

---

2018

---

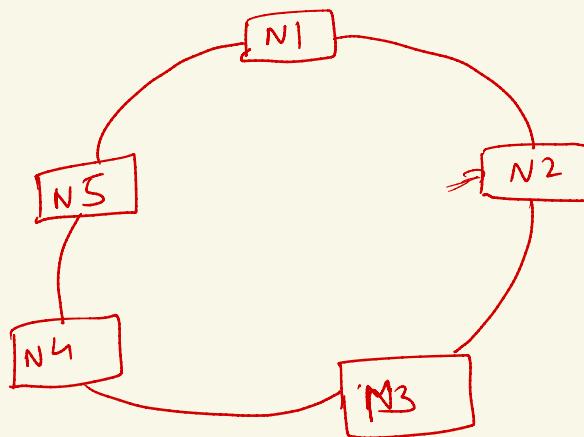
---

---

---



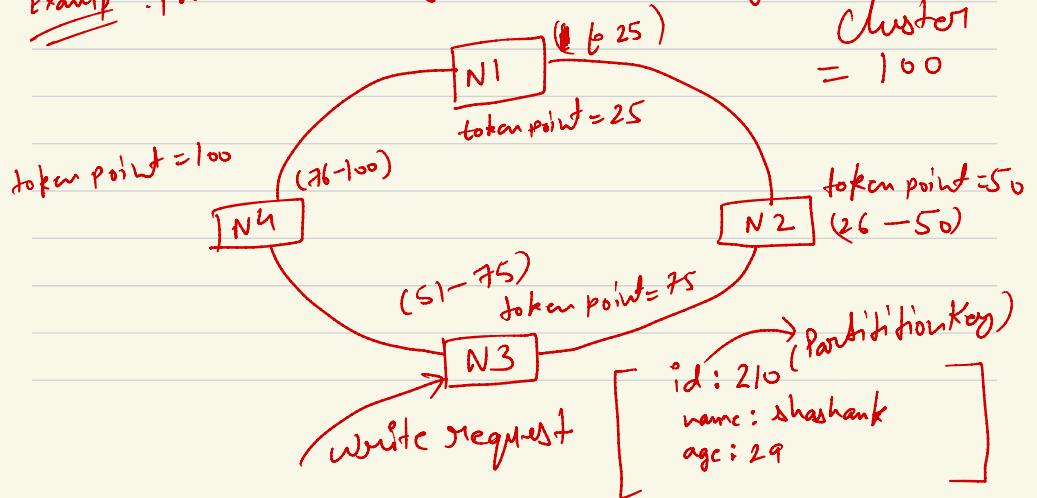
## How data is partitioned in Cassandra



Token  $\Rightarrow$  A hash value generated for a given partition key.

Token Range  $\Rightarrow$  Each node in Cassandra will be given a range of token values.

Example: Total number of token served by Cassandra



$\text{Token} = \text{hash}(\text{Partition-Key})$

$\hookrightarrow (\text{Partition-Key}) \% \text{ total\_tokens}$

$\text{Token} = \{0 \dots 99\}$

$id = 210$ , Token?

$$\text{Token} = 210 \% 100 = 10$$

Cassandra follows MurMur3 hash algorithm to generate tokens.

Hive

$\text{hash}(\text{id}) = \text{hash}(\text{Partition-Key}) \% 5$

1 Country  
2 INDIA

$$= \text{hash}(\text{INDIA}) \% 5 = 5 \% 5 = 0$$

2 USA

$$= \text{hash}(\text{USA}) \% 5 = 0$$

3 INDIAN

$$= \text{hash}(\text{INDIAN}) \% 5 = 3$$

4 INDIA

5 CHINA

$id = 1, \text{Partition\_Value} = 0$

6 CHINA

$id = 3$

B1  
0

B2  
1

B3  
2

B4  
3

B5  
4

Collision  $\Rightarrow$  if same hash value  
got generated for  
different values

In Cassandra Token Range

between

$-2^{63}$  to  $(2^{63}-1)$

8 Bytes long Integer  
64 bits Value  
 $(0-2^{64})$

$[-10259642382 \text{ to } 10259642381]$

Total Token Range

Let say we have 4 Node in cluster

for each Node =  $\{0 - 99\} / \text{Number of Nodes}$

token Range

OR

=  $(\text{Total Token Range}) / \text{num\_tokens}$

for each

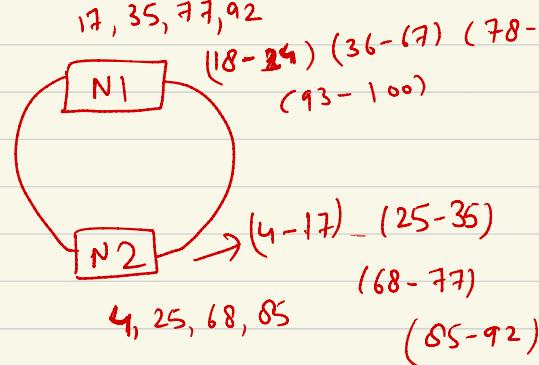
node

num\_tokens = Number of Ranges each Node can handle

Node = 2

num\_tokens = 4

token\_range = 0-100



## Consistency Levels in Cassandra

Q.) based on CAP theorem Cassandra lies where?

↳ AP (Available, Partition Tolerance)

Q.) in Cassandra <sup>(Yes)</sup> Can we impose consistency?

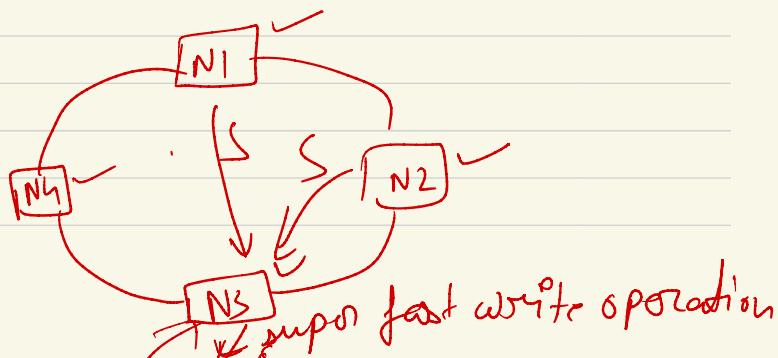
Q.) What will be the impact if we apply consistency? (Speed or performance)

What happen for consistency?

① Write → if consistency not maintained the write operation will fail

② Read → if consistency not maintained then Read will fail

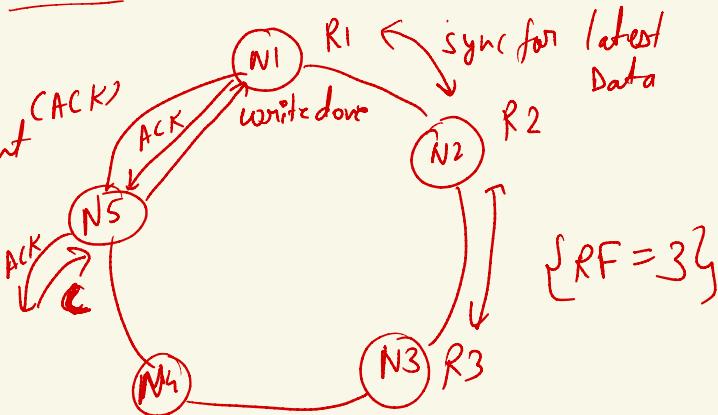
RBMS is there any intermediate state?



Write request (S)  
(forcefully write it for ALL)

### Consistency for write operations

↳ Based on acknowledgement from replicas



→ ONE: it needs acknowledgement from only one Replica Node.

In this case data will be synched asynchronously (from backend) on remaining replicas.

→ TWO: it needs ~~one~~ acknowledgement from 2 replica nodes.

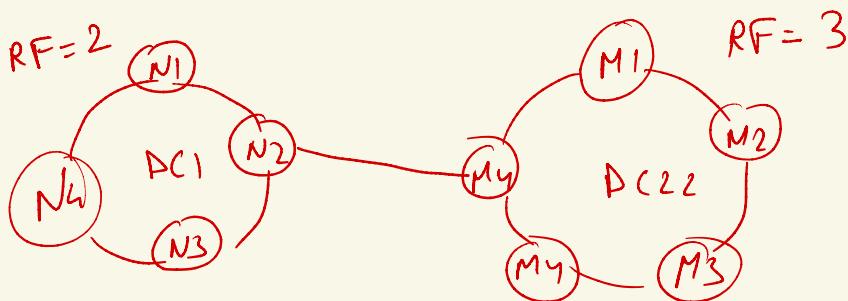
(majority)  
Quorum Calculation formula =  $\frac{\text{sum of Replication factor}}{2} + 1$

if we have 100% of something then what least value we need for majority? (51%)

ans 6+ 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12

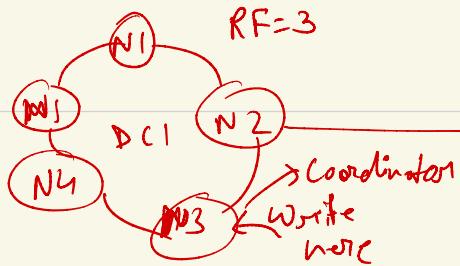
(10)  $= \left(\frac{n}{2} + 1\right) = \left(\frac{10}{2}\right) + 1 = 5 + 1 = 6$

: QUORUM  $\rightarrow$  needs acknowledgement from  
51% or majority of replica nodes  
across all the data centers.



$$\begin{aligned} \text{Value of QUORUM} &= (2+3)/2 + 1 \\ &= \left(\frac{5}{2}\right) + 1 \\ &= 2 + 1 \\ &= 3 \end{aligned}$$

: LOCAL - QUORUM  $\rightarrow$  need acknowledgement from  
51% or a majority of replica nodes  
just within the same data center  
as the coordinator.



LOCAL QUORUM

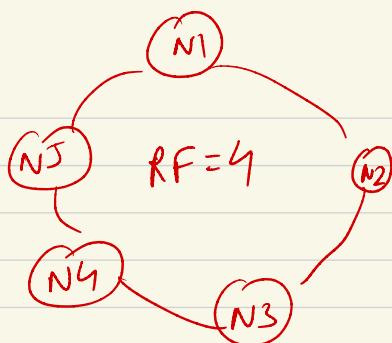
$$\text{for } DC_1 = \frac{3}{2} + 1 \\ = 2$$

LOCAL QUORUM

$$\text{for } DC_2 = \frac{4}{2} + 1 \\ = 3$$

: ALL  $\rightarrow$  needs acknowledgement from all replica nodes

so it need ACK  
from all  
4 replicas

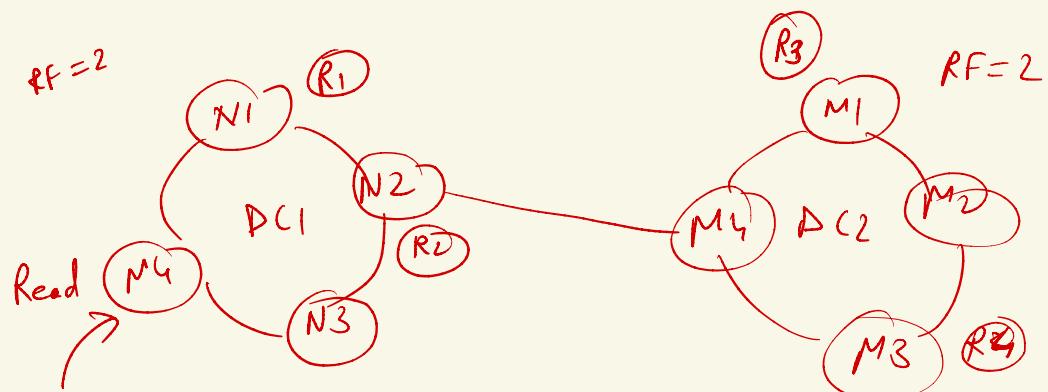


Ex:- Updates done using Consistency Quorum  
Set total\_purhase = 100 where cust\_id = 4;

Consistency for Read Operation

: ONE  $\rightarrow$  Only one replica node returns the data.

: QUORUM  $\rightarrow$  returns the data from a quorum with the most up-to-date data.



$$\{ \text{last update} = 200 \} \quad \text{Quorum} = 3$$

- $\hookrightarrow$  Output from R1  $\rightarrow$  Last update happened on 20-NOV-2022 10:00PM
- $\hookrightarrow$  Output from R2  $\rightarrow$  20-NOV-2022 9:00PM
- $\hookrightarrow$  Output from R3  $\rightarrow$  20-NOV-2022 11:00 PM

: LOCAL\_QUORUM  $\rightarrow$  Returns a result from

LOCAL\_QUORUM with latest data  
and follow coordinator node for local quorum value.

: ALL  $\rightarrow$  reads data from all Replica nodes and

Returns latest record based on timestamp.

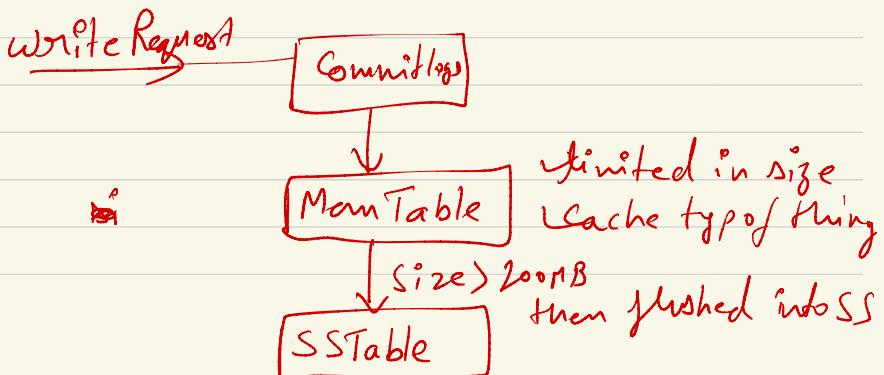
Ex: Select total purchase from Sales  
using consistency ALL where customer-id=5;

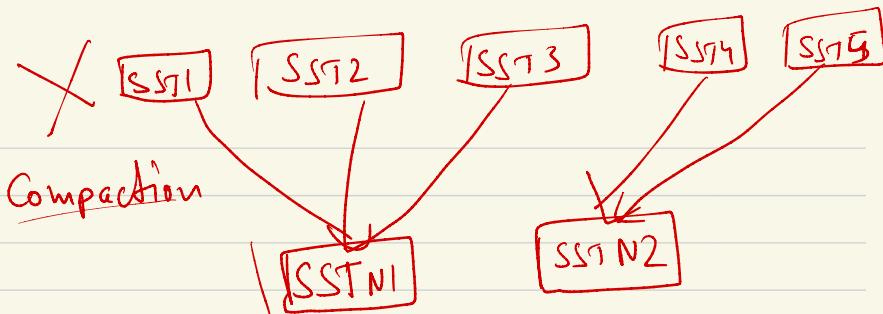
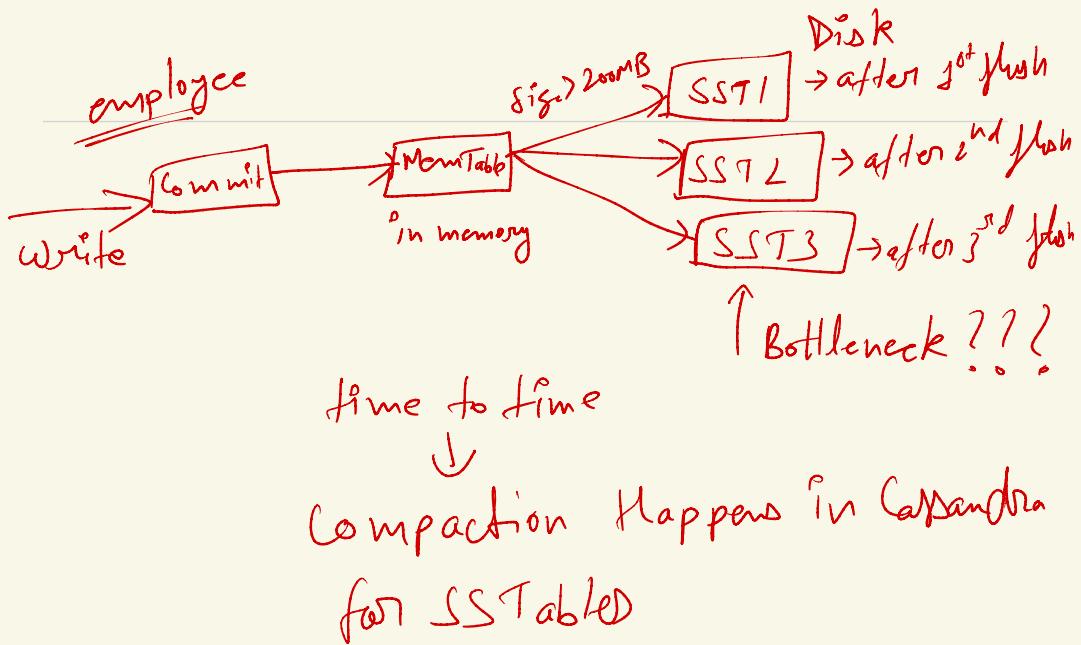
### Write Path in Cassandra

→ Commit log ⇒ transaction logs will be written here for crash recovery. Non-volatile, so if power failure happens still we can recover.

→ Memtable ⇒ memory cache to store the copy of the data.

→ SSTable ⇒ the final ~~destination~~ destination where actual data gets stored in the disk.





## Write Path

- ① Cassandra appends writes in commit log for durability.
- ② Data then will be written in a in memory Cache known as Memtable
- ③ The memtable store writes in a sorted order until reaches the

memory threshold value.

- ⑤ Then it is flushed to sorted string table called as SSTable. Write operations are atomic on Row Level, means values for each column of that row will be written/updated or none will be written.

### Read path in Cassandra

① Cassandra first checks if the data is available in the memtable

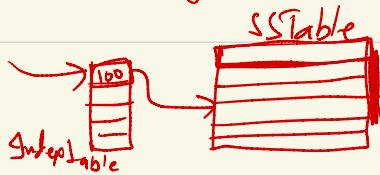
② if not found then data will be read from SSTable.

③ To optimize the reads

    → Cassandra uses Bloom filter technique to minimize the data scan.

    → Cassandra uses indexing on SSTable

        ↓  
Metadata for Actual Data



→ Compaction also helps to improve reads

Types of Reads coordinator can send to Replica?

↳ Direct Read Request

↳ Digest Request

↳ Background Read Repair Request

Direct Req:

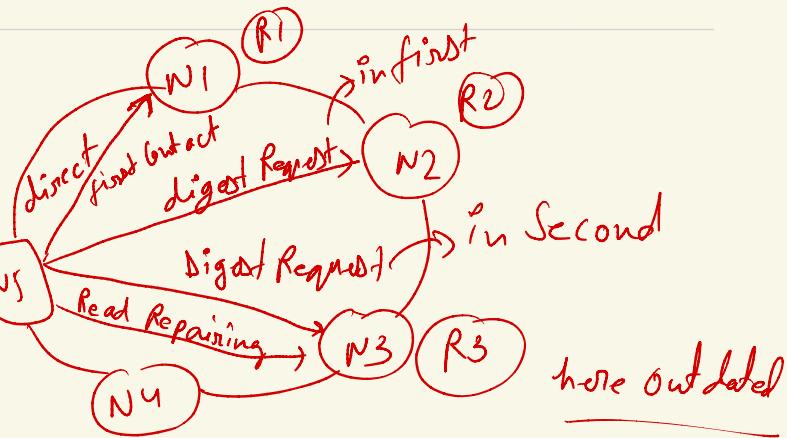
Coordinator node contacts one Replica node, after that digest request will be sent to the remaining replica nodes based on consistency level. The digest request will check if data on the replicas are up-to-date or not. Then Coordinator will send the digest request to remaining nodes.

If any replica node have out of date data, a background read repair request will be sent. Read repair request ensure that the requested row(data) is made consistent on all replicas involved in a read query.

$$RF = 3$$

$$\begin{aligned} RF - CLR \\ 3 - 2 = 1 \end{aligned}$$

$$\begin{aligned} (curr\_id = 200) \\ \text{Read } 3 \text{ opn} \end{aligned}$$



Select with consistency Quorum

= 2