

Database Systems

Lecture11 – Chapter 12: Physical Storage Systems

Beomseok Nam (남범석)

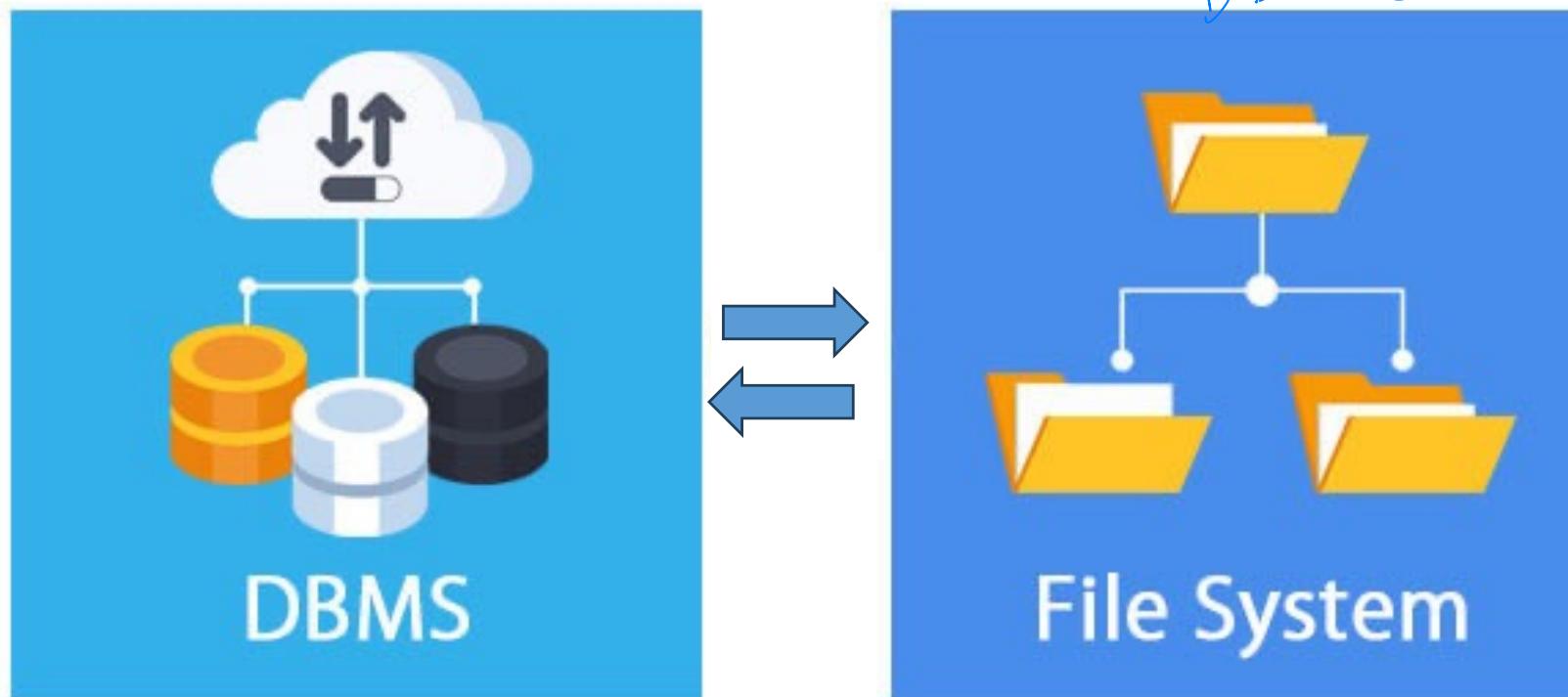
bnam@skku.edu

Table of Contents We Skip For Now

- Chap 8: Complex Data Types
 - XML, JSON, etc
- Chap 9: Application Development
 - HTML, Servlet, JavaScript, etc
- Chap 10: Big Data
 - Hadoop
- Chap 11: Data Analytics
 - Data Warehousing, OLAP, Data Mining

Part III: Storage Management

- DBMS is a system S/W designed to manage databases
- DBMS interacts with file systems for **durability**

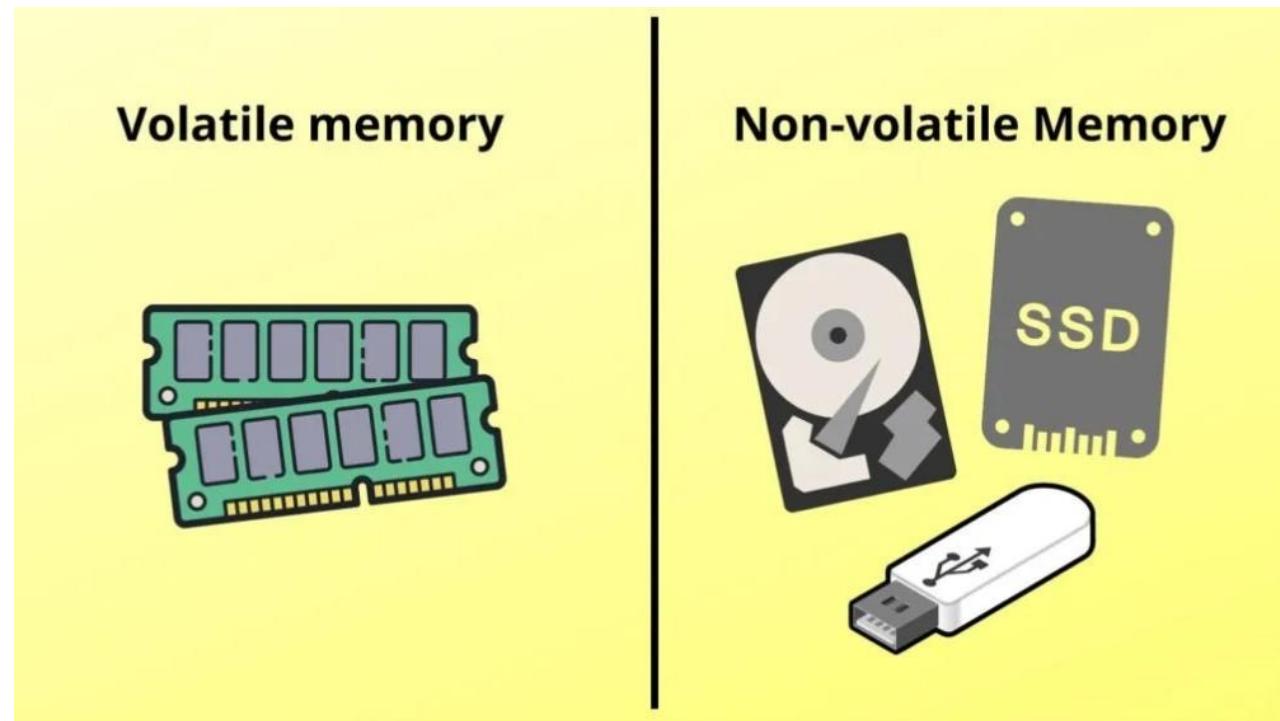


내적성

DBMS는 F.S 와 함께 실행됨

Classification of Storage Media

- Can differentiate storage into:
 - **volatile storage:** loses contents when power is switched off
(회발성기)
 - **non-volatile storage:**
(비-회발성)
 - Contents persist even when power is switched off.



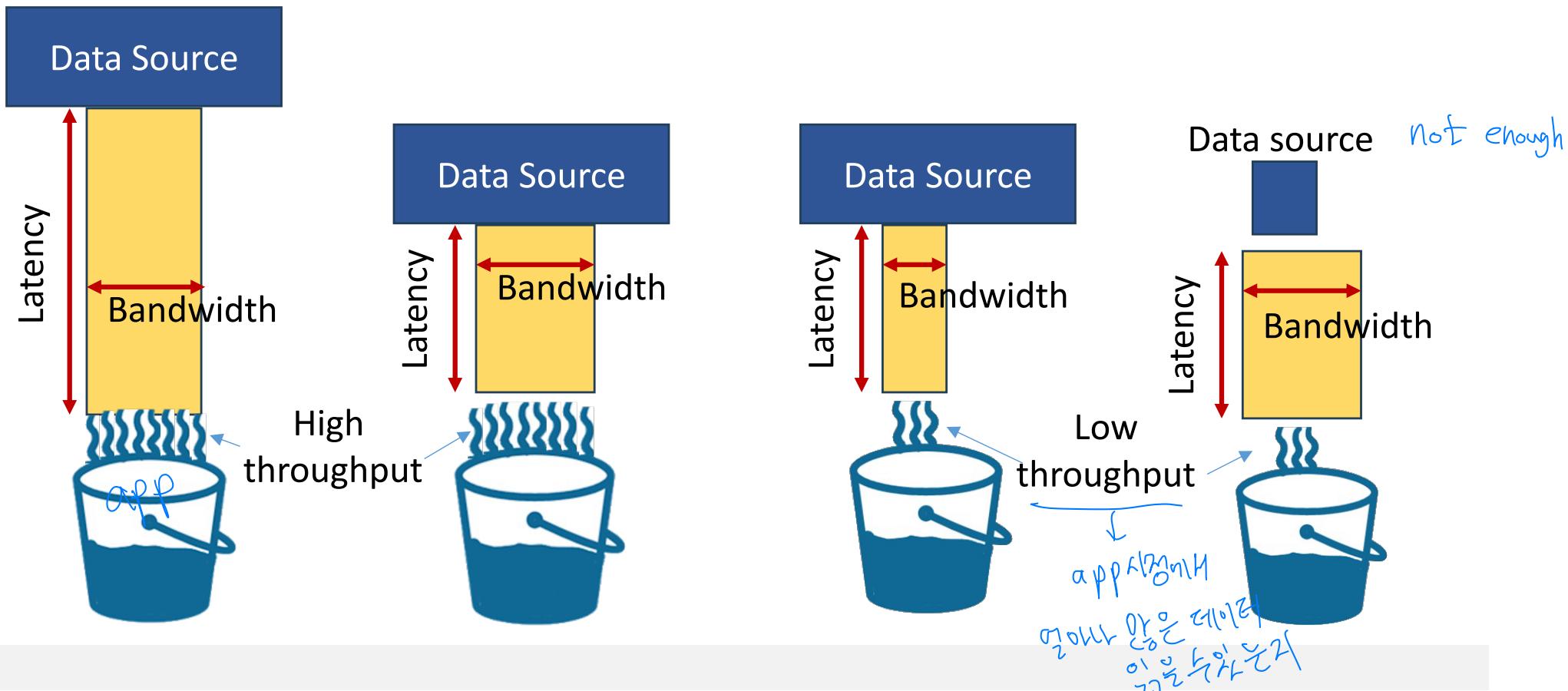
Performance of Storage Media

■ Factors affecting choice of storage media

- Latency → 시간(시간지연)

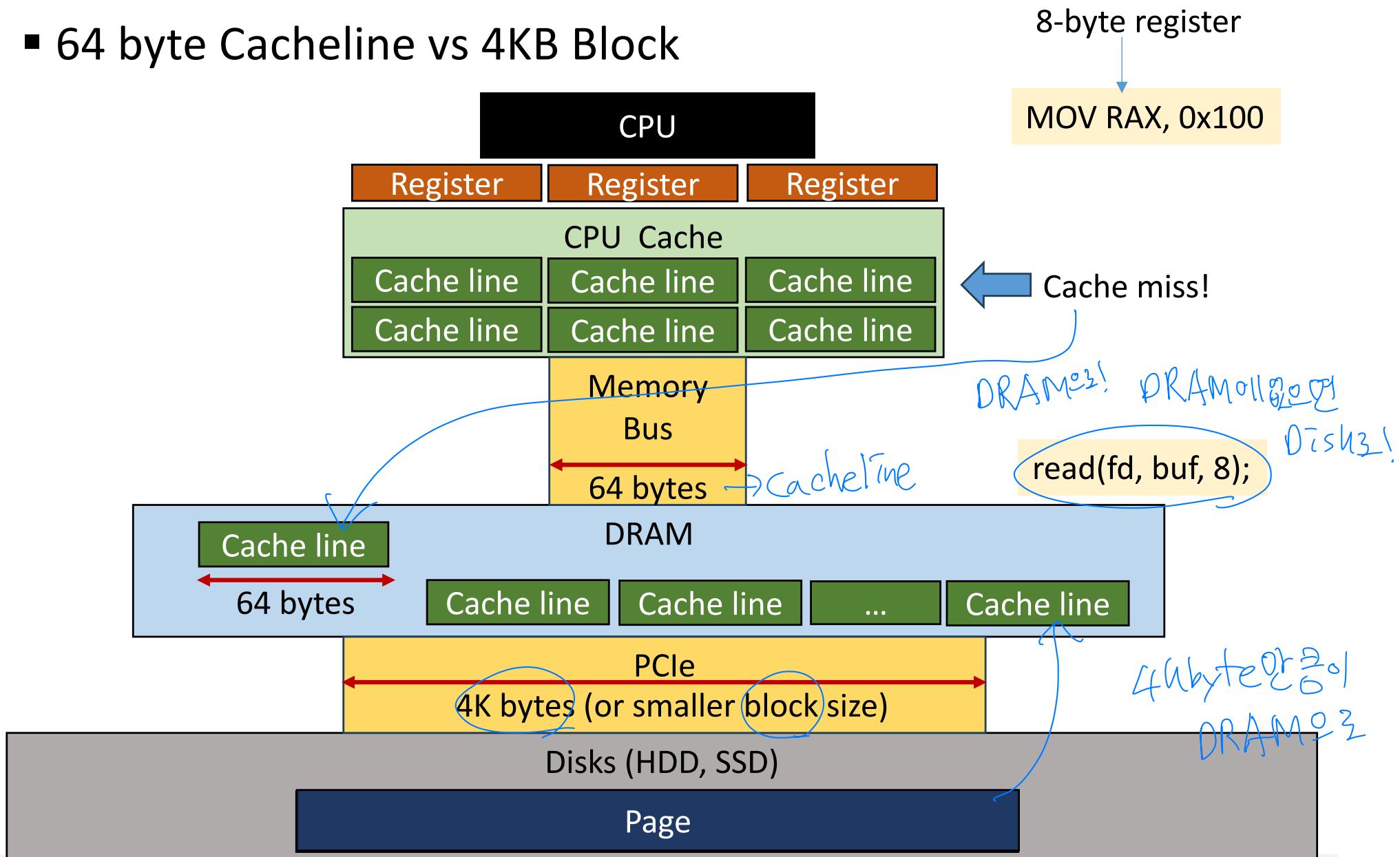
- Bandwidth \hookrightarrow I/O吞吐
Bandwidth metric: I/O operations per second (IOPS)
한 번에 처리 가능한 데이터 전송 가능성이? (device 시장)

- Cost



IO Granularity of Storage Media

- 64 byte Cacheline vs 4KB Block



Before HDD, there was Magnetic Tapes

최전기연기 운차집근

Tape Drive for Sequential Access



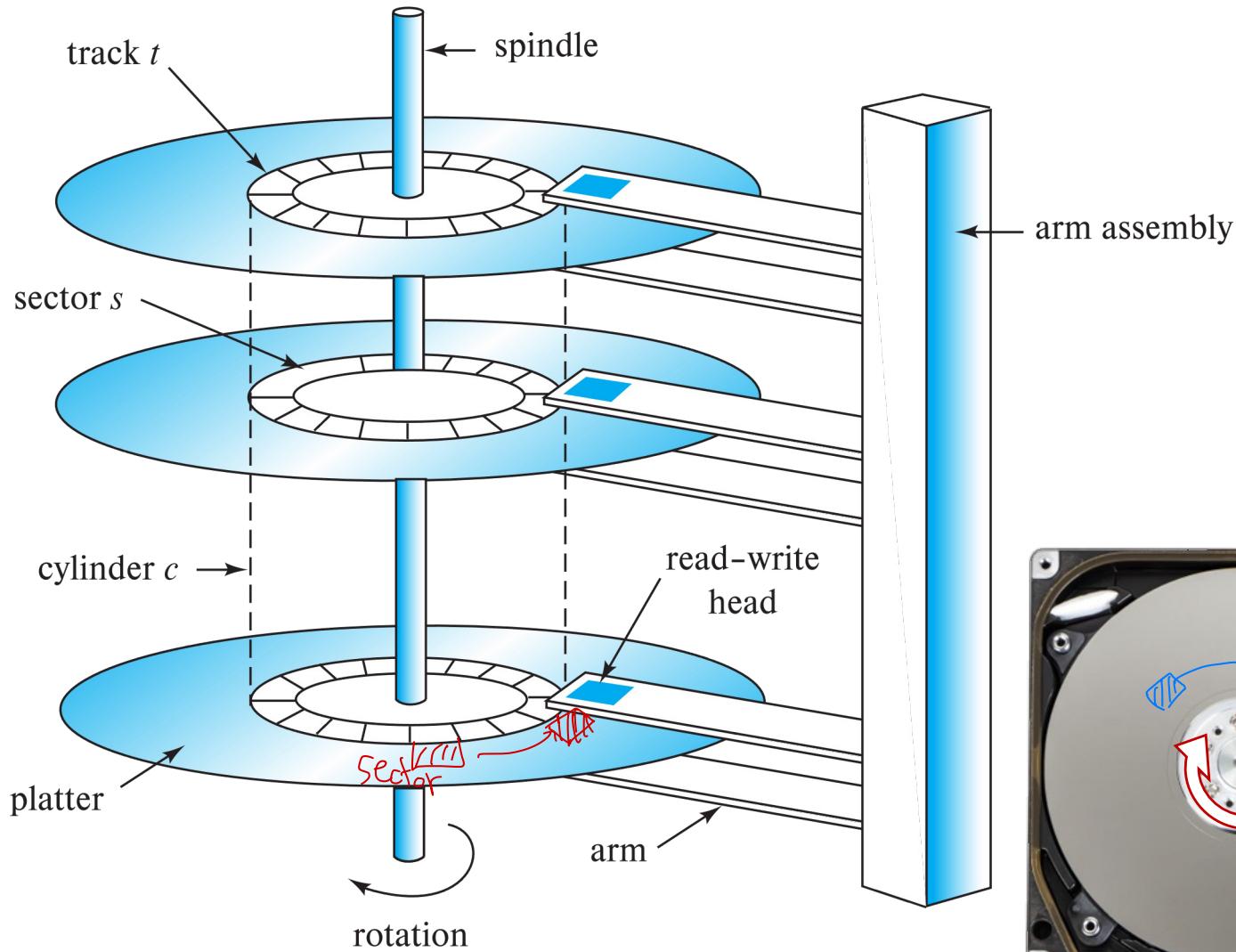
Floppy Disks for Random Access



- 1-Dimensional movement
 - read through the tape to find the desired data
 - Tape head does not move
- Still popular as backup devices in large data centers

- 2-Dimensional movement
 - Circular magnetic disk spins
 - Tape head moves

Magnetic Hard Disk Mechanism



Schematic diagram of magnetic disk drive

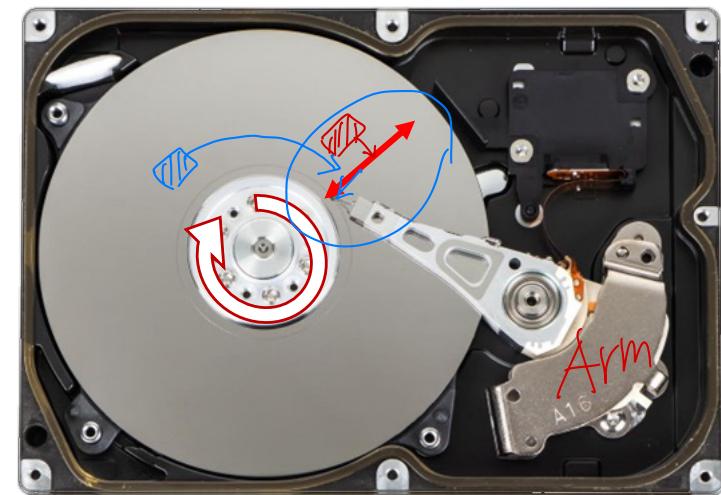


Photo of magnetic disk drive

Performance Measures of Disks

읽고쓰는 sector 를 접근하는 시간

- **Access time** – the time it takes from when an I/O request is issued to when data transfer begins. Consists of:

- **Seek time** – time it takes to reposition the arm over the correct track.
 - 4 to 10 milliseconds on typical disks
- **Rotational latency** – time it takes for the sector to be accessed to appear under the head.
 - 4 to 11 milliseconds on typical disks (5400 to 15000 r.p.m.)
- Overall latency is 5 to 20 msec depending on disk model

arm이 움직이는

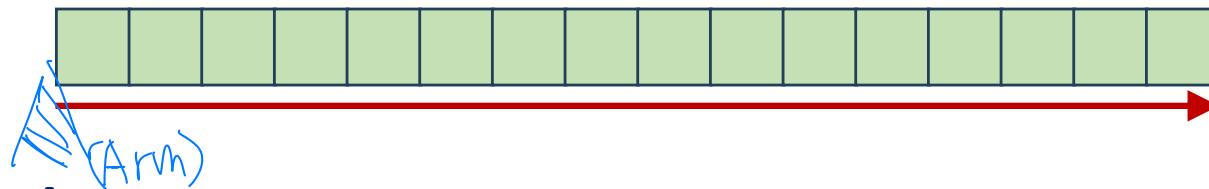
- **Data-transfer rate** – the rate at which data can be retrieved from or stored to the disk.
 - 25 to 200 MB per second max rate, lower for inner tracks

Performance Measures (Cont.)

- **Disk block** is a logical IO unit for storage
 - 4 to 16 kilobytes typically
 - Problem of Large block: unnecessary data transfer

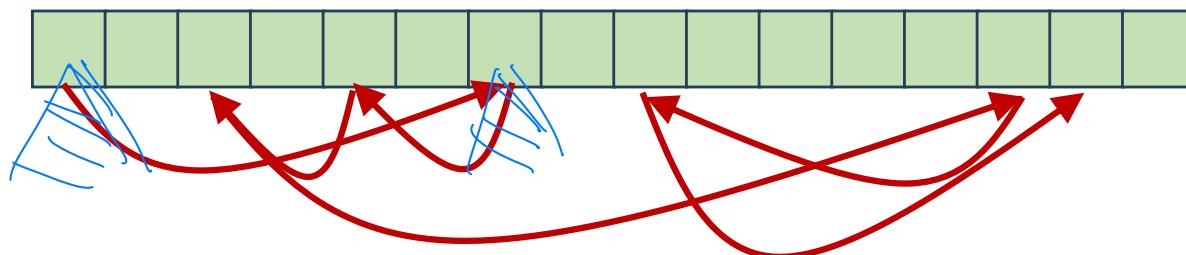
▪ Sequential access pattern

- Disk seek required only for first block



▪ Random access pattern

- Each access requires a seek
- Transfer rates are low since a lot of time is wasted in seeks



Flash Storage

■ NAND flash

- page: 512 bytes to 4 KB
 - 20 to 100 microseconds for a page read
- Page can only be written once
 - Must be erased to allow rewrite



■ Solid state disks (SSD)

- Use standard block-oriented disk interfaces, but store data on multiple flash storage devices internally
- Transfer rate of up to 500 MB/sec using SATA, and up to 3 GB/sec using NVMe PCIe



↑
HDD → 직렬
SSD → 전기적
↓
randomaccess
↓
Samsung 840 EVO

Flash Storage (Cont.)

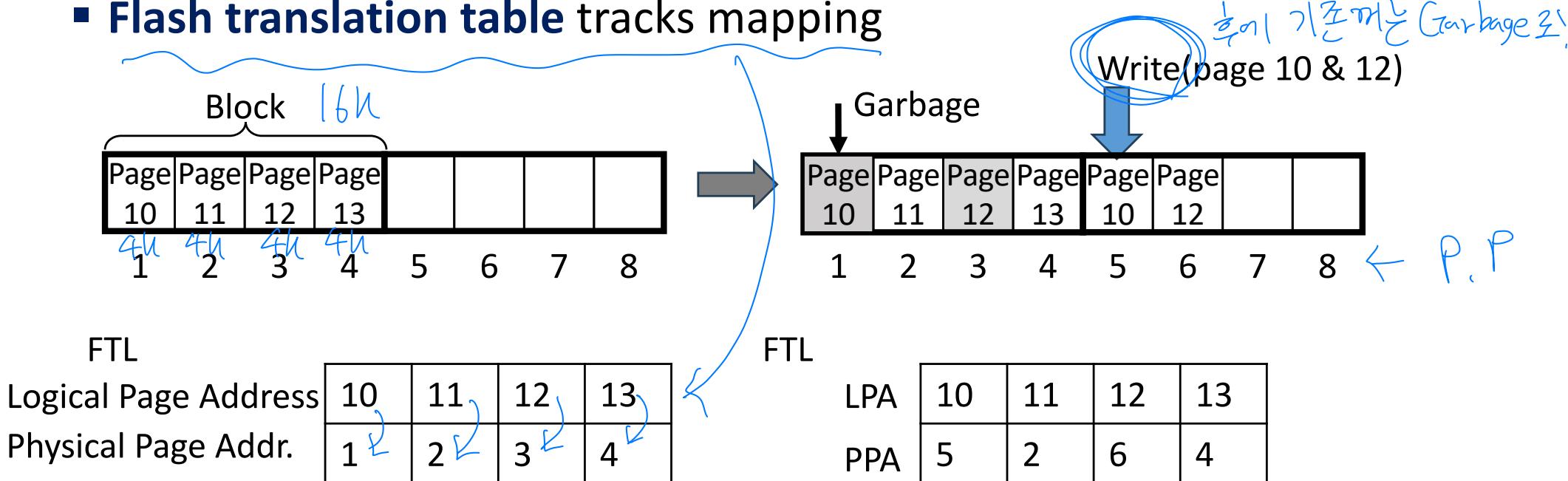
- Erase happens in units of erase block

- Erase block typically 256 KB to 1 MB (128 to 256 pages)

block 단위로 지울

- Remapping of logical page addresses to physical page addresses

- Flash translation table tracks mapping

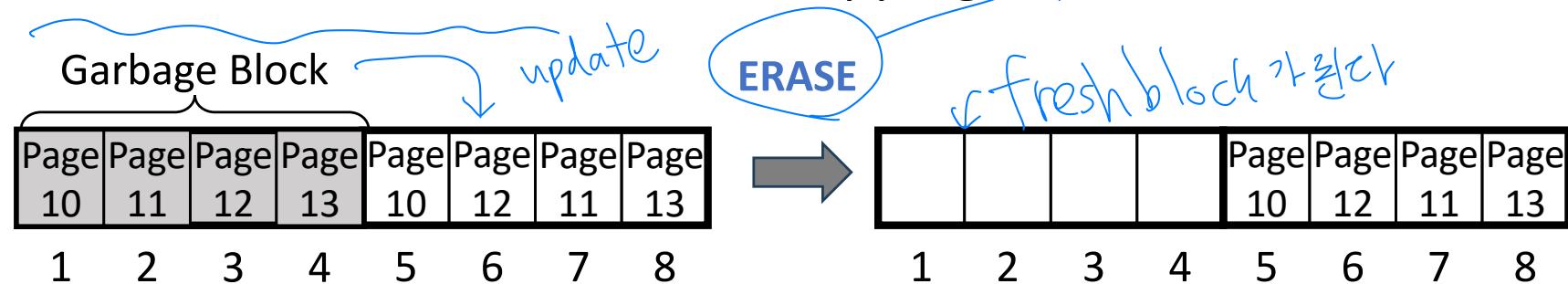


Flash Storage (Cont.)

-read 4h
-write 4h
-erase 4h X N

여기 page가 많아

- Erase happens in units of **erase block**
 - Erase block typically 256 KB to 1 MB (128 to 256 pages)
- Remapping** of logical page addresses to physical page addresses
- Flash translation table** tracks mapping



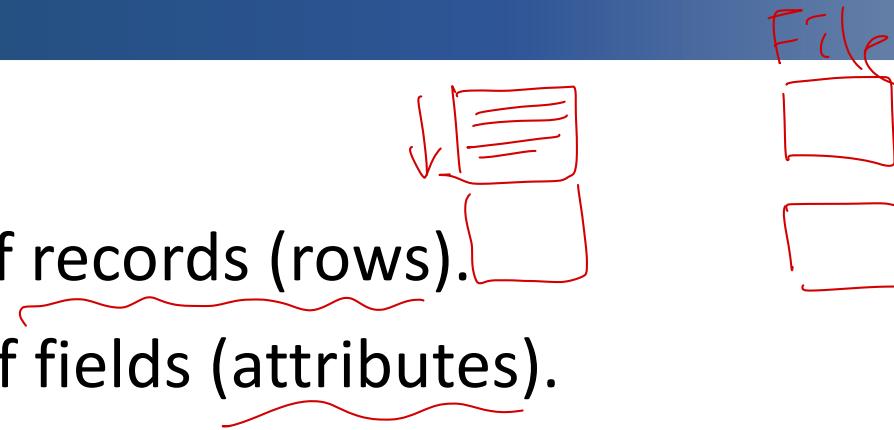
FTL				
LPA	10	11	12	13
PPA	5	7	6	8

FTL				
LPA	10	11	12	13
PPA	5	7	6	8

- After 100,000 to 1,000,000 erases, block becomes unusable
 - wear leveling** becomes necessary

File Organization

- Each Table = File(s)
- Each file is a sequence of records (rows).
- A record is a sequence of fields (attributes).



- We assume that records are smaller than a disk block.
- We will consider variable length records later

record << 4 kB

Fixed-Length Records

- Simple approach:

- File offset for record k is $n * (k - 1)$, where n is the size of record.

fixed size!

7개의 Tuple \rightarrow 28byte

$28 \times 6 = 168$

168

± 4
172

Katz의 시작주소: 172byte

	^{int} 4 bytes	^{name} 10 bytes	10 bytes	^{salary (int)} 4 bytes
record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califeri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000

Q: What's the file offset for the name of record 6?

Fixed-Length Records

- Let's delete record k :

- Option 1) move records $k + 1, \dots, n$ to $k, \dots, n - 1$**
- Option 2)
- Option 3)

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califeri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000

Fixed-Length Records

- Let's delete record k :

- Option 1) move records $k + 1, \dots, n$ to $k, \dots, n - 1$**
- Option 2)
- Option 3)

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califeri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000

Fixed-Length Records

- Let's delete record k :
 - Option 2) move record n to k**
 - Option 3)
- E.g. Record 3 deleted and replaced by record 9

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 9	76766	Crick	Biology	72000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califeri	History	62000
record 8	76543	Singh	Finance	80000

Fixed-Length Records

- Let's delete record k :

- Option 3) do not move records, but link all free records on a *free list*

header
record 0
record 1
record 2
record 3
record 4
record 5
record 6
record 7
record 8
record 9
record 10
record 11

			Free
10101	Srinivasan	Comp. Sci.	65000
			Free
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
33456	Gold	Physics	87000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000

FreeSpace[1|2|3|4|5] linked list update

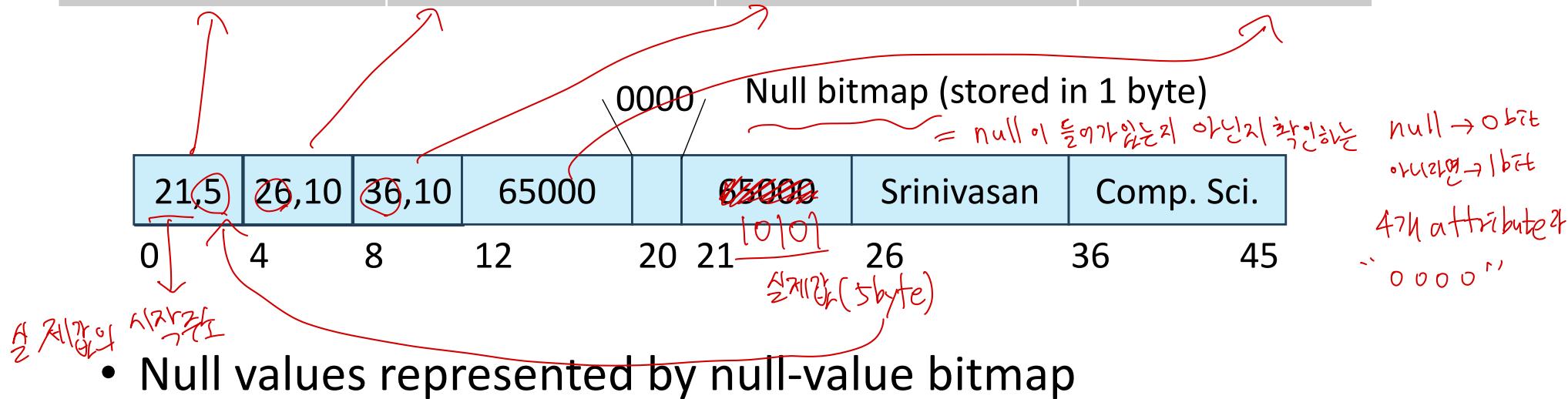
Variable-Length Records

■ Variable-length record ^{"VLR"}

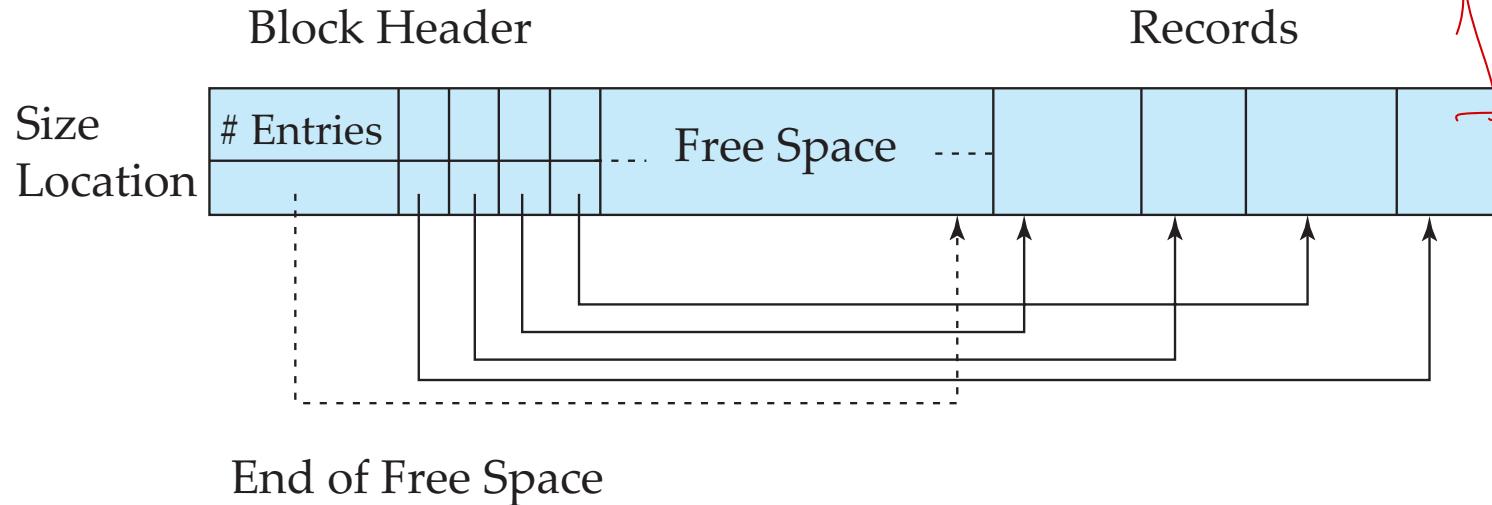
- A tuple with column(s) of varchar type
- An array of <offset, length> pairs points to actual data stored after the fixed-length array

record을 만드는 과정

ID (varchar)	name (varchar)	dept_name (varchar)	salary (8)
10101	Srinivasan	Comp. Sci.	65000



Page Layout with Variable-Length Records



- Slotted page header contains: *meta data*
 - number of record entries
 - end of free space in the block
 - location and size of each record

Slotted Page Structure

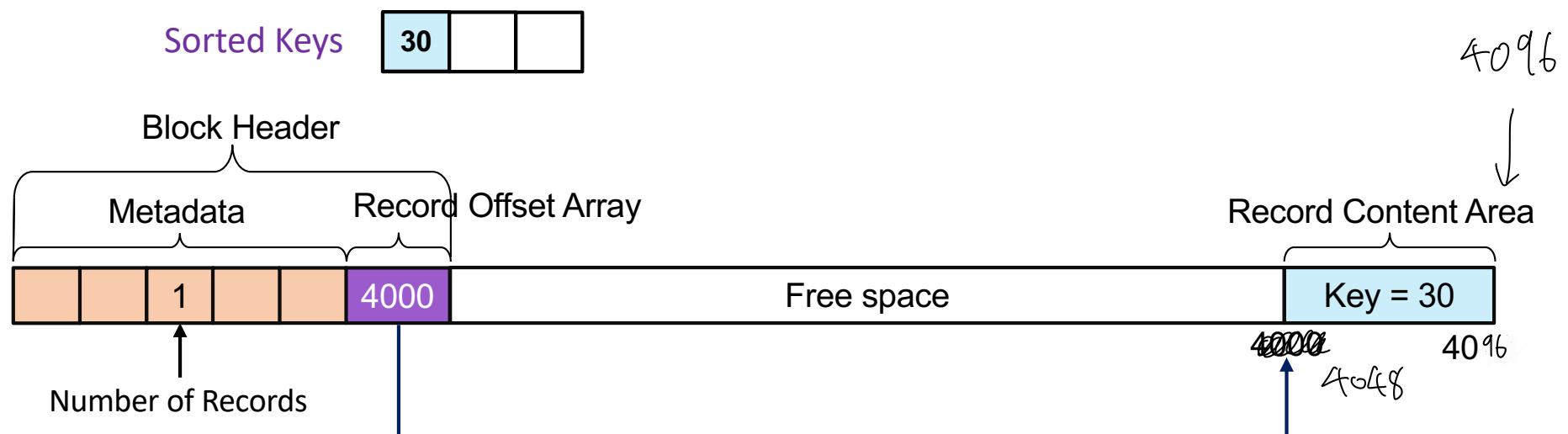
- **Slotted page** can sort records without moving them

- Eg. Insert variable length records:

- primary key 30 (48 bytes) →

- primary key 50 (100 bytes) →

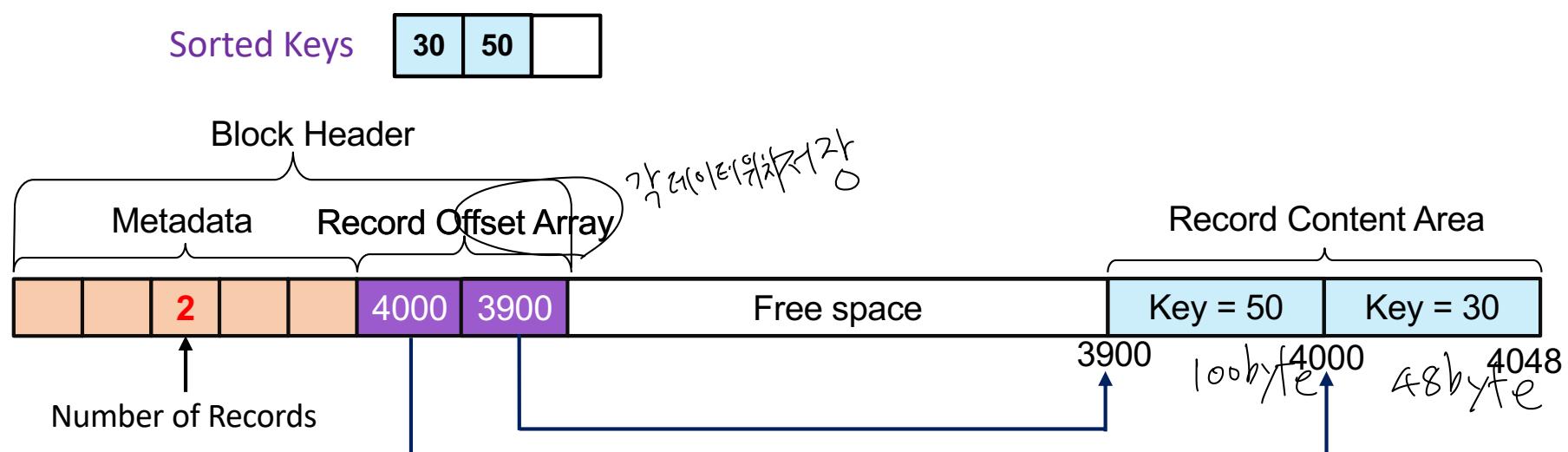
- primary key 40 (50 bytes)



Slotted Page Structure

- **Slotted page** can sort records without moving them

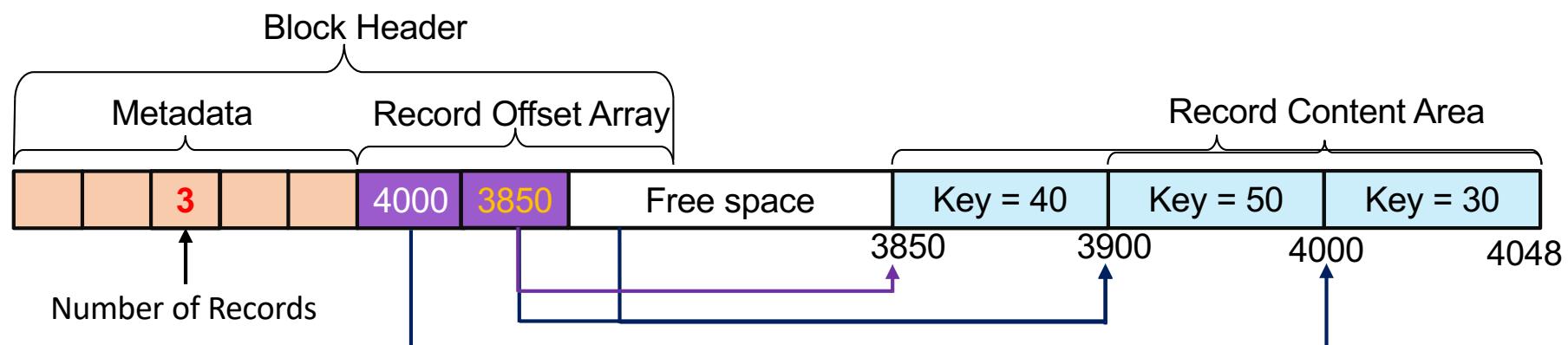
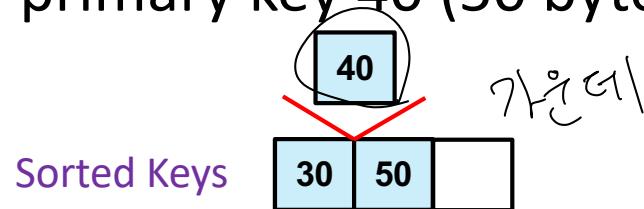
- Eg. Insert variable length records with
 - primary key 30 (48 bytes) →
 - primary key 50 (100 bytes) →
 - primary key 40 (50 bytes)



Slotted Page Structure

- **Slotted page** can sort records without moving them

- Eg. Insert variable length records with
 - primary key 30 (48 bytes) →
 - primary key 50 (100 bytes) →
 - primary key 40 (50 bytes)

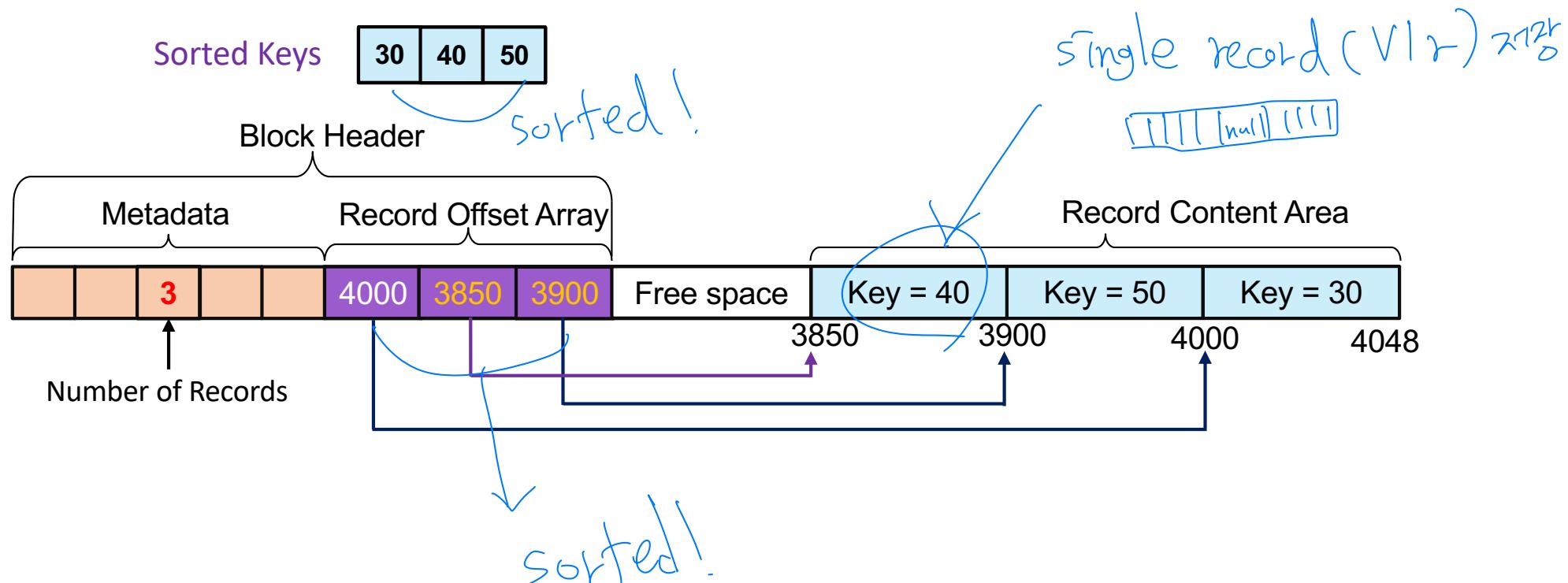


Slotted Page Structure

- Slotted page can sort records without moving them

- Eg. Insert variable length records with
 - primary key 30 (48 bytes) →
 - primary key 50 (100 bytes) →
 - primary key 40 (50 bytes)

page²
page는 고정



Organization of Records in Files

⇒ 파일은 반드시 고정
(연속적 page)

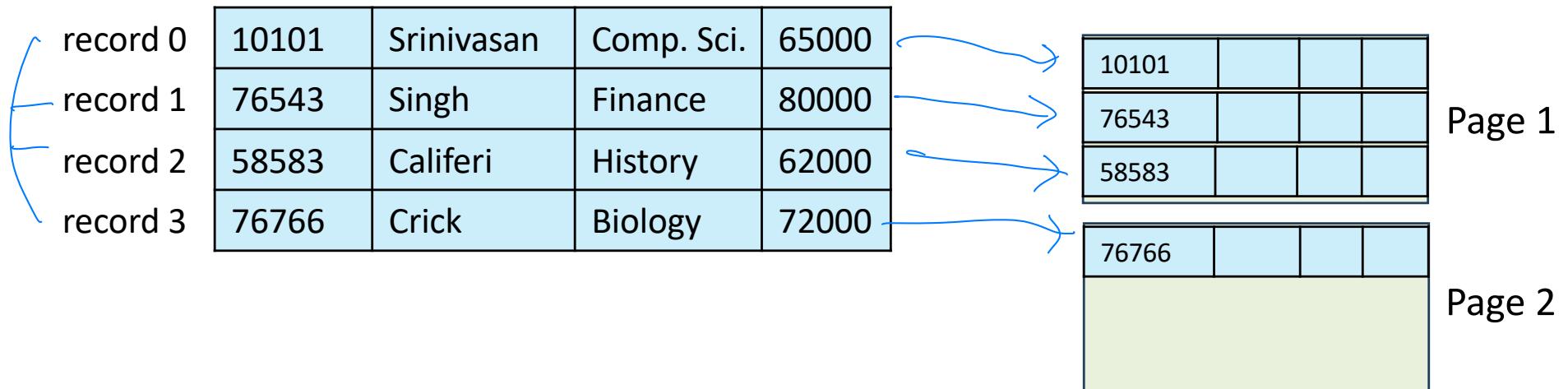
- **Heap** – record can be placed in any free space in the file
- **Sequential** – store records in sorted order by key
- **B⁺-tree file organization**
 - Dynamic ordered index (i.e., sorted by key)
 - More on this in Chapter 14
- **Hashing**
 - Unordered but fast index
 - More on this in Chapter 14

Heap File

bunch of random obj

- **Similar to Record Content Area of Slotted Page Structure**

- Records can be placed anywhere in the file
- Records do not move once allocated



Heap File

- Similar to Record Content Area of Slotted Page Structure
 - Records can be placed anywhere in the file
 - Records do not move once allocated

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	76543	Singh	Finance	80000
record 3	76766	Crick	Biology	72000

10101				
76543				
remove				

Page 1

76766				

Page 2

Heap File

- Similar to Record Content Area of Slotted Page Structure
 - Records can be placed anywhere in the file
 - Records do not move once allocated

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	76543	Singh	Finance	80000
record 3	76766	Crick	Biology	72000
record 4	22222	Einstein	Physics	95000

free space
available

10101				
76543				
22222				

Page 1

76766				

Page 2

Heap File

- Similar to Record Content Area of Slotted Page Structure
 - Records can be placed anywhere in the file
 - Records do not move once allocated

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	76543	Singh	Finance	80000
record 3	76766	Crick	Biology	72000
record 4	22222	Einstein	Physics	95000
record 5	33456	Gold	Physics	87000

10101				
76543				
22222				

Page 1

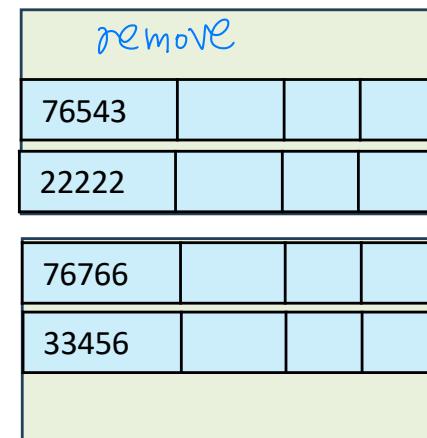
76766				
33456				

Page 2

Heap File

- Similar to Record Content Area of Slotted Page Structure
 - Records can be placed anywhere in the file
 - Records do not move once allocated

record 1	76543	Singh	Finance	80000
record 3	76766	Crick	Biology	72000
record 4	22222	Einstein	Physics	95000
record 5	33456	Gold	Physics	87000



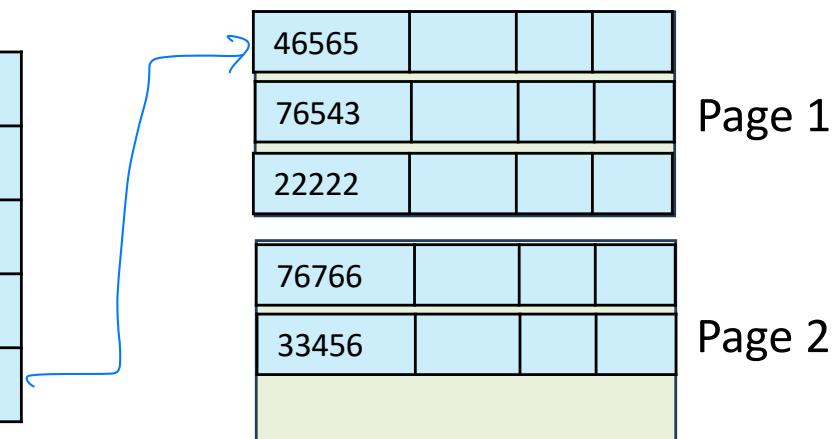
Page 1

Page 2

Heap File

- Similar to Record Content Area of Slotted Page Structure
 - Records can be placed anywhere in the file
 - Records do not move once allocated

record 1	76543	Singh	Finance	80000
record 3	76766	Crick	Biology	72000
record 4	22222	Einstein	Physics	95000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000

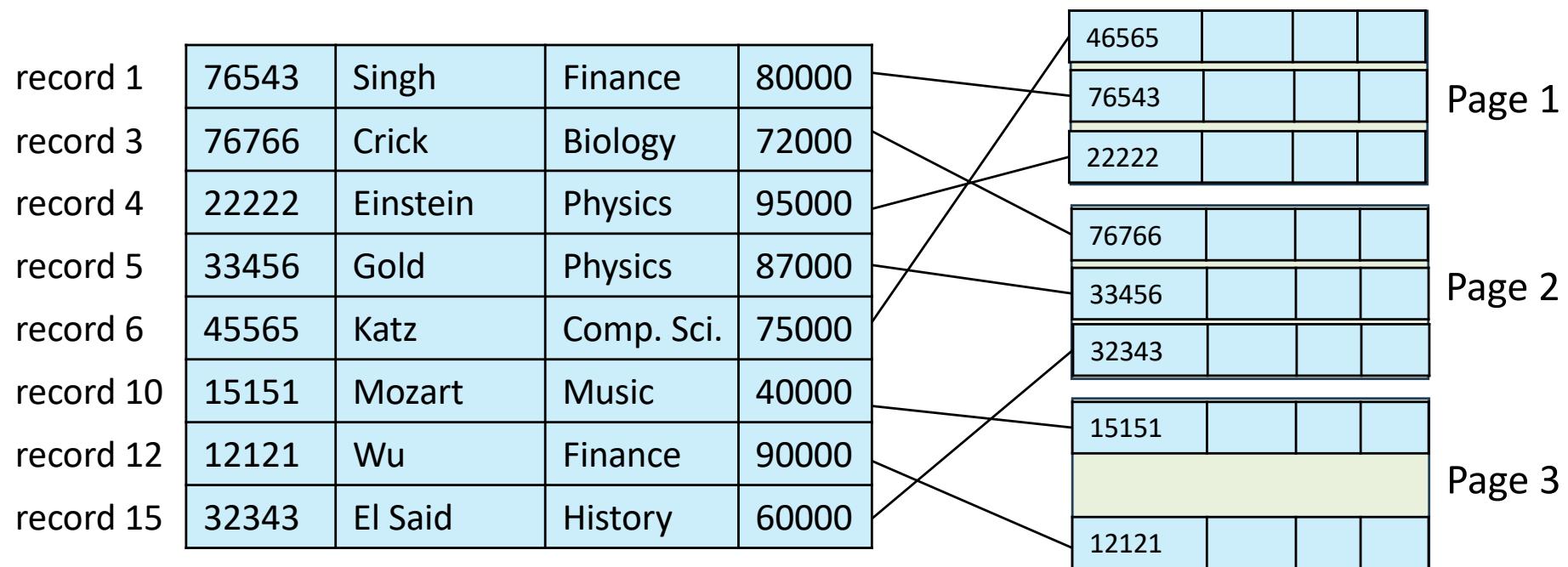


Heap File

- Similar to Record Content Area of Slotted Page Structure
 - Records can be placed anywhere in the file
 - Records do not move once allocated

Heap file는 Table을 같은 듯

결국 Random position이 된다
(순서가 없음)



Heap File

- Q: How to find free space in a heap file?

- A: Free-space map

- Each entry indicates fraction of free blocks

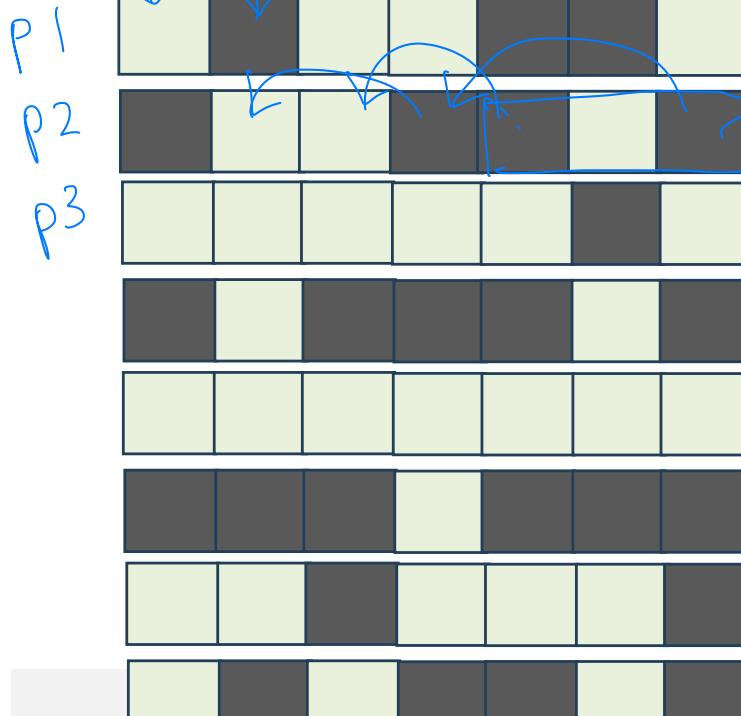
- Each entry size is a few bits ~ a byte

3bit → 8개 : 0~7

- E.g., 3 bits per block, value divided by $2^3=8$ indicates

fraction of free blocks out of 7

Freespace
0~7개 중 몇개는



연속도(개) 0~7
Page 1 (4: 100)

Page 2 (3: 011)

Page 3 (6: 110)

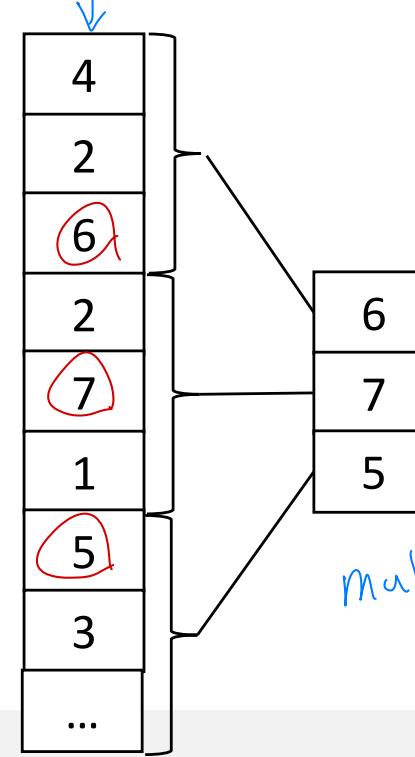
Page 4 (2: 010)

Page 5 (7: 111)

Page 6 (1: 001)

Page 7 (5: 101)

Page 8 (3: 011)



Sequential File Organization

Tape Drive (순서적)

- Suitable for applications that perform sequential writes and sequential reads of the entire file
 - E.g., LSM tree-based key-value stores (RocksDB)
- The records in the file are ordered by search-key

free list

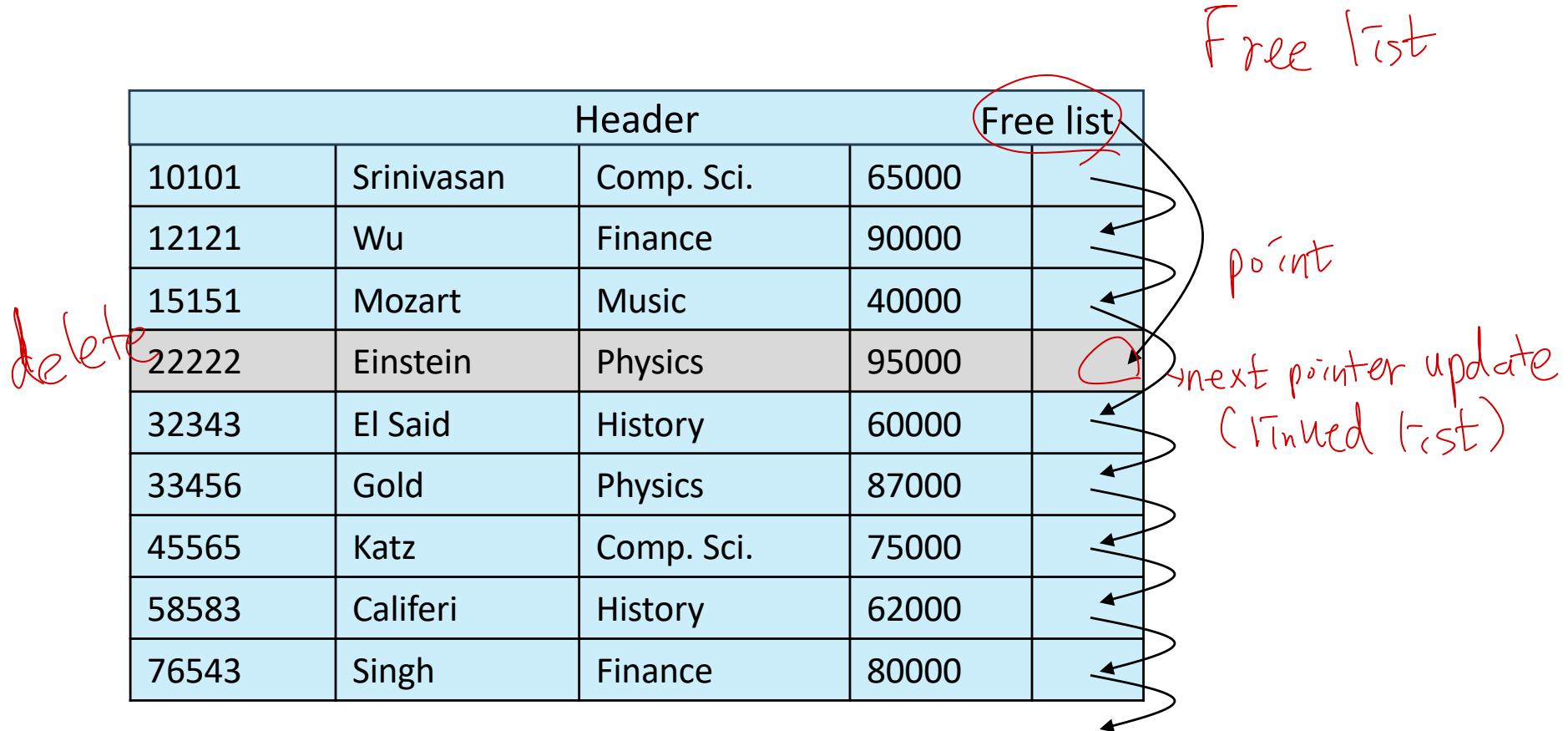
Header				
10101	Srinivasan	Comp. Sci.	65000	
12121	Wu	Finance	90000	
15151	Mozart	Music	40000	
22222	Einstein	Physics	95000	
32343	El Said	History	60000	
33456	Gold	Physics	87000	
45565	Katz	Comp. Sci.	75000	
58583	Califeri	History	62000	
76543	Singh	Finance	80000	

delete 10101
cost 12121
for (de shift
for
Free list

→ Free list

Sequential File Organization (Cont.)

- Deletion – use pointer chains



Sequential File Organization (Cont.)

■ Insertion

- if there is free space, insert there

Header				
10101	Srinivasan	Comp. Sci.	65000	
12121	Wu	Finance	90000	
15151	Mozart	Music	40000	
76766	Crick	Biology	72000	
32343	El Said	History	60000	
33456	Gold	Physics	87000	
45565	Katz	Comp. Sci.	75000	
58583	Califeri	History	62000	
76543	Singh	Finance	80000	

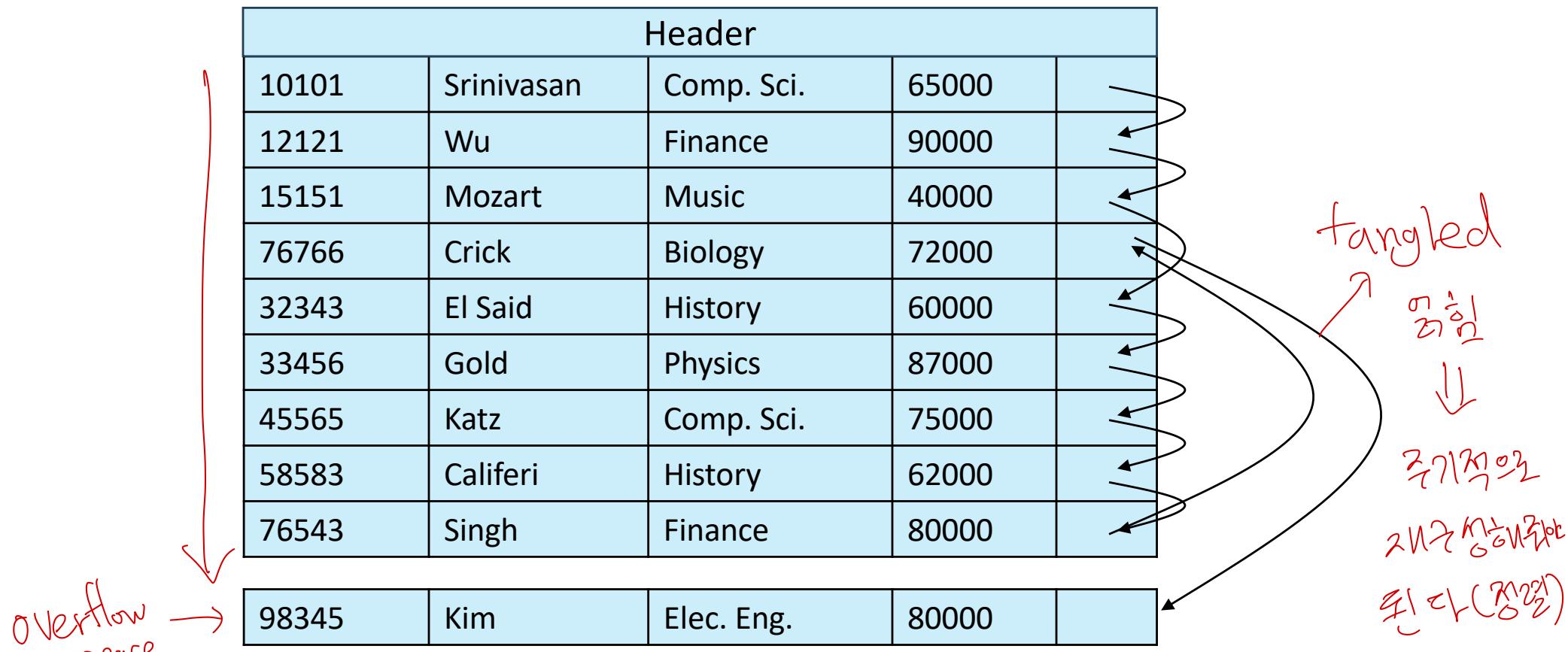
Insert →

76543 < 76766
017120101
next pointer
01713 point
∴ 정렬 완료

Sequential File Organization (Cont.)

■ Insertion

- if no free space, insert the record in an overflow block



- Need to completely reorganize the file from time to time to restore sequential order

Multitable Clustering File

- Store several tables in one file

- good for queries involving \bowtie
- bad for queries involving only table

department

instructor

multitable clustering
of *department* and
instructor

dept_name	building	budget
Comp. Sci.	Taylor	100000
Physics	Watson	70000

ID	name	dept_name	salary
10101	Srinivasan	Comp. Sci.	65000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
83821	Brandt	Comp. Sci.	92000

Comp. Sci.	Taylor	100000	
10101	Srinivasan	Comp. Sci.	65000
45565	Katz	Comp. Sci.	75000
83821	Brandt	Comp. Sci.	92000
Physics	Watson	70000	
33456	Gold	Physics	87000

파일을
▷ 툴을 쓰면 성능이 좋다
dep 테이블은 고(Comp)
Inst 테이블은 고(Physics)
(Srin → Katz → Brandt)
dep 알고(Physics)
inst 알고 --
join 풍경 임이 좋다
(Seek Time ↑
→ 성능 저하)

Column-Oriented Storage

- Also known as columnar representation
- Store each attribute of a relation separately
 - Commonly used in NoSQL
- Example

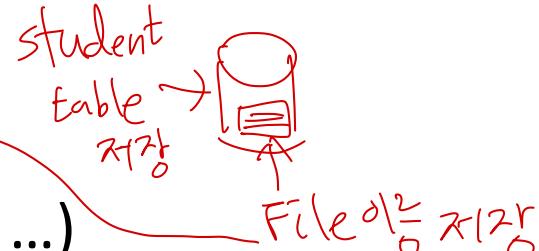
10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000



Data Dictionary

- Data dictionary (aka. system catalog) stores metadata

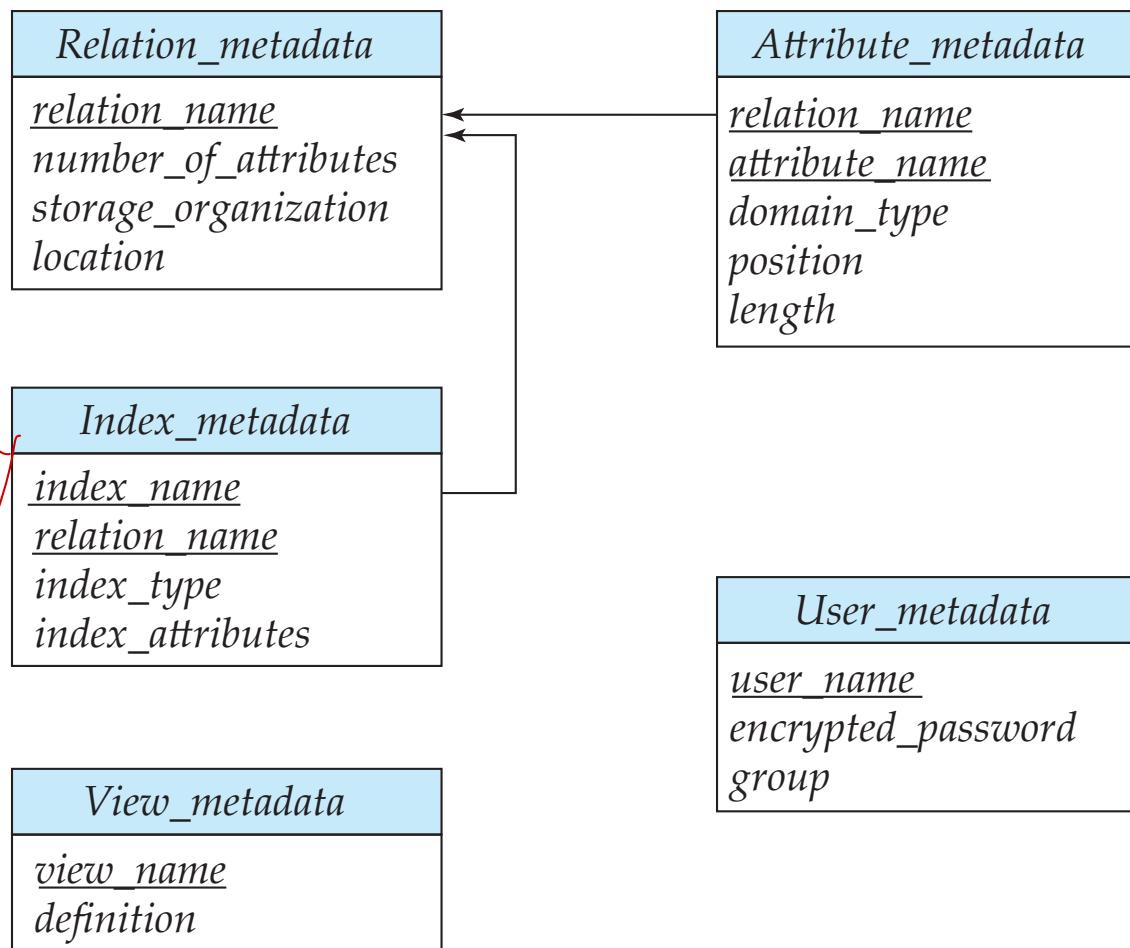
- Information about tables
 - File names of tables
 - File formats (Sequential, Heap, ...)
 - names, types and lengths of attributes of each table
 - Statistical and descriptive data
 - number of records in each table
 - Key distribution, etc
 - Information about indexes (Chapter 14)



Data Dictionary as a Database

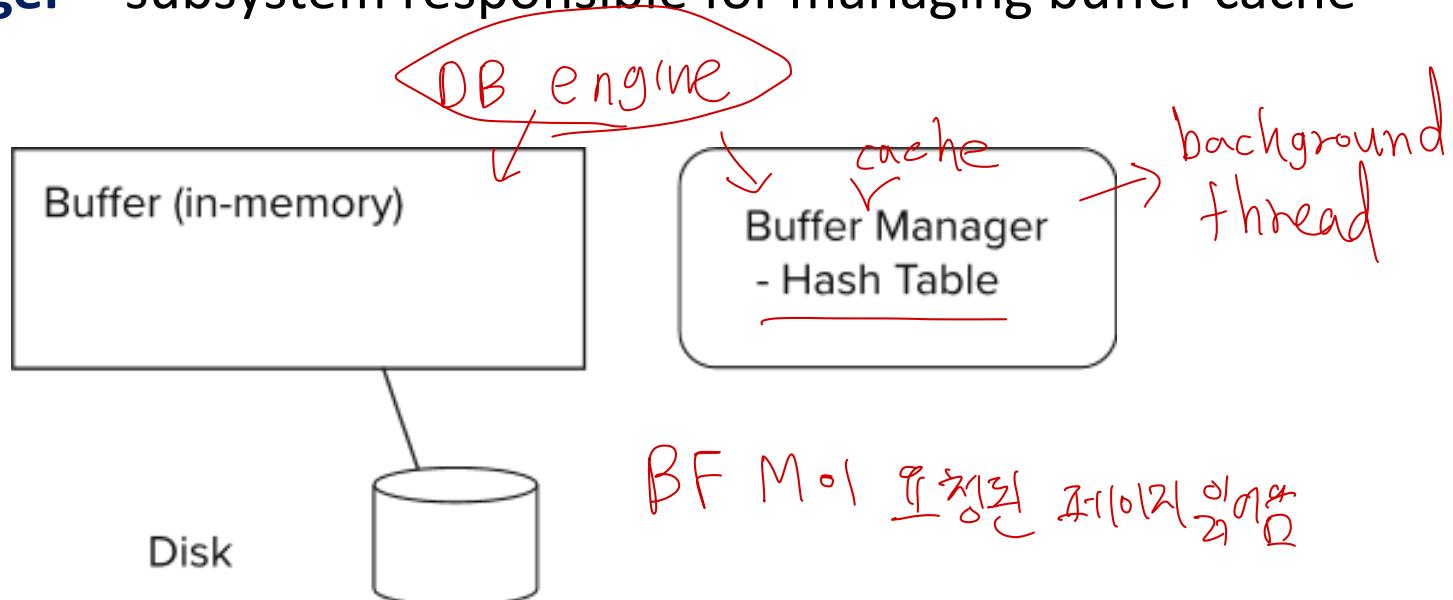
- Data Dictionary is also a **Database**
- Relational representation

E-R
Diagram
of
Dictionary



Storage Access

- Block (Page) = a unit of storage allocation and data transfer (I/O). Disk에 대한 접근 줄이자!
- DBMS tries to minimize the number of block transfers
 - by keeping as many blocks as possible in main memory.
- **Buffer (or Buffer Cache)** – portion of main memory available to store copies of disk blocks.
- **Buffer manager** – subsystem responsible for managing buffer cache

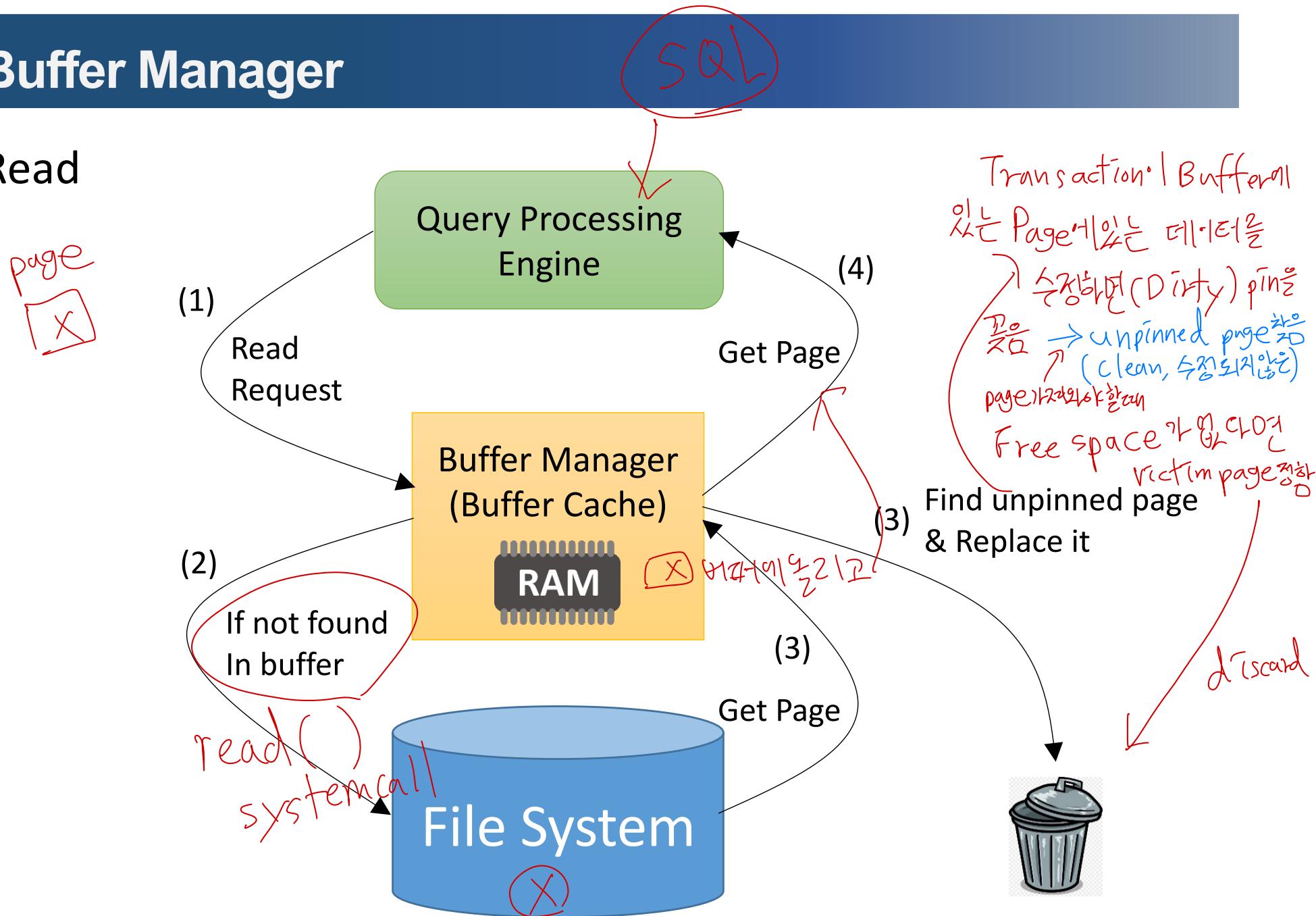


Buffer Manager

- DBMS calls on the buffer manager when it needs a block
(= page)
 - If the block is already in the buffer, buffer manager
 - returns the address of the block in main memory
 - If the block is not in the buffer, the buffer manager
 - Allocates space in the buffer for the block
 - Replacing some other block to make space for the new block.
 - Replaced block written back to disk if it was modified (i.e., dirty)
수정되었으면
 - Reads the block from the disk to the buffer, and returns the memory address of the block to the requester.

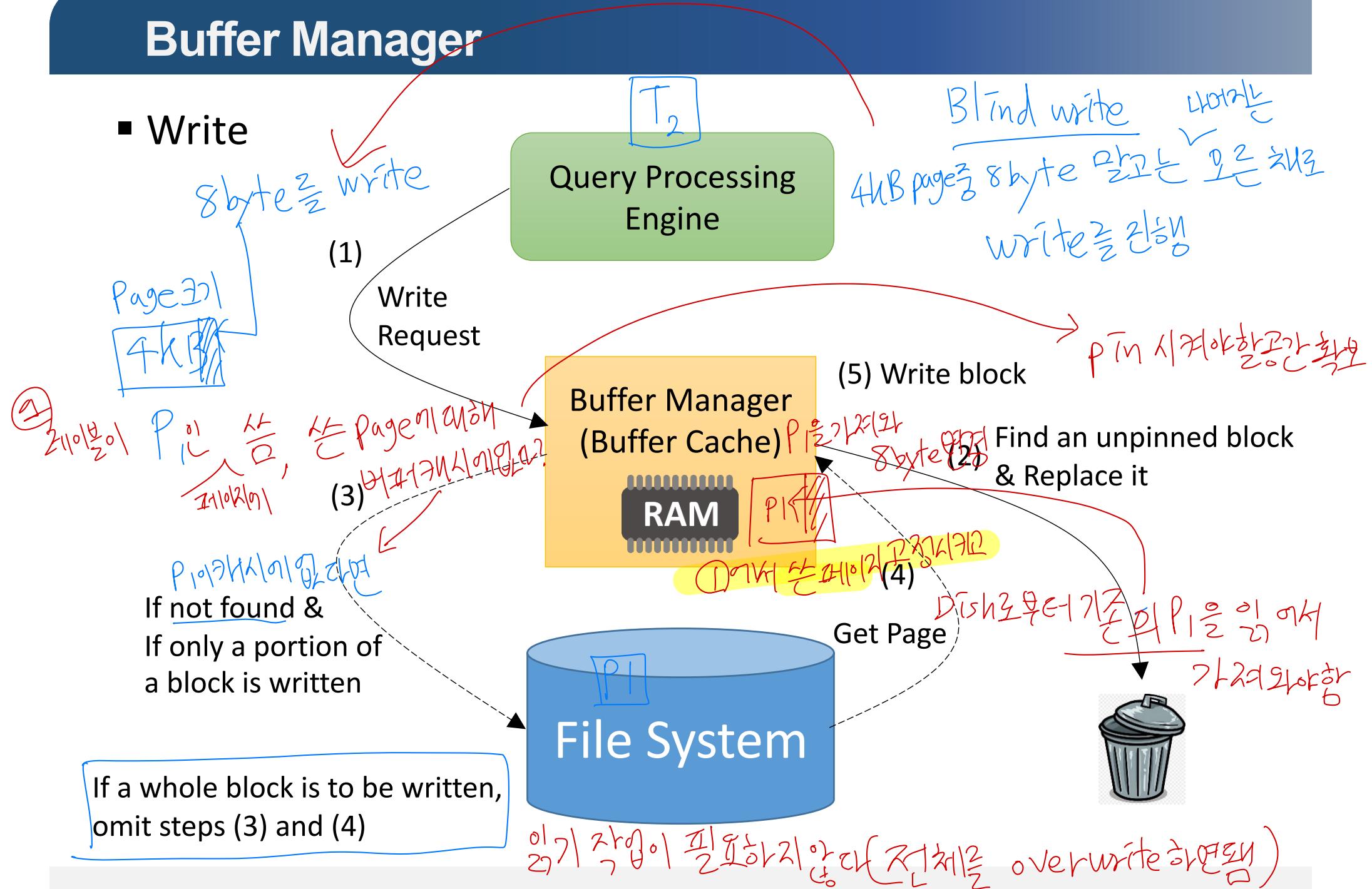
Buffer Manager

■ Read



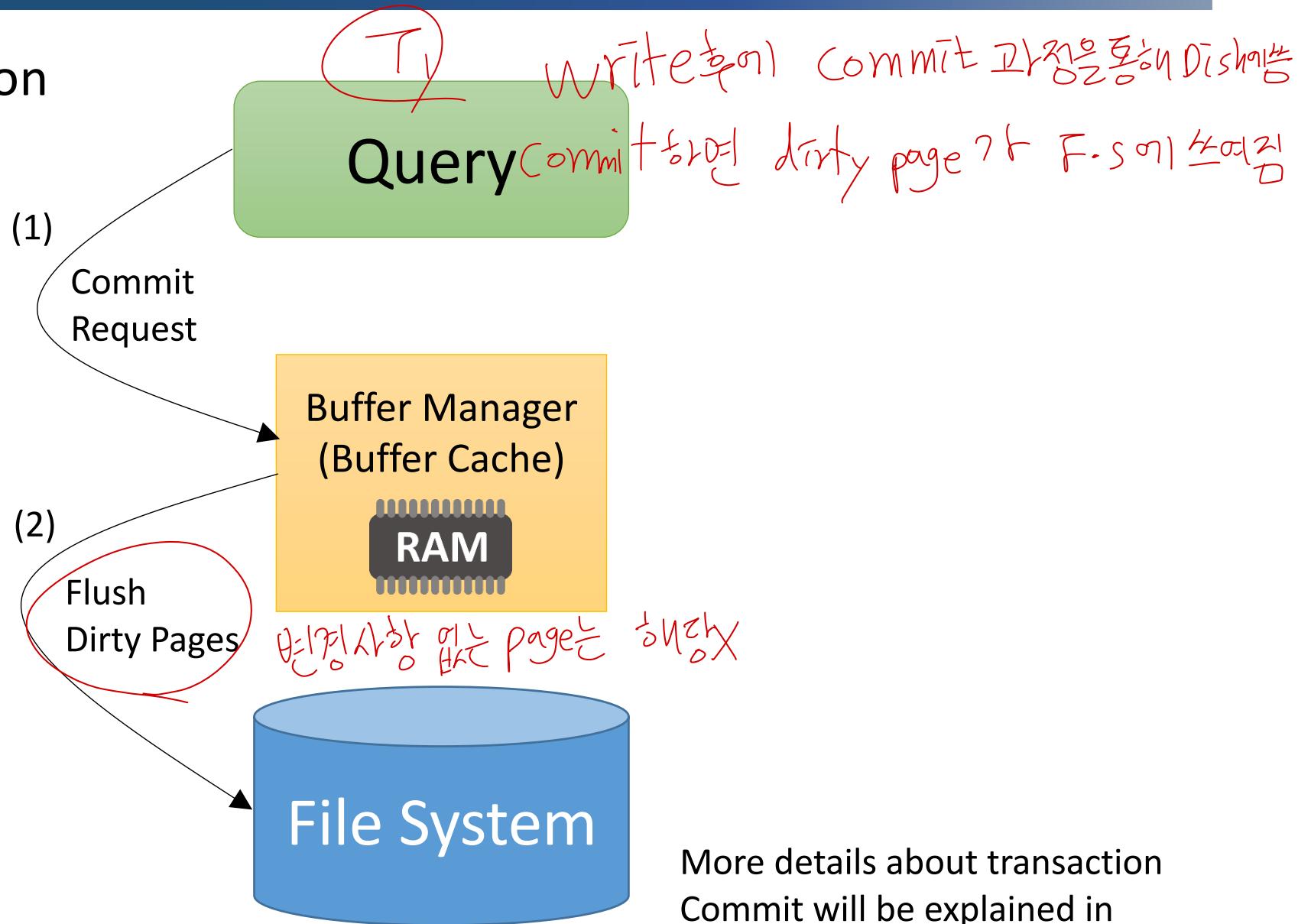
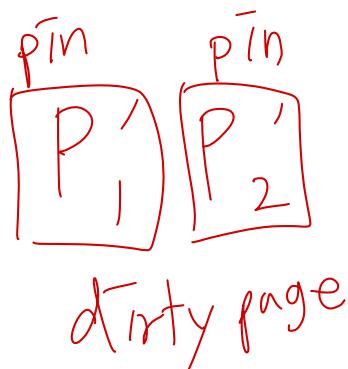
Buffer Manager

■ Write



Buffer Manager

■ Transaction Commit

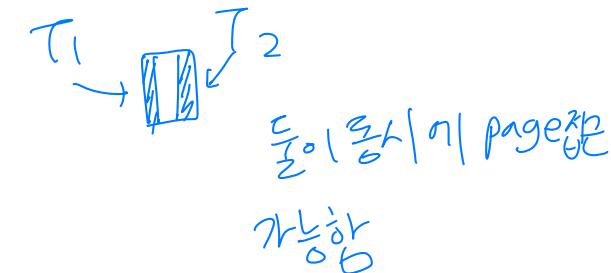


More details about transaction Commit will be explained in Chapter 19.

Buffer Manager & Replacement Policy

■ Pinned block is not allowed to be written to disks

- Pin before reading/writing data from a block *수정*
- Unpin when read /write is complete *수정X*
- Multiple concurrent pin/unpin operations possible
 - Keep a pin count, *pin count로 추적*
 - Can be evicted only if pin count = 0



■ Shared and exclusive locks on buffer

- Readers: need **shared lock** (read lock)
- Writers: need **exclusive lock** (write lock)
- **Locking rules:**
 - Only one process can get exclusive lock at a time
 - Shared lock cannot be given if exclusive lock is held
 - Multiple processes may be given shared lock concurrently

Buffer-Replacement Policies

- Most popular policy in OS: **Least Recently Used (LRU)**

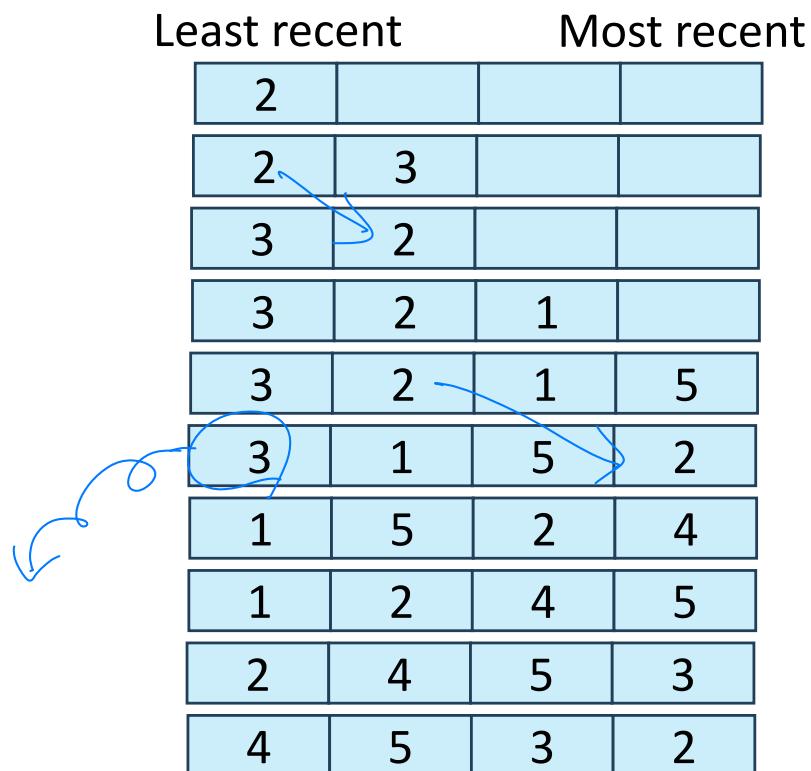
- Idea behind LRU

- use past block reference patterns to predict future references → Heuristic policy 

Block Reference Sequence

Page Request	Hit/Miss	Evict
2	Miss	
3	Miss	
2	Hit	
1	Miss	
5	Miss	
2	Hit	
4	Miss	3
5	Hit	
3	Miss	1
2	Hit	

Buffer Cache



Buffer-Replacement Policies

- Most popular policy in OS: **Least Recently Used (LRU)**

- LRU can be bad for database queries

- Sequential Flooding ← Join 때문에 LRU 안 쓸 뿐 (DBMS는!)

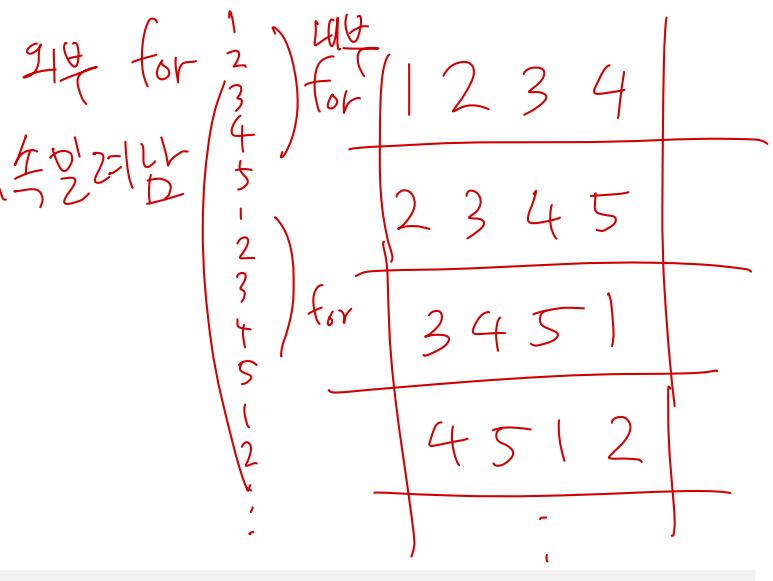
- Example: when computing the join
for each tuple tr of r do
 for each tuple ts of s do
 if the tuples tr and ts match

- If ts is larger than the buffer cache, no cache hits

수학적으로 이해가

캐시가 빠르게 채워지고 이전에 읽은 데이터가 계속 일려남

나중에 재사용되는게



Buffer-Replacement Policies in DBMS

- DBMS access disks in well-defined patterns
 - DBMS is aware of the disk access patterns for each query processing algorithm, e.g., sequential scan, join, etc.

metadata 이용 → 각 query가 맞는 효율적인 정책을 선택할 수 있다

- Data Dictionary manages statistics

- Probability that a request will reference a particular table
- Size of tables
- Key distribution of tables
- Etc.

- DBMS can do better than LRU

Buffer-Replacement Policies in DBMS

▪ **Toss-immediate** strategy

- Evicts a block as soon as the final record of that block has been processed even if free space is not required

→ avoid cache update

→ 자주 사용하지 않는 테이블에 대해 주로 사용
join

Block Reference Sequence

Page Request	Hit/Miss	Evict
...	...	
2	Hit	
4	Miss	DISK & GROWL
5	Hit	
3	Hit	
2	Hit	

Buffer Cache

Buffer-Replacement Policies in DBMS

▪ Most recently used (MRU) strategy

- After a block is processed, the block becomes the candidate for replacement, i.e., most recently used block.

지나간 temporal locality의 경향에 따른 것을 알면 MRU 사용

Block Reference Sequence

Page Request	Hit/Miss	Evict
...	...	
2	Hit	
4	Miss	2
5	Hit	
3	Hit	
2	Miss	

Buffer Cache

Least recent		Most recent	
3	2	1	5
3	1	5	2
3	1	5	4
3	1	4	5
1	4	5	3
1	4	5	2