# HLD - NoSQL internals

Tanmay Kacker

# Agenda

- NoSQL storage
- Motivation
- Key-value store
- WAL
- LSM (SSTables + Memtable)
- Compaction

$$\underbrace{\qquad\qquad\qquad} \times \underbrace{\qquad\qquad\qquad} \times \underbrace{\qquad\qquad}$$

| SQL | NoSQL |
|---|---|
| Structured | Unstructured |

| id | name | email |

Two

8   20
1   Tan tia Tope
—
2   Moriarty

28

56

20

1

→ ID    NAME    EMAIL
   ID      NAME

---

Key-value
   - Redis

   -   get ( key )
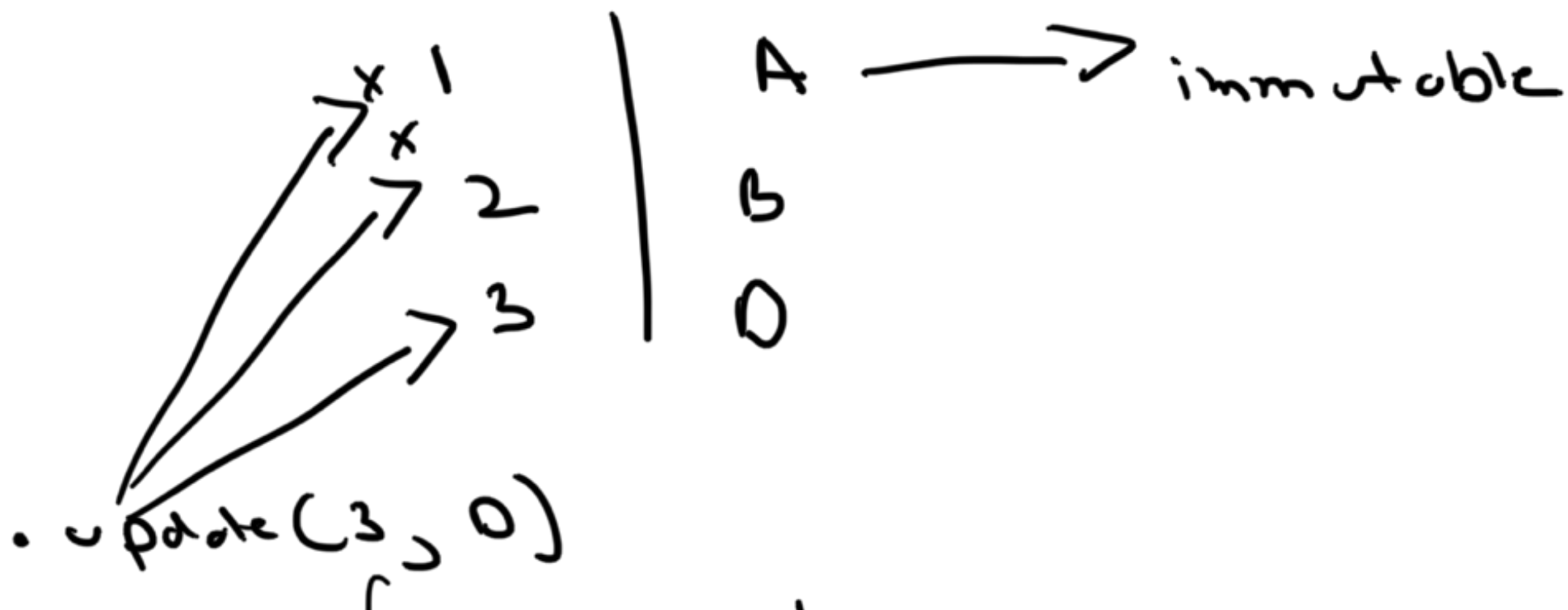   -   update ( key , value )

---

Solution 1   =   File

       S

M

$g e(1) \sim$

"1" : "Tantia Tope"

"2" : "Sherlock"

Key

①

2

Value

Tantia

Sherlock

D

read (x) → O (N)

update (x, y)

$$\longrightarrow O(n)$$

Solution 2 — $O(1)$

Write Ahead logging
 — bin log

x 1
x 2
  3

A $\longrightarrow$ immutable

B

0

• update(3, 0)

write( )

$O(1)$

1
2
3
3

A
B
e
D

1 → F

② 

1
3

A
B
C
F

$O(N)$

②

① Read complexity — $O(N)$

(2) Duplicate data

Solution 3 — read $O(1)$

WAL + index

WAL

| | |
|---|---|
| 1 | F |
| 2 | B |
| | C |
| 3 | D |
| 2 | F |

1000
1001
1002
..3
4

$O(1)$

Index

$\xi$        O(1)

1 : 1000

2 : 1001

③ → 1002

4 : 1003

(Memory)

y

get ( 3 )        O(1)

1000

Write

update ( 1, F)

get
→ Index → Address → Disk
→ O(1)

update
→ update ( 1, F)
→ Append

$\rightarrow$ Update the index

RAM - $\dfrac{8\ GB}{256\ GB}$ $\dfrac{32}{5\ TB}$

Disk -

1 Trillion - $8b + 8b$

$\boxed{16}$ -

① Flat file

read + write ~ O(N)

② WAL

- append to the end
- immutable          X1   A⟵
- write - O(1)        2    B
- read - O(N) ' - F    b    C
                          1  F⟵

③ WAL + index
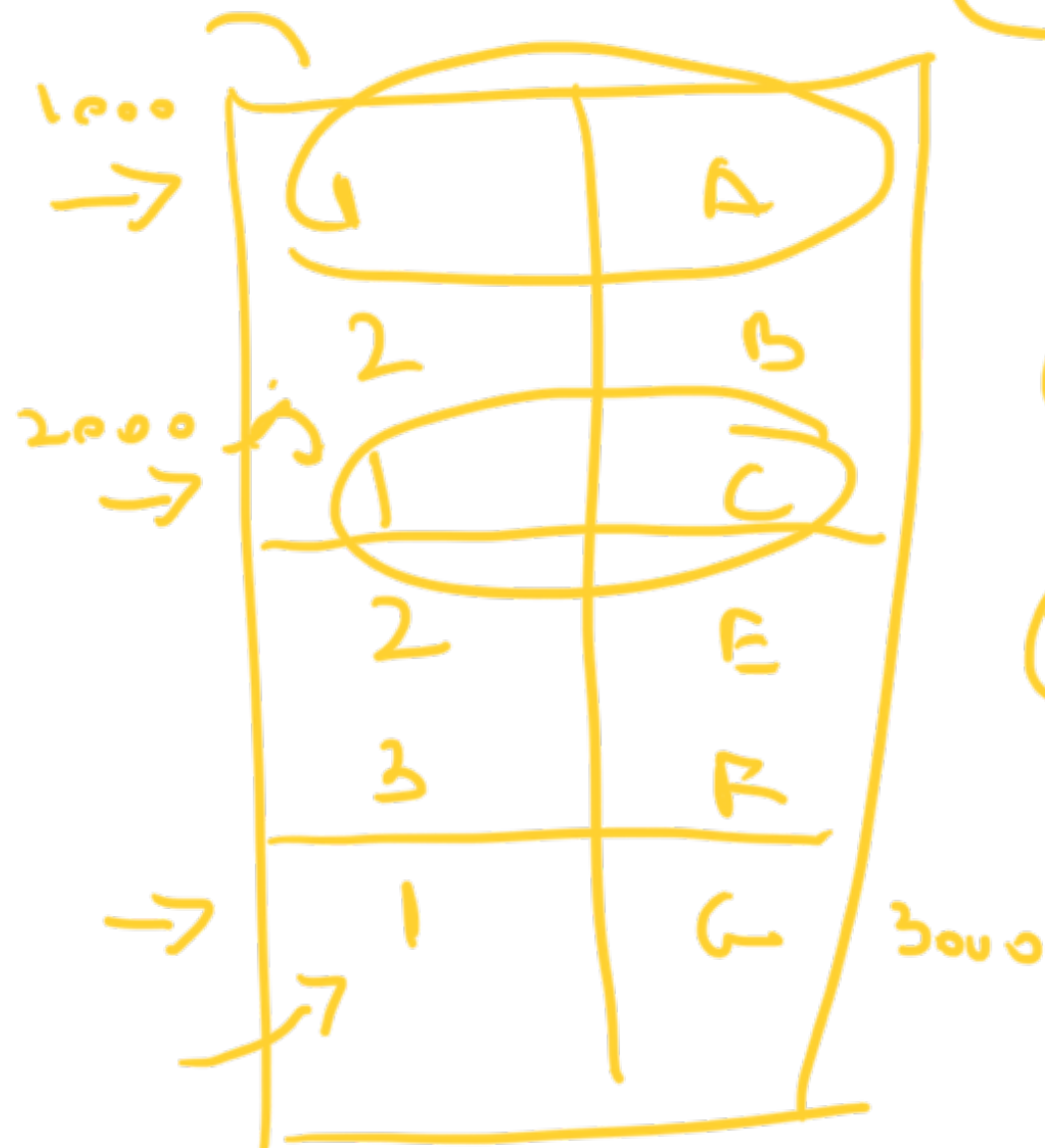
- memory index
  - key to address
- get() - index
  - address
  - seek
- update

- append to the end
- update index

1 → 1003

$\{$
1 : 1000
2 : 1001
3 : 1002

$\}$

get (1)

| 1000 | 1 | A |
|------|---|---|
| 1001 | 2 | B |
| 1002 | 3 | C |
| 1003 | 1 | F |

update (1, F)

Compaction

1000 →

| 1 | A |
| 2 | B |

2000 →

| 1 | C |
| 2 | E |
| 3 | F |
| 1 | G | 3000 |
| 7 | | |

① Periodic

② Read the file
   — Chunkify

③ De-duplication
   if (address != index)

   Store

   ⟵ | 1 : A |
       |       |

{ 1 : 3000

$1000$ ①

$2000$ ①

$3000$ ①

A

B

ε

①

$1:$ ③⓪⓪⓪

ε

②

$1 = 1000$ ✗

$1 = 2000$ ✗

$1 = 3000$

5000

1 | C

$\{$ 1 : 5000

$\}$

1 | A
2 | A
—

Fragmented

$\}$

BREAK

5 : 54 — 6 : 00
— 10 : 30

Solution S

WAL
10

Roohy

10

VAL
a

| 1 | A |
|---|---|
| 2 | b |
| 3 | c |

VAL
b

⋮

memory

100M

| | | |
|---|---|---|
| 2 | | |
| | 2: | A |
| | 3: | B |
| | 1: | C |
| 3 | | |

push →

| | |
|---|---|
| 1: | C |
| 2: | A |
| 3: | B |

Tree Map

( 1   b )

Update
→ $O(\log n)$

Get

= if key in memory
   - log(N)

= if the key is not in
      memory

MEM

DISC



S

1: A

2: B

3: C

3

log(m)

get(1)

1: D

2: B

4: A

update

get (4)

10    10    10    10    10

50    10    20

1    2

**LSM** — Log structured merge tree

① Tree Map = Latest data
  — Memory (RAM)
  ( MemTable ) — 100 MB

② Chunks of WAL file
  ( SS Table )
  ↓

# Sorted String

③ Index — Memory

ID — Start of that

Mem Table

SS Table 1

1000

$\{$

100 MB

1 : A    flush

2 : Ⓐ

3 : B

| 1 | A |
| 2 | A |
| 3 | B |

update (1, A)

→ 4 | 0

$(2, A)$

$(3, b)$

S    E

get $(2)$

get $(4)$

→ binary search

→ {
   1 : 1000
   4 : 200
}

Volatile Memory - Mem Table?
→ WAL file

(2) Delete
→ Sof delete

→ Tomb storing

SST able



del (1)

get(1)

not prexnt

③ Can a key be point of multiple SS Tables?

100

{
1

2
3

3

update (1)

(2)

(3)

Compaction
→

1

1

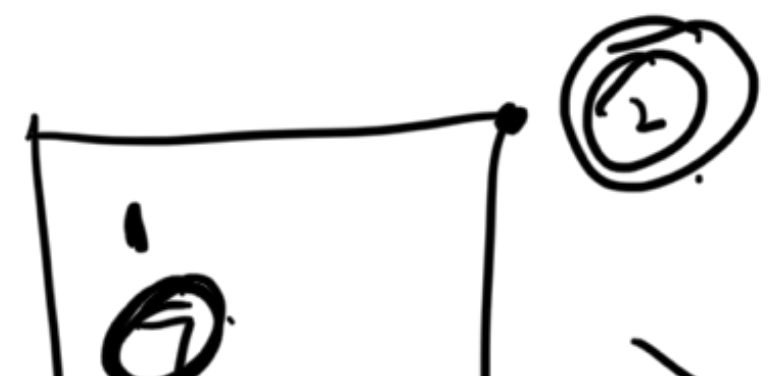Search fun bey across chunks

Bloom filter

does A Key Exist ()

False → Does not exist
True → May exist

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 0 |

hash 1 - 0

hash 2 — ②

④ h1 — 1

h2 — 3

③ h1 ① ✓

h2 ④ O ✗

②

1 — 1000

2 — 2000

1000

get()

2
5

2

2
3
6

3

7

1

5

2

2
4

6

5

3
5

SS2

100

2
4
6

200

1
2
3

index

→ 4

200

5

2

100

is itprst