Agenda.

⇒ NoSQL Internals.

↳ How NoSQL DBs stores data internally.

# SQL DB.

↳ Structured Data

↳ fixed Schema.

Users.

| | id | name | email | — | — |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |
| | | | | | |

6B → id

50B (name), 40B (email), 12B, 20B

140B.

@100

$R_1 \doteq 140B$

| $R_1 \doteq 140B$ | $R_2 \equiv 140B$ | $R_3 \doteq 140B$ |
|---|---|---|
| ←40→ .... | | |
| | | |

(int) a[] :

108, 116, 124

| 4 | 10 | 20 | 7 | 8 | 74 | — |
|---|---|---|---|---|---|---|

@104, 112, 120, 128

4B.

$$a[i] \leftarrow \underbrace{a[0] + i \times 4}_{O(1)}$$

# NoSQL DB.

↳ Unstructured Data

↳ No fixed Schema.



⇒ WAL

↳ Write Ahead Log.

↳ Append Only file.

C
R
U
D
⟩ Write

WAL

```
A = 12
B = 40
C = 20
A = 24
D = 100
B = 41
E = 500
C = -1
A = 30
```
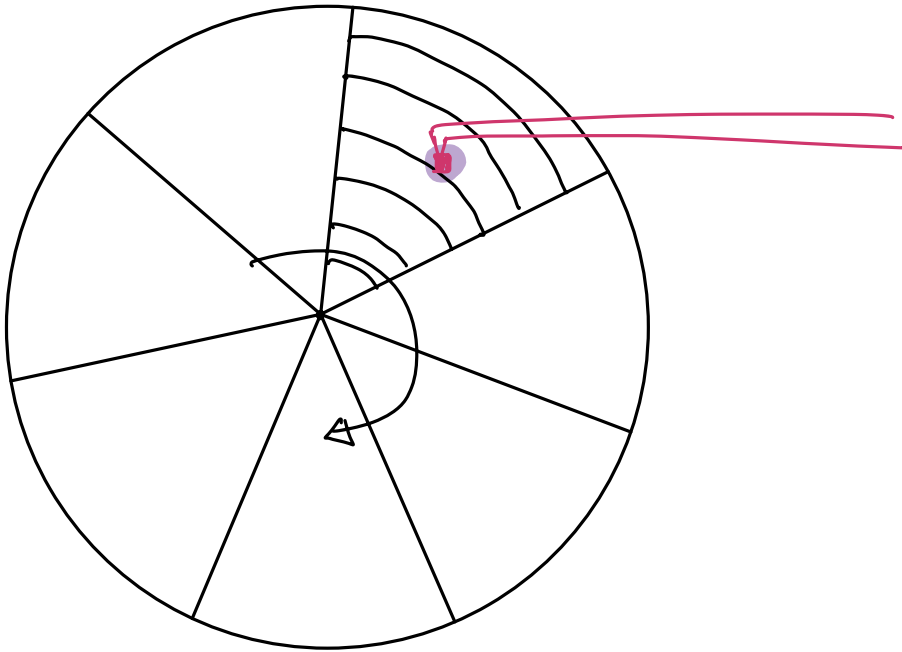EOF →

K, V

WAL file

→ Maintains the Complete history of Operations.

→ Reconstruct our DB from WAL.

→ Used to sync Replicas

Write op^n TC in WAL ⇒ $O(1)$

Read op^n TC in WAL ⇒ $O(N)$ → # of Write op^n

⤷ Start reading the file from starting till the end and get the latest value of the key.



SQL DB

⤷ Read TC : $O(logN)$

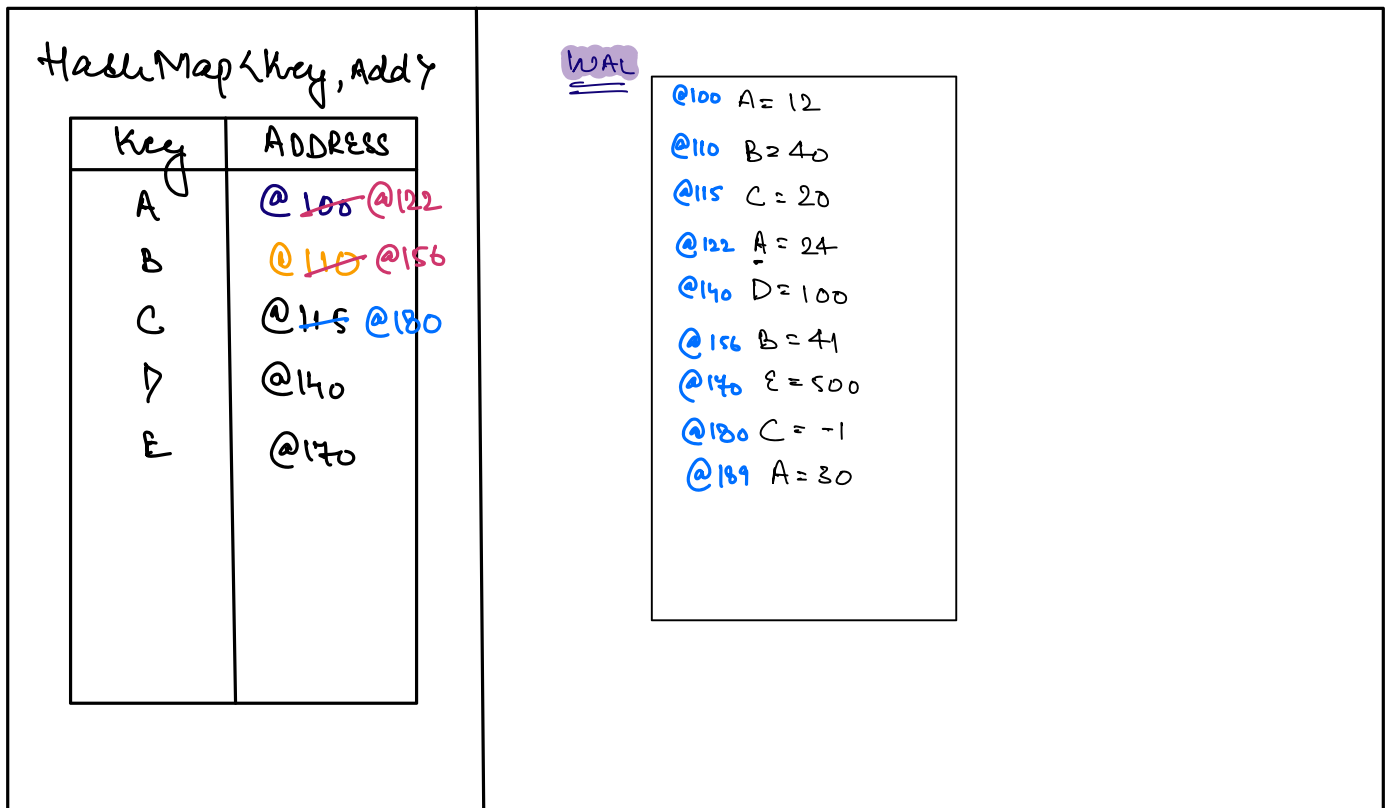⤷ Write TC : $O(logN)$

# Approach #1

→ Only WAL.

Read TC : $O(N)$

Write TC : $O(1)$

# Approach #2

→ WAL + HashMap.

| RAM | HDD. |
|---|---|

**RAM**

HashMap <Key, Addr>

| Key | ADDRESS |
|---|---|
| A | @~~100~~ @122 |
| B | @~~110~~ @156 |
| C | @~~115~~ @180 |
| D | @140 |
| E | @170 |

**HDD.**

WAL

@100  A = 12
@110  B = 40
@115  C = 20
@122  A = 24
@140  D = 100
@156  B = 41
@170  E = 500
@180  C = -1
@189  A = 30

# Write Opⁿ

↳ Append in WAL ⇒ $O(1)$

↳ Update HM ⇒ $O(1)$

TC of write $\Rightarrow$ $O(1)$

Read $Op^n$ $\Rightarrow$ $O(1)$

$\rightarrow$ Get the address of the key from HM and read the value from this address in WAL

Cons.

$\rightarrow$ HM is present in the RAM (which is volatile storage), in case our m/c restarts, we'll have to rebuild the complete HM from scratch.

Size of WAL = 10 TB

Size of 1 entry = 10 B.

No of entries = $\dfrac{10 TB}{10 B}$ = $\dfrac{10 \times 10^{12} B}{10 B}$

$= 10^{12}$ entries.

Assumptions : Unique keys = $\dfrac{10^{12}}{10}$ = $10^{11}$

Size of 1 K,V pair in HM $\approx$ 16 B.
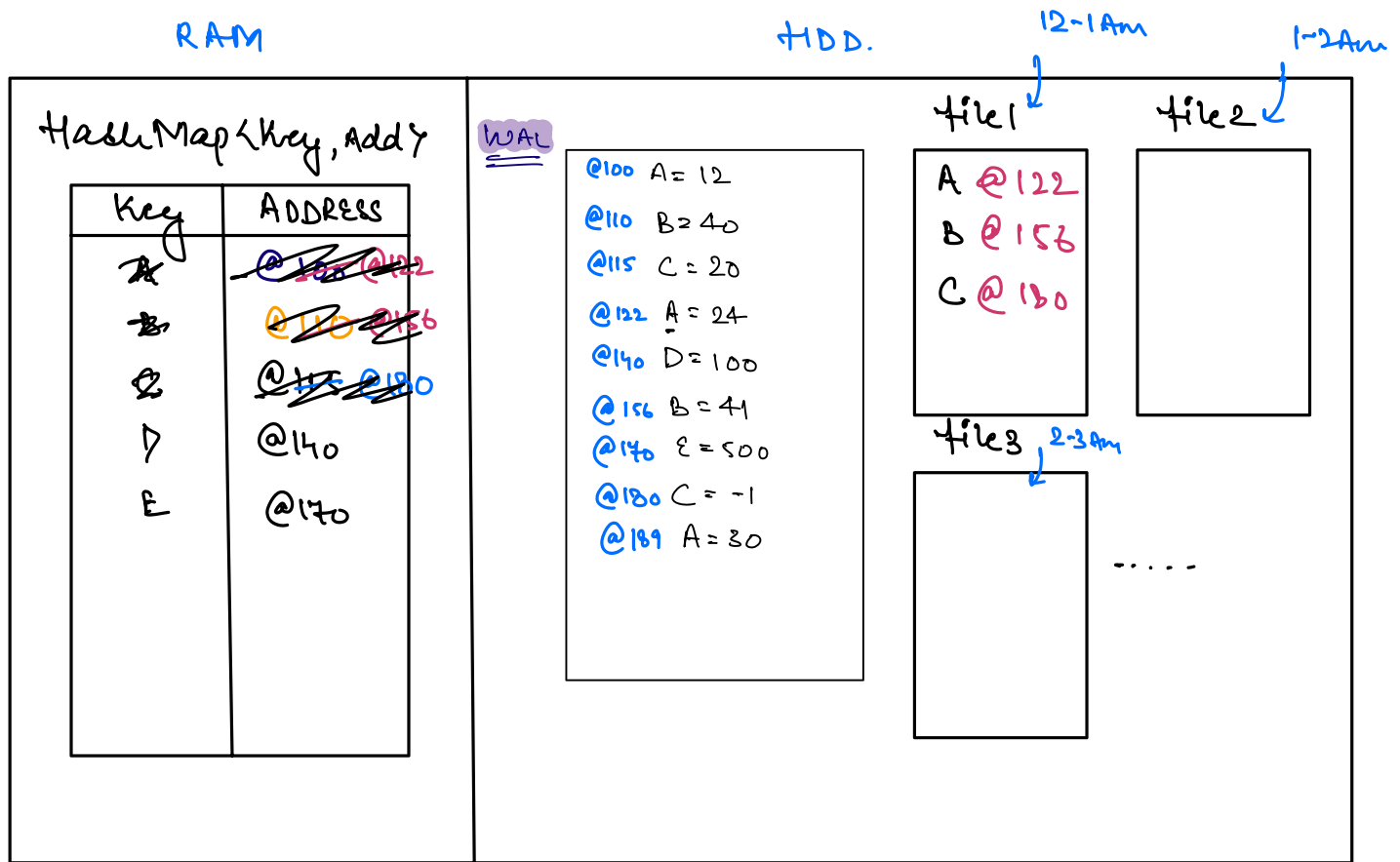
Size of the HM $= 10^{11} \times 16$ B.

$$= 1.6 \times 10^{12} \text{ B.}$$

$$= 1.6 \text{ TB.}$$

$\Rightarrow$ Size of the HM can be huge, it will cost huge amount of money.

# Background Script.

$\Rightarrow$ Every 1 hour, take the Data from HM and put it in the file in HDD & reset the HM.

## RAM

HashMap<Key, Add>

| Key | ADDRESS |
|-----|---------|
| A~~A~~ | ~~@100~~ ~~@122~~ |
| ~~B~~ | ~~@110~~ ~~@156~~ |
| ~~C~~ | ~~@115~~ ~~@180~~ |
| D | @140 |
| E | @170 |

## HDD.

WAL

@100  A = 12
@110  B = 40
@115  C = 20
@122  A = 24
@140  D = 100
@156  B = 41
@170  E = 500
@180  C = -1
@189  A = 30

**file1**  *12-1Am*

A @122
B @156
C @180

**file2**  *1-2Am*

**file3**  *2-3Am*

. . . . .

---

## Write

→ Write in WAL ⟹ O(1)

↘ Write in HM ⟹ O(1)

## Read.

↘ find the key in the Map

**if found** → Get the Address from HM & read from WAL.

**if NOT found** → Read all the files from latest to oldest (Max 24 files)

TC : $O(1) + (\text{\# of files}) \times \boxed{N}$

$\downarrow$
HM

→ Iterate every hourly file linearly.

$: (\text{\# of files}) \times \underline{N}.$

Summary.

|  | Read TC | Write TC |
|---|---|---|
| Only WAL | $O(N)$ | $O(1)$ |
| WAL + HM | $O(1)$ | $O(1)$ |
| WAL + HM + Hourly files | $(\text{\# of files}) \times \underline{N}.$ | $O(1)$ |

# Binary Search

└→ Sorted Dataset + Equal Size Data.

@100                  116                              136                    166

| 1 | 4 | 12 | 19 | 24 | 30 |

4B    12B      8B      12B      10B    20B

104                  124      133      146

$$\frac{100 + 166}{2} = 133$$

⇒ Instead of HashMap, use Tree Map.

**RAM**

Tree Map ⟨Key, Add⟩

| Key | ADDRESS |
|-----|---------|
| A | @100 @122 |
| B | @110 @156 |
| C | @115 @180 |
| D | @140 |
| E | @170 |

**HDD.**

WAL

@100  A = 12
@110  B = 40
@115  C = 20
@122  A = 24
@140  D = 100
@156  B = 41
@170  E = 500
@180  C = -1
@189  A = 30

**12-1Am**  file1

A @122
B @156
C @180

**1-2Am**  file2

A

**2-3Am**  file3

A

.....

A

## Write

- → Write in WAL ⟹ $O(1)$    ⎫
- → Write in Tree Map ⟹ $O(\log N)$    ⎬ $O(\log N)$
                                              ⎭

## Read.
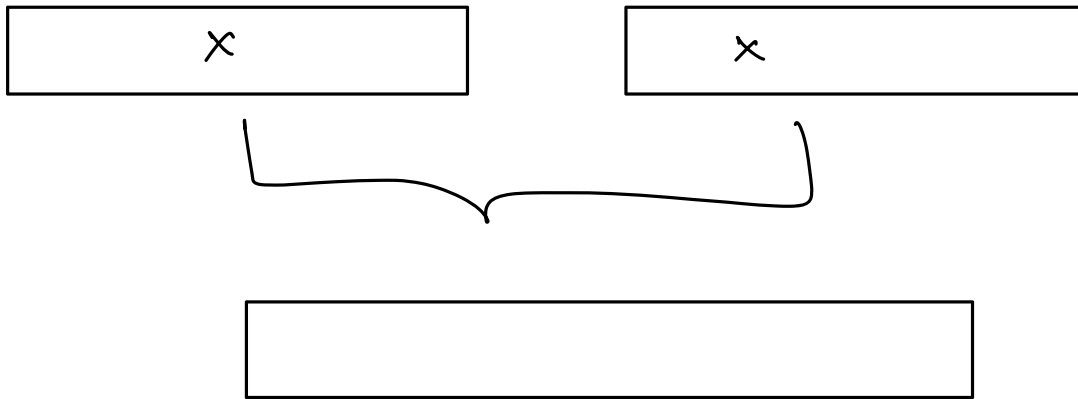
- → find the key in the TreeMap

  ↙ if found          ↘ if NOT found
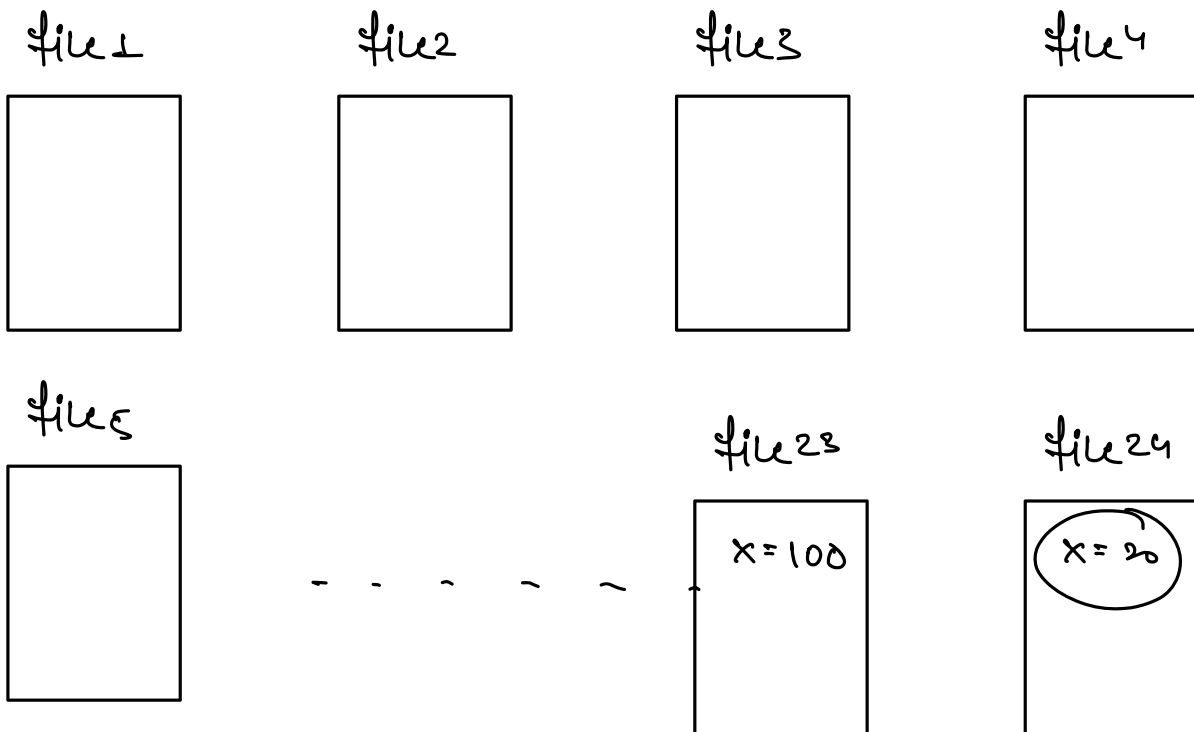
## Read TC : $O(\log N) + (\# \text{ of files}) \times \boxed{N}$

### Summary.

|               | Read TC                       | Write TC   |
|---------------|-------------------------------|------------|
| Only WAL      | $O(N)$                        | $O(1)$     |
| WAL + HM      | $O(1)$                        | $O(1)$     |
| WAL + HM + Hourly files | $(\# \text{ of files}) \times N$ | $O(1)$ |
| WAL + TM + Hourly files | $(\# \text{ of files}) \times N$ | $O(\log N)$ |

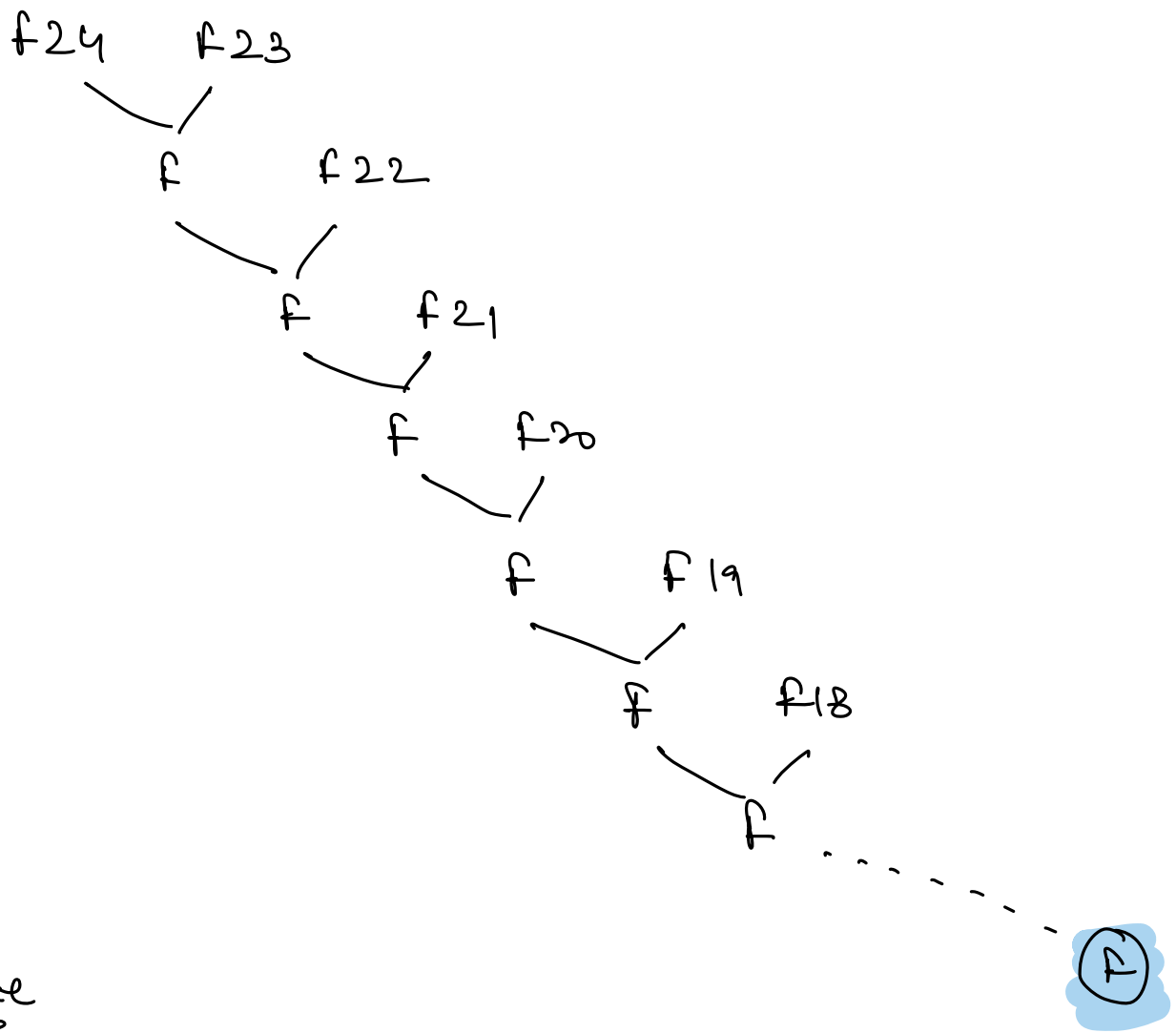## 2 Sorted Array.



# Background Script.

→ Every 24 hours, merge all the files (from file 1 to file 24) into 1 sorted file.



file 1    file 2    file 3    file 4

file 5    - - - - - - ↑    file 23    file 24

file 23: $x = 100$

file 24: $x = 20$

f24   f23

f

f22

f

f21

f

f20

f

f19

f

f18

f

⋯⋯⋯

(f)

**Write**

↳ Write in WAL ⟹ $O(1)$  ⎫
                          ⎬ $O(\log N)$
↳ Write in Tree ⟹ $O(\log N)$  ⎭
        Map

**Read.**

↳ Read from TreeMap.
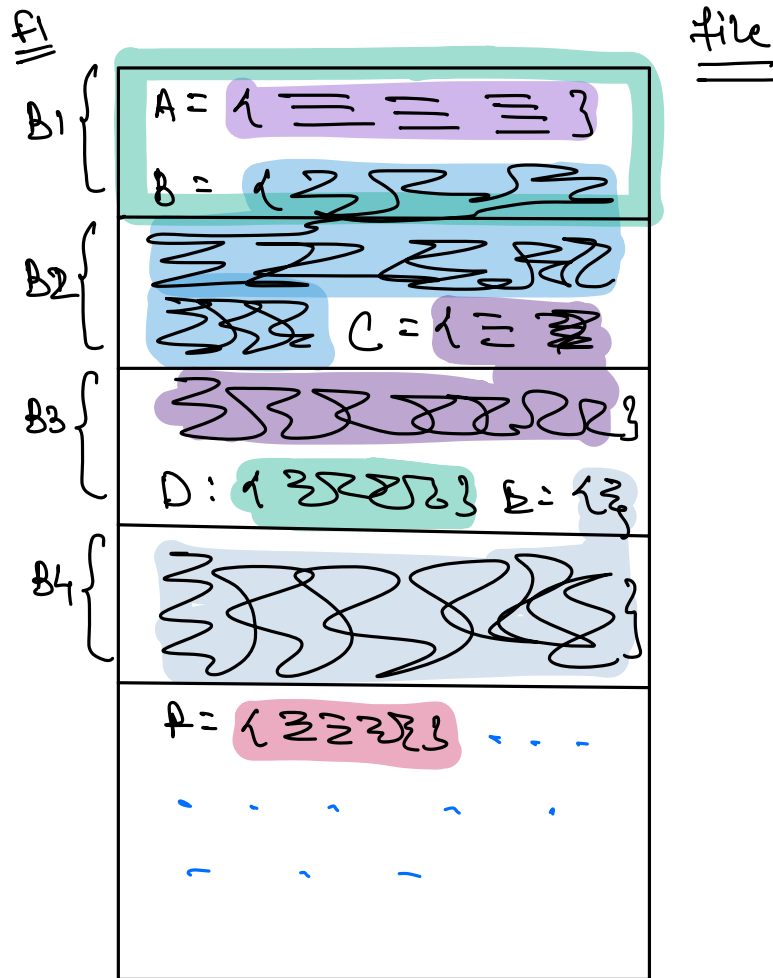
if found                    if NOT found
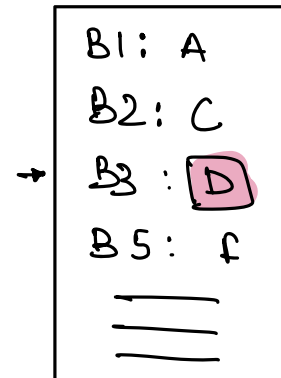
↳ Get from WAL       Iterate through all the
                     files from latest to oldest.

Read TC : $O(\log N) + (\#\ of\ files) \times N.$

Sol$^n$ : Some how if we can apply Binary Search on the files, then we are done.



f1

file

B1 { A = { ≡ ≡ ≡ }
B = { ≈ ≈ ≈ ≈ }

B2 { ≈ ≈ ≈ ≈
C = { ≡ ≈ }

B3 { ≈ ≈ ≈ ≈ ≈ }
D : { ≈ ≈ } E = { ≈ }

B4 { ≈ ≈ ≈ ≈ }

f = { ≡ ≡ ≡ } - - -

64 B

(Map)

MetaData for f1

B1 : A
B2 : C
→ B3 : D
B5 : f

final Sol$^n$

# Write

- → Write in WAL ⟹ $O(1)$ ⎫
- → Write in **Tree Map** ⟹ $O(\log N)$ ⎬ $O(\log N)$
  ⎭

# Read.

→ Read from **TreeMap.**

if found → Get from WAL

if NOT found →

⟹ Apply Binary search on the files from latest to oldest.

$$f24 \to f23 \to f22 \ldots\ldots f1 \to Ⓕ$$

$\log N \quad \log N \quad \log N \quad \log N \quad \log N$

Read TC ⟹ $\log N + (\# \text{ of files}) \times \log N$.

$= \log N.$

**Log file** + **Sorted** + **Merge** + **Trees**
(WAL)           Strings                  (Tree Maps)

$$\equiv \text{LSM}$$
$$\text{Trees}$$