

# Distributed Systems

CS425/ECE428

*Instructor: Radhika Mittal*

*Acknowledgements for the materials: Indy Gupta*

# Logistics

- HW5 and MP3 due Friday.
- Final exam starts next week!
  - See Campuswire post #420 for cheatsheet and syllabus.
- Further delay in releasing midterm 2 grades. Thank you for your patience!

# Logistics

- Participation score: directly taken from Campuswire
  - If reported score  $> 100$ , you get full 1%
  - Else you get  $(\text{reported score} / 100)\%$
- Upto 100% bonus for active participation in class.
  - Watch out for an email from me in the next 2-3 weeks.
  - I'll send a Campuswire post after sending all emails.
    - If you do not receive such email, but believe you have actively participated in class, send me an email with your latest headshot!

# DRF Algorithm

## Algorithm 1 DRF pseudo-code

---

$R = \langle r_1, \dots, r_m \rangle$       ▷ total resource capacities  
 $C = \langle c_1, \dots, c_m \rangle$     ▷ consumed resources, initially 0  
 $s_i$  ( $i = 1..n$ )    ▷ user  $i$ 's dominant shares, initially 0  
 $U_i = \langle u_{i,1}, \dots, u_{i,m} \rangle$  ( $i = 1..n$ ) ▷ resources given to user  $i$ , initially 0

**pick** user  $i$  with lowest dominant share  $s_i$

$D_i \leftarrow$  demand of user  $i$ 's next task

**if**  $C + D_i \leq R$  **then**

$C = C + D_i$       ▷ update consumed vector

$U_i = U_i + D_i$     ▷ update  $i$ 's allocation vector

$s_i = \max_{j=1}^m \{u_{i,j}/r_j\}$

**else**

    remove user  $i$  from consideration

**end if**

(terminates when no more users can be accommodated)

Our example

Job 1's tasks: 2 CPUs, 8 GB

    => Job 1's resource vector = <2 CPUs, 8 GB>

Job 2's tasks: 6 CPUs, 2 GB

    => Job 2's resource vector = <6 CPUs, 2 GB>

Consider a cloud with <18 CPUs, 36 GB RAM>

# DRF Algorithm

## Algorithm 1 DRF pseudo-code

---

$R = \langle r_1, \dots, r_m \rangle$       ▷ total resource capacities  
 $C = \langle c_1, \dots, c_m \rangle$     ▷ consumed resources, initially 0  
 $s_i$  ( $i = 1..n$ )    ▷ user  $i$ 's dominant shares, initially 0  
 $U_i = \langle u_{i,1}, \dots, u_{i,m} \rangle$  ( $i = 1..n$ ) ▷ resources given to user  $i$ , initially 0

**pick** user  $i$  with lowest dominant share  $s_i$

$D_i \leftarrow$  demand of user  $i$ 's next task

**if**  $C + D_i \leq R$  **then**

$C = C + D_i$       ▷ update consumed vector

$U_i = U_i + D_i$     ▷ update  $i$ 's allocation vector

$s_i = \max_{j=1}^m \{u_{i,j}/r_j\}$

**else**

    remove user  $i$  from consideration

**end if**

(terminates when no more users can be accommodated)

Our example

Job 1's tasks: 2 CPUs, 8 GB

    => Job 1's resource vector = <2 CPUs, 8 GB>

Job 2's tasks: 6 CPUs, 2 GB

    => Job 2's resource vector = <6 CPUs, 2 GB>

Consider a cloud with <18 CPUs, 36 GB RAM>

# Our agenda

- Brief overview of key-value stores
- Distributed Hash Tables
  - Peer-to-peer protocol for efficient insertion and retrieval of key-value pairs.
- Key-value stores in the cloud
  - How to run large-scale distributed computations over key-value stores?
    - Map-Reduce Programming Abstraction
    - Cloud scheduling
  - How to design a large-scale distributed key-value store?
    - Case-study: Facebook's Cassandra

# Distributed datastores

- Distributed datastores
  - Service for managing distributed storage.
- Distributed NoSQL key-value stores
  - BigTable by Google
  - HBase open-sourced by Yahoo and used by Hadoop.
  - DynamoDB by Amazon
  - Cassandra by Facebook
  - Voldemort by LinkedIn
  - MongoDB,
  - ...
- *Spanner is not a NoSQL datastore. It's more like a distributed relational database.*

How to design a distributed  
key-value datastore?

# Design Requirements

- High performance, low cost, and scalability.
  - Speed (high throughput and low latency for read/write)
  - Low TCO (total cost of operation)
  - Fewer system administrators
  - Incremental scalability
    - Scale out: add more machines.
    - Scale up: upgrade to powerful machines.
    - *Cheaper to scale out than to scale up.*

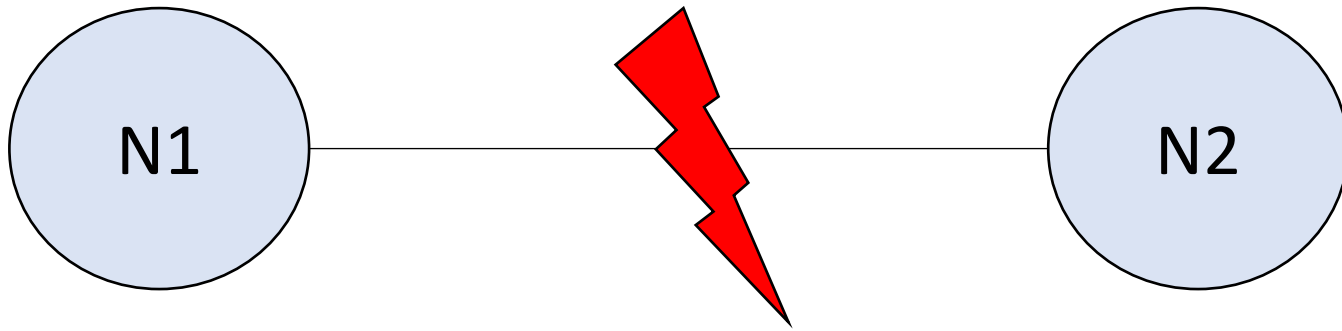
# Design Requirements

- High performance, low cost, and scalability.
- Avoid single-point of failure
  - Replication across multiple nodes.
- Consistency: reads return latest written value by any client (all nodes see same data at any time).
  - *Different from the C of ACID properties for transaction semantics!*
- Availability: every request received by a non-failing node in the system must result in a response (quickly).
  - Follows from requirement for high performance.
- Partition-tolerance: the system continues to work in spite of network partitions.

# CAP Theorem

- **C**onsistency: reads return latest written value by any client (all nodes see same data at any time).
- **A**vailability: every request received by a non-failing node in the system must result in a response (quickly).
- **P**artition-tolerance: the system continues to work in spite of network partitions.
- **In a distributed system you can only guarantee at most 2 out of the above 3 properties.**
  - Proposed by Eric Brewer (UC Berkeley)
  - Subsequently proved by Gilbert and Lynch (NUS and MIT)

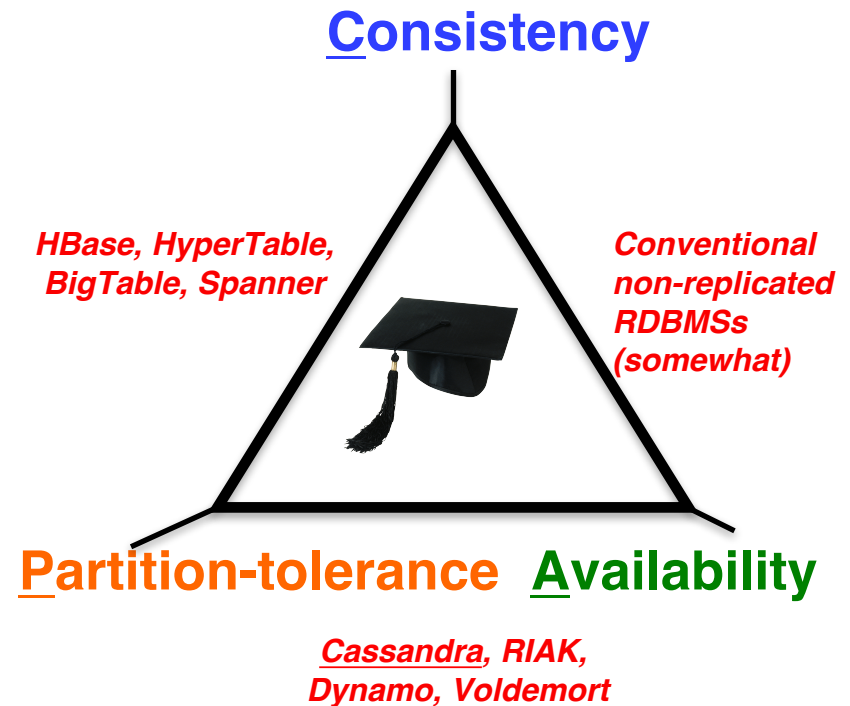
# CAP Theorem



- Data replicated across both N1 and N2.
- If network is partitioned, N1 can no longer talk to N2.
- Consistency + availability
  - N1 and N2 must talk (no partition-tolerance).
- Partition-tolerance + consistency:
  - only respond to requests received at N1 (no availability).
- Partition-tolerance + availability:
  - write at N1 will not be captured by a read at N2 (no consistency).

# CAP Tradeoff

- Starting point for NoSQL Revolution
- A distributed storage system can achieve **at most two of C, A, and P.**
- When partition-tolerance is important, you have to choose between consistency and availability



# Case Study: Cassandra

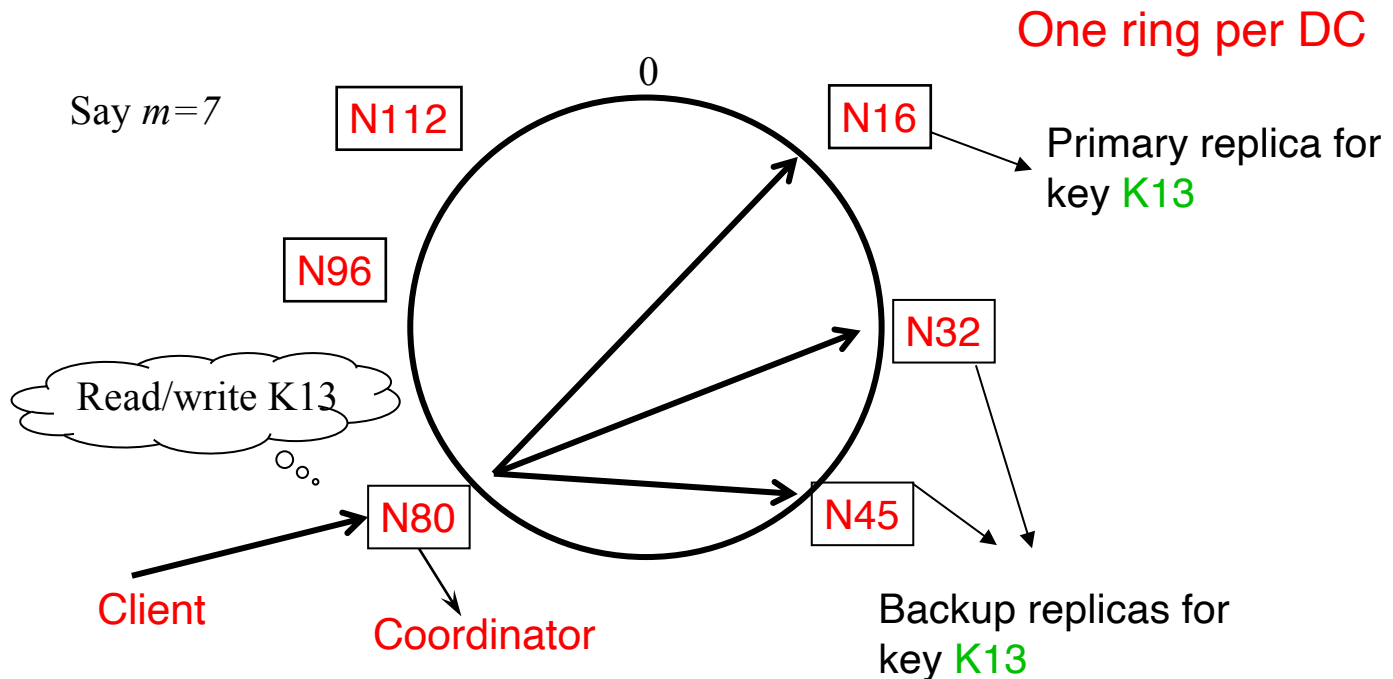
# Cassandra

- A distributed key-value store.
- Intended to run in a datacenter (and also across DCs).
- Originally designed at Facebook.
- Open-sourced later, today an Apache project.
- Some of the companies that use Cassandra in their production clusters.
  - IBM, Adobe, HP, eBay, Ericsson, Symantec
  - Twitter, Spotify
  - PBS Kids
  - Netflix

# Data Partitioning: Key to Server Mapping

- How do you decide which server(s) a key-value resides on?

Cassandra uses a ring-based DHT but without finger or routing tables.



# Partitioner

- Component responsible for key to server mapping (hash function).
- Determines the primary replica for a key.
- Two types:
  - *Chord-like hash partitioning*
    - *Murmur3Partitioner* (default): uses *murmur3* hash function.
    - *RandomPartitioner*: uses MD5 hash function.
  - *ByteOrderedPartitioner*: Assigns ranges of keys to servers.
    - Easier for range queries (e.g., get me all twitter users starting with [a-b])

# Replication Policies

Two options for replication strategy:

## 1. SimpleStrategy:

- First replica placed based on the partitioner.
- Remaining replicas clockwise in relation to the primary replica.

## 2. NetworkTopologyStrategy: for multi-DC deployments

- Two or three replicas per DC.
- Per DC
  - First replica placed according to Partitioner.
  - Then go clockwise around ring until you hit a different rack.

# Writes

- Need to be lock-free and fast (no reads or disk seeks).
- Client sends write to one coordinator node in Cassandra cluster.
  - Coordinator may be per-key, or per-client, or per-query.
- Coordinator uses Partitioner to send query to all replica nodes responsible for key.
- When  $X$  replicas respond, coordinator returns an acknowledgement to the client
  - $X =$  any one, majority, all....(consistency spectrum)
  - More details later!

# Writes: Hinted Handoff

- Always writable: Hinted Handoff mechanism
  - If any replica is down, the coordinator writes to all other replicas, and keeps the write locally until down replica comes back up.
  - When all replicas are down, the Coordinator (front end) buffers writes (for up to a few hours).

# Writes at a replica node

On receiving a write

1. Log it in disk commit log (for failure recovery)
2. Make changes to appropriate memtables
  - **Memtable** = In-memory representation of multiple key-value pairs
  - Cache that can be searched by key
  - Write-back cache as opposed to write-through
3. Later, when memtable is full or old, flush to disk
  - Data File: An **SSTable** (Sorted String Table) – list of key-value pairs, sorted by key.
  - Each SSTable accompanied by some other data structures (index tables, Bloom filters) for efficient look-ups.

# Compaction

- Data updates accumulate over time and over multiple SSTables.
- Need to be compacted.
- The process of compaction merges SSTables, i.e., by merging updates for a key.
- Run periodically and locally at each server.

# Deletes

Delete: don't delete item right away

- Write a **tombstone** for the key.
- Eventually, when compaction encounters tombstone it will delete item

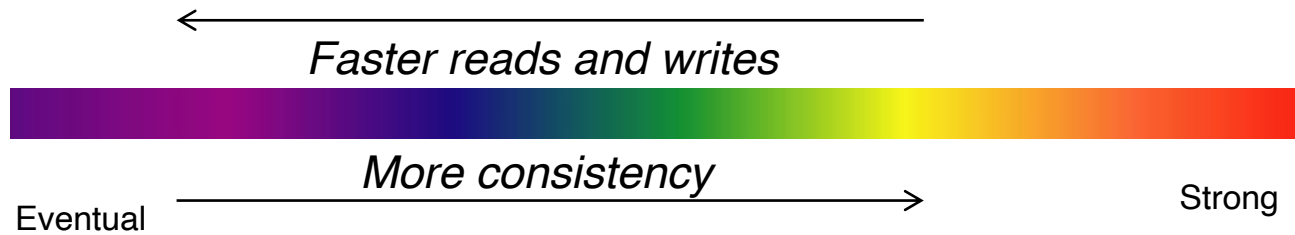
# Reads

- Coordinator contacts  $X$  replicas (e.g., in same rack)
  - Coordinator sends read to replicas that have responded quickest in past.
  - When  $X$  replicas respond, coordinator returns the latest-timestamped value from among those  $X$ .
  - $X$  = based on consistency spectrum (more later).
- Coordinator also fetches value from other replicas
  - Checks consistency in the background, initiating a **read repair** if any two values are different.
  - This mechanism seeks to eventually bring all replicas up to date.
- At a replica
  - Read looks at Memtables first, and then SSTables.
  - A row may be split across multiple SSTables => reads need to touch multiple SSTables => reads slower than writes (but still fast).

# Cross-DC coordination

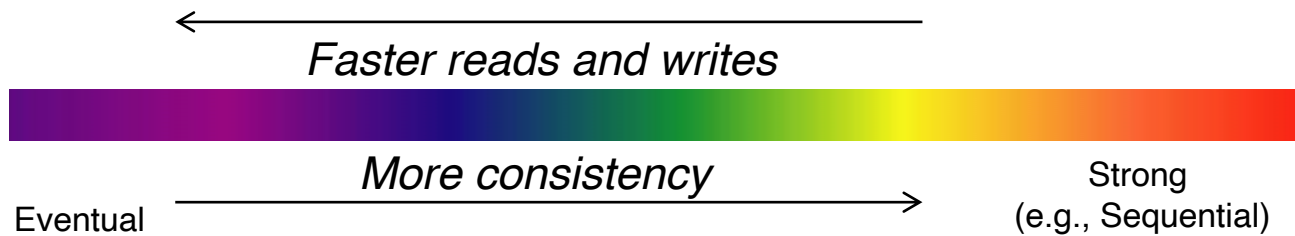
- Replicas may span multiple datacenters.
- Per-DC coordinator elected to coordinate with other DCs.
- Election done via Zookeeper which runs a Bully algorithm variant.

# Consistency Spectrum



# Eventual Consistency

- Cassandra offers **Eventual Consistency**
  - If writes to a key stop, all replicas of key will converge.
  - Originally from Amazon's Dynamo and LinkedIn's Voldemort systems



# Cassandra write and read recap

- Writes
  - Client sends write request to a *coordinator*.
  - Coordinator writes to all replicas.
  - Waits for **X** replicas to respond before returning acknowledgement to the client.
  - Hinted handoff: if a replica is down, it receives the write request once it comes back up.
- Reads
  - Client sends read request to a *coordinator*.
  - Coordinator contacts **X** replicas, and returns the latest returned value.
  - Read repair: After returning a response, coordinator continues with fetching values from other replicas, and initiates repairs to outdated values.

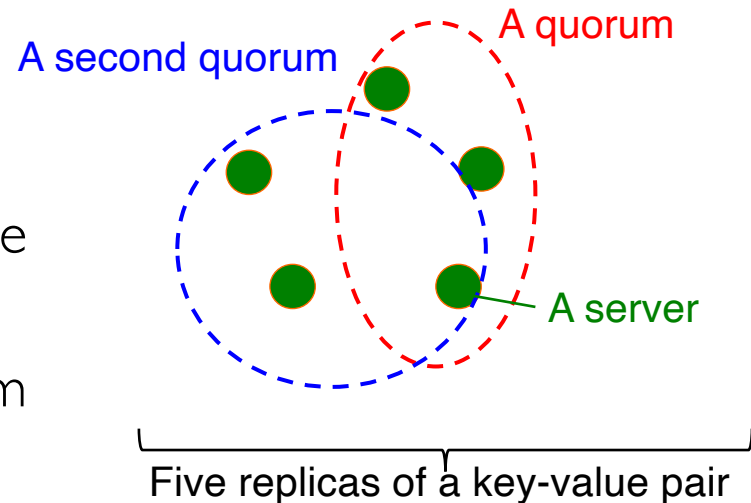
# Consistency levels: value of X

- Cassandra has [consistency levels](#).
- Client is allowed to choose a consistency level for each operation (read/write)
  - ANY: any server (may not be replica)
    - Fastest: coordinator caches write and replies quickly to client
  - ALL: all replicas
    - Ensures strong consistency, but slowest
  - ONE: at least one replica
    - Faster than ALL, but cannot tolerate a failure
  - QUORUM: quorum across all replicas in all datacenters (DCs)

# Quorums?

In a nutshell:

- Quorum = (typically) majority
- Any two quorums intersect
  - Client 1 does a write in red quorum
  - Then client 2 does read in blue quorum
- At least one server in blue quorum returns latest write
- Quorums faster than ALL, but still ensure strong consistency
- Several key-value/NoSQL stores (e.g., Riak and Cassandra) use quorums.



# Read Quorums

- Reads
  - Client specifies value of  $R$  ( $\leq N$  = total number of replicas of that key).
  - $R$  = read consistency level.
  - Coordinator waits for  $R$  replicas to respond before sending result to client.
  - In background, coordinator checks for consistency of remaining  $(N-R)$  replicas, and initiates read repair if needed.

# Write Quorums

- Client specifies  $W$  ( $\leq N$ )
- $W$  = write consistency level.
- Client writes new value to  $W$  replicas and returns when it hears back from all.
  - Default strategy.

# Quorums in Detail (Contd.)

- $R$  = read replica count,  $W$  = write replica count
- Necessary conditions for consistency:
  1.  $W+R > N$ 
    - Write and read intersect at a replica. Read returns latest write.
  2.  $W > N/2$ 
    - Two conflicting writes on a data item don't occur at the same time.
- Select values based on application:
  - To be continued in next class.....