

→ power 절약 .
→ 성능이 어떻게 되느냐
→ "SRAM!!" 이라고 끝낼 것이 아님.

Introduction to Cache Memory – Part#1

지금까지 배웠던 컴퓨터 속의
cache memory에 대해 알아볼 것임.

Yoonjin Kim

Full Professor

**Dept. of Computer Science
Sookmyung Women's University**

Outline

- Before beginning
- What is **cache memory**? (⇒ cache memory 개념)
- Why do **CPUs need cache memory**? (⇒ 때 CPU는 캐시메모리를 필요로 하니까)
- What **data is stored into cache memory**?
(⇒ 캐시메모리에 어떤 데이터가 저장되는가) 이론의 시작점.
↑ 이벤트수 분량
- How does cache memory work?
 - Basic concepts
 - Cache read-operation
 - Cache write-operation
 - Cache memory space allocation
↳ 메모리 공간 할당

Before Beginning

(구글에서 찾아본 결과)

- Let's search for commercial CPU specification.

(기본 상용 CPU의 spec)

Cache가 어떤지 알면

(해당 CPU의 clock rate)를 명확히 알 수 있음 - 성능의 valence

clock rate (클럭 주파수)

MANUFACTURER	Intel
MODEL	Core i7-6700
PART #	BX80662I76700
DATA WIDTH	64-bit
SOCKET	LGA1151
OPERATING FREQUENCY	3.4GHz
MAX TURBO FREQUENCY	4GHz
CORES	4

L1 CACHE
4 x 32KB Instruction
4 x 32KB Data
L2 CACHE
4 x 256KB
L3 CACHE
1 x 8MB

같이 읽어야 해!

왜 이렇게 구성되어 있는가?
L1, L2, L3는 뭐지?
size 차이는 왜 있는 거지?

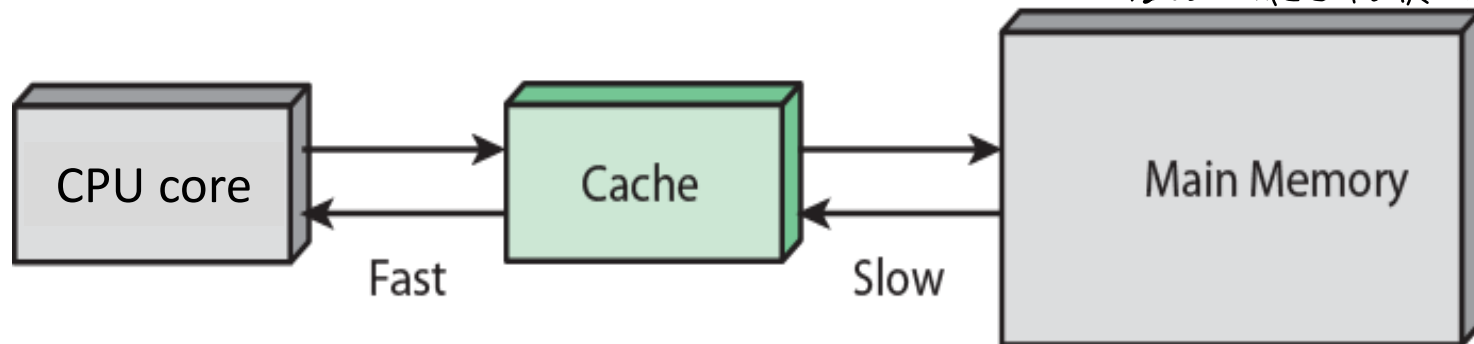
Cache?

Apple A10 Fusion

	
Produced	From September 7, 2016 to Present
Designed by	Apple Inc.
Common manufacturer(s)	TSMC ^[1]
Max. CPU clock rate	to 2.34 GHz ^[2]
Min. feature size	16 nm
Instruction set	A64, A32, T32
Microarchitecture	Hurricane and Zephyr both ARMv8-A-Compatible
Product code	APL1W24
Cores	Quad-core (2× Hurricane + 2× Zephyr)
L1 cache	Per core: 64 KB instruction + 64 KB data
L2 cache	3 MB shared
L3 cache	4 MB shared

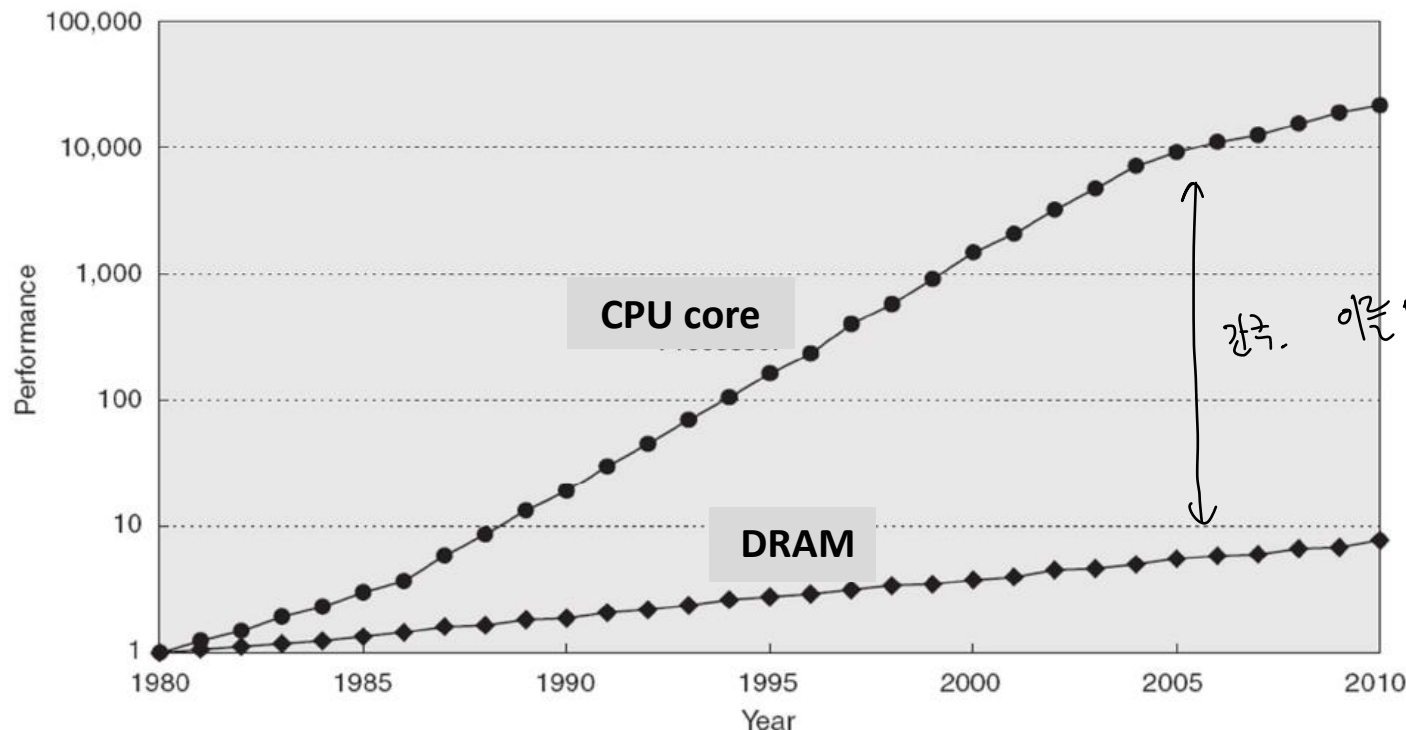
What is Cache Memory?

- Cache memory is an intermediate buffer between CPU Core and main memory (DRAM).
↳ 캐시메모리는 DRAM과 CPU 사이에 “중간 buffer”.
 - CPU = (CPU core + Cache memory) on a chip
- It is small amount of fast memory (SRAM). ↳ 작은 크기의 빠른 메모리.
- It may be located on a chip with CPU core or off the chip.
(↳ CPU 안에 위치할 수도 있고, 밖에 위치할 수도 있음.)
↳ 보통은 칩셋 안에 들어있음.



Why do CPUs need Cache Memory?

- **CPU core - DRAM performance gap** (CPU와 DRAM의 성능 gap)
 - Main memory (DRAM) access time is higher than CPU core clock period.
 - Therefore, high speed core's time is wasted during memory access.
 - So it results in significant delay.



이론에 따르면
이렇게 빠른 CPU에서
기타는
시간을 최소화하기 위해
SRAM을 쓴다!
이것이?
가속.

Why do CPUs need Cache Memory?

- To minimize the waiting time for the CPU core

→ DRAM 대신 SRAM을 메인메모리로 사용하면
어떨까?

- We can consider the fastest memory (SRAM) as main memory.
- However, unfortunately, SRAM is too expensive for most people to buy a lot of.
- There is a tradeoff between speed, cost and capacity as follows.

↳ 비쌌

↳ 중간에 뭘 넣어야 할까? (CPU와 DRAM)

Storage	Speed (Delay)	Cost (Price/MB)	Capacity (Size)
SRAM (Static RAM)	Fastest (1-10 cycles)	Expensive (~\$5)	Smallest (24KB-12MB)
DRAM (Dynamic RAM)	Slow (100-200 cycles)	Cheap (~\$0.10)	Large (512MB-64GB)

이게 빠지

사서

대부분
off-chip ←

(pg. 92가서
revisit)

↑
CPU 클럭을 cycle
(앞에서 봤을때 클럭은 클럭)

Why do CPUs need Cache Memory?

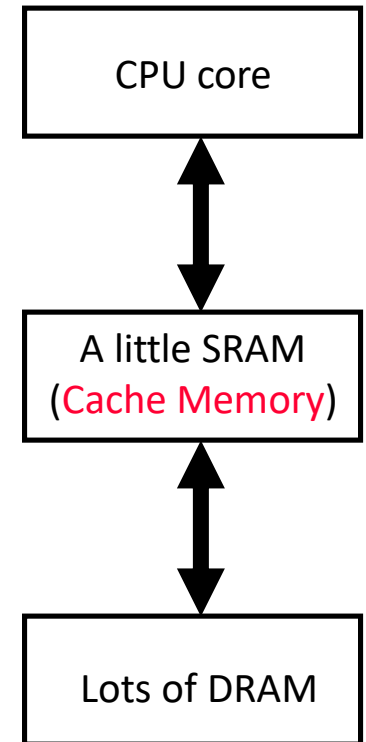
- To minimize the waiting time for the CPU core

→ SRAM과 DRAM의 valence를 발견하면 그게 더 좋지 않을까?

- Wouldn't it be nice if we could find a balance between fast & expensive memory (SRAM) and slow & cheap memory (DRAM)?
- We do this by introducing a **cache memory**, which is a small amount of fast, expensive memory (SRAM).
- The cache memory goes between the CPU core and main memory (DRAM). → CPU와 DRAM 사이에 있는 캐시메모리
- **The cache memory keeps a copy of the most frequently used data from the main memory.**

메인메모리에서 가장 많이 사용되는 데이터를 캐시메모리에 keep 하는 것

⇓
캐시메모리의 역할



CPU와 DRAM의 gap을
cache memory를 통해 메꾸는 것.

Why do CPUs need Cache Memory?

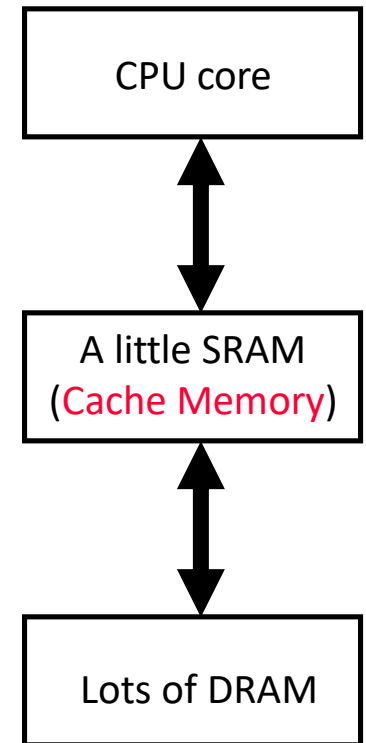
- The waiting time for the CPU core is minimized because (⇒ CPU의 waiting time이 줄어든다)

- Reads and writes to the most frequently used addresses will be serviced by the cache memory.

⇒ 캐시메모리가 자주 사용하는 데이터를 읽고 쓸 수 있게 하죠.

- We only need to access the slower main memory for less frequently used data.

⇒ 덜 사용하는 데이터는 메인메모리라는 이용해서 (느리게) 접근하죠.



어떤 성능 개선을 원하느냐.

어떤 컴퓨터가 발전할 수 있는 이유 (성능 개선)

- 1) CPU 클럭 증가
- 2) core 개수 증가
- 3) cache 메모리

Why do CPUs need Cache Memory?

- **Recall** - Basic Architecture of Programmable Processor

(Memory Hierarchy Aspect)

⇒ 같은 clk이 안들어갈수 있음.
(SRAM, DRAM, 레지스터에)

- **Register File** – 38 transistors for 1-bit

- Fastest
- But biggest size
- Very expensive

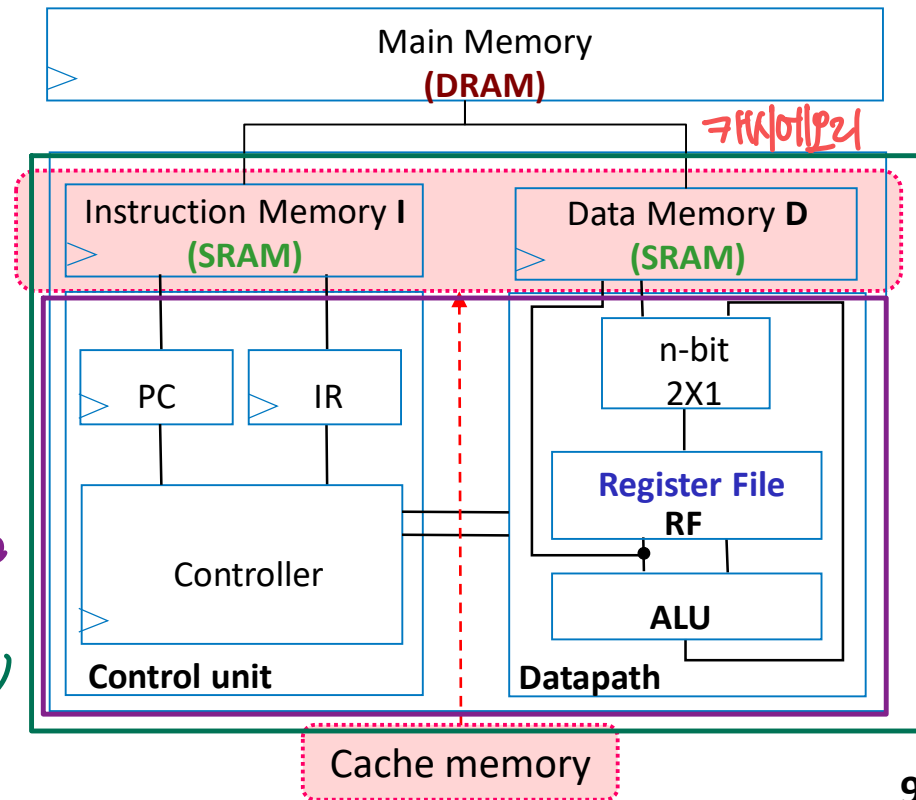
- **SRAM** – 6 transistors for 1-bit

- Fast
- More compact than register file
- Expensive

- **DRAM** – 1 transistor for 1-bit

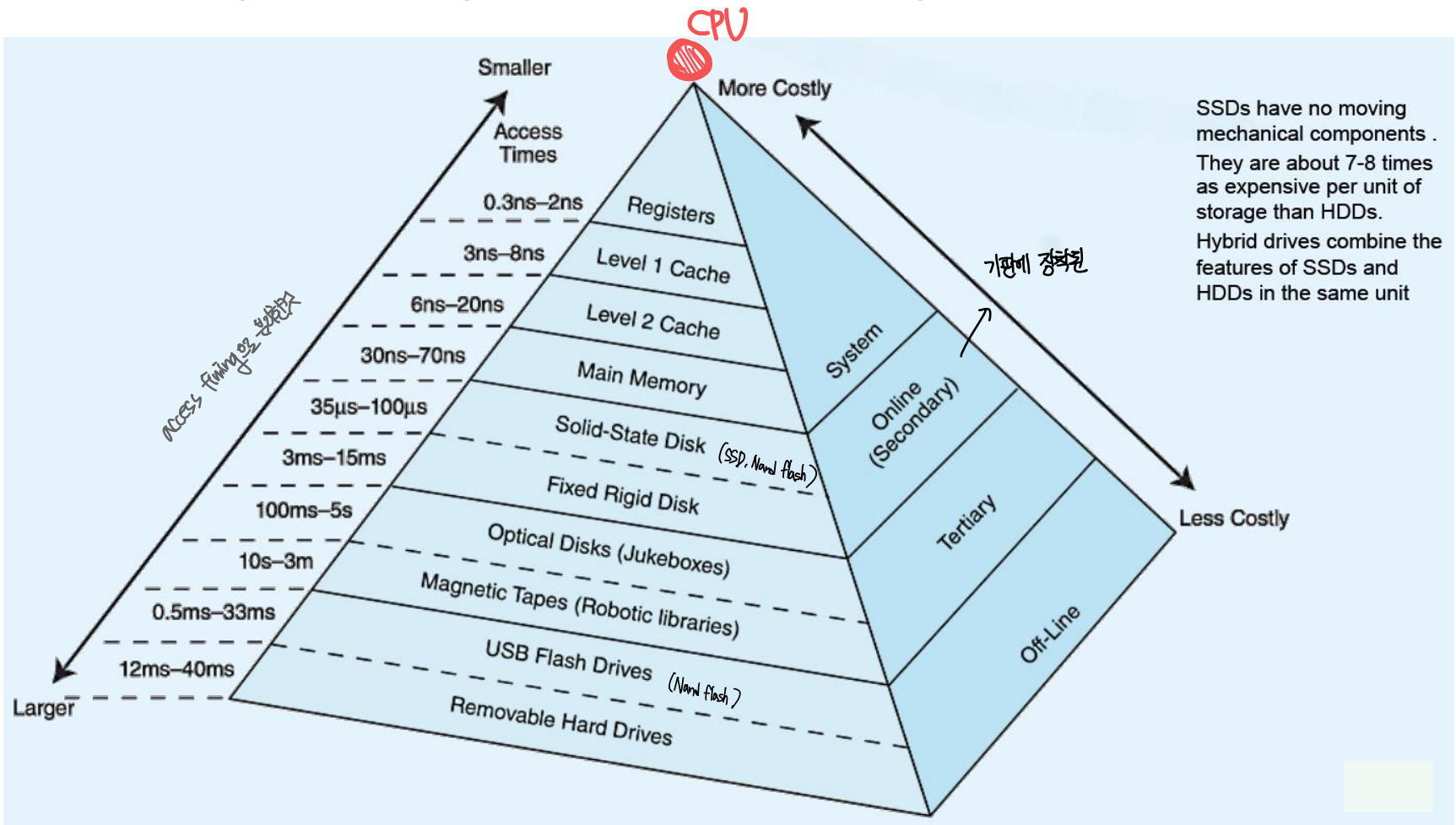
- Slowest
 - ❖ And refreshing takes time
- But very compact
- Cheap

CPU
core
CPU



Why do CPUs need Cache Memory?

- Memory hierarchy with cache memory



→ 캐시메모리가 어떤 데이터를 저장하는가?

What Data is stored into Cache Memory?

- **How to know the most frequently used data?** *** → 어떻게 가장 많이 사용되는 데이터인지를 캐시에 데이터를 넣는 것인가?
 - The cache memory keeps a copy of **the most frequently used data** from the main memory.
 - However, It's usually difficult or impossible to figure out what data will be "**most frequently used**" before a program actually runs. (프로그램의 동작전에)
- **The principle of locality** *** (정규성의 특징)
 - In practice, most programs exhibit locality, (which the cache memory can take advantage of.) → 모든 프로그램은 locality를 가지고 있다. 이를 캐시 메모리에서 이용할 수 있다.
 - The principle of **temporal locality** ① says that if a program accesses one memory address, there is a good chance that **it will access the same address again**. → 프로그램이 한 메모리 주소에 접근하면, 또 같은 주소에 접근할 가능성이 높음 (같은 cycles 근처)
 - The principle of **spatial locality** ② says that if a program accesses one memory address, there is a good chance that **it will also access other nearby addresses**. → 프로그램이 한 메모리 주소에 접근하면, 근처의 다른 주소에 접근할 가능성이 높음. (프로그램은 순차적으로 동작하기)

아직에
이해하지
못해서 캐시메모리에
저장될지 걱정됨.

What Data is stored into Cache Memory?

① Temporal locality in programs (⇒ 여기서 program은 "명령어"를 의미함)

- The principle of **temporal locality** says that if a program accesses one memory address, there is a good chance that **it will access the same address again**.
- Loops are excellent examples of temporal locality in programs.
→ Loop은 temporal locality를 잘 보여주는 예시임.

- The loop body will be executed many times. ⇒ 계속 반복됨.
- The computer will need to access those same few locations of the instruction memory repeatedly.
- For example:

assembly 코드의 loop
예시

☆ **Loop:** lw \$t0, 0(\$s1)
add \$t0, \$t0, \$s2
sw \$t0, 0(\$s1)
addi \$s1, \$s1, -4
bne \$s1, \$0, **Loop**

4개의 명령어가 계속 반복됨.

한번 실행된 명령어 시간이 주기를 주기에
또 반복 ⇒ temporal locality

조건을 만족할 때까지
loop를 하는 것.

- Each instruction will be fetched over and over again, once on every loop iteration.

What Data is stored into Cache Memory?

① Temporal locality in data (⇒ 여기서 data는 명령어를 제외한 data를 의미함)

- Programs often access the same variables over and over, especially within loops. Below, `sum` and `i` are repeatedly read and written.

①

```
sum = 0;
Loop for (i = 0; i < MAX; i++)
    sum = sum + a[i];
```

for loop이 끝날 때까지

sum, i가 반복적으로
(계속 사용되는 data)

⇒ data가 자주 사용되므로 temporal locality.

- Commonly-accessed variables can sometimes be kept in registers, but this is not always possible because of a limited number of registers.

레지스터는
크기가 작아서
loop 변수를 다 넣을 수 없음.

아주 작은 레지스터
SRAM에
넣는 것.

What Data is stored into Cache Memory?

② Spatial locality in programs (⇒ 여기서 program은 "명령어"를 의미함)

- The principle of spatial locality says that if a program accesses one memory address, there is a good chance that **it will also access other nearby addresses**.

공간에 조건에 의해
뛰어넘더라도
조건에 해당하지 않는 것들은
연속된 명령어를 가져옴.
(연속성)

```
sub $sp, $sp, 16
sw $ra, 0($sp)
sw $s0, 4($sp)
sw $a0, 8($sp)
sw $a1, 12($sp)
```

→ 0번지
→ 1번지
→ 2번지
→ 3번지 ↓
순서대로 가져옴

⇒ 명령어는 순서대로 실행이 되기 때문.

- **Nearly every program exhibits spatial locality**, because instructions are usually executed in sequence—if we execute an instruction at memory location i then we will probably also execute the next instruction, at memory location $i+1$.
 ⇒ loop과 같은 경우, temporal과 spatial locality를 둘 다 가지고 있음.
- **Code fragments such as loops exhibit both temporal and spatial locality.**

What Data is stored into Cache Memory?

②

- **Spatial locality in data** (→ 여기서 data는 명령어를 제외한 data를 의미함)

- Programs often access data that is stored contiguously.

“배열”
특성

- **Arrays**, like **a** in the right code, are stored in memory contiguously.

- The **individual fields** of a record or object like **employee** are also kept contiguously in memory.

- Can data have both spatial and temporal locality?

- Yes, array like **b** in the right code

배열 → spatial locality

+

for loop → temporal locality

```
sum = 0;  
for (i = 0; i < MAX; i++)  
    sum = sum + a[i];
```

배열은 순서대로 메모리 저장

배열

```
employee.name = "Petter";  
employee.boss = "Linda";  
employee.age = 45;
```

field → 각각의 인스턴스 한 번에 하나씩

```
sum = 0;  
for (i = 0; i < MAX; i++)  
    for (j = 0; j < MAX; j++)  
        sum = a[i] + b[j];
```

What Data is stored into Cache Memory?

- How caches take advantage of temporal locality

⇒ 캐시 메모리가 어떻게 "temporal locality"를 이용하는가!?

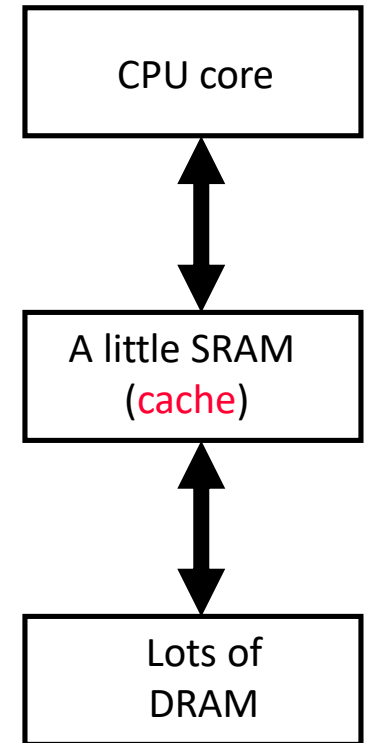
- The first time the CPU core reads from an address in main memory, a copy of that data is also stored in the cache. → 캐시에 메인메모리 데이터를 읽어들, 그리고 복제 저장.

- The next time that same address is read, we can use the copy of the data in the cache *instead* of accessing the slower dynamic memory.

↳ 짧은 주기로 읽으면, DRAM에서 묻지 않고 캐시에서 읽어들.

- So the first read is a little slower than before since it goes through both main memory and the cache, but subsequent reads are much faster.

↳ 첫번째는 느리지만 이후 것은 메인메모리에서 매우 빨라짐.



What Data is stored into Cache Memory?

- How caches take advantage of spatial locality

⇒ 캐시메모리가 어떻게

"spatial locality"를 이용하는가?

- When the CPU core reads location i from main memory, a copy of that data is placed in the cache.

- But instead of just copying the contents of location i , we can copy *several* values into the cache at once, such as the four bytes from locations i through $i+3$.

한번에 2만
가져오는 대신

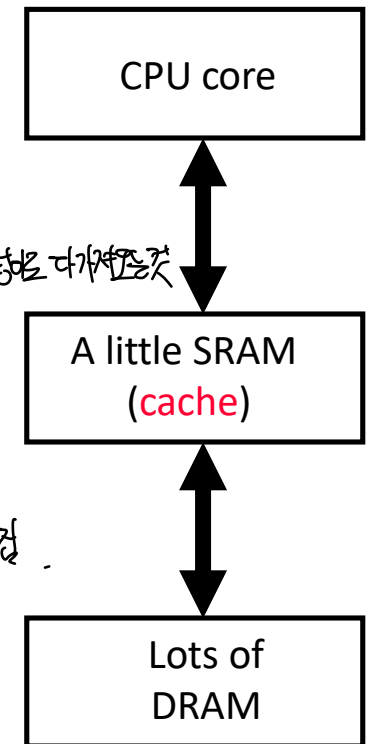
→ 용량을 다 가져오는 것

- If the CPU core later does need to read from locations $i+1$, $i+2$ or $i+3$, it can access that data from the cache and not the slower main memory.

→ 이 데이터는 캐시에 가져다와서 다시 메모리 접근의 필요가 없음.

- Again, the initial load incurs a performance penalty, but we're gambling on spatial locality and the chance that the CPU core will need the contiguous data.

→ CPU가 연산한 데이터는 필요할 것 없고 생략하고 다음 진행하는 것.



How does Cache Memory Work?

⇒ 캐시 메모리가 어떻게 작동하는가?

- Basic Concepts

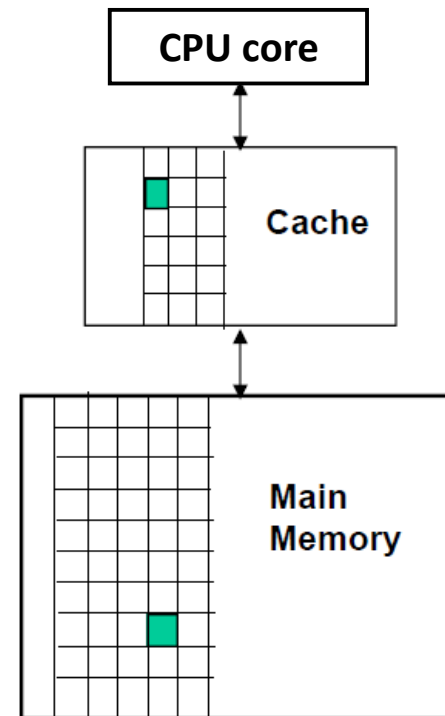
Cache Hit

- If the requested data is found in one of the cache blocks (upper level)
⇒ there is a **hit** in the cache access

CPU core 가 요청한 데이터가
캐시 메모리 안에 들어 있는 상황

⇒

캐시 메모리가
히트 했다!!!



How does Cache Memory Work?

- Basic Concepts

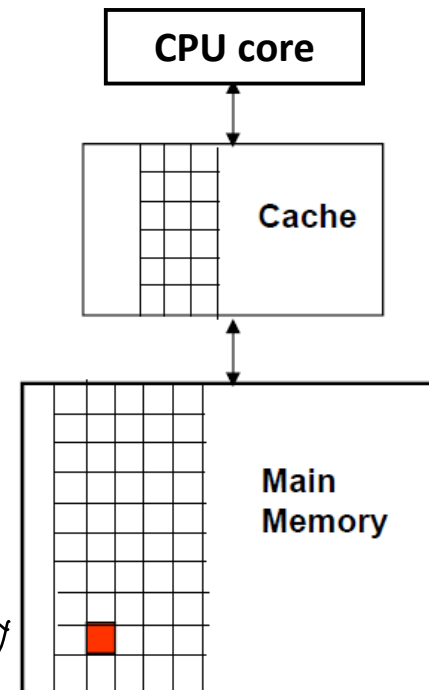
Cache Miss

- If the requested data is not found in in one of the cache blocks (upper level) \Rightarrow there is a **miss** in the cache access
 \Rightarrow to find the block, we need to access the lower level of the memory hierarchy

- In case of a data miss, we need:

- To stall** the CPU core; \Rightarrow CPU core 멈추고
- To require block from the main memory \Rightarrow 메인 메모리에 block을 요청함
- To copy the block in cache; \Rightarrow 캐시에 block copy
- To repeat the cache access (hit).

\Rightarrow 다시 캐시에 접근함.



CPU에서 요청한 데이터가
캐시에 없으면 미스 상태

\downarrow
캐시 메모리가
MISS !!

일단(처음)
캐시 메모리를
활용하는 방안
(Miss가 있을 때)

How does Cache Memory Work?

- Basic Concepts

⇒ 캐시에어라는 hit, miss를 가지고 성능을 평가할 수 있음.

- **There are two basic measurements of cache performance.**

- The **hit rate** is the percentage of memory accesses that are handled by the cache memory.
- The **miss rate** (1 – hit rate) is the percentage of accesses that must be handled by the slower main memory.

→ 보면 좋지 않은 것임.

$$(\text{hit-rate}) + (\text{miss-rate}) = 1$$

- **Typical caches have a hit rate of 95% or higher, so in fact most memory accesses will be handled by the cache memory and will be dramatically faster.**

⇒ 말만 듣고 hit-rate은 95%를 넘음.

→ 예를 들어 성능 개선이 가능함.

→ 기본선은 95%임.

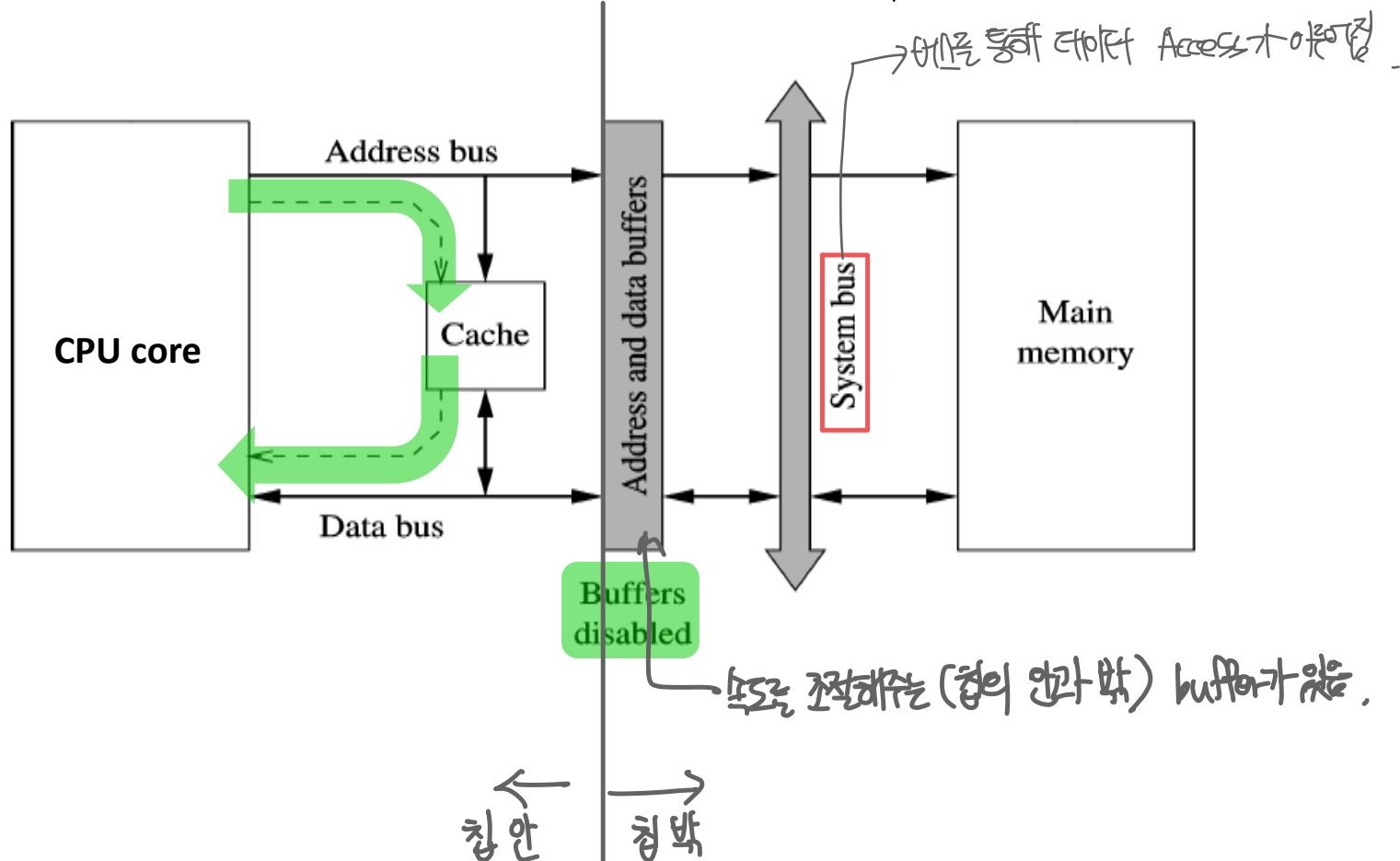
How does Cache Memory Work?

HIT

- Cache-Read Operation

- Read hit

– CPU core reads data in cache. ⇒ 빠른 속도로 캐시에서 데이터를 읽어들임.

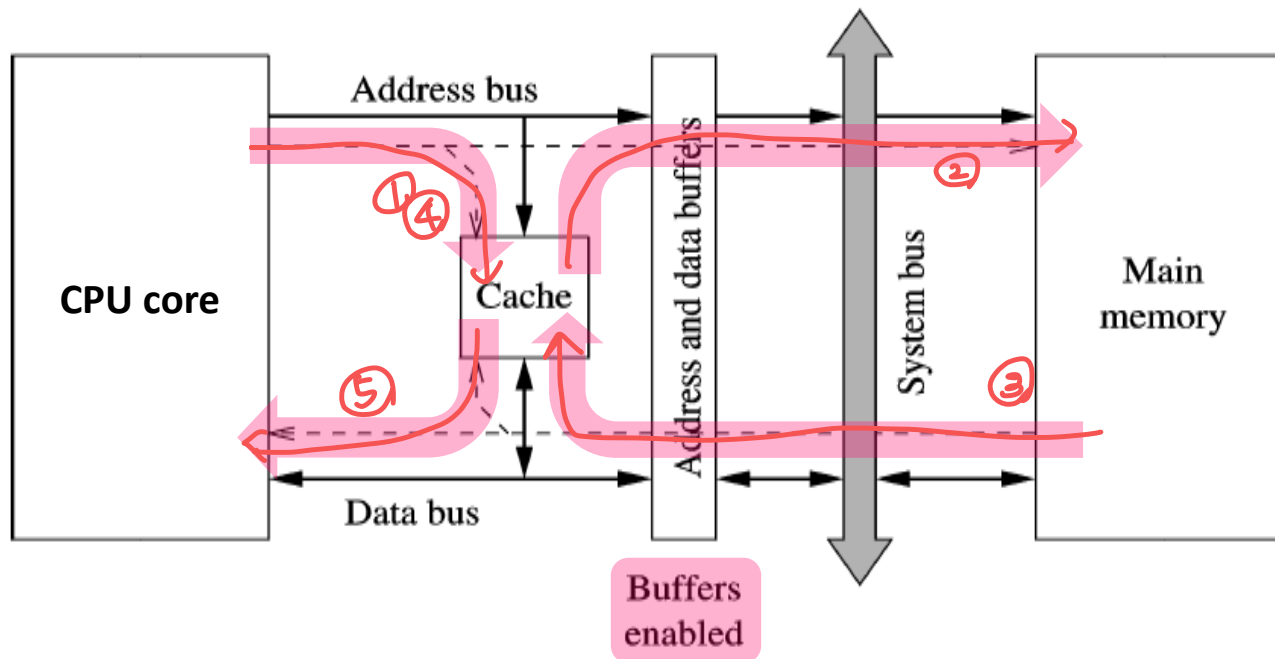


How does Cache Memory Work?

MISS

- Cache-Read Operation

- **Read miss** ⇒ 캐시에 필요한 데이터가 없는 상황
 - CPU core stalls and requests the data to main memory. core가 필요한 데이터에 요청.
 - The data is copied in cache memory. 캐시에 데이터 복사
 - CPU core repeats cache-read operation. CPU core가 캐시를 다시 찾음.

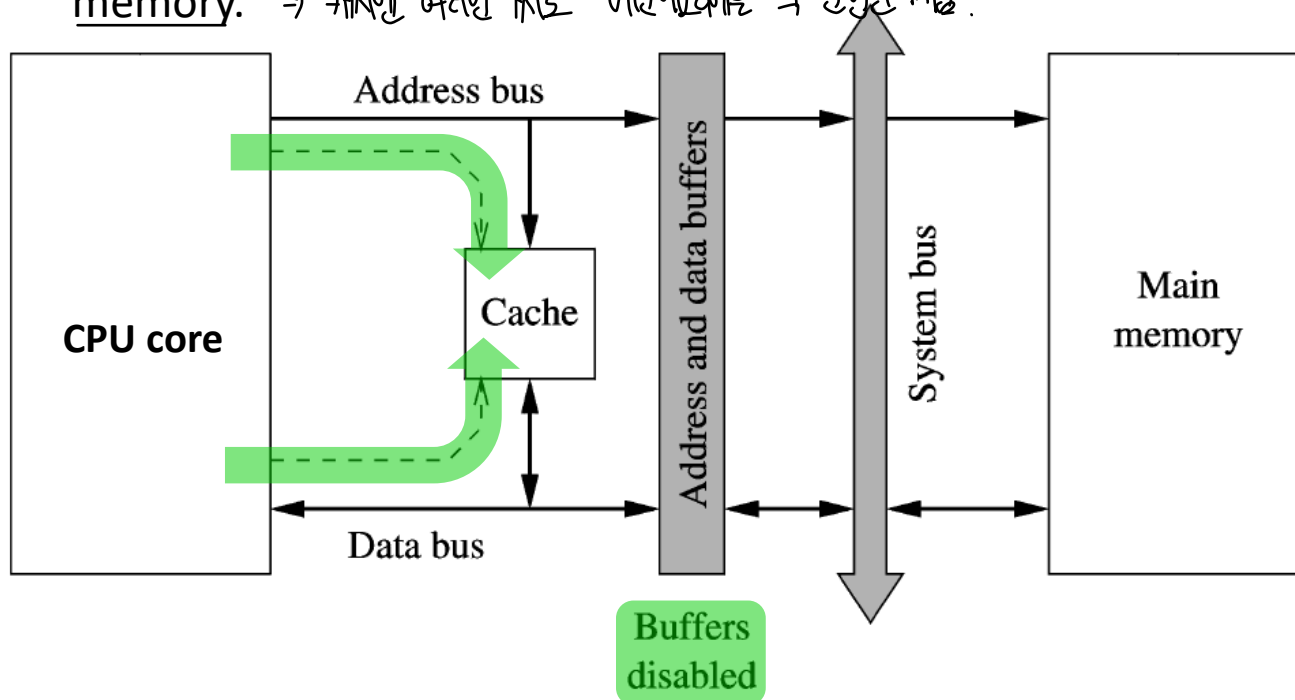


How does Cache Memory Work?

HIT

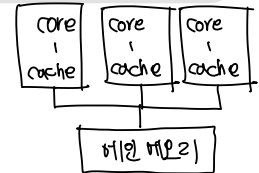
write는 hit 방식 2가지, miss 방식 2가지 있음. - **Cache-Write Operation**

- **Write hit**: ① write-back ⇒ 변경된 데이터를 바로 메인메모리에 쓰지 않음.
 - CPU core writes data in cache memory only. ⇒ 캐시에요 CPU core가 데이터를 저장.
 - The modified data is written to the main memory when it is replaced due to a miss ⇒ miss가 난 경우에만 메인메모리에 변경된 데이터를 쓰줌. (1번만) → 그렇지 않은 경우 계속 캐시와 core만 주고받음.
 - Multiple writes to the same data require only a single write to the main memory. ⇒ 캐시에 여러번 쓰고 메인메모리는 딱 한번만 쓰줌.



메인메모리에
접근하는 횟수를 줄임

* CPU가 multicore인 경우



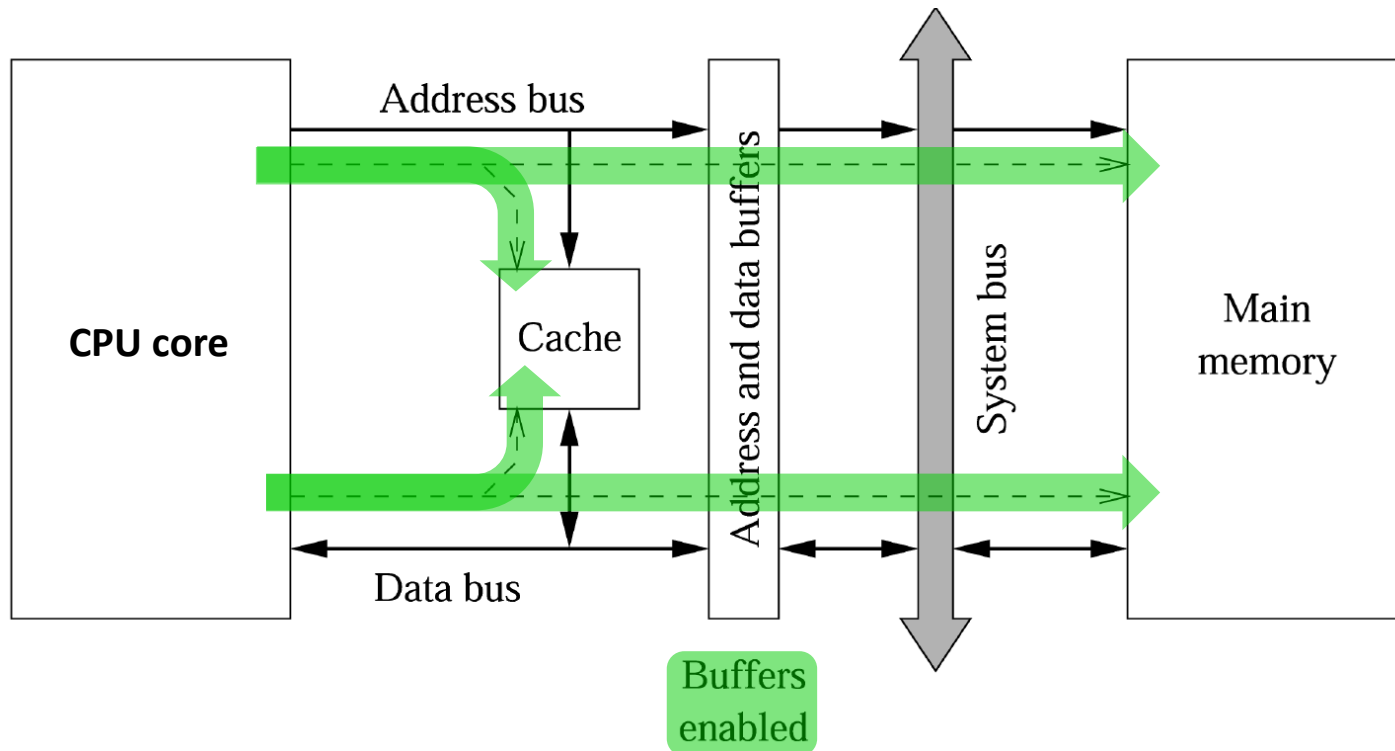
메인메모리를 공유하므로
write back 방식이 적용될 수 없음.

How does Cache Memory Work?

HIT

- Cache-Write Operation

- **Write hit**^②: write-through \rightarrow 캐시에 데이터도 쓰고 동시에 메인메모리에도 쓰는 것.
 - CPU core writes data both in cache memory and in main memory.
 - \rightarrow 둘다 업데이트가 되므로 multi-core에서도 사용가능.

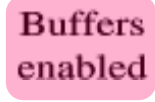


MISS

→ 7월 15일에 새로 해설 공이 없 경우

①

- CPU core stalls and requests the data to main memory. \Rightarrow **데이터를 캐시에서 가져오기 위해 대기**
- The data is copied in cache memory. \Rightarrow **캐시 메모리에 데이터를 복사함**
- CPU core repeats cache-write operation. \Rightarrow **그후에 CPU core가 다시 캐시 메모리에 데이터를 쓴다**



How does Cache Memory Work?

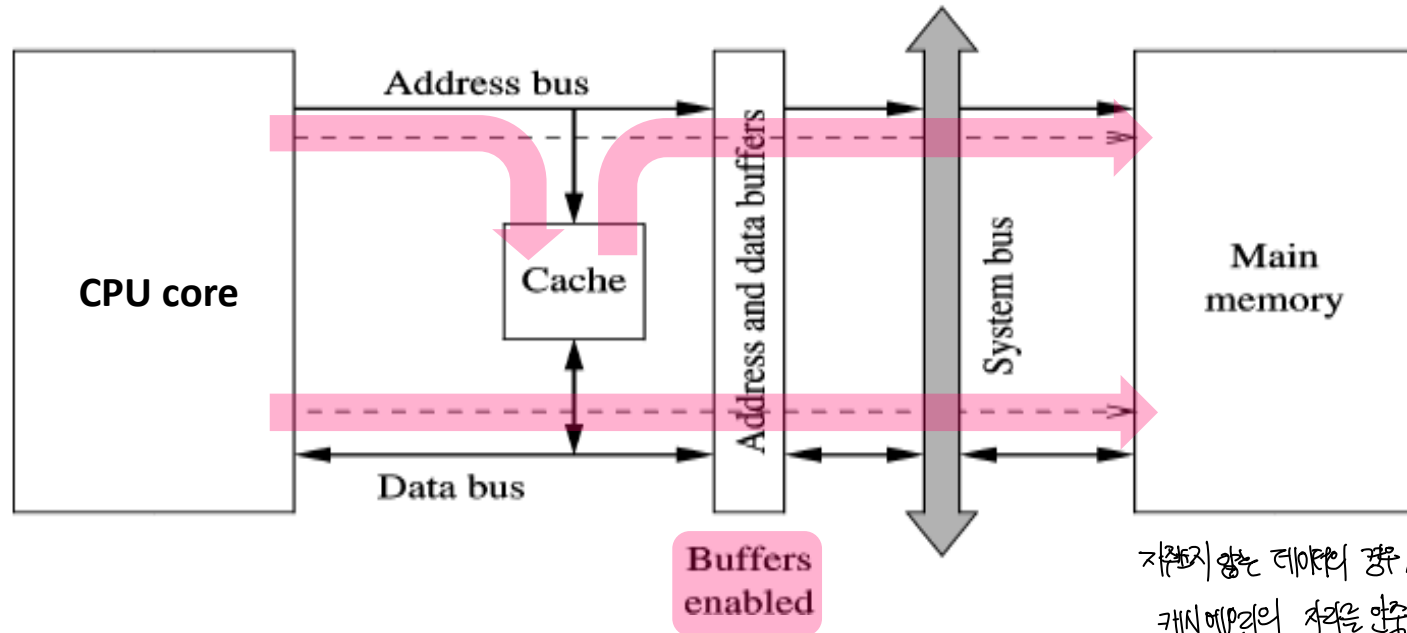
MISS

- Cache-Write Operation

- **Write miss**^②: no write-allocate ⇒ 캐시에 없기에 해당 address의 공간을 내어주지 않는 것

- CPU core stalls and simply send write-data to main memory.

CPU 멈추고 메인메모리 address에 보내주고 바로 CPU에서 메인메모리로 데이터를 보내버림.



저장하지 않는 데이터의 경우,
캐시에 없기에 저장할 수 없음. ⇒ no write-allocate

수많은 SW/HW 문제가 존재함.

↓
필요치 않으면 100% 해결할 수 있는 대안으로 개선하며
① write-allocate or ② no write-allocate 인지를 결정함.

How does Cache Memory Work?

- Cache Memory Space Allocation

- How to allocate cache memory space for data from main memory? \Rightarrow 어떻게 캐시메모리에서 allocate 하니까?

– When?

어제

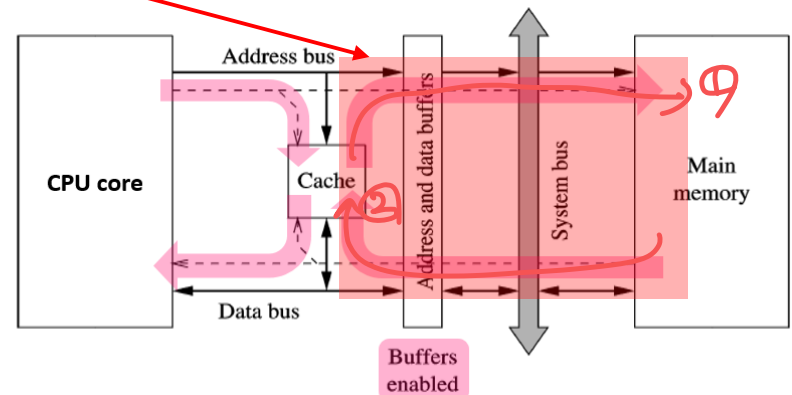
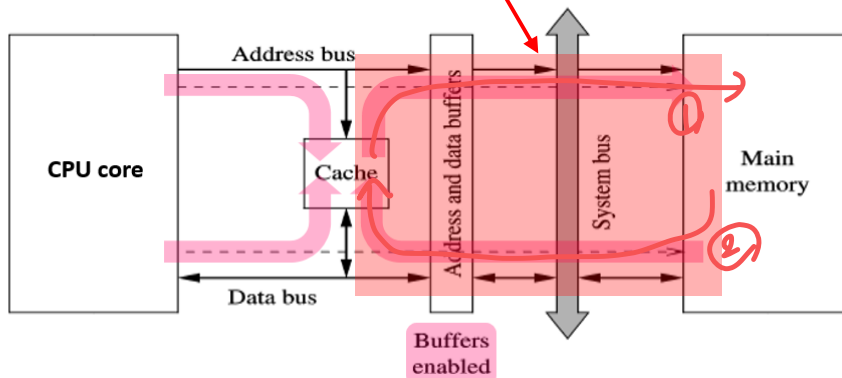
과거 상황

- **Read miss**

- CPU core stalls and requests the data to main memory.
- The data is copied in cache memory.
- CPU core repeats cache-read operation.

- **Write miss** : write-allocate

- CPU core stalls and requests the data to main memory.
- The data is copied in cache memory.
- CPU core repeats cache-write operation.



How does Cache Memory Work?

캐시관련 중요한 (문제)가
아직 안 나왔어!
이것은 중요하네.

- Cache Memory Space Allocation

- How to allocate cache memory space for data from main memory?

↳ miss 경우, 어떻게 캐시메모리에 allocate 할 것인가?

- In the aspect of data-placement & identification

↳ address gap의 차이를 정확하게 identify 하주는 필요함.

이것으로 데이터를 배치시켜줘야 함.

- Directed mapped
- Fully associative
- N-Way Set associative

⇒ 설계자들에게 의미있는 Cache,
사용자들은 걱정하지 않을 것
이런
상황에서는
하는 것

- In the aspect of data-replacement

- LRU (Least Recently Used) ⇒ 가장 안 쓴 것 (뜻 반대로)
- FIFO (First In First Out) ⇒ 선입선출
- Random ⇒ 랜덤

어떤 이유로 캐시메모리에서
데이터가 삭제 되었는가?
(MISS)

대신하는 거임.

Cache

- 1) 캐시는 데이터가 메모리에서
언제에 교체되었는지를 알아야 함.
구간 관리법이 다름에
이름은 같지만 "구간"이 필요함.
- 2) 기존 캐시메모리에 있는 데이터를
어떤 데이터를 뺄 것인가?