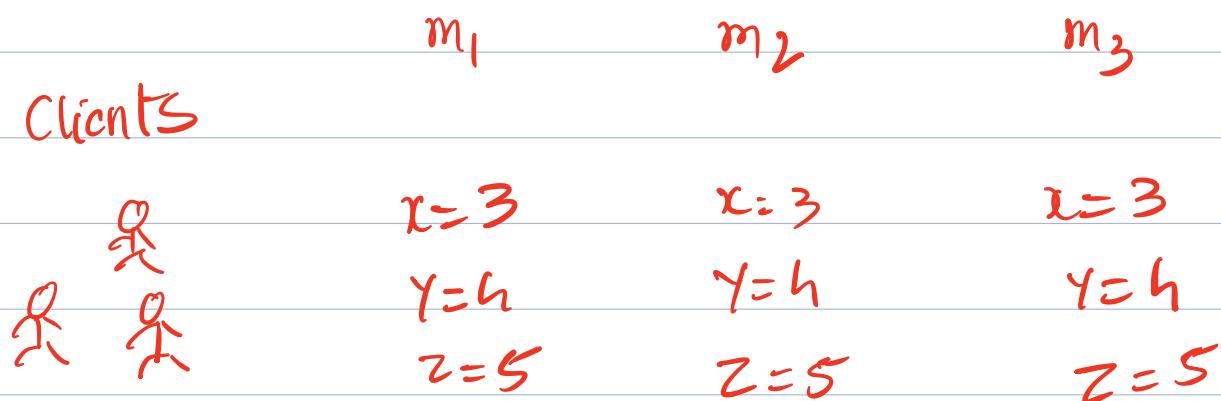


Agenda:

- Intro to sharding (data partitioning)
- consistent hashing

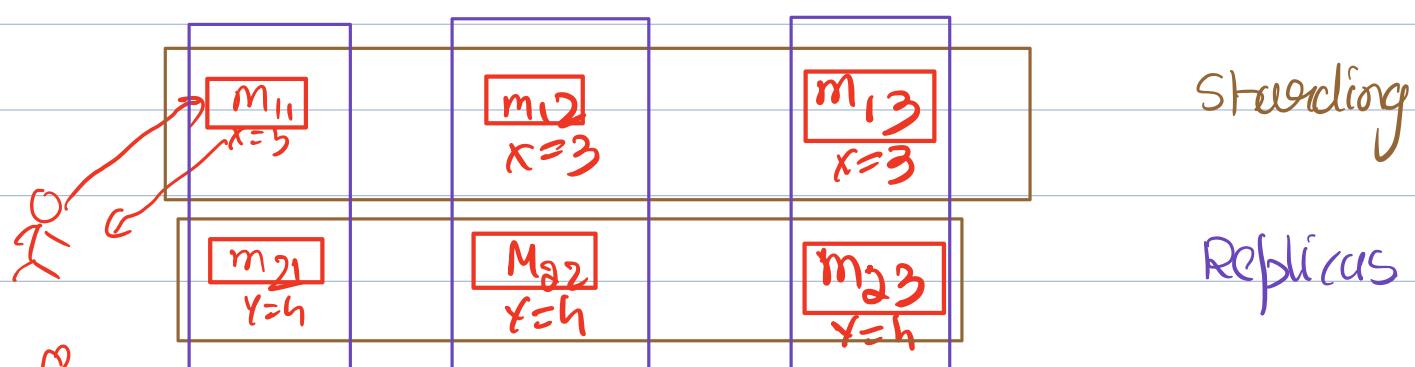
Riak

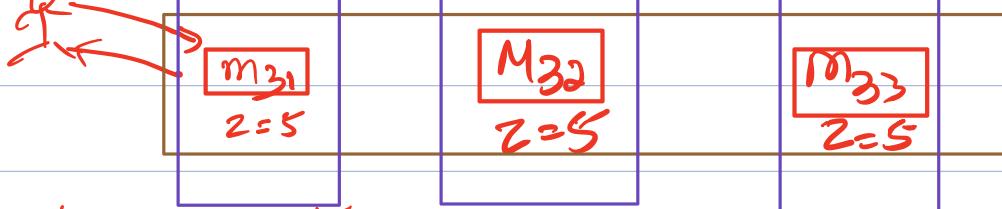
- Inspired by Dynamo DB
- Open source clone of Dynamo



What's wrong with an approach where everyone stores all of the data?

- What if you have more data than can reasonably fit on one machine?
- If everyone stores all data, consistency is more expensive to maintain.





- Lesser cost
- Split data across machines
- Save effort of keeping consistent
- Distribute load

What did you lose by changing design?

- No fault tolerance.

Sharding & Replication are orthogonal to each other

Also called data partitioning.

Why would you do this?

- capacity increases.
- throughput increases

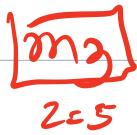
Sharding

How do we decide how to split up our data among shards?



What makes a sharding strategy good?

gzh



Randomly assign data to nodes?

✓ - evenly distribute data

✗ - hard to find what we want

Put it all in one place?

✓ - easy to find it

✗ - the worst at even distribution

Every partition stores a range of keys

✗ - Distribution might be uneven

How do you take data that is not uniformly distributed & make it uniformly distributed?

Hashing!

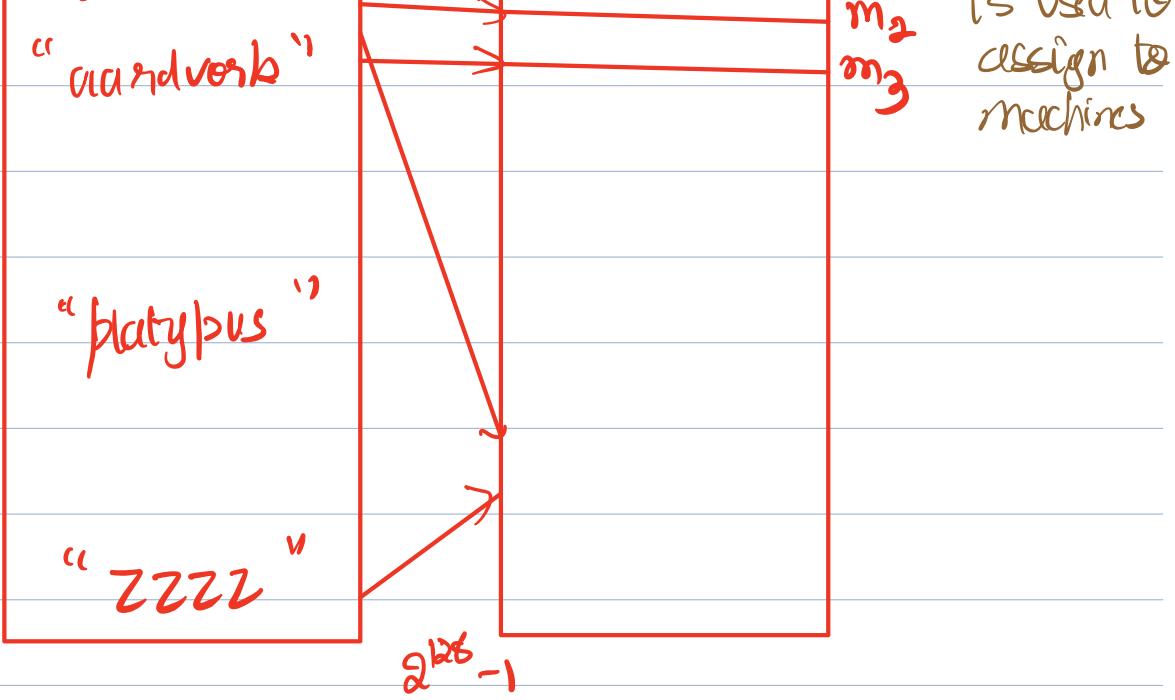
Hashing

MD5: input → output

"aaaaa"

0

Assuming hash % 3 is result
m1 m2 m3



m_1 "aaaa"

m_2 "aardvark"
"apple"
"zzzz"

m_3 "awkward"
"platypus"

Positioning by hash
of key $\% N$,
where $N = \# \text{ of}$
partitions.

What is the catch with the above approach?

- Assumes that N is fixed!
- If N changes, most of the data needs to be moved around!

Assume, $MDS(\text{"aardvark"}) \% 3 = 1$

$MDS(\text{"aardvark"}) \% 4 = 2$

is used
design to
machines

Even if there is a new machine, the key is moved to another old machine !!

∴ Too much data movement when N changes

What is the minimum acceptable movement?

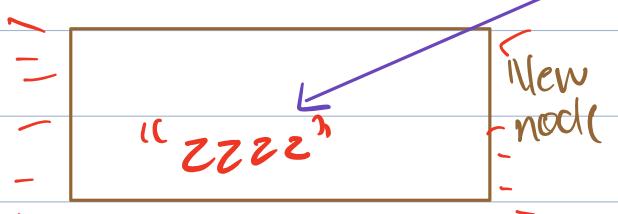
If there are 6 keys split across 3 machines, on average, each machine should have 2 keys
 $\frac{6}{3}$ machines $\Rightarrow 1.5$ keys

"aaaaa"
"platypus"

"apple"
"awkward"

"aaandwork")
"zzzzz"

$\frac{6}{3}$ about 1.5 per node

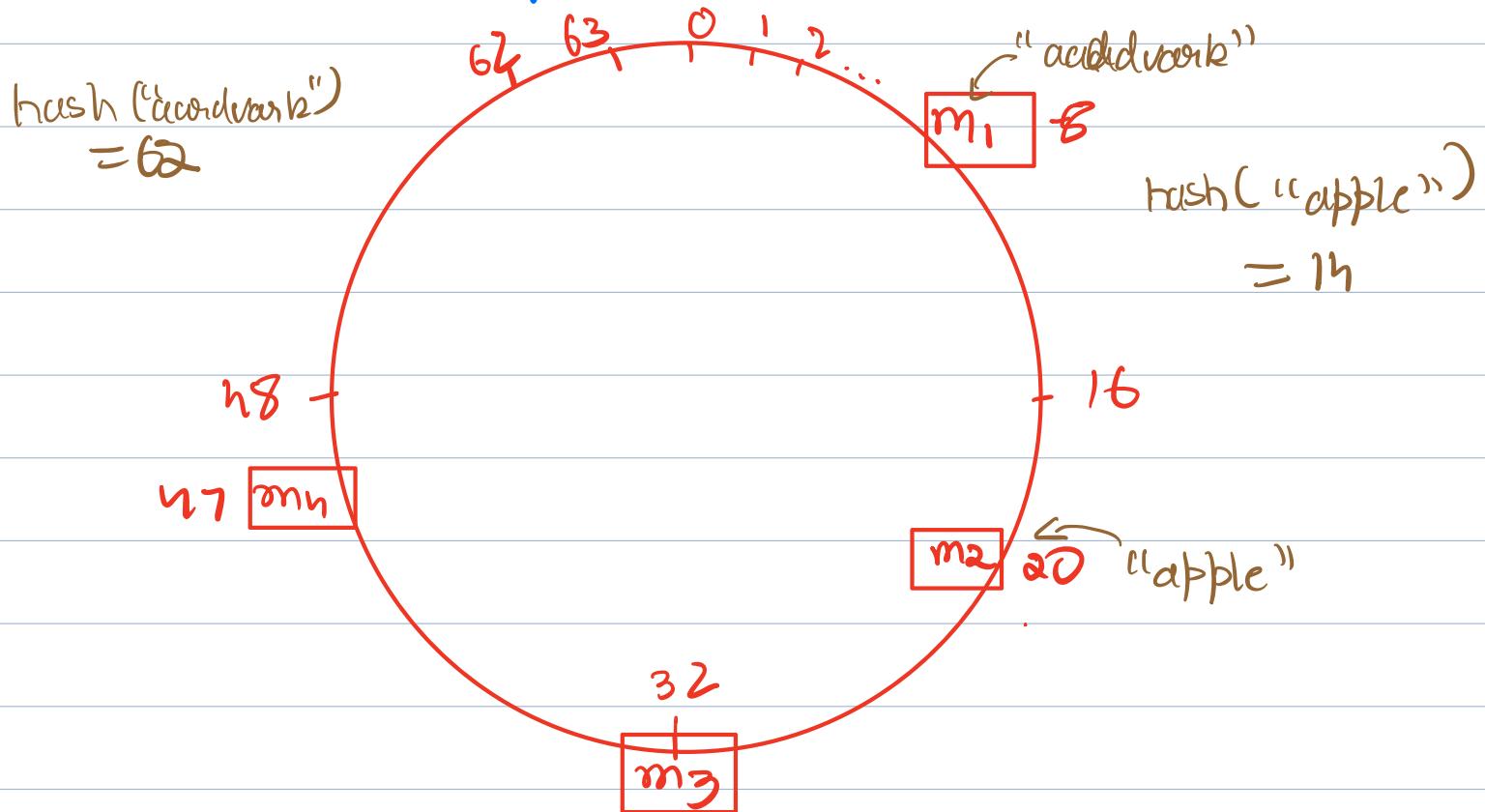


For K keys & N nodes, $\frac{K}{N}$ is the minimum

movement possible to get an even split

"Consistent Hashing" gives us this minimum movement

Consistent Hashing $2^b - 1$



- Nodes are organized around a ring
- Assuming max no. of slots on the ring is $2^b - 1$ i.e. The hash function used gives $2^b - 1$ values, each node exists on a particular slot
- How do you determine the slot for each node?
 - Maybe you can hash the IP address of the node
- Suppose you want to store the key "apple". After hashing, we get the value 14
- Since there is no node on 14, we will walk the ring clockwise & store it on the first node we see. Hence, "apple" is stored on

node M₂

How does Dynamo use consistent hashing for replication?

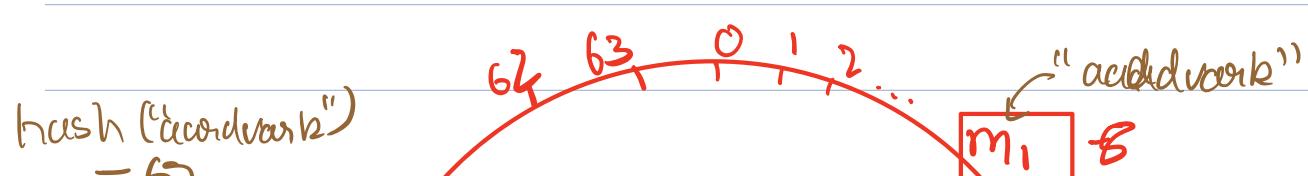
- Assuming the number of replicas is 3.
- The original hash value of the key is used to store the data in a node called the Home Node.
- Since the replicas is 3, we will walk the ring clockwise & store the key in the next 2 nodes. This will be m₃ & m₁

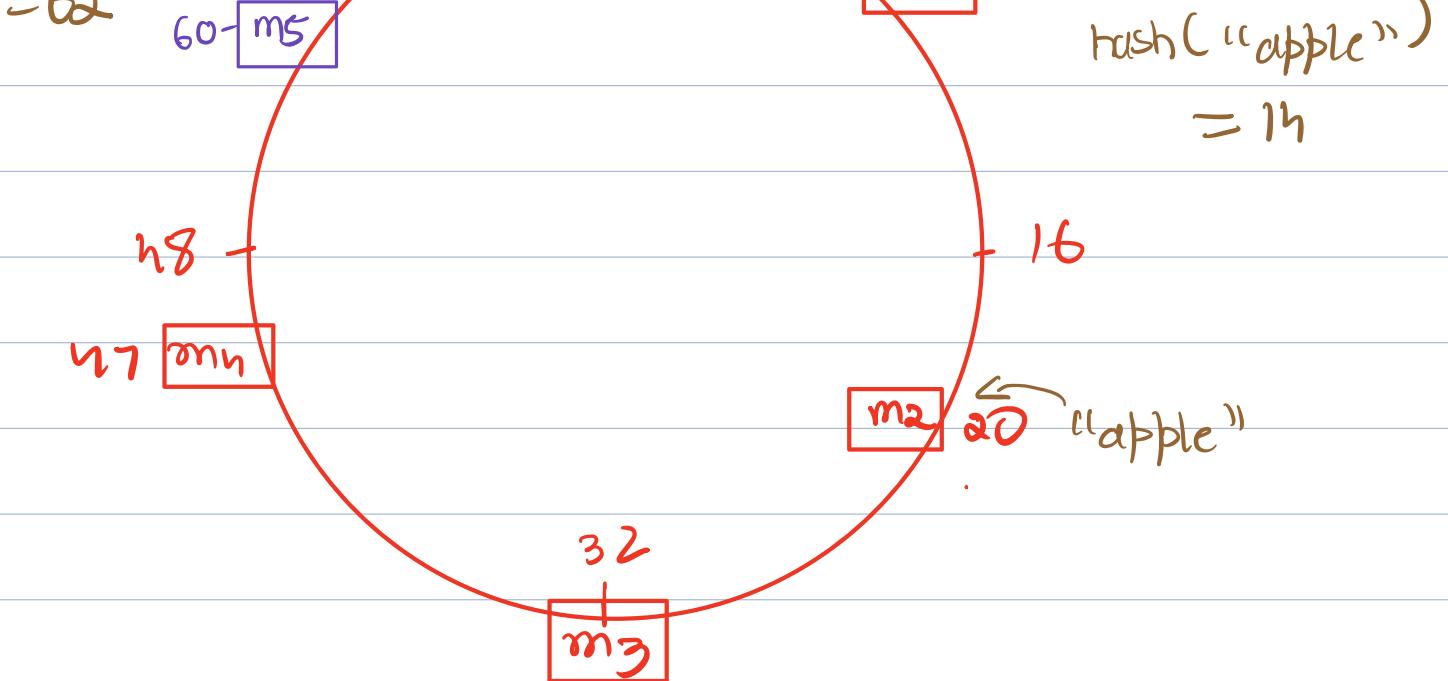
In dynamo, every node has a preference list i.e. list of nodes where the data should be replicated.

Generally the preference list is longer than the replication factor why?

Some nodes could be unavailable for some reason.

How does consistent hashing help reduce data movement?





New node m_5 at 60. Which keys do we have to move?

If $\text{hash}("zzzz") = 50$, before m_5 , it would go to m_1 ,

After m_5 , m_1 is responsible for smaller range.

\therefore Anything that would have hashed to m_1 initially when the value fell in the space for m_5 ($n_8 - 60$) would have to be moved to m_5

However, on average, $\frac{k}{N}$ keys will be moved

All other nodes don't care about this data movement!!

What happens if a node crashes?

Suppose M_2 crashes

M_2 was responsible for range (9-20).
Since M_2 is gone, M_3 needs to store more keys

However, all other nodes are unaffected.

Hopefully, M_3 already has one backup of these keys, and can act as temporary home node for these keys before M_2 replacement is up.

How/When can consistent hashing go wrong?

- Human needs to monitor system health & ensure not too many nodes crash.
- If couple of nodes crash, neighbor has to pickup lot of slack.
- In the case where nodes are not evenly spaced, the OIP space may not be evenly covered.

Solution?

Hash a node to many spots)
Same node, but it is available on many
spots. These nodes are called virtual nodes.

Advantages

- This makes it more likely that nodes are evenly distributed on the ring.
- Virtual node count can be proportional to the capacity of each node.

M_1 has 10 TB space \Rightarrow no virtual nodes

M_2 has 1 TB space \Rightarrow 16 virtual nodes

Helps categorize heterogeneous machines

Disadvantages:

- If M_1 goes down, virtual nodes are down too. Lots of other nodes need to pickup slack from other hosts on the ring.
In Dynamo, this is stated as an advantage as data movement happens between different machines.