

# Today

- Dynamo
- End class: what have we learned?  
(wrap up)
- Admin
  - class presentations
  - Midterm #2
  - other

## Dynamo

Problem: DBs of time  $\Rightarrow$  didn't scale

Goals:

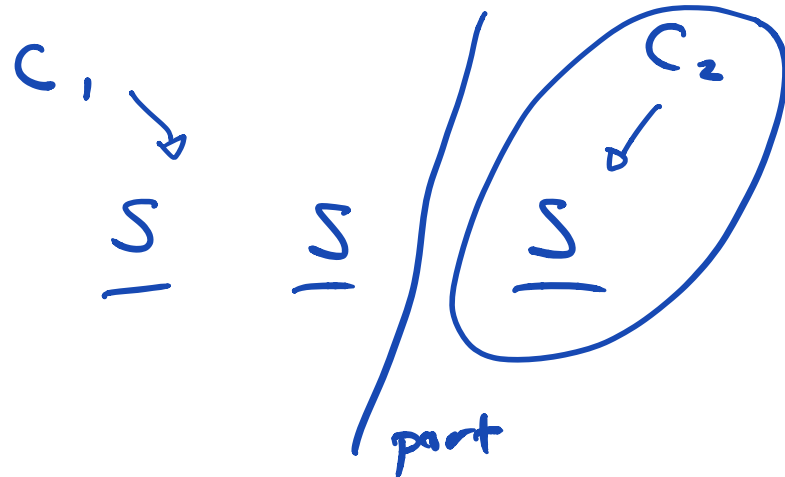
scale (incremental)

Availability: at all costs

$\Rightarrow$  Brewer: "CAP" theorem  $\Rightarrow$  truism

avail  
consistency — partition tolerance

$\Rightarrow$  pick any two



## Dynamo : availability

always writeable

→ conflict res on reads (later)

conflicts:

app-specific (Bayou)

↳ just accept writes

## Performance :

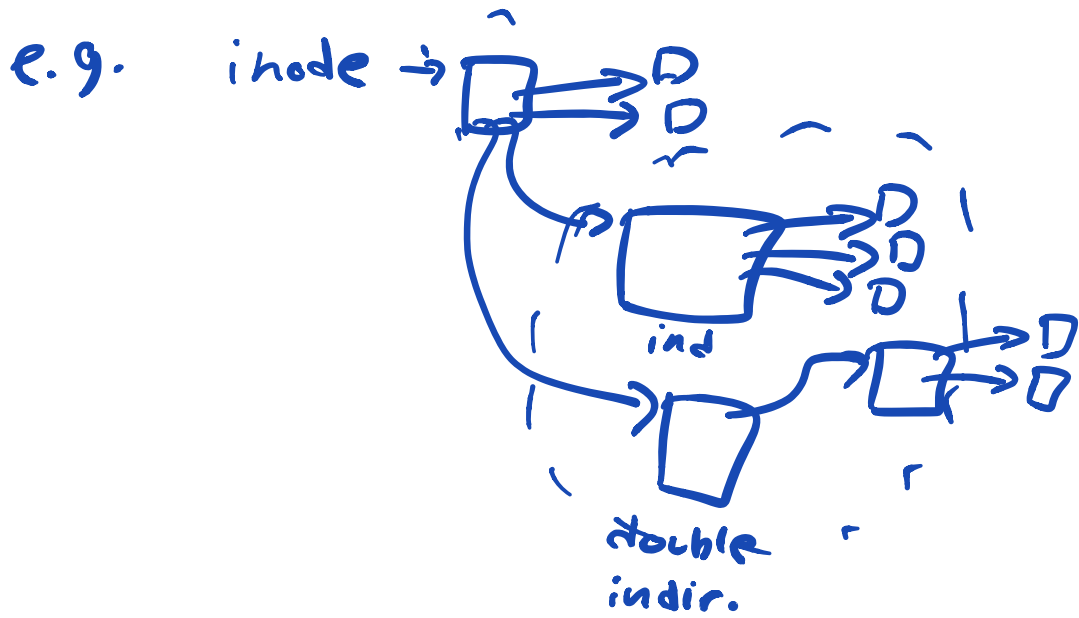
latency ⇒ tail latencies (99.9<sup>th</sup> percentile)

(predictable)

{ Service Level Agreements }

⇒ known h/w

⇒ admission control



## System Arch

API: get/put "key value"  
 hash(key) ⇒ lookup

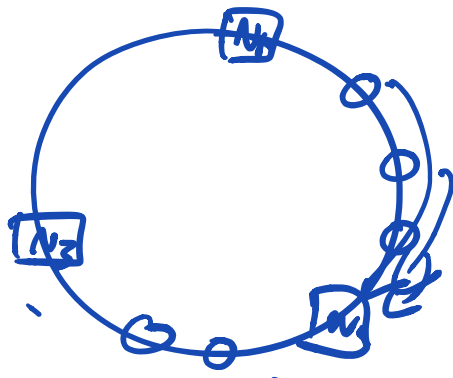
Partitioning:

"consistent hashing"

virtual nodes:

→ load balance  
 under join/leave

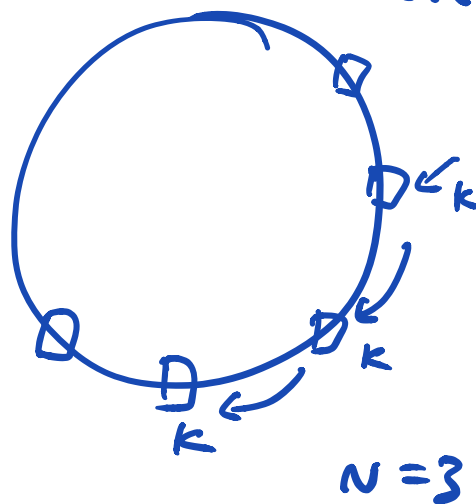
→ hetero.



Replication : N-way

writes: => coordinator

↓  
replicate key/value => successors  
careful: (N distinct phys nodes)



Version S

"eventual consistency"

W . . . . .

R  
R  
R

}  
}

may see new or old

---

e.g. shopping cart

how to map shop cart ops  
onto get/put?

=> op => put    add cart => put  
del cart => put

~ versions

=> vector clocks

list of (node, counter) ...  
allows detection of  
conflicts, non-conflicts

concern: size => lists can get  
big

limit: 10

Get/Put:

=> routing

2 options:

→ Load Balancing tier

→ client

-v<sub>1</sub>

-v<sub>2</sub>

-v<sub>3</sub>

-v<sub>4</sub>

-v<sub>5</sub>

writes: coordinator =>  
 replicate to first  
 N-1 "healthy" nodes

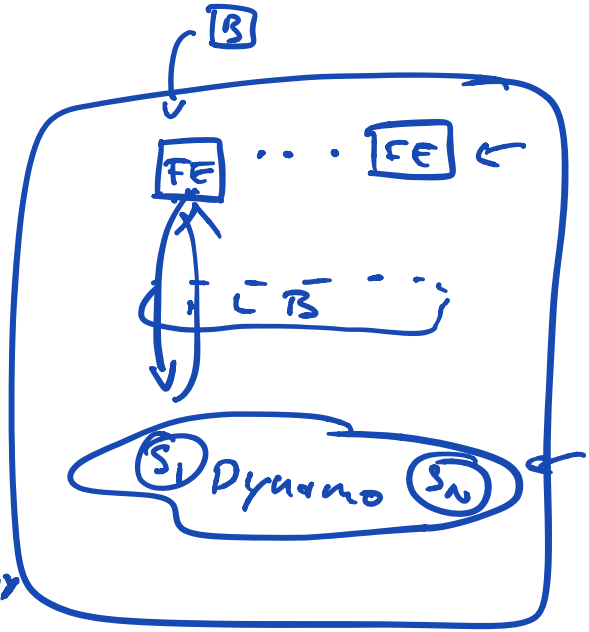
quorum: N, R, W  
 read write

"sloppy" quorums:  
 => avail.

replicate to N healthy  
 nodes

-> ~~N<sub>1</sub>~~ N<sub>2</sub> N<sub>3</sub> N<sub>4</sub>

recovers



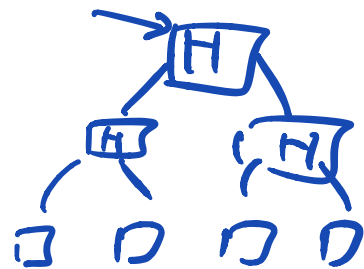
list: eventually  
 consistent  
 N<sub>1</sub> => replica  
 "hinted handoff"

## Replicas Sync.

periodic nodes sync w/ each other

merkle trees over key/values

=> efficient sync

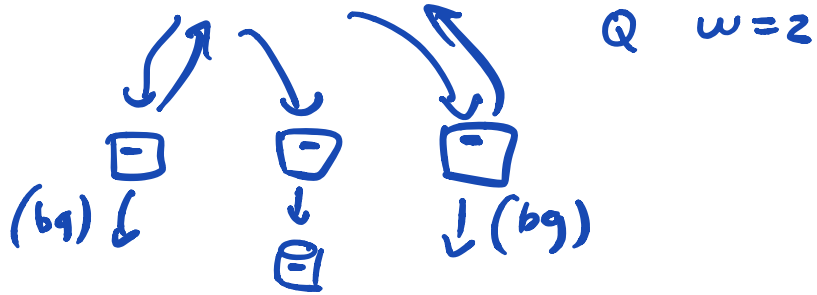


# Experience + Lessons

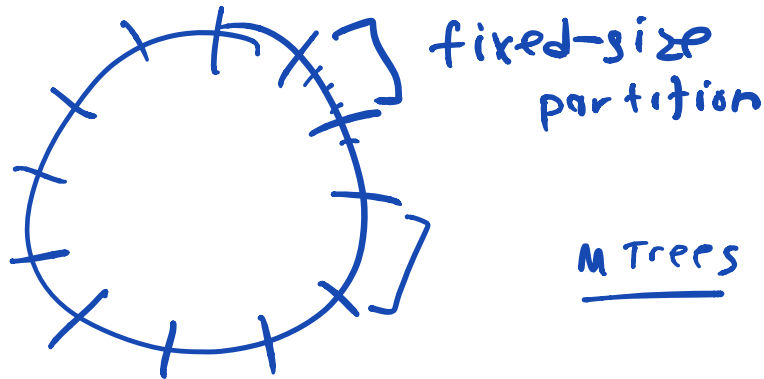
## Perf vs. Durability:

write buffer: 5x tail latency

write to N nodes



## Load Distribution



## Background Tasks:

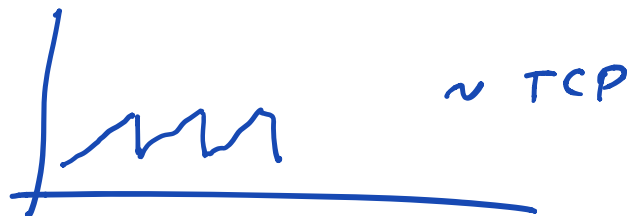
replica sync, etc.

⇒ bg tasks can slow fg tasks

Schedule bg work:

→ driven by measurement  
fg tail latency

→ throttle bg work accordingly



## Conclude:

Early NoSQL

→ Cassandra  
Risk  
⋮

Interesting:

wasn't super popular

inside amazon (@first)

→ system, not service  
(includes op.)

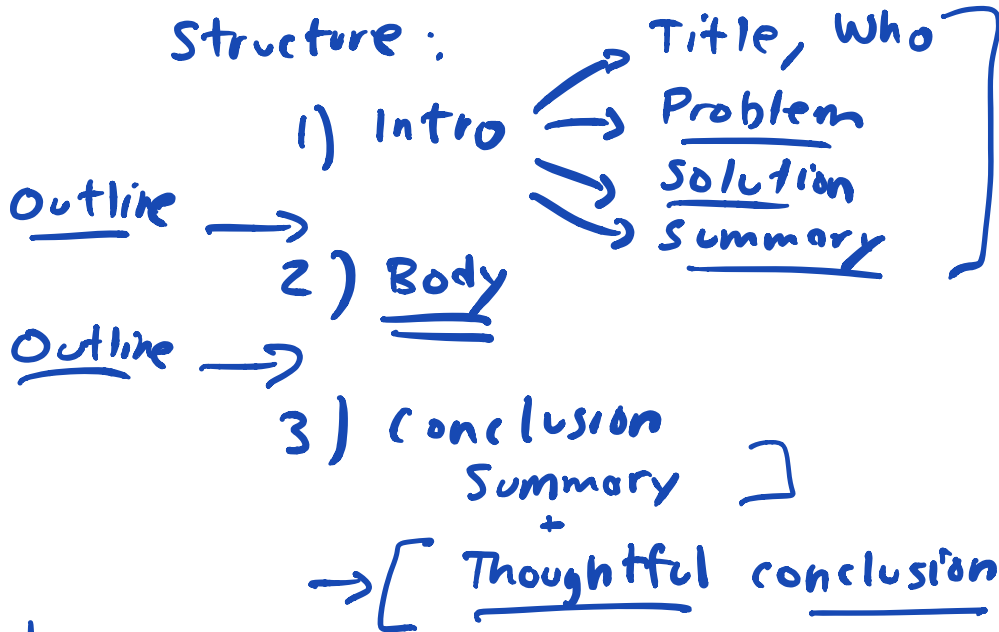


Admin :

⇒ Presentations :

20 minute talk  
slides, bring laptop

Structure :



2<sup>nd</sup> Midterm : no thanksgiving questions(!)

⇒ take home  
typing ok

Sun a.m. →  
Mon p.m (6)



Project write up:

---

{ short  
5-pages (+) }

)