

Using Rust to Build a Distributed Transactional Key-Value Database

LiuTang | tl@pingcap.com

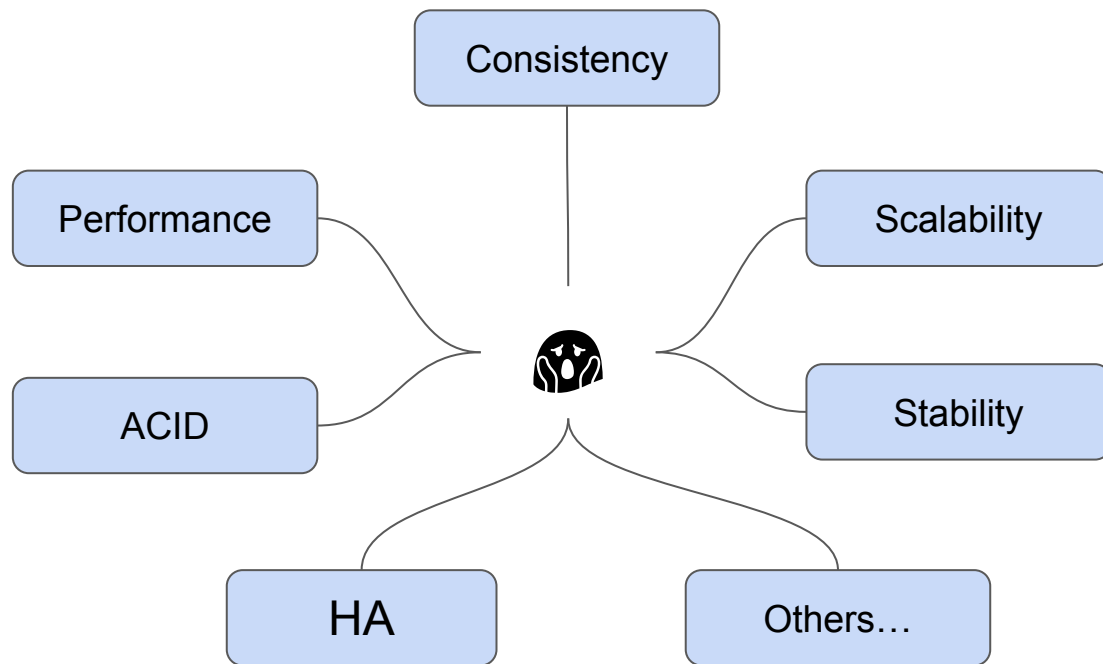
About me

- Chief Architect at PingCAP
- TiDB and TiKV
- Open source projects
 - LedisDB
 - go-mysql
 - go-mysql-elasticsearch
 - rust-prometheus
 - ...

Agenda

- Introduction
- Hierarchy
 - Storage
 - Raft
 - Transaction
 - RPC Framework
 - Monitor
 - Test
- Combine them all

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



A High Building, A Low Foundation

Language

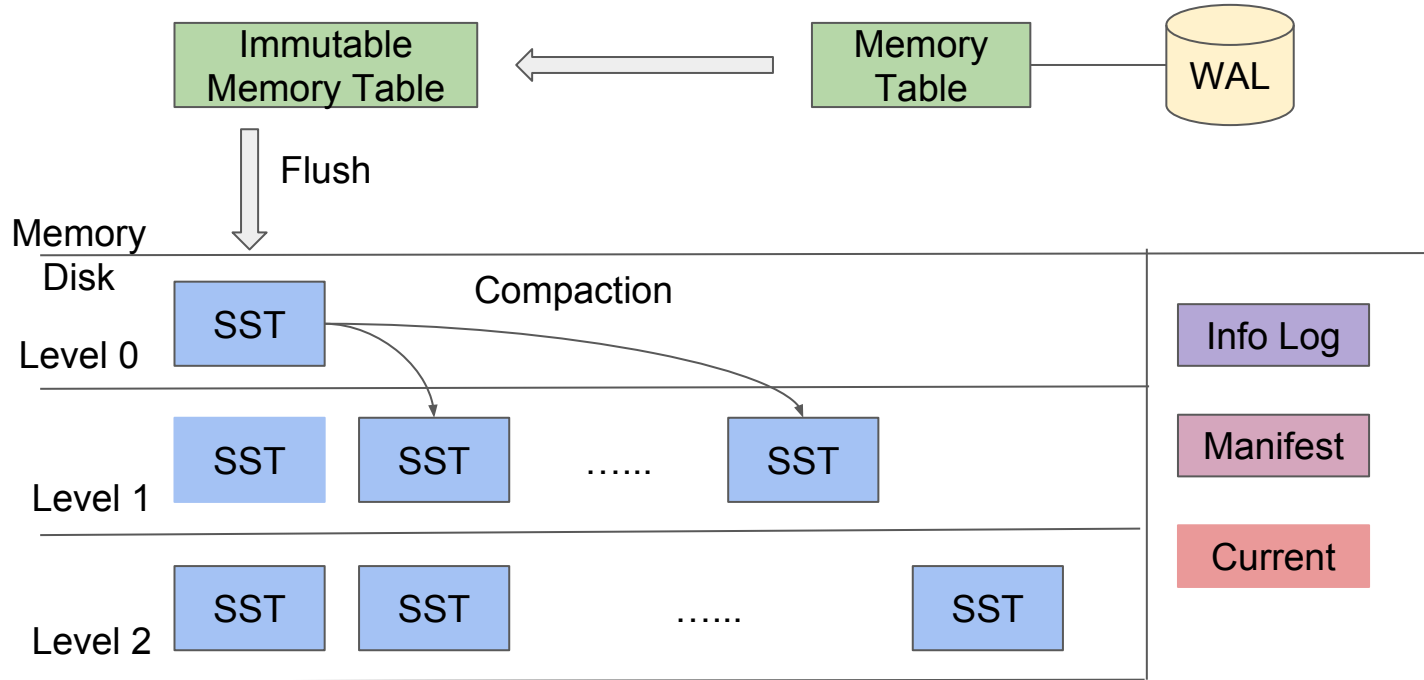


Let's start from scratch!!!



RocksDB

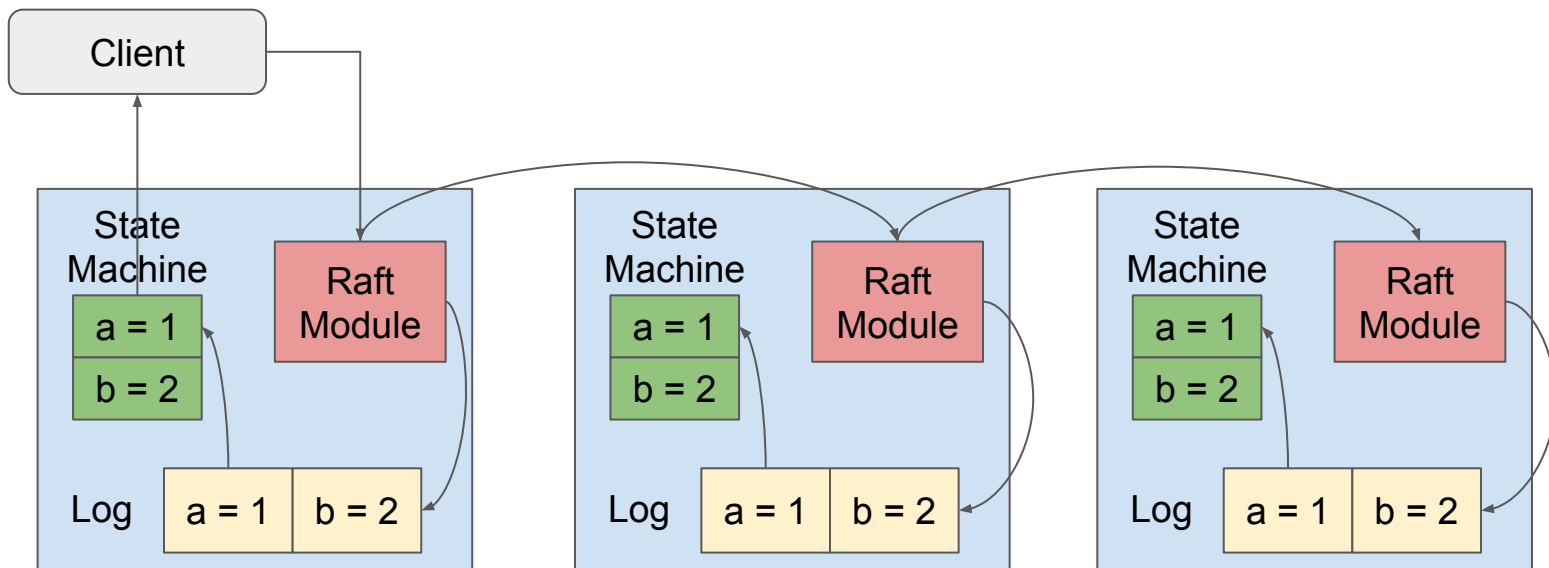
RocksDB



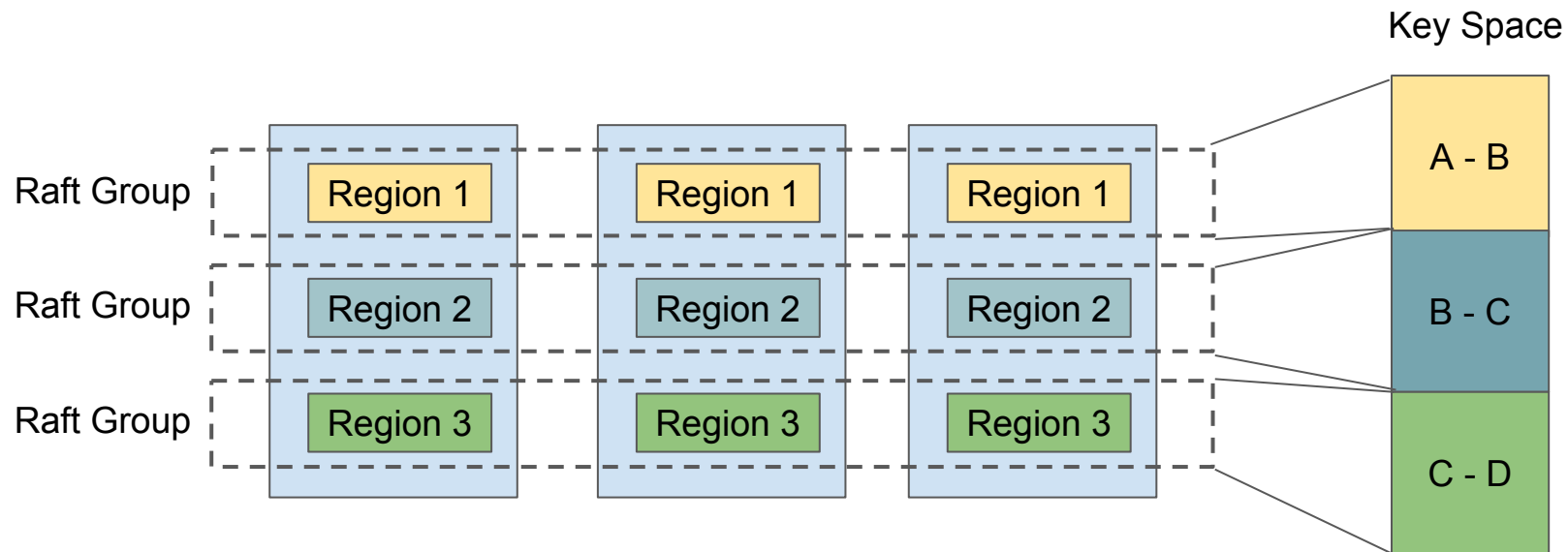
<https://github.com/pingcap/rust-rocksdb>



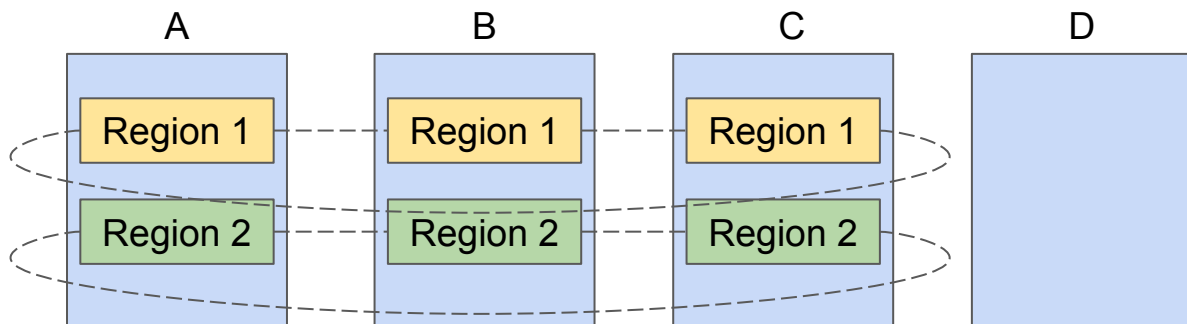
Raft



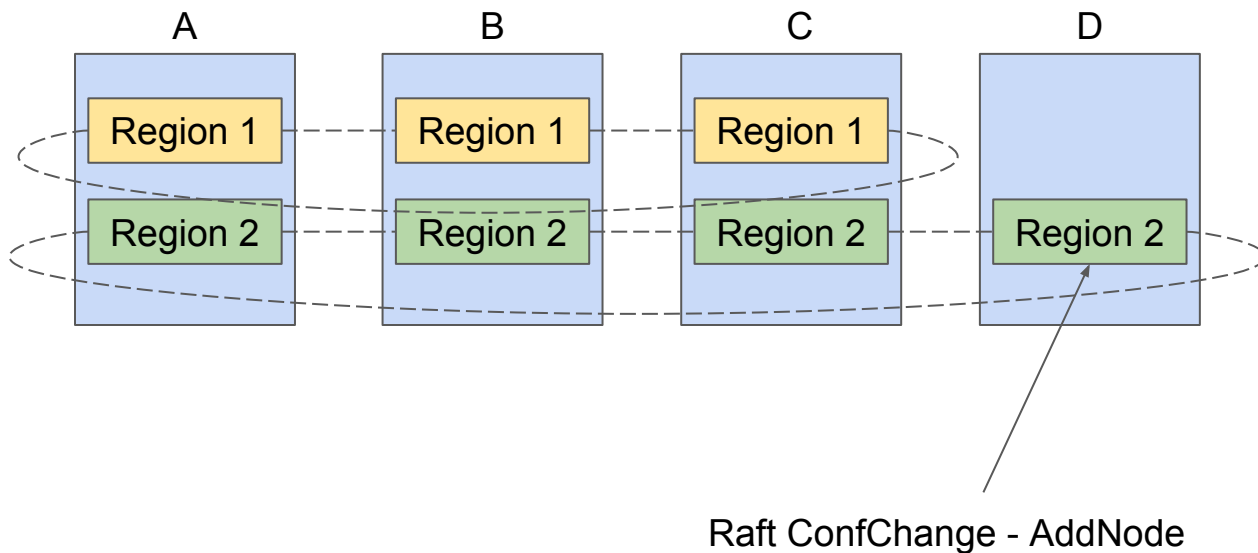
Multi-Raft



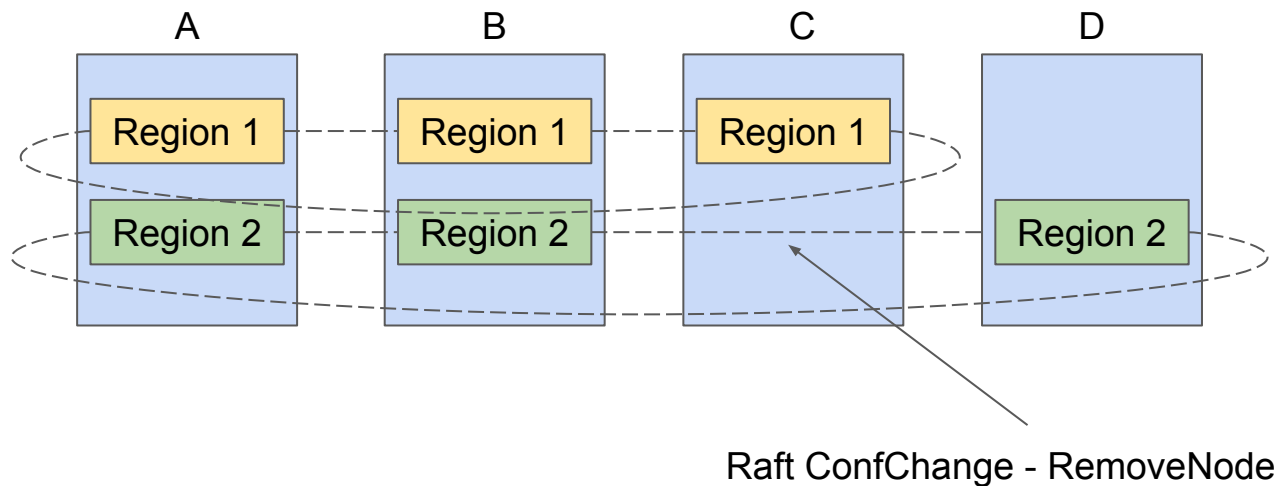
Multi-Raft - Scalability



Multi-Raft - Scalability



Multi-Raft - Scalability




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

Transaction

Transaction

```
let mut txn = store.begin()
let value1 = txn.get(region1_key)
let value2 = txn.get(region2_key)
// do something with value
txn.set(region1_key, new_value1)
txn.set(region2_key, new_value2)
txn.commit()
// or txn.rollback()
```

How to keep consistency
crossing multi-Raft Groups?



Transaction

1. Inspired by Google Percolator
2. Optimized Two Phase Commit (2 PC)
3. Multiversion Concurrency Control (MVCC)
4. Snapshot Isolation
5. Optimistic Transaction



↑GRPC↓

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>



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>

Testing

Testing - Failure Injection

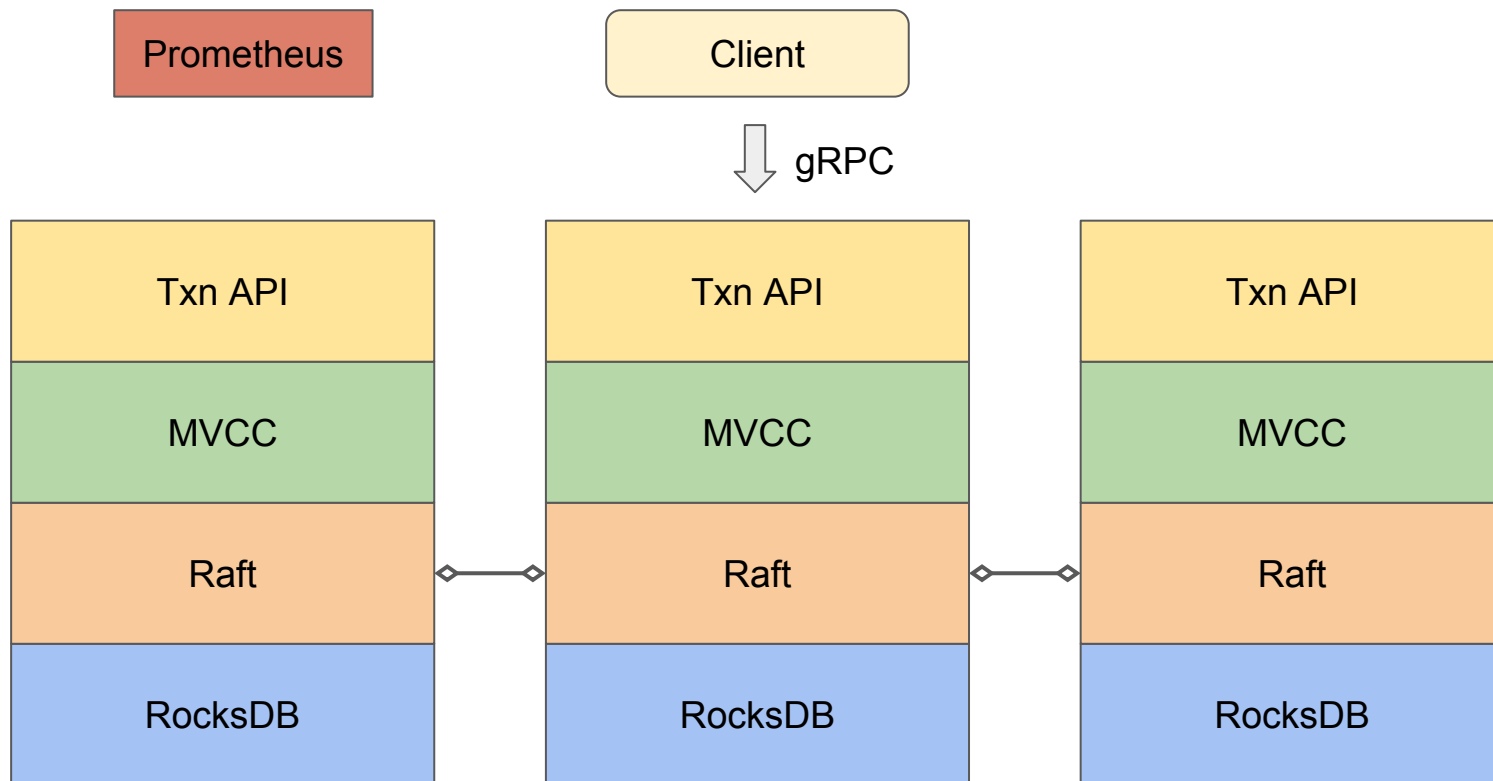
```
// Ingest a failure
fn function_foo() {
    fail_point!("foo");
}
```

```
// Run and Trigger the failure
FAILPOINTS=foo=panic cargo run
```

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

Architecture

Architecture

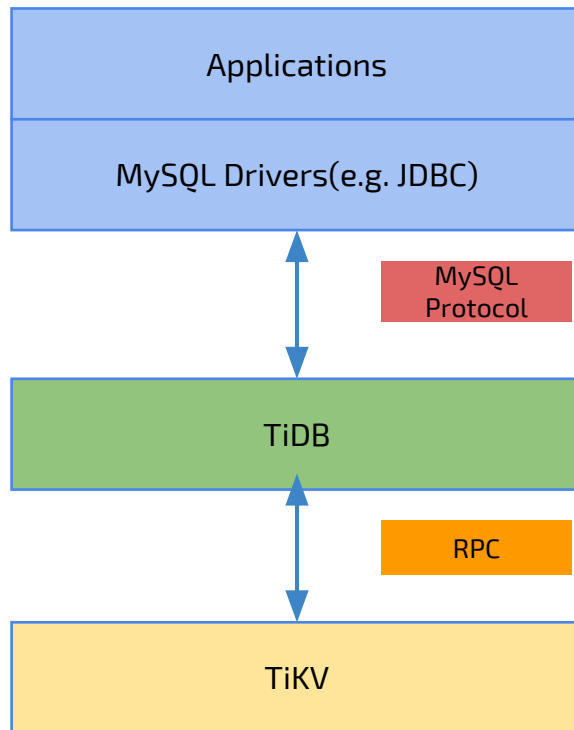


<https://github.com/pingcap/tikv>

Beyond TiKV

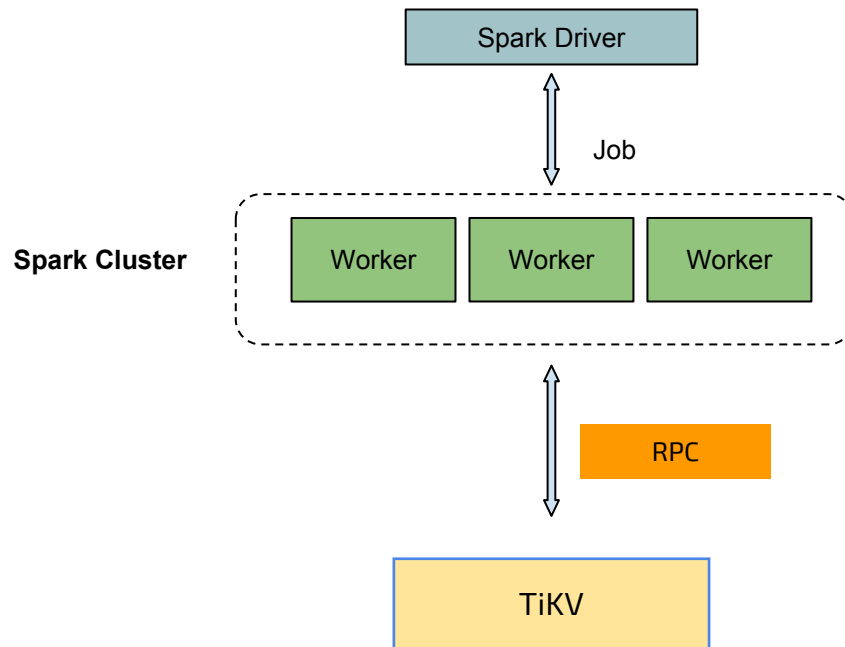
A Distributed Relational Database

TiDB

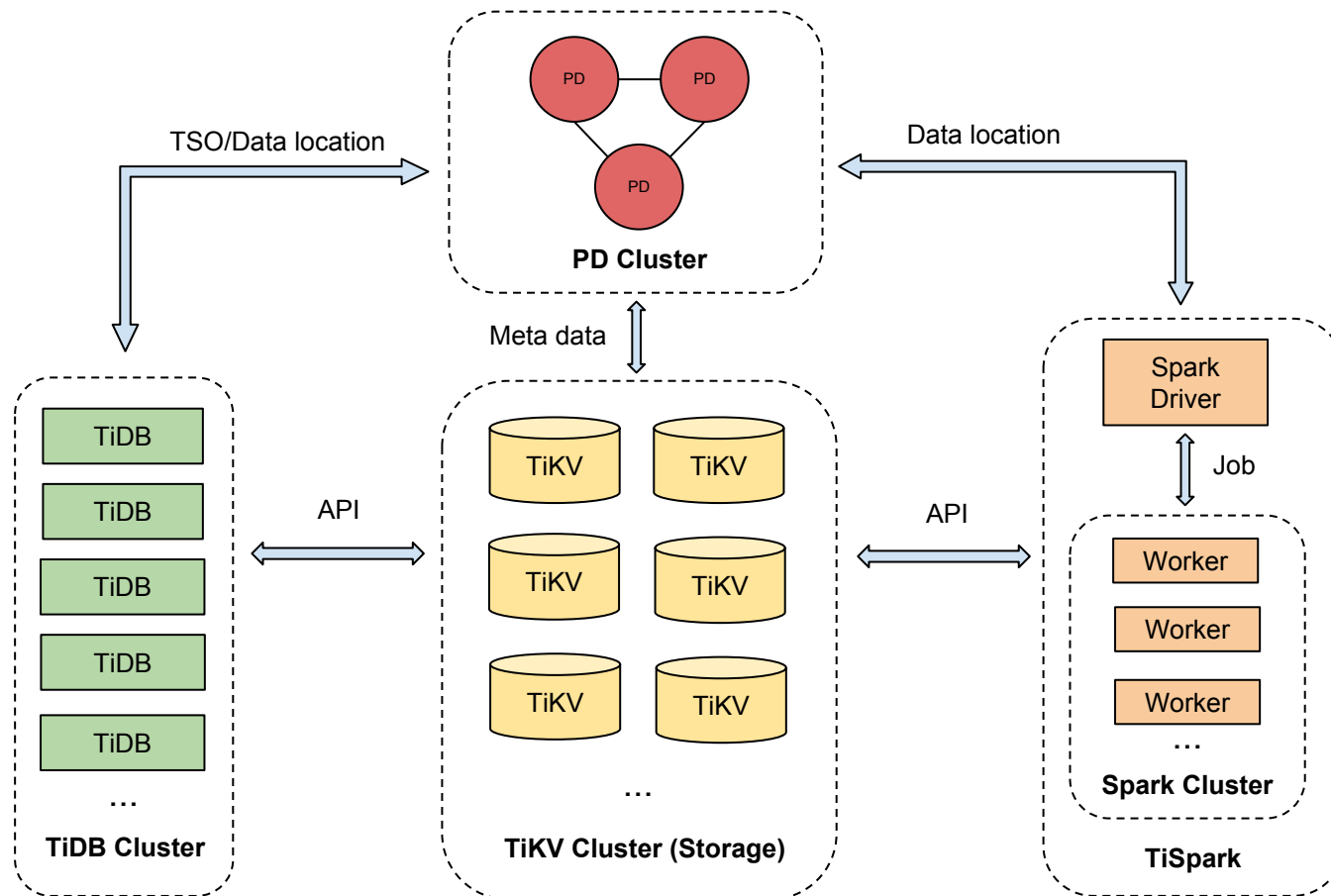


A Distributed Analytical Database

TiSpark



Hybrid Transactional/Analytical Processing Database



Thank you!

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

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

We are hiring...

@China @Silicon Valley @Home