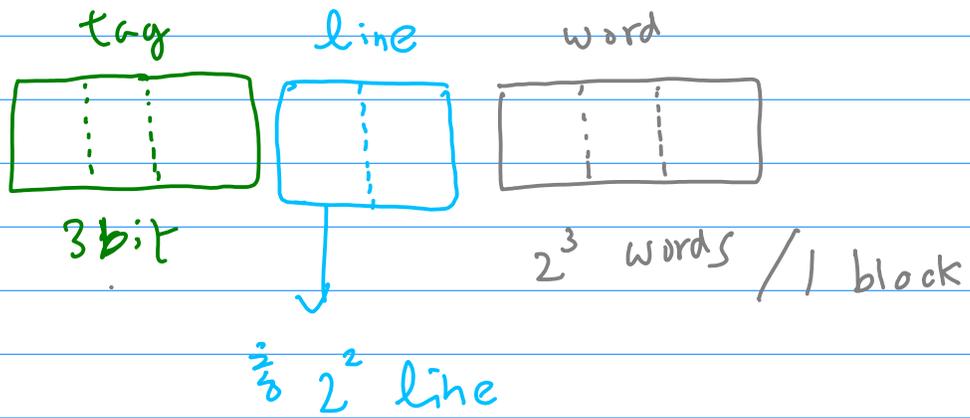
tag가 1-bit 이면
2-bit

2^1 개 block들이 같은 line에 ...
 2^2 개 block들이 같은 line에 ...

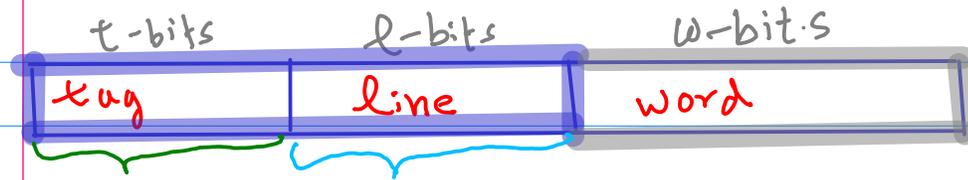
각 cache line은

{ t-bit tag
2^w개 word로 구성된 block은
cache memory의 data로 저장한다





* m.m address



Cache Memory의 data를 취급함.

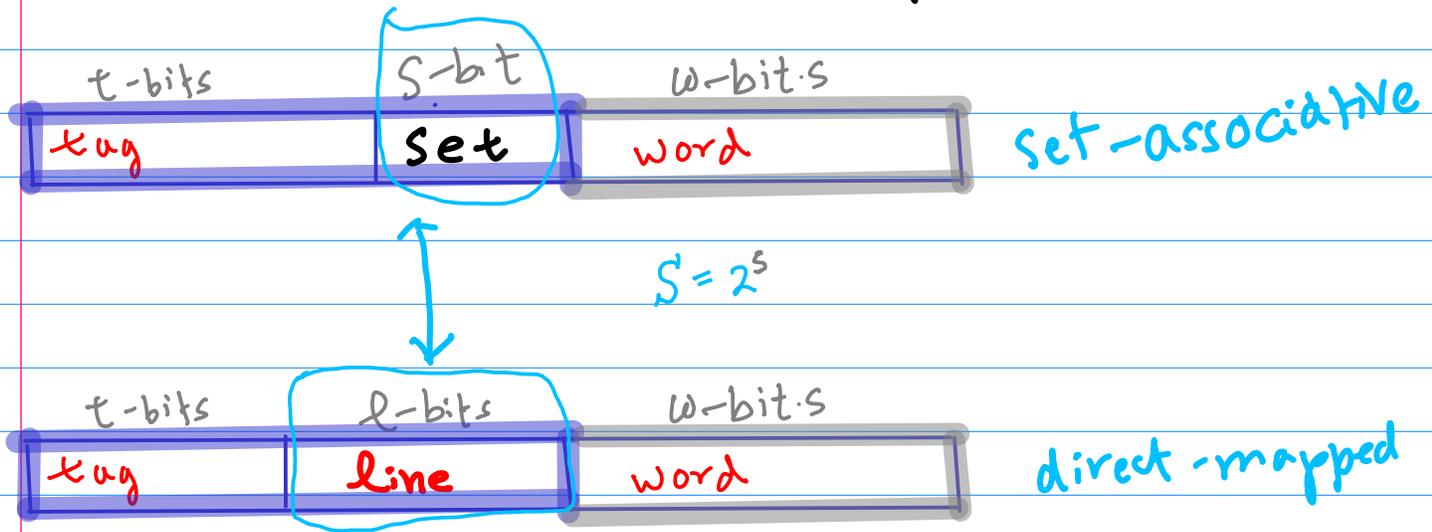
Cache Memory의 address를 사용.

* Cache

2^l 개의 line이 있다.

k-way

* Set - associative mapping



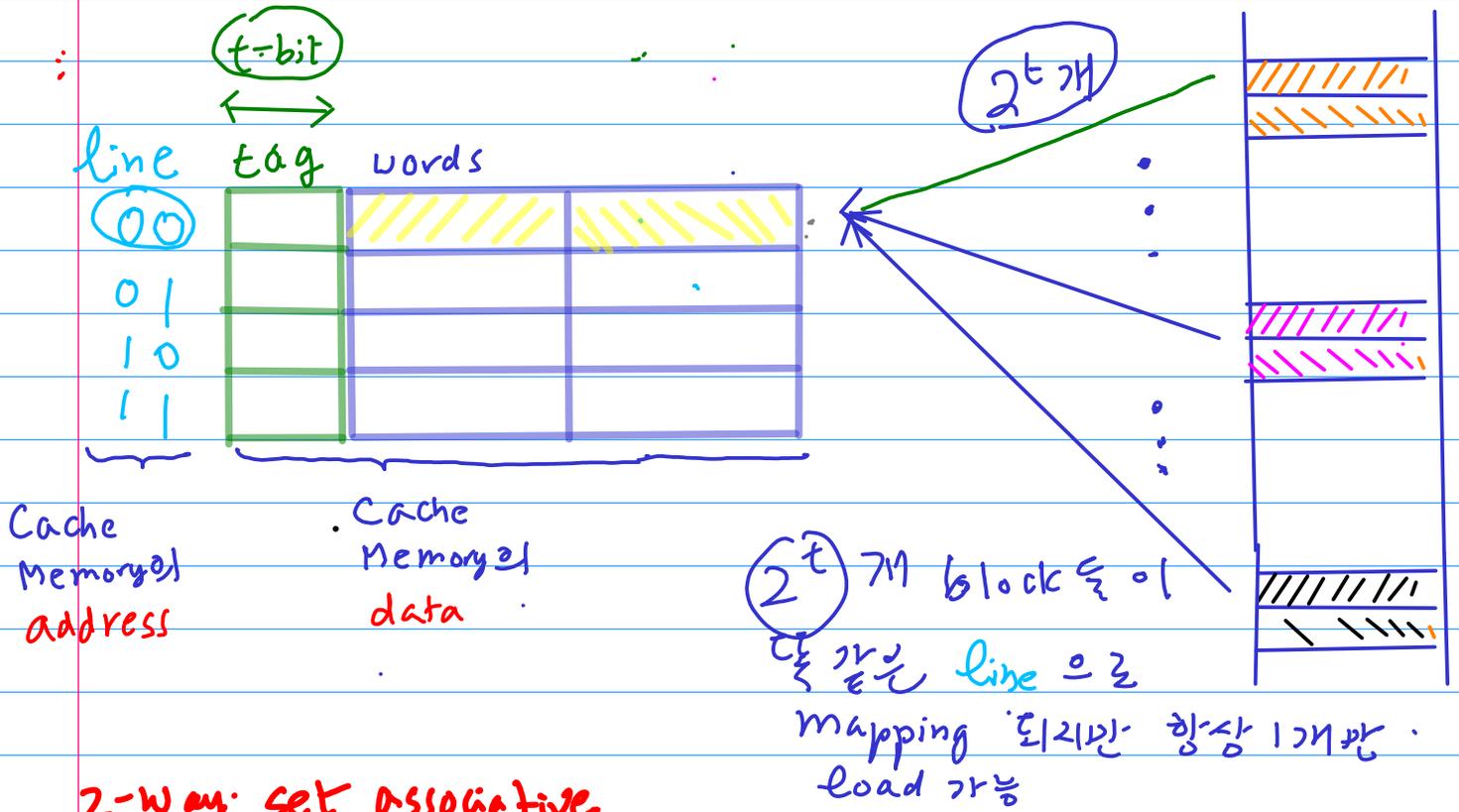
{ direct mapping
fully-associative mapping }

$$2^k \times k \Rightarrow$$

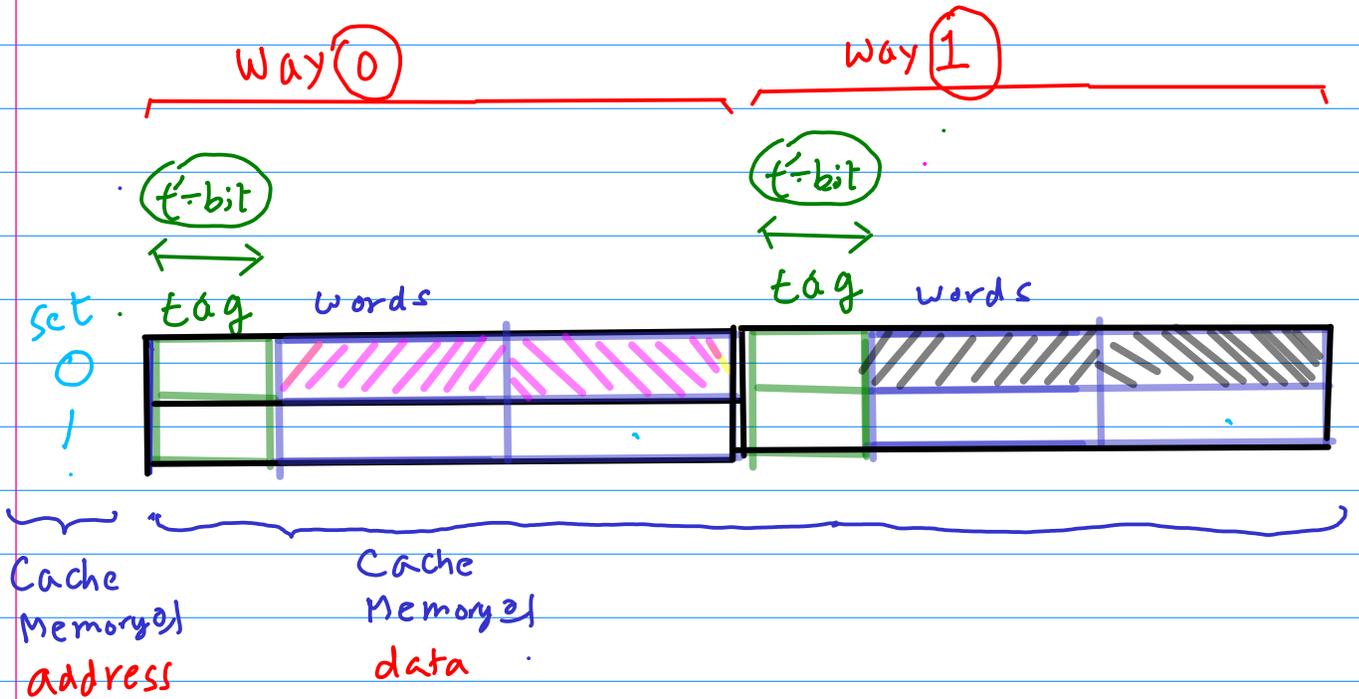
Set-associative mapping

direct mapped

각 line당 1개의 block만 load



2-way set associative

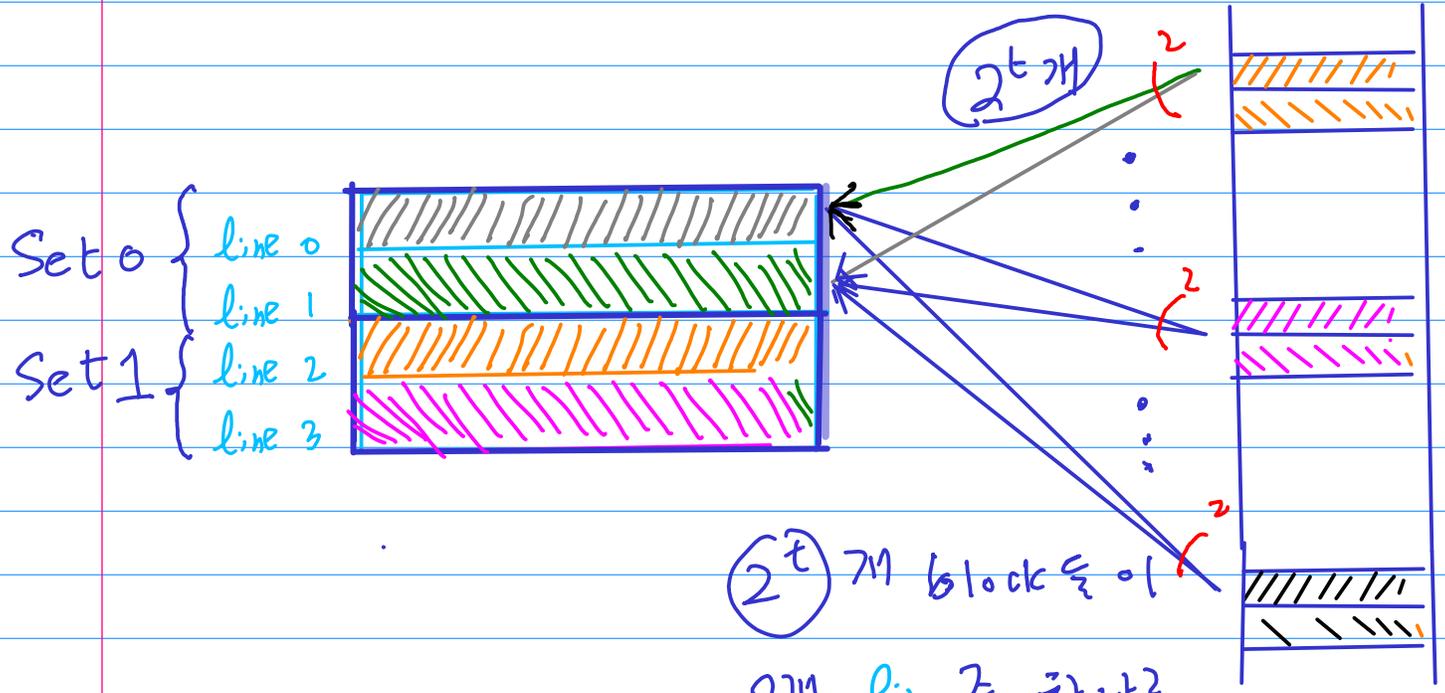


2-way set associative에서는

각 set당 2개의 block이 load 될수 있음

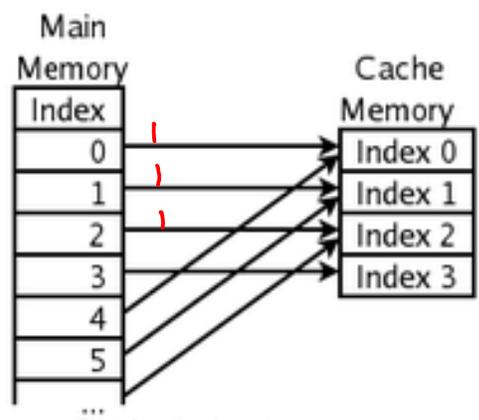
2-Way Set Associative

Set 당 2개의 block load 가능



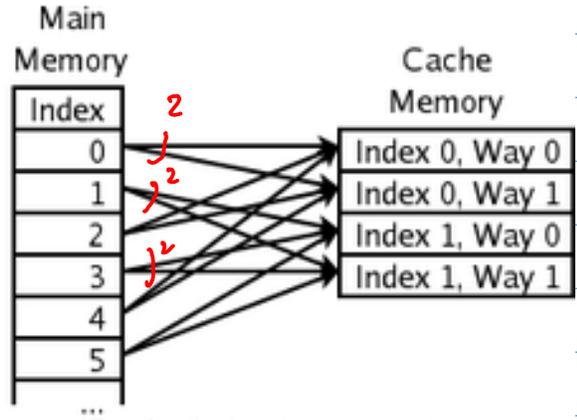
2개 block들이
2개 line 중 하나로
mapping 되고 항상 2개
load 가능

Direct Mapped Cache Fill



Each location in main memory can be cached by just one cache location.

2-Way Associative Cache Fill



Each location in main memory can be cached by one of two cache locations.

원래 cache에

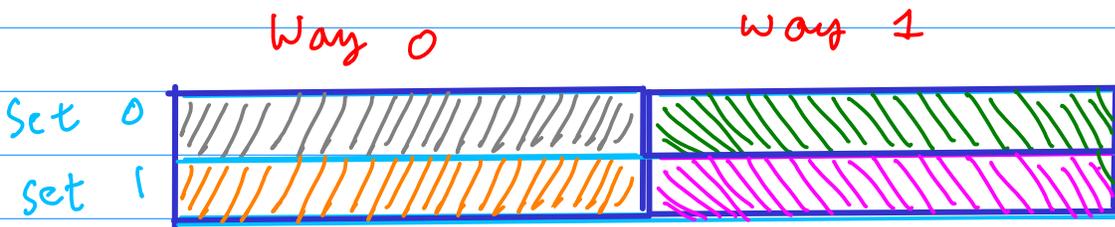
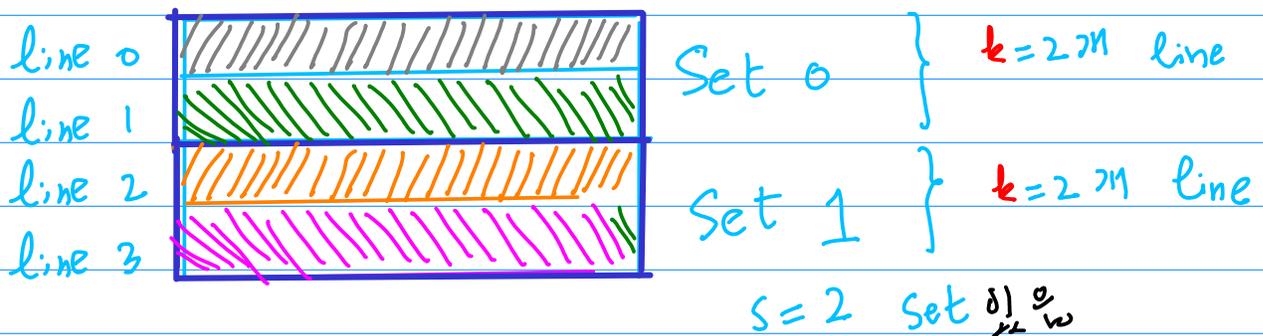
l 개의 line 있다고 가정.

S 개의 set 있음

각 set는 k 개의 line을 저장할 수 있음.

총 $S \cdot k$ 개의 line들이 있다.
 $= l$

원래 cache에 $l = 4$ lines 있음



2-way set associative mapping

전체 Cache size가 32 bytes }
 line size는 4 bytes }

무슨 줄 $32/4 = 8$ 개 line이 있다.

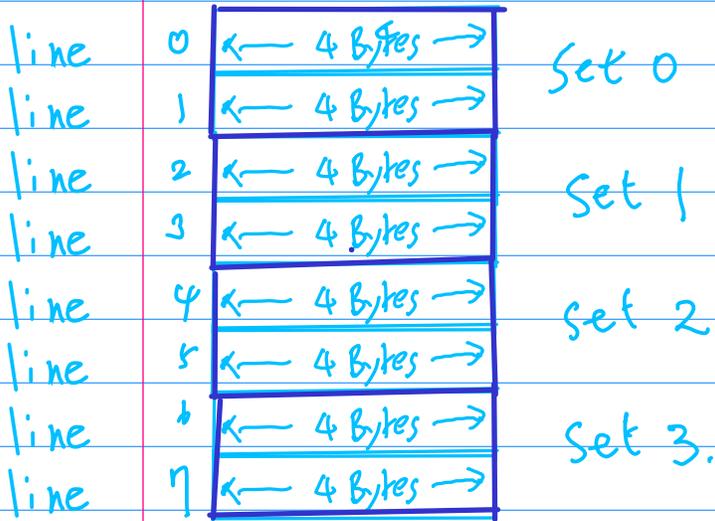
만약 1 set 당 2개의 line을 저장하면

$$S \cdot k = 4 \cdot 2 = 8 = L$$

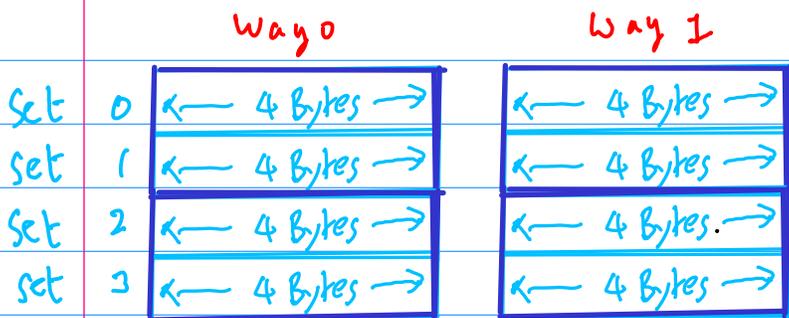
$S = 4$ 개의 set가 존재

Cache는 $S=4$ 개의 set들로 구성되고

각 set는 2개의 line이 있다

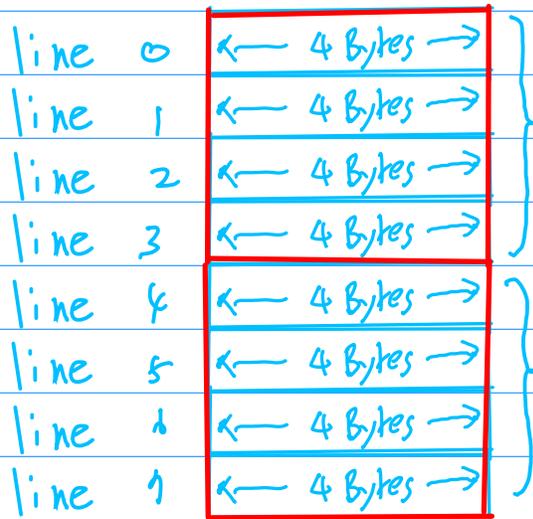


2-way
 set-associative
 mapping



* tag도 block과 함께 cache memory의 data로 저장된

2개의 tag 저장



$k=4$

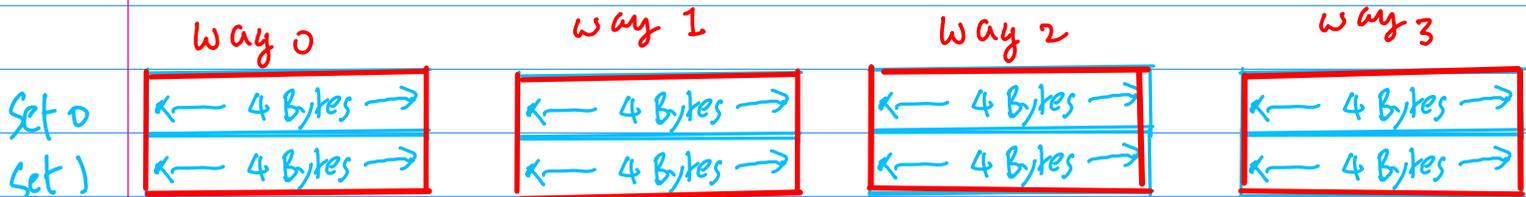
4-way
set associative
mapping

Set 1

$s=2$

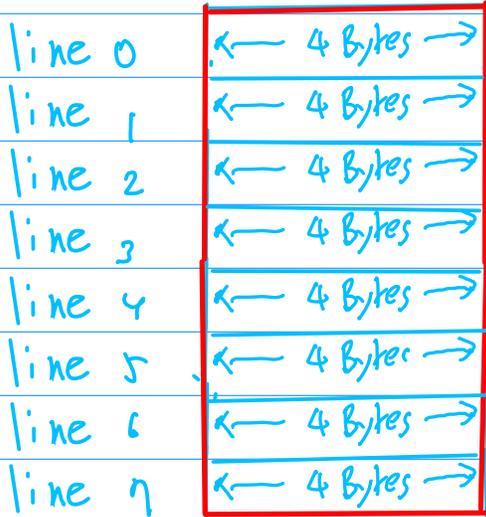
각 Set에 k 개씩 line이 있으면

k -way set-associative mapping



* tag도 block과 함께 cache memory의 data로 저장됨.
2way는 tag 저장

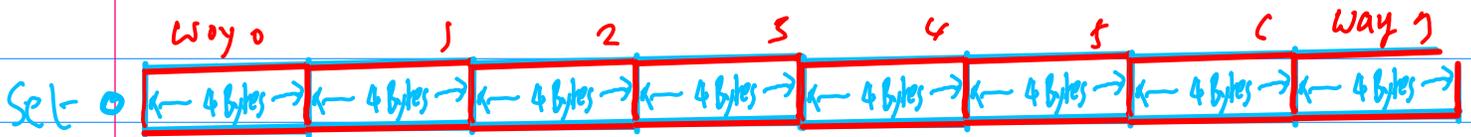
Fully Associative Mapping



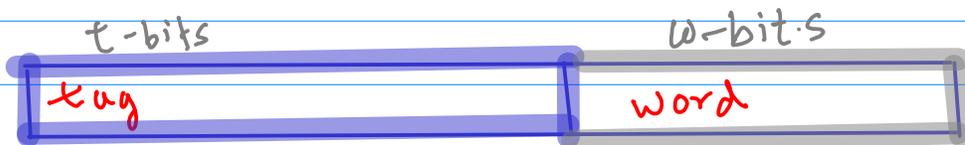
Set 0

only 1 set

$k = l$ lines



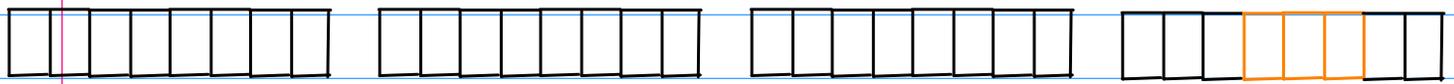
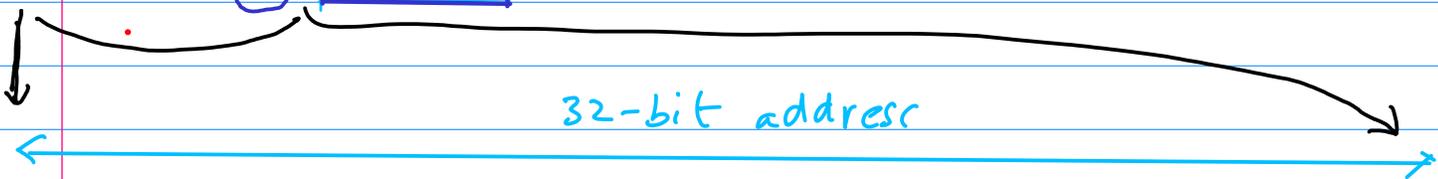
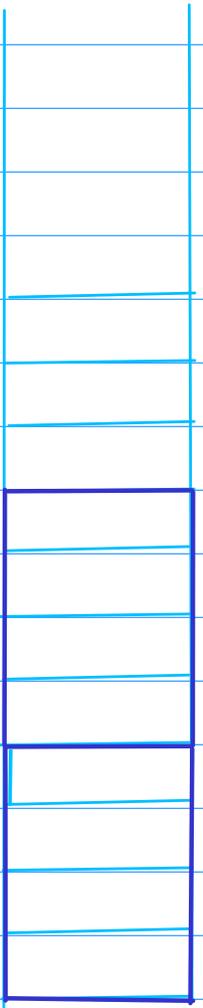
* Fully associative mapping에서 는 tag가 address에서 하위 w -bit를 제외한 전체가 된다



* tag도 block과 함께 cache memory의 data로 저장된다.
 2^w에서 tag 저장

1 Byte
↔

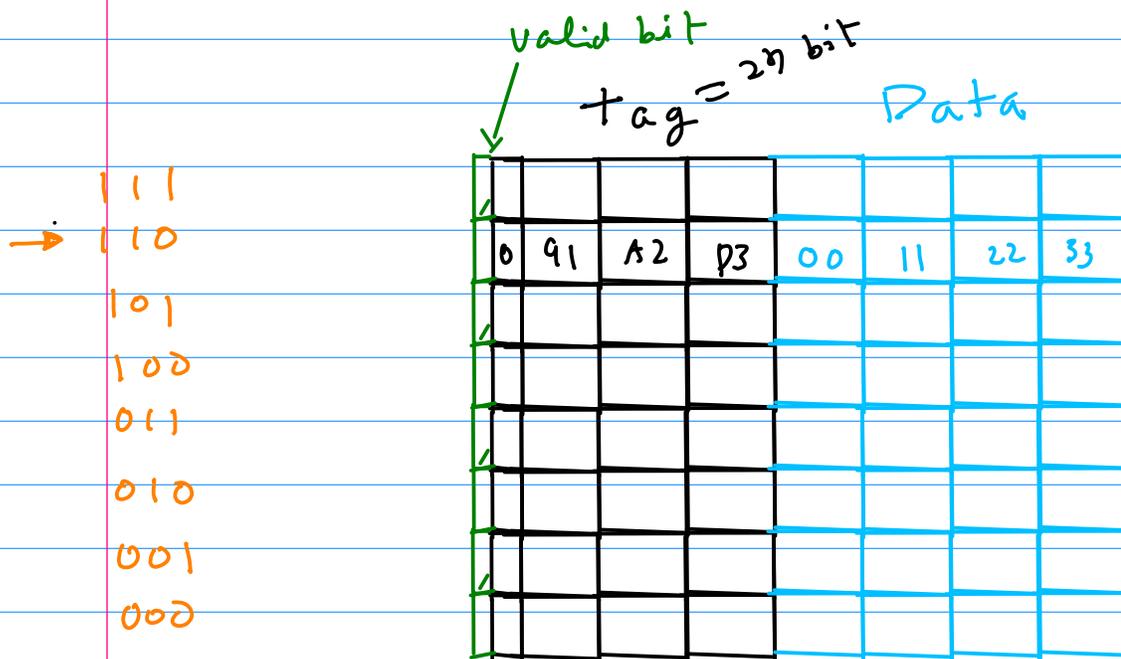
- 0 ... 0 0 1 1
- 0 ... 0 0 1 0
- 0 ... 0 0 1 0 1
- 0 ... 0 0 1 0 0
- 0 ... 0 0 0 1 1
- 0 ... 0 0 0 1 0
- 0 ... 0 0 0 0 1
- 0 ... 0 0 0 0 0



$$\text{tag} = (32 - 2 - 3) = 27$$

- 0 0 0
- 0 0 1
- 0 1 0
- 0 1 1
- 1 0 0
- 1 0 1
- 1 1 0
- 1 1 1

Direct Mapping



a 32-bit address

data

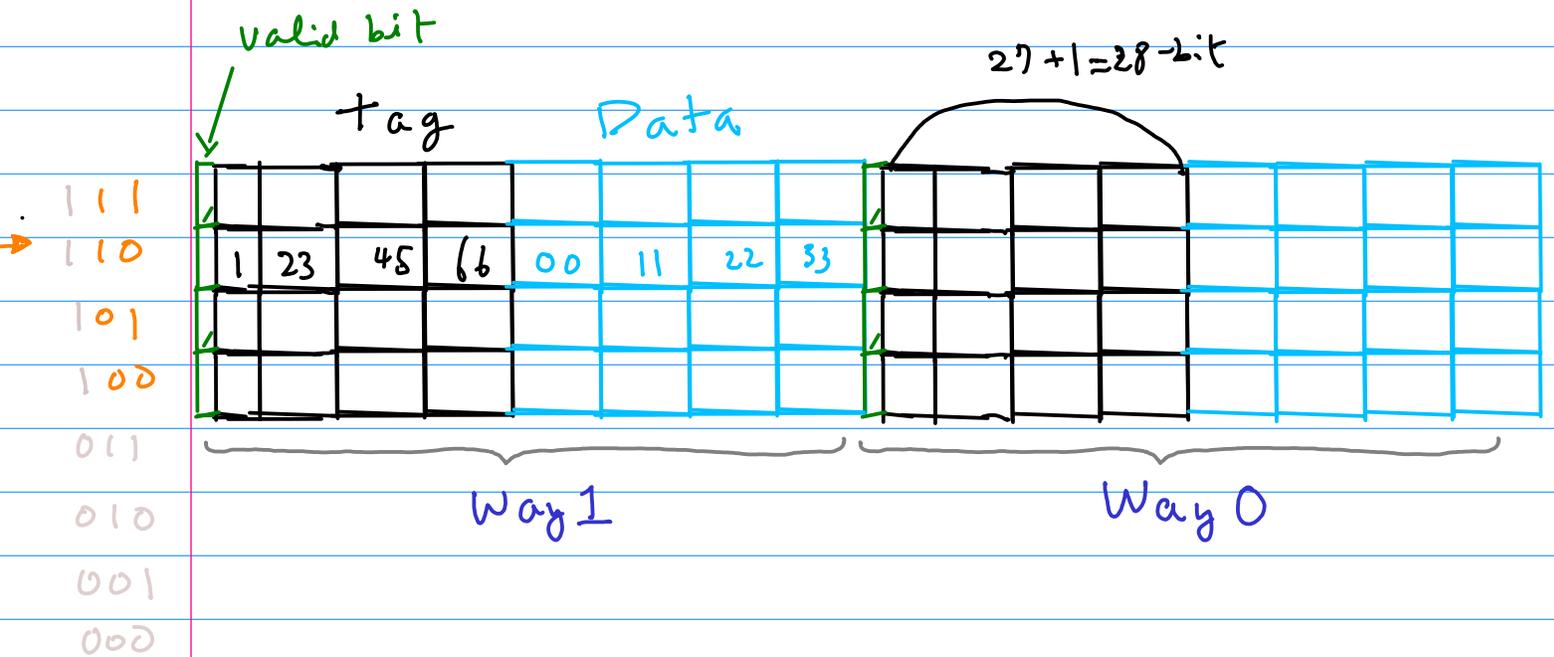
0x 12 34 56 18 00 11 22 33

01111000

0001 0010 0011 0100 0101 0110 0111

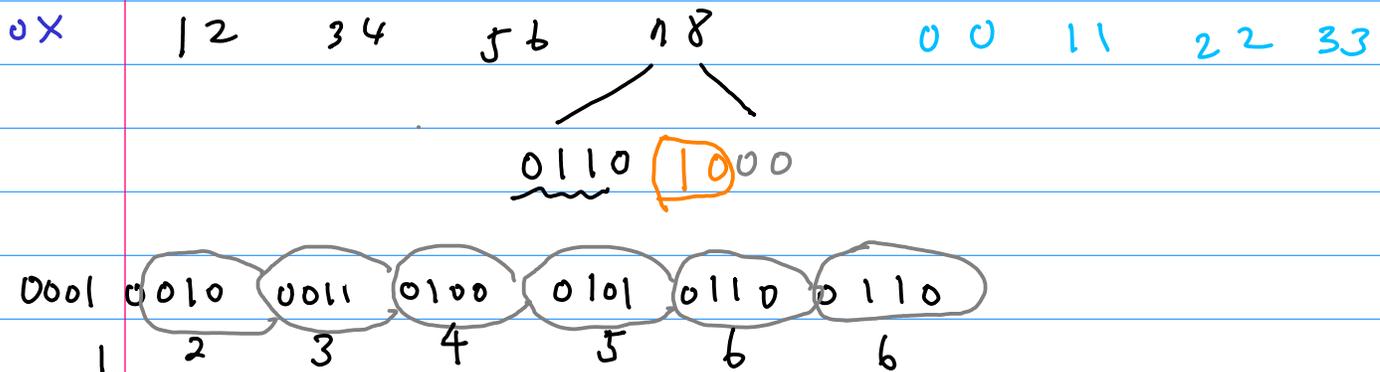
0 9 1 A 2 D 3

2-Way Set associative mapping



a 32-bit address

data



Direct

tag 27-bit

line 3-bit

word 2-bit

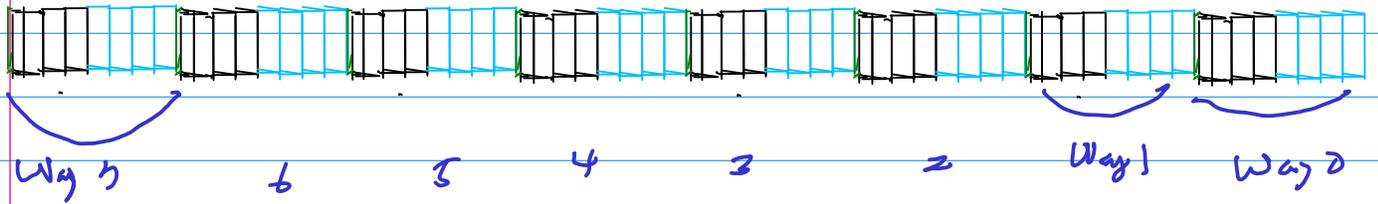
2-way

tag 28-bit

set 2-bit

word 2-bit

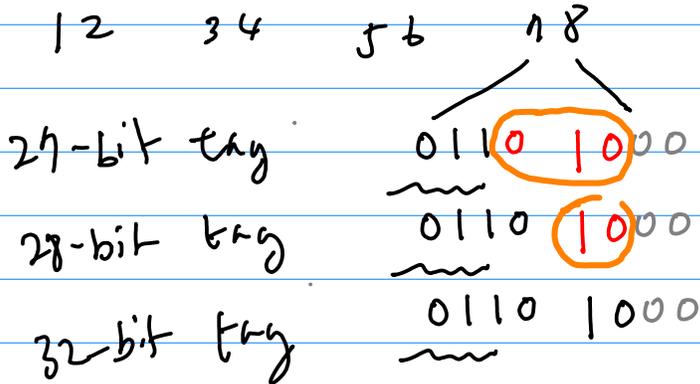
Fully Associative



a 32-bit address

data

0x

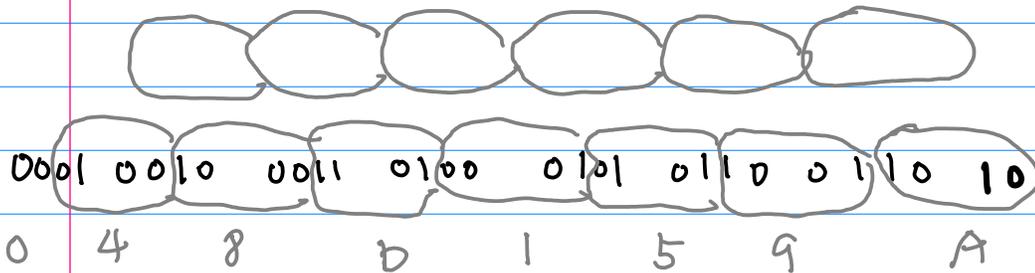


00 11 22 33

Direct Mapping

2-Way Set-Assoc

Fully Assoc

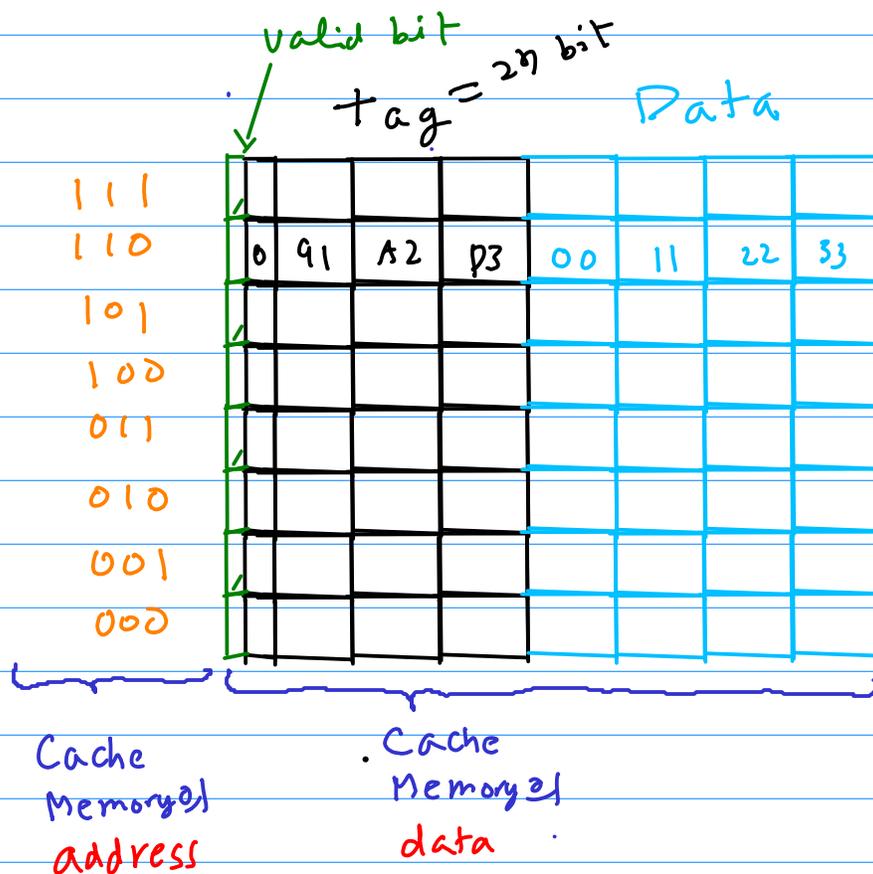


$$27 + 1 + 2 = 30$$

Address Matching Logic

no address decoding X
instead, address matching

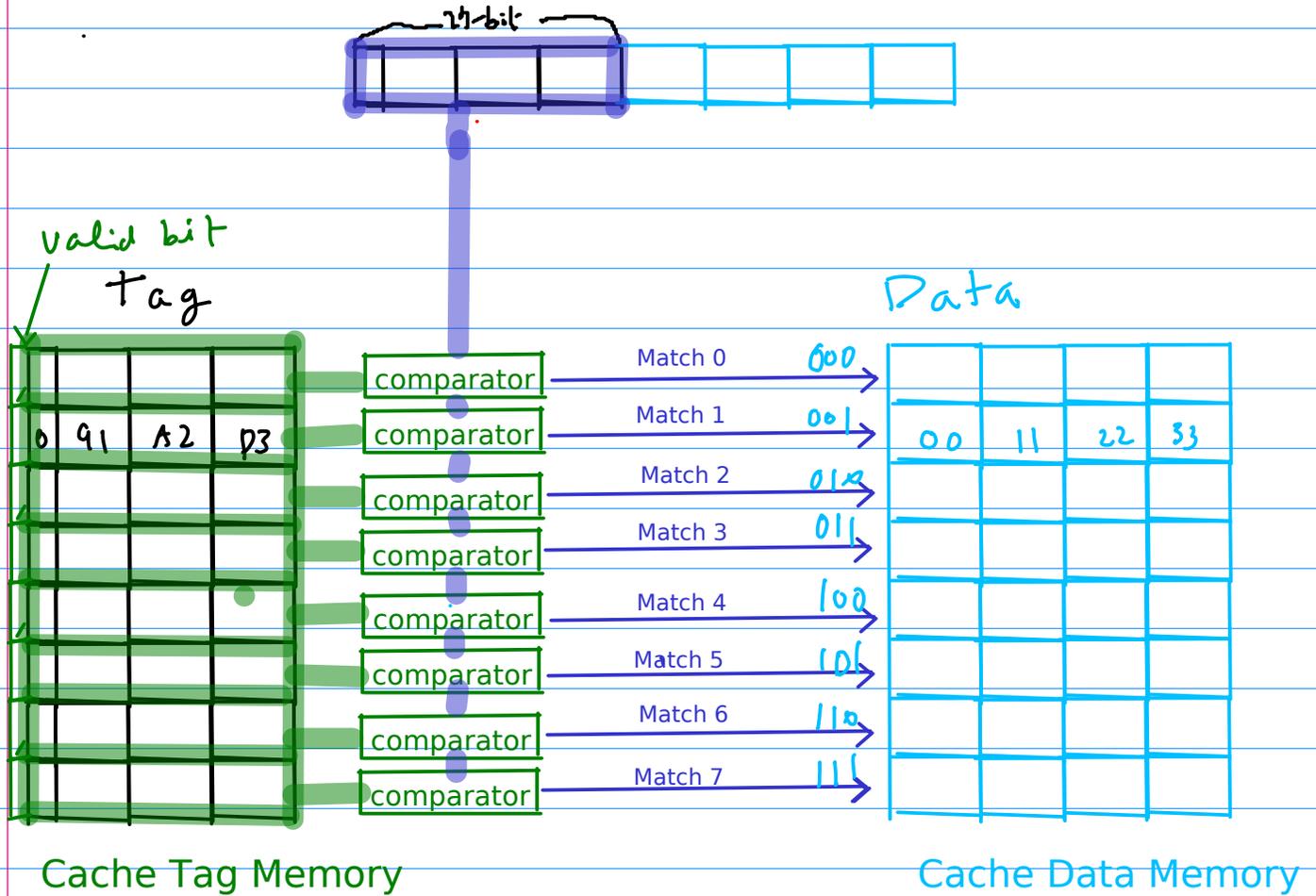
⇒ Content Addressable Memory (CAM)
Associative Memory



↓
one of these are selected
not by address decoding logic
but by address matching logic

Address Matching

CPU's address

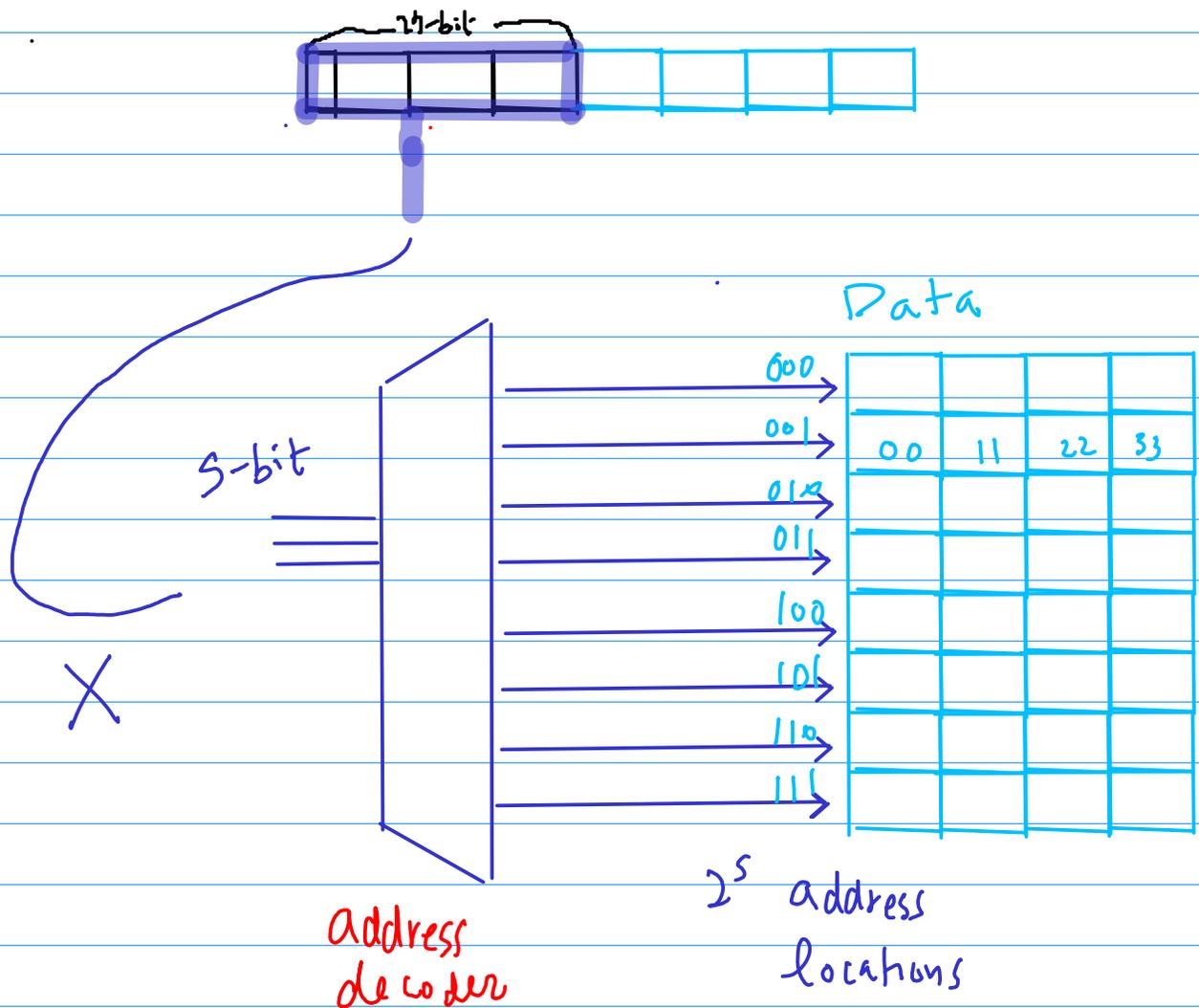


the tag field of a CPU address is matched with each content of tag memories.

Content addressable

Address Decoder

CPU's address



In cache memories, address decoding is not used, address matching is used.

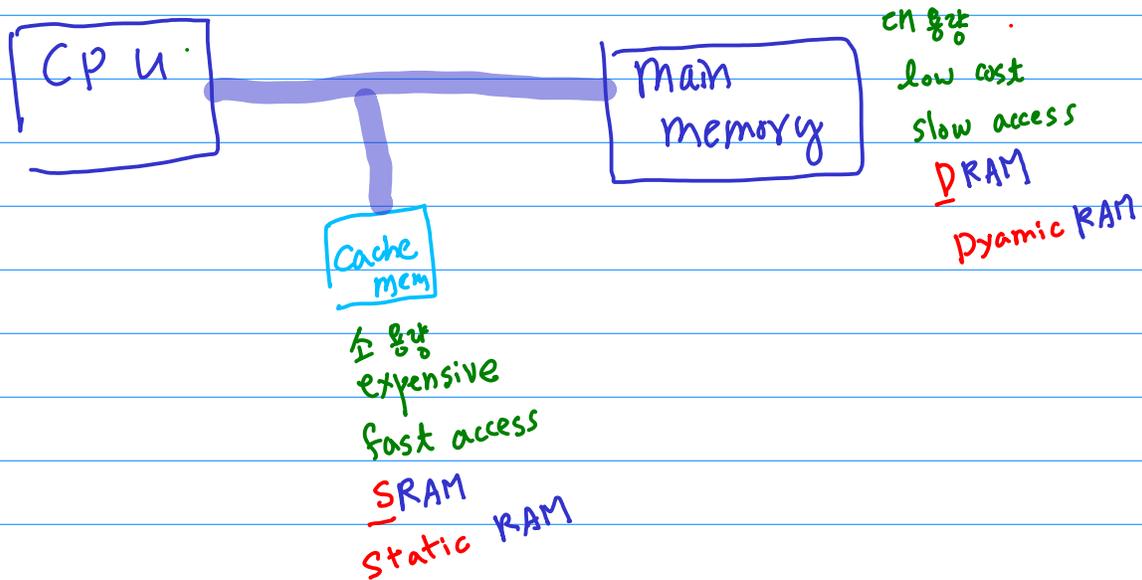
Cache Memory

Speed of CPU $>$ Speed of Main Memory
빠르다

instruction 수행시간 $<$ memory access time 길다

CPU가 연산결과를 memory에 쓰기 위해서
CPU가 idle하면서 기다리는 상황이 있을수 있다 \rightarrow performance 성능저하

CPU와 main memory 간의 속도차이 인함
성능저하를 방지하기 위해서 cache를 사용한다.



CPU가 Data를 Read 하는 경우

먼저 Cache를 찾아보고

- 원하는 data가 cache에 있으면 (Cache hit)

cache에 있는 data를 읽어본다

느린 main memory access 없이

- 원하는 data가 cache에 없는 경우 (Cache miss)

main memory를 access해서

CPU가 읽어오려 했던 data 뿐만 아니라

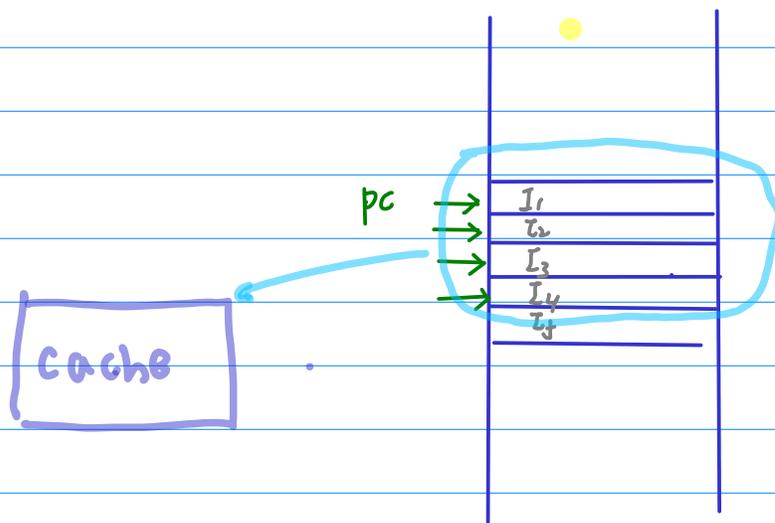
그 주변 data block을 cache에 Load 한다

Locality (지역성)

CPU가 main memory의 특정 부분 (지역) 이 위치한

program code나 data를 빈번히 / 자주 / 집중적으로

access 하는 현상을 의미



temporal locality (시간적 지역성)

Spatial locality (공간적 지역성)

Sequential locality (순차적 지역성)

temporal locality (시간적 지역성)

최근에 access한 data가 가까운 미래에 다시 access 될 가능성이 높다.

for loop, subroutine, common variables.

Spatial locality (공간적 지역성)

memory에서 서로 인접되어 있는 data들이 연속적으로 access 될 가능성이 많다.

· table, array

* Sequential locality (순차적 지역성)

branch가 발생하지 않는 한

명령어들은 기억장치의 저장된 순서대로 인출되어 실행된다.

program이 수행되는 동안에

이런 지역성으로 인하여

cache에 이미 load된 data를

cpu가 access 하는 경우가 많다.

(cache hit ratio 증가)

→ 평균 memory access time을 단축시킬 수 있다.

cache design 시 목표.

- Cache hit ratio를 극대화.

CPU가 원하는 data가 cache에 존재할 확률을 높여야 한다.

- cache access time을 최소화.

cache hit인 경우 CPU가 cache로부터 읽어오는 시간을 줄인다.

- cache miss인 경우 delay time을 최소화

main memory로부터 cache로 data를 읽어오는 시간을 최소화.

- main memory와 cache의 data consistency (일관성)

유리 / overhead 최소화

CPU가 cache의 내용을 변경하였을 때

main memory에 그 내용을 update하는

결과로 인한 delay를 최소화

{ Write-Through
Write-Back

Fetch Policy Cache ← M.M

- demand fetch (요청 인출)
- prefetch (선인출)

Write Policy Cache → M.M

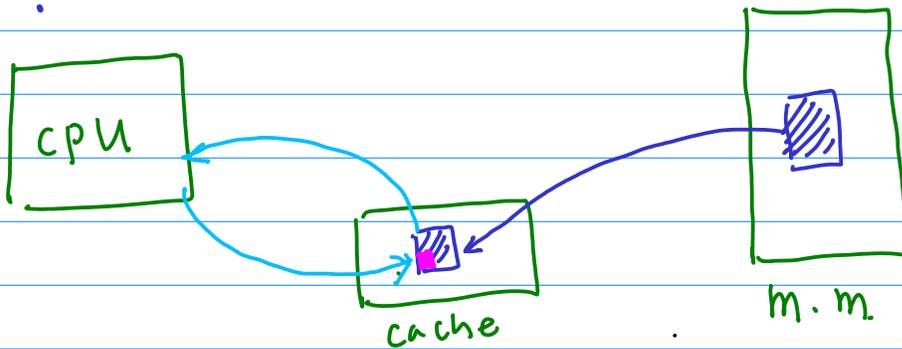
- Write-Through
- Write-Back.

Replace Policy Cache ↔

- Least Recently Used
- FIFO (First In First Out)
- Least Frequently Used
- Random

Write Policy

(쓰기 정책)



cache에 load된 data ← 변경 ← program 수정
↓ 변경

main memory에 있는 상응하는 data

일관성 consistency 유지 cache와 m.mem가 같은 data를 갖도록 유지.

cache write 시간 < main memory write 시간

Write-Through

- cache가 변경될 때 마다 main memory의 내용도 변경하는 방식

Write-Back

- cache가 변경되더라도 main memory를 update하지 않고 변경된 data가 속한 라인, replace 할 때 update한다

Modified bit Cache의 내용이 변경되면
main memory로 update할 필요가
있을 때 나타내는 bit

Write-through

Simple.

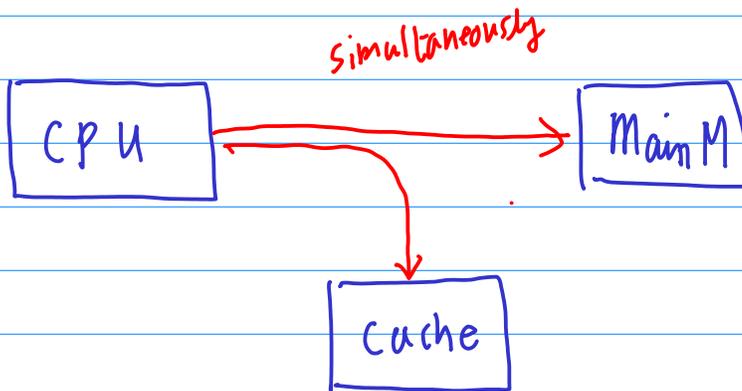
모든 write operation이

(cache 와
main memory 에)

동시에 수행된다 (simultaneously)

항상 main memory에 있는 data들은 valid하다
consistent한 data. (유효)

평균 access time이 W-Back에 비해서 많이 걸린다.



Write-back policy

어떤 block이 cache에서 replace 되어 cache에서 나갈 때

- 변경된 내용이 없는 경우 ($M=0$)

다른 block으로 replace 하면 된다

- 변경된 내용이 있는 경우 ($M=1$)

replace 되는 block은 main memory에 update 한다

다른 block으로 replace 하면 된다

Write-Back에서 Modified Bit set된

block을 교체할 때 (cache → main memory)

block 전체를 main-memory write → update

장점: main memory write operation 횟수가 감소

단점: cache에서 수정된 내용이 main memory에서 update 될 때까지 main memory의 해당 블록이 무효화 상태 (valid x) data inconsistency
→ control 회로가 복잡하다.

cache miss가 발생하여 원하는 data가 있는 block을 cache에 load 해야 하는 replace 되는 block이 변경된 경우는 write-back을 하기 때문에 추가의 시간이 소요된다.

$$\left(\begin{array}{l} \text{replace 되는 block을 m.m 쓰는 시간} \\ + \text{새로운 block을 m.m에서 cache로 load하는 시간} \end{array} \right) = \text{Cache miss 일 때 걸리는 시간}$$

Cache hit 인 access time < cache miss access time <<

Replacement Policy (교체정책)

Cache miss 발생하여 (CPU가 원하는 data가 cache에 없을 때)
새로운 block을 m. mem 에 cache로 load 할 때
· cache에 빈 자리가 없으면 cache에 있는 block을
교체 (replace)해야 한다

LRU (Least Recently Used)

사용되지 않은 블록으로 가장 오래 동안 load 되어 있던 block을 교체
(가장 오래 동안 access 되지 않은 block)

FIFO (First In First Out)

cache에 load 된지 가장 오래된 block을 교체
(가장 먼저 load 된 block)

LFU (Least Frequently Used)

cache에 load된 이후 access된 횟수가 가장 적은 block을 교체

Random

임의의 block을 교체.

Fetch Policy (인출 방식)

main memory → cache

demand fetch (요구 인출)

필요한 정보만 인출.

prefetch (선 인출)

필요한 정보뿐만 아니라 앞으로 필요할 것으로 예측되는 정보까지
fetch

CPU가 원하는 data를 fetch할 때
그 data와 근방의 data까지 함께
fetch해서 cache에 load한다.

block

