

Storage System Internals

Ankur Shrivastava (Pricing at GoJek)

January 17, 2022



Who am I?

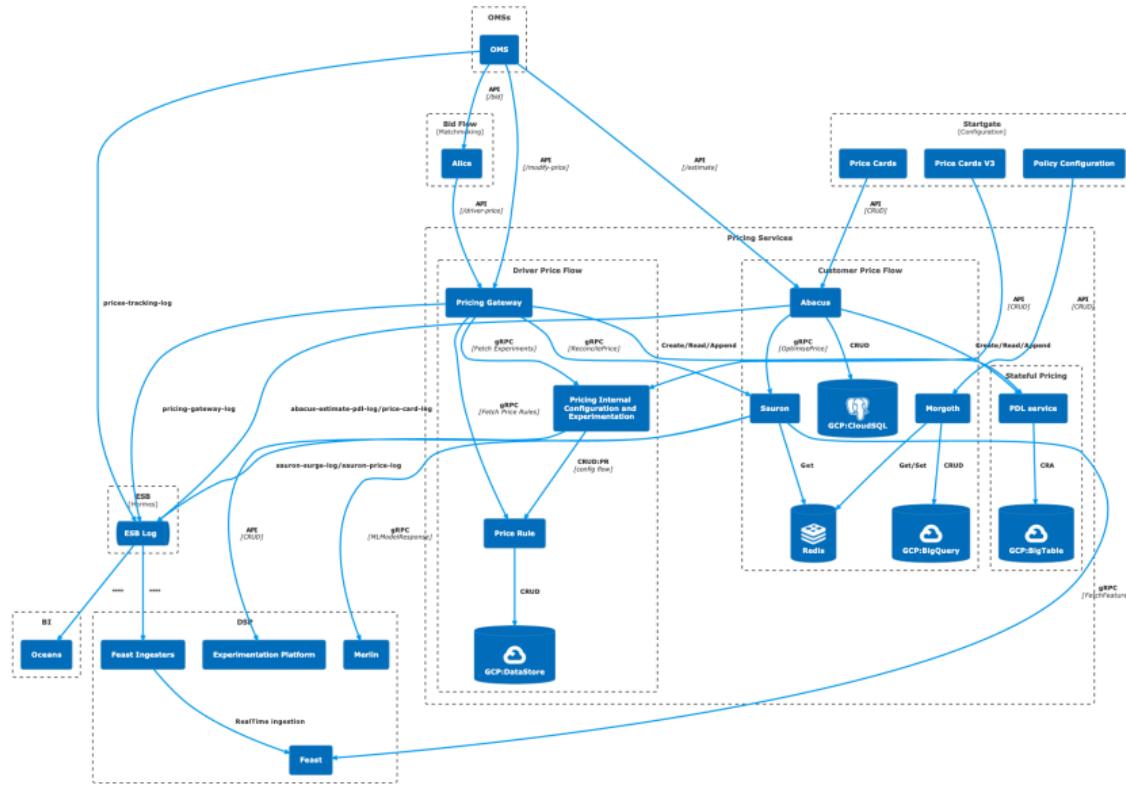
- Ankur Shrivastava
 - Principal Software Engineer
 - Team Lead of Pricing (GoJek)
 - Worked at Skyscanner, Carousell, Flipkart, Amazon, Zynga, etc.



Context



Pricing Architecture (Simplified)

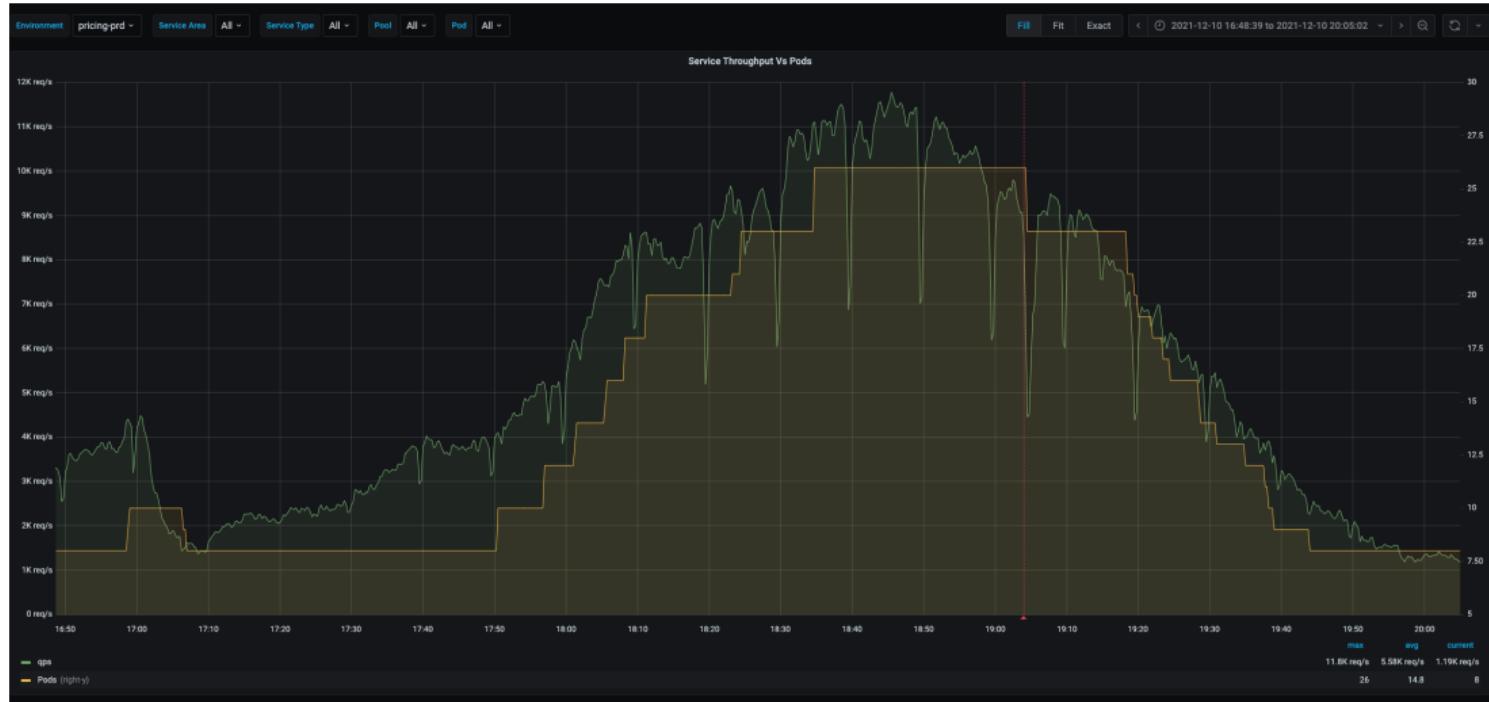


Pricing Architecture

- Pricing stack is entirely on K8S (PDG Managed)
- We have no VMs (other than GKE)
- No state in K8S
- All services autoscale with traffic
- Deployments triggered by CI and take ~5 mins max.
- All service to service calls over gRPC
- No Single Point of Failure (DB, Cache, LB, Pods have automated recovery)
- Read More at <https://go-jek.atlassian.net/l/c/z0wxUr9R>



Pricing Gateway Autoscaling



Dynamic Driver Pricing system processes 11.8k individual driver price requests per second (708,000 rpm) and can handle sudden 5x increase in traffic.



gojek

Zero downtime cluster/project movement



We can migrate K8S cluster / GCP Project without any impact on latency.



New Stateful Pricing Service

Requirements

Based on <https://tools.ietf.org/html/rfc2119>

- Need to store Pricing PDL

Must Have

- Read and Append API for PDL data exposed by Pricing
 - PDL can only be created by pricing
 - PDL components can never be updated only appended
- OMSs should communicate price changes back to pricing
- Ability to lookup PDL data from a given PDL ID



New Stateful Pricing Service

- Get PDL (read)
 - Food: 300 qps * number of service types (=1) * number of payment types (=1) => 300qps
 - Transport: 350qps * number of service types (=1) * number of payment types (=2) => 700 qps
- Append PDL (write)
 - Food: 600 qps * number of service types (=1) * number of payment types (=5) => 3k qps
 - Transport: Need to be done (@Amit Kumar Rai - Jan 2022)
- Estimate API (write)
 - Food: 600 qps * number of service types (=1) * number of payment types (=1) => 600 qps
 - Transport: 350qps * number of service types (=10) * number of payment types (=2) => 7k qps

- Need to store Pricing PDL
- Writes \Rightarrow 11,000+ per second
- Reads \Rightarrow 1,000+ per second
- Writes > Reads
(by a factor of 11)



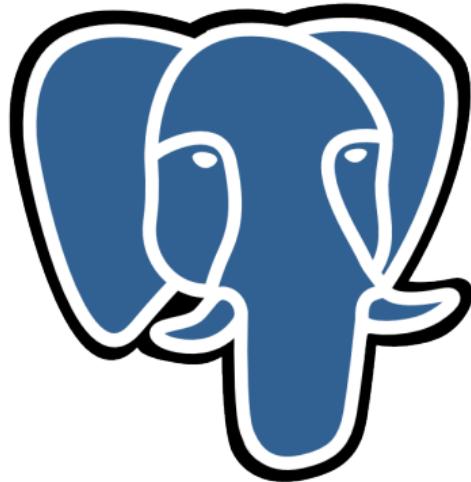
New Stateful Pricing Service

- Get PDL (read)
 - Food: 300 qps * number of service types (=1) * number of payment types (=1) => 300qps
 - Transport: 350qps * number of service types (=1) * number of payment types (=2) => 700 qps
- Append PDL (write)
 - Food: 600 qps * number of service types (=1) * number of payment types (=5) => 3k qps
 - Transport: Need to be done (@Amit Kumar Rai - Jan 2022)
- Estimate API (write)
 - Food: 600 qps * number of service types (=1) * number of payment types (=1) => 600 qps
 - Transport: 350qps * number of service types (=10) * number of payment types (=2) => 7k qps

- Need to store Pricing PDL
- Writes \Rightarrow 11,000+ per second
- Reads \Rightarrow 1,000+ per second
- Writes $>$ Reads
(by a factor of 11)
- Ingestion: \sim 400M keys per hr, \sim 40GB data per hour
- QPS will continue to grow
- Each price will only be accessed \sim 10 times over its lifetime
- We will never access this data after a few hours



Stateful Pricing - Storage



- PostgreSQL+Redis for the win!

PostgreSQL

Stateful Pricing - Storage

 **Severan Rault** Feb 4th, 2021 at 7:18 PM
I wrote this earlier today in a slack channel. What is your guys view ? Is it time to phase out all our Postgres(PG) prod DBs for caching / transient datastore scenarios ? (edited)

26 2 8 19 3 6

122 replies

 **Meghajit Mazumdar** 11 months ago
Does it mean we make our services serverless ? (functions) ?

 **Severan Rault** 11 months ago
I think my point is more around the fact that a relationalDB is never a hard requirement, a KV store is enough to implement anything, and its a much lower burden to scale and make operationally robust

 **Severan Rault** 11 months ago
I think we should aspire to simplify and standardize along few , well mastered patterns that scale well:
Lambda/ziggurat , KV (redis) , Kafka ...

 **Darshan Chaudhary** 11 months ago
IMO few things will miss with redis - enforcement of data integrity, data view from multiple tables (joins), concurrent access

 **Severan Rault** 11 months ago
u can design around this "limitations" 😊

 **Severan Rault** 11 months ago
thing is eventually, you have to give up on the relational datastore because it won't scale

 **Severan Rault** 11 months ago
and by scale, I am not just referring to performance but also engineers working concurrently on functionality linked to it

 **Mitesh Patel** 11 months ago
Just curious, when you say PG is a ticking timebomb 🕒 . Can you please elaborate a bit.

 **Severan Rault** 11 months ago
it will eventually explode in your face 😊 usually a few years down the road as the disc space slowly dwindles because the cron job that clean up the DB , was paused as it was slowing down perf due to re-indexing... type of issues

 **Ashar Fuadi** 11 months ago
| eventually, you have to give up on the relational datastore because it won't scale

are you saying at the current Gojek's scale, PG is definitely not a good storage solution? is there any references we could read / look up to back up this claim? 😊

 **Severan Rault** 11 months ago
PG is not designed for gojek scale, heck no relational DB is

 **Severan Rault** 11 months ago
PG is suitable for fast prototyping and quick iterations , but once you want speed and reliability , you flatten the data structure and go KV. We doing that for a while now. I guess I am trying to get all of you to think long term

 **Severan Rault** 11 months ago
and also to be clear, I am referring of the usage of PG for caches and transient informations that are in the critical path of our users flow.

1 3 6

 **Severan Rault** 11 months ago
PG for long term storage of seldom accessed data is appropriate I suppose (edited)

● PostgreSQL + Redis for the win!

Stateful Pricing - Storage

```
23 $ bash misc/loadtest/grpc-write.sh
24 stateful-pricing-grpc-int.pricing.golabs.io
25 Summary:
26   Count: 696028
27   Total: 300.00 s
28   Slowest: 659.14 ms
29   Fastest: 8.95 ms
30   Average: 20.99 ms
31   Requests/sec: 2320.09
32 Response time histogram:
33   8.953 [1] |
34   73.972 [690581] |#####
35   138.991 [4796] |
36   204.010 [333] |
37   269.028 [140] |
38   334.047 [71] |
39   399.066 [44] |
40   464.085 [8] |
41   529.104 [8] |
42   594.122 [8] |
43   659.141 [12] |
44 Latency distribution:
45   10 % in 13.94 ms
46   25 % in 15.55 ms
47   50 % in 17.95 ms
48   75 % in 21.85 ms
49   90 % in 29.48 ms
50   95 % in 40.77 ms
51   99 % in 67.58 ms
52 Status code distribution:
53   [OK]          695978 responses
54   [Unavailable] 50 responses
55 Error distribution:
56   [50]  rpc error: code = Unavailable desc = transport is closing
58 Cleaning up file based variables
60 Job succeeded
```

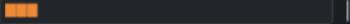
- PostgreSQL + Redis for the win!
- BigTable might work better
- Lets build a service and load test (integration)

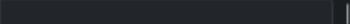
Stateful Pricing - Storage

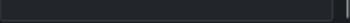
 **Polly Results** ✖

... **Results**

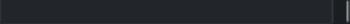
Have you worked with HBase/Cassandra/BigTable before ?

I have used them in Production
 | 10% (1)
@Ankur

I have used them Personally (non production)
 | 0% (0)

I am aware of how they work but haven't used them
 | 0% (0)

I have heard about them but haven't used them
 | 90% (9)
@Mauro, @siddhant.tiwari, @Akshay, @nsp, @leyi, @kevin.fan, @Qing, @derek, @andre

Wait, what is HBase/Cassandra/BigTable
 | 0% (0)

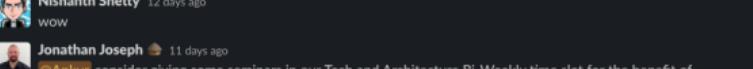
Total Votes: 10

- PostgreSQL + Redis for the win!
- BigTable might work better
- Lets build a service and load test (integration)
- Lets ask if the rest of the team have used it

Stateful Pricing - Storage

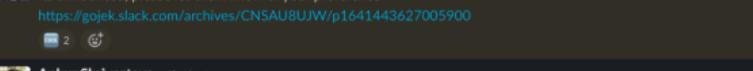
Also sent to the channel

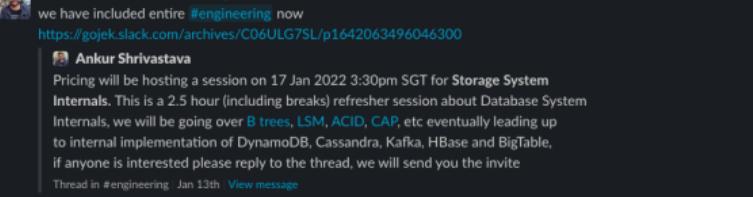
Ankur Shrivastava 12 days ago
^ the results are a bit scary
how will you guys feel if I do ~2hr session in the next few week going over B+ trees, LSM, ACID, CAP, etc eventually leading up to BigTable internal implementation ?


Nishanth Shetty 12 days ago
WOW


Jonathan Joseph 11 days ago
@Ankur consider giving some seminars in our Tech and Architecture Bi-Weekly time slot for the benefit of Marketplace? 

Ankur Shrivastava 11 days ago
I will send the invite out to Marketplace (with priority to pricing) and since I would need at-least 2 hours i was thinking of running it on a Friday
Also sent to the channel

Ankur Shrivastava 10 days ago
Its announced, please let them know of your preference -
<https://gojek.slack.com/archives/CNSAU8UJW/p1641443627005900>


Ankur Shrivastava < 1 minute ago
we have included entire #engineering now
<https://gojek.slack.com/archives/C06ULG7SL/p1642063496046300>


- PostgreSQL + Redis for the win!
- BigTable might work better
- Lets build a service and load test (integration)
- Lets ask if the rest of the team have used it
- That is a risk we need to mitigate

Before we begin

- We will Simplify things
- Real world is more nuanced



Before we begin

- We will Simplify things
- Real world is more nuanced
- After this talk you will be able to understand
 - How Databases differ from one another
 - Identify which DB works best for a given use case
 - How to identify possible bottlenecks



Before we begin

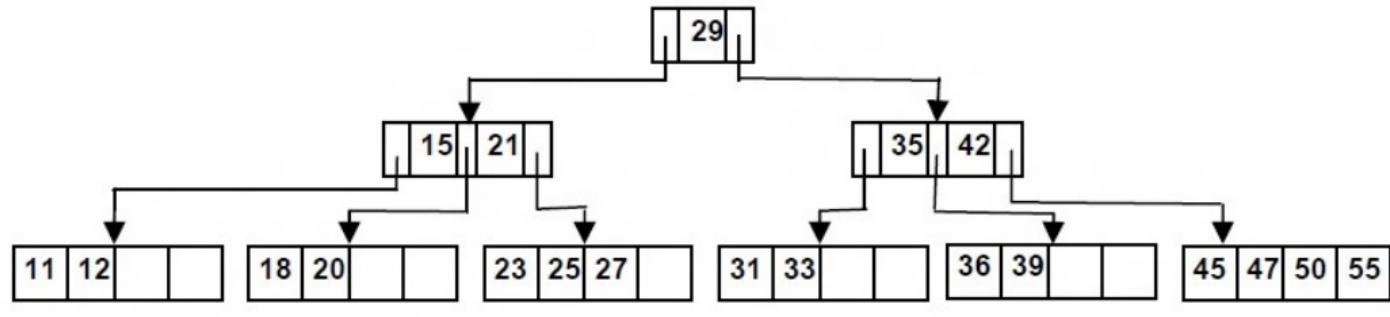
- We will Simplify things
- Real world is more nuanced
- After this talk you will be able to understand
 - How Databases differ from one another
 - Identify which DB works best for a given use case
 - How to identify possible bottlenecks
- You will still have to
 - Read the Manual/Documentation
 - Do your own load/failure tests



B–Tree



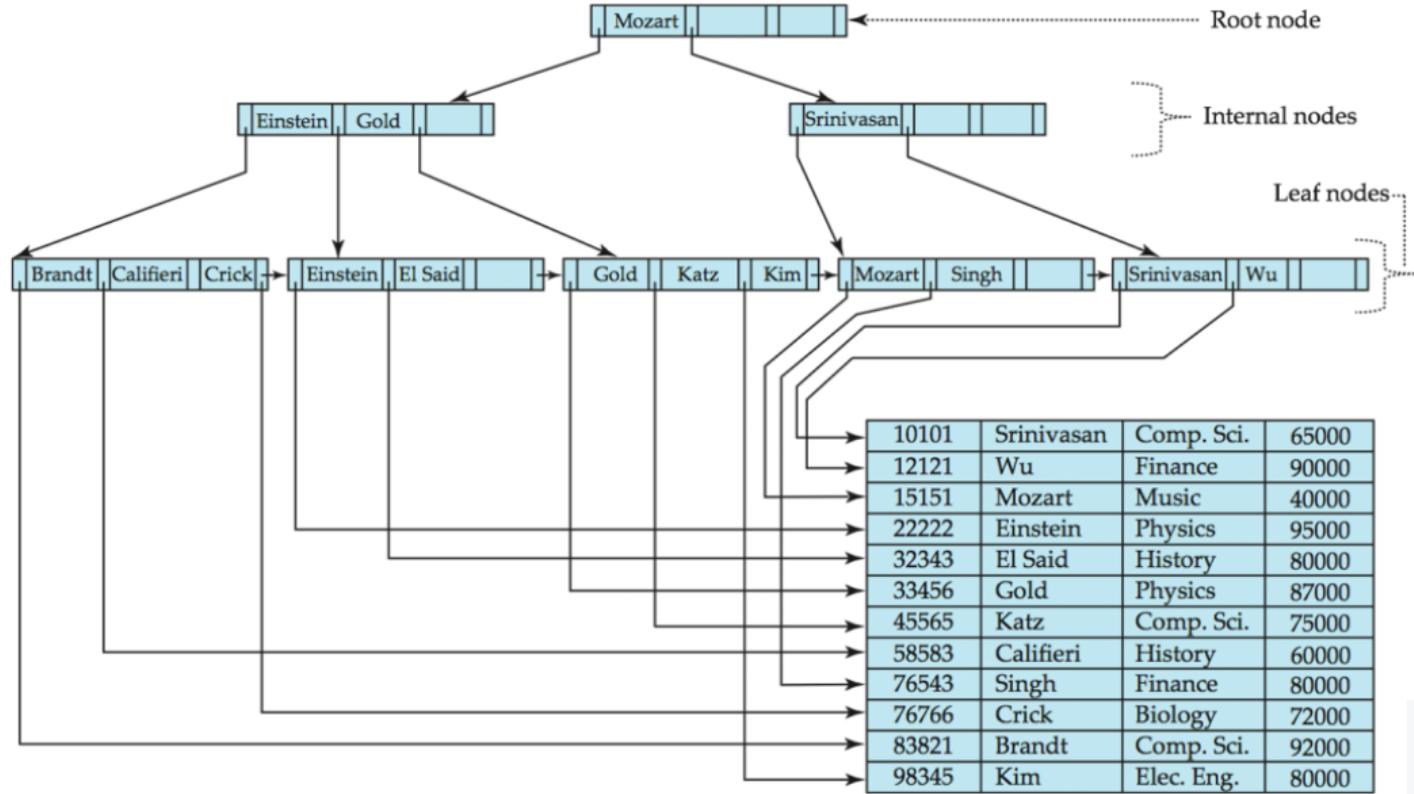
B-Trees



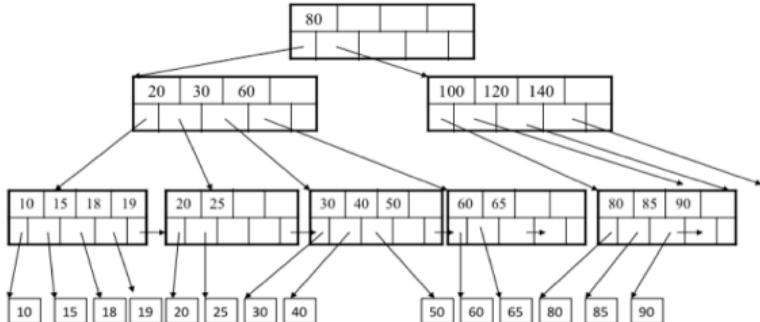
B-Tree of order 5



B+ Tree - index



B-Trees



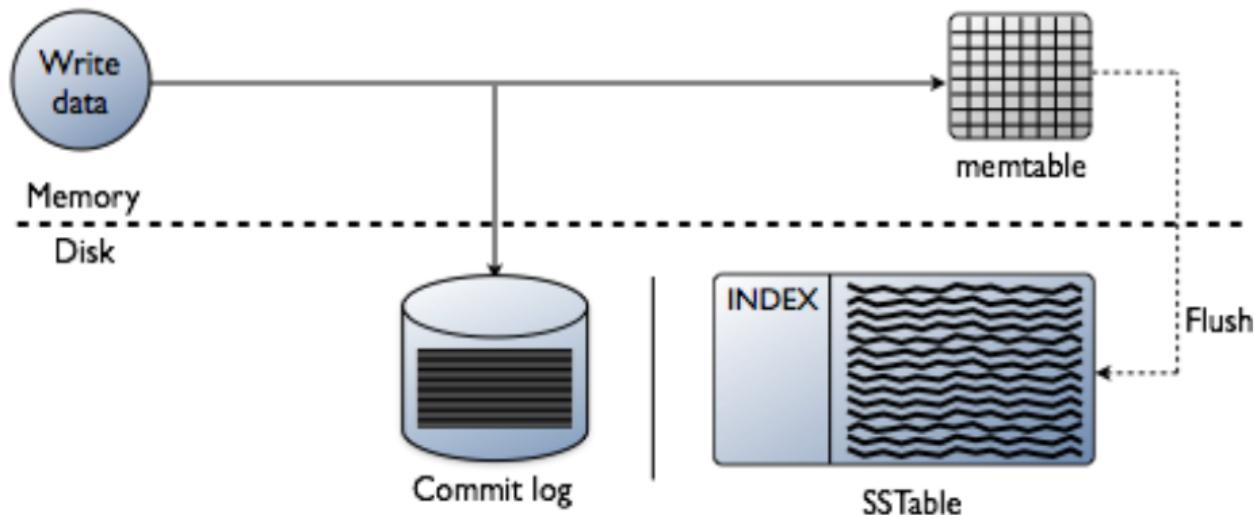
- B Tree variations ⇒ B+ Trees, B* Trees, etc
- Typical reads much higher than writes
- Keeps keys sorted for sequential traversing
- Uses a hierarchical index to minimize the number of disk reads
- Self-Balancing
- No additional maintenance step required
- No Mechanism to delete stale/expired data



Log Structured Merge (a.k.a LSM)



LSM - Write



LSM - SSTable

Index

key	offset
key	offset
...	...

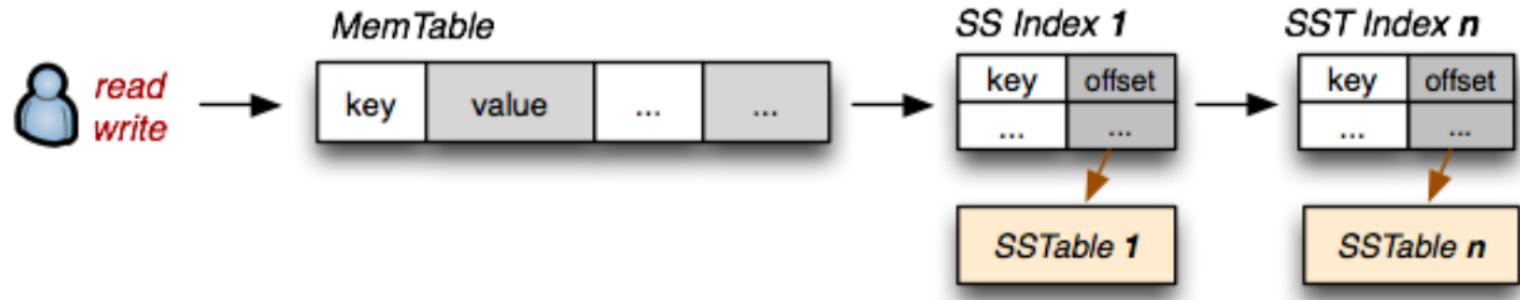
SSTable file

key	value	key	value	key	value
-----	-------	-----	-------	-----	-------	-----	-----

- Sorted String Table
- 2 sections - Index and Data
- SSTables also contain a Bloomfilter
- Append only storage
- Persisted on Disk



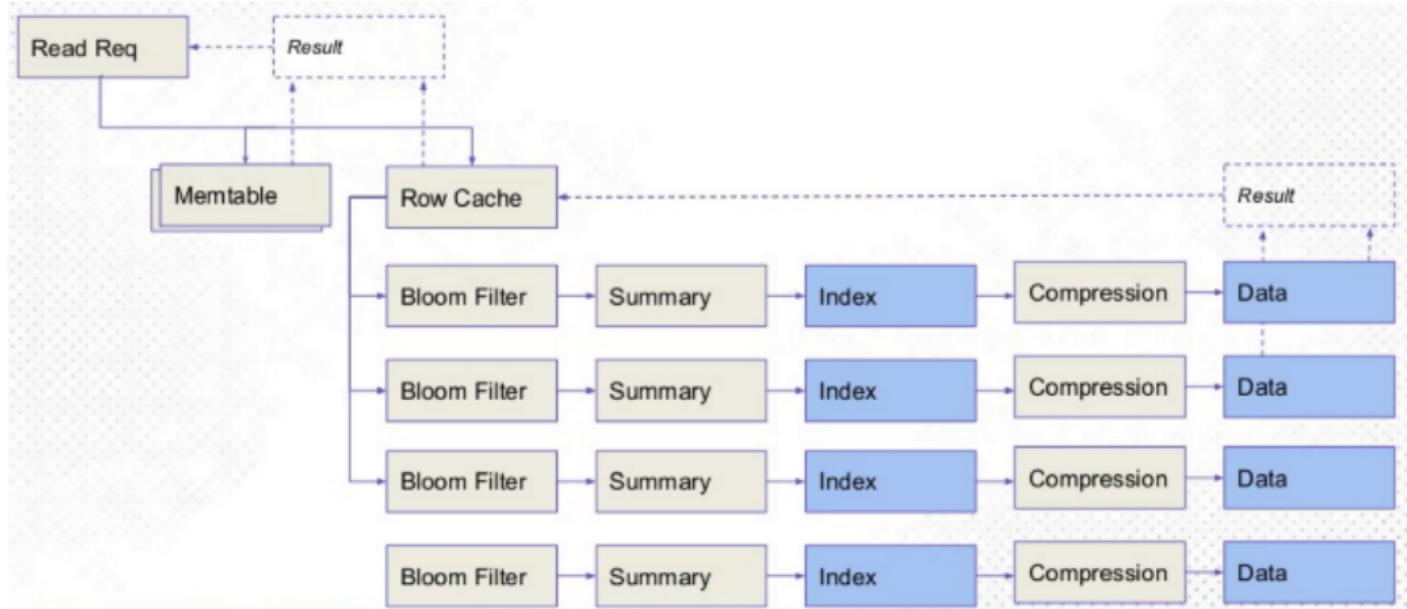
LSM - Memtable



- Memtable helps speed up the requests
- Memtable is flushed to disk into SSTables (time/size)
- Allows Batching of writes to disk



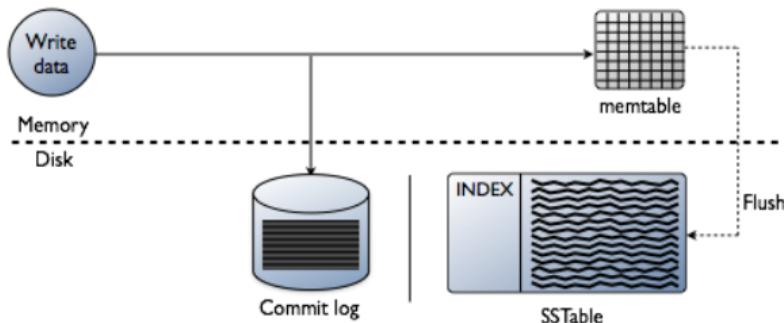
LSM - Read



gojek



Log Structured Merge (LSM)



- Typical writes much higher than reads
- Deletion of values results in tombstones
- Updates across results in versions (by-design / memtable flush)
- Append only storage
- Extra maintenance required (SS table compaction)



What is used where?

DATABASE STORAGE ENGINES

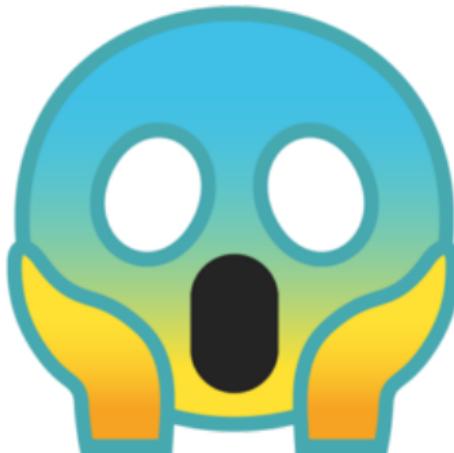
B-TREE



LSM TREE



Time for ACID Trip



ACID defines properties of a Transaction.

- **A**tomicity
- **C**onsistency
- **I**solation
- **D**urability



Issues with Transaction

Primary focus is Isolation.



Issues with Transaction

Primary focus is Isolation.

- **Dirty Read** – Transaction reads data that's not committed.



Issues with Transaction

Primary focus is Isolation.

- **Dirty Read** – Transaction reads data that's not committed.
- **Non Repeatable Read** – Different values on multiple reads of same row.



Issues with Transaction

Primary focus is Isolation.

- **Dirty Read** – Transaction reads data that's not committed.
- **Non Repeatable Read** – Different values on multiple reads of same row.
- **Phantom Read** – Different number of rows returned when same query executed multiple times.



Isolation Levels

SQL defines 4 levels of Isolation.



Isolation Levels

SQL defines 4 levels of Isolation.

- **Read Uncommitted** – Transaction can read uncommitted data.
- **Read Committed** – Transaction only read committed data (Lock on current row).
- **Repeatable Read** – Transactions don't have *Non Repeatable Reads* (Locks on all referenced rows).
- **Serializable** – Concurrently executing Transaction appear to be serially executing (No defined order)



ACID – how does it come together?

Isolation Level	Dirty reads	Non-repeatable reads	Phantoms
Read Uncommitted	May occur	May occur	May occur
Read Committed	Don't occur	May occur	May occur
Repeatable Read	Don't occur	Don't occur	May occur
Serializable	Don't occur	Don't occur	Don't occur



Let's put our CAP on!



CAP Theorem

A **Distributed Database** can only provide two out of the following three guarantees.



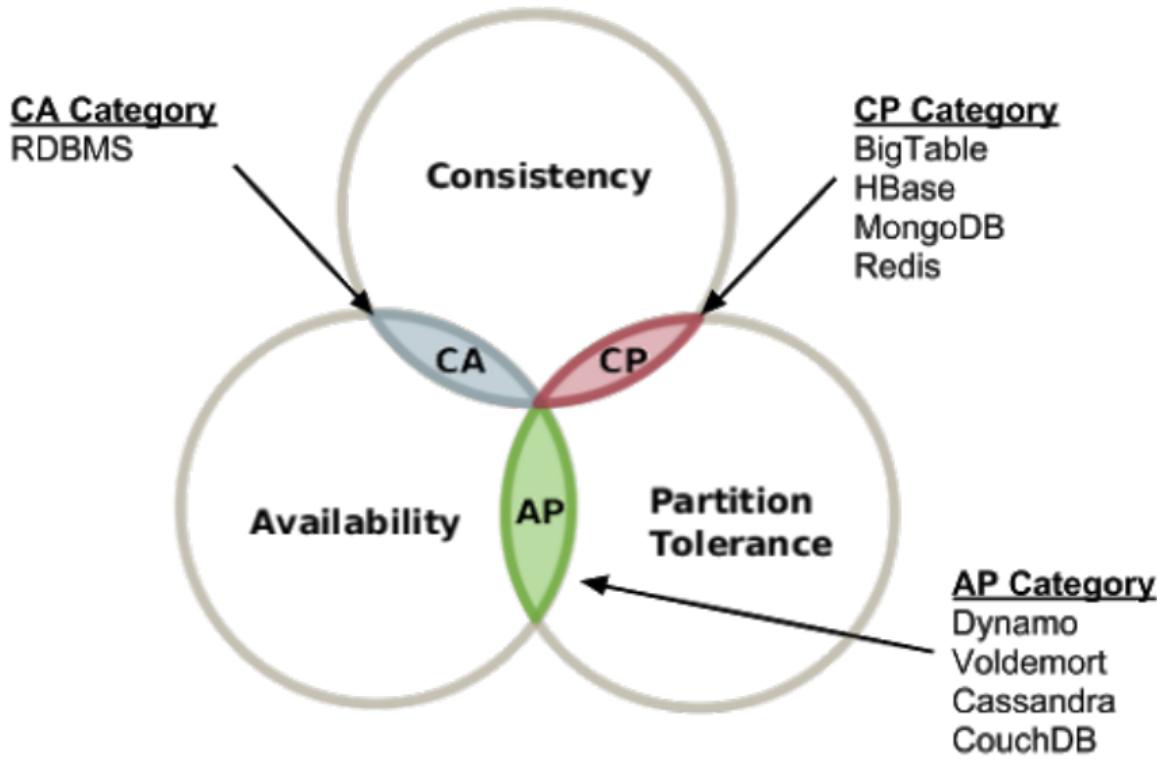
CAP Theorem

A **Distributed Database** can only provide two out of the following three guarantees.

- **Consistency** – Every read receives the most recent write or an error.
- **Availability** – Every request receives a (*non-error*) response.
- **Partition Tolerance** – System continues to operate in case of network partition.



CAP Theorem



CAP Theorem Issues

- CAP Theorem is **only** valid for partition events.
- CAP Theorem does not define normal working.



CAP Theorem Issues

- CAP Theorem is **only** valid for partition events.
- CAP Theorem does not define normal working.
- PACELC Theorem goes one level extra.



PACELC Theorem



PACELC Theorem

- Extension of CAP Theorem.
- Describes normal case tradeoff between
 - Consistency
 - Latency



PACELC Theorem

- Extension of CAP Theorem.
- Describes normal case tradeoff between
 - Consistency
 - Latency
- **PAC** – When there is **P**artition a system is either **A**vailable or **C**onsistent.



PACELC Theorem

- Extension of CAP Theorem.
- Describes normal case tradeoff between
 - Consistency
 - Latency
- **PAC** – When there is **P**artition a system is either **A**vailable or **C**onsistent.
- **ELC** – **E**lse the system focuses on either **L**atency or **C**onsistency.



Where do modern DBs fit?

DDBS	P+A	P+C	E+L	E+C
DynamoDB	✓		✓ ^[a]	
Cassandra	✓		✓ ^[a]	
Cosmos DB		✓	✓ ^[b]	
Couchbase		✓	✓	✓
Riak	✓		✓ ^[a]	
VoltDB/H-Store		✓		✓
Megastore		✓		✓
BigTable/HBase		✓		✓
MySQL Cluster		✓		✓
MongoDB	✓			✓
PNUTS		✓	✓	
Hazelcast IMDG ^{[8][6]}	✓	✓	✓	✓
FaunaDB ^[9]		✓	✓	✓

↑ is default behaviour, some can be tuned



Sharding



Sharding

- Shard is partition of data.
- Each shard acts as the single source for this subset of data



Sharding

- Shard is partition of data.
- Each shard acts as the single source for this subset of data
- Each shard can be kept in different physical location.
- Allows databases to scale linearly.
- Multiple copies of same shard are normally kept on different servers.



Modulo Hashing

Modulo operation finds the remainder after division of one number by another

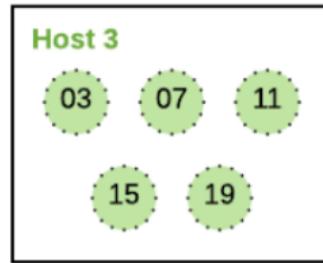
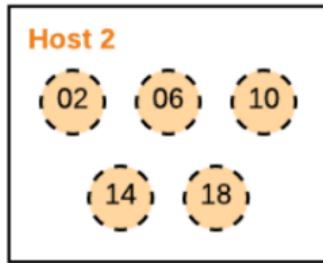
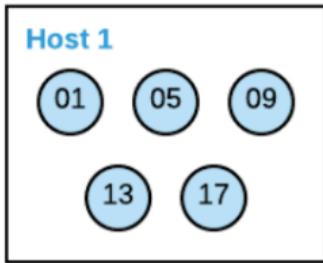
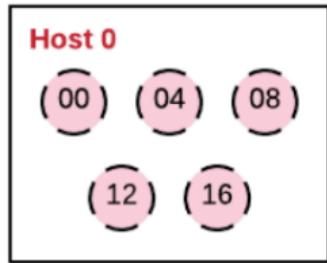


Modulo Hashing

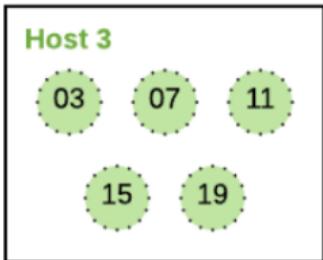
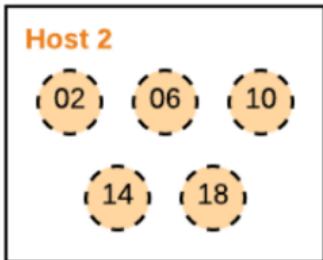
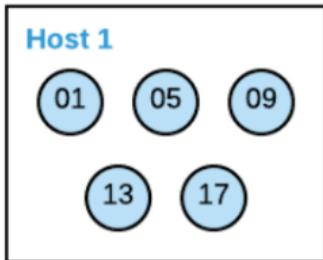
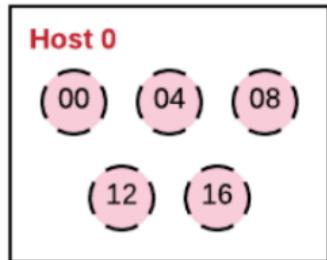
- Choose a hash function (Murmur3, xxHash)
- Hash the incoming key and take modulo
- Select the server based on that modulo

$$\text{server} := \text{servers}[\text{hash}(\text{key}) \% \text{len}(\text{servers})]$$

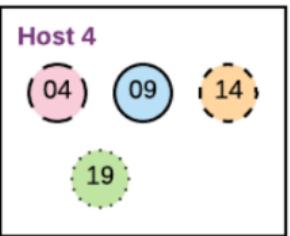
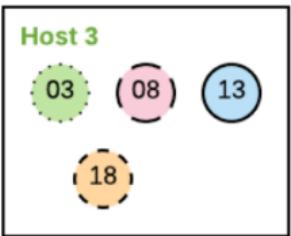
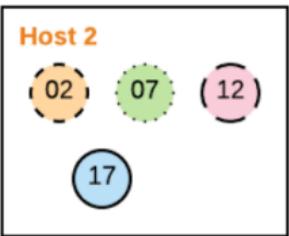
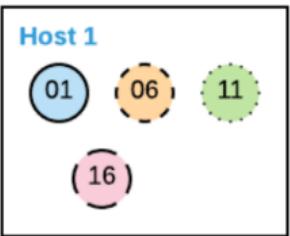
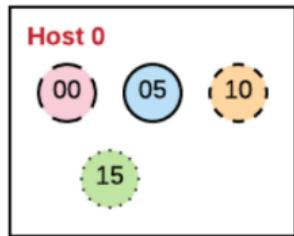

Problems with Modulo



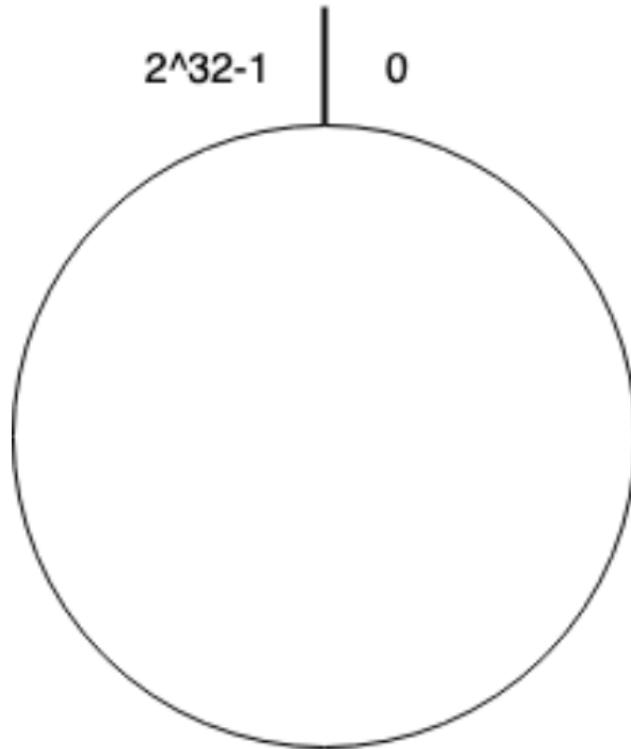
Problems with Modulo



On adding new server



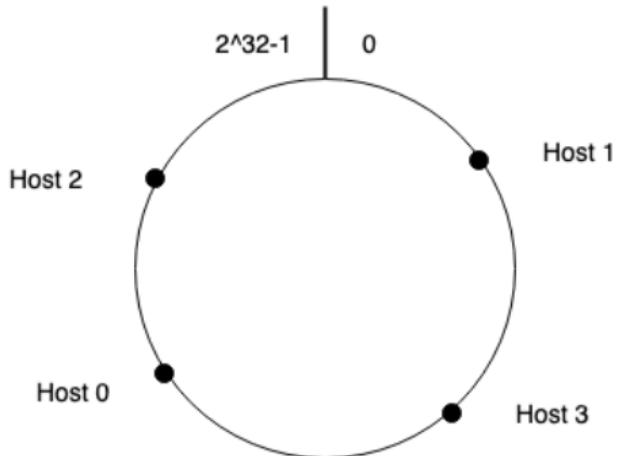
Consistent Hashing



- Take key space from $0 \dots 2^{32} - 1$
- Each point on circle becomes a shard (vbuckets)



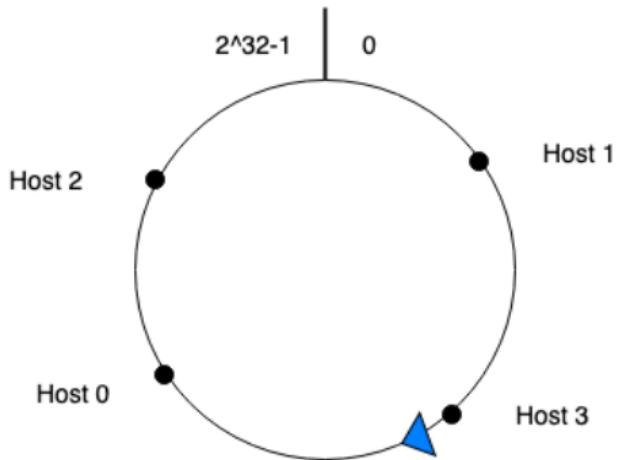
Consistent Hashing



- Take key space from $0 \dots 2^{32} - 1$
- Each point on circle becomes a shard (vbuckets)
- Hash and put each server on the circle.



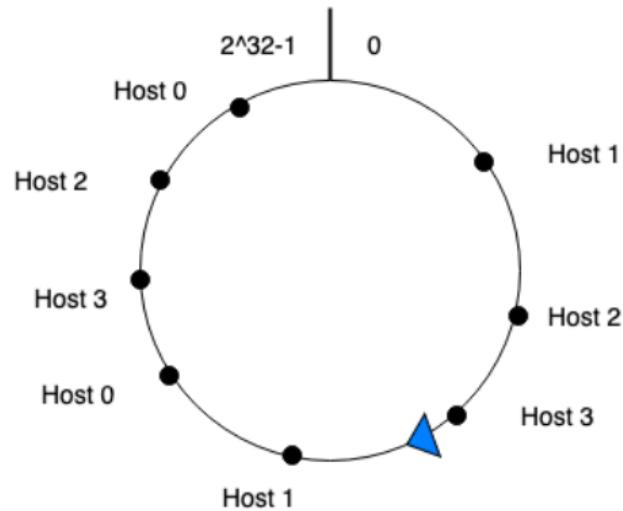
Consistent Hashing



- Take key space from $0 \dots 2^{32} - 1$
- Each point on circle becomes a shard (vbuckets)
- Hash and put each server on the circle.
- Host owns shard (vbuckets) starting from itself till the next host



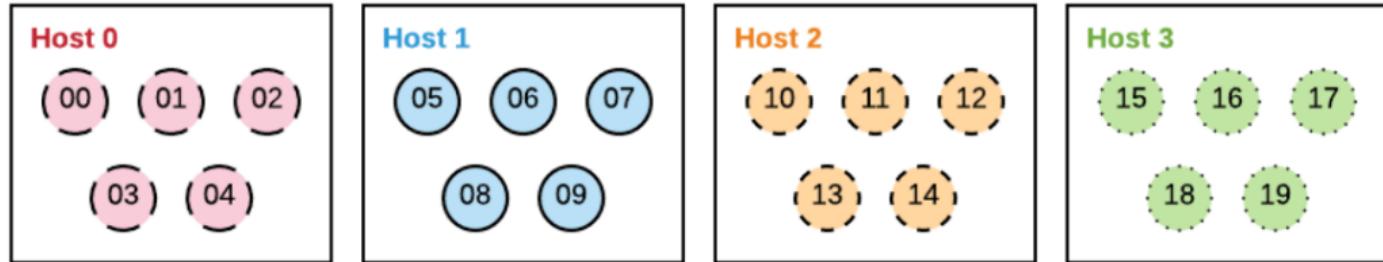
Consistent Hashing



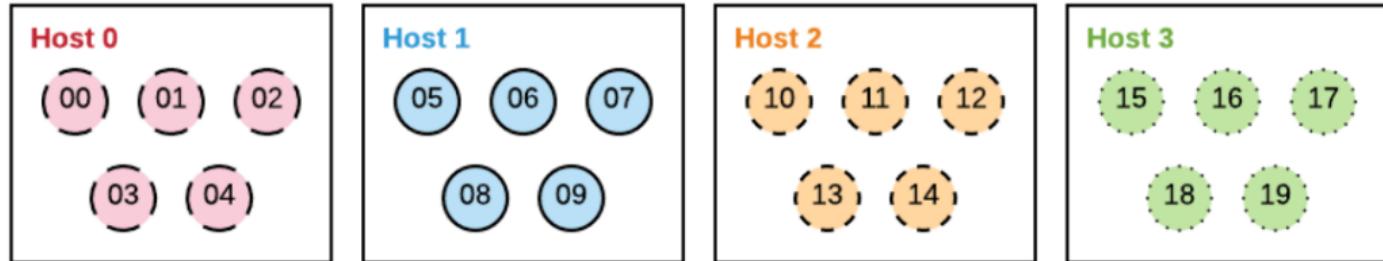
- Take key space from $0 \dots 2^{32} - 1$
- Each point on circle becomes a shard (vbuckets)
- Hash and put each server on the circle.
- Host owns shard (vbuckets) starting from itself till the next host
- Do this multiple times



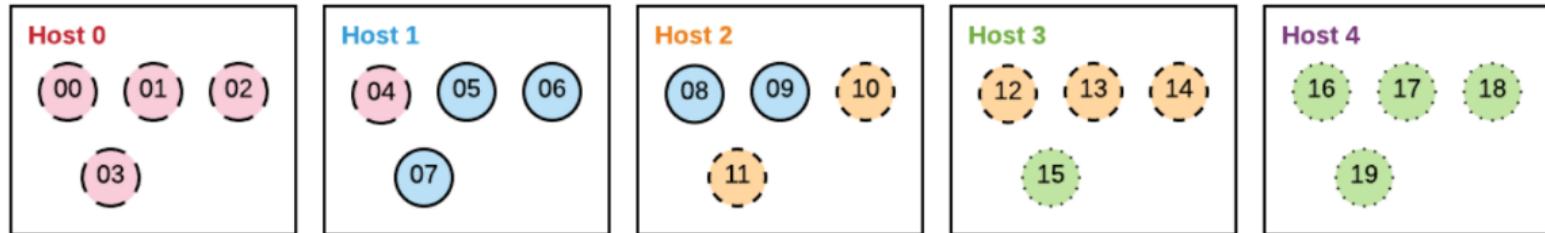
Using Consistent Hashing - ketama



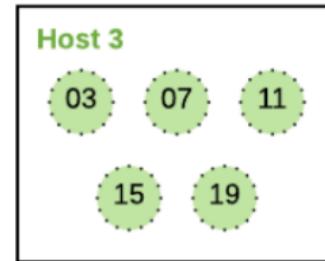
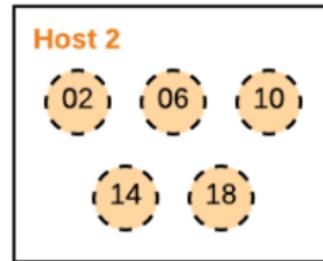
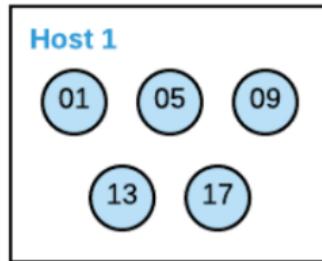
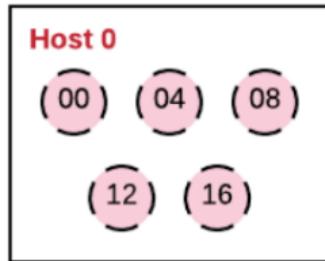
Using Consistent Hashing - ketama



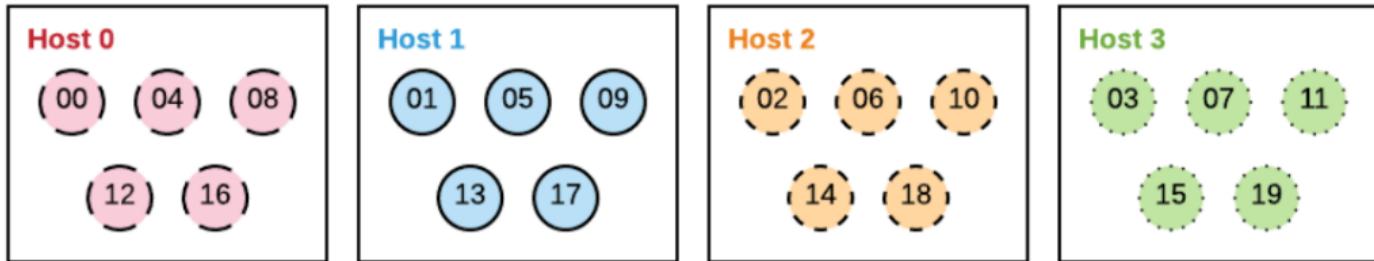
On adding new server



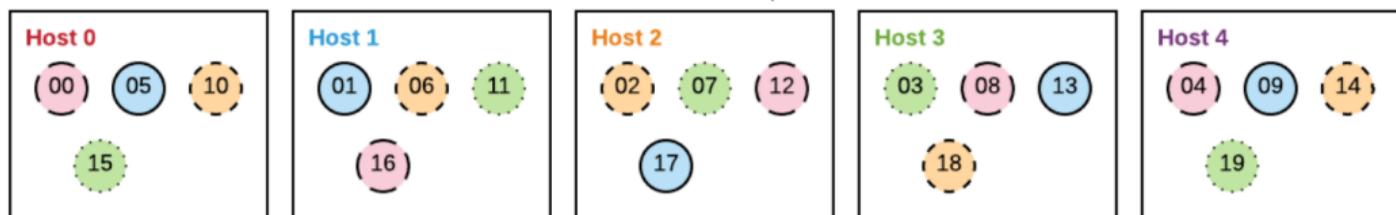
Initial ↓



Initial ↓

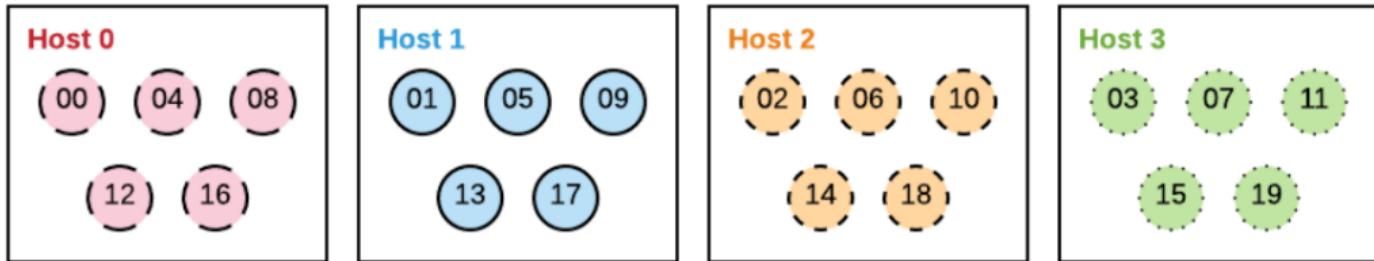


Modulo ↓

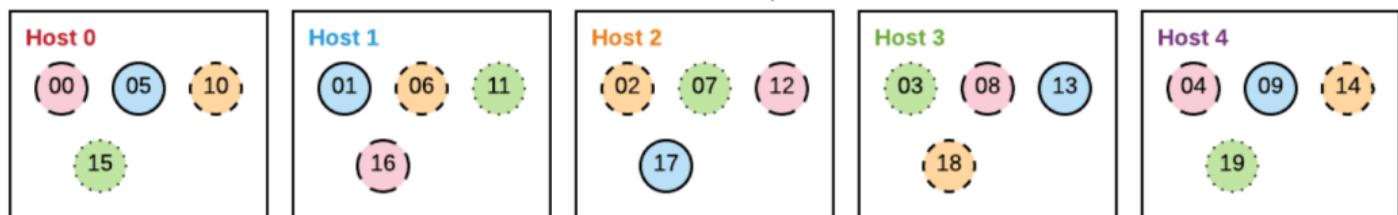


gojek

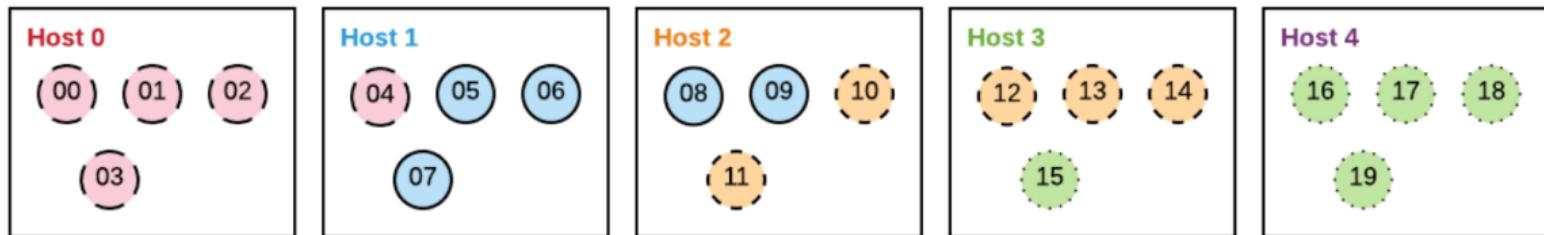
Initial ↓



Modulo ↓



Consistent (ketama) ↓



Quorums and Consensus Algorithms

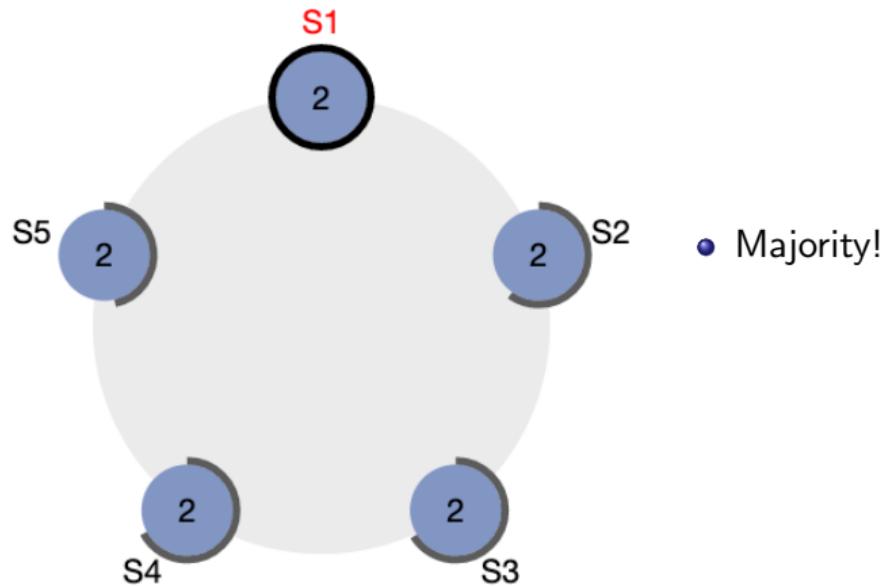


Quorums

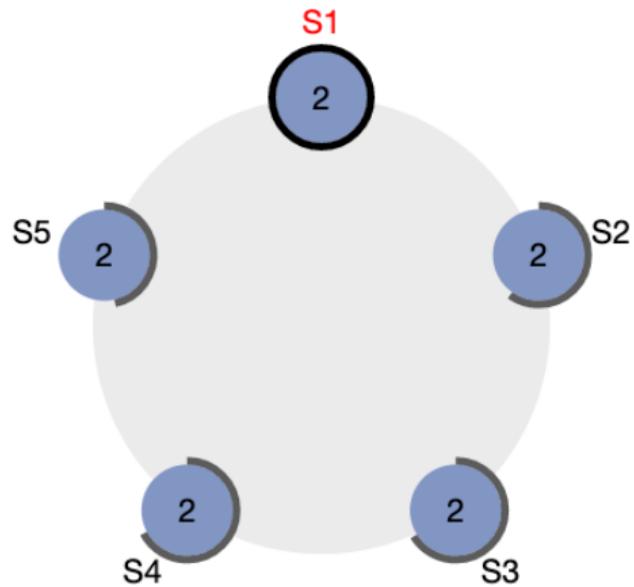
The minimum number of members of an assembly or society that must be present at any of its meetings to make the proceedings of that meeting valid.



Quorums

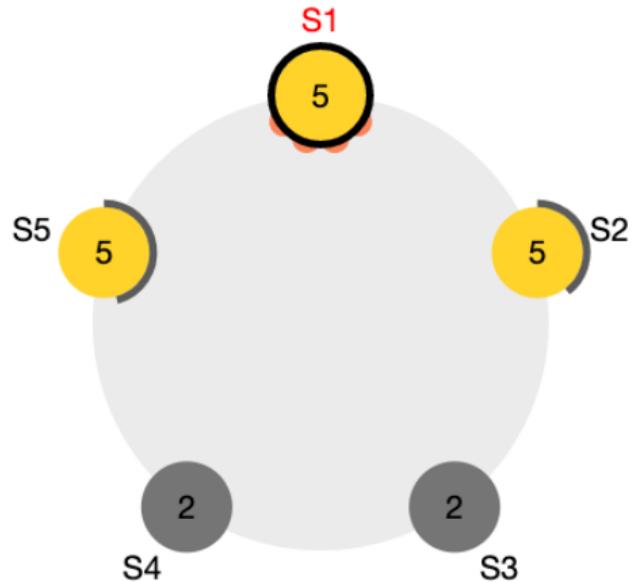


Quorums



- Majority!
- Defined as $\lceil \frac{n+1}{2} \rceil$ where n = number of nodes

Quorums



- Majority!
- Defined as $\lceil \frac{n+1}{2} \rceil$ where $n =$ number of nodes
- Can work with failure of $n - \lceil \frac{n+1}{2} \rceil$ nodes

Consensus Algorithms

- Consensus is a fundamental problem in fault-tolerant distributed systems.
- Consensus involves multiple servers agreeing on values.



Consensus Algorithms

- Consensus is a fundamental problem in fault-tolerant distributed systems.
- Consensus involves multiple servers agreeing on values.
- Consensus Algorithms can help maintain Quorums.



Consensus Algorithms

- Consensus is a fundamental problem in fault-tolerant distributed systems.
- Consensus involves multiple servers agreeing on values.
- Consensus Algorithms can help maintain Quorums.
- Some of the most popular Consensus Algorithms are
 - Paxos
 - Raft



Consensus Algorithms

- Consensus is a fundamental problem in fault-tolerant distributed systems.
- Consensus involves multiple servers agreeing on values.
- Consensus Algorithms can help maintain Quorums.
- Some of the most popular Consensus Algorithms are
 - Paxos
 - Raft
- Zookeeper Atomic Broadcast Protocol (ZAB) another Consensus Algorithm but uses Two Phase commit.



How can nodes maintain data consistency without complex consensus algorithms?

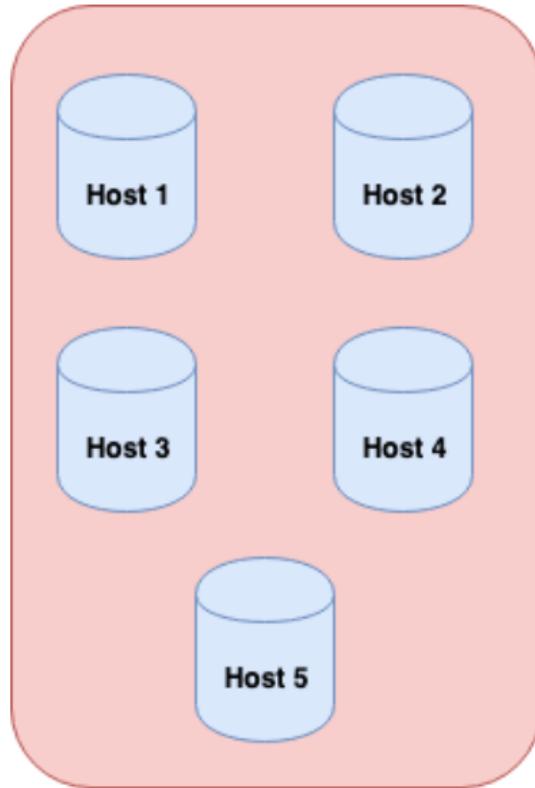


How do multiple nodes maintain data consistency?

$$R + W > N$$

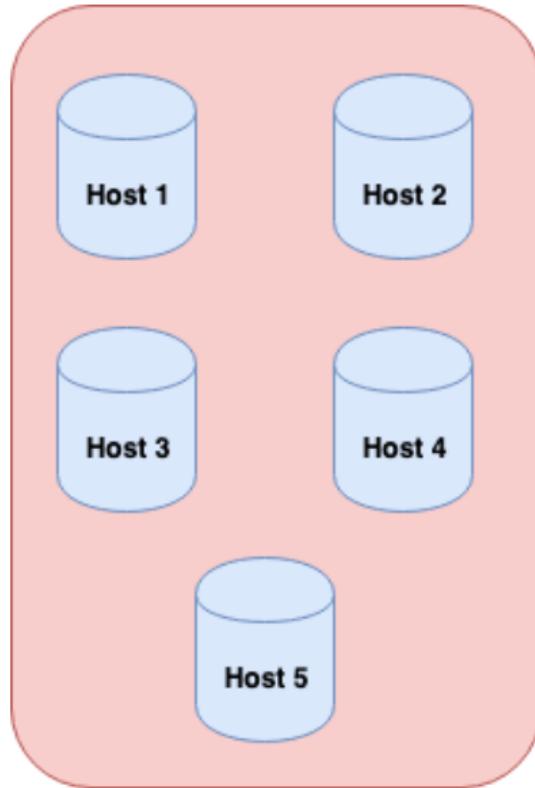


How do multiple nodes maintain data consistency?



- Reads Ack + Writes Ack > Total Number of Replicas

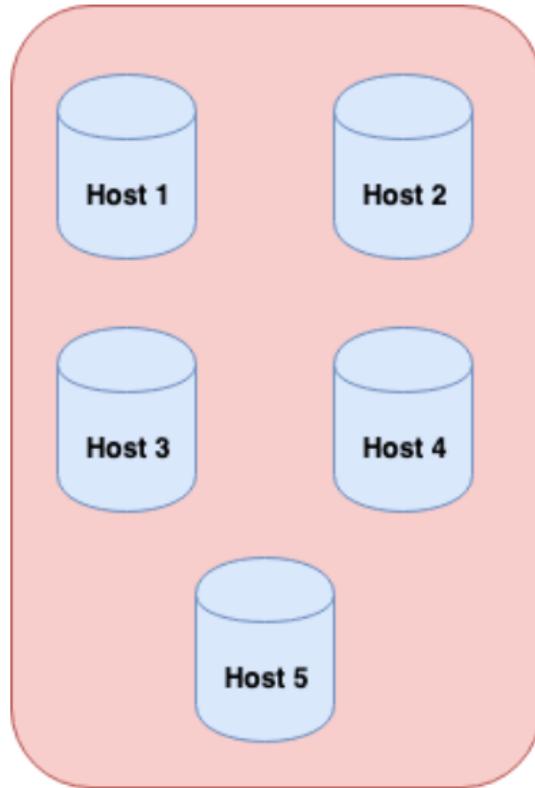
How do multiple nodes maintain data consistency?



- Reads Ack + Writes Ack > Total Number of Replicas
- Lets walk through some example
- R:1 + W:5 > N:5



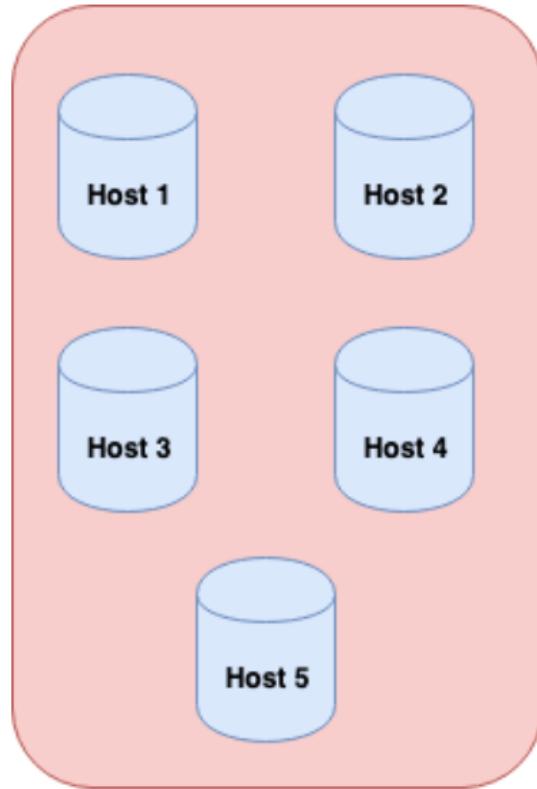
How do multiple nodes maintain data consistency?



- Reads Ack + Writes Ack > Total Number of Replicas
- Lets walk through some example
- R:1 + W:5 > N:5
- R:5 + W:1 > N:5



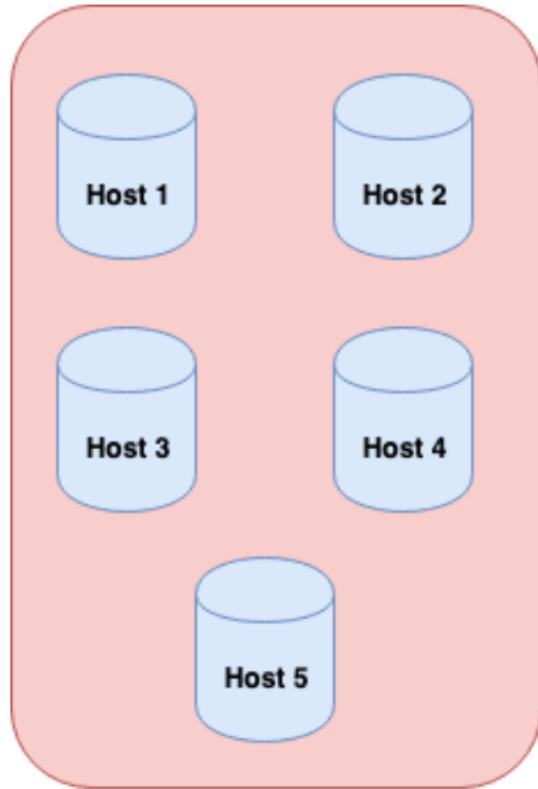
How do multiple nodes maintain data consistency?



- Reads Ack + Writes Ack > Total Number of Replicas
 - Lets walk through some example
 - R:1 + W:5 > N:5
 - R:5 + W:1 > N:5 (read repair)
 - R:3 + W:3 > N:5



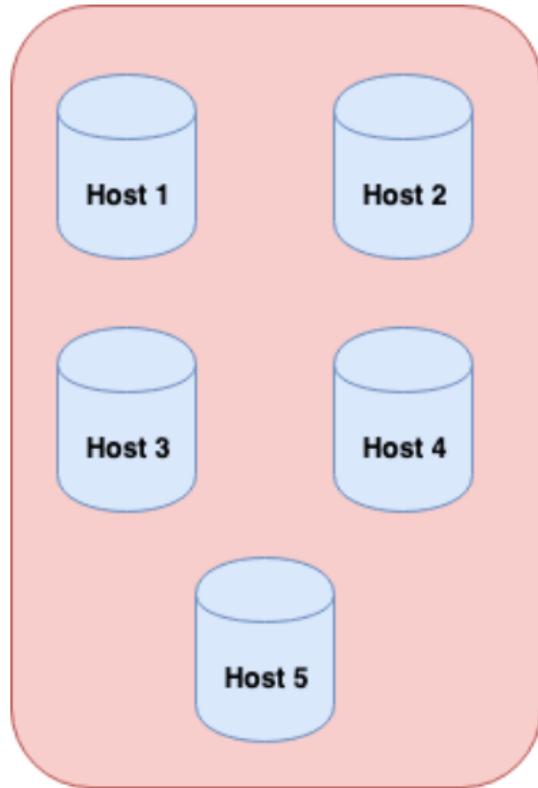
How do multiple nodes maintain data consistency?



- Reads Ack + Writes Ack > Total Number of Replicas
- Lets walk through some example
- R:1 + W:5 > N:5
- R:5 + W:1 > N:5 (read repair)
- R:3 + W:3 > N:5
- R:2 + W:3 = N:5 ??



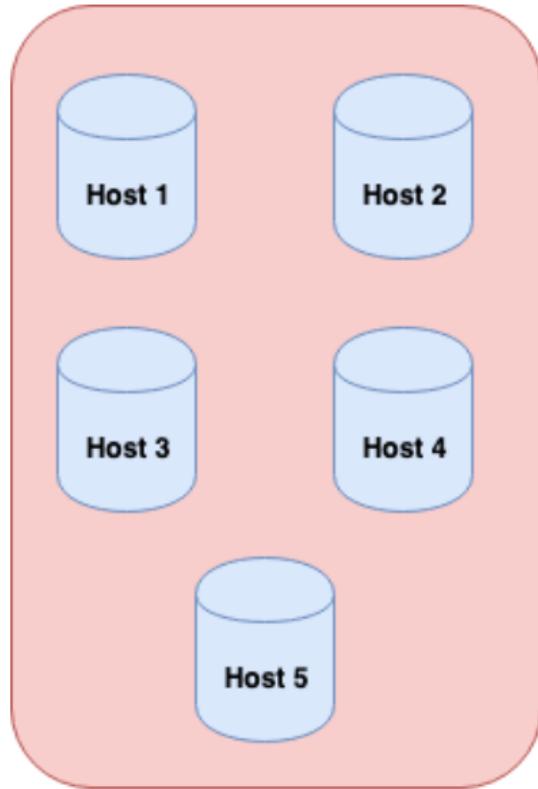
How do multiple nodes maintain data consistency?



- Reads Ack + Writes Ack > Total Number of Replicas
- Lets walk through some example
- R:1 + W:5 > N:5
- R:5 + W:1 > N:5 (read repair)
- R:3 + W:3 > N:5
- R:2 + W:3 = N:5 ??
(eventual consistency)



How do multiple nodes maintain data consistency?



- Reads Ack + Writes Ack > Total Number of Replicas
- Lets walk through some example
- R:1 + W:5 > N:5
- R:5 + W:1 > N:5 (read repair)
- R:3 + W:3 > N:5
- R:2 + W:3 = N:5 ??
(eventual consistency)
- R:2 + W:2 < N:5



How do multiple nodes maintain data consistency?

$$R + W > N \text{ (Strong)}$$
$$R + W \leq N \text{ (Eventual)}$$




Cassandra

- LSM Based





- LSM Based
- Tunable Consistency (request time)
 - ONE / LOCAL_ONE / TWO / THREE
 - QUORUM / LOCAL_QUOROM / EACH_QUOROM
(only for writes)
 - ANY (only for writes), ALL ...

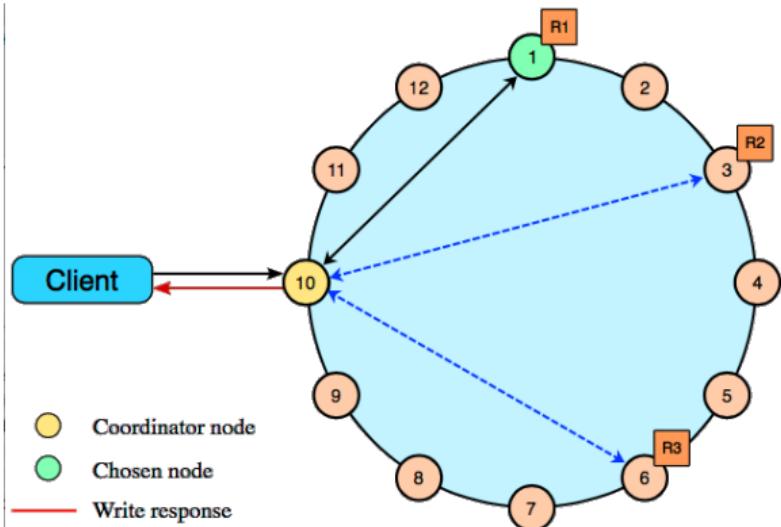




- LSM Based
- Tunable Consistency (request time)
 - ONE / LOCAL_ONE / TWO / THREE
 - QUORUM / LOCAL_QUORUM / EACH_QUORUM
(only for writes)
 - ANY (only for writes), ALL ...
- Uses repair mechanism to fix inconsistency
 - hinted handoff (node is down)
 - read repair (fix data when reading)
 - anti-entropy node repair (manual/background repair)
- Consistent Hashing
- Built-in support for multi-master multi-DC clusters
- Gossip for cluster membership



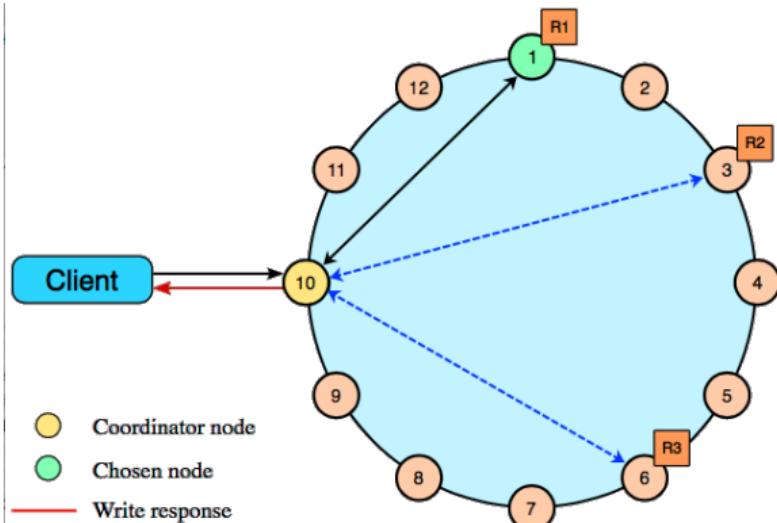
Cassandra - ONE Write (W:1 N:3)



- Writes can go to any node
- Node receiving the request becomes Coordinator



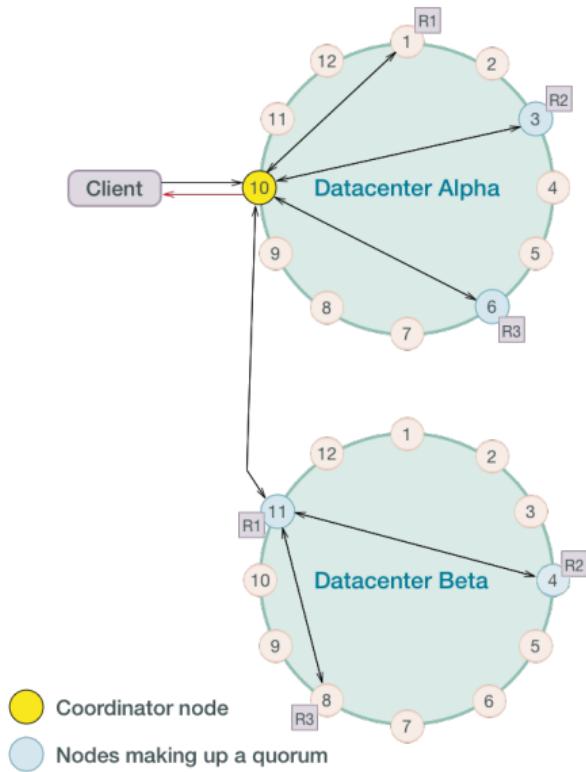
Cassandra - ONE Write (W:1 N:3)



- Writes can go to any node
- Node receiving the request becomes Coordinator
- Coordinator writes to all replicas
- When one replica write is successful, user gets a success
- Num Ack depends on Consistency level
 - ONE / LOCAL_ONE / TWO / THREE
 - QUORUM / LOCAL_QUORUM / EACH_QUORUM (only for writes)
 - ANY (only for writes), ALL ...



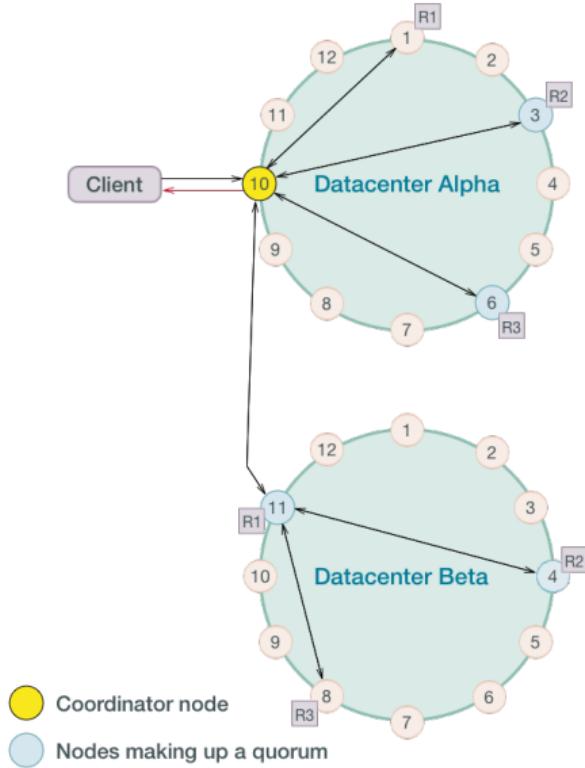
Cassandra - Multi DC Write



- Writes can go to any node
- Node receiving the request becomes Coordinator



Cassandra - Multi DC Write

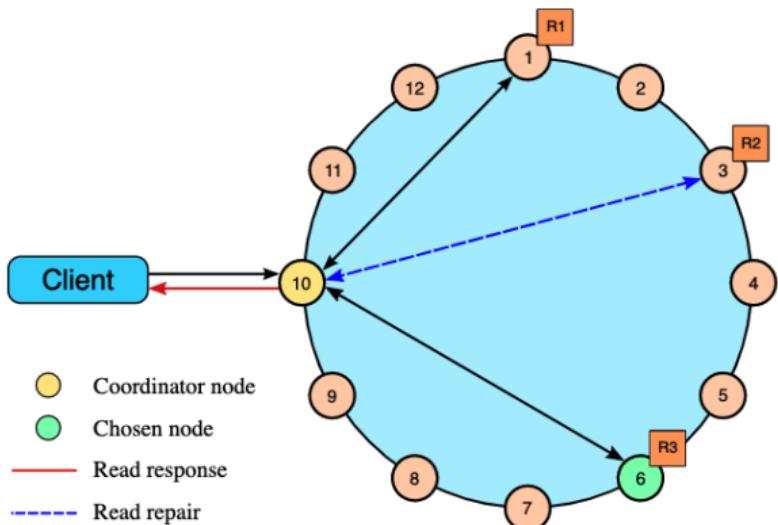


- Writes can go to any node
 - Node receiving the request becomes Coordinator
 - Coordinator writes to all replicas
 - When one replica write is successful, user gets a success
 - Num Ack depends on Consistency level
 - ONE / LOCAL_ONE / TWO / THREE
 - QUORUM / LOCAL_QUORUM / EACH_QUORUM (only for writes)
 - ANY (only for writes), ALL ...

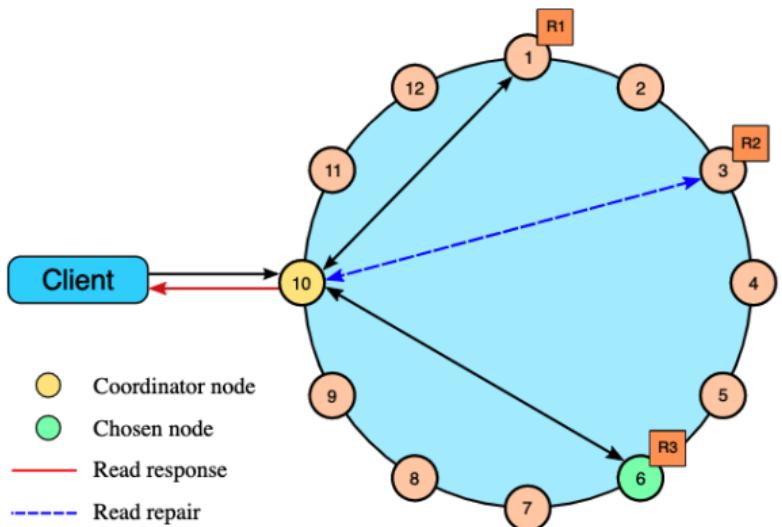


Cassandra - QUORUM Read (R:2 N:3)

- Reads can go to any node
 - Node receiving the request becomes Coordinator



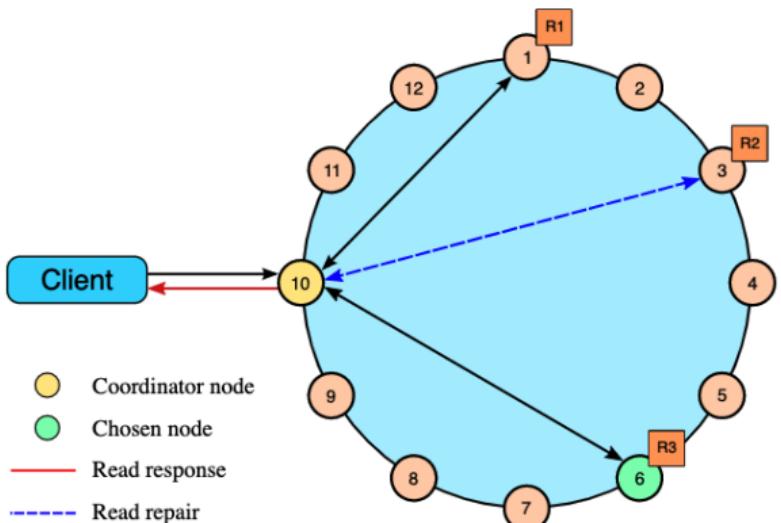
Cassandra - QUORUM Read (R:2 N:3)



- Reads can go to any node
 - Node receiving the request becomes Coordinator
 - Coordinator reads from all replicas
 - When Quorum (2) replicas return up to date values, user get a response



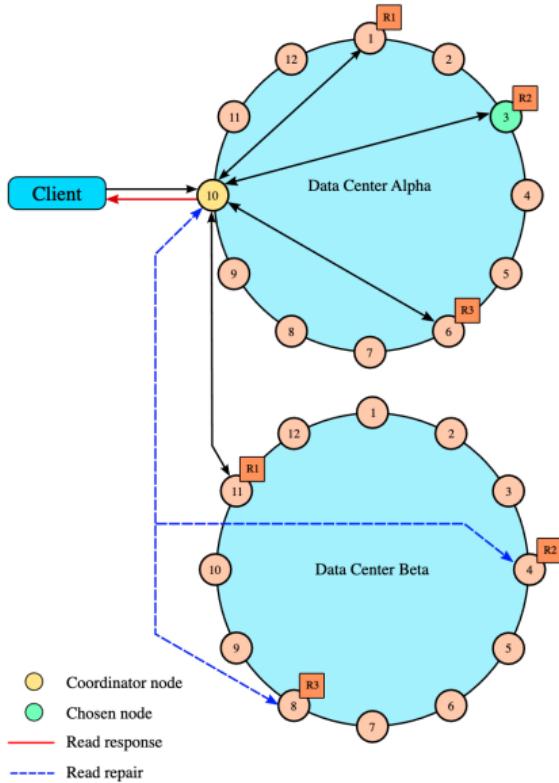
Cassandra - QUORUM Read (R:2 N:3)



- Reads can go to any node
- Node receiving the request becomes Coordinator
- Coordinator reads from all replicas
- When Quorum (2) replicas return up to date values, user get a response
- When replicas have different versions of the row, the replica with the most recent version will return the requested data
- a read repair may be initiated for the disagreeing replicas
- Last Write Wins (timestamp)

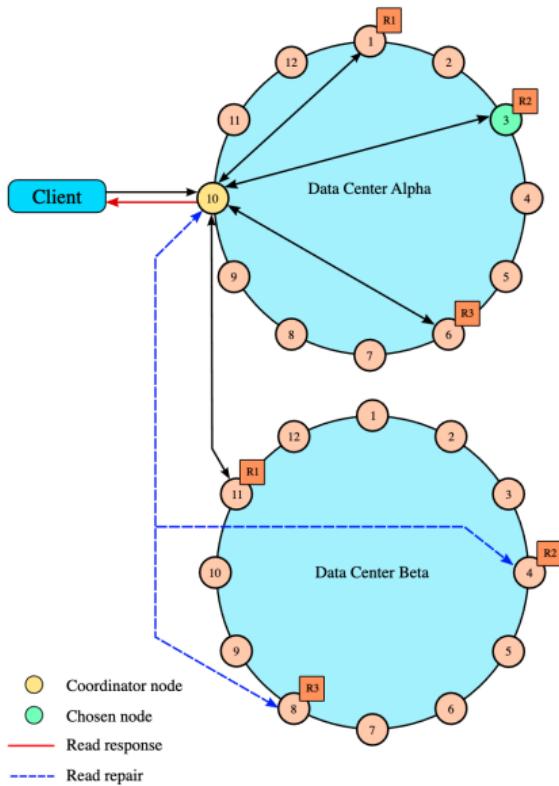


Cassandra - Multi-DC QUORUM Read (R:4 N:6)



- Reads can go to nodes in any DC
- Node receiving the request becomes Coordinator

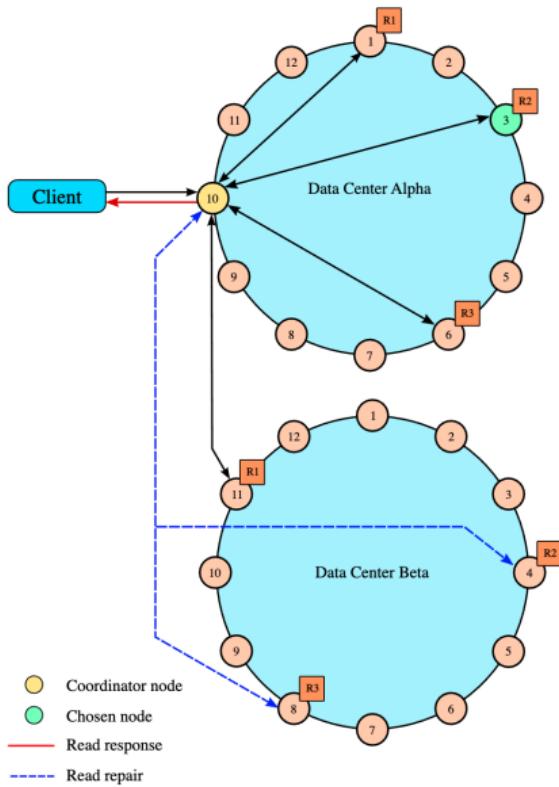
Cassandra - Multi-DC QUORUM Read (R:4 N:6)



- Reads can go to nodes in any DC
- Node receiving the request becomes Coordinator
- Coordinator reads from all replicas across DC
- When Quorum (4 out of 6 total) replicas return up to date values, user get a response



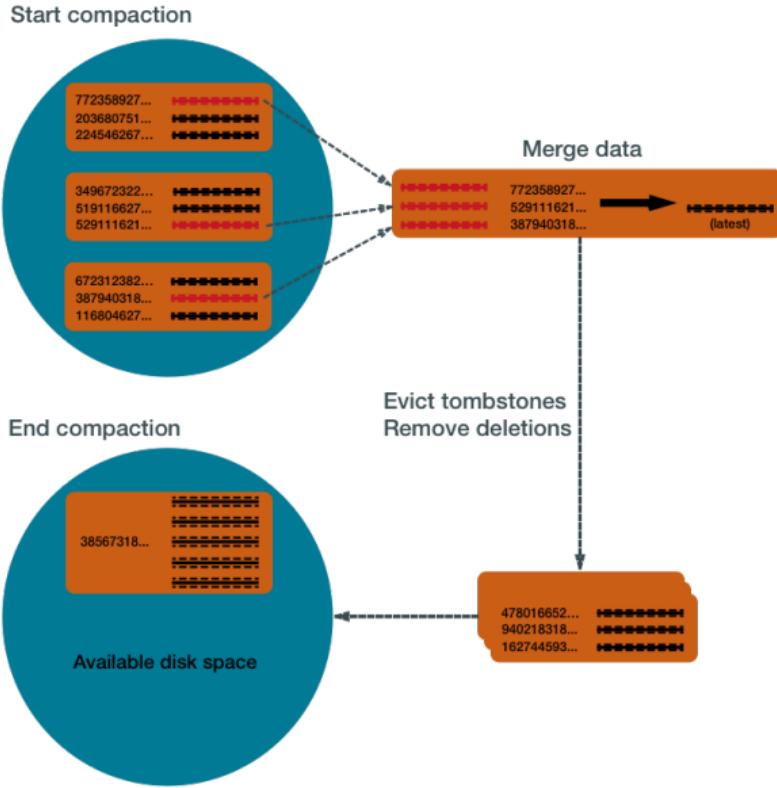
Cassandra - Multi-DC QUORUM Read (R:4 N:6)



- Reads can go to nodes in any DC
- Node receiving the request becomes Coordinator
- Coordinator reads from all replicas across DC
- When Quorum (4 out of 6 total) replicas return up to date values, user get a response
- When replicas have different versions of the row, the replica with the most recent version will return the requested data
- a read repair may be initiated for the disagreeing replicas
- Last Write Wins (timestamp)



Cassandra - Compaction



- Deletes / Updates can result in multiple versions of data in SSTables
 - SSTables are compacted based on defined strategy
 - Size-Tiered Compaction (STCS)
 - Date-Tiered Compaction(DTCS)
 - Time-Window Compaction(TWCS)
 - Leveled Compaction (LCS)
 - Compaction takes care of clean up as well (TTL)



For more please refer to the Excellent Cassandra Documentation!
<https://docs.datastax.com/en/cassandra-oss/3.x/cassandra/dml/dmlIntro.html>

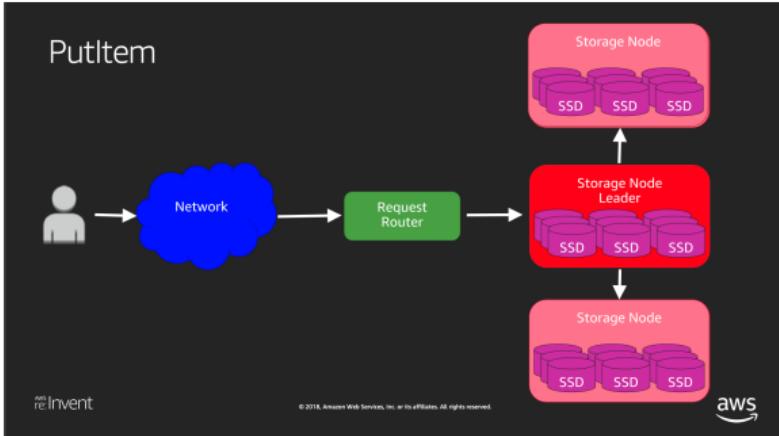




DynamoDB

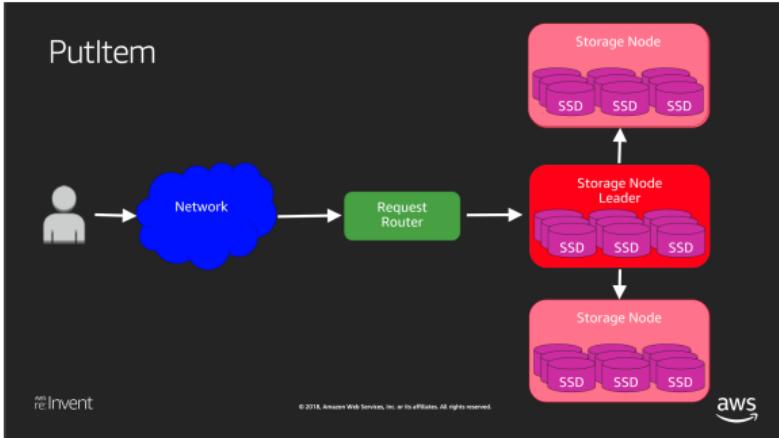


DynamoDB – Writes



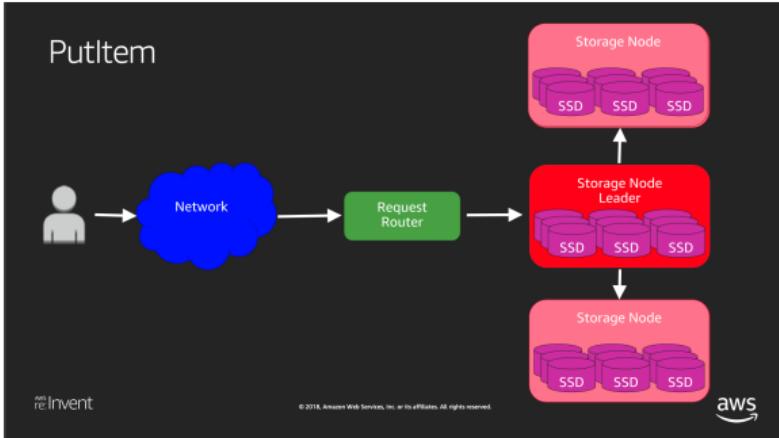
- B-tree based structure with replication log.

DynamoDB – Writes



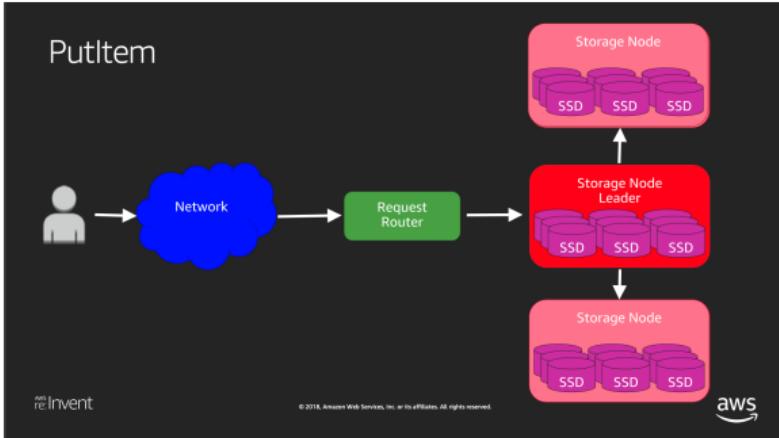
- B-tree based structure with replication log.
- Requests flow through a request router

DynamoDB – Writes



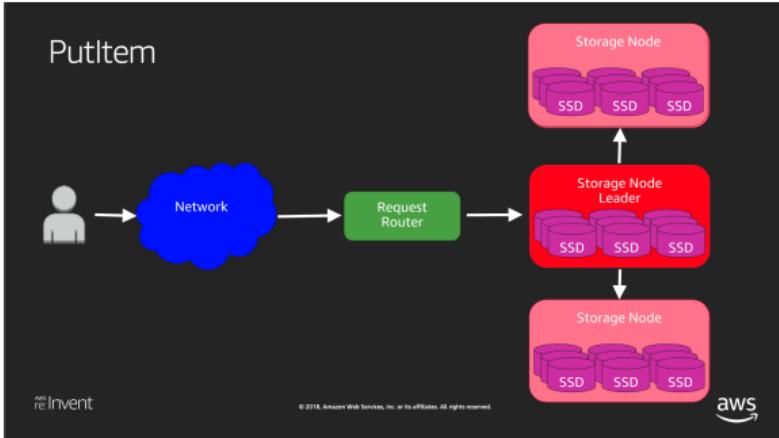
- B-tree based structure with replication log.
- Requests flow through a request router
- Consistent Hashing for partitioning

DynamoDB – Writes



- B-tree based structure with replication log.
- Requests flow through a request router
- Consistent Hashing for partitioning
- Paxos for Leader Election.

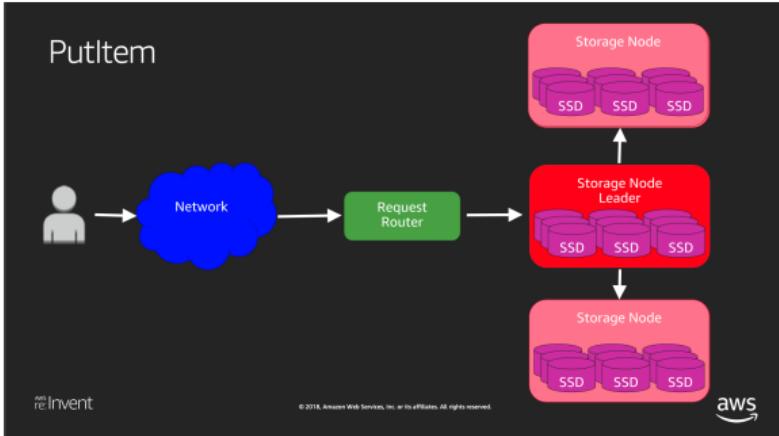
DynamoDB – Writes



- B-tree based structure with replication log.
- Requests flow through a request router
- Consistent Hashing for partitioning
- Paxos for Leader Election.
- Num Ack of 2
- Replica Count of 3
- Replicas across AZ



DynamoDB – Writes

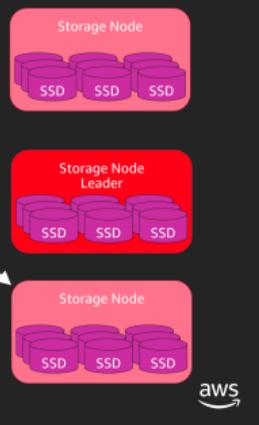


- B-tree based structure with replication log.
- Requests flow through a request router
- Consistent Hashing for partitioning
- Paxos for Leader Election.
- Num Ack of 2
- Replica Count of 3
- Replicas across AZ
- HeartBeat for peer failure detection.



DynamoDB – Reads

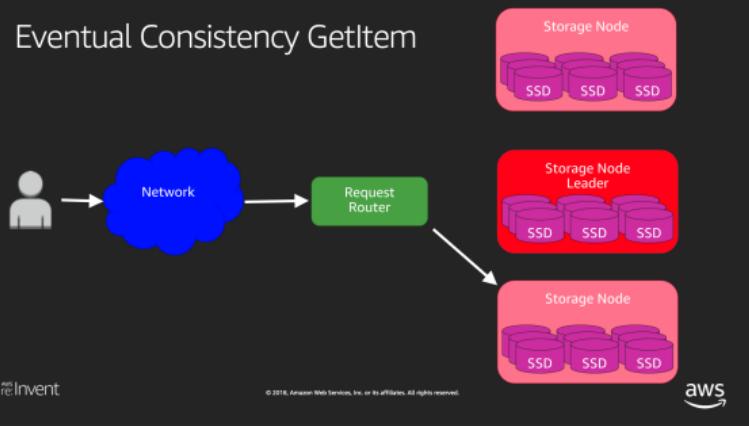
Eventual Consistency GetItem



- Consistent Hashing for partitioning

DynamoDB – Reads

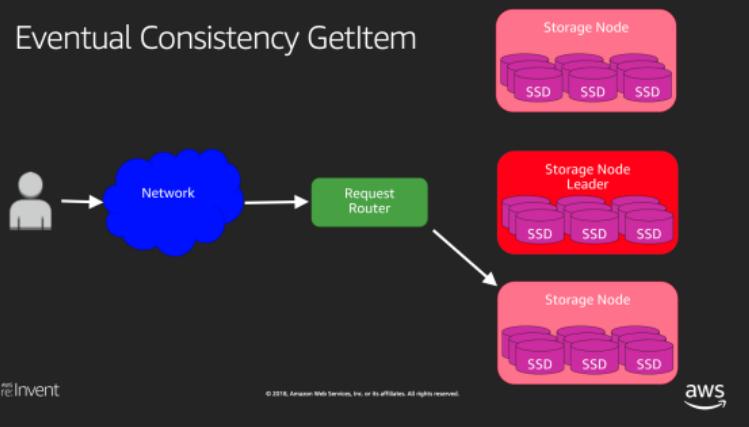
Eventual Consistency GetItem



- Consistent Hashing for partitioning
- Read from single replica.

DynamoDB – Reads

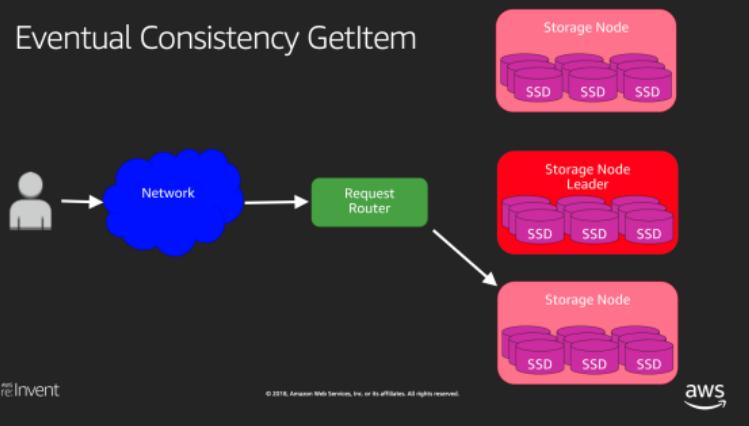
Eventual Consistency GetItem



- Consistent Hashing for partitioning
- Read from single replica.
- $R = 1, W = 2, N = 3$
- $R + W == N$ (eventual consistency)

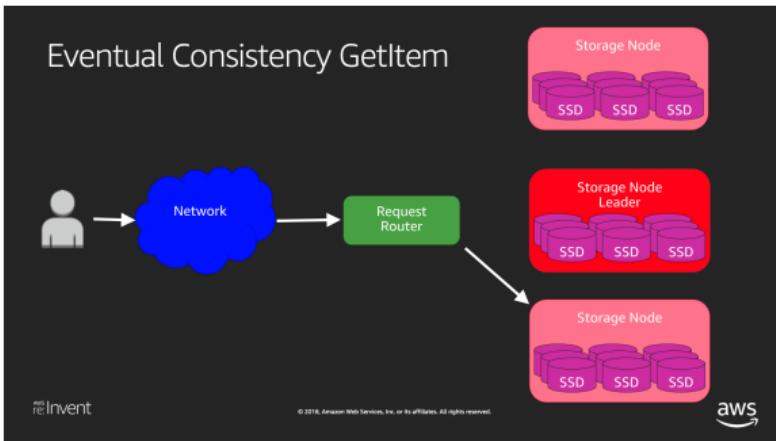
DynamoDB – Reads

Eventual Consistency GetItem



- Consistent Hashing for partitioning
- Read from single replica.
- $R = 1, W = 2, N = 3$
- $R + W == N$ (eventual consistency)
- Reads by default are eventually consistent (flag available in API – ConsistentRead)

DynamoDB – Reads

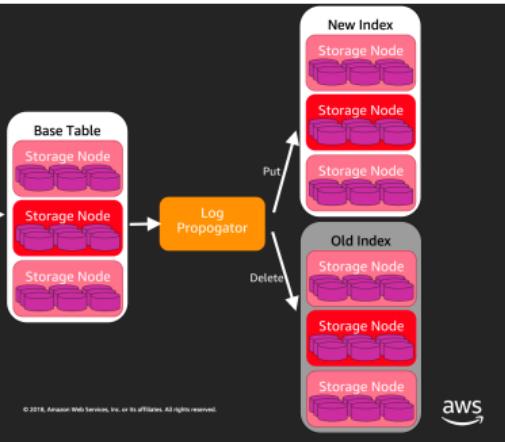


- Consistent Hashing for partitioning
- Read from single replica.
- $R = 1, W = 2, N = 3$
- $R + W == N$ (eventual consistency)
- Reads by default are eventually consistent (flag available in API – ConsistentRead)
- In case of conflicts use Vector Clocks (Last Write Wins / Custom)



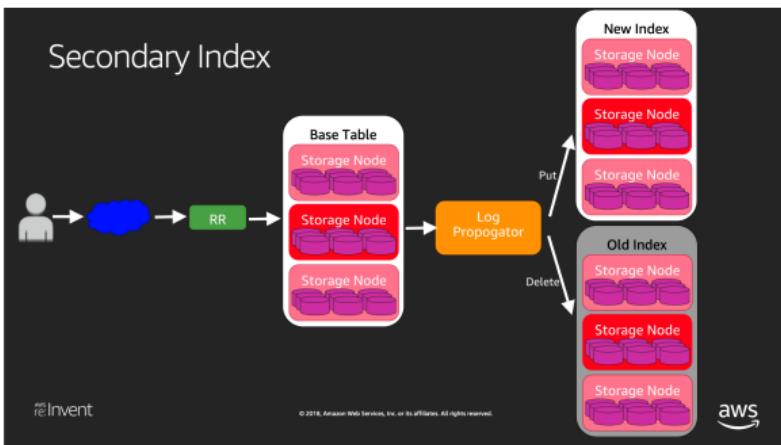
DynamoDB – Secondary Index

Secondary Index



- Independent Table.

DynamoDB – Secondary Index

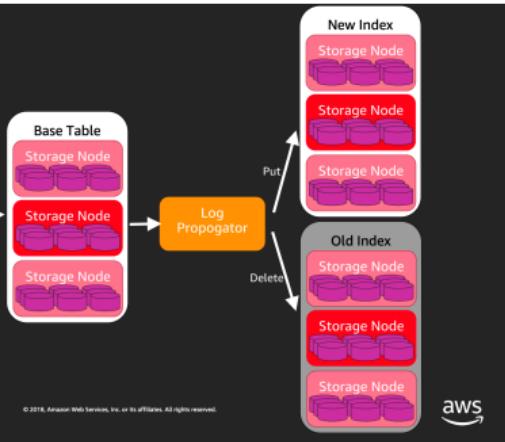


- Independent Table.
- Consistent Hashing for partitioning



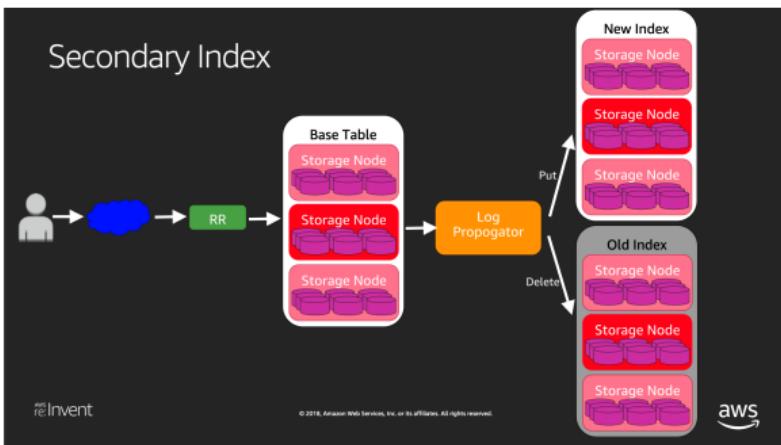
DynamoDB – Secondary Index

Secondary Index



- Independent Table.
- Consistent Hashing for partitioning
- Build from replication log.

DynamoDB – Secondary Index

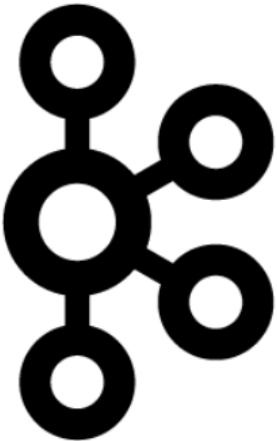


- Independent Table.
- Consistent Hashing for partitioning
- Build from replication log.
- Write Amplification (Specially for updates)

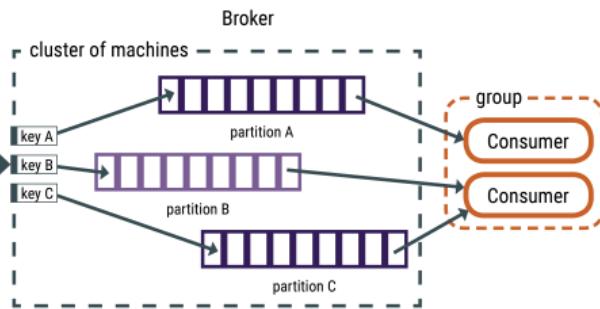


For more please refer to **DAT321** talk from AWS re:Invent!
<https://www.youtube.com/watch?v=yvBR71D0nAQ>



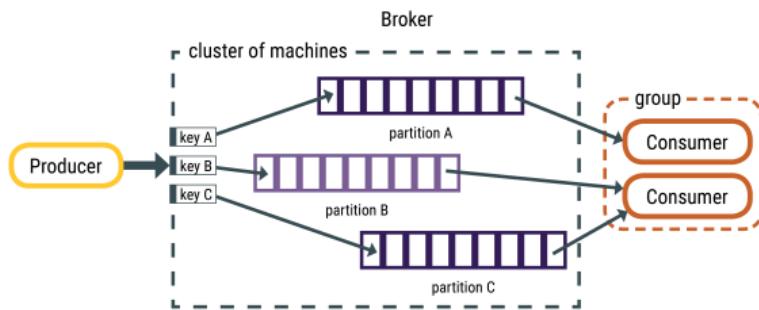


kafka



- Uses Replicated Logs
- Configurable replica count (topic creation / updatable)
- Configuration partition (topic creation / updatable)
- Configurable write Acknowledge (producer configuration)

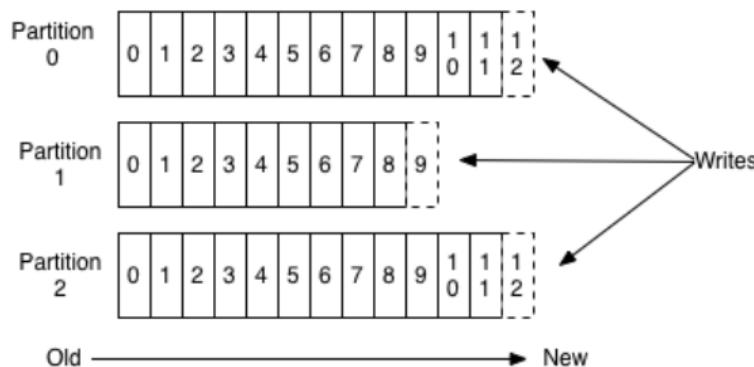




- Uses Replicated Logs
- Configurable replica count (topic creation / updatable)
- Configuration partition (topic creation / updatable)
- Configurable write Acknowledge (producer configuration)
- Typically topic has a retention value (TTL)
- We can store values indefinitely by using Log Compaction.

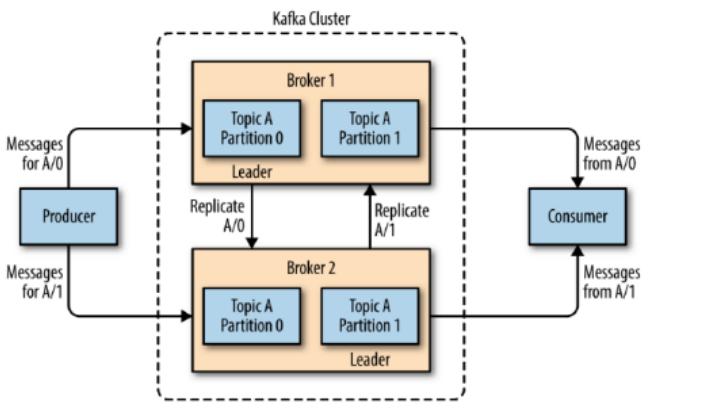


Anatomy of a Topic



- Each Topic (Queue) comprises of multiple Partitions

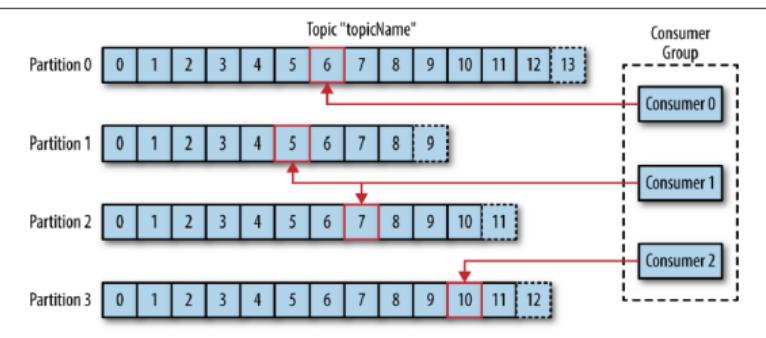
Kafka - Topics



- Each Topic (Queue) comprises of multiple Partitions
- Each Broker is responsible for multiple partitions.
- Each partition has 1 leader
- Leader election (Zookeeper/raft quorum)
- All reads/writes go to the leader
- Non leader replicate data from leader



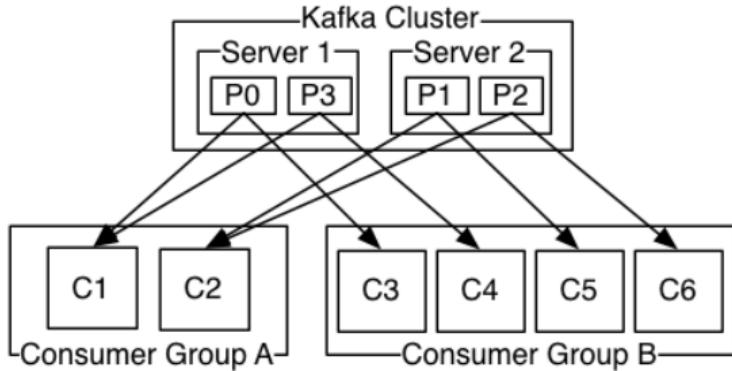
Kafka - Topics



- Number of Consumers \leq Number of Partitions
- Ordering guarantees apply only within the partition



Kafka - Topics



- Number of Consumers \leq Number of Partitions
- Ordering guarantees apply only within the partition
- Multiple consumer groups can read same topic



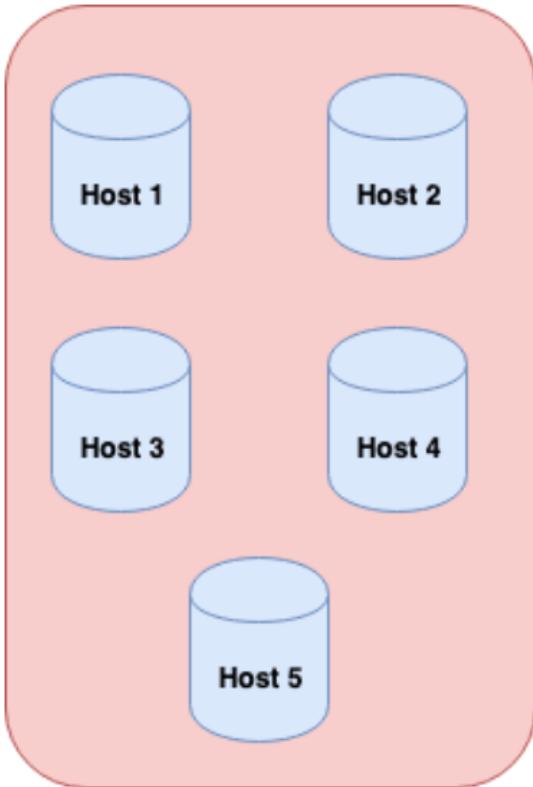
For more please refer to the Excellent Kafka Documentation!
<https://kafka.apache.org/documentation/#design>



Is there another way multiple nodes can maintain data consistency?



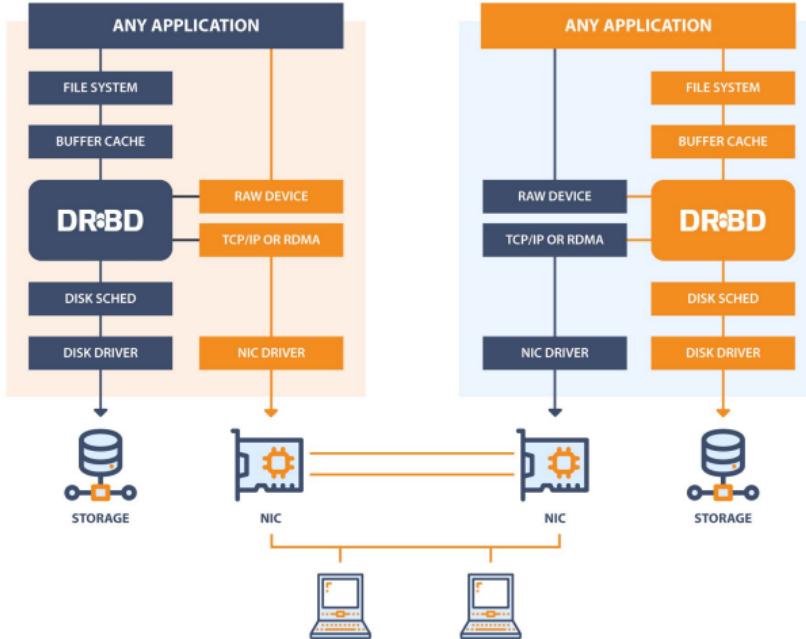
Replicated Storage



- Replication of writes is needed for data consistency
- Can we replicate stored data instead ?

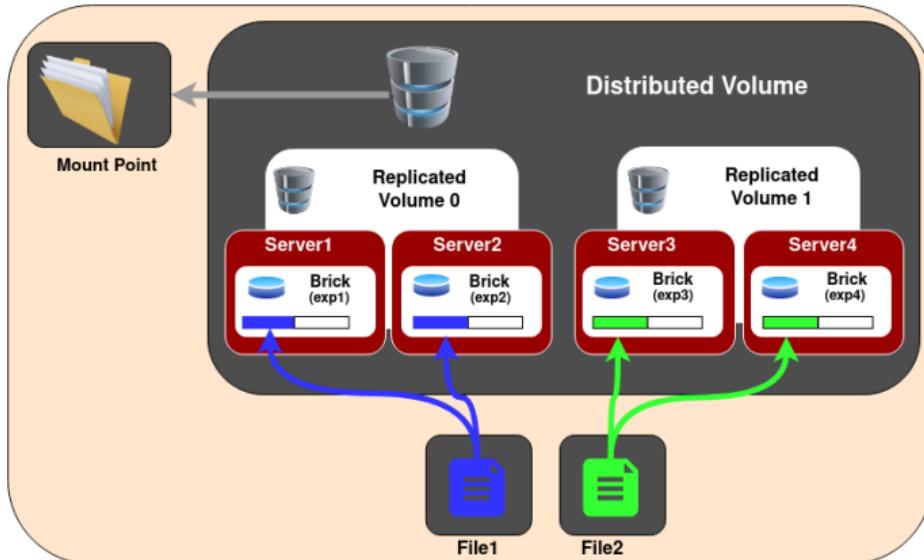


Replicated Storage



- Replication of writes is needed for data consistency
- Can we replicate stored data instead ?
- DRBD

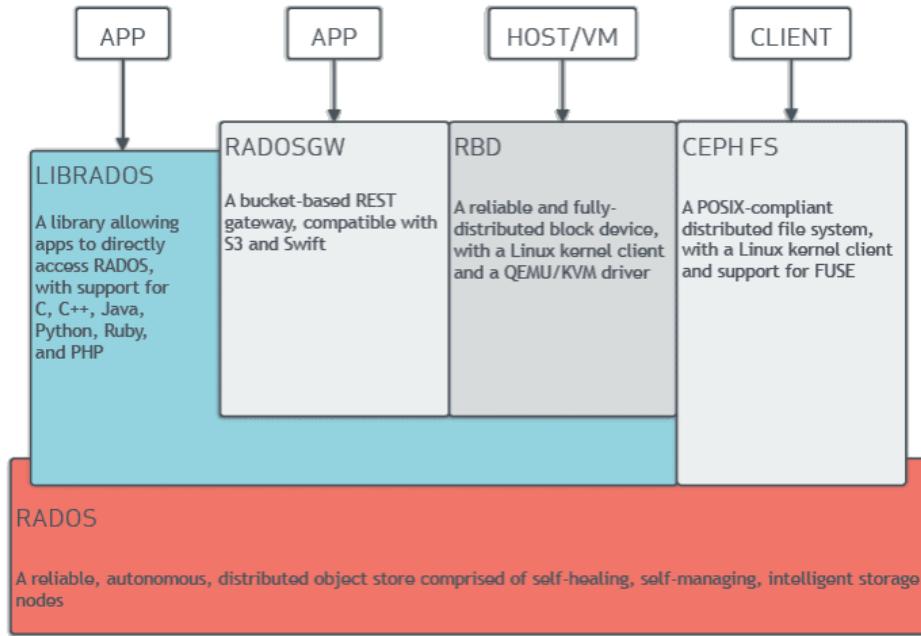
Replicated Storage



- Replication of writes is needed for data consistency
- Can we replicate stored data instead ?
- DRBD
- GlusterFS



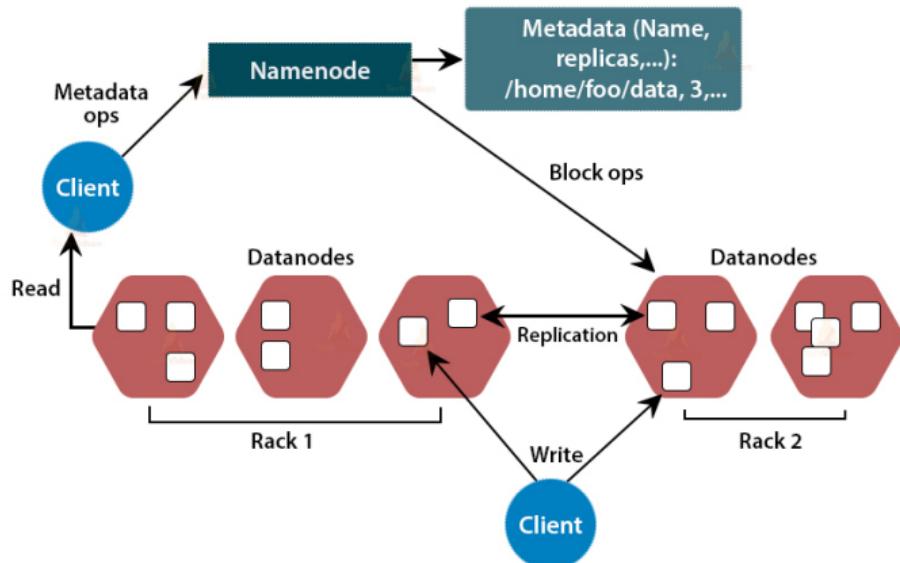
Replicated Storage



- Replication of writes is needed for data consistency
- Can we replicate stored data instead ?
 - DRBD
 - GlusterFS
 - Ceph.io



HDFS Architecture



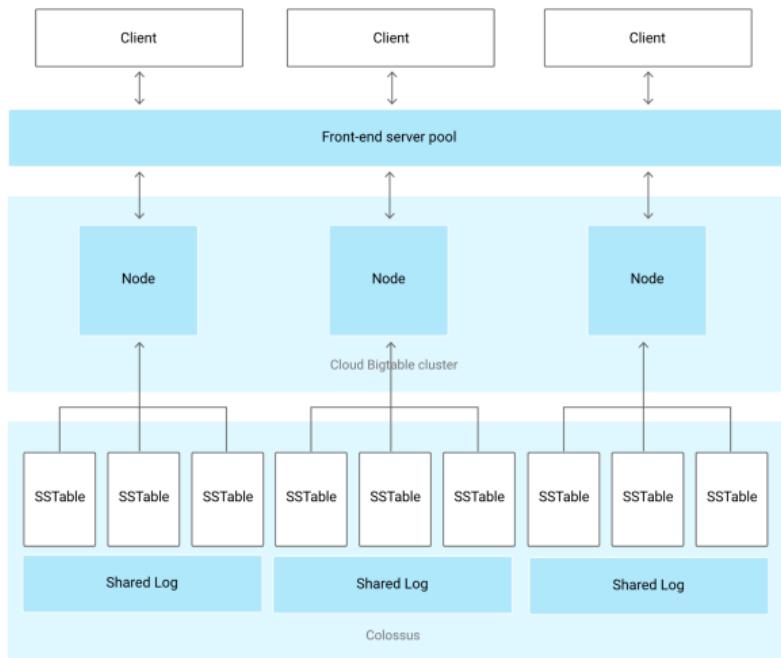
- Replication of writes is needed for data consistency
- Can we replicate stored data instead ?
- DRBD
- GlusterFS
- Ceph.io
- HDFS





BigTable

BigTable

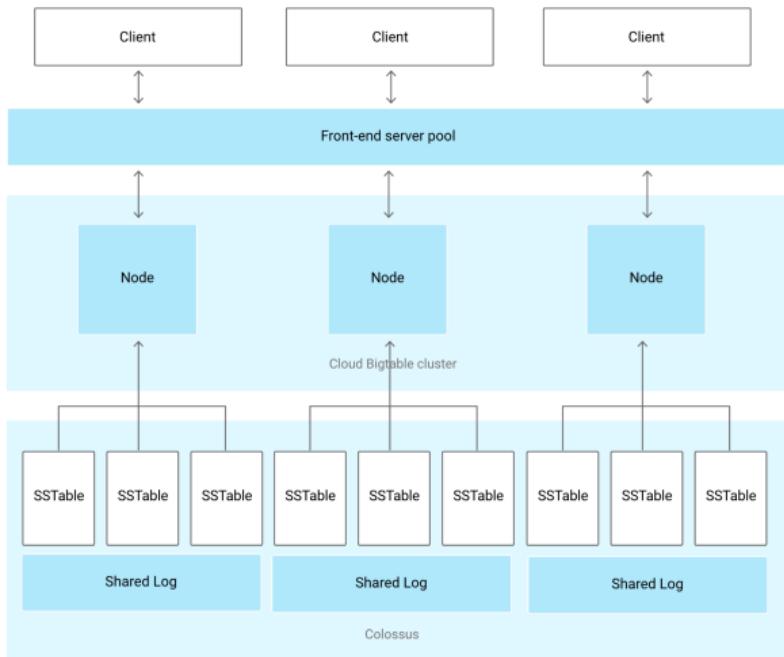


- LSM Based
- Uses dedicated request routers



gojek

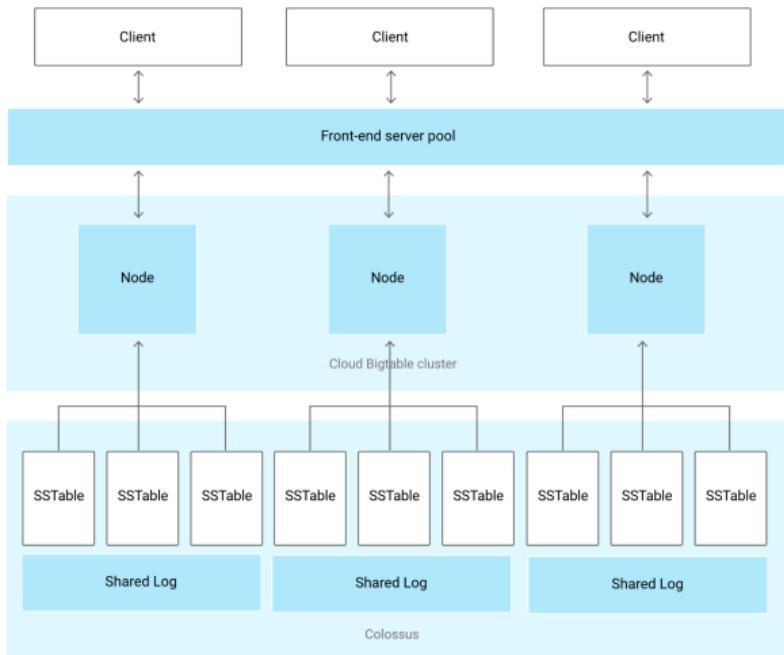
BigTable



- LSM Based
- Uses dedicated request routers
- Key ranges are partitioned across nodes (Tablet)
- Each Node is responsible for multiple Tablets
- No replica Nodes/Tablets (strong consistency)



BigTable

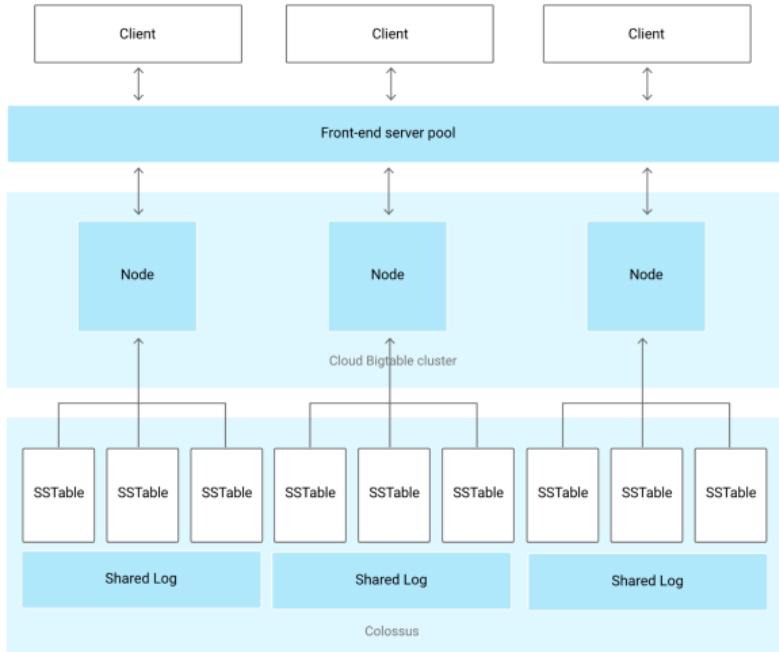


- LSM Based
- Uses dedicated request routers
- Key ranges are partitioned across nodes (Tablet)
- Each Node is responsible for multiple Tablets
- No replica Nodes/Tablets (strong consistency)
- Clusters in single AZ
- multi-cluster can be AZ across regions (multi-master)
- Uses Google's Internal distributed file system - Colossus (replaces GFS)



gojek

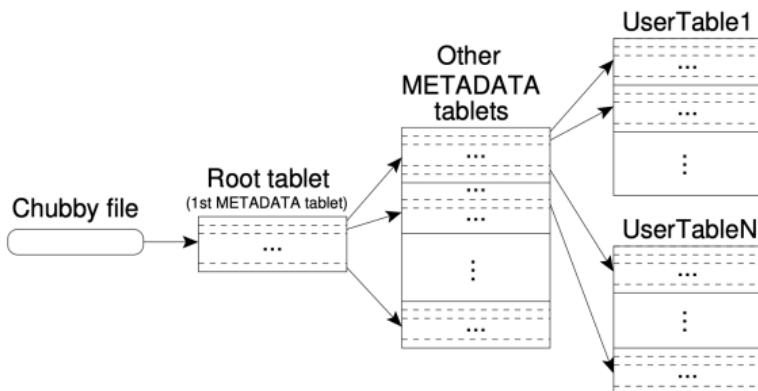
BigTable - Master



- Bigtable master manages nodes/tablets (does not serve traffic)
- **BigTable Master** is responsible for
 - Tablet assignment to server
 - detecting the addition and expiration of tablet
 - balancing tablet-server load
 - garbage collection of files
 - table and column family creations



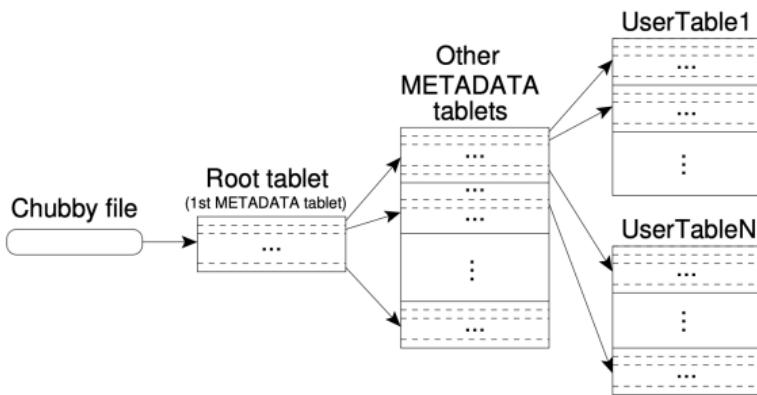
BigTable - Tablets



- Multiple METADATA tablets (B+ Tree style)
- Tablets contain associated row range (e.g. aardvark...apple)
- Multiple tablets may map to a single SSTable



BigTable - Tablets



- Multiple METADATA tablets (B+ Tree style)
- Tablets contain associated row range (e.g. aardvark...apple)
- Multiple tablets may map to a single SSTable
- Initially, each table consists of just one tablet.
- As a table grows, it is split into multiple tablets
- Each tablet is 100-200 MB in size by default.
- tablets can be pre-split at table creation



BigTable - App Profile

[Edit application profile stateful-pricing-multi](#)

Instance ID	stateful-pricing
Application profile ID	stateful-pricing-multi

Description

Cluster routing

Manage how incoming requests are routed.

Single-cluster

Routes to one cluster

Multi-cluster

Routes to nearest available cluster

Cluster group

Select all

stateful-pricing-asia-east1-b
asia-east1-b

stateful-pricing-asia-east1-a
asia-east1-a

Any cluster

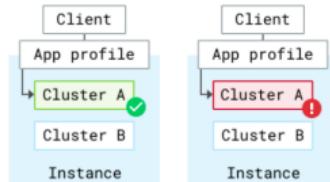
[CANCEL](#) [OK](#)

- Application profile defines how incoming requests are routed to clusters



BigTable - App Profile

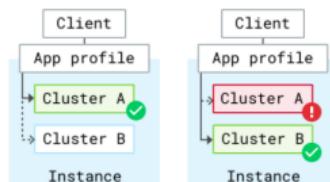
Single-cluster



1. Request goes to Cluster A, the specified cluster.

2. If Cluster A is unavailable, the request fails. The app profile must be manually re-routed.

Multi-cluster



1. Request goes to the nearest available cluster. Here, it's Cluster A.

2. If Cluster A is unavailable or a request fails or is delayed, the request is resent to Cluster B

- Application profile defines how incoming requests are routed to clusters
- Single-cluster profiles always route to the selected cluster
- Multi-cluster profiles can route to any available cluster



BigTable - App Profile

Routing comparison

Single-cluster routing	Multi-cluster routing
Routes to one cluster	Routes to nearest available cluster
⚠ Manual failover	Automatic failover
Helps isolate CPU-intensive loads	Higher availability
Read-your-writes consistency	⚠ No read-your-writes consistency
Single-row transactions supported	⚠ No single-row transactions

- Application profile defines how incoming requests are routed to clusters
- Single-cluster profiles always route to the selected cluster
- Multi-cluster profiles can route to any available cluster
- Single-cluster profiles unlock strong consistency
- Multi-cluster profiles have higher availability



BigTable

Read the BigTable Paper at <https://research.google/pubs/pub27898/>

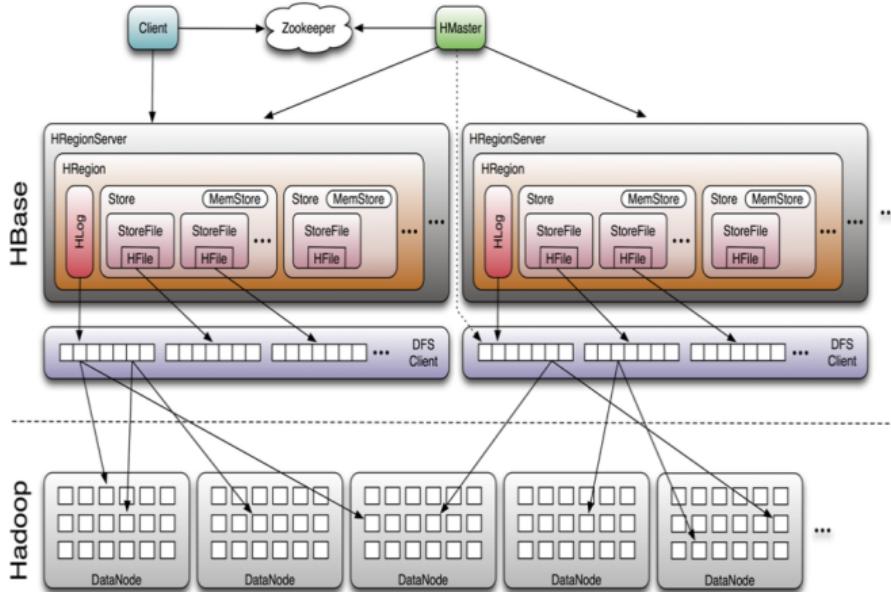
You can also read the BigTable Documentation!

<https://cloud.google.com/bigtable/docs/overview>





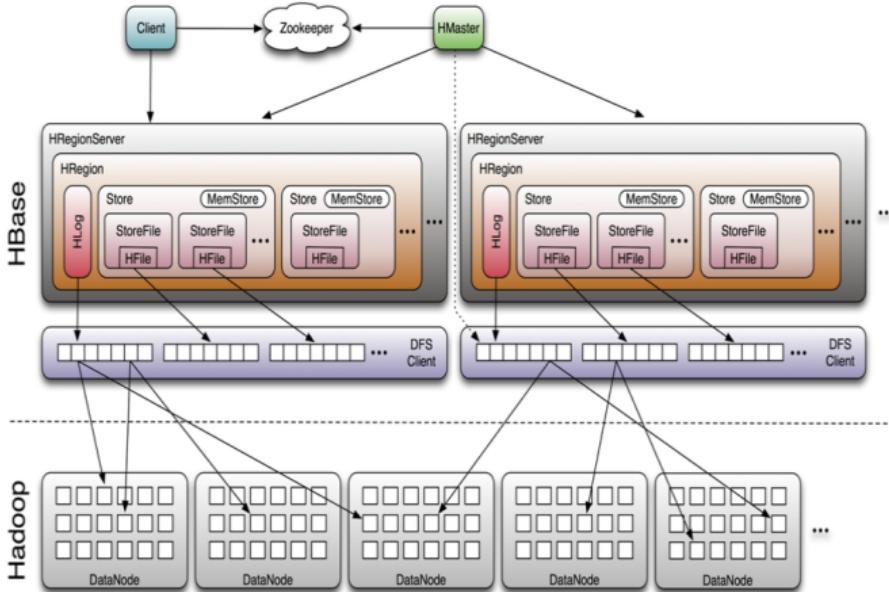
HBase



- Inspired by BigTable Paper
- LSM Based
- Uses HDFS for storage
- Open Source



HBase



- Inspired by BigTable Paper
- LSM Based
- Uses HDFS for storage
- Open Source
- HMaster - similar to BigTable master node
- RegionServer - similar to BigTable node, responsible for multiple regions
- Region - similar to BigTable tablet



For more please refer to the Excellent HBase Documentation!

<https://hbase.apache.org/book.html>



Q&A



Pricing is hiring L2/L3/L4
Internal Mobility and Referrals

