CS 61: Database Systems

Distributed systems

Adapted mongodb.com unless otherwise noted

Agenda

1. Centralized systems

- 2. Distributed systems
 - High availability
 - Scalability
- 3. MongoDB

A single database can handle many thousands of transactions per second



Your start up that you're certain will be a smashing success in the market is unlikely to overwhelm a database running on a single reasonable server for quite some time

-- Pierson

Scale vertically – get a bigger box Scale horizontally – get more boxes

Source: https://www.mysql.com/why-mysql/benchmarks/

Let's estimate performance

Assume:

Source: percona.com

Each user interaction takes 10 queries on average (normalization) Average of 30 user interactions/visitor Database can handle 100,000 queries per second If you exceed these numbers, you'll need some help from someone who took more than Max user interactions/second: an introductory database class! 100,000 queries x 1 interaction = 10,000 interactions 10 queries second second Max user interactions/day: 10,000 interactions x 60*60*24 seconds = 864M interactions day day second Max user visits/day: 864M interactions x 1 visitor 28.8M visitors day 30 interactions day

4

Agenda

1. Centralized systems

- 2. Distributed systems
 - High availability
 - Scalability
- 3. MongoDB

With one database, you've put all you eggs in one basket!



With SANs you can have real-time, blocklevel replication to another database



7

Log shipping is another, often more costeffective high availability solution



Network speed/reliability important

Partitioning can help with scalability when data become large

Partitioning

Horizontal partitioning

(aka sharding) slices data by rows and spreads across multiple nodes

Vertical partitioning slices

data by columns



ID	Name	Salary
100	Alice	100,000
200	Bob	90,000
300	Charlie	85,000

Partitioning increases capacity

- Each machine may be small, but only handles a subset of overall data
- Tradeoff: increased complexity



Global distributed schema keeps track of data locations





Sharding horizontally partitions data based on an attribute chosen as the shard key

Partitioning

Choose an attribute to serve as the shard key

- Shard key difficult change once database is partitioned
- Want cardinality > number of shards

	Logica	l data	view
--	--------	--------	------

ID	Name	Salary
100	Alice	100,000
200	Bob	90,000
300	Charlie	85,000

Hashing:

 Hash shard key to get replica number like CS10 hash table

Range partitioning:

Distribute based on a partition key (Names A-E go to database 1,

F-J go to database 2, ...







Two-phase commit (2PC) protocol

UPDATE Employee

SET Salary = Salary * 1.05

One node chosen

Two-phase commit: DO-UNDO-REDO protocol

- **DO:** Record before and after values in write ahead transaction log
- UNDO: Reverses operation using transaction log
- **REDO:** redoes an operation written by DO

11



Two-phase commit (2PC) protocol

UPDATE Employee

SET Salary = Salary * 1.05

Phase 1: Preparation

- Coordinator send Prepare to Commit message to all subordinate nodes
- Subordinates write transaction log and send acknowledgement to coordinator
- Coordinator ensures all nodes ready to commit or aborts



One node chosen

Two-phase commit (2PC) protocol

UPDATE Employee

SET Salary = Salary * 1.05

Phase 1: Preparation

- Coordinator send Prepare to Commit message to all subordinate nodes
- Subordinates write transaction log and send acknowledgement to coordinator
- Coordinator ensures all nodes ready to commit or aborts



One node chosen

Two-phase commit (2PC) protocol

UPDATE Employee

SET Salary = Salary * 1.05

One node chosen

Phase 2: Commit

- Coordinator broadcasts commit message
- Each subordinate updates with DO
- Subordinates reply with COMMITTED or NOT COMMITTED
- If any nodes reply NOT COMMITTED, then UNDO followed by REDO



Two-phase commit (2PC) protocol

UPDATE Employee

SET Salary = Salary * 1.05

One node chosen

Phase 2: Commit

- Coordinator broadcasts commit message
- Each subordinate updates with DO
- Subordinates reply with COMMITTED or NOT COMMITTED
- If any nodes reply NOT COMMITTED, then UNDO followed by REDO



Two-phase commit (2PC) protocol

UPDATE Employee

SET Salary = Salary * 1.05

One node chosen

Phase 2: Commit

- Coordinator broadcasts commit message
- Each subordinate updates with DO
- Subordinates reply with COMMITTED or NOT COMMITTED
- If any nodes reply NOT COMMITTED, then UNDO followed by REDO



Data might be replicated to several nodes located across the globe



Data replication scenarios

Fully replicated: multiple copies of each database partition at multiple sites **Partially replicated:** multiple copies of some database partitions at multiple sites **Unreplicated:** stores each database partition at a single site

All replicated nodes should have same data, but network latency raises issues



A potential problem: consistency rule says all copies must be identical when data changes are made

- **Push replication** (focus on consistency)
 - After data update, send changes to all replicas
 - Data unavailable until changes to propagate across all copies, but data is always consistent across copies
- **Pull replication** (focus on availability)
 - Send message to all replicas, they decide when to apply change
 - Data is available, but not consistent until changes propagate

Read operations just query the nearest replica

The network may be partitioned by communication breaks



The network may have multiple communication links between each node If one link fails, other nodes will still be reachable Multiple link failure, however, may separate some nodes – called a network partition

The network may be partitioned by communication breaks



The network may have multiple communication links between each node If one link fails, other nodes will still be reachable Multiple link failure, however, may separate some nodes – called a network partition

The network may be partitioned by communication breaks



The network may have multiple communication links between each node If one link fails, other nodes will still be reachable Multiple link failure, however, may separate some nodes – called a network partition

It impossible to be consistent, available, and partition tolerant simultaneously

CAP theorem

The CAP theorem showed that it is impossible to have three desirable properties at the same time in distributed systems

Consistency

- All nodes should return the same data at the same time
- Replicas should be immediately updated
- Network latency means this cannot happen

Availability

- A request is always fulfilled by the system
- No request is ever lost

Partition tolerant

- The system will operate even if nodes fail
- Operations that are lost due to node failure are picked up by other nodes
- The system will only fail if all nodes fail

Trade-off between consistency and availability

BASE rather than ACID: Basically Available, Soft state, Eventually consistent (BASE)

Data changes are not immediate but propagate slowly through the system until all replicas are eventually consistent

Agenda

- 1. Centralized systems
- 2. Distributed systems
 - High availability
 - Scalability



MongoDB is a NoSQL database designed for high availability and scalability

mongoDB

MongoDB is a document database designed for scalability and flexibility

- MongoDB is "schemaless"
 - Stores data in JSON-like documents
 - Fields can vary from document to document
 - Data structure can change over time
- MongoDB is a distributed database at its core
 - High availability (replication)
 - Horizontal scaling (sharding)
 - Geographic distribution built in
- MongoDB is free to use
 - Cloud-based version at MongoDB Atlas (<u>https://www.mongodb.com/cloud/atlas</u>)



Mongo solves the high availability problem using two different types of nodes

Mongo replication

Replica set: Group of database servers that provide high availability and redundancy using two different types of nodes

- **Primary:** Receives all write operations
- Secondary: replicate operations from primary to maintain an identical data set

Mongo recommends at least a threemember replica set (more even better)

All members of replica set can accept read operations

Replica sets store the same data in all nodes

Purpose: redundancy for high availability



If the primary nodes fails, other nodes hold an election to pick a new primary

Mongo replication

Replica sets use elections to determine which set member will be primary

Elections triggered:

- Start up
- New node added to set
- Secondary member looses connectivity to primary for 10 seconds (default)
- Called automatic failover when new primary takes over



Mongo shards data across multiple nodes, each shard can be replicated

Replication with sharding



Source: https://severalnines.com/blog/turning-mongodb-replica-set-sharded-cluster

If network is partitioned, behavior depends on primary's location

Replication with sharding



If network is partitioned, behavior depends on primary's location

Replication with sharding



Homework: get access to a MongoDB database

Create a cloud-based MongoDB database

- 1. Create a free cloud-based account at Mongo Atlas: <u>https://www.mongodb.com/cloud/atlas</u>
- Install mongo shell to interact with your database <u>Mac</u> (assumes you have brew installed): brew tap mongodb/brew brew install mongodb-community-shell

<u>Windows</u>: <u>https://www.mongodb.com/download-center/community</u> (install only shell)

3. Optional: install Compass (like MySQL Workbench): <u>https://www.mongodb.com/products/compass</u>

OR

Install MongoDB locally on your machine

https://docs.mongodb.com/manual/installation/ (includes shell)