# CSE138   Lecture 16

- intro to sharding
  (aka data partitioning)
- consistent hashing

### Clients

| $M_1$ | $M_2$ | $M_3$ |
|-------|-------|-------|
| $x=3$ | $x=3$ | $x=3$ |
| $y=4$ | $y=4$ | $y=4$ |
| $z=5$ | $z=5$ | $z=5$ |

One issue with storing all data on all replicas:

- What if there's more data than can fit on one machine?

### Clients

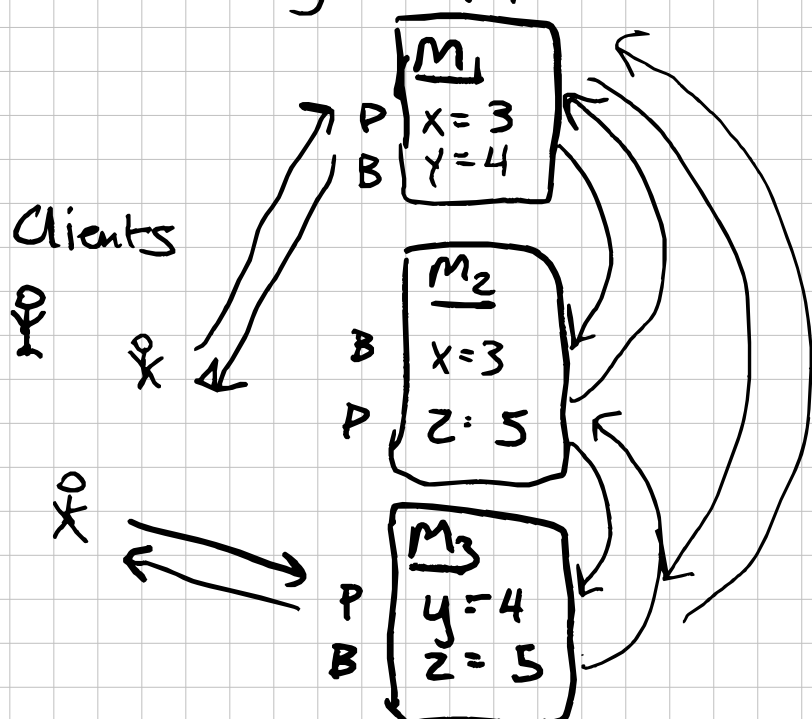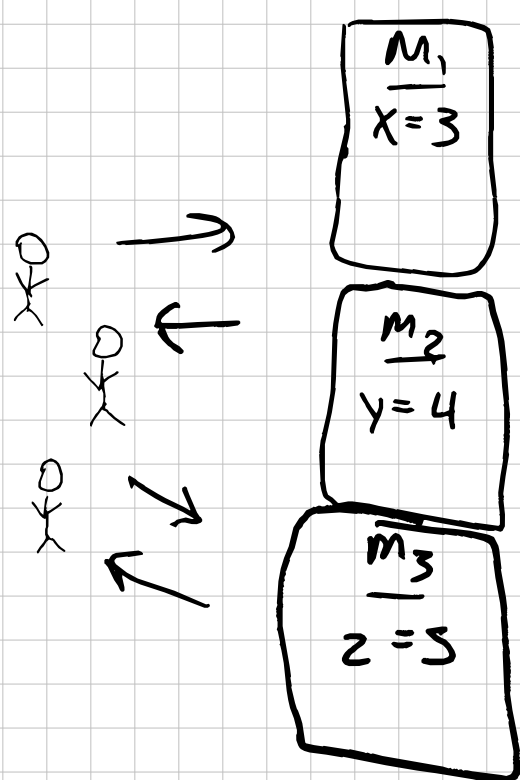| $M_{11}$ | $M_{12}$ | $M_{13}$ | } shard 1 |
|----------|----------|----------|-----------|
| $x=3$ | $x=3$ | $x=3$ | |
| $M_{21}$ | $M_{22}$ | $M_{23}$ | } shard 2 |
| $y=4$ | $y=4$ | $y=4$ | |
| $M_{31}$ | $M_{32}$ | $M_{33}$ | } shard 3 |
| $z=5$ | $z=5$ | $z=5$ | |

Shards (different data on each) are rows

replicas (same data on each) are columns.

Reasons to do Sharding:

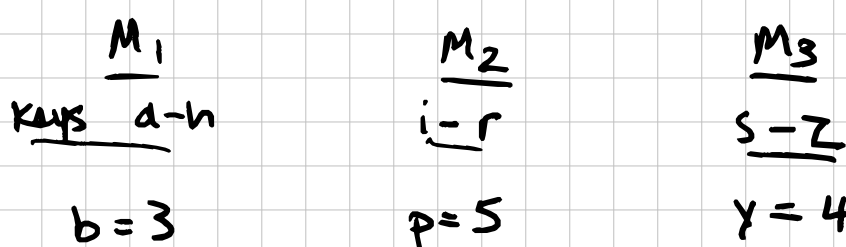* — Store more data than can fit on one machine
* — improving throughput

| $M_1$ | |
|---|---|
| P | $x=3$ |
| B | $y=4$ |

### Clients

| $M_2$ | |
|---|---|
| B | $x=3$ |
| P | $z=5$ |

| $M_3$ | |
|---|---|
| P | $y=4$ |
| B | $z=5$ |

Focus on partitions:

| $M_1$ |
|-------|
| $x=3$ |

| $M_2$ |
|-------|
| $y=4$ |

| $M_3$ |
|-------|
| $z=5$ |

Goals for our partitioning strategy:

- Evenly spread data across the nodes.
- Make it fast and easy to find data we want!

Partitioning by key range:

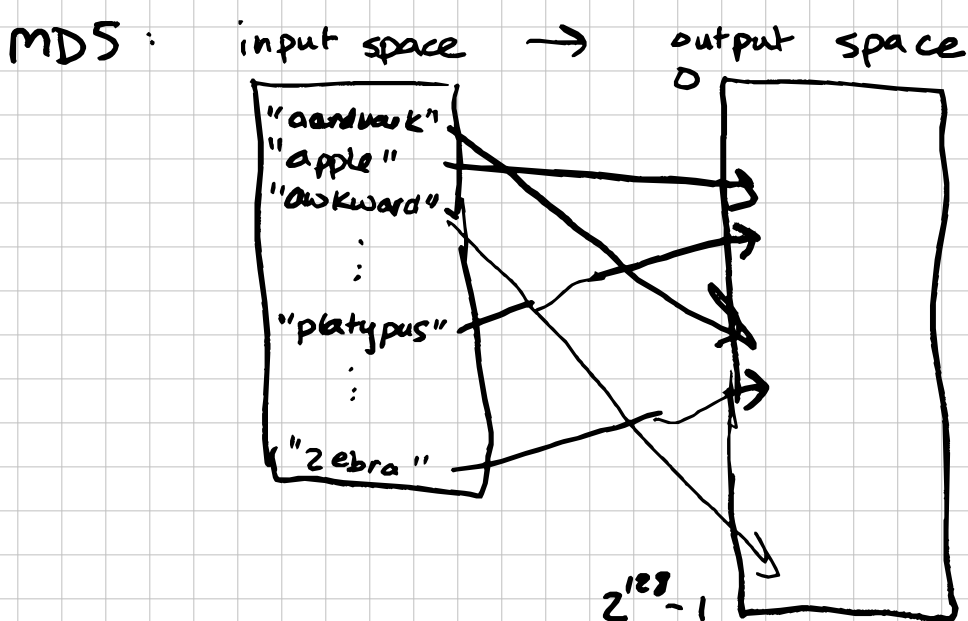| $M_1$ | $M_2$ | $M_3$ |
|-------|-------|-------|
| Keys a-h | i-r | s-z |
| $b=3$ | $p=5$ | $y=4$ |

$$\begin{cases} s = 10 \\ z = 11 \\ x = 12 \end{cases}$$

Probably bad unless you know that your keys are uniformly distributed.

But there's hope... hashing!

MD5 : input space → output space



"aardvark"
"apple"
"awkward"
⋮
"platypus"
⋮
"zebra"

0

$2^{128} - 1$

$$hash(key) \mod N,$$

where N is the number of
Nodes you're partitioning onto.

"aardvark" = [ data blob ]

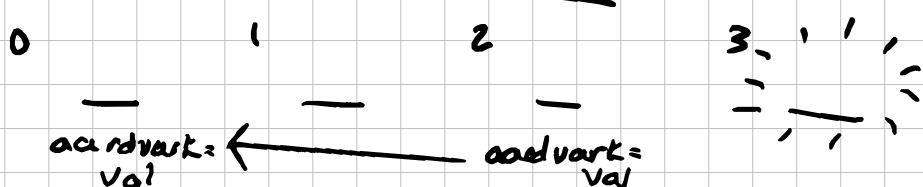$$hash(\text{"aardvark"}) \mod N = 2$$

This lets us split up the data
more or less evenly.

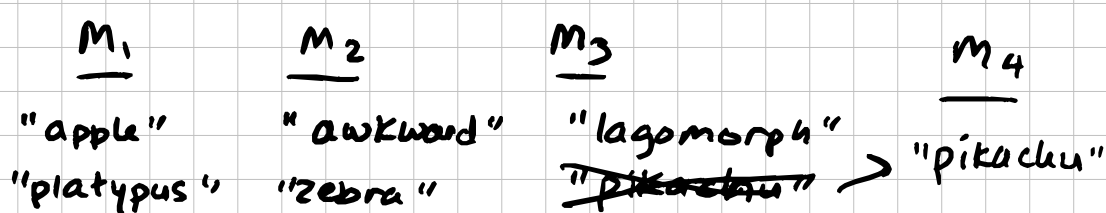But there's a catch...
What if the number of nodes (N) changes?

Data that doesn't need to move
will get moved around! "

$$hash(\text{"aardvark"}) \mod 3 = 2$$
$$hash(\text{"aardvark"}) \mod 4 = 0$$

0         1         2         3 ' ' '
—         —         —        =/ /
aardvark =  ←————————  aardvark =
val                    val

e.g.,
Some data will move from an old
machine to a different old machine
when we add a new machine.
    Seems bad.

M₁              M₂              M₃                       M₄

"apple"        "awkward"       "lagomorph"
"platypus"     "zebra"         ~~"pikachu"~~  >  "pikachu"

$\frac{6}{4}$ or 1.5 keys per node. 1 should move.
                              (or 2)

in general, if you have
        K keys and N nodes,

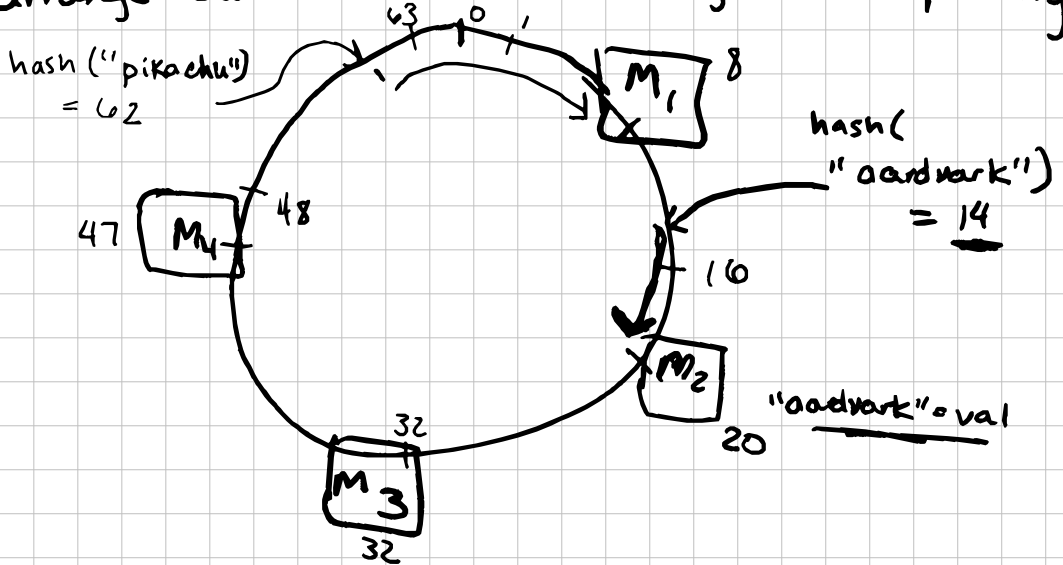$\boxed{\frac{K}{N}}$ is the number of keys that
                should move when
                a node joins or leaves.

                              invented for
To accomplish this,
we use [ consistent ] hashing! ← CDNs
            ↑                        in the
    is totally different from        90's
        consistency models e.g.    (1997 paper)
            causal consistency;    Akamai
        totally different from
            consistent snapshots.

How does consistent hashing work?
arrange our nodes in a "ring" conceptually.

hash ("pikachu")
= 62



hash(
    "aardvark")
    = 14

"aardvark" → val

output space of our hash function
    goes from 0 to 63 in this example.

both the keys and the node names get hashed,
    and nodes get located at a point
    on the ring.