

NoSQL Internals

→ SSTables
→ LSM Trees

Structures that are used by most of NoSQL databases.

B Trees / BTrees

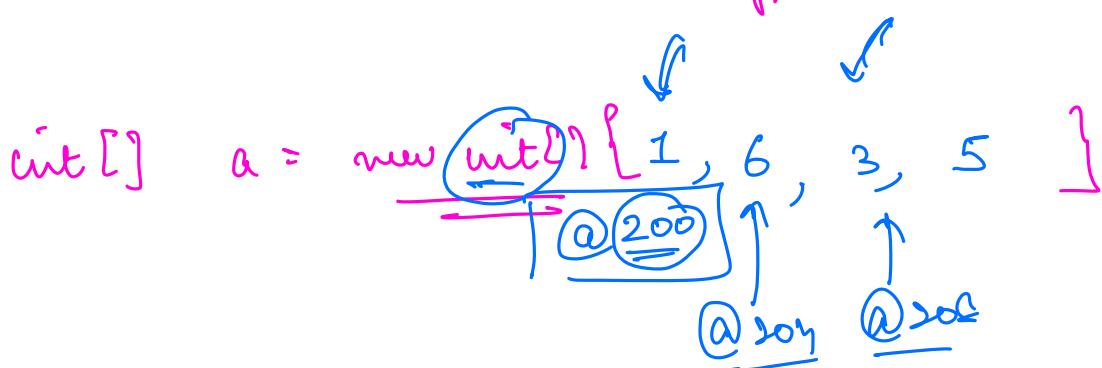
{ Optimized SQL Queries }

⇒ No class next Saturday

10th Feb

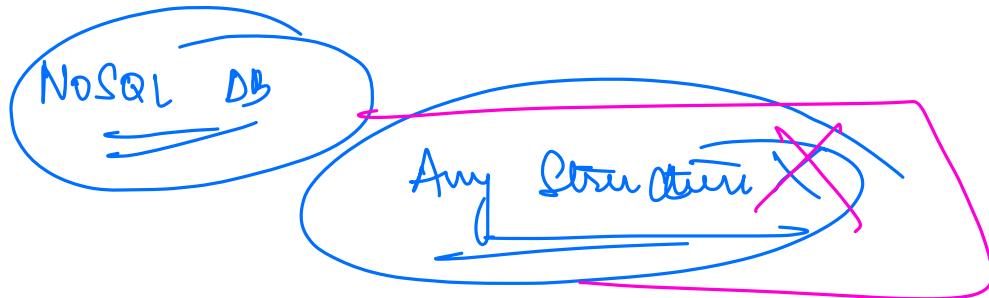
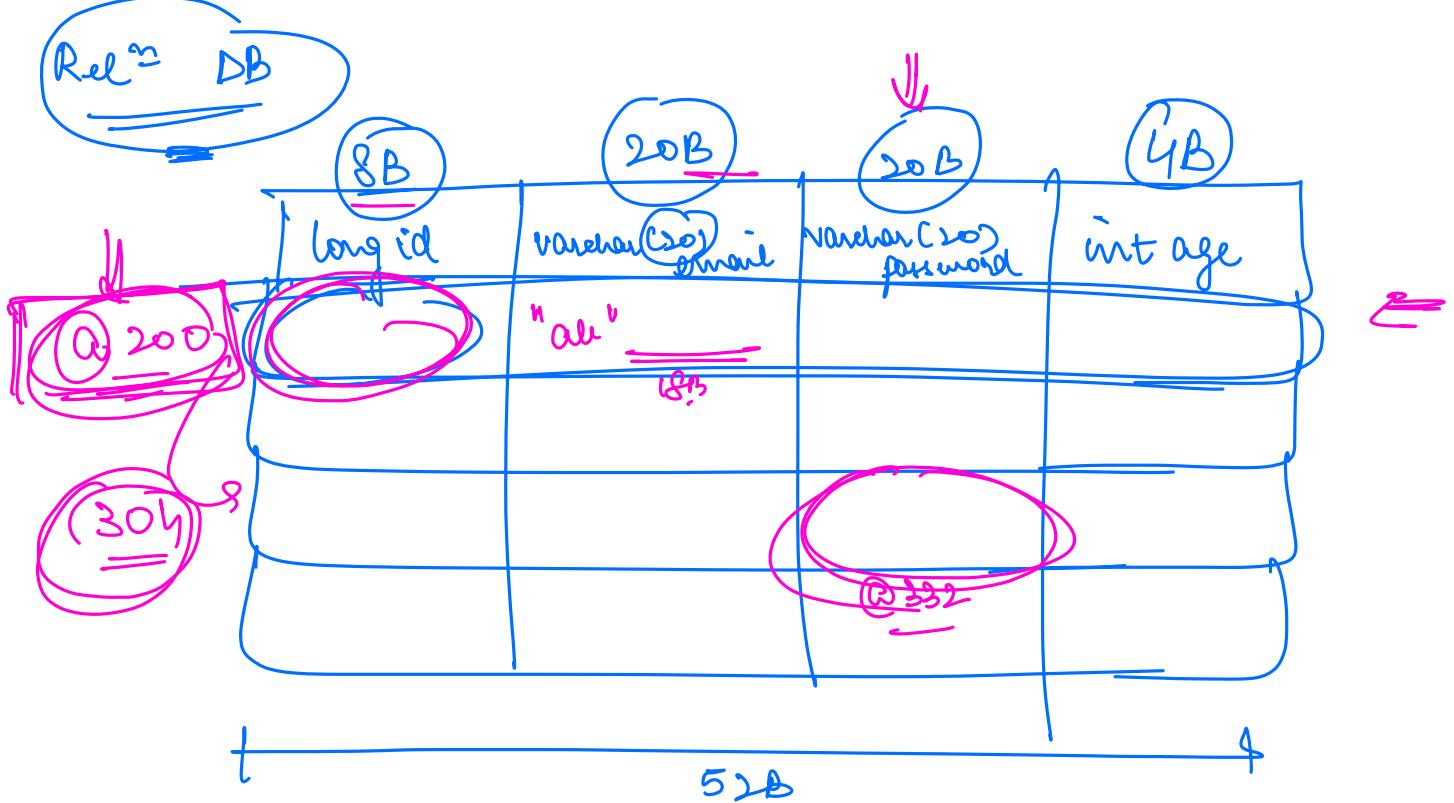
Creating own SQL DB is still easier
↳ Structure

↳ every column has a fixed data type.



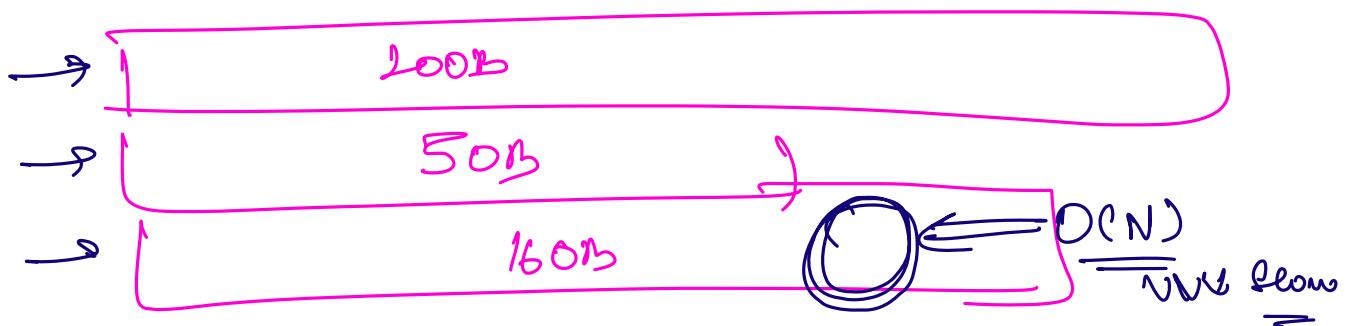
⇒ Binary Search on an array

↳ ~~Worst Case~~



↳ How will I store data internally
↗ How will I query data?

?



As every entity / attribute can have its own segi

isn't possible for me to find value of a attr of iter entity in $O(1)$

"TOP10 — how is india" \Rightarrow [" , — ,

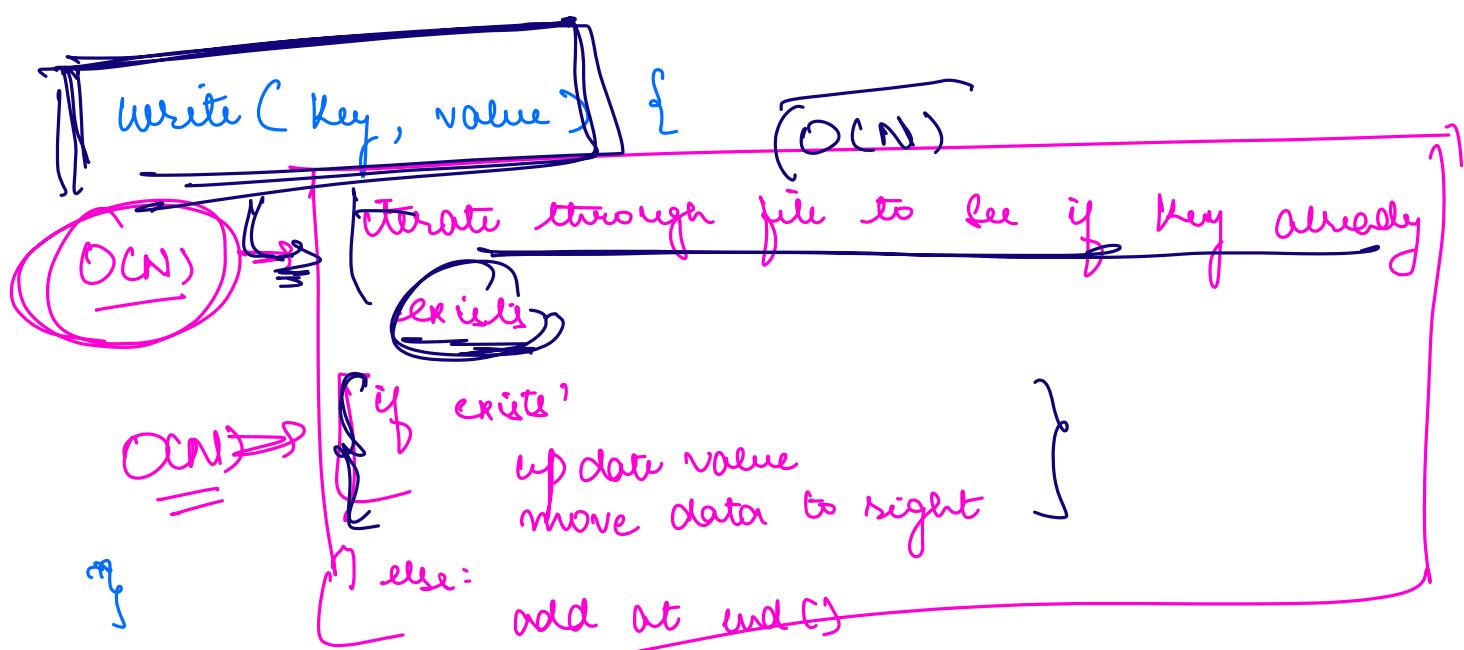
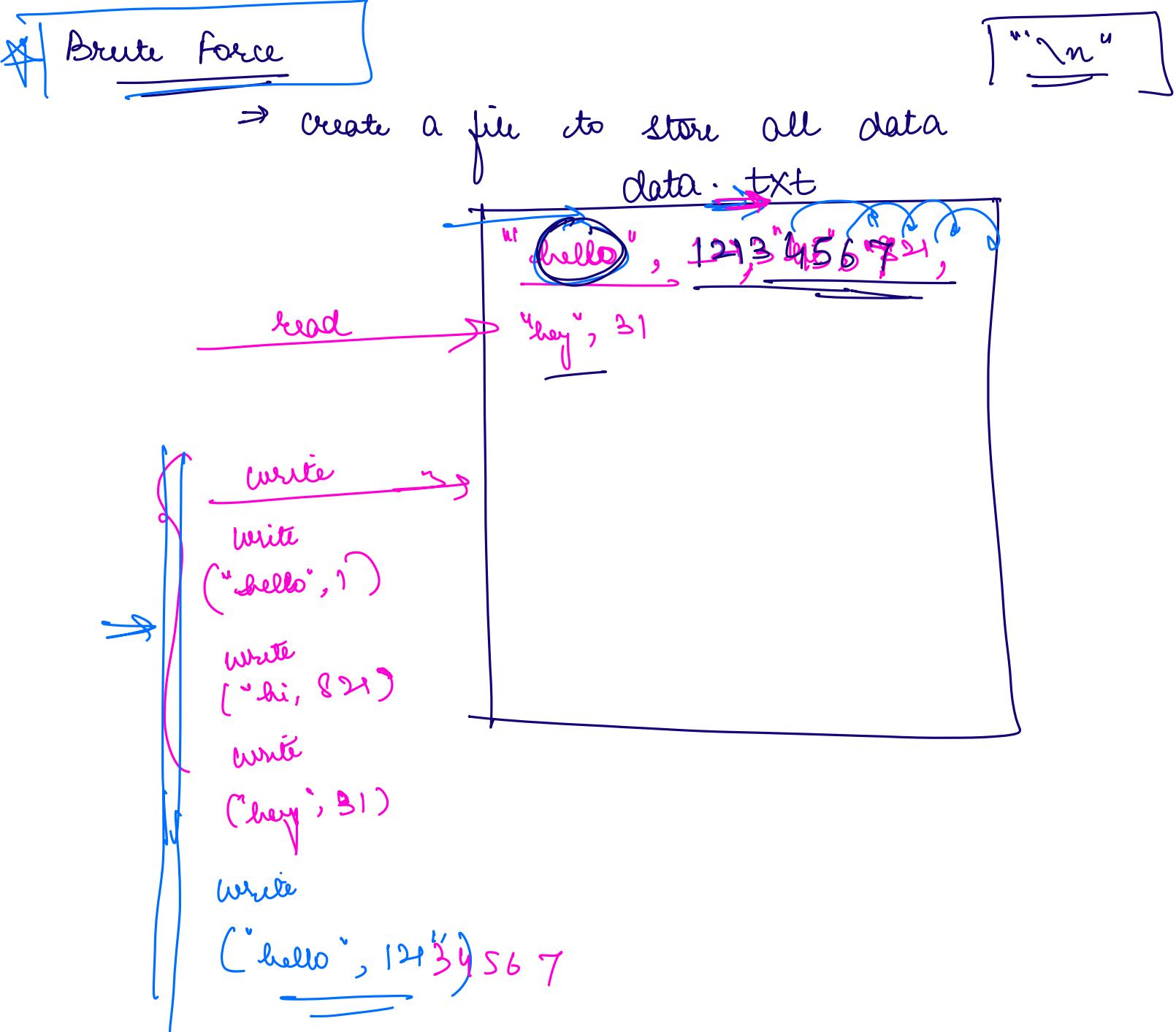
K-V, Document DB

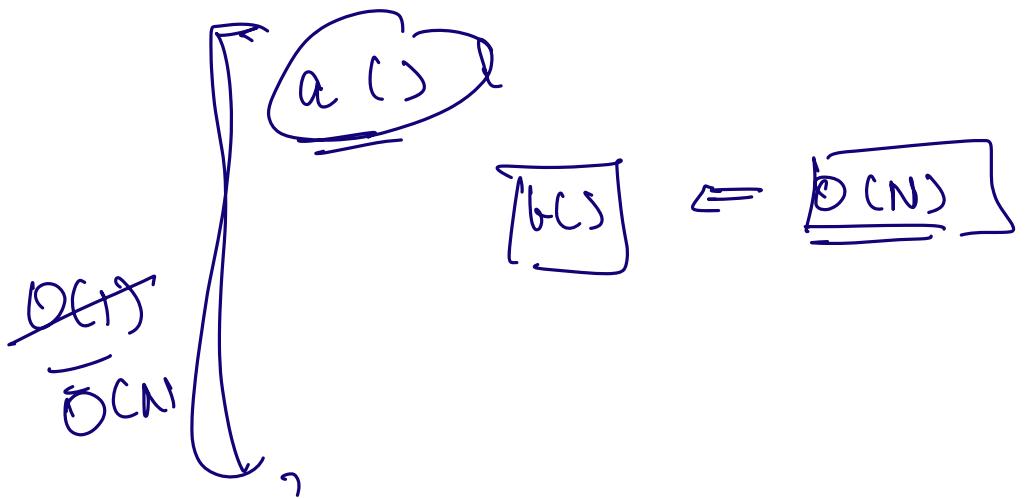
$\Rightarrow O(N)$

\Rightarrow NWVV B ad

How to store data in NoSQL DB

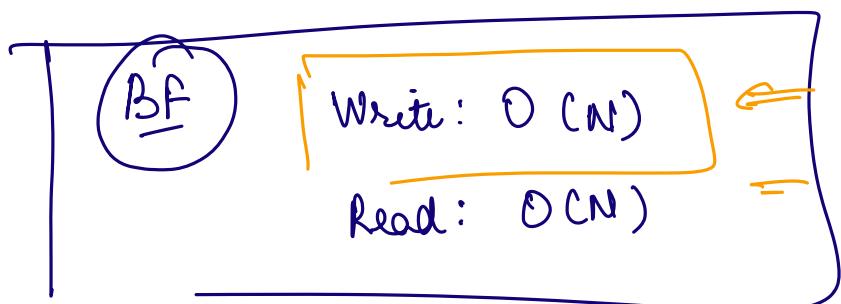
On disk





read(Key)
 iterate through file
 find key
 return if exists

$O(N)$



TC II

Only append on write (no update on write)

Read ("hello")

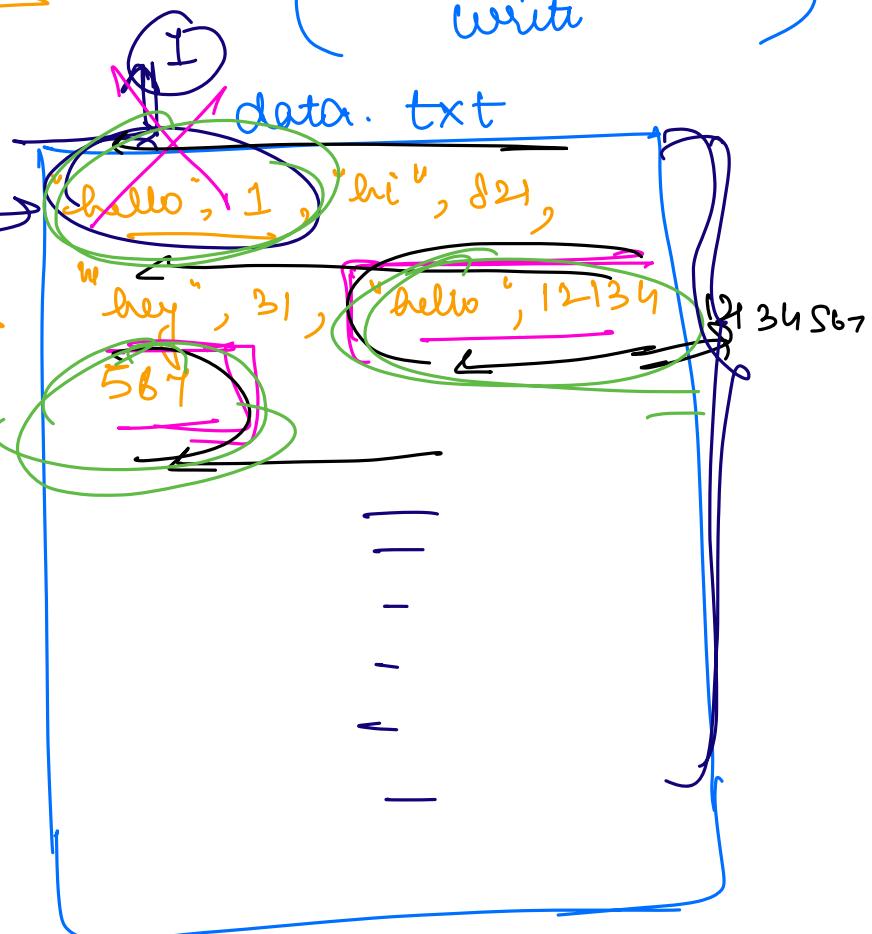
Write ("hello", 1)

Write ("hi", 821)

Write ("hey", 31)

Write

("hello", 1234567)



Now, on write, we are just appending to file and not checking if key already exists

TC

Write $\Rightarrow \mathcal{O}(1)$

Read $\Rightarrow \mathcal{O}(N)$ in worst case

write (key, value)

append to file

$\Rightarrow O(1)$

→ read (key)

start reading file from end

to find key

if exists:

return

else: error

$O(N)$

Problems

① $\Rightarrow O(N)$ is still bad

VVV Slow

② A lot of duplicates
 \Rightarrow Wasted Storage

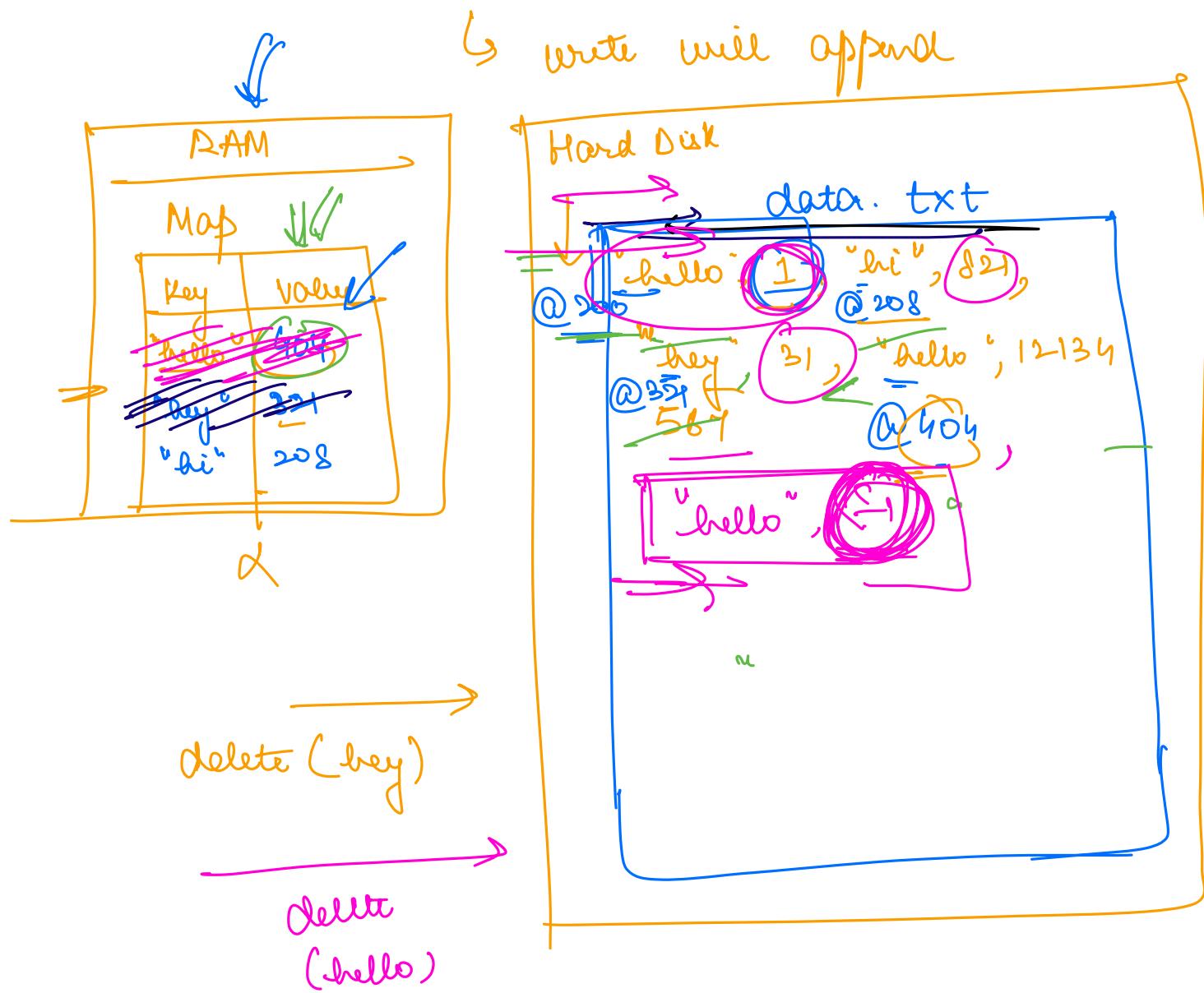
?

$O(N^2 + N) \Rightarrow \underline{O(N^2)}$

13

Also try to optimize reads somewhat

⇒ In disk we have same thing as in V2



⇒ In memory, → create a Map

$O(1) \Rightarrow$ write (key, value)

append to file $\Rightarrow O(1)$

update hm with the latest = $O(1)$
address of that key

read (key) $\Rightarrow O(1)$

if hm.contains (key) = $O(1)$

go to address
return value

else: error

Write $\Rightarrow O(1)$

Read $\Rightarrow O(1)$

What if computer restarts?

- all will go through the file
- rebuild the map
- db will start

Q how to handle delete?

Sol¹:

just delete from HM

⇒ But that will not work because when computer restarts, while rebuilding HM, it will find key in file and bring it back.

Sol²:

⇒ delete from file

⇒ will have to iterate over complete file ⇒ O(N)

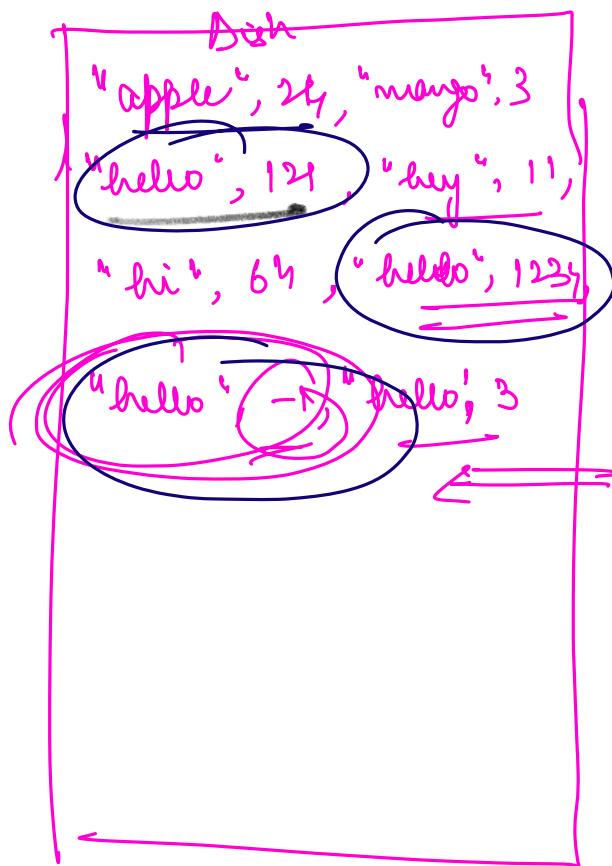
users			
id	name	email	is deleted
1			true false

Tombstone Method

sol²¹³

Tombstone

delete (key) {
remove key from hm \Rightarrow O(1)
add tombstone (key, (-1)) at
null
O(1)
the end of file \Rightarrow O(1)}



Memory

key	Add on Disk
apple	
hey	11
hi	64
hello	3

hm. put (K, (-1))

	write	read
$v_0 \Rightarrow$	updated value $O(n)$ on disk	read the complete file $O(n)$
$v_1 \Rightarrow$	append to file on disk $O(1)$	read the file from end $O(n)$
v_2	append to disk + update HM in Mem $O(1)$	get add from HM + go to add on disk and read $O(1)$

restores

deletes

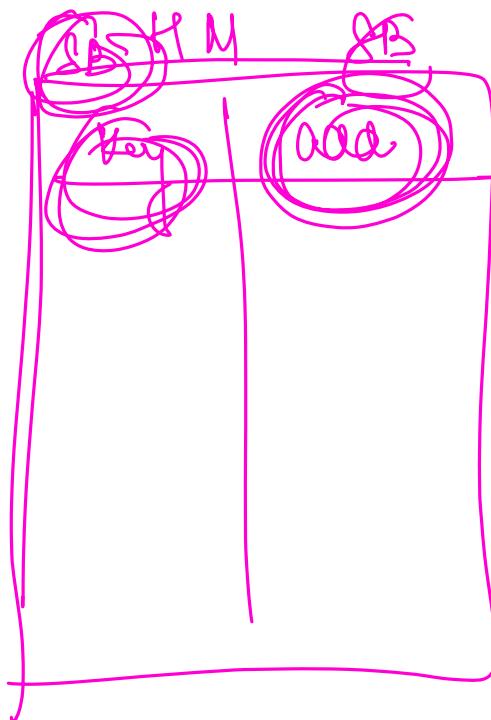
The problem still remaining is

→ a lot extra disk storage is being used.

1 TB disk \Rightarrow

every entry on avg takes 50B

$\Rightarrow \# \text{Keys} = \frac{10^3 \times 10^9 \text{ B}}{50}$



$$\Rightarrow \frac{10^2 \times 10^{10} \text{ B}}{50\%} \rightarrow [10 \times 10^9 \text{ keys}]$$

$$\Rightarrow 16 \times 20 \times 10^9 \text{ B} \rightarrow 320 \text{ GB}$$

- Our RAM might not be able to store
- the value of all keys

Problems

①

Too much storage

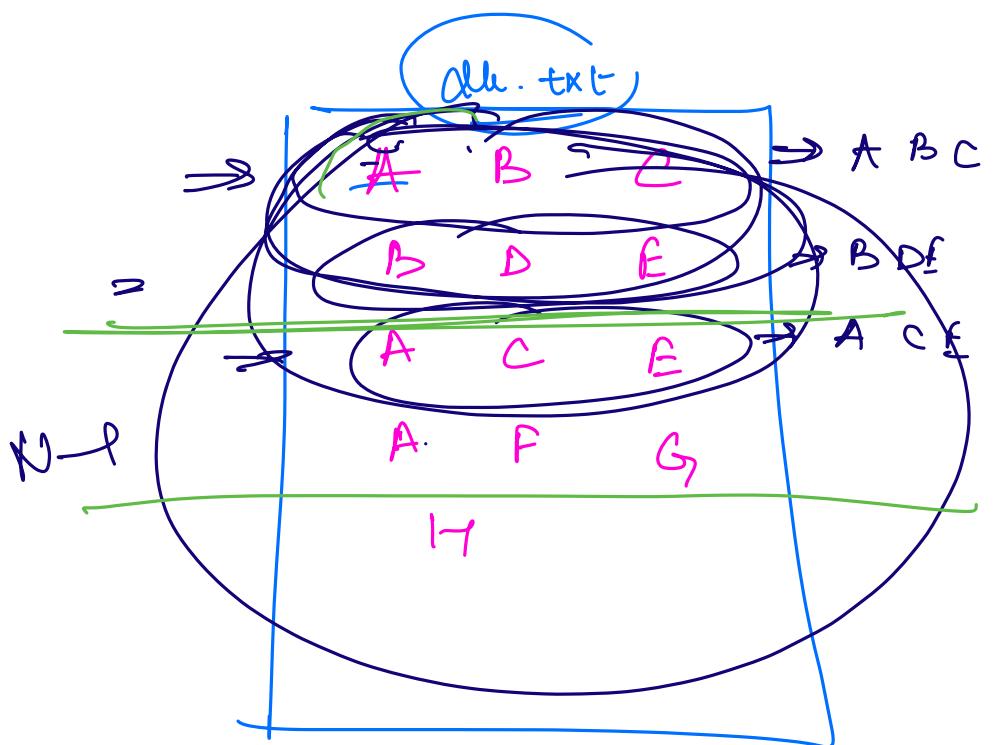
~~⇒ duplicates~~

②

RAM won't be able to fit complete

HM





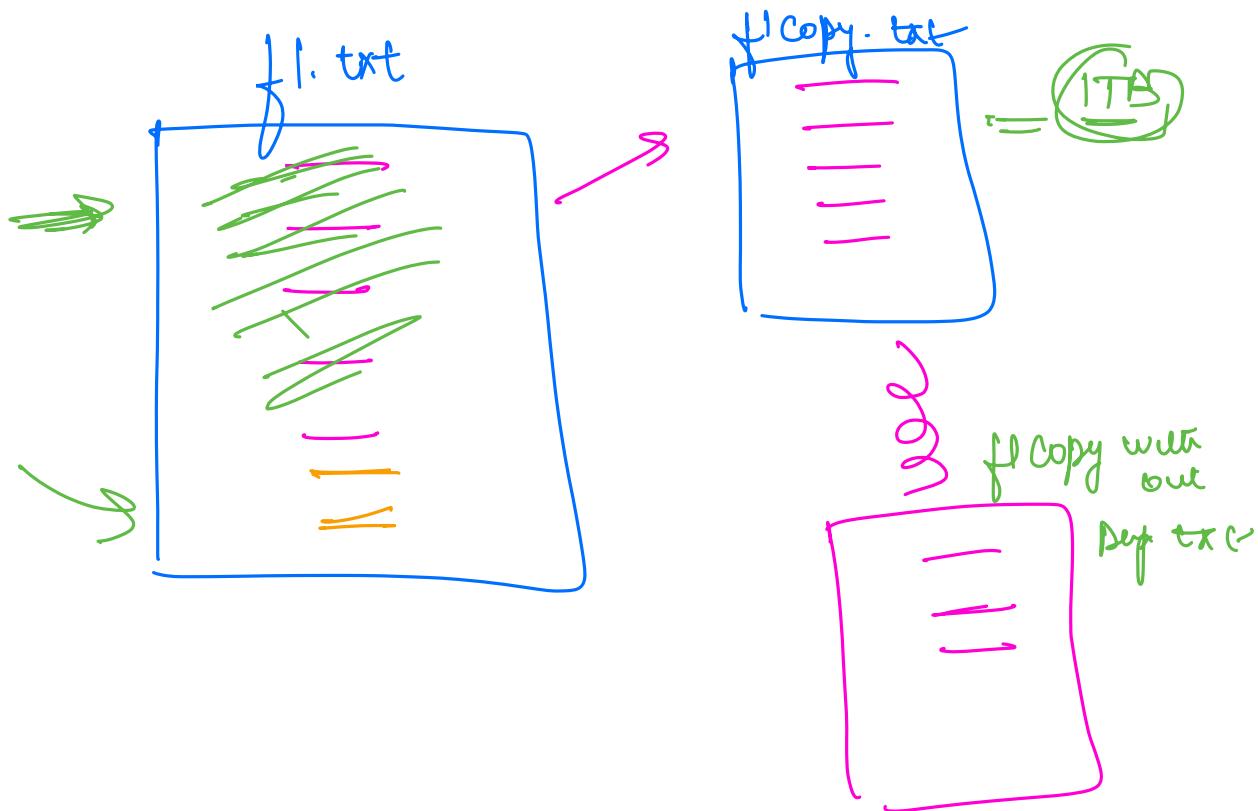
How to remove duplicate

Solⁿ1: at the time of insertion \rightarrow writes will be slow

Solⁿ2: at restart

Solⁿ3: To background job running every X hrs.

{ (i) create a copy of file
 (ii) iterate over copy, find duplicates
 (iii) create a new version with all dup removed



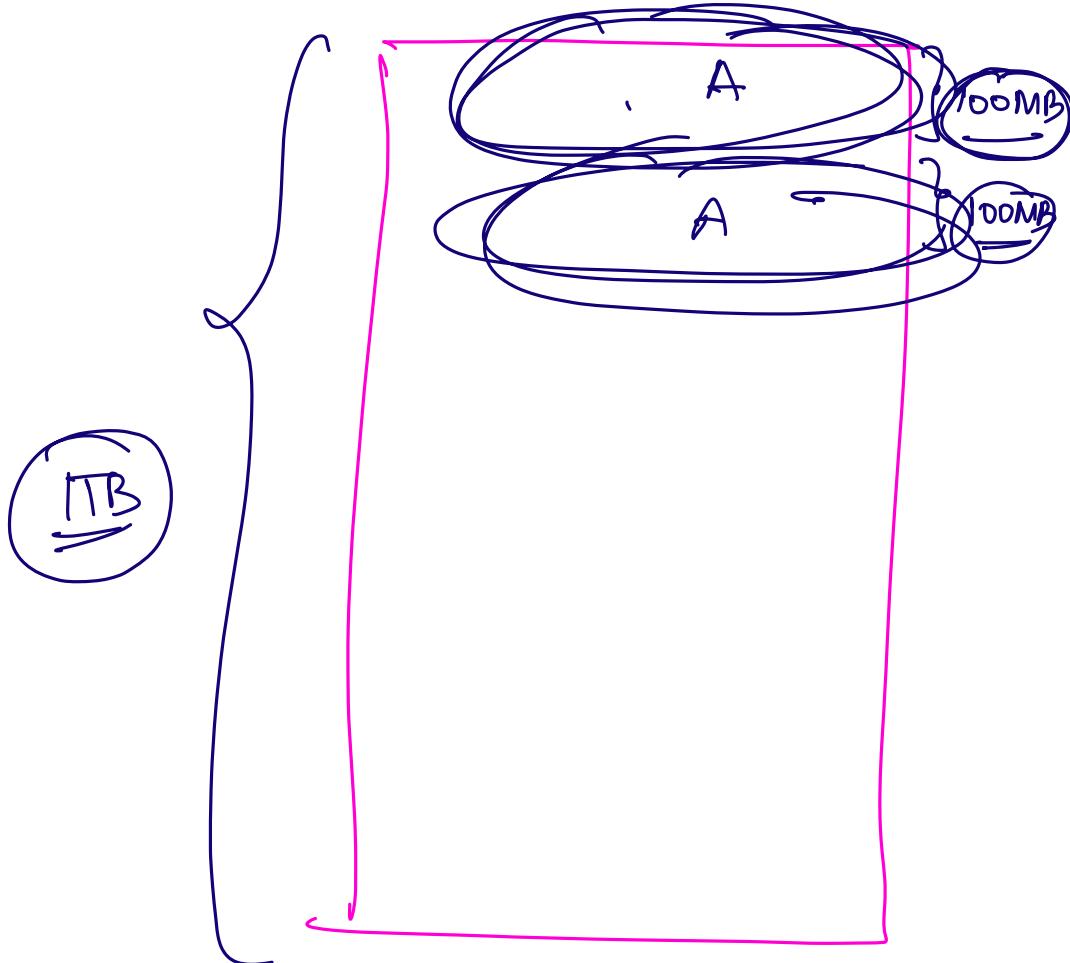
⇒ Implementation is very troublesome.

```
removeDuplicates(list < String >) {
```

HSet

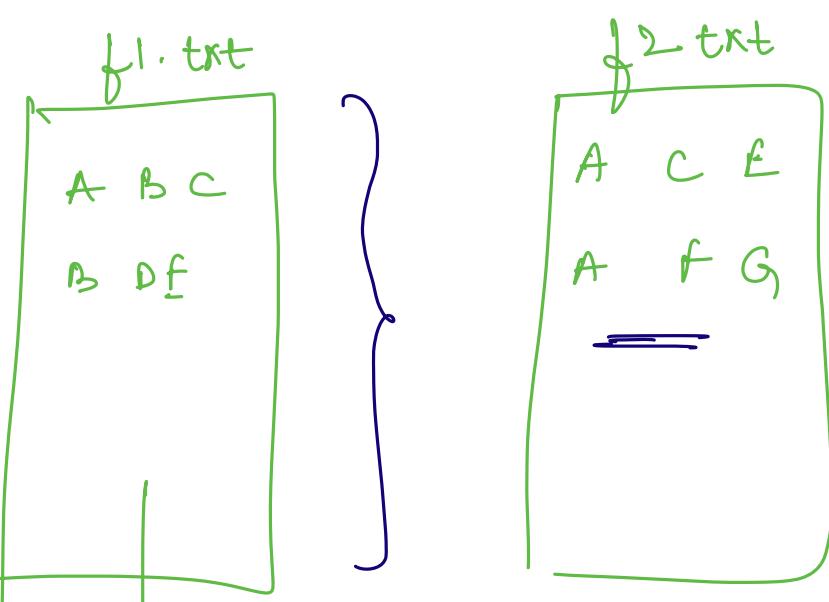
}

⇒ for a 1TB file putting all keys in a Hash Set is not scalable.



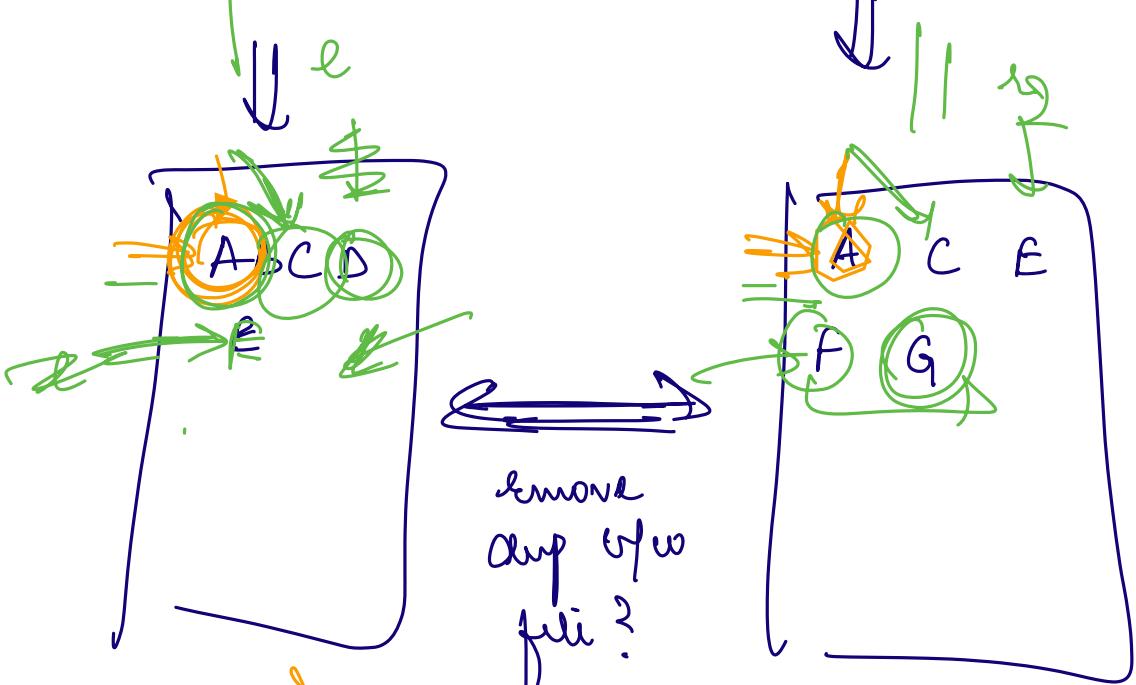
Solⁿ
→

Split the data amongst multiple files

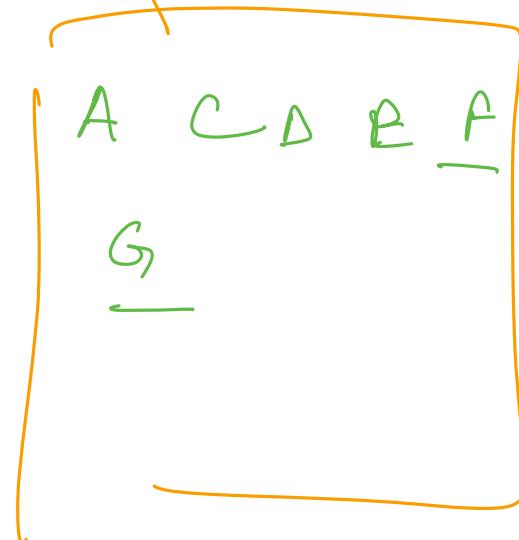


each file is small enough to be brought into



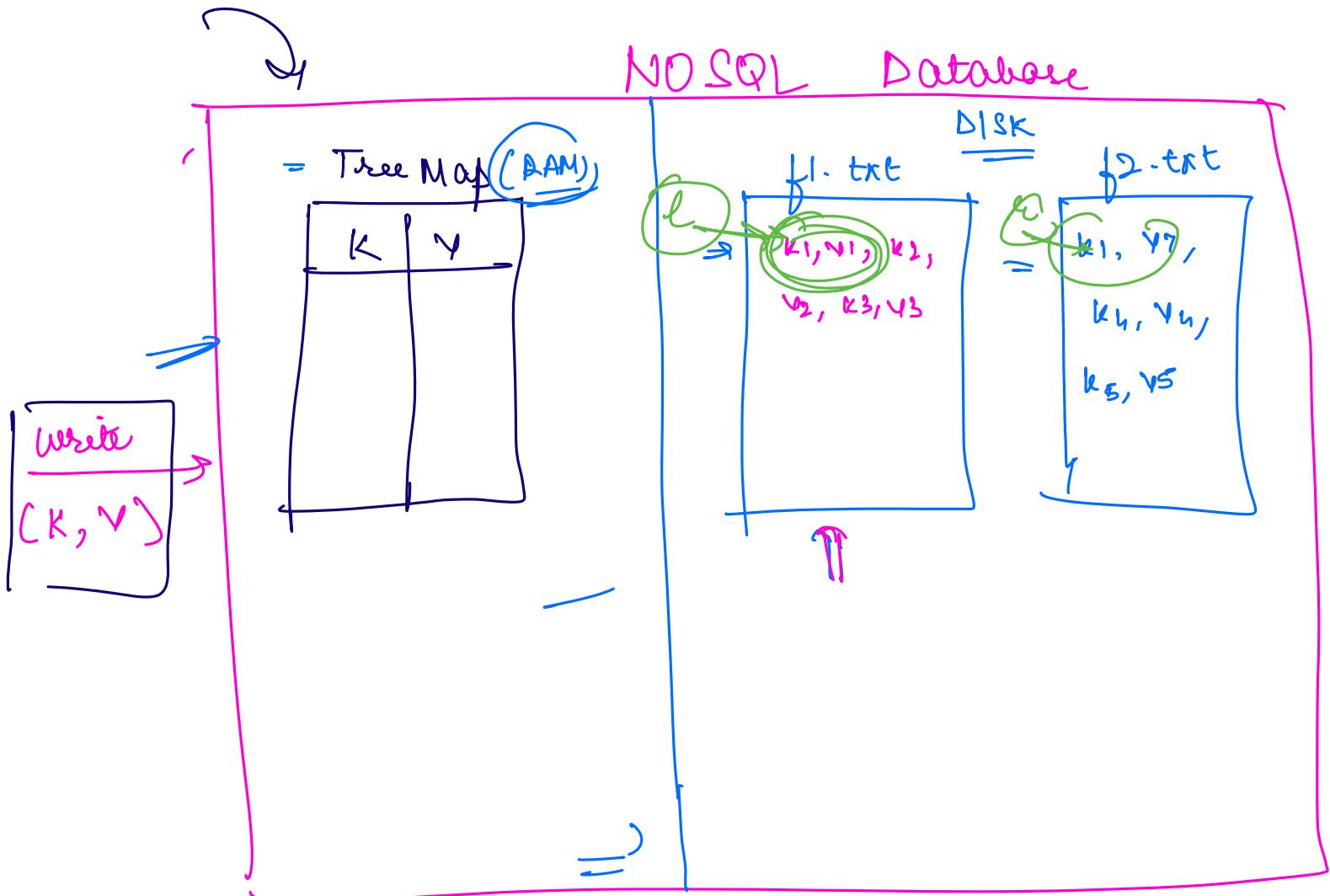


is there a technique using which
 I don't bring the whole file into RAM
 but still can find dup w fili?



=> if somehow I am able to get
 my data into multiple file, and
 each file is sorted, I can

remove dups will eat



⇒ Write (K, V) {
 treeMap.put(K, V);
}

⇒ $O(\log N)$

= job running every X Min () {
 create a copy of treeMap
 clear current treeMap.
 → tc treeMap

Within a Map
no dups are
possible

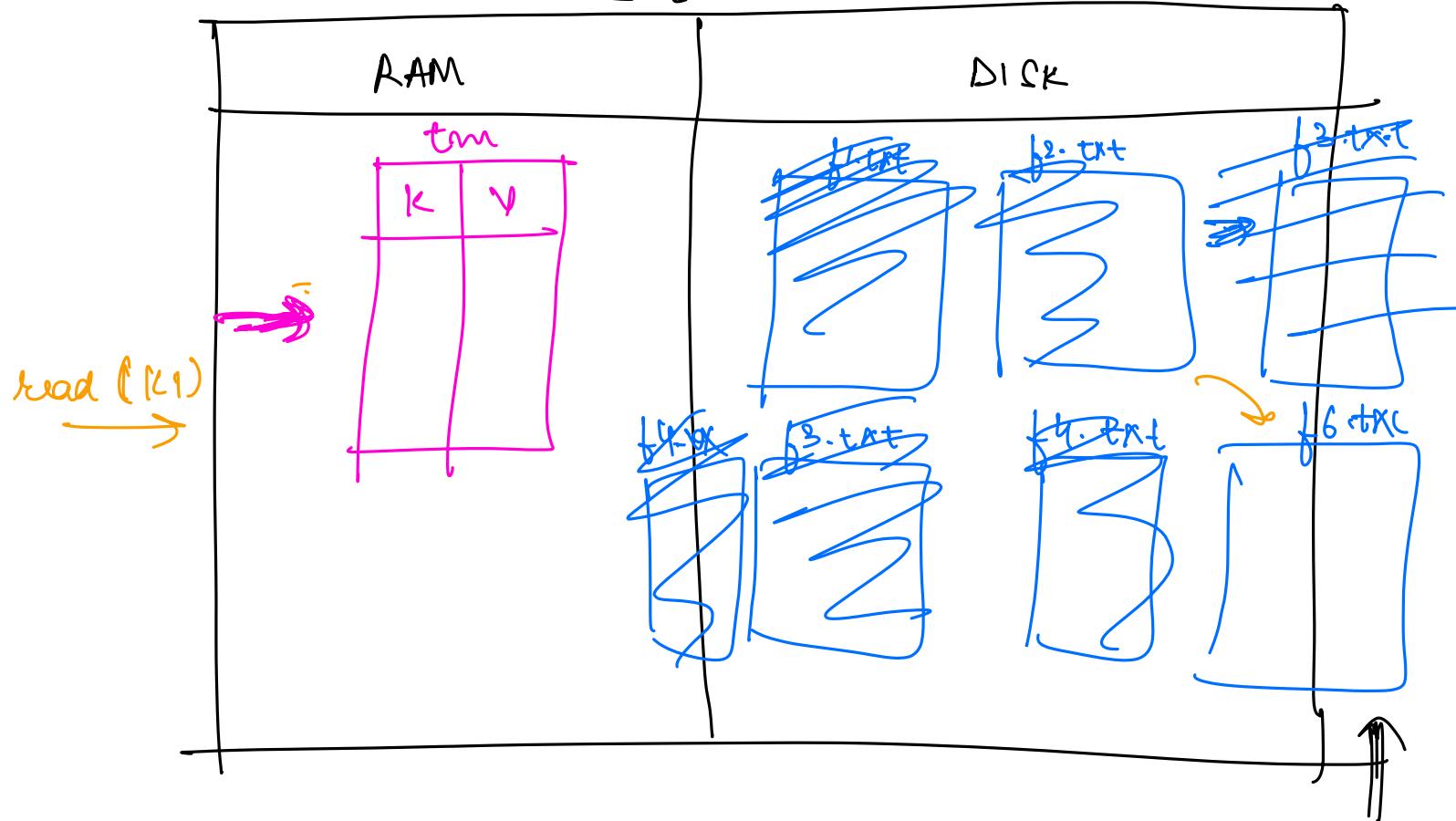
```
treeMap = new TreeMap< > }
```

V iterate through copy of tc:
put the k, v on a new file

}

- ① within a file \Rightarrow no dups
- ② each file will be sorted.
 - \Rightarrow if dups are now there on my disk
they are b/w 2 files.
 - \Rightarrow As both files are sorted I can remove
dup without needing to bring everything
in mem.

NO SQL Database



Write (K, V) {

 treemap.put (K, V)

}

VVV fast

Convert TreeMap To New File ()

Create a copy of tm and empty current tm.

Create a new file:

Iterate over the copy of TreeMap:

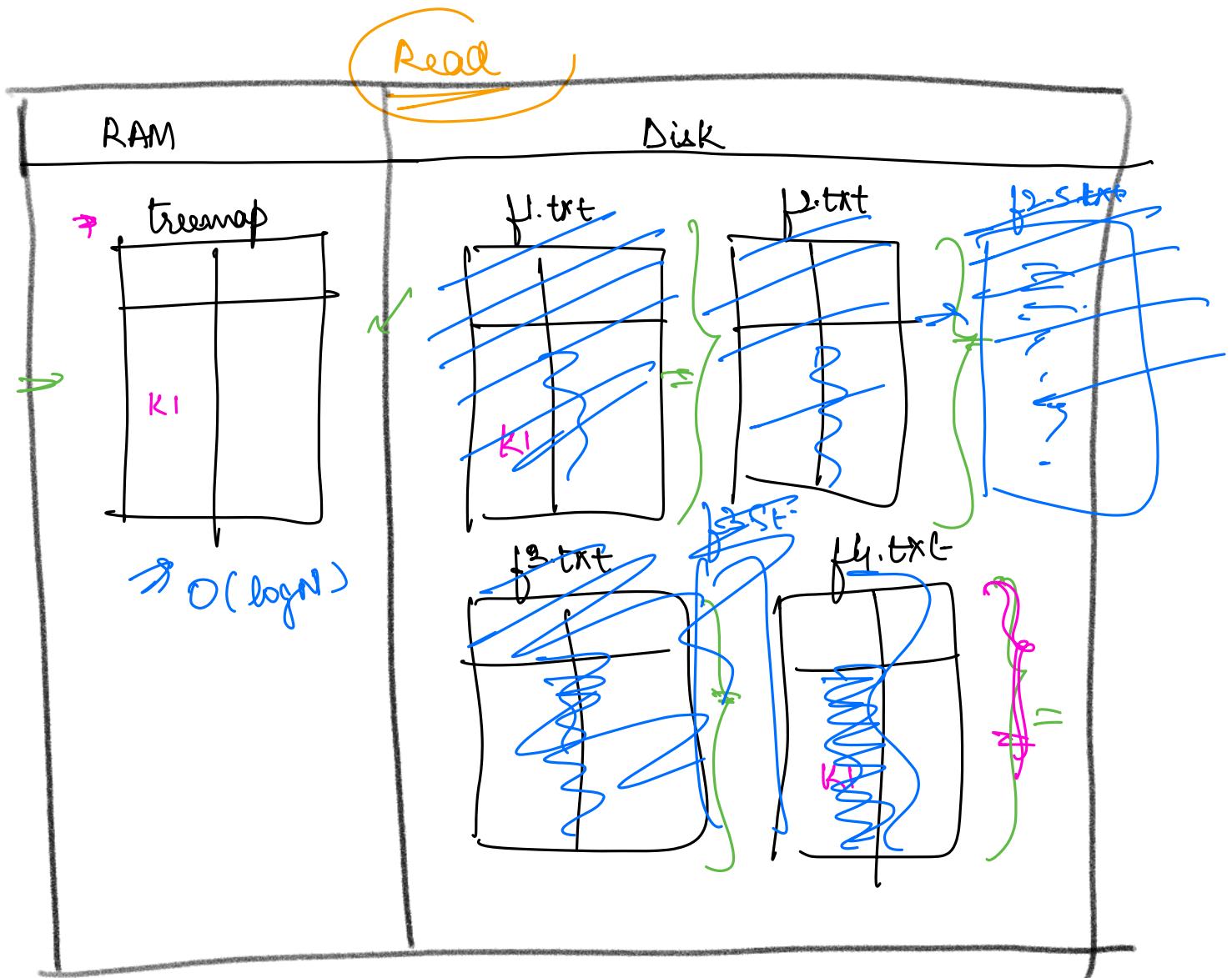
 put key, value on file

} \Rightarrow within each file, no dupes and sorted data

```
TreeMap< > copy = treeMap;  
treeMap = new TreeMap< >();
```

{ merge files () { // also called ~~compaction~~
while there is more than file:
 pick 2 files
 merge them
}
} ⇒ at end 1 file with no dup.

BUT!!! How will I read data ???



```
read (key) {
```

```
if (treemap. contains (key) ) {  
    return treemap.get(key);  
}
```

三

read each file from latest to oldest:

if file contains key:

return key

return -1

{}

// for which key read will be very slow:

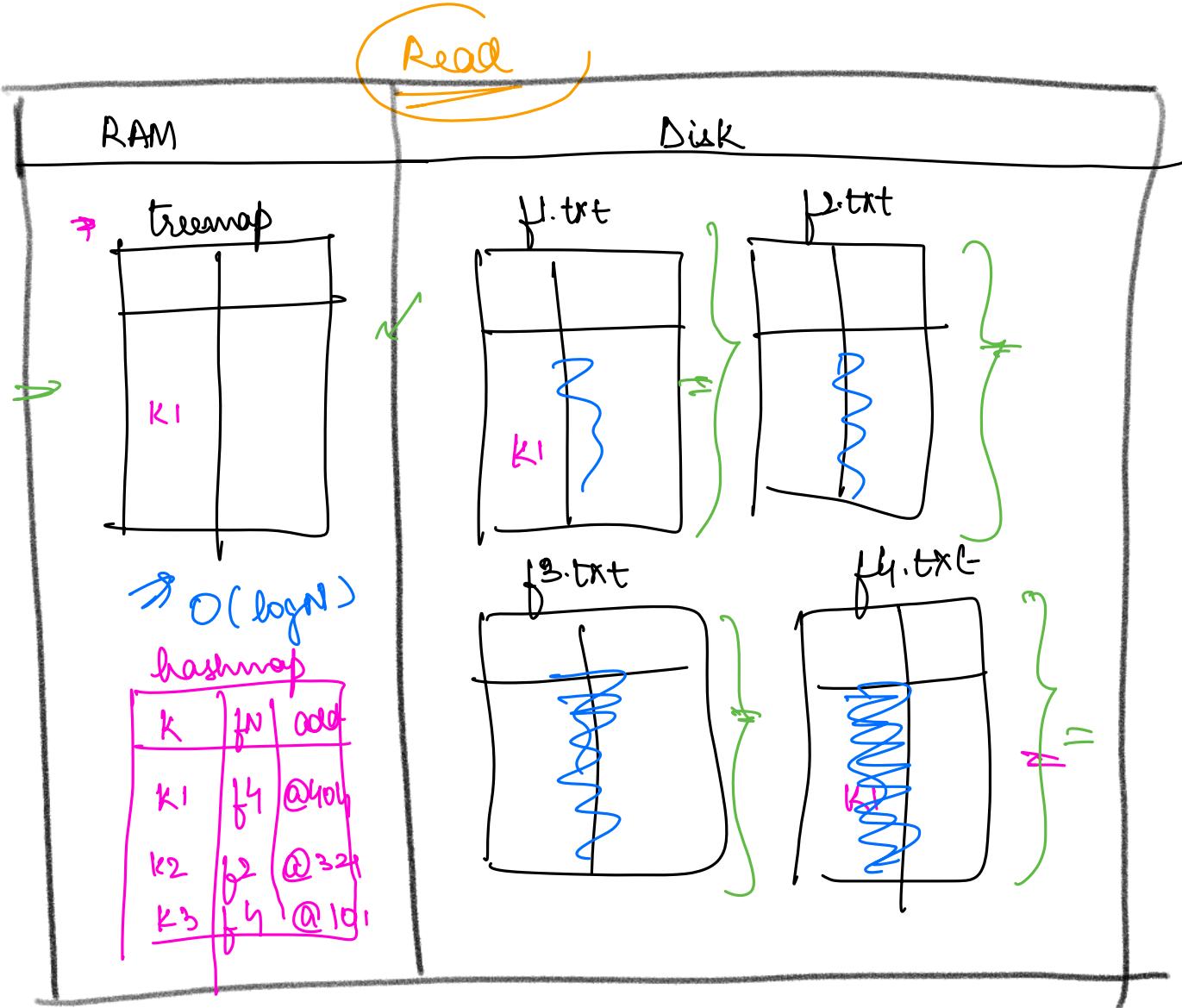
⇒ reading non existing key will be slower

⇒ worst case \Rightarrow have to read all files -

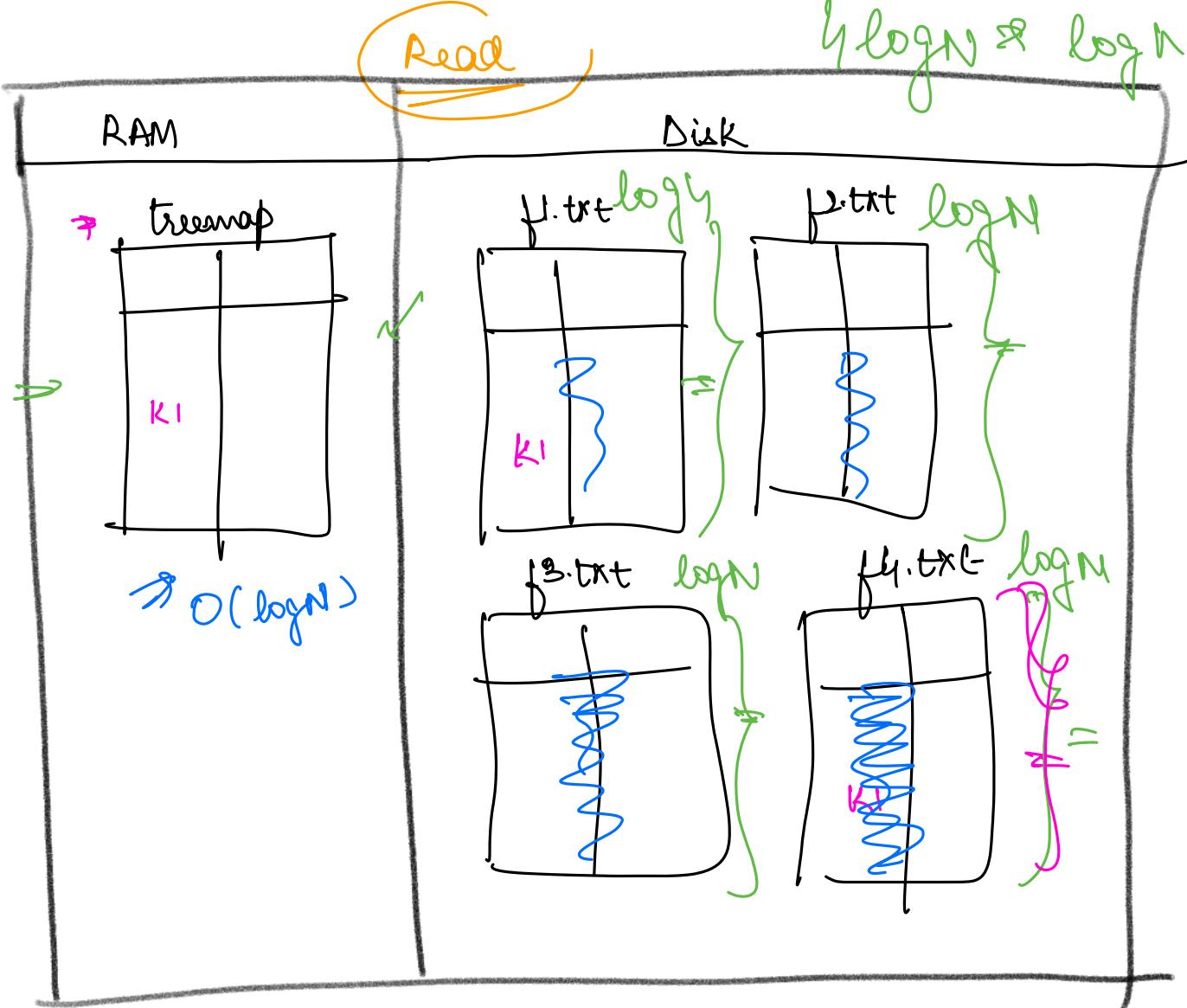
⇒ PROBLEM: reads have now again become slow.

What if I again create a HM, which tells for every key, what file and writing that file at what address that key exists

Slow Reads v/s I can't store HM in Mem.



→ BUT!!! Data within each file is sorted.
Can I make use of this observation

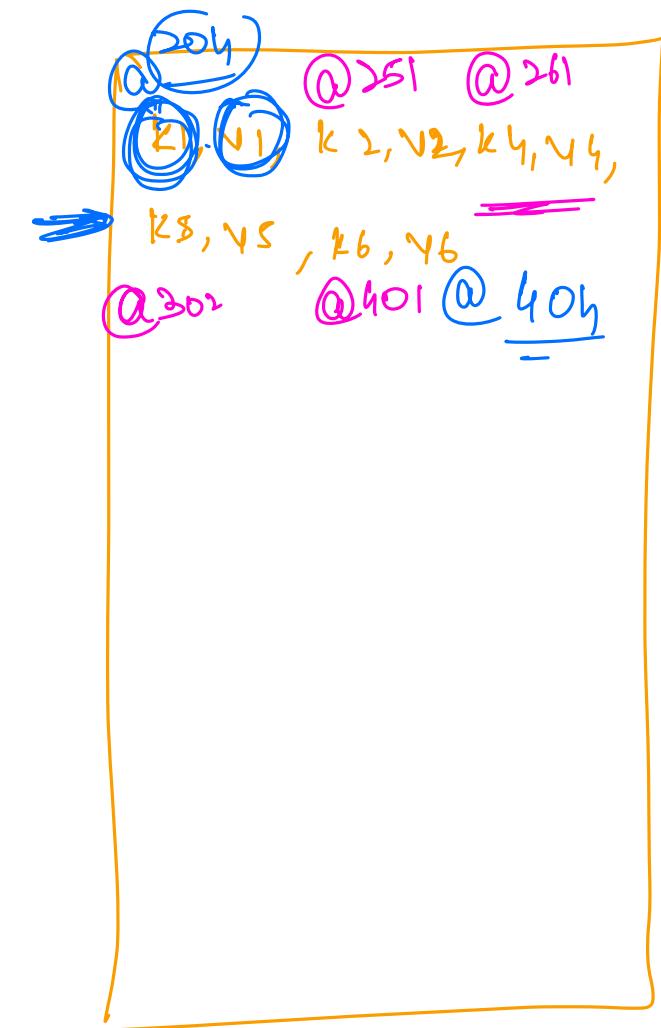


```

read (key) {
    if (TreeMap.contains (key)) {
        return TreeMap.get (key)
    }
    read each file from latest to oldest:
    if file contains key:
        find the key within file via BS
        return key
    return -1
}

```

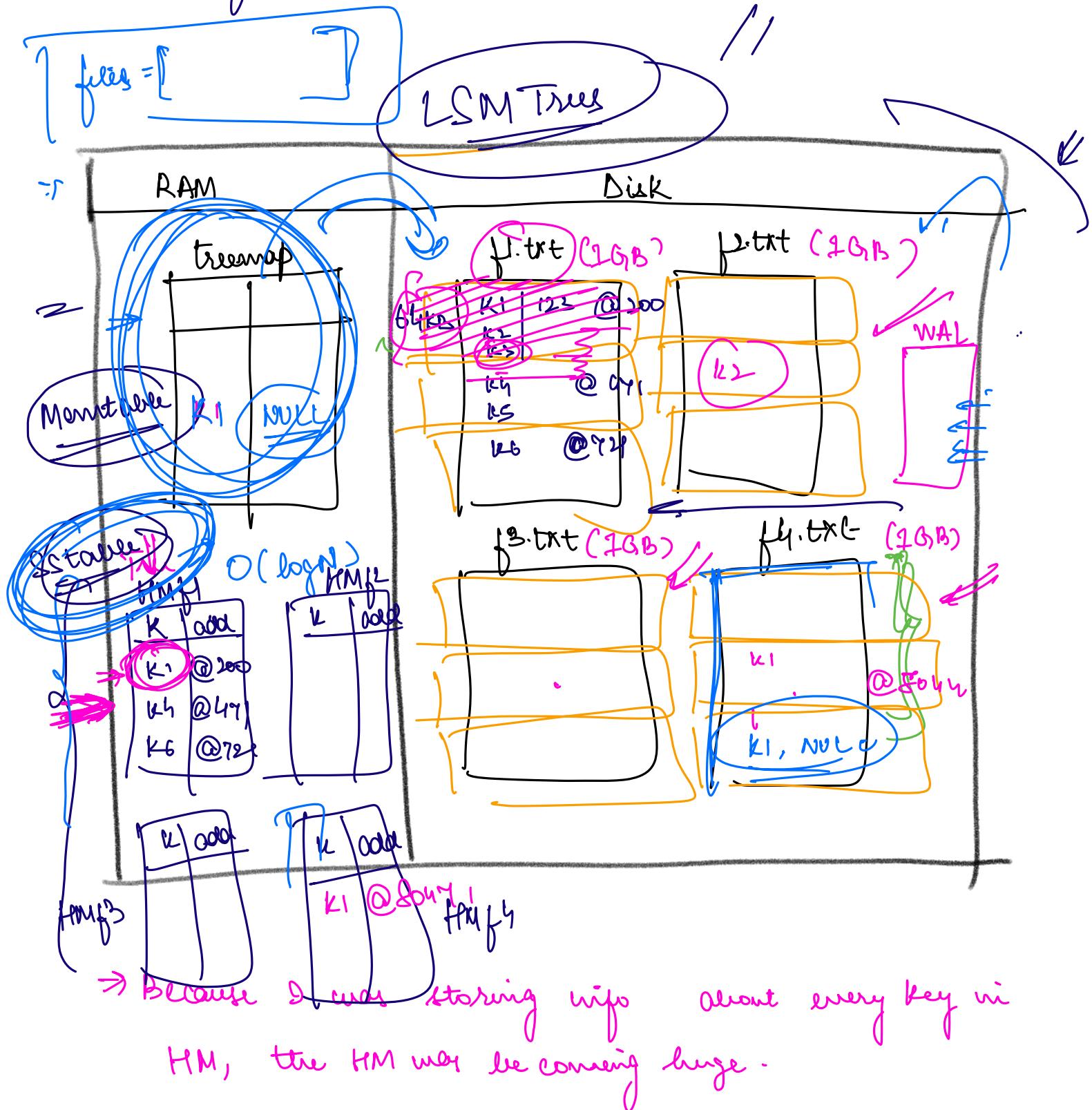
$\Rightarrow O(N)$



$$\frac{404 + 204}{\Rightarrow \boxed{304}}$$

Only Problem: Find a key in logN circuit file

$\text{Sol}^n \rightarrow$ Create another additional Hashmap
for every file.



to find a key in a file

(i) ~~Binary Search~~

(ii) ~~Linear Search~~

→ v.v. slow.

→ 1 GB

Divide file into multiple blocks,

(i) Create a HM associated to every file

(ii) In HM, dont store info about every
key but some keys (first key of every
block)

Now to find a key in a file:

e.g. find K3 in f1:

→ I will go to the boshmap of f1

$O(\log N)$ → I will find the key just smaller than
(K3)

→ I will just read that block

→ 64 KB

```

read (Key) {
    if (TreeMap. contains (Key)) {
        return TreeMap. get (Key);
    }
    for each file from latest to oldest:
        find first key smaller than k
        in TreeMap of that file:
        read that complete block linea
        if k exist
            return
}
return -1

```

{ Writes are fast b/c \Rightarrow put in "memtable" }
 reads are slower \Rightarrow you may have to read
 some part of multiple files.

But can we lose writes

Why: I am only writing to memtable
(treemap in mem)

Ans: No

Every DB has 2 things that it stores -

(i) Data

(ii) Write Ahead Log (WAL)

Commit log
Bin log



a file where we keep a note
of every write query that is
coming to my system

WAL

T0 \Rightarrow put("Hello", 3)
T1 \Rightarrow put("Hello", 2)
T2 \Rightarrow put("Hello", 1)

Write queries coming from the

\Rightarrow file that has
history of
everything
in DB.

