

11/Dec/2023

NoSQL Internals

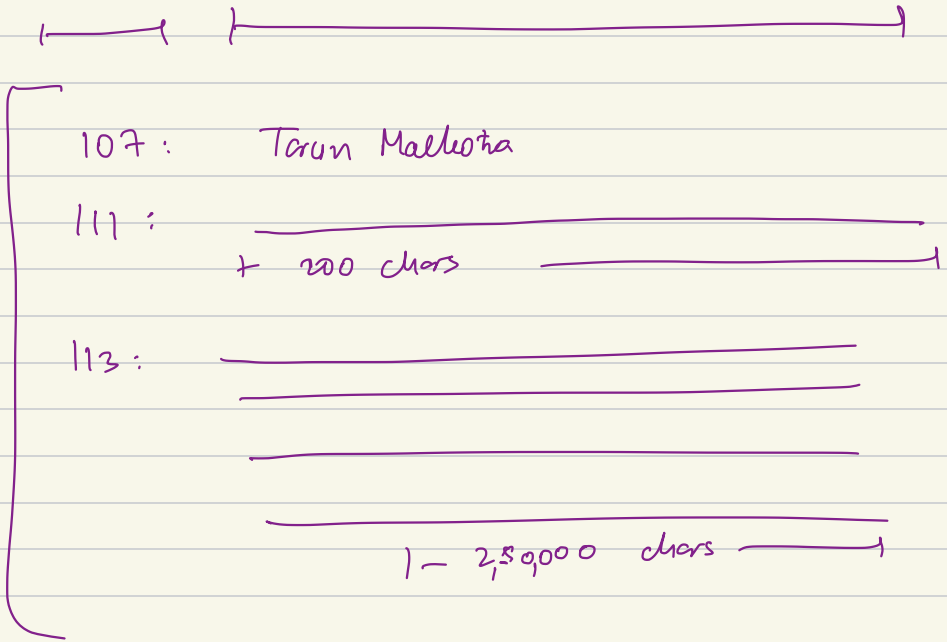


10th 

B+ Trees

[SQL → optimized for Reads + Writes
Because of fixed sized contiguous
rows 😊]

Key value



Update C (107)

107:

Torun Malluota

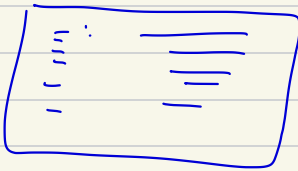
1 - 10,000 chars

① ~~WAL~~ (Write Ahead log)

AOF

WAL \equiv AOF
(Append only File)

② Snapshots of data



snapshot is saved
in the disk

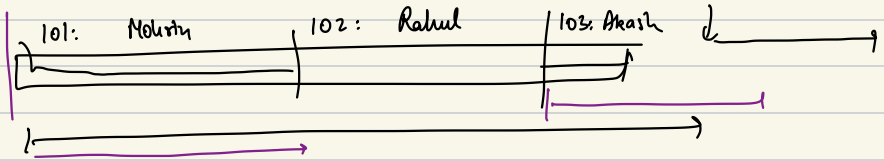
Iteration 1 :

Search : $O(n)$
Read : $O(n)$
Write : $O(1) + O(x)$

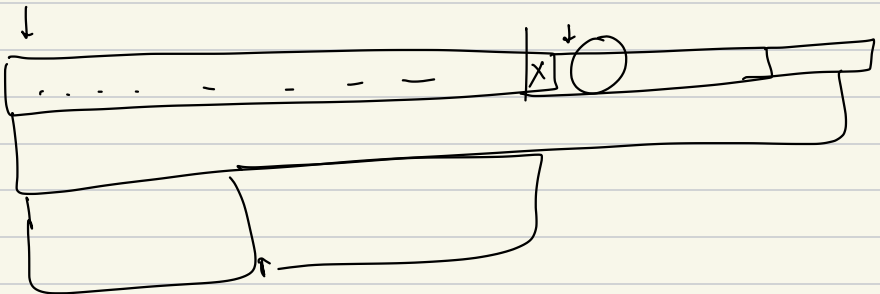
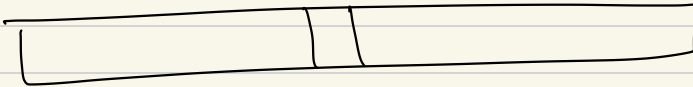
101 : Mohsin

102 : Rahul

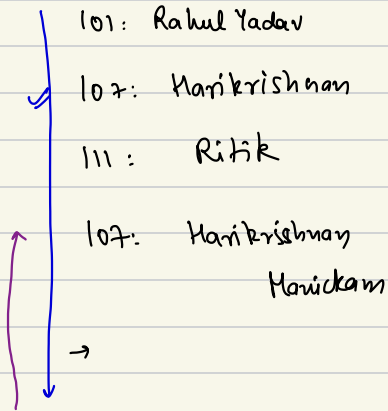
103 : Akash



Update : $O(N) + O(N) = O(N)$



Iteration 2:



T1: 101 → Rahul Yadav

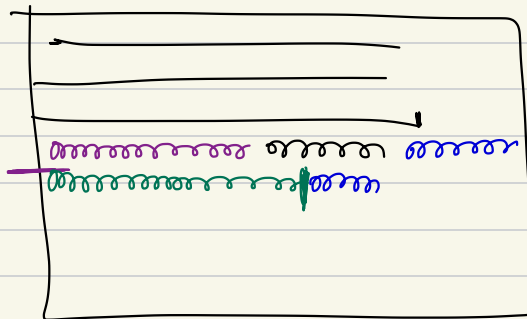
T2: 107 → Harikrishnan

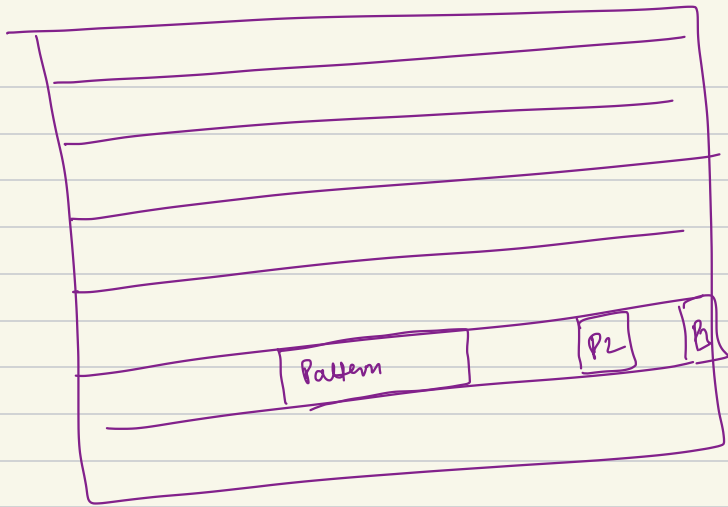
T3: 111 → Ritik

T4: update 107 →
Harikrishnan
Harickam

WRITE | UPDATE → 0(1) 😊

T5: Read (107)



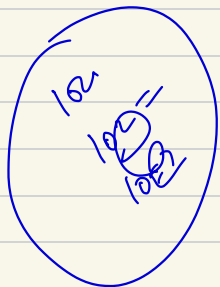


$N \rightarrow$ no of unique keys

$O(1)$ Write \rightarrow at the EOF

$O(1)$ update \rightarrow No search; simply write at EOF

$\Rightarrow O(1)$ Read \rightarrow start reading from end
and stop at the first
key match



↑ 101: Tarun
102:
103:
104:

110:
120:
110:
102:
102:
103:
1

delete 102

Writes are great

$O(1)$

Updates are great

$O(1)$

Reads are super bad

$O(N \times \text{avg duplicates})$

Storage space is super bad

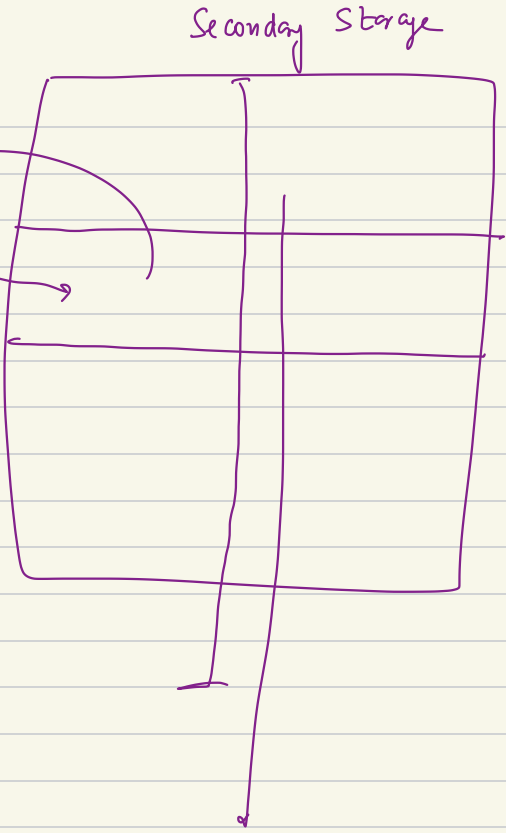
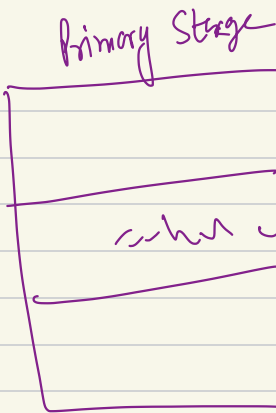
$O(N \times \text{avg duplicates})$

✓ 102: —

102: —

102: deleted

CPU

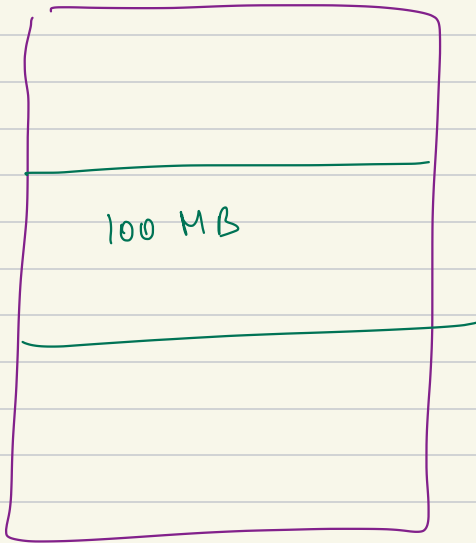


1GB
100GB

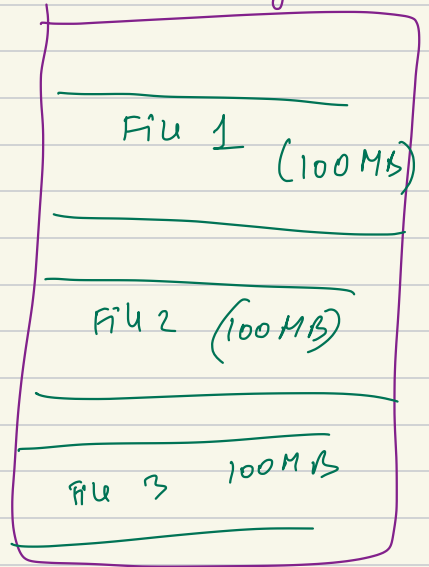
Two curved arrows pointing from the text '1GB' and '100GB' towards the right side of the page.

Given that your file which stores your
database will keep on becoming bigger and bigger
within no time it will become big enough that
it can't be loaded to main memory in 1
chunk 😊

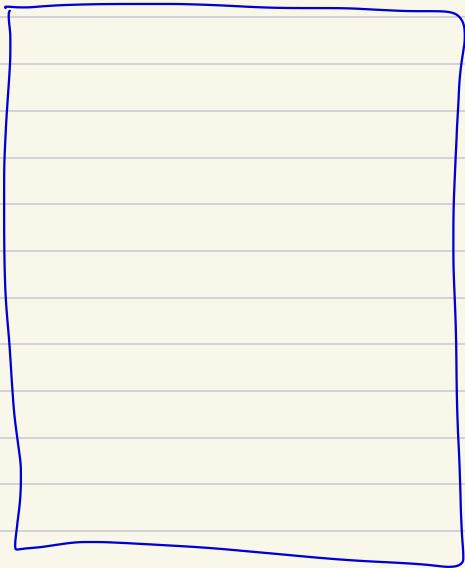
Main Memory



Secondary Memory



Main Memory



Secondary Memory

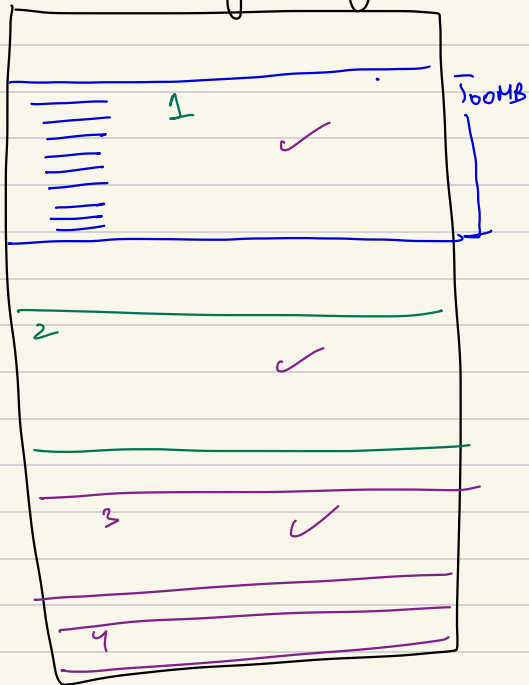


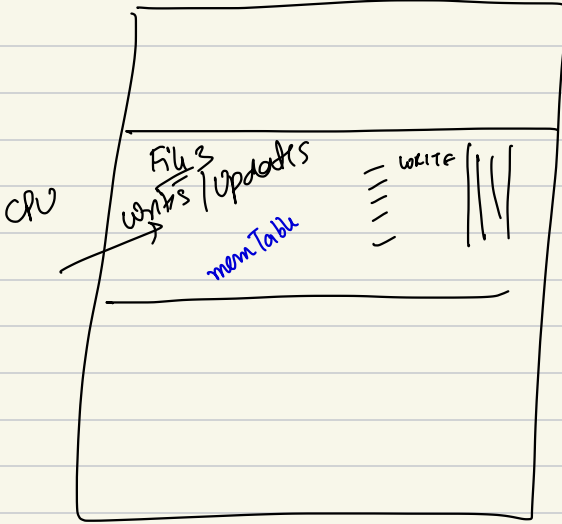
Fig 1] 100MB

Fig 2] 100M

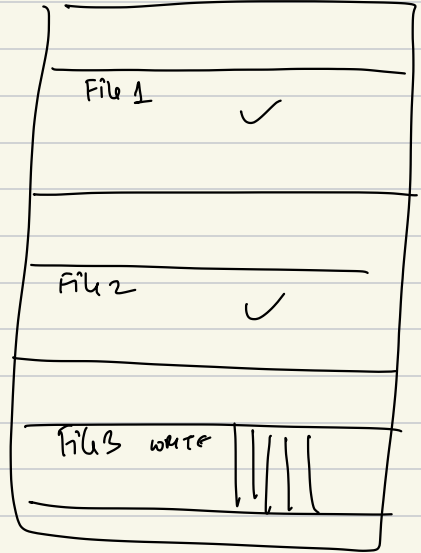
Fig 3]

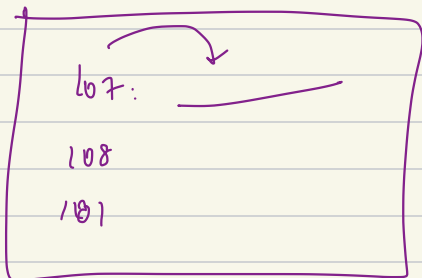
ITERATION 3

Primary



Secondary

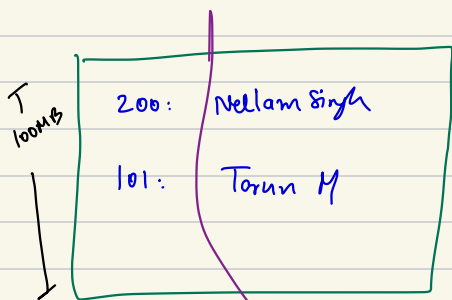




107: _____
 108: _____
 109: _____
 107: _____

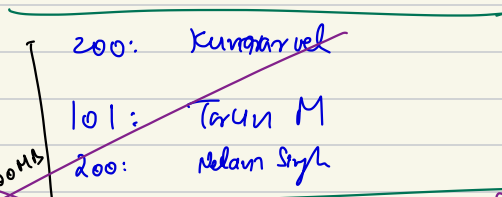
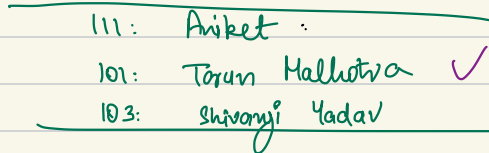
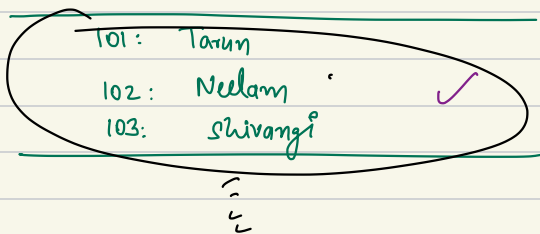
Iteration 4

Primary Memory



memTable

Secondary Memory



②

safely delete

①

deletes last

Primary



memTable



Why NOT have this

as a sorted map 😊 😊

Secondary

③

delete backup Bk

File 1 Backup

101: Bhagwan

111: Narendra

101: Bhagwan Singh

Serialize this memTable

and dump it

②

File 1

≡

101: Bhagwan Singh

111: Narendra

File 2 Bkp

Sorted String Table

Iteration 4:

SSTable

All the chunk files stored on the
secondary memory are

- ① without duplicates inside the same file
- ② They are sorted on keys 😊
- ③ All of them are ≈ 100 MB in size 😊

① Duplicates can still exist between
SSTables 😐 😐

② Only within a SSTable keys are sorted,
not b/w 2 different SSTables 😊
😐

Break.

10:40 PM - 10:48 PM

😊

MemTable

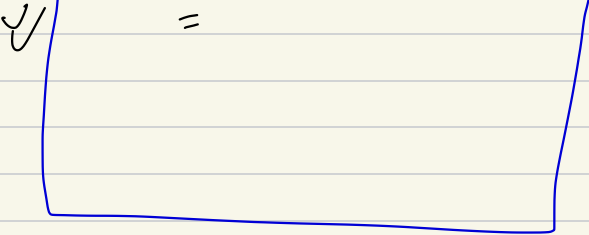
SSTables

How Reads will happen?

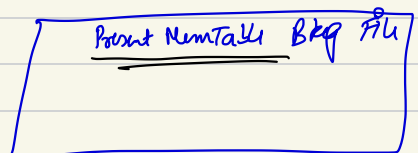
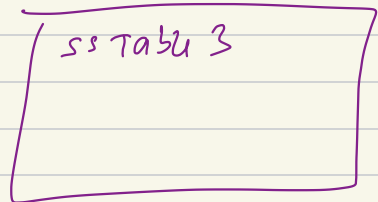
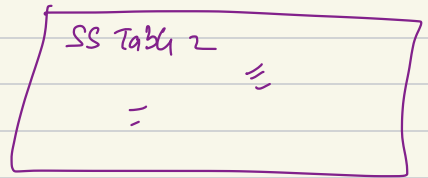
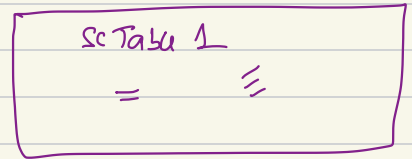
Read (107) \Rightarrow

Primary Memory Table
sorted hashmap

memTable \rightarrow



Secondary Memory



Read (key)

① Check if memTable contains the key

Yes 😊😊

No 😞

o(1) / ^{o(1)} you can reply
with the value

② Look into the latest SSTable [ex SSTable 3]
and see if key is present there.

Yes

No

You can return that
value

③ Penultimate SSTable (SSTable 2)

STEP 2 might have to be performed recursively.

① Updates are great

② Writes are great

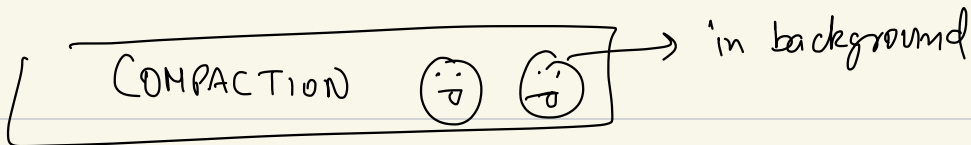
③ Reads - Best case scenario is great.

- Avg case scenario is GOOD

- Worst Case scenario is Very Bad 😞 😞

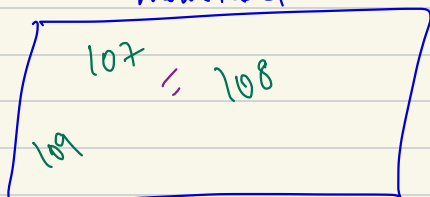


④ Space utilization is also Bad 😞 😞



Primary

memTable



Secondary

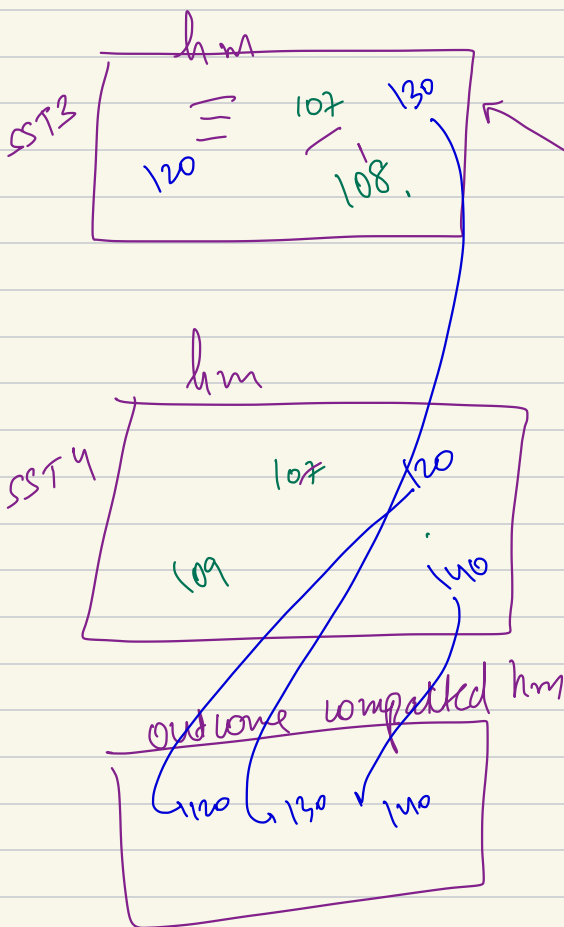
SST 1

SST 2

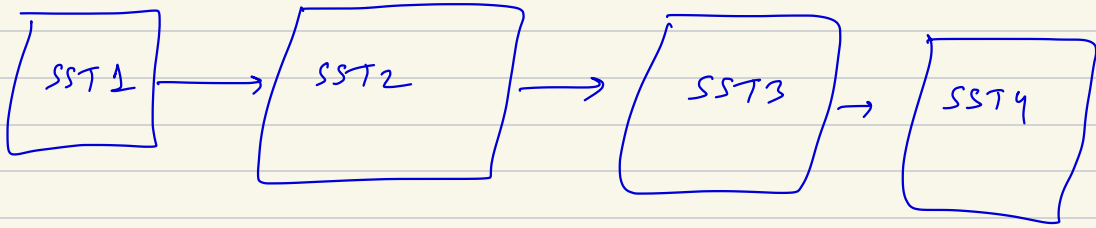
SST 3

SST 4

Step MemTable

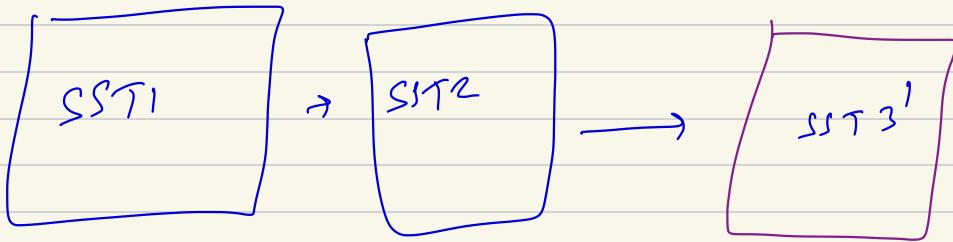


Imagine SSTables as a LL



Pick 3 + 4

Compact



ITERATIONS: includes compaction

① Updates are great ✓✓

② Writes are great ✓✓

③ Reads - Best case scenario is great. ✓✓

- Avg case scenario is GOOD ✓✓

- Worst case scenario is ~~Very Bad~~ ☹️☹️

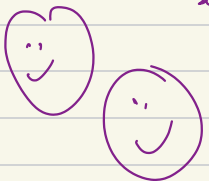
slightly less bad

④ space utilization is also ~~Bad~~ ☹️☹️

improved this a lot

because now we are removing duplicates all throughout

using the background process of compaction.



Indexing - Using Binary search to make Read way better....

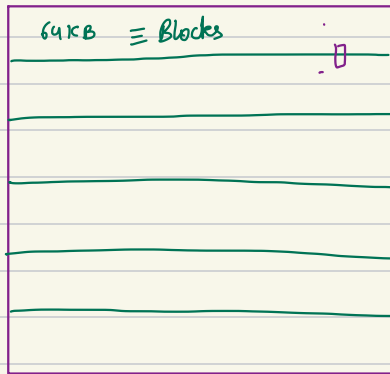
Why binary search can't be applied directly?

SSTable 2	
101:	Kumar vel
111:	Soumyakanti
200:	Aditya
400:	Shubham

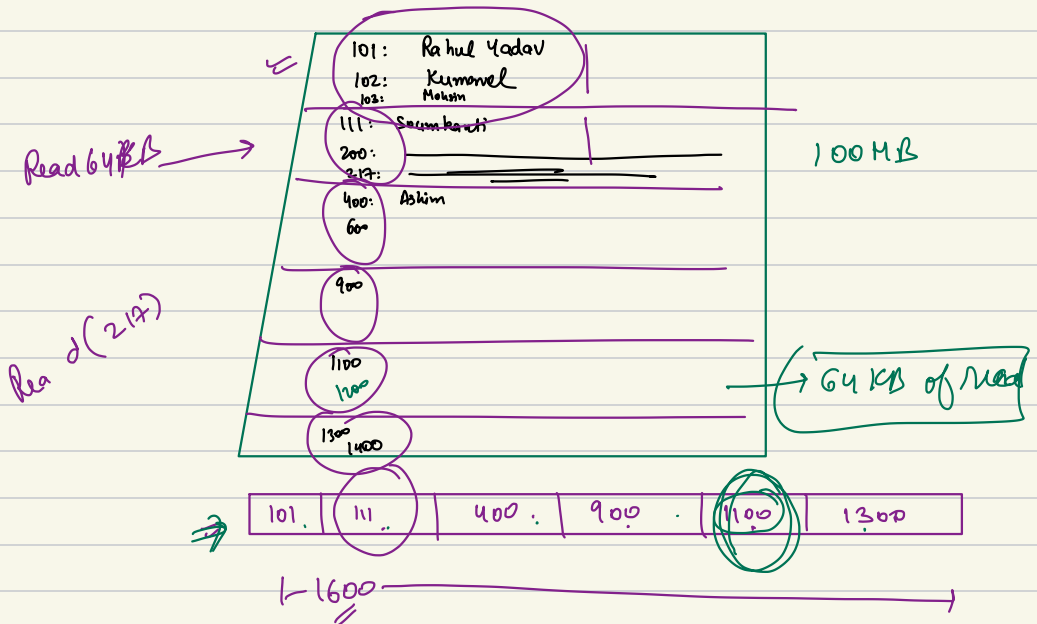


★ we can't apply B.S directly
as elements are NOT of equal size

T
1600
Blocks
inside
the SSTable



$$\# \text{ Blocks} = \frac{100 \text{ MB}}{64 \text{ KB}} = 1600 \text{ Blocks}$$



Read Algorithm

Read (700)

① Mem Table

$O(\log N)$

X
Repeated

SS Table last

$O(100MB)$

→ $O(64KB)$



1600 Times faster 😊

Block Size

✓ 64 KB

✓ 32 KB

✓ 128 KB

✓ 1 MB

1600 Times faster
 $O(100MB)$ → $O(64KB)$

3200x
 $O(100MB)$ → $O(32KB)$

800x
 $O(100MB)$ → $O(128KB)$

100x
 $O(100MB)$ → $O(1MB)$

Finally 😊 → 😊

① Updates are great

② Writes are great



③ Reads - Best case scenario is great. 😊

- Avg case scenario is GOOD 😊

- Worst case scenario is ~~Very Bad~~



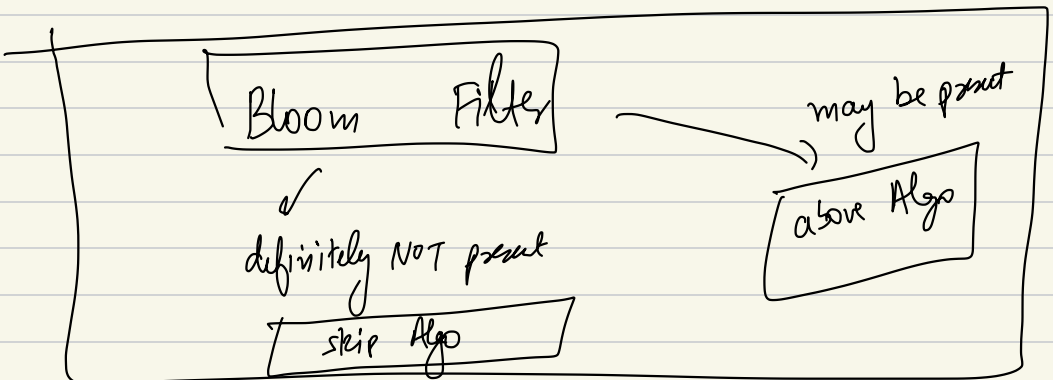
Good 😊



④ space utilization is also ~~Bad~~



Removed all duplicates due to comparison.



LSM Trees



↳ log Structured Merge Trees 😊