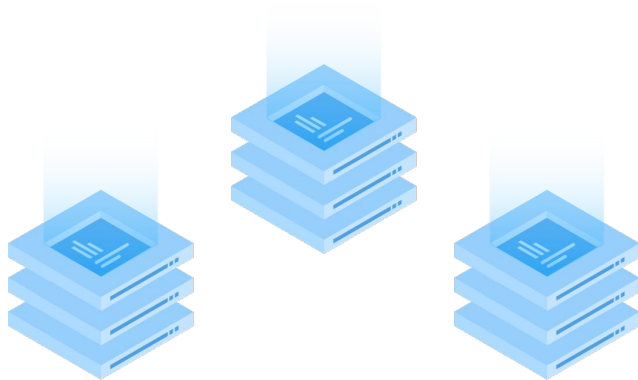


Using Rust to Build a Distributed Transactional Key-Value Database

LiuTang



About Me

- Chief Engineer of PingCAP
- Leader of TiKV
- Open source lovers: go-mysql, raft-rs, rust-prometheus, etc.
- tl@pingcap.com
- siddontang (Github, twitter)

Agenda

- Problems
- How do we build TiKV?
- Beyond TiKV

When we want to build a distributed transactional key-value database...



Problems

1. How to save the data in the machine?
2. How to support fault-tolerance?
3. How to provide ACID features?
4. How to communicate with the servers?
5. How to test the database?
6. ...



A High Building, A Low Foundation



Language





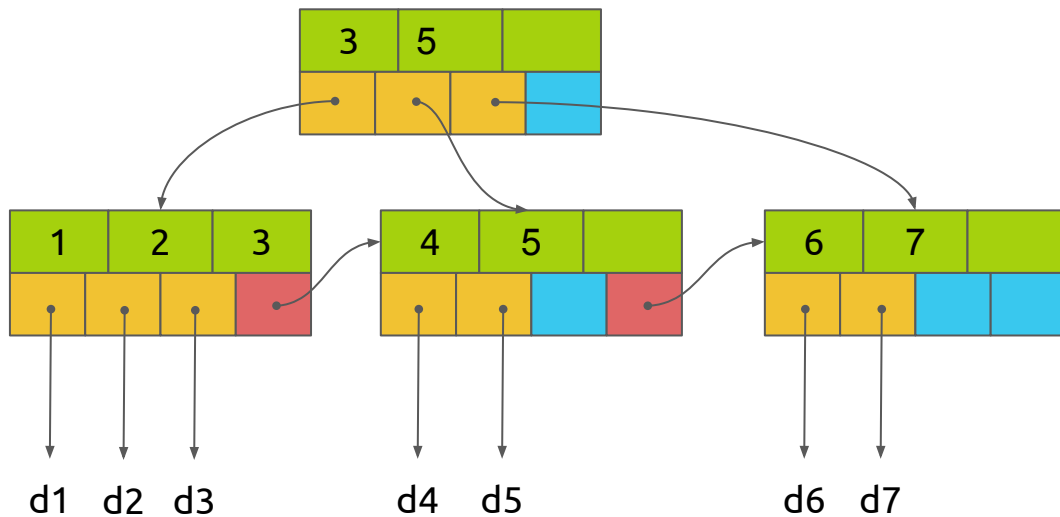
Let's start from scratch!!!



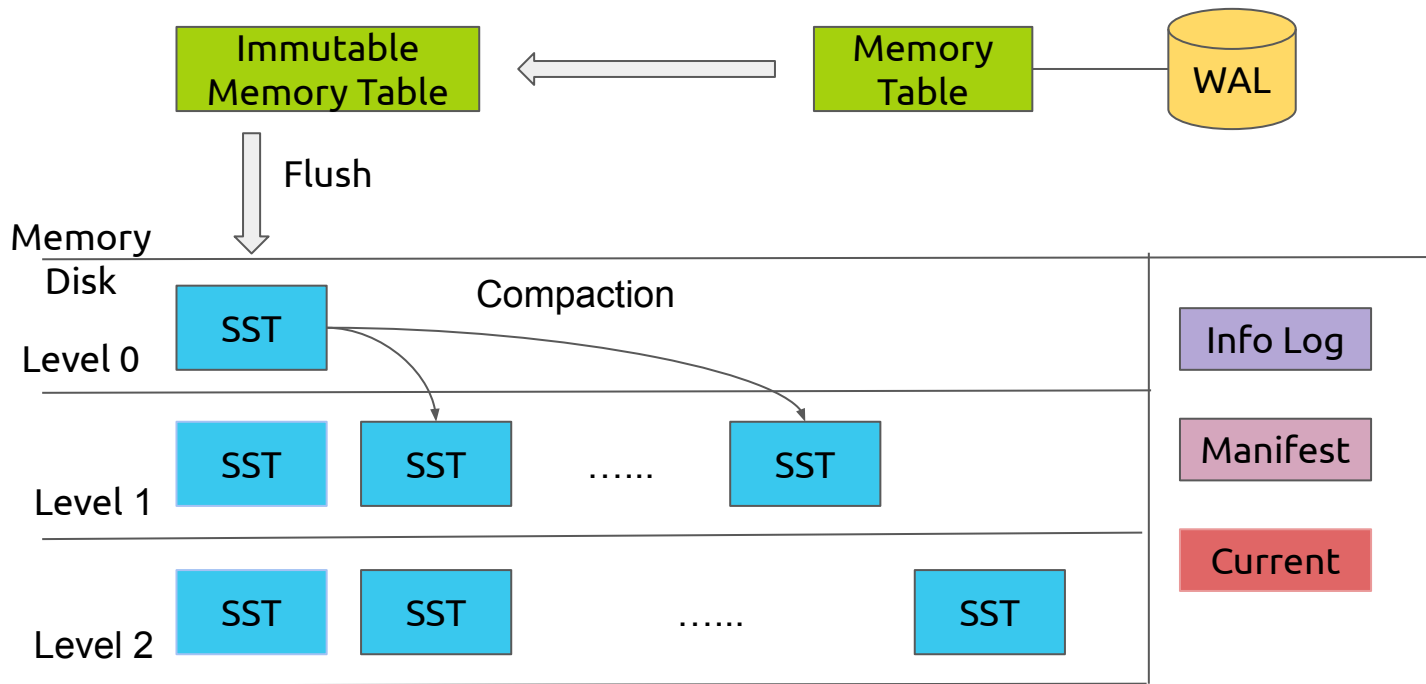
Storage



B+ Tree



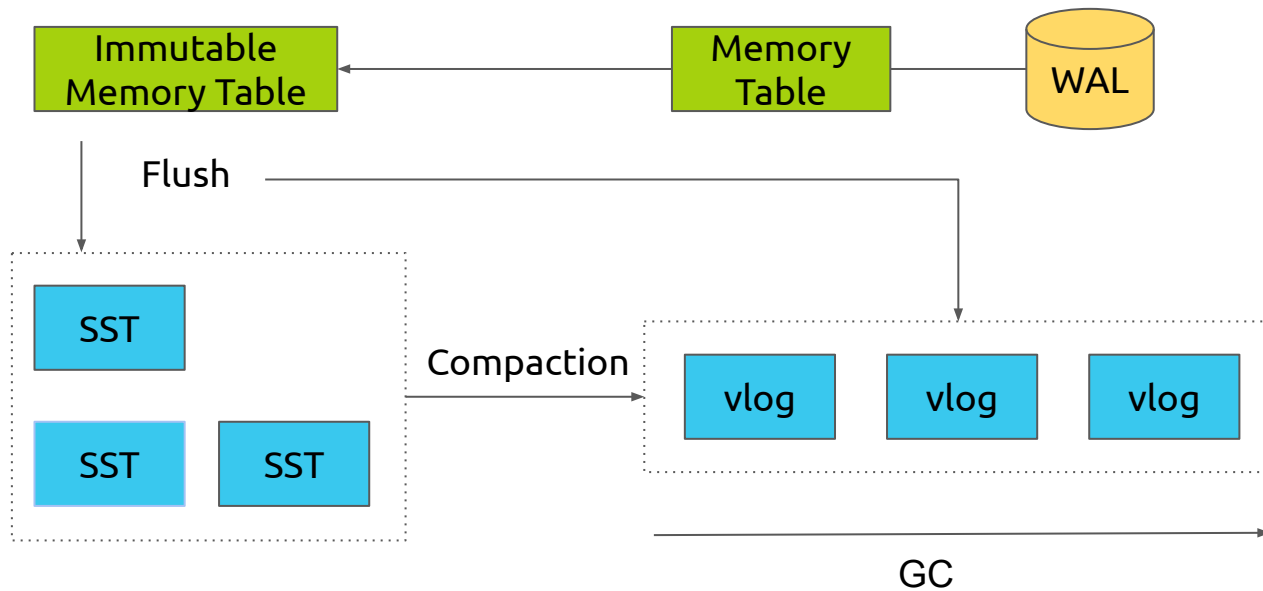
LSM





RocksDB

Titan



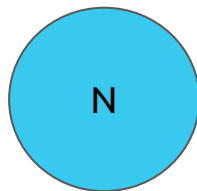
<https://github.com/pingcap/rust-rocksdb>



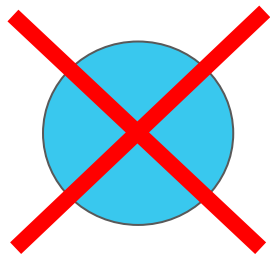
Raft



One Node



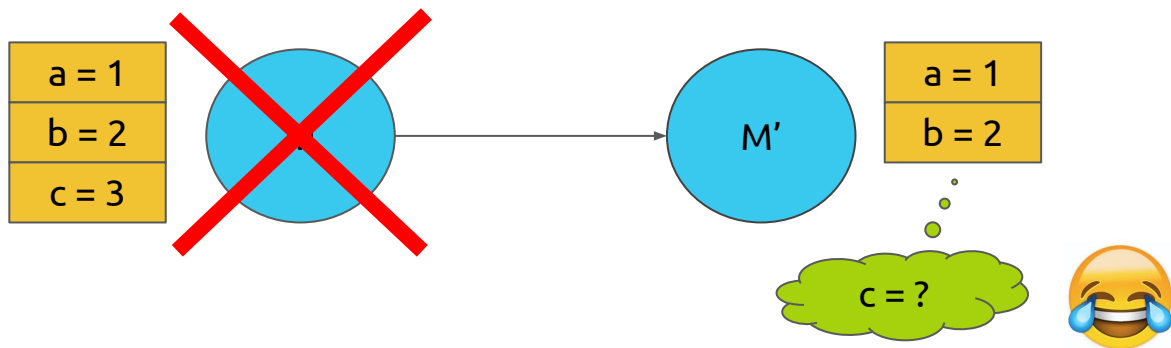
One Node



Async Replication



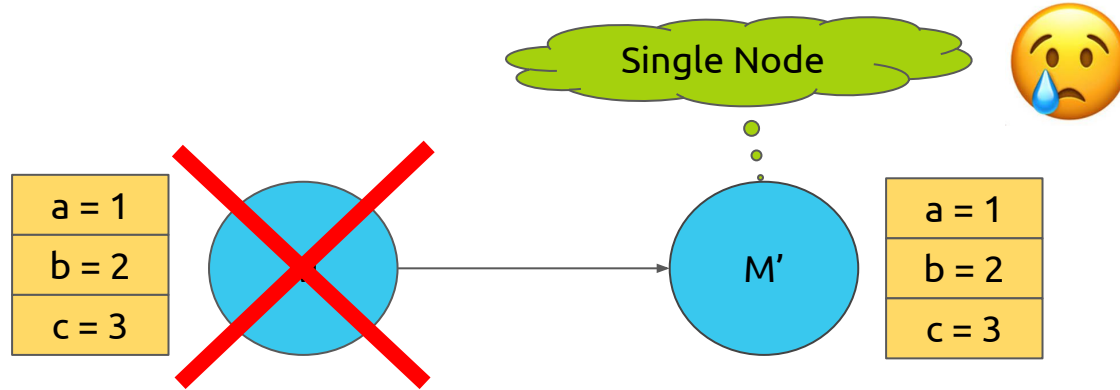
Async Replication



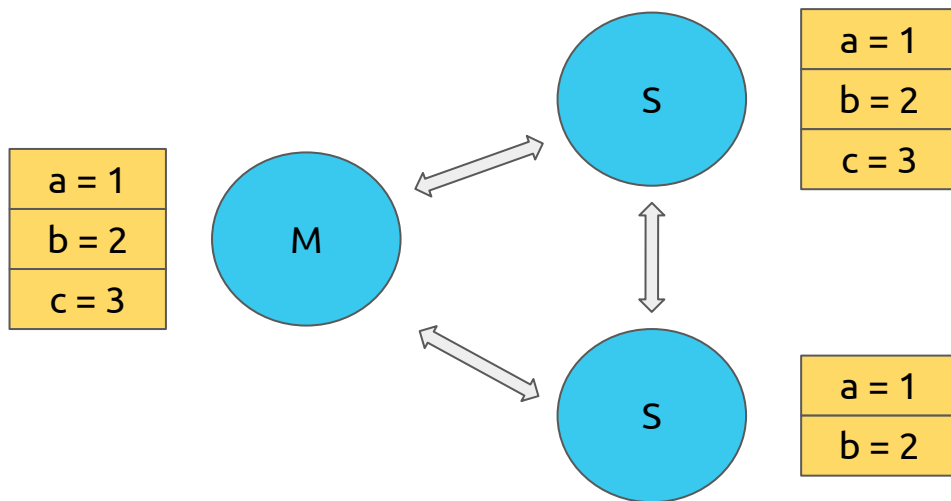
Sync Replication



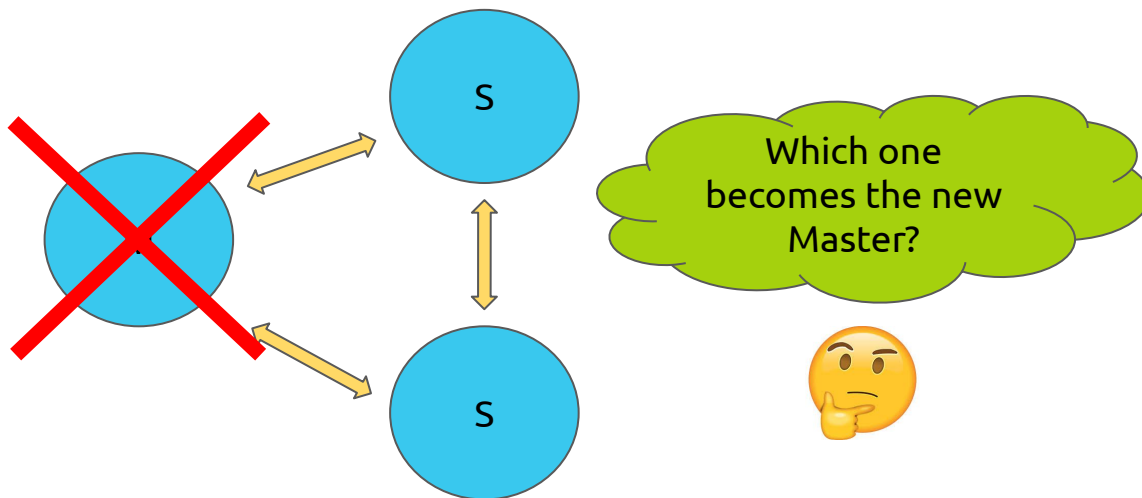
Sync Replication



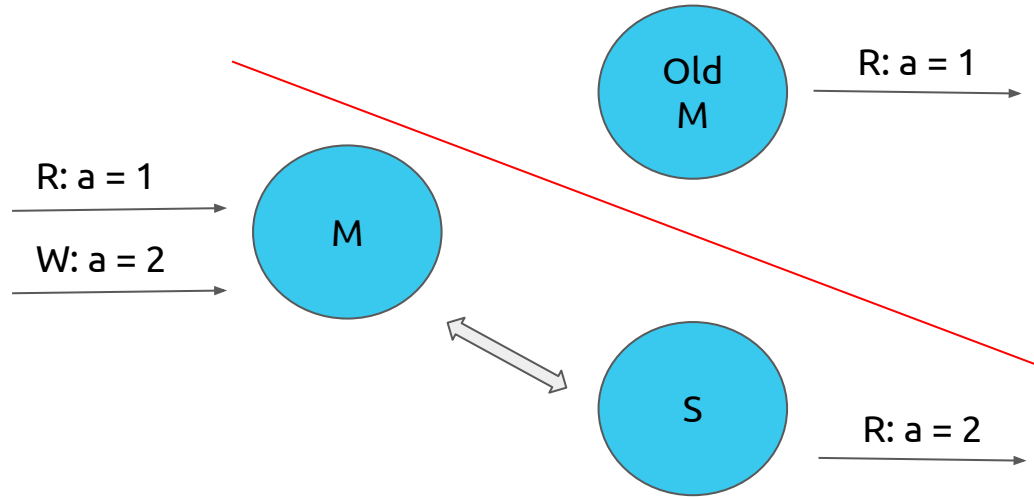
Quorum, $2N + 1$



Failover



Brain Split





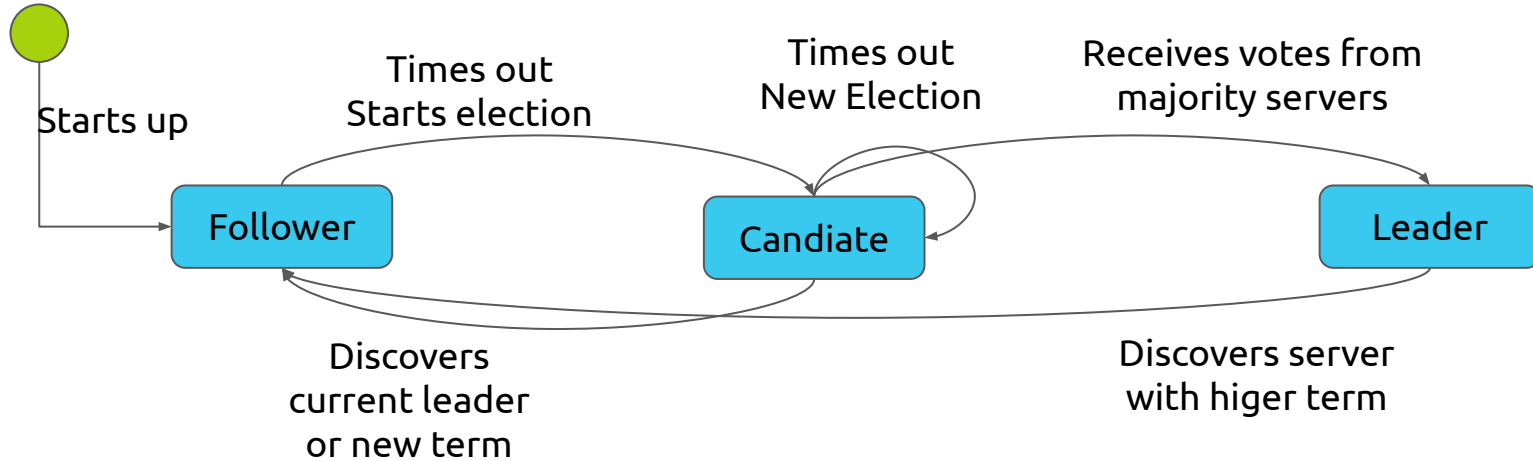
Key Points

- Election
- Replication

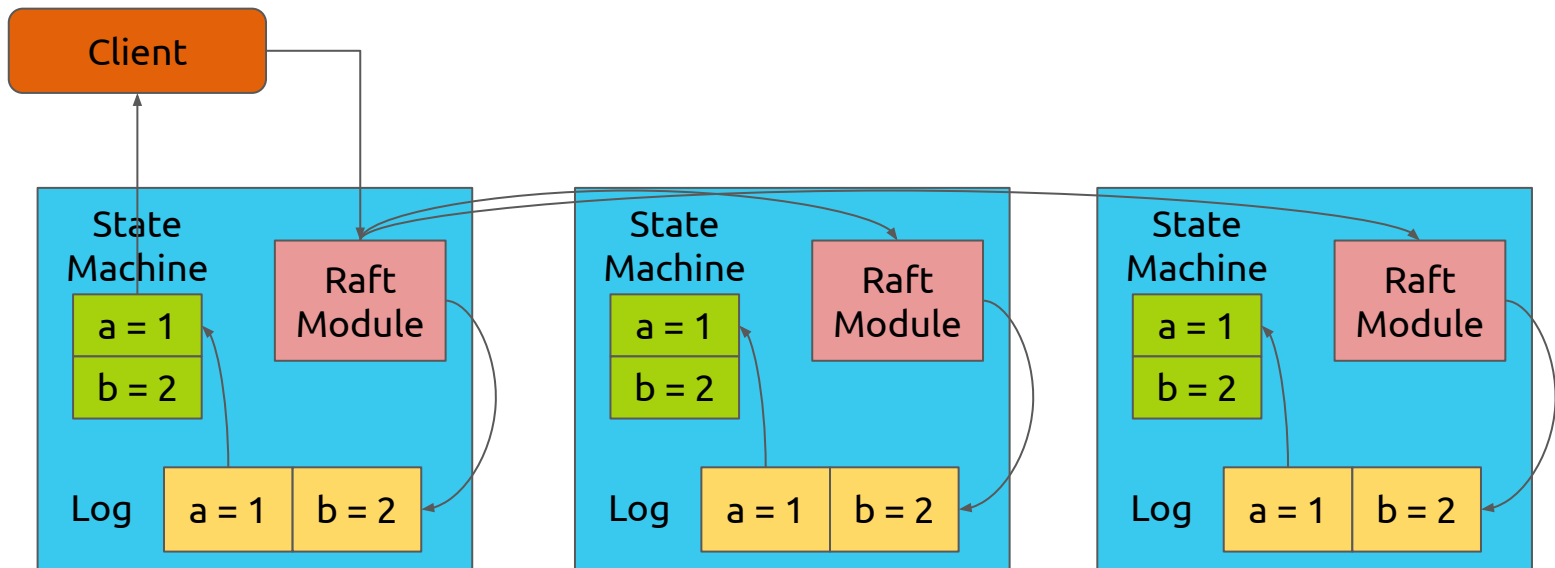
Role

- Leader
 - Only one Leader
 - Elected by the majority of the peers
 - Handles all the client requests
- Follower
 - Receives Replicated Logs from the Leader
- Candidate
 - Campaigns to become leader

Election



Log Replication

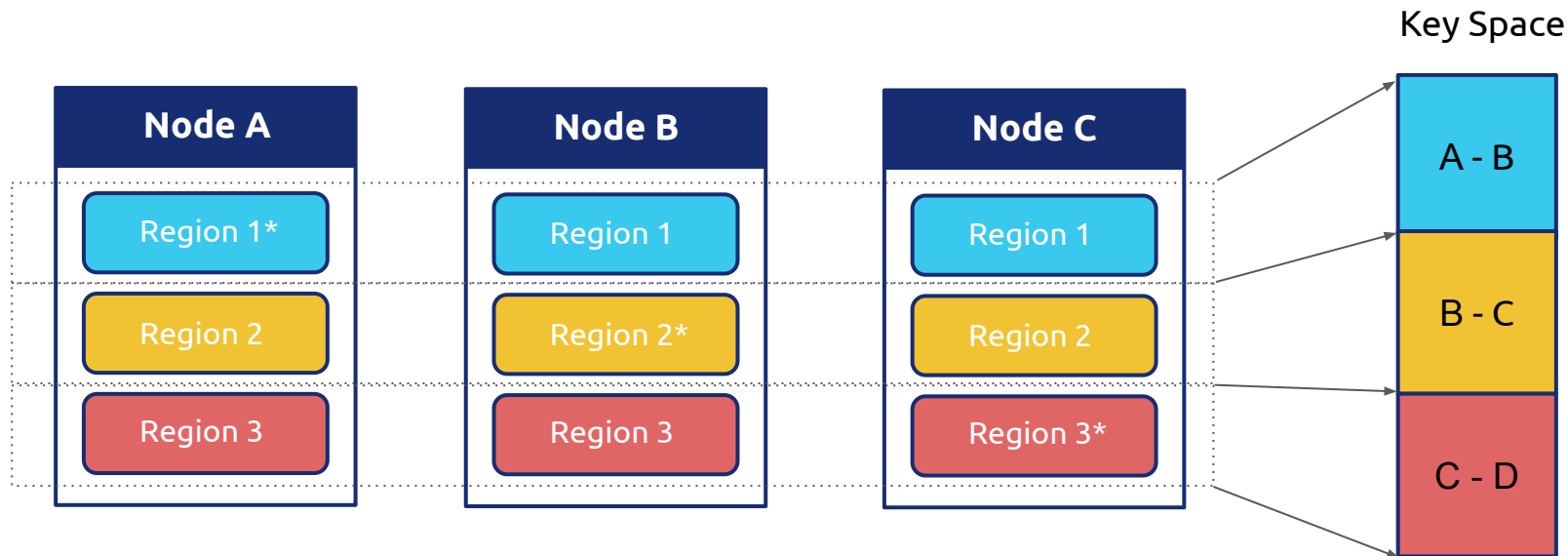


<https://github.com/pingcap/raft-rs>

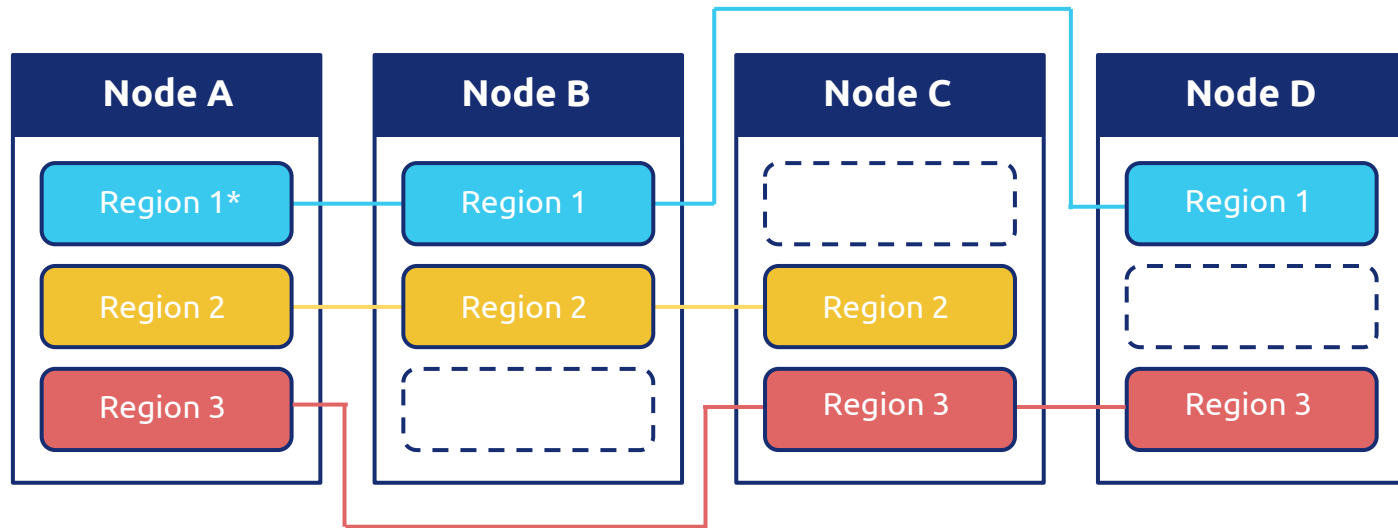
Multi Raft



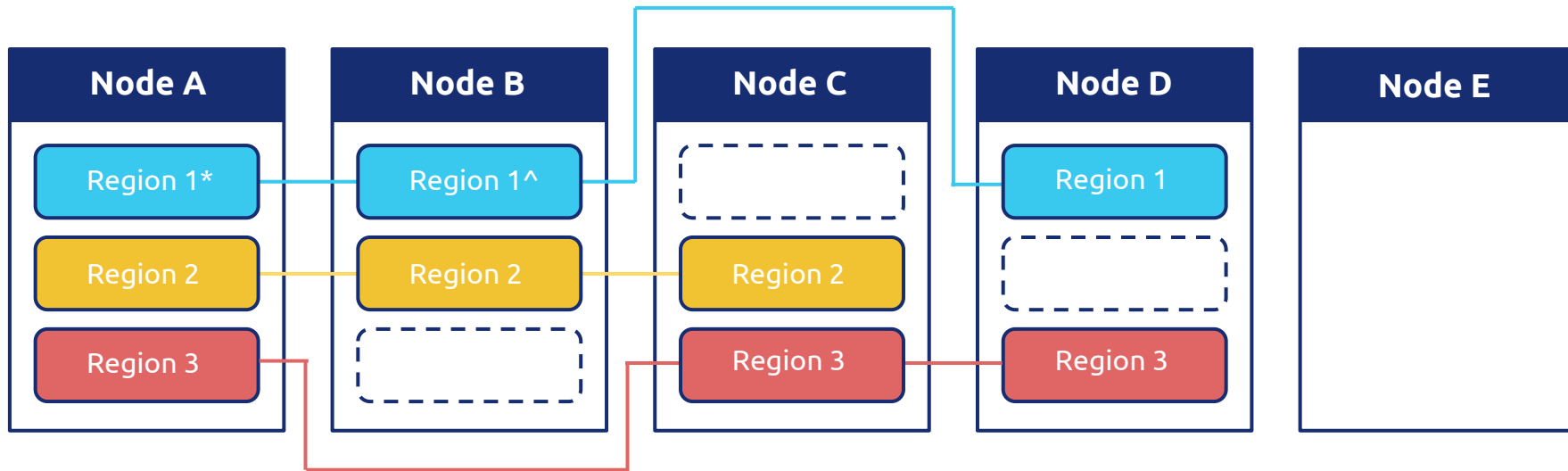
Multi-Raft



Scale-out

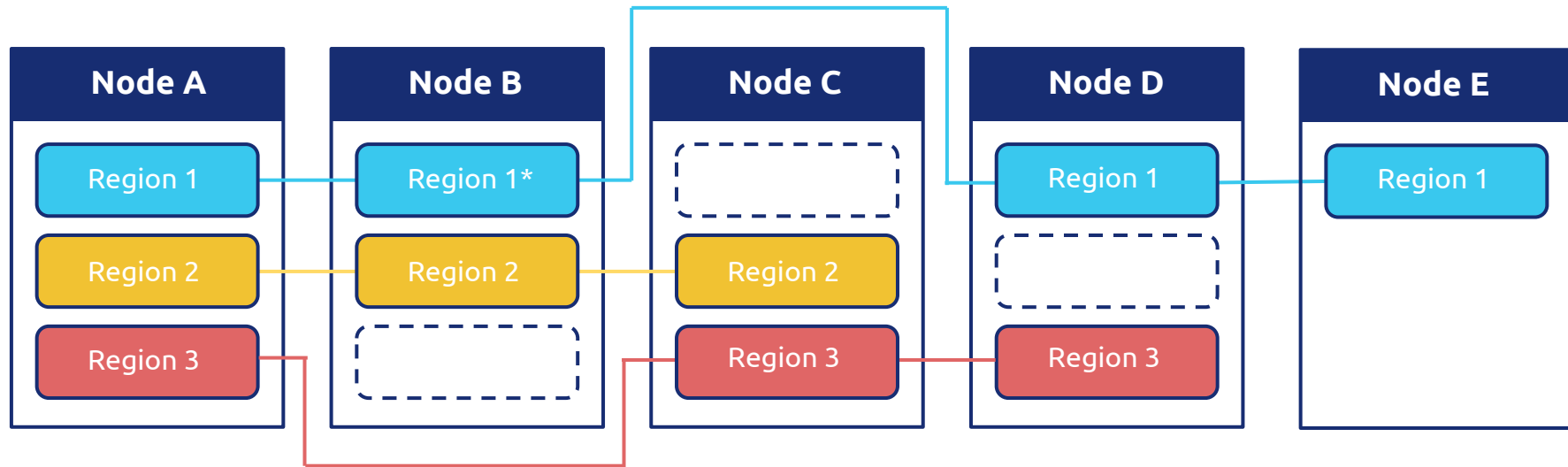


Scale-out



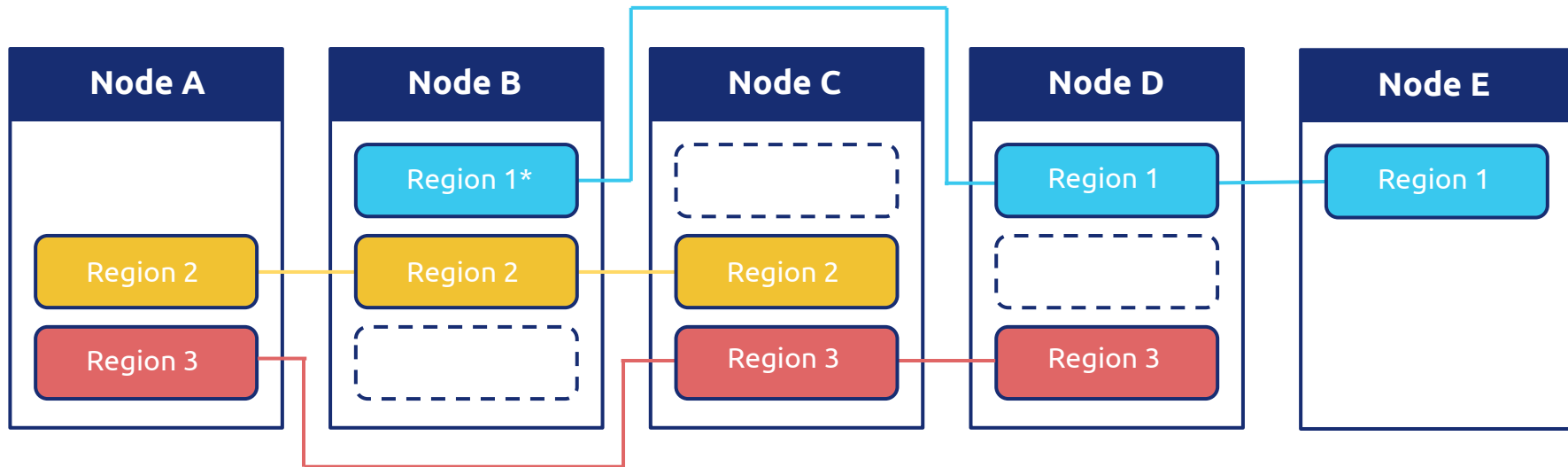
1. Transfer leadership of region 1 from Node A to Node B.

Scale-out



2. Add Replica on Node E.

Scale-out

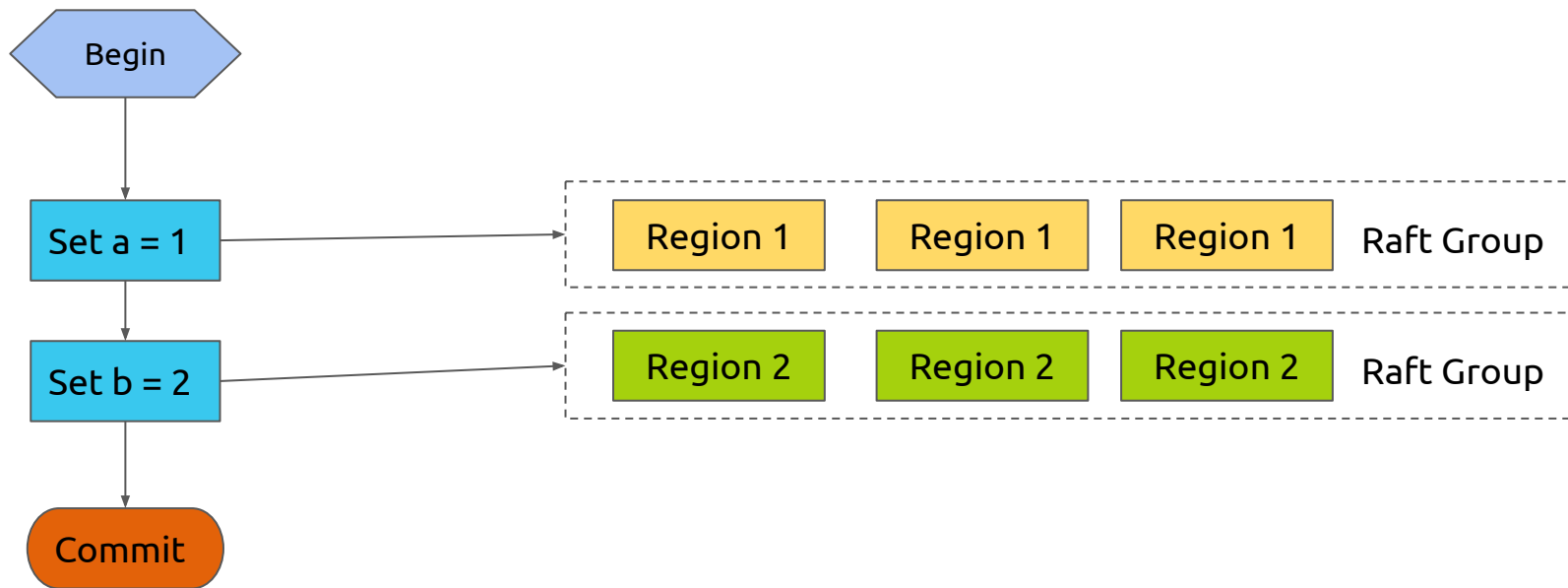


3. Remove replica from Node A.

Distributed Transaction



Why?



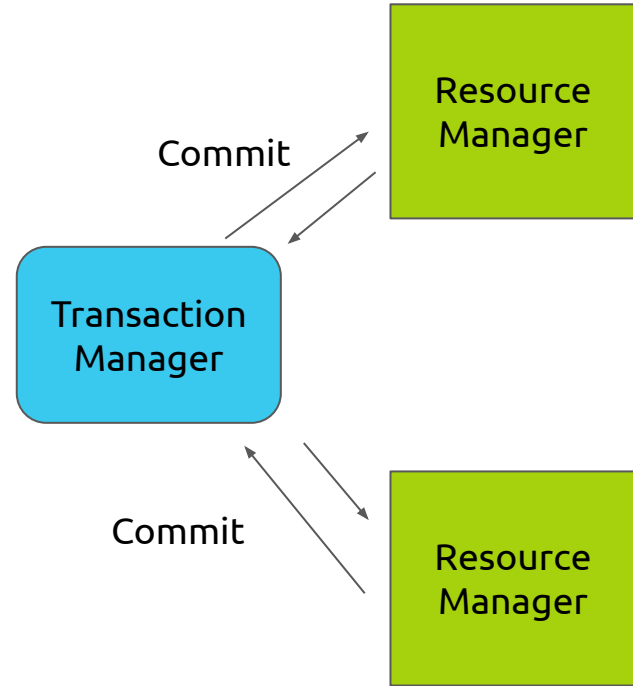
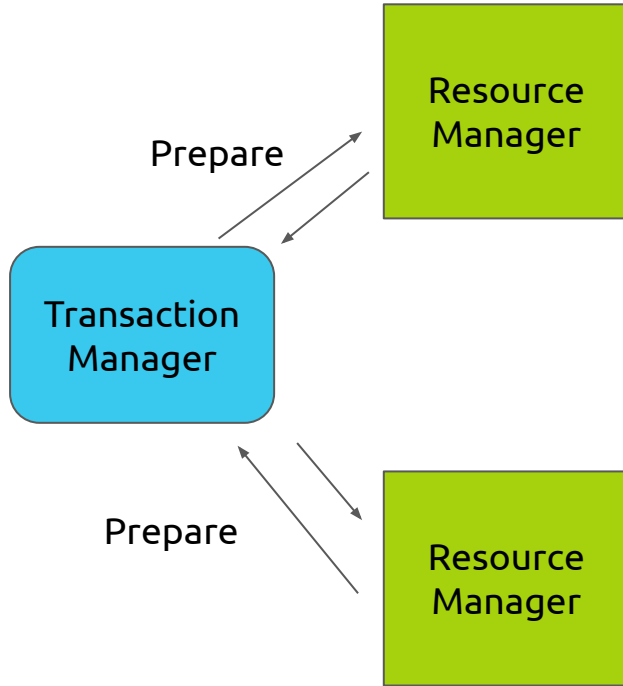
ACID



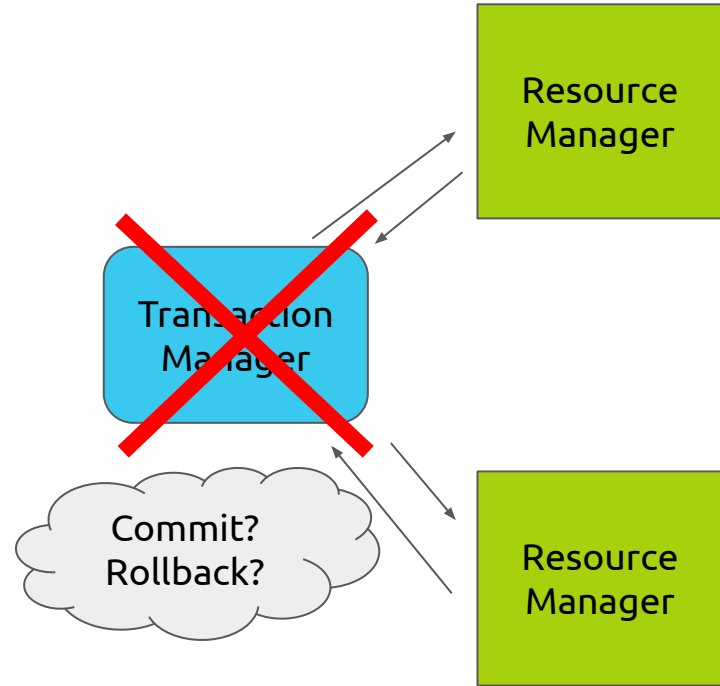
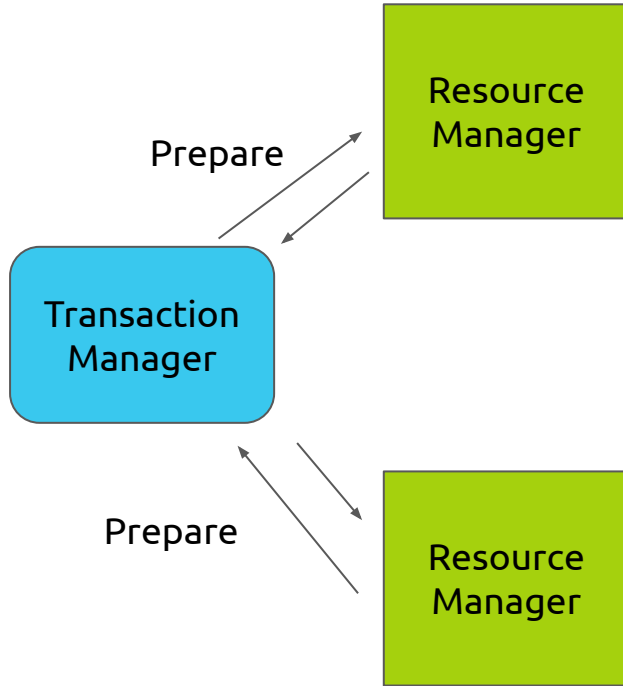
ACID

- Atomicity
- Consistency
- Isolation
- Durability

Two-Phase Commit



Two-Phase Commit



Communication



↑GRPG↓

gRPC

- Mode
 - Unary
 - Client streaming
 - Server streaming
 - Duplex streaming
- Using Futures to wrap the asynchronous C gRPC API

```
let f = unary(service, method, request);  
let resp = f.wait();
```

<https://github.com/pingcap/grpc-rs>

Monitoring





Prometheus

Prometheus

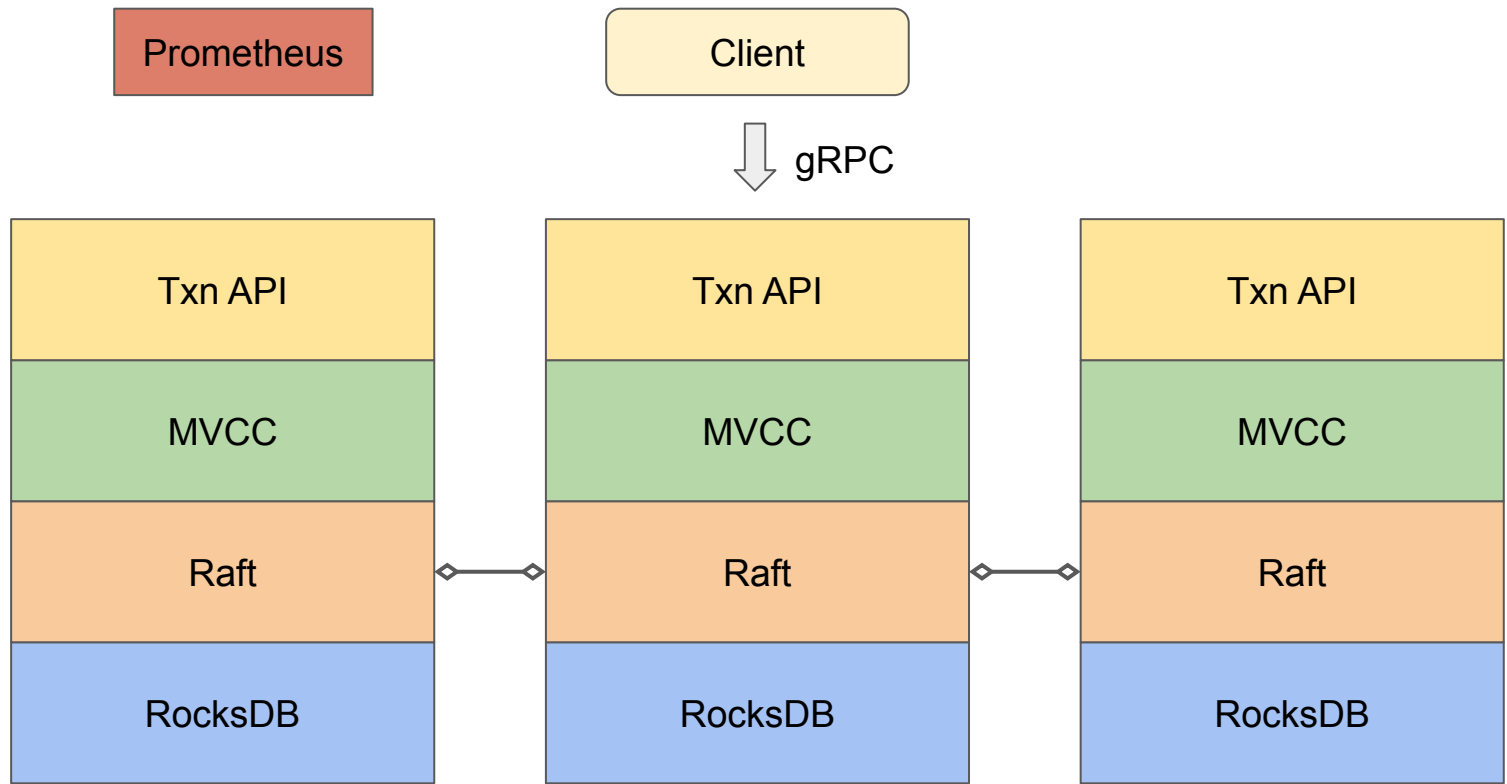
- Type
 - Counter
 - Gauge
 - Histogram

```
lazy_static! {  
    static ref HTTP_COUNTER: Counter = register_counter!(  
        "http_request_total",  
        "Total number of HTTP request."  
    ).unwrap();  
}  
  
HTTP_COUNTER.inc();
```

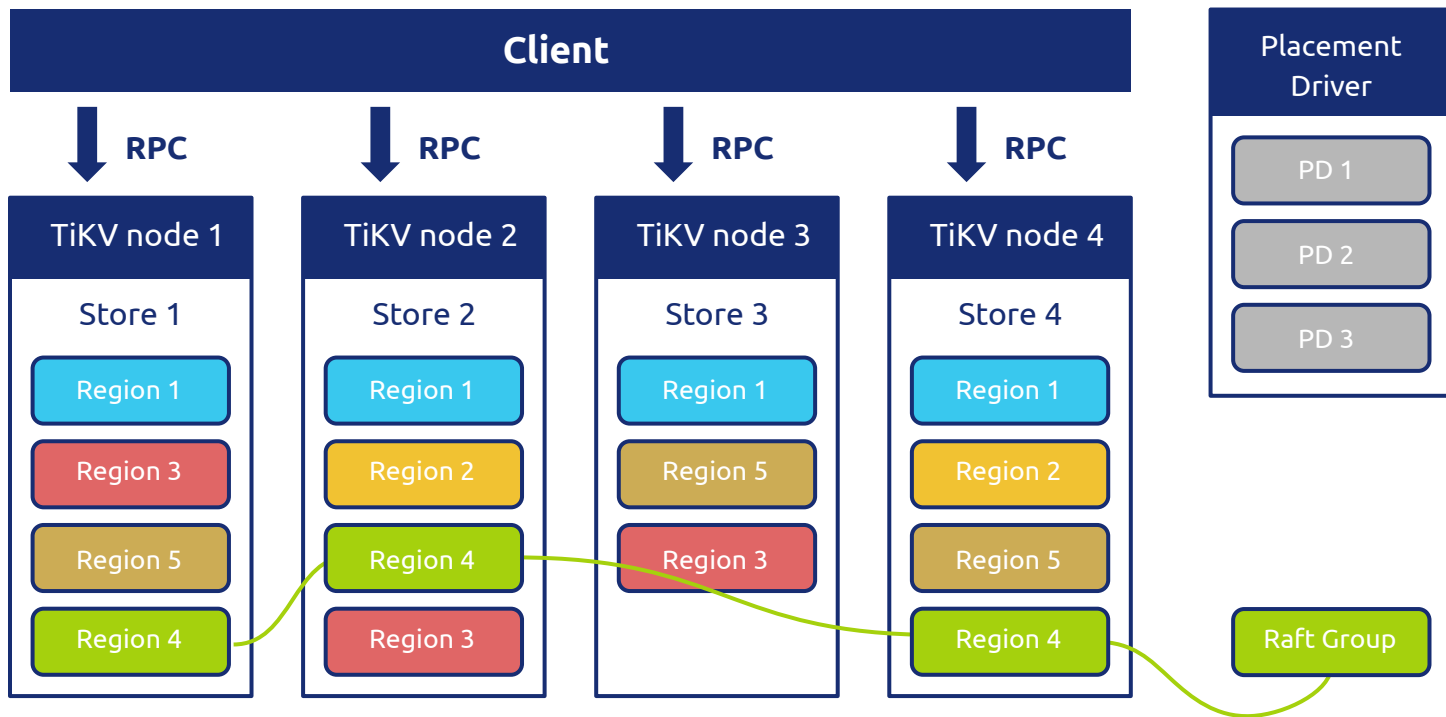
<https://github.com/pingcap/rust-prometheus>

TiKV





TiKV: The whole picture





CLOUD NATIVE
COMPUTING FOUNDATION

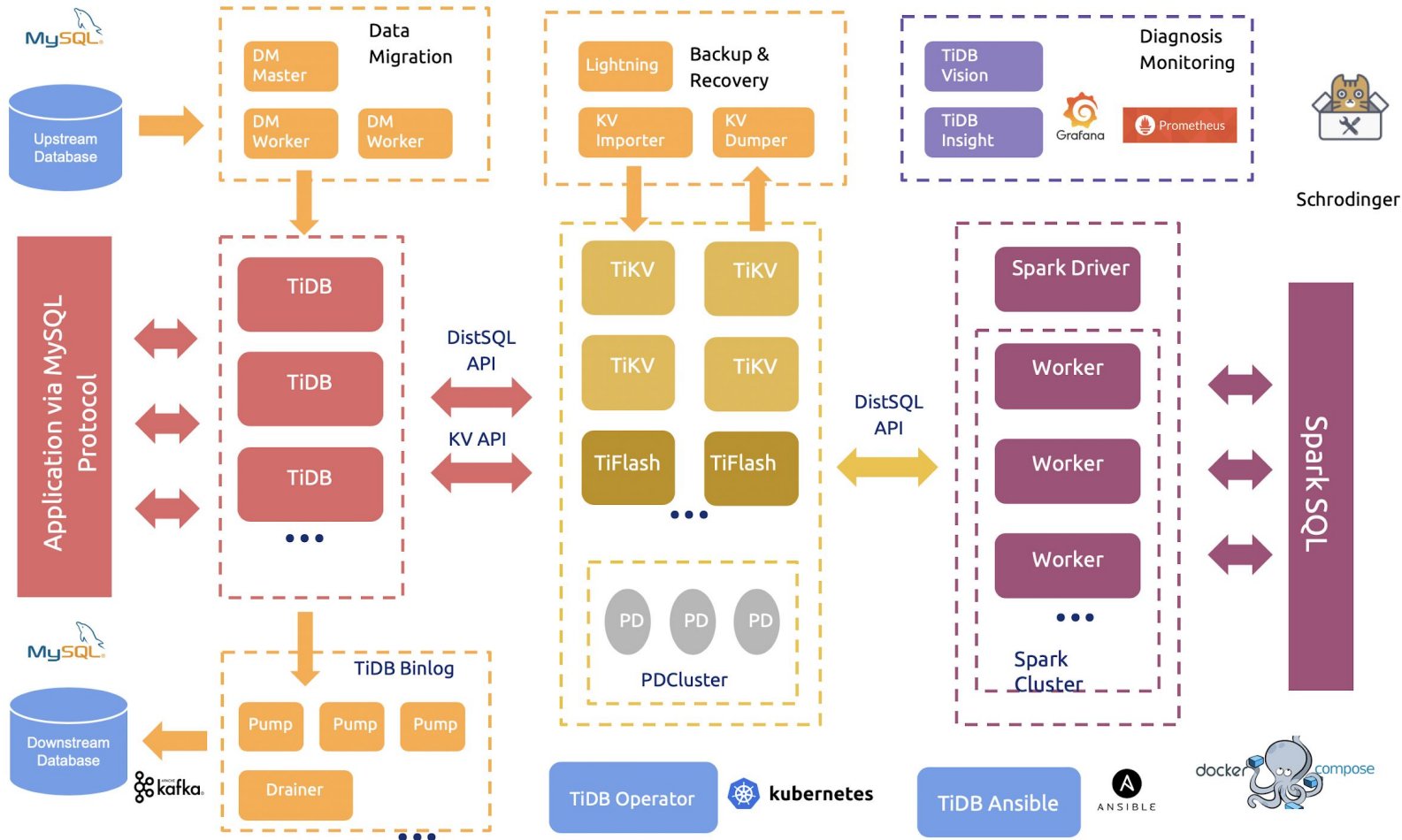


KV

<https://github.com/pingcap/tikv>

Beyond TiKV





Thank You !

<https://github.com/pingcap/tidb>

<https://github.com/pingcap/tikv>

