# Dynamo

=> early
key-value
storage

⎧ Goals
⎨ _____
⎩ System Arch
  (most)

# Dist Sys

Theory
Academic ———> Practice

'00 ——> '10

↗web
services

Goals:                  ↪ shopping carts
⎧ —> **Scale** ←
hard ⎨ —> Efficiency/Latency
⎩                    (tail)
     —> Availability  (over all other things)

⎧ —> Simple API : get/put
⎩                    keys/value
   => rise of key/value stores
~"NoSQL movement" ⟷ "SQL"
                         ⎡ two ⎤
                         ⎣ religions ⎦

# Background:

$\rightarrow$ "eventual consistency"

## Availability:

"always writeable" $\rightarrow$ make writes highly available on reads

$\Rightarrow$ allow writes conflict resolution

$\Rightarrow$ how to resolve conflicts?

$\rightarrow$ system
e.g. last
writer
wins
(overly
general)

$\rightarrow$ app can
$\rightarrow$ & be
smarter

# System Architecture:

## API: get/put
$\Rightarrow$ shopping cart
$\Rightarrow$ is this interesting? (adv.)

## Partitioning
$\Rightarrow$ consistent hashing
(but not full thing a la Chord paper)



key

(use virtual nodes)
+ ~~rate~~

Replication: N hosts
1 of N: coordinator
put -> make sure key properly replicated

"preference list": nodes resp.
for key k
N-way replication,
pref list > N: why?
-> N distinct phys nodes
-> may be some down

"Eventual" Consistency: replicas
not always
in sync.

put new
$c_1$: get -> new
$c_1$: get -> old }

support for "shopping cart"
-> never lose "add to cart"
adds/deletes => puts (writes)

how to track versions/conflicts:
version vectors

↳ ⇒ easy conflict detection

↳ standard problem: large (can be)
limit size to 10

---

Executing get/put :

1) load balancing

→ generic load balancer (might need extra hop)

→ partition-aware client library

↳ why better?

2) coordinator :

first N healthy nodes
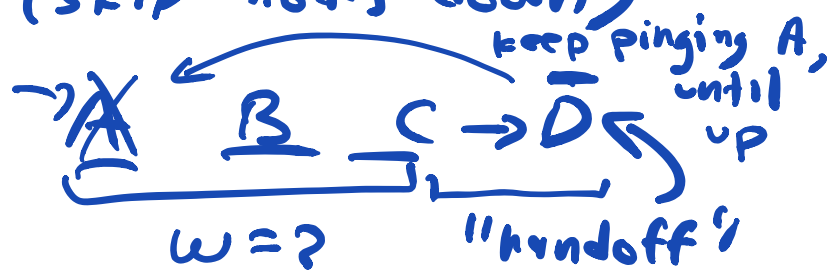(w/ failures → go further
down pref
list)

Quorum-based : R, W
if R+W > N, intersection
but: availability
instead "sloppy" quorum:

=> first N healthy nodes
(skip nodes down)



keep pinging A, until up

A ⤫  B  C → D

w=?     "handoff"

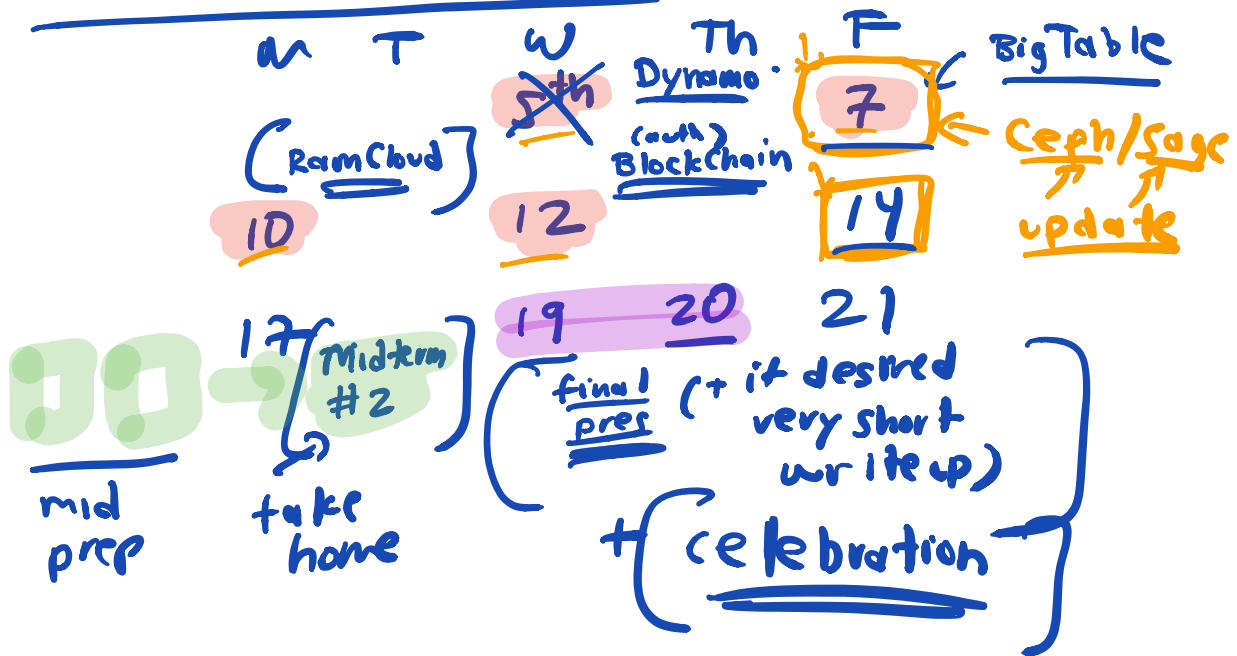## Replica Synchronization

=> Anti-entropy { Random pairwise sync }

=> Detect/fix differences:
[ Merkle trees ]

problem: classic "consistent hashing"

=> key ranges keep changing w/
node addition/deletion

(Break)

# Ending the class:

| M | T | W | Th | F | |
|---|---|---|---|---|---|
| | | 5th | Dynamo | 7 | BigTable |
| (RamCloud) | | | (auth) Blockchain | 14 | Ceph/Sage update |
| 10 | | 12 | | | |

17 [Midterm #2]   19   20   21

mid prep   take home   final pres (+ if desired very short write up)

+ (celebration)

# Dynamo: last tech piece

## membership:
join/leave

worried partitioned: solve by using "seed" nodes

## failure detection:
mostly to "local" failure detection

e.g., coordinator → replicate write

$$x_1, \ \cancel{x_2} \ldots \to x_N$$

historical: p2p          not p2p:
                            smaller sak

Experience / Lessons :

Performance vs. Durability :
       (trade-off)    (latency)
   e.g. in Unix file system
        disk / OS cache :  write()
                            → mem

Dynamo:
   write buffer :   batch locally
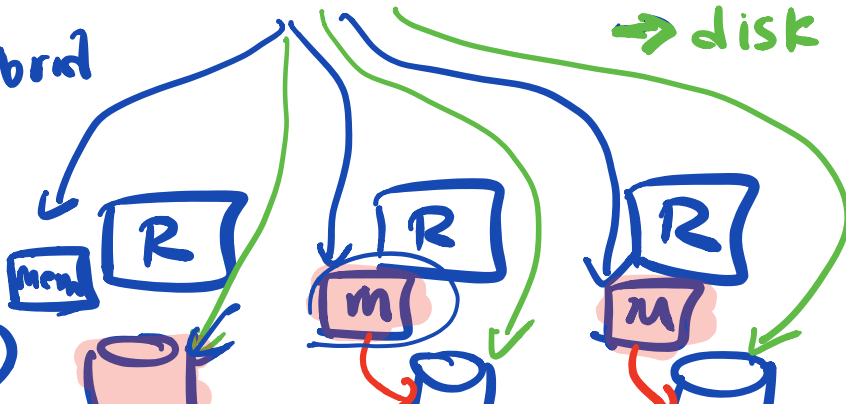                    ( perf =>
      => can lose      reduced tail
         data          latency by 5x)

   writes to N nodes :
                          → memory
                          ⇒ disk

Dynamo: hybrid
→ N node
   rep.
→ 1 durable [mem]  [R]    [R]    [R]
   rest (w)        [m]    [m]    [m]

memory:
(eventually persisted)
memory:
    perf but
    not as
      "reliable"

disk:
    other
    way

→ but: make more reliable by
    → placement ( across racks,
    (but bugs...)    datacenters,
                         etc.)

## Load Distribution

=> Partitioning:
  Consistent Hashing
  => really simplified
    (as compared to
        Chord)

why?
1) => [ Chord: variable size
                ranges ]
  → hard to sync.

## Background Tasks

=> scheduling
  of bg
  tasks

(details:
not shown
here)

instead: fixed-size partitions

2) knowledge of key-range
responsibility

$\rightarrow$ Chord: distributed

$\rightarrow$ Dynamo: all nodes have
this mapping info

Conclude:
$\rightarrow$ key/value  (NoSQL)