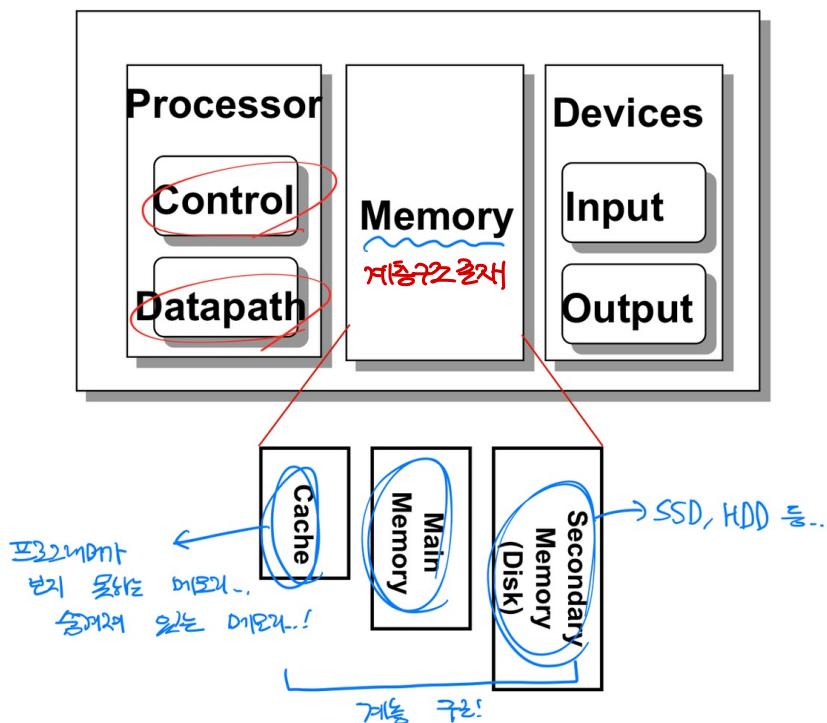




Ch.5-1 Large and Fast: Exploiting Memory Hierarchy Part 1 Cache Memory

* Computer의 주요 Component



- processor \Rightarrow 빠르고 정정 빨라짐
But, PRAM \Rightarrow " 느려진다면? \rightarrow 빠르고 정정
 \Rightarrow cache memory 개념 탄생
- Memory Hierarchy
 - i) size $\uparrow \rightarrow$ speed \uparrow
size $\downarrow \rightarrow$ speed \downarrow \rightarrow memory마다 특징 다른.
 \Rightarrow 적재적출에 알맞게 이용하여 memory 활용도 높여 주고 빠르게 처리?
 \Rightarrow 계획적, 병렬적으로 구현 가능

ii) Hierarchy가 유용한 이유 \Rightarrow Locality

• Locality : access하는 부분을 다른 access하는 경향 존재

[Temporal Locality] 최근에 참조된 Item이 또 참조될 가능성 큼

[Spatial Locality] " 같은 공간에 있는 Item "

iii) Hierarchy 관련 terms

• Block : 정부의 흐름 단위 (multiple words) \rightarrow copy할 때 Block 단위로

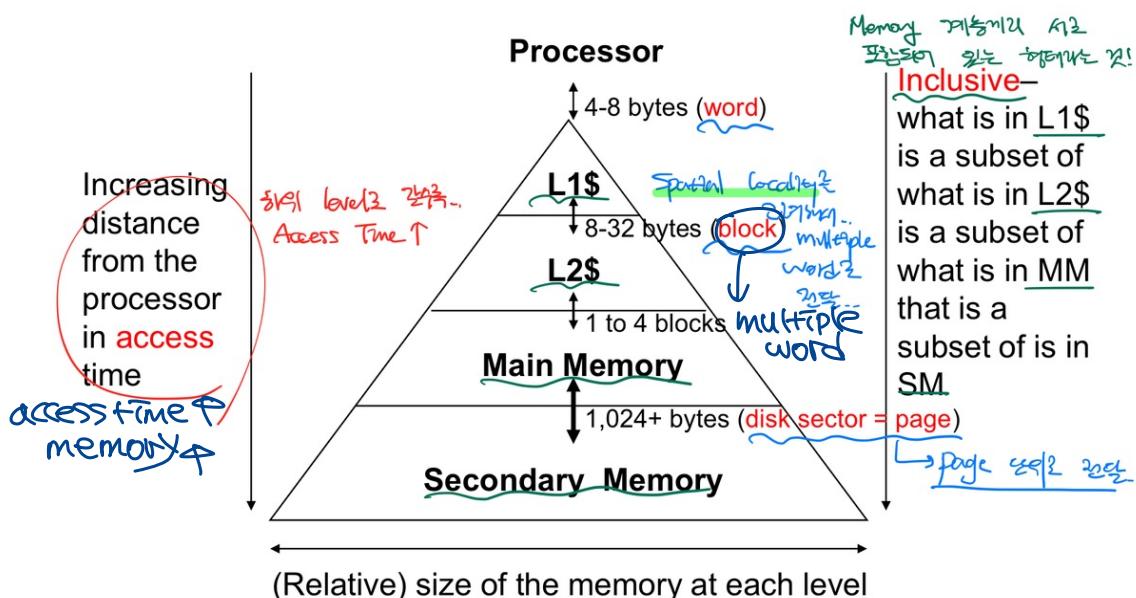
• Hit Rate : memory access를 통해 원하는 Data 찾는 비율

↳ Hit Time : access block + determine hit/miss

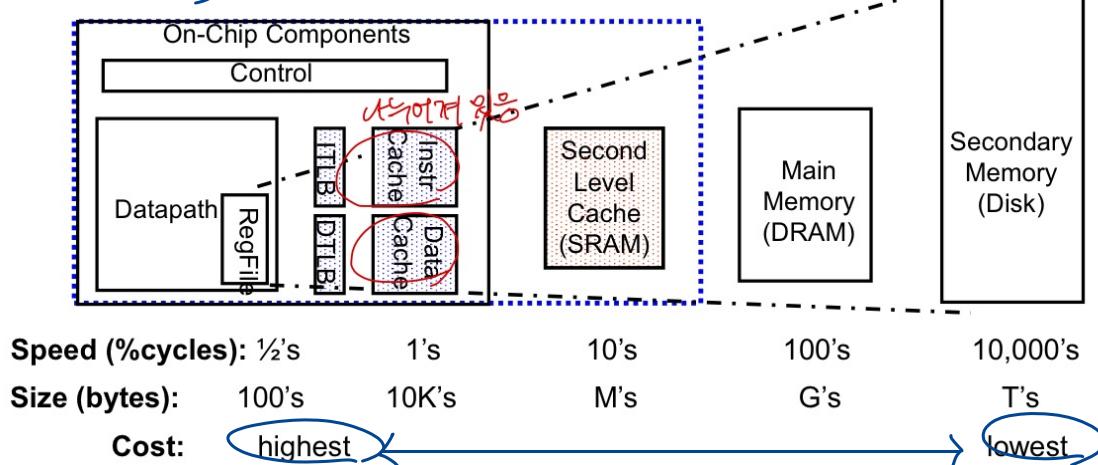
• Miss Rate : " 찾지 못하는 " (1 - HitRate)

↳ Miss penalty : 특정 단위의 대상으로 Block \rightarrow 더 low level로 replace

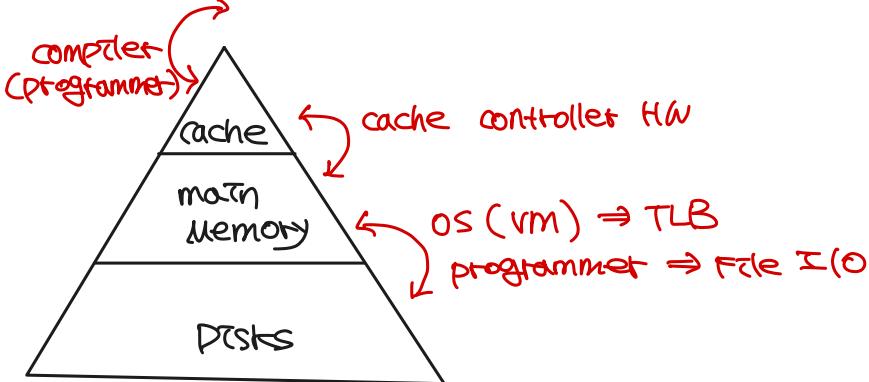
\therefore Hit Time < Miss Penalty (\because 찾기 쉬운 단위가 더 높아)



\Rightarrow locality 흐름의 적용 내용



IV) Memory Hierarchy Tech



① Cache : SRAM

- 속도↑ (access time : 0.9 ~ 2.5 ns)
- 집적도↑ ⇒ power 끊임이 필요, 가격↑
- Static : 비휘발성 !!

② Main memory : DRAM

- 속도↓
- 집적도↑ ⇒ power가 비효율적일 수 ↑, 가격↓
- Dynamic : 주기적으로 refreshed (every 8~16 ms)
 - ↳ DRAM은 1~2% 사용
- Address : row, column 구조화 활용
 - RAS : row detector triggering
 - CAS : column "

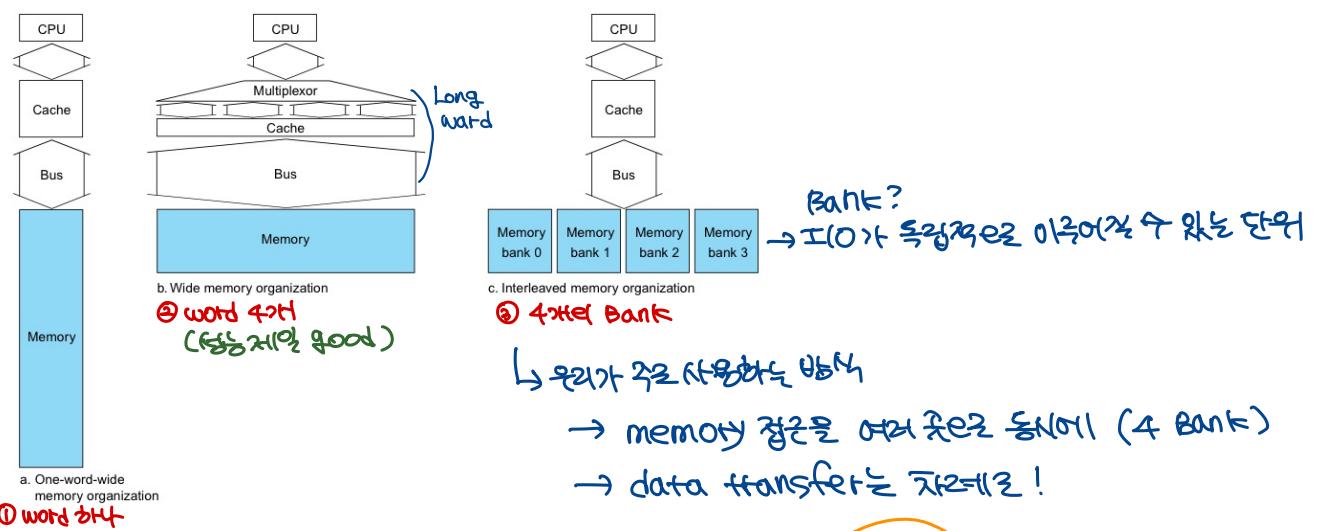
Advanced DRAM organization

- DRAM에 2진수의 모든 bit ⇒ rectangular array처럼 접근
 - ⇒ Burst Mode : 연속적인 정보를 Main Memory에서 한 번에读后
 - latency ↓

Double Data Rate (DDR) : Rising edge에서 전송

Quad " (QDR) : DDR의 INPUT, OUTPUT 병기

③ Increasing Memory Bandwidth → memory bank 사용 → multiple word read 가능!



$$③ = ① \text{el H/w} + ② \text{의 성능}$$

XOR'd after

→ example Addr + DRAM + Data *

memory system

cache block : 4 words per block

1 memory bus clock cycle : address 전송

15 " : DRAM 주기 상태 접근

1 " : data(word 단위) 전송

• One-word-wide

(1단위 address 전송 + 4개의 DRAM 접근) $\Rightarrow 1 + 4 \times 15 + 4 \times 1 = 69$ bus cycles

• Wide-memory bank

1단위 " + 2개의 " $\Rightarrow 1 + 2 \times 15 + 2 \times 1 = 33$ "

• Interleave

1단위 " + 1단위 " + 4개의 cycle $\Rightarrow 1 + 1 \times 15 + 4 \times 1 = 20$ "

↳ 효율적

* Caching

Q1 : cache 안에 data item이 있는지 어떻게 찾을까?

Q2 : 만약 없다고 판斷 어떤 방식?

① Direct mapped

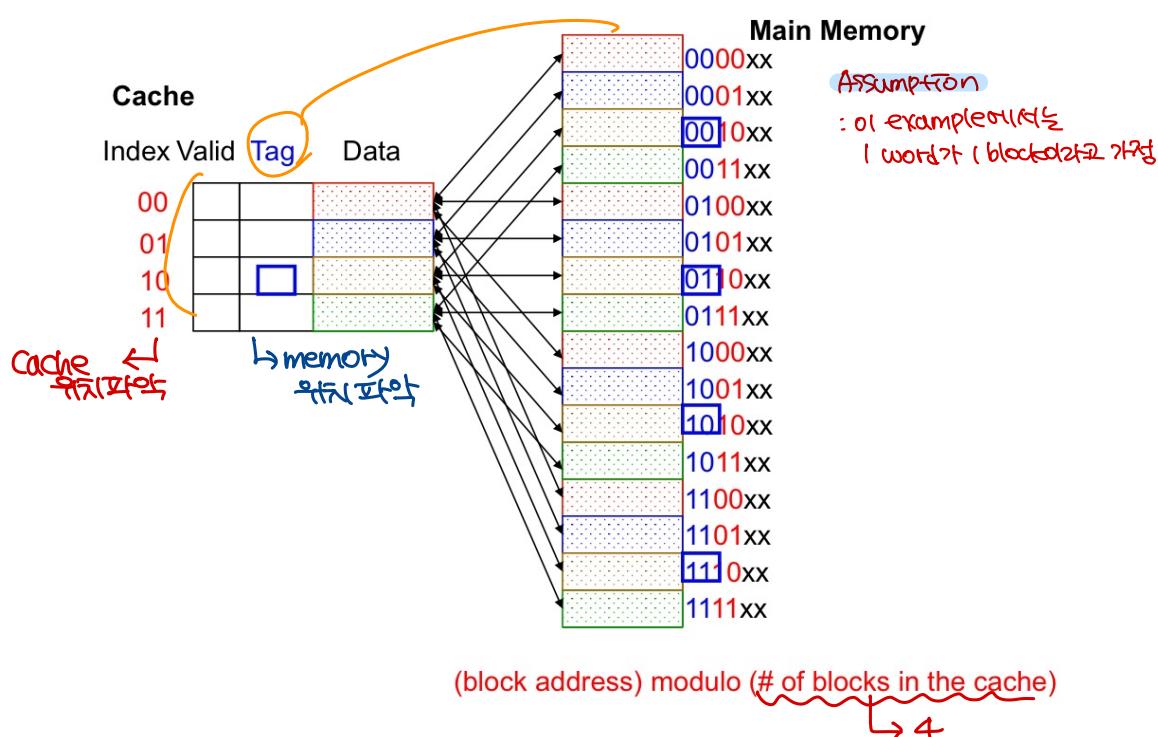
memory의 각 word에 Cache 블록 위치를 memory 주소에 기반을 두어 촬영

↳ 각 word가 한 장소 내에서만 있을 수 있음

→ ① Address mapping : (Block address) modulo (Cache Block 개수)

→ ② tag : Memory Block 구분함 \Rightarrow 상위 2비트

\neq index : Cache Block 구분함 \Rightarrow tag 다음 2비트 (modulo 연산값) *



② HIT / MISSES \Rightarrow Cache 3차원 여부

Read HTs : 읽었을 때 cacheonl 해당 Data 3차

Read MISSES : CPU 명령 \rightarrow memory로 블록 Data 가져온다 \rightarrow 3차원적

WRITE HTs

write-through : Cache, main memory 동시에 write \Rightarrow 불일치 해결 방지 가능

write-back : 무선 Cache에 있는 것만 write \rightarrow 이후에 memory update

WRITE MISSES : cache 3 Block 전체 Read \rightarrow word write

\hookrightarrow Read misses \rightarrow data \rightarrow write hits

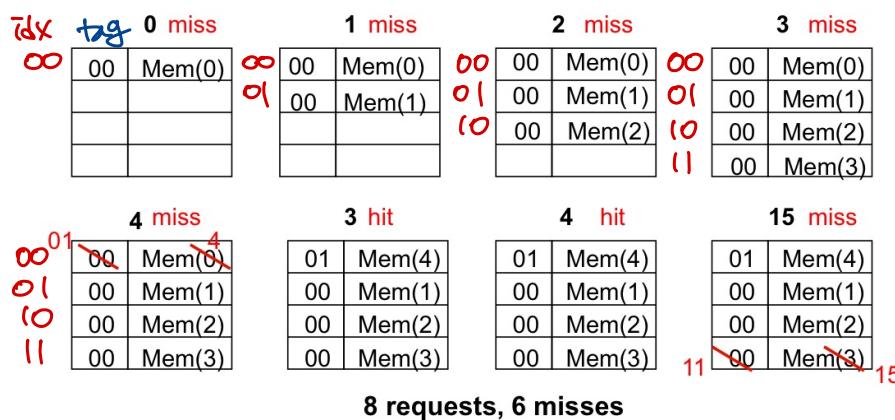
Cache pt

속도 경비 느낌

이관성 알고리즘 필요
다

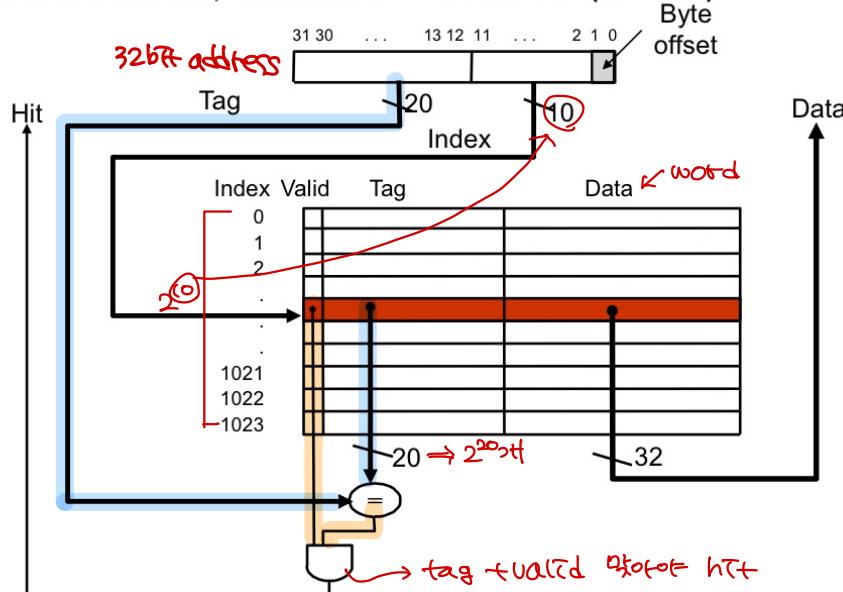
Start with an empty cache - all blocks initially marked as not valid

0 1 2 3 4 3 4 15



③ Direct Mapped Example

→ One word blocks, cache size = 1K words (or 4KB)



\Rightarrow spatial locality 장점

(한 번 access한 data를 끌어온다) \rightarrow (word block 끝까지 저장)

ex) address 1200 \rightarrow In 64thel 16byte Block

Block address = 1200 / 16 = 75 \rightarrow 75번 cache

Block number = 75 modulo 64 = 11
↓
Index

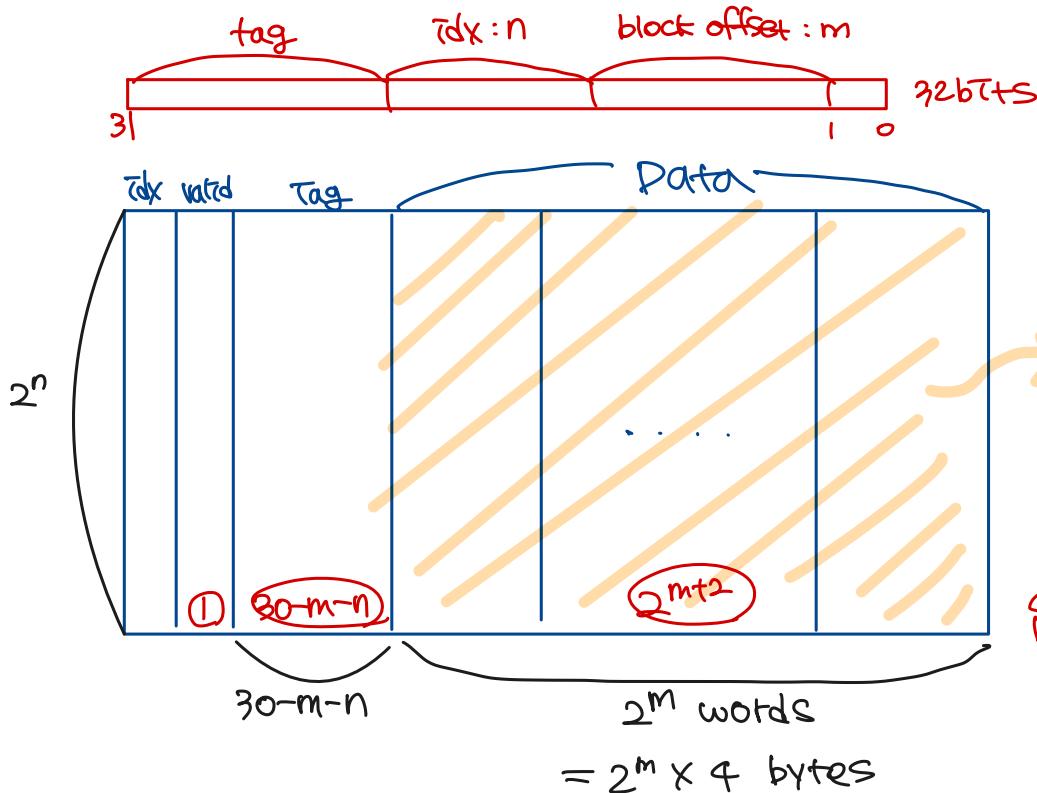
Cache size : 64 * 6 = 1KB
Block당 4개의 word ((6B))

④ Block size Consideration

적절한 크기로 정하는 것이 중요

- Block size ↑ → miss rate ↓ But, pollution ↑
- miss penalty ↑

⑤ Cache field size 구조체 설계



⑥ Cache hits handling

i) Read hits : 읽는다는데 결과 (I\$, D\$ 모두!)

ii) Write hits : (D\$에만)

- cache, memory ⇒ 일치성 → But, 성능 Bad
- " ⇒ 일치성 X 사용 ⇒ write-back

⑦ Cache Miss handling

i) Read misses (I\$, D\$ 모두!) → main memory로부터 data 가져오기

- ⇒ Block size ↑ miss penalty ↑
- [Early Restart
Requested word first]

ii) Write misses (D\$ only!)

- [Write allocate] : 이미 할당한 블록에 쓰기 (memory Read)
- [No-write allocate] : Block 자체를 개선 (memory Read X)

→ memory에 block 가져온 후 word를 block에 쓴다

⑧ 성능 측정 with Direct mapped cache

* Simplified model

$$\text{execution time} (\text{실행시간}) = (\text{execution cycles} + \text{memory stall cycles}) \times \text{cycle time}$$

$$\text{memory stall cycles} = \# \text{ of inst} \times \text{miss rate} \times \text{miss penalty}$$

\Rightarrow miss rate 혹은 miss penalty 조여보자

* example

- Cache miss rate \rightarrow instruction : 2%
- " \rightarrow Data : 4%
- CPI (stall 횟수) \rightarrow 2 (perfect cache ver)
- Miss penalty \rightarrow 100 cycles
- Total memory reference \rightarrow 40% (LW/GW)

$$\Rightarrow \text{miss cycles} = I \times 2\% \times 100 = 2I$$

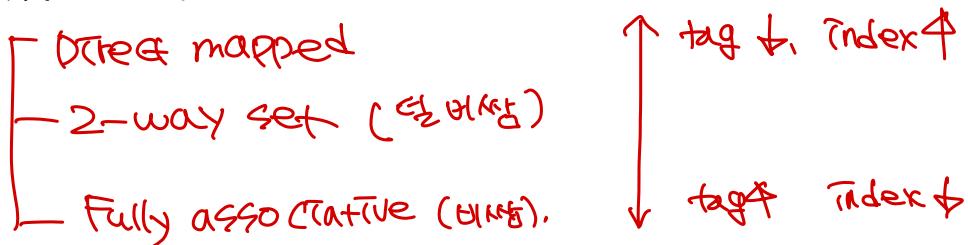
$$\text{Data } .. = I \times 40\% \times 4\% \times 100 = 1.6I$$

$$\text{CPI-stalls} = 2 + 1.6 + 2 = 5.6, \star \quad \text{perfect} \quad \rightsquigarrow 2.8x \text{ speed up}$$

* average access time

$$\text{AMAT} = \text{Hit Time} \times \text{miss rate} \times \text{miss penalty}$$

⑨ Associative Caches



⑩ Replacement Policy

- Direct mapped \rightarrow parallel하게 실행
- set associative Cache \rightarrow 정렬된 구조 내에서 block 훈련
- LRU \rightarrow HW 필요
- Random

\Rightarrow associativity 높일 때마다 miss rate 향상하기 쉬워

* Multi-Level Caches

- ① L1 Cache : 주로 빠름
- ② L2 Cache : 크지 않은 느림

example) CPU base CPI = 1

$$\text{clock rate} = 4\text{GHz} \Rightarrow \text{cycle time} = 0.25\text{ns}$$

instruction miss rate = 2%

memory access time = 100ns

- 1) L1 cache만 있을 때

$$\text{miss penalty} = 100\text{ns} / 0.25\text{ns} = 400 \text{ cycles}$$

$$\text{effective CPI} = 1 + 0.02 \times 400 = 9$$

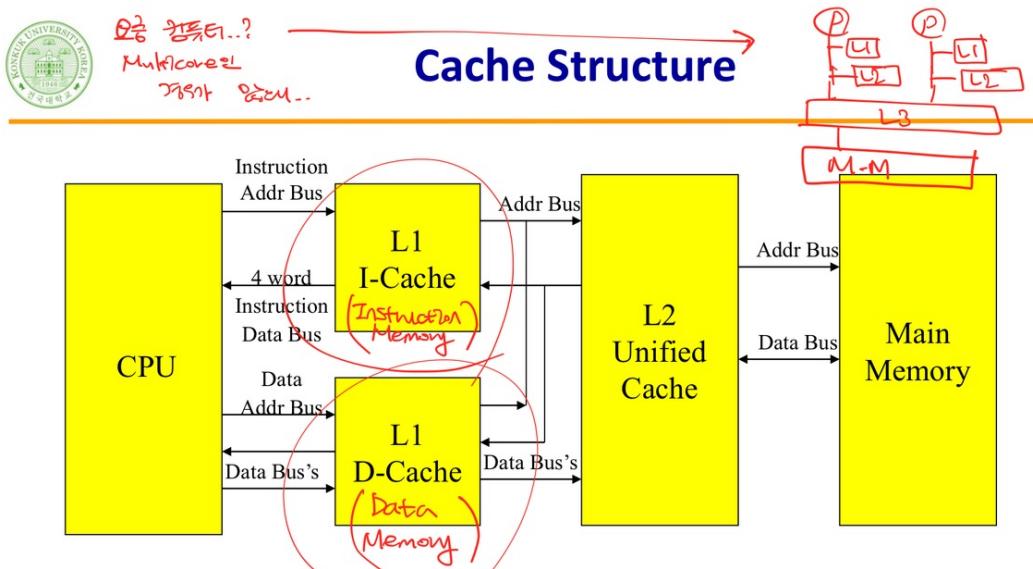
- 2) L2 cache 추가

$$\text{access time} = 5\text{ns}$$

$$L1 \rightarrow L2 \text{ miss rate} = 0.5\%$$

$$\Rightarrow L2 \text{ miss penalty} = 5\text{ns} / 0.25\text{ns} = 20 \text{ cycles}$$

$$\Rightarrow \text{effective CPI} = 1 + 0.02 \times 20 + 0.005 \times 4000 = 3.4$$



- ❑ The L1 cache is usually separated into separate Data and Instruction caches (L1 cache는 주로 Data/Instruction Cache로 나누어놓는다!)
- ❑ The L2 and higher caches are normally unified caches, i.e. both instructions and data are in the same cache