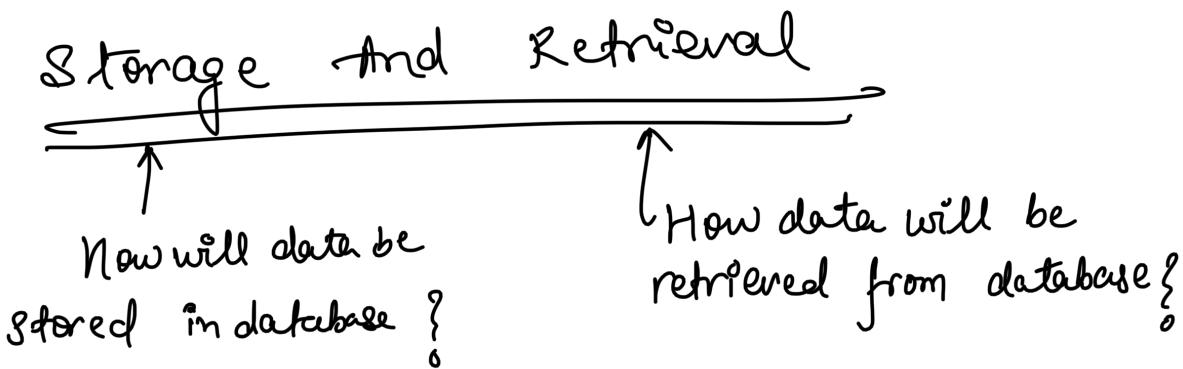


# Designing Data Intensive Applications

## CHAPTER - 3



### Date Structures that Power your Database

Why do I care what's happening in backend?  $\Rightarrow$  It helps in selecting a data storage for your system. There are lot of options these days 😊

Log File  $\rightarrow$  Assume you're keeping a key value pair and store in text file in append mode.

For any new write query, you'll just append an entry.

How will read value for given key?

$\hookrightarrow$  Scan entire file for occurrence of key?  $\rightarrow$  that will be  $O(n)$  operation

$\hookrightarrow$  It's not efficient if data grows to large scale

$\hookrightarrow$  You'll need indexes (watch youtube video on Indexes on channel)  
create indexes for data written.

$\hookrightarrow$  Should I create indexes for all the data? It will make read faster

$\rightarrow \underline{\text{NO}}$ , it will make write slower. Along with append operation, you also need to write data to indexes.

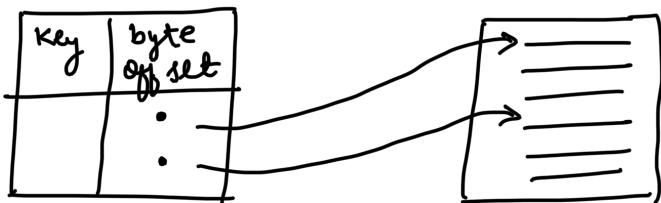
Hash Indexes  $\rightarrow$  How will you index data on disk?

Assumption - Data storage provides only append operation on file

Indexing Strategy  $\rightarrow$

$\hookrightarrow$  Create in-memory hashMap.

$\hookrightarrow$  Every key is mapped to byte offset in file.



$\hookrightarrow$  gt is used by Bitcask (Storage engine in Riak)

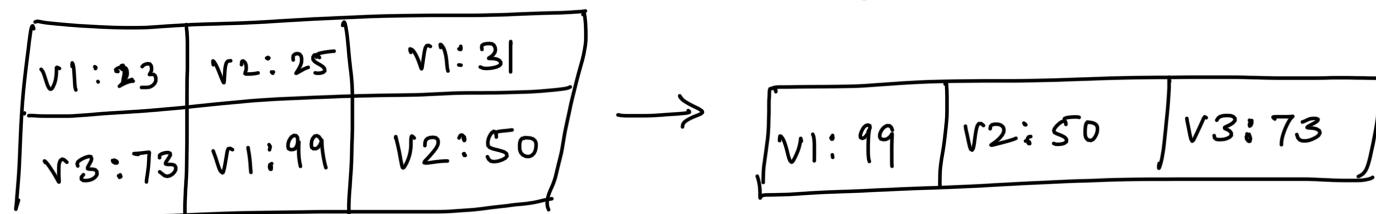
$\hookrightarrow$  In-memory Hashmap helps if there are limited no of keys.

$\hookrightarrow$  Example - maintain a counter how many times video is played.

$\hookrightarrow$  How will you ensure you're not running out of disk space (since you're only appending)

$\hookrightarrow$  ① break the log into segments of certain size.

$\hookrightarrow$  ② perform compaction - remove duplicate keys in log and only keep recent update.



$\hookrightarrow$  Each segment has its own in-memory hash Table.

$\hookrightarrow$  Find value for key : check most recent segment's hash map.

$\hookrightarrow$  If not found, check 2nd most recent.

## Implementation Details

- ① file format - Binary → encodes the length of string in bytes, followed by raw string.
- ② Record Deletion - append special deletion record to data file.  
↳ called Tombstone {watch Discard msg storage video}
- ③ Crash Recovery - Maps in memory will be lost in events of crash.  
↳ keep snapshot of segment's hash map on disk.
- ④ Partially written records - use checksums, helps in identifying corrupted parts and delete them.
- ⑤ Concurrency - Data files are append only → immutable.  
↳ can be read concurrently.

Why not update the key instead of appending?

- ① much faster than random writes, especially on magnetic spinning disks.
- ② concurrency and crash recovery is simple.

## Limitations

- ① Hash Table must fit in memory.
- ② Range queries are not efficient.

## SSTables and LSM-Trees

optimizing the limitations mentioned above ↑

↳ change the format of how segment files are stored

↳ make sequence of key-value pair "sorted by key"  
↑

This is called Sorted String Table (SSTable)

↳ Each key should appear only once in merged segment.  
(handled in compaction process)

## A) Merging of segments → algorithm similar to merge sort.

① Read input file side by side.

② look for first key in each file.

③ Copy lowest key to output file repeat again.

④ If same key in multiple segments → pick value from most recent segment.

S1	V1: 5   V2: 3   V3: 3   V4: 7
----	-------------------------------

S2	V5: 9   V6: 11   V2: 5   V3: 9
----	--------------------------------

S3	V1: 11   V2: 9   V7: 23   V3: 21
----	----------------------------------

↓	V1: 11   V2: 9   V3: 21   V4: 7   V5: 9   V6: 11   V7: 23
---	---

## Search for a Key

- ↳ ① Inmemory index is not required for maintaining offset for key.
- ② Maintain offset for some of keys (sparse index) and they traverse between these offsets to find particular key.
- ③ Group the records for keys you're maintaining in Inmemory index.  
↳ saves disk space and reduce I/O.

## Construction & maintenance of SSTables

for any write query -

- ① Add it to in-memory balanced tree structure. — called memtable.
- ② Write memtable to disk after certain threshold (say 1MB) — called SSTable
- ③ For read query, Search in memtable, then most recent segment of SST.
- ④ Run merge & compaction time to time.

## Performance of Database built on LSM Tree

- ① Reads can be slow as you traverse from one segment to another.  
 ↳ Bloom filters are used → data structure helpful in approximating contents of set.  
 Example → will tell if key is present in DB or not saving traversing DB.
- ② Compaction & merging → how to determine order & timing.

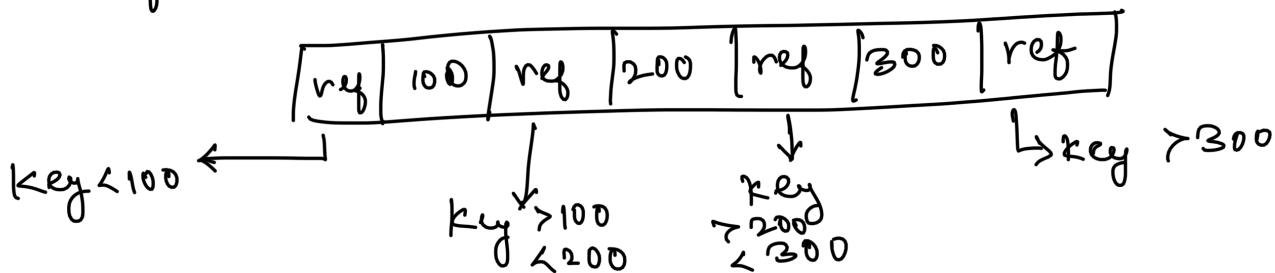
Size-Tiered	Leveled Compaction
HBase. Cassandra uses both. smaller SSTables are successively merged into larger SSTables.	LevelDB, RocksDB key range is split up into small SSTable and older data moved to separate levels.

B Tree → it is not "Binary" Tree.

↳ key value pairs sorted by key. → on disk

↳ BTree break database into fixed size blocks.

↳ Each Block can be identified by address, using which one page refer to another.



Handling Crashes in Btree → what if hardware crashed at time of any update?

It maintains additional data structure → WAL (Write ahead Log)

↓  
 append only file - every modification to Btree is written here before actual update.

Write amplification → one write to DB resulting in multiple writes on disk over course of DB's lifetime.

## LSM Tree

- ① Typically faster for writes
- ② It can sustain higher write throughput due to less write amplification
- ③ Can be compressed better
- ④ Compaction process can affect ongoing read/write.
  - ↳ more observable at high throughput.

## BTree

- ① Typically faster for reads.
- ② More write overhead, update existing page.
- ③ Leaves some space unused due to fragmentation.
- ④ Each key is present a single place.
  - ↳ offer strong transactional support

## Full Text Search and fuzzy indexes

↳ Lucene example Elasticsearch

Ex. to search for mis-spelled words or synonyms.

## In-memory Databases

memcached → intended for caching purpose, it's ok to lose data if machine is restarted.

↳ Inmemory databases for durability → dump periodic snapshots to disk.

↳ Redis and couchbase provide weak durability by writing to disk asynchronously.

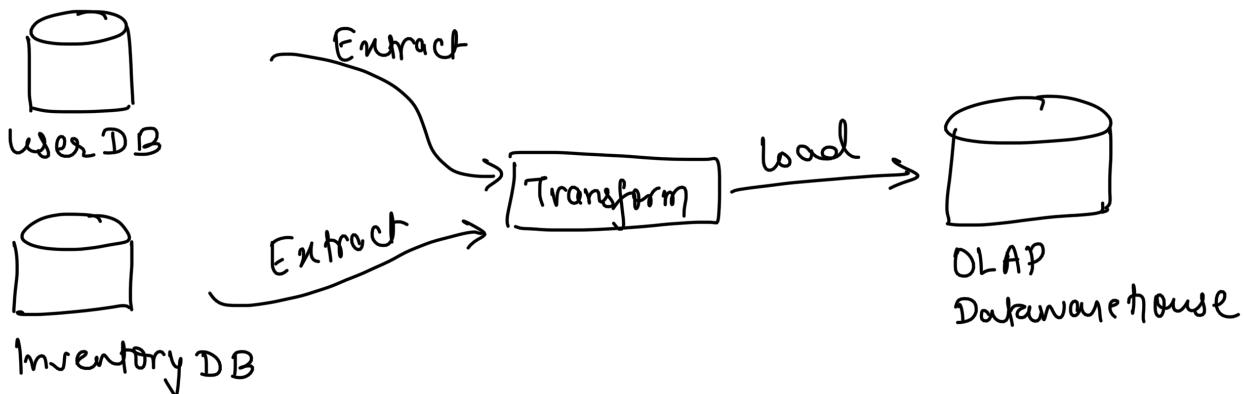
## Your use-case → Transaction Processing or Analytics ?

### OLTP

- ① Read - small number of records per query fetched by key
- ② Write - Random access, low latency writes
- ③ GB to TB

### OLAP

- ① Aggregate over large set of data.
- ② ETL
- ③ TB to PB



## OLTP databases

### Column Oriented Data Storage

- ① The table can be  $\geq 100$  columns wide but we might be accessing only 4 to 5 columns at a time.
- ② To get the data, you need to load all rows matching criteria to these columns.
- ③ Column Oriented  $\rightarrow$  Store values from each column together instead of rows.
- ④ Column compression is easier  $\rightarrow$  there are chances that values in column are same for different rows.
- ⑤ Example - Cassandra, HBase
  - $\hookrightarrow$  not strongly column oriented.
  - $\rightarrow$  they store all columns for row together along with rowkey.
  - $\rightarrow$  don't use column compression.

Subscribe for more such content.

YouTube channel- Ms Deep Singh

Happy Learning 😊