

05/22

DBMS : Non-volatile System with volatile / Heap file

Lecture #04: Database Storage (Part II)

15-445/645 Database Systems (Fall 2023)
<https://15445.courses.cs.cmu.edu/fall2023/>
Carnegie Mellon University
Andy Pavlo + Jignesh Patel

what is heap file

1 Log-Structured Storage

Tuple보다 page 바르다!

Some problems associated with the Slotted-Page Design are:

- Fragmentation: Deletion of tuples can leave gaps in the pages, making them not fully utilized. *조작남*
- Useless Disk I/O: Due to the block-oriented nature of non-volatile storage, the whole block needs to be fetched to update a tuple. *메모리 없으면 디스크에서 가져와야함. 1개를 읽어도*
- Random Disk I/O: The disk reader could have to jump to 20 different places to update 20 different tuples, which can be very slow. *DBMS에 따라서 달라질 수 있다. / 20개가 다른 페이지에 있다면 20개를 읽어야함*

What if we were working on a system which only allows creation of new pages and no overwrites? The log-structured storage model works with this assumption and addresses some of the problems listed above.

Log-Structured Storage: Instead of storing tuples, the DBMS only stores the log records of changes to the tuples. The DBMS appends new log entries to an in-memory buffer without checking previous records and then writes out the changes sequentially to disk.

→ HDFS 같은 경우 Overwrite 안됨.

- Records contain the tuple's unique identifier, the type of operation (PUT/DELETE), and, for put, the contents of the tuple. *메모리 버퍼에 key-value pair tuple*
- To read a record, the DBMS scans the log file backwards from newest to oldest to find the most recent contents of the tuple. *Update x / 2개 밖에 없음*
- Fast writes, potentially slow reads. Disk writes are sequential and existing pages are immutable which leads to reduced random disk I/O. Good for append-only storage. *old ← new Sequential I/O. 20개가 다른 페이지에 있다면?*
- To avoid long reads, the DBMS can have indexes (for bookkeeping) to allow it to jump to specific locations in the log. *read가 빠르다. / id를 통해서 (index) 시간을 조금 줄일 수 있다.*
- The log will eventually get quite big. The DBMS can periodically compact the log by taking only the most recent change for each tuple across several pages. *뒤에서부터 압축하면 최신의 것만 남을 수 있다. ⇒ Compaction*
- After compaction, the ordering is no longer needed since there's only one of each tuple, so the DBMS can sort by id for faster lookup. These are called **Sorted String Tables (SSTables)**. *tuple을 compact한 것 같은데.*
- In Universal Compaction, any log files can be compacted together. In Level Compaction, the smallest files are level 0. Level 0 files can be compacted to create a bigger level 1 file, level 1 files can be compacted to a level 2 file, etc. *수업 시간 X But Key를 기준으로 정렬*
- The downside is that **compaction is expensive** and **also leads to write amplification** (for each logical write, there could be potentially multiple physical writes).

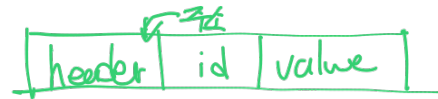
2 Index-Organized Storage

Observe that both (page-oriented storage) and (log-structured storage) rely on additional index to find individual tuples because the tables are **inherently** unsorted. In the **index-organized storage** scheme, the DBMS directly stores a table's tuples as the value of an index data structure. The DBMS would use a page layout that looks like a slotted page, and tuples are typically sorted in page based on key.

tuple은 정렬되지 않음 / index는 정렬 index-Organized Storage



3 Data Representation



The data in a tuple is essentially just byte arrays and doesn't keep track of what kinds of values the attributes are. It is up to the DBMS to know how to keep track of that and interpret those bytes. A **data representation scheme** is how a DBMS stores the bytes for a value.

DBMSs want to make sure the tuples are **word-aligned** so that the CPU to access it without any unexpected behavior or additional work. Two approaches are usually taken: *→ CPU가 읽을 수 있는 단위(워드)를*

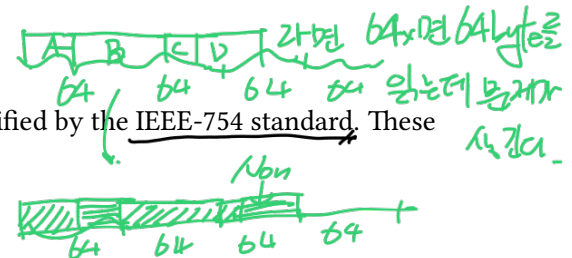
- **padding:** add empty bits after attributes to ensure that tuple is word aligned.
- **reordering:** switch the order of attributes in the physical layout to make sure they are aligned. *→ 다시 정렬해서 64비트 맞추기만 하지 않을 거 같은데(?)*

There are five high level datatypes that can be stored in tuples: integers, variable-precision numbers, fixed-point precision numbers, variable length values, and dates/times.

Integers

Most DBMSs store integers using their “native” C/C++ types as specified by the IEEE-754 standard. These values are fixed length.

Examples: INTEGER, BIGINT, SMALLINT, TINYINT.



Variable Precision Numbers *→ Level of detail → floating number.*

These are inexact, variable-precision numeric types that use the “native” C/C++ types specified by IEEE-754 standard. These values are also fixed length.

Operations on variable-precision numbers are faster to compute than arbitrary precision numbers because the CPU can execute instructions on them directly. However, there may be rounding errors when performing computations due to the fact that some numbers cannot be represented precisely.

Examples: FLOAT, REAL.

Fixed-Point Precision Numbers

These are numeric data types with arbitrary precision and scale. They are typically stored in exact, variable-length binary representation (almost like a string) with additional meta-data that will tell the system things like the length of the data and where the decimal should be.

These data types are used when rounding errors are unacceptable, but the DBMS pays a performance penalty to get this accuracy.

Examples: NUMERIC, DECIMAL.

Variable-Length Data

These represent data types of arbitrary length. They are typically stored with a header that keeps track of the length of the string to make it easy to jump to the next value. It may also contain a checksum for the data.

Most DBMSs do not allow a tuple to exceed the size of a single page. The ones that do store the data on a special “overflow” page and have the tuple contain a reference to that page. These overflow pages can contain pointers to additional overflow pages until all the data can be stored. *→ 커번 데이터의 경우 허용*

Some systems will let you store these large values in an external file, and then the tuple will contain a pointer to that file. For example, if the database is storing photo information, the DBMS can store the

photos in the external files rather than having them take up large amounts of space in the DBMS. One downside of this is that the DBMS cannot manipulate the contents of this file. Thus, there are no durability or transaction protections.

Examples: VARCHAR, VARBINARY, TEXT, BLOB → Object File

Dates and Times

Representations for date/time vary for different systems. Typically, these are represented as some unit time (micro/milli)seconds since the unix epoch.

Examples: TIME, DATE, TIMESTAMP.

Null Data Types

There are three common approaches to represent nulls in a DBMS.

- **Null Column Bitmap Header:** Store a bitmap in a centralized header that specifies what attributes are null. This is the most common approach.
- **Special Values:** Designate a value to represent NULL for a data type (e.g., INT32_MIN).
- **Per Attribute Null Flag:** Store a flag that marks that a value is null. This approach is not recommended because it is not memory-efficient. For each value, the DBMS has to use more than just a single bit to avoid messing up with word alignment. → 하나의 bytes를 할당하는건 그냥 낭비다.

4 System Catalogs

→ 디코딩 (decode something → 해석하기)

In order for the DBMS to be able to decipher the contents of tuples, it maintains an internal catalog to tell it meta-data about the databases.

Metadata Contents:

- The tables and columns the database has as well as any indexes on those tables.
- Users of the database and what permissions they have.
- Statistics about the table and what contents are contained within them (i.e., max value of an attribute).

Most DBMSs store their catalog inside of themselves in the format that they use for their tables. They use special code to “bootstrap” these catalog tables.

