

design Key Value DB : (Dynamo DB)

Goals : ① Scalability

② Decentralization

③ Eventual consistency.

- Steps:
- ① Partition
 - ② Replication
 - ③ Get & Put ops
 - ④ Data Versioning
 - ⑤ Gossip Protocol
 - ⑥ Merge tree

① Partition: we introduce n nodes and manage data with the help of consistent Hashing, this would lead to Scalability.

decentralization: even if there are failures in the nodes, the system should not get down, even in failures system should fulfill the request for all incoming keys.

② Replication: Keeping a copy of the data to other servers in order to fallback in case of any node failures

N = Replication count

1 is the main node called coordinator
 $N-1$ nodes in clockwise dir. where the copy of data would be made.

each Key Range

Maintains

Preference List

[1-50]

Key Range

S1

S2

S3

Preference

3 Replic

\uparrow

N

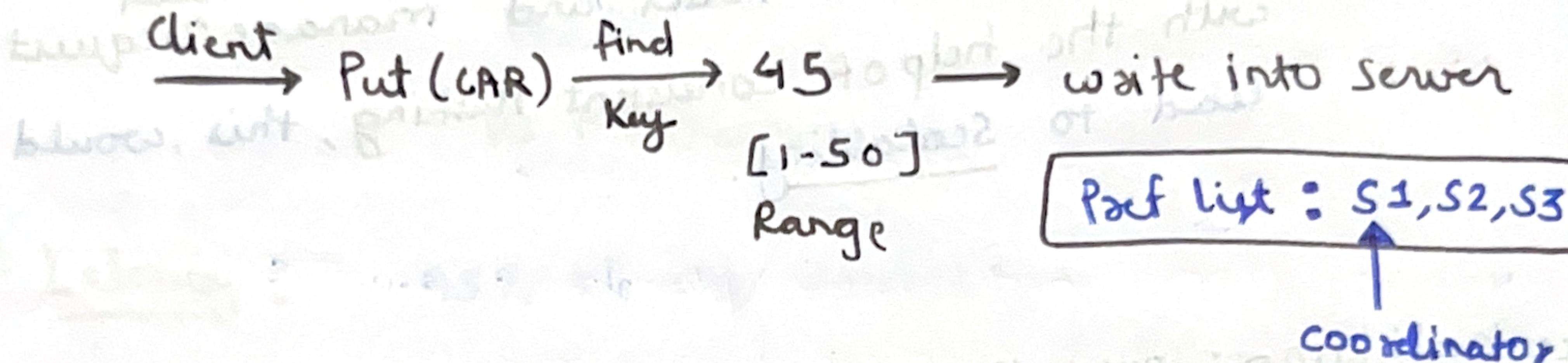
coordinator

⇒ it may not choose the next $N-1$ servers immediately meant to coordinator, Replication nodes are also chosen based on regional preferences as well.

& skips virtual nodes.

③ Get & put ops:

Put



We send success response to client when;

no. of servers to read from while reading (get req.)

$$R + W > N \quad \text{no. of replicas}$$

no. of servers to wait for the async copy to complete

get we already have prof. list for each key Range

there is a LB

request



C.H.

partition aware LB

generic LB

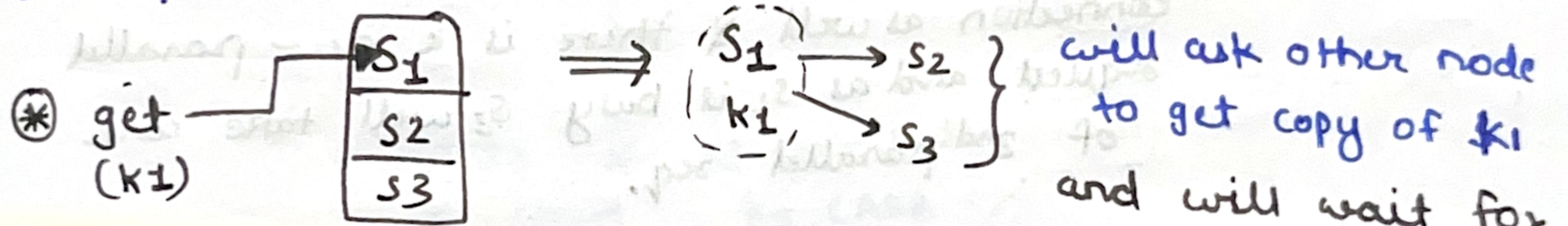
generic LB: get or put request comes to any node, now each and every node is aware of all Key Range's Pref. List, so it will take a hop and send request to coordinator. if coordinator is down next one will take the charge from assigned node. (same for get)

- hopping is the main issue, high latency

partition aware : here the LB is aware of the Pref. List and partition so request will go to coordinator directly.

- low latency because of minimized hopping.

* control will come to coordinator in case of all the type of LB.



\Rightarrow Once S_1 (coordinator) get all the replicas (R); it would respond back to client for get Req.

why do we need R & W for above operation; before declaring the operation as success?

\Rightarrow Data Versioning will answer that.

④ Data Versioning :

These might be a chance of failures, right?

* So it might be possible that due to network failure or any other failure, the keys can get inconsistent as they might not be (yet) updated in other nodes.

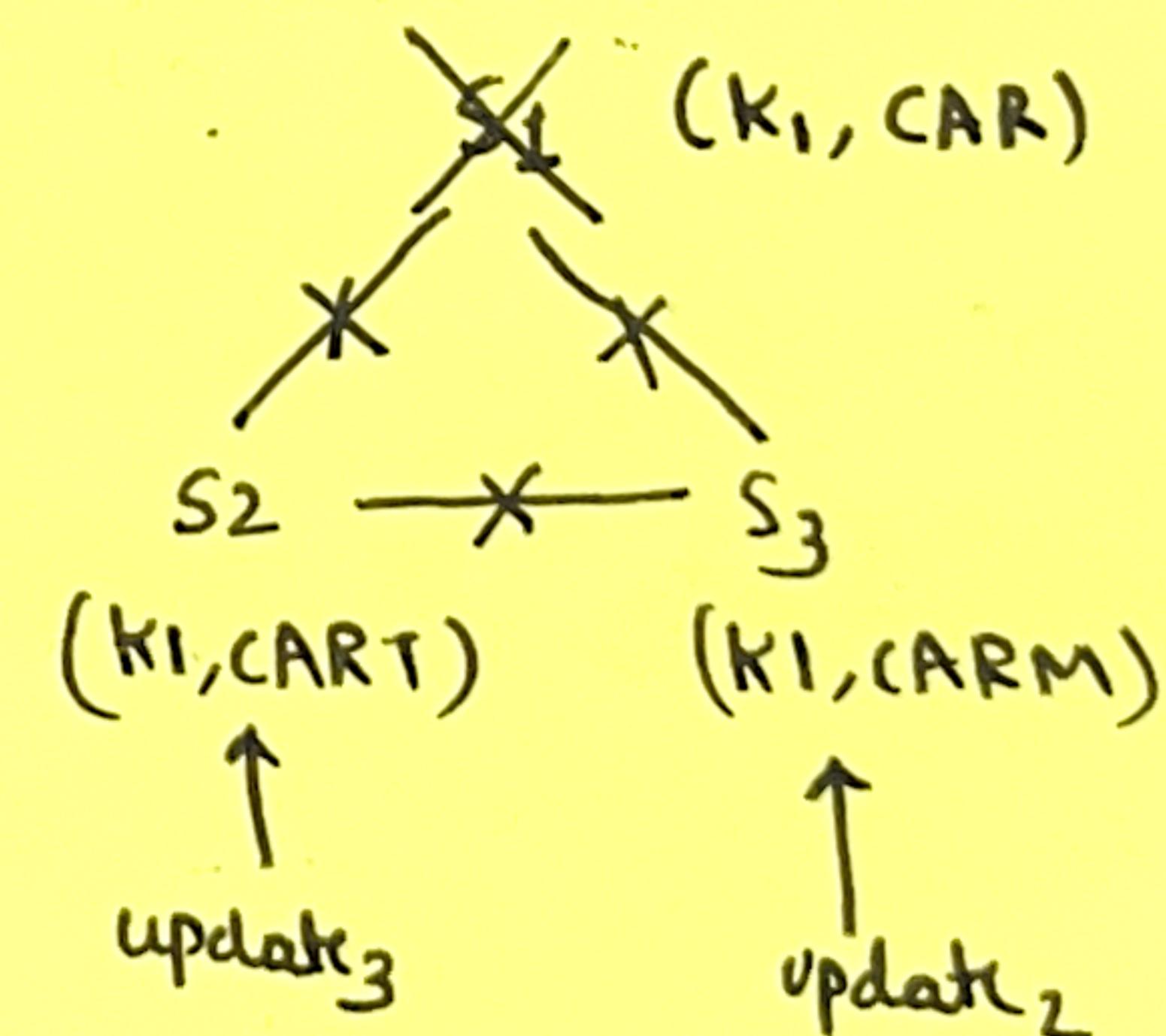
[For example, on second update of same key, if our main coordinator S_1 goes down, go it would contain non-updated entries.]

→ more complexity arises when S_2 and S_3 lose connection as well & there is a new parallel request and as S_2 is busy S_3 will take care of 2nd parallel seq.

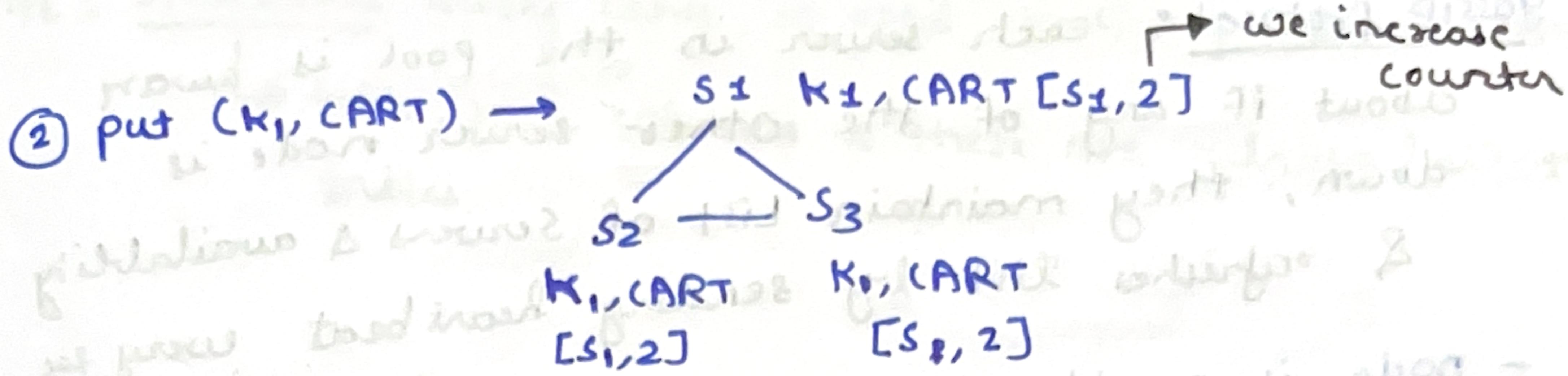
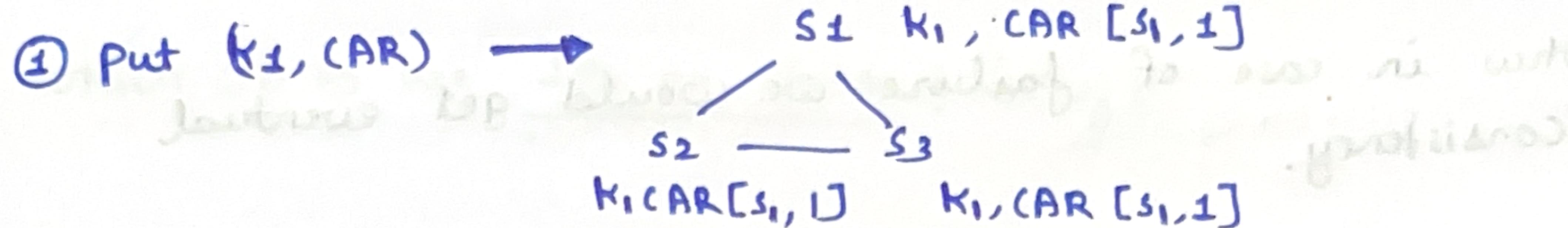
To resolve this, we use something called Vector clock.

Each entry in server for key contains additional information

List of [server, counter]

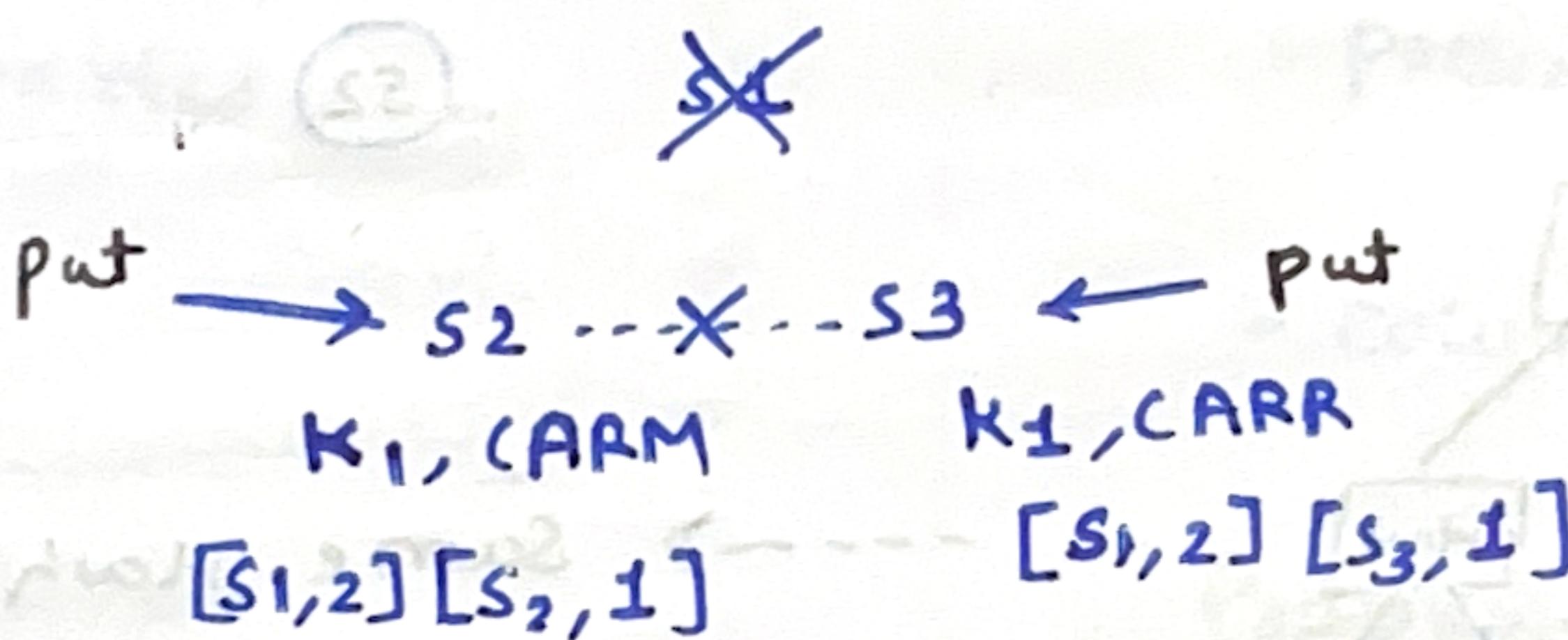


all 3 are inconsistent.

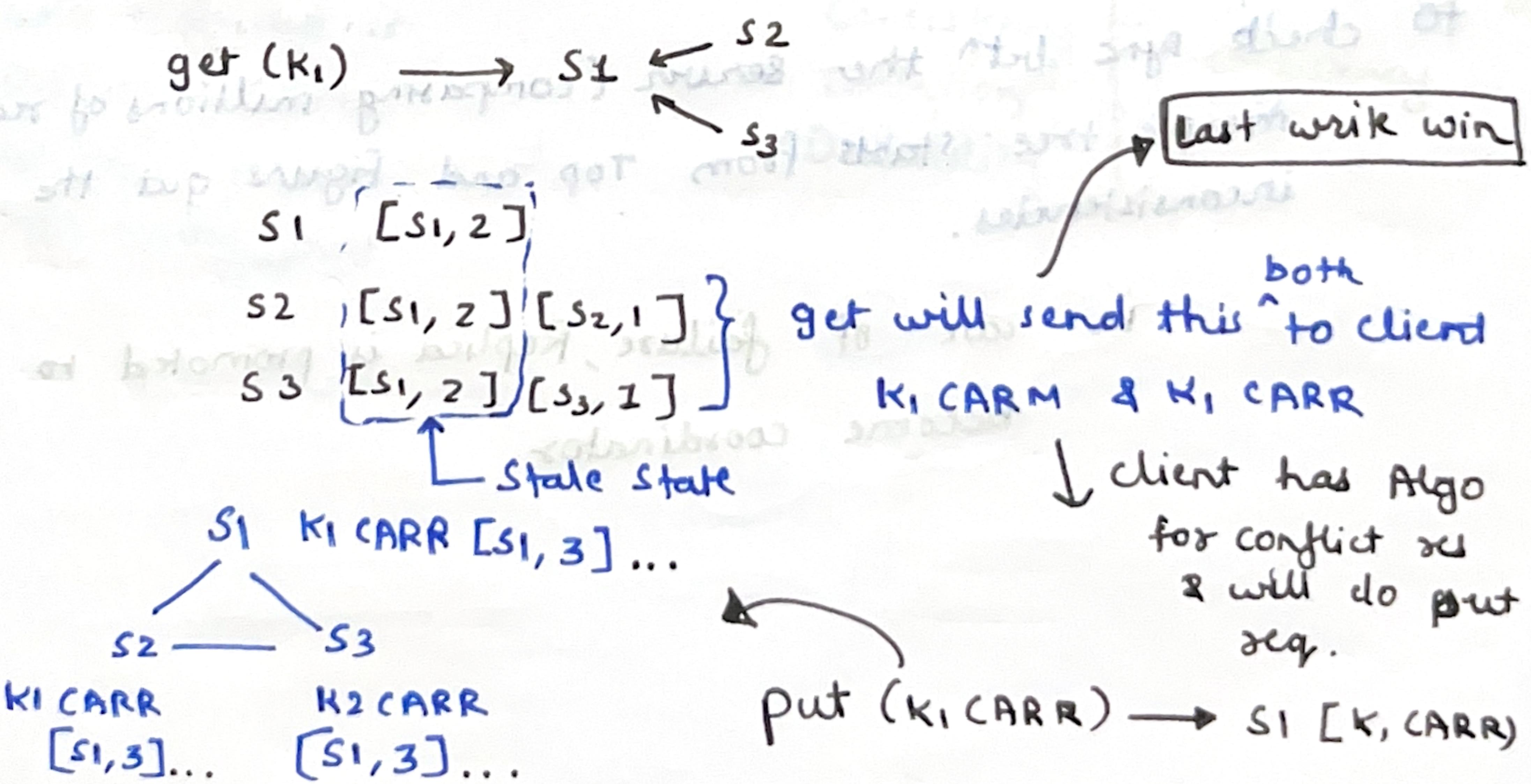


Now node S1 goes down & therefore partition b/w S2 & S3.

③ put (k, CARM) & $\text{put}(k_1, \text{CARR})$ in parallel:



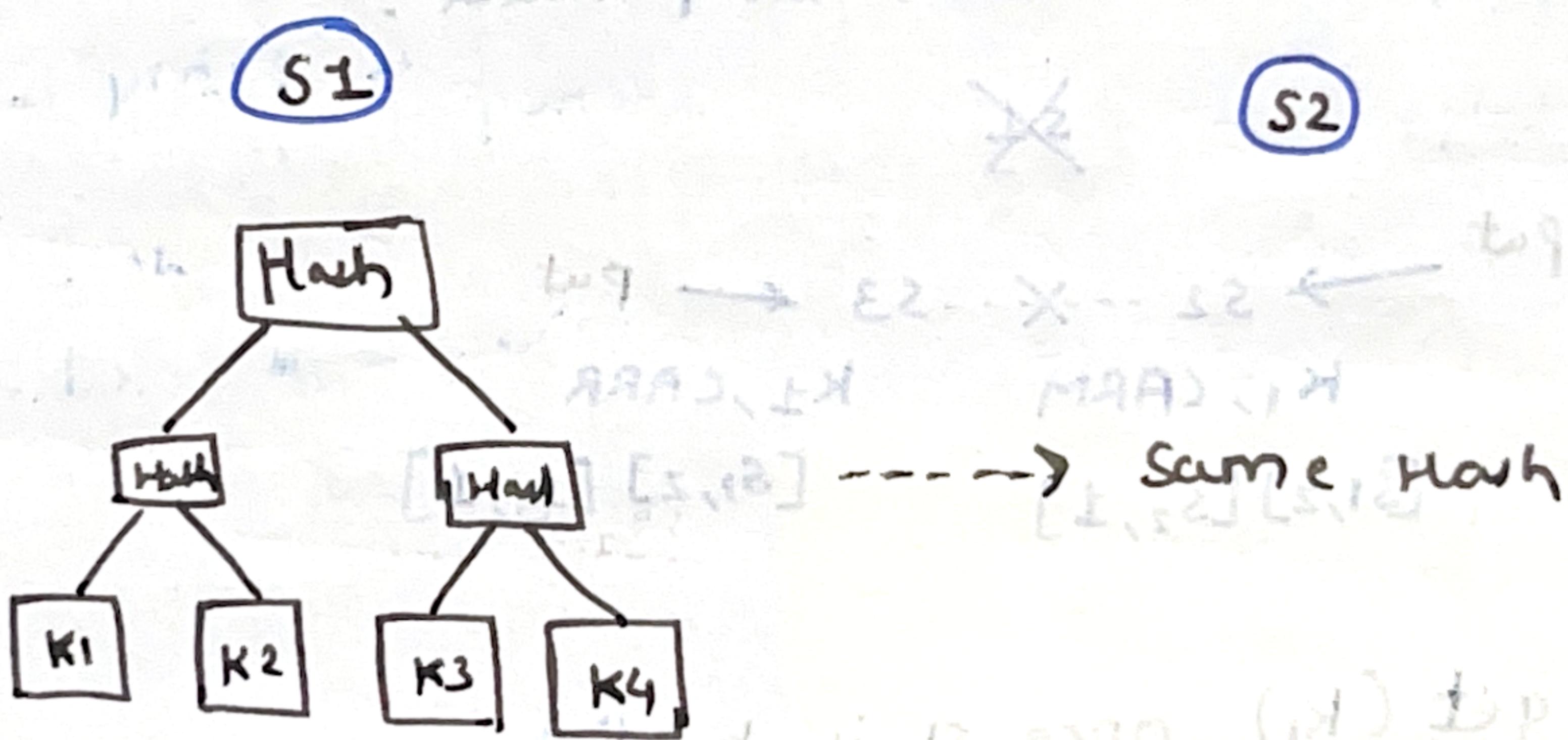
Now if we get (k_1) once S1 is back as well;



thus in case of failure we would get eventual consistency.

- ⑤ Gossip protocol: each server in the pool is aware about if any of the other server node is down, they maintain list of servers & availability & refreshes list by sending heart beat every sec.
- Node is considered down if more than one servers notices & reports failure heartbeat report.

⑥ Merkle Tree:



To check sync betw the servers (comparing millions of record) Merkle tree starts from Top and figures out the inconsistencies.

In case of failure, Replica is promoted to become coordinator.