### Scaling Data Lucas Dohmen Senior Consultant @ INNOQ





### Scaling Reads

### Scaling Writes

**INOQ** / Motivation

### Failure Resistance

### Why?

### Geographical Distribution

# Big Data Sets





- Senior Consultant at INNOQ
- Everything Web & Databases
- Previously worked at ArangoDB
- <u>http://faucet-pipeline.org</u>









- 1. Consistency
- 2. Scaling
- 3. Trouble











## Part 1: Consistency











W



INOQ / Consistency





W



**INOQ** / Consistency





W



**INOQ** / Consistency



W



**INOQ** / Consistency



W



**INOQ** / Consistency



### Consistency Models: Which histories are valid?

Read() => a Write(b) Read() => b Read() => b Read() => a Write(b) Read() => a Read() => a

**INOQ** / Consistency

Read() => b Write(b) Read() => b Read() => b





## https://aphyr.com/posts/313strong-consistency-models

**INOQ** / Consistency







**INOQ** / Consistency

<u>Highly Available Transactions: Virtues and Limitations – Bailis et al.</u>



## Part 2: Scaling





- Share nothing between application servers
- Put behind a load balancer
- Add servers

### **INOQ** / Scaling Applications



- Share nothing between application servers
- Put behind a load balancer
- Add servers

### **INOQ** / Scaling Applications



## Share Nothing for Databases?

- Possible & Underused
- Separate databases for separate data
- If we need to join data, we need to join in the application



**MySQL** 





Replication





## Replication

# Same data on multiple nodes

**INOQ** / Scaling / Replication

- Failover
- Read scaling
- No write scaling

**INOQ** / Scaling / Replication





## Sync or Async Replication?

- Trade-off between consistency & speed
- Sync: Every follower we add decreases performance
- Async: If our leader dies and the replication is not done, we have lost acknowledged data



- Redis
- MariaDB
- PostgeSQL
- MongoDB



## Examples

- Failover
- Read & write scaling

**INOQ** / Scaling / Replication

## Multi Leader





- Two leaders can accept a conflicting write
- We usually resolve them when reading
- Do we have all information we need to resolve a conflict at read time?

**INOQ** / Scaling / Replication









- CouchDB
- Percona Server for MySQL
- ArangoDB



## Examples





- Failover
- Read & write scaling

**INOQ** / Scaling / Replication







- Clients write to multiple nodes at once
- When more than n nodes acknowledged the write, the write is successful (n is the write quorum)
- When we read, we read from m nodes (m is the read quorum)

### **INOQ** / Scaling / Replication

### Guorum





- riak
- Cassandra
- aerospike



## Examples





**INOQ** / Scaling / Sharding

Sharding



## Sharding

### Each node only has part of the data

**INOQ** / Scaling / Sharding



## Sharding by Primary Key

H-L



**INOQ** / Scaling / Sharding





### Sharding by Hashed Primary Key

### • Equal distribution to all shards

**INOQ** / Scaling / Sharding



## **Combining Replication & Sharding**





Shards







## Part 3: Trouble









## Clocks are monotonic & synchronized

**INOQ** / Trouble



## Clocks are monotonic & synchronized

**INOQ** / Trouble



## Clocks are monotonic & synchronized

**INOQ** / Trouble





## Clocks dre monotonic & synchronized

**NTP Sync**  $\Rightarrow$  **Going back in time** 

**INOQ** / Trouble





## Clocks dre monotonic & synchronized

**NTP Sync**  $\Rightarrow$  **Going back in time** 

**INOQ** / Trouble

### **NTP fails**

### NTP is an estimation



## Clocks dre monetoric & synchronized

### **NTP Sync** $\Rightarrow$ **Going back in time**

**INOQ** / Trouble

### **NTP fails**

### NTP is an estimation



# DO NOT USE WALL CLOCKS FOR ORDERING

**INOQ** / Trouble



## Solution: Vector Clocks

**INOQ** / Trouble



The network is reliable



### Problem

**Reordered Messages** 

Lost Messages

Duplicated Messages

**Delayed Messages** 









The network is reliable



## The network is reliable

**INOQ** / Trouble



### packages can take a loooong time

## The network is reliable

**INOQ** / Trouble



### packages can take a loooong time

## The network is reliable

**INOQ** / Trouble

### the network can fail partially/entirely







Node Failure

### Node Recovery

### Crash Amnesia









Availability vs. Consistency

# YOUR NODES

**INOQ** / Trouble

# YOUR NETWORK























- Stop taking requests
- Not available, but consistent



- Continue taking requests
- Available, but not consistent







<u>Highly Available Transactions: Virtues and Limitations – Bailis et al.</u>

**INOQ** / Trouble











Wrap-Up





- Nodes will fail
- The network will fail
- Clocks aren't reliable

**INOQ** / Wrap Up

## Remember!





## What are your requirements?



**INOQ** / Wrap Up







- @moonbeamlabs on Twitter
- Photo Credit
  - Slide 5: <u>Shoot N' Design</u> on <u>Unsplash</u>
  - Slide 11: <u>Andy Hall</u> on <u>Unsplash</u>
  - Slide 30: <u>Hermes Rivera</u> on <u>Unsplash</u>





