

Agenda

- ✓ Problem Statement
- ✓ Write Ahead Log (WAL)
- ✓ LSM Tree } Overview / intuition
- ✓ Bloom Filter

NoSQL → handling automatically

→ managing servers

Today → how is data stored on disk.

Problem Statement

↳ SQL → well defined schema

I know the exact size of each row

col → datatype

User

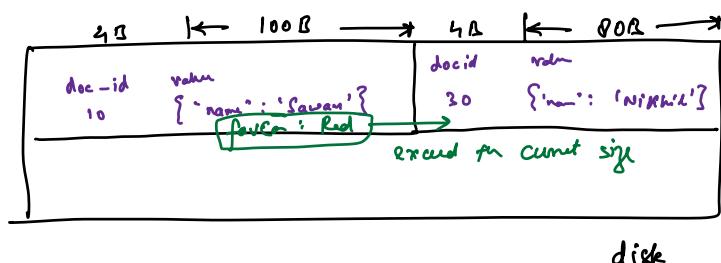
id	name
8 byte	varchar(50)

20 30	Rohan
5	Nikhil

Nikhil Agrawal

Size of row remains same in changes

→ NoSQL → different



db.update(10, { 'value': 'Sawan', 'favColor': 'Red' })

① Can lead → overwriting

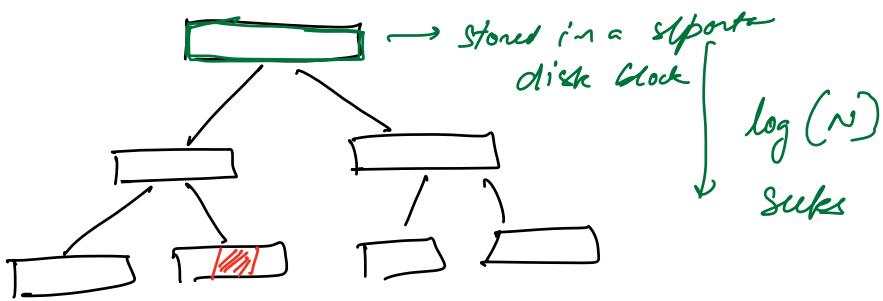
② Split the document across disk → fragmentation
a lot of disk seeks → very expensive



even for SSD

↳ sequential I/O
is 100x faster
than random I/O

SQL → Indexes → \mathbb{R}^+ Trees

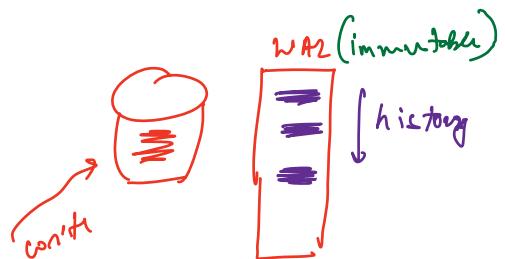


Goal No SQL → optimized for both read & write loads

for any database (Especially NoSQL)
in SQL you can turn it off

① Write Ahead Log (WAL)

- ↳ append only
- ↳ every commit → appended to WAL
- ↳ allows you to replay the entire history of DR.



- ↳ power failure recovery is done via WAL

- ↳ replication
- ↳ db Backups
- ↳ Point of Time Recovery

WAL can also be read from
 read the tail of the file
 look at the last 2mb / 200 entries
 of the file

② Current state of data
 discuss how to implement efficiently.

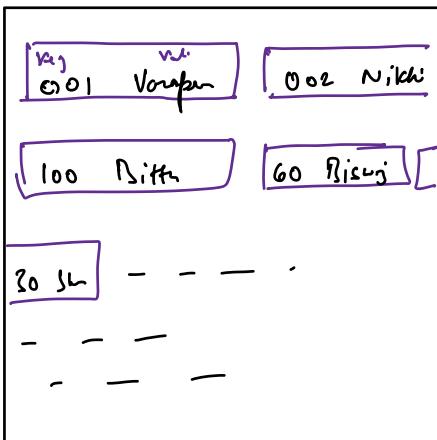
Key-Value Store → Design how the data is
 written on disk

if each entry was fixed size → R^+ true

↳ not the case for NoSQL dbs



Brute force



write values to a file

- ① Read key (60)

Linear Scan entire file

- ② Write key (60) → Bisujit

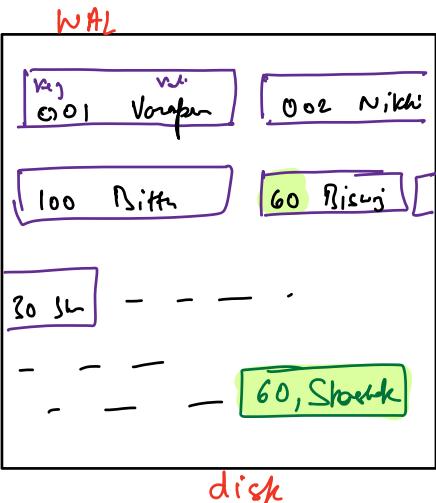
Linear Scan & update

key :

value can be
of varied size

Un WAL

every write → appended to the file



- ① read key (60)

Linear scan
from last to first
(take most recent write)

- ② write key (60)

↳ Snapshot

append a new
entry to file

③ (i) writes
extremely
fast

There will be duplicate entries!

writing → super fast

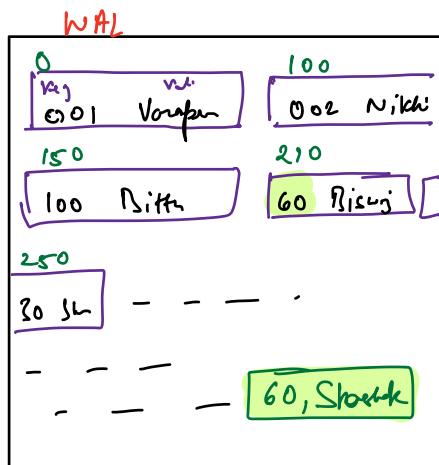
read → still super slow.

Option \rightarrow B-tree \rightarrow sorted slate \rightarrow can't scan
 ↳ varied sized entries

HashMap \rightarrow given a key its location in file
 (memory)

HashMap	
key	offset
001	0
002	100
100	150
60	210
30	250

in-memory



def read(key):

0.1μs {
 ↓ 10,000x }
 10ms }
 offset = Hmap.find(key)
 buffer = read the next n bytes from file
 Start from 'offset'
 $O(1)$
 key, value = parse(buffer)
 return value

```

def write(key, value):
    offset = WAL.size()
    append(filem, [key, value])
    Hmap.set(key, offset)

```

Still duplicates in WAL

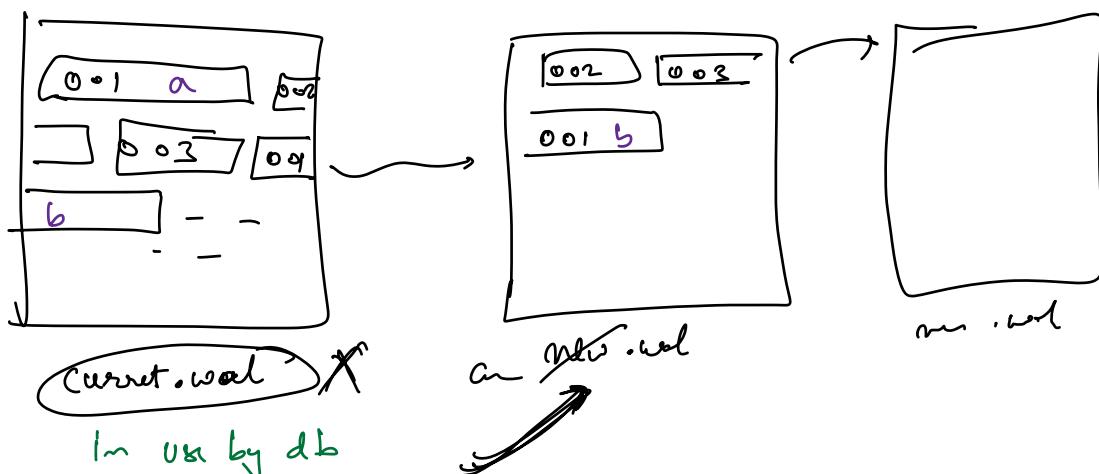
The reason we're still appending in WAL for writes is because data has varint size.

If I try to modify existing value → overflow

① Still duplicate entries → try to remove duplicates periodically

② In a DB, you might have billion / trillion keys
Key - offset Hashmap → might not fit in RAM

We have a cron job / background process



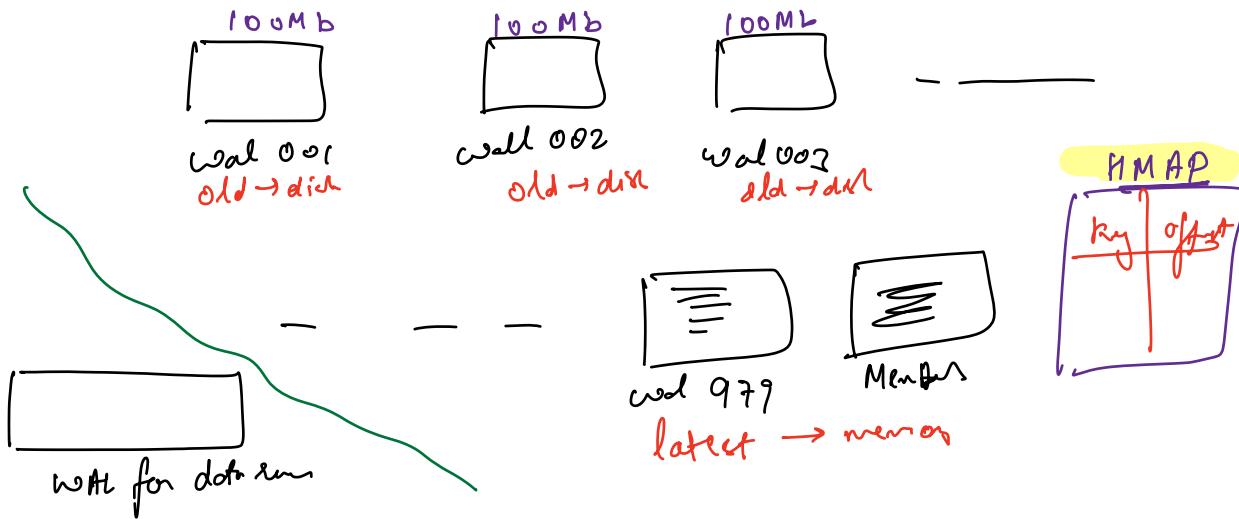
WAL are pretty large

↳ ① can we find duplicates efficiently?

↳ ② we can't read entire WAL in one go

Can we process the WAL in chunks? ✓

Can we also just store WAL in chunks?



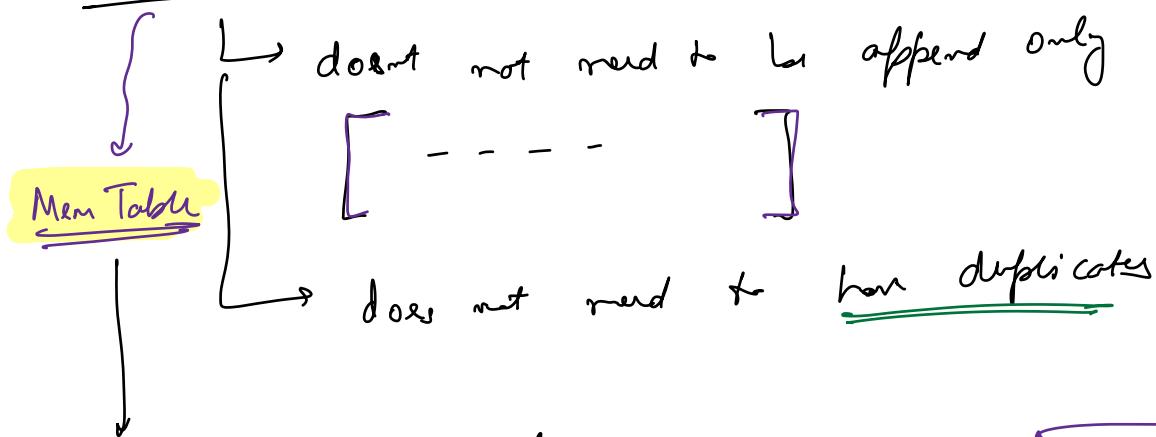
- ① Any write that happens → latest WAL file
 - ② each chunk of WAL is small → stored latest (current) chunk in memory instead of Disk??
- writes are now even faster!!

to ensure that data is not lost in case of power failure

Any write →

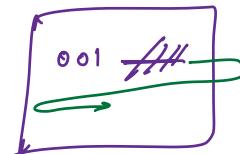
- ↳ changed in memory
- ↳ appended to separate WAL on disk → very small ~100 MB

latest WAL → in-memory



assume: this is a hashing

this MemTable will also have
latest data



↑ count 001

any read for which key found in WAL

↳ read can get from memory

↳ doesn't have to touch disk!!

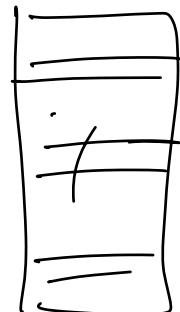
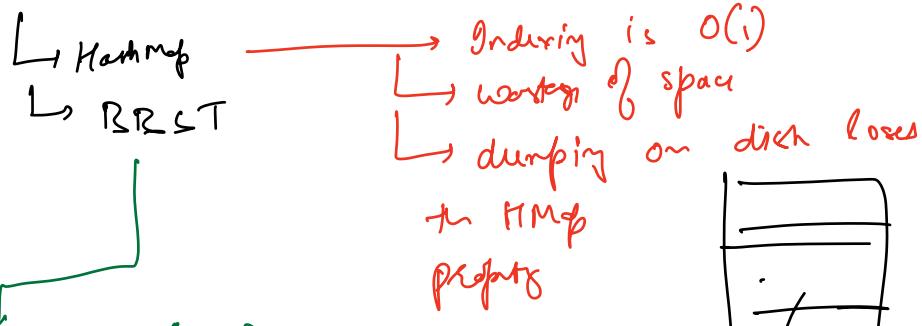
After some counts the MemTable will get large
 $> 100 \text{ Mb}$ (threshold)

↳ dump it on disk as a WAL file

- Any wAI chunk will not have any duplicates
- Across different wAI chunks → duplicates

So we still need a background process
to compact data

① Memtable → in memory → random order



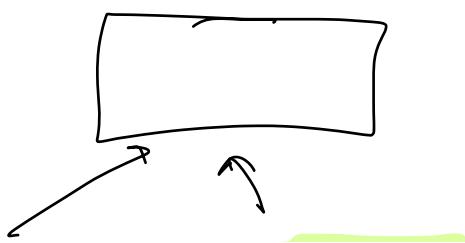
→ Operation is $O(\lg n)$

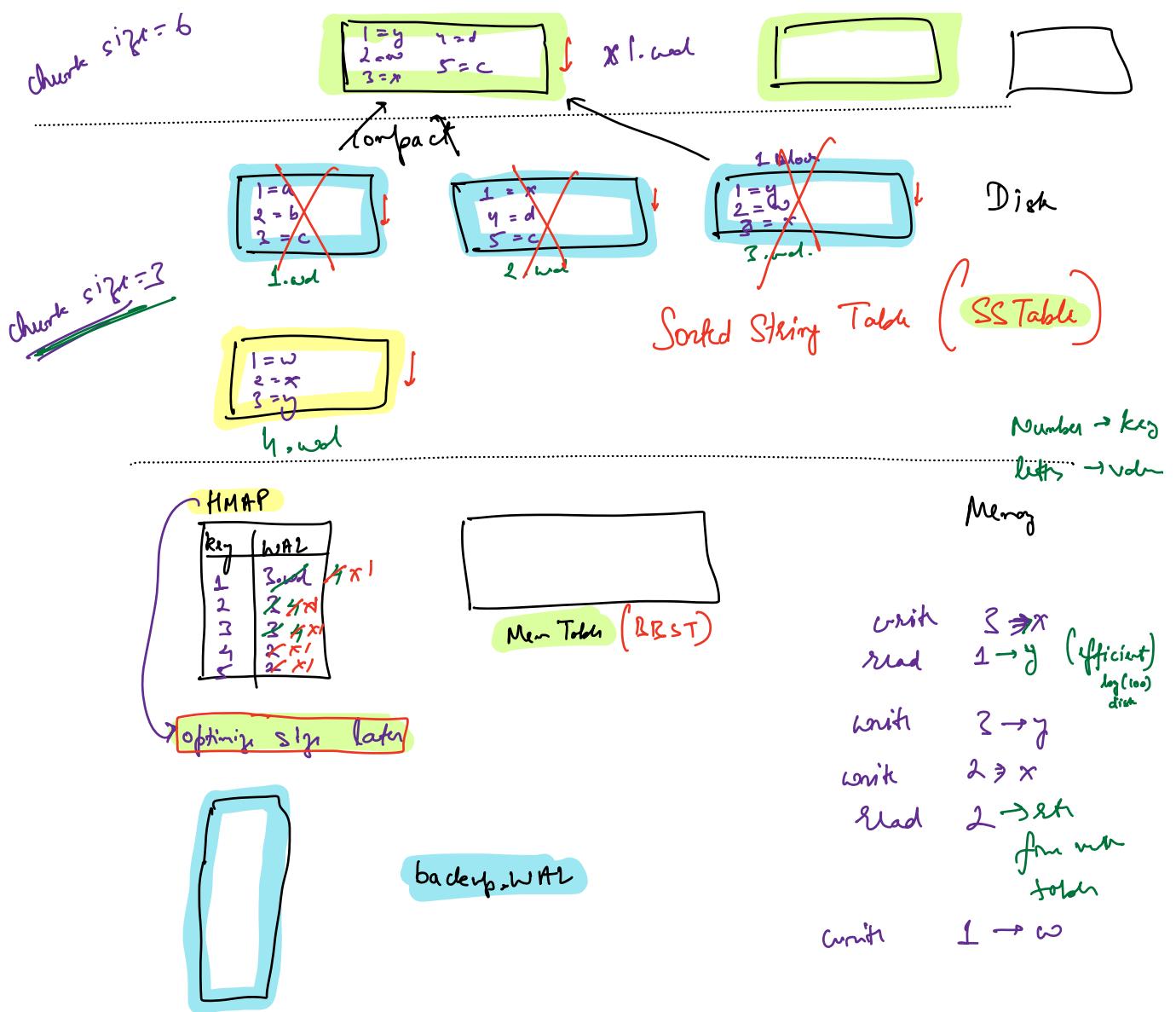
→ no space waste

→ dump it on disk → sorted order

even wAI chunk will be automatically sorted!!

→ allows us to avoid the key → (val, offset)
hashmap partially





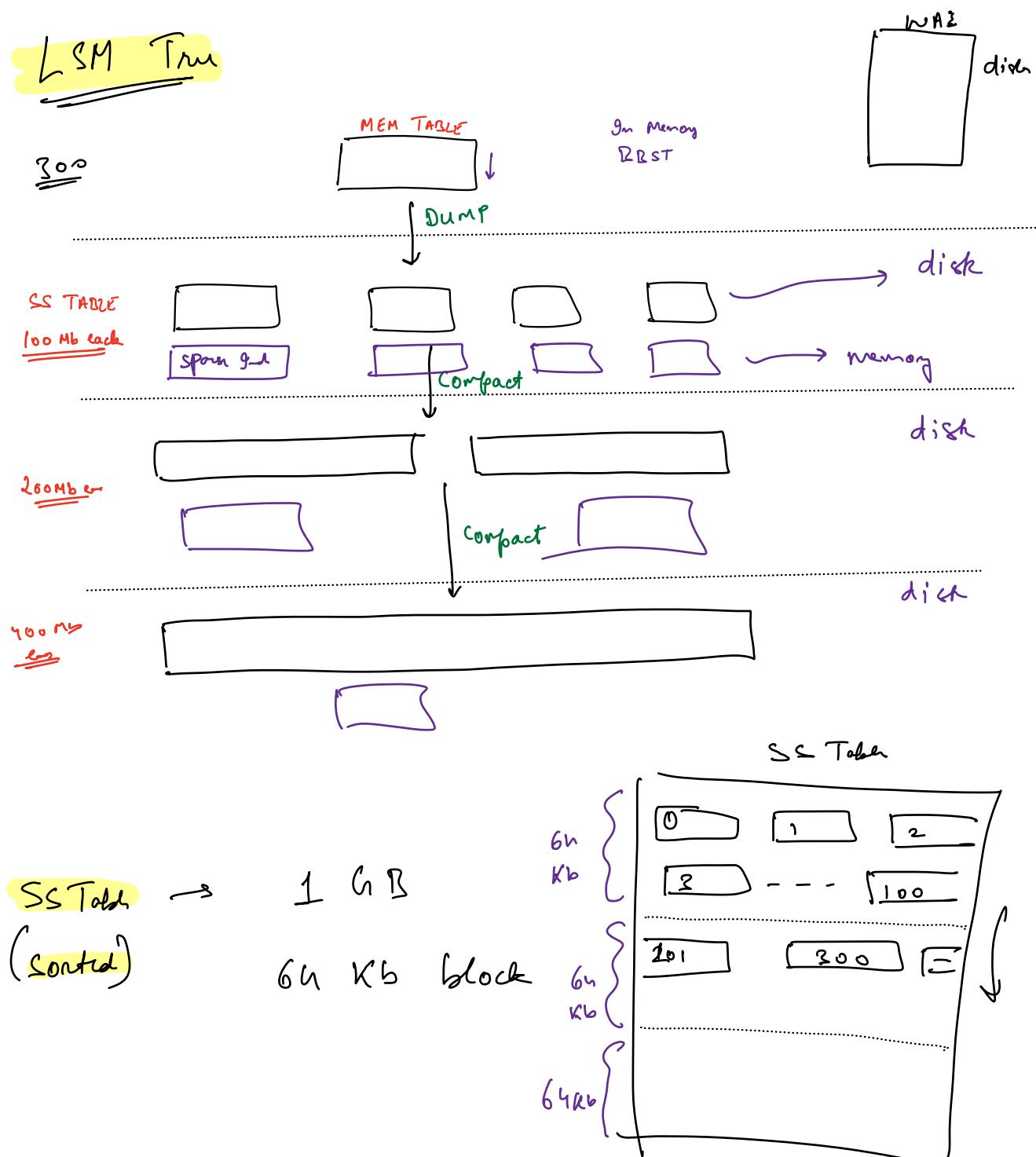
① Power failure → evict from backups_WATL

② Mem table is full

③ Backgr compaction process

LSM Tree

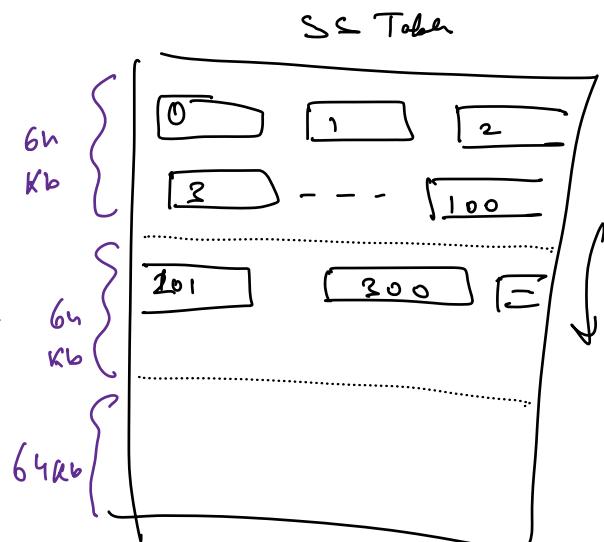
LSM Tree



SS Table \rightarrow 1 GB

(sorted)

64 Kb block

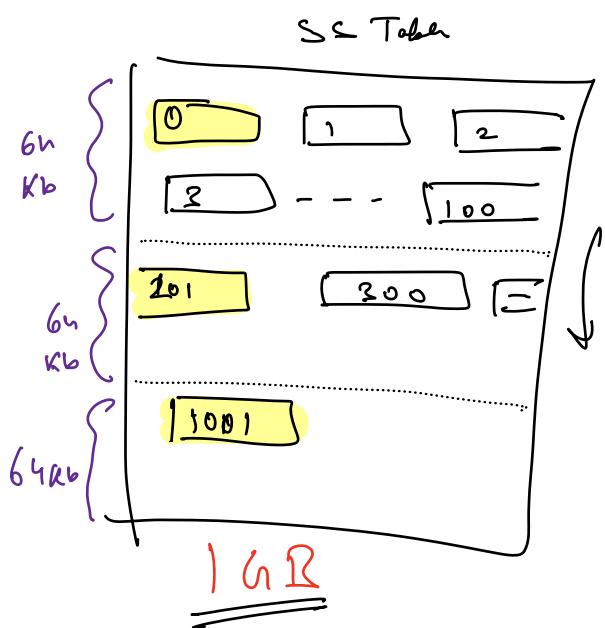


how many chunks?

$$\frac{1 \times 2^{30}}{64 \times 2^{10}} = 2^{12} \text{ chunks}$$

↳ 16 K chunks

Sparse



~~Index~~

0, offset = 0
201, offset = 2000
1001, offset = 5000

1 MB 16 K entries
 → no values (only keys)

$$\begin{aligned} &\rightarrow \text{on avg each key} \\ &\approx 64 \text{ bytes} \\ &16 \times 2^{10} \times 64 = 2^{20} \text{ B} \\ &\approx 2^6 \qquad \qquad \qquad = 1 \text{ MB} \end{aligned}$$

I problem if key is not there → not slow
 → but still we do it
 or burst of work
 to return nothing !!

Bloom Filter → probabilistic data structure

insert(key) → does Not increase size at all!

contains(key) → true/false

if key is present → guaranteed to return true

is absent → return false with

high probability
true false ←
probability. (false positive)

insert → Nikhil → 2, 3, 9
contains(Nikhil) → if all one set
H1 H2 H3 → 2, 3, 9

insert → Visha → 2, 5, 11

contains → Visha → 1, 3, 7

contains → Akhil → 3, 9, 11 → False true

Bloom Filter

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

M = 12
(size of bloom filter)

A = 3

Hash functions (k)

H1(s) → INDEX(S0) % 12

H2(s) → ≠ (S1) % 12

⋮

$K \rightarrow$ Hash func

$N \rightarrow$ # of values inserted

$M \rightarrow$ # bits in Bloom filter

$$\text{Probability of false pos} \approx \left(1 - e^{-kn/m}\right)^k$$

$$\underline{M = 10 \text{ MB}} = 10 \times 2^{20}$$

$$N = 1 \text{ billion keys} = 10^9$$

$$k = 10$$

NoSQL \rightarrow Bloom filter to store which keys are present

Akamai (Cloudflare \rightarrow CDNs)

(70% of URLs are only visited once)

Cache a URL only if it has been visited 2 times

Trillion URLs \rightarrow each URL path by $\sim 500B$

500 TB of data \rightarrow store

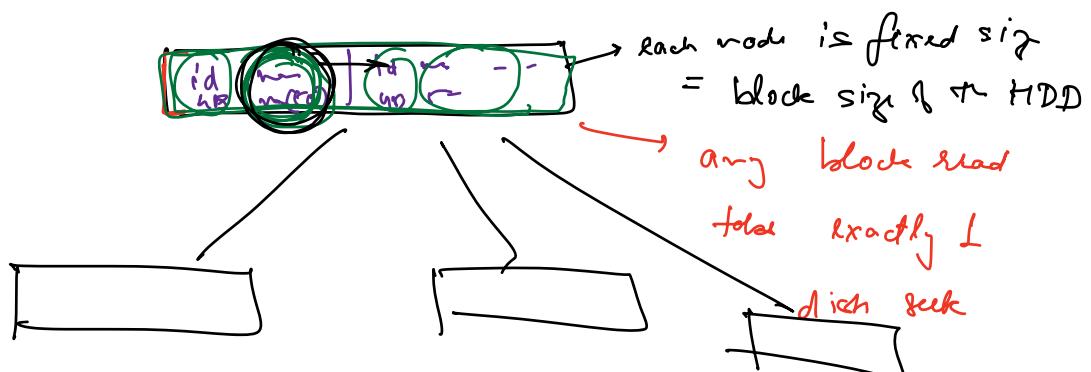
Bloom Filter

Delete (key) $\xrightarrow{\text{LSM}}$ write (key, Tombstone)
 ↓
 mark for a
 value special
 significance

Read (key) $\xrightarrow{\text{if val = Tombst}}$
 get \rightarrow not found
 or
 get value

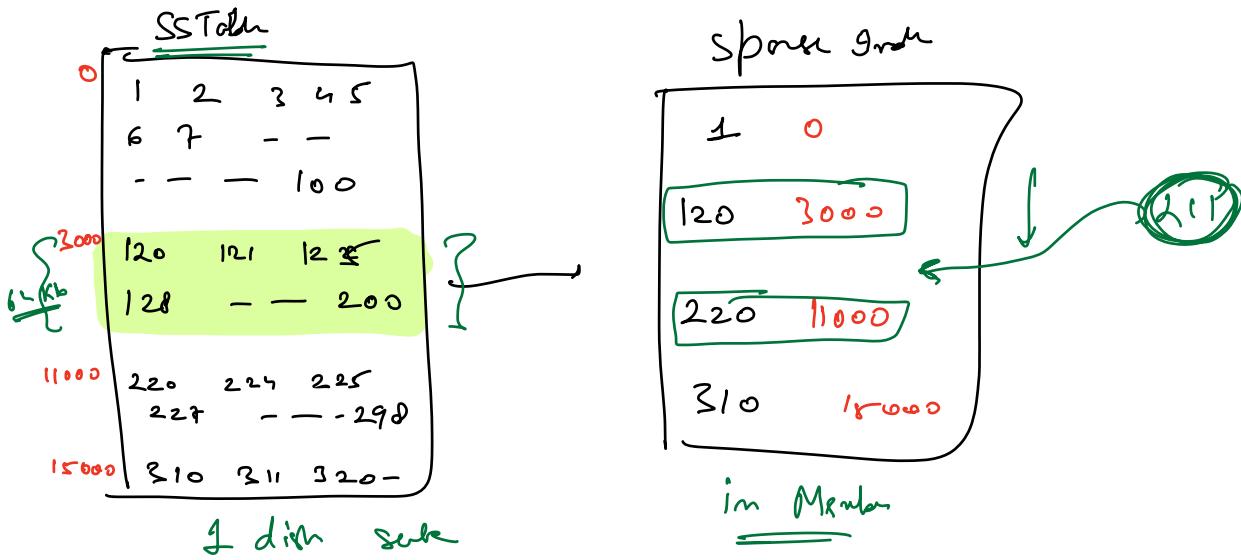
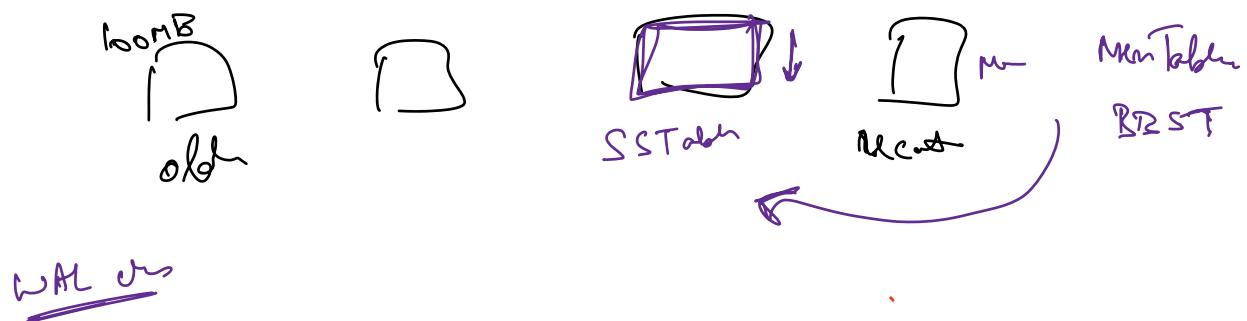
T1 k: Vanquishes LSN Then Deleted Value $\xrightarrow{\text{?}}$

T2 k: ~~OK~~





- ① duplisch
② $\text{Hn} \rightarrow \text{PV.bn}$
- remove duplicats.



graph \rightarrow $R^T + \alpha$ update
 \curvearrowleft ω_{AI} file