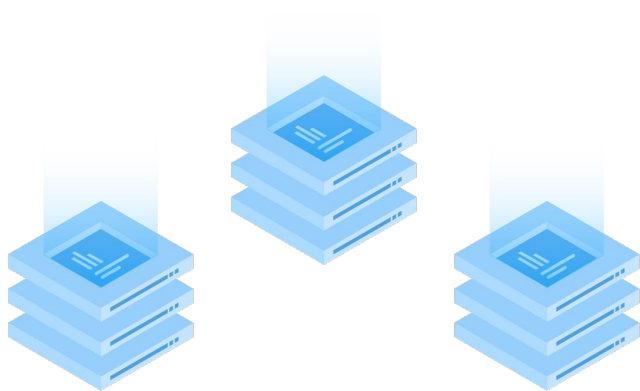


# How Does Geo-replication Work in TiDB

Presented by Jay Lee



# About Me



**Jay Lee (李建俊)**

Senior Engineer of PingCAP

Distributed system engineer / Database engineer/

Open- source advocator

TiKV Maintainer

Rustacean: raft-rs / grpc-rs

Github @busyjay



# Overview

- What is TiDB
- How does replication work
- Deployment
- Q&A



# Part I - Intro to TiDB

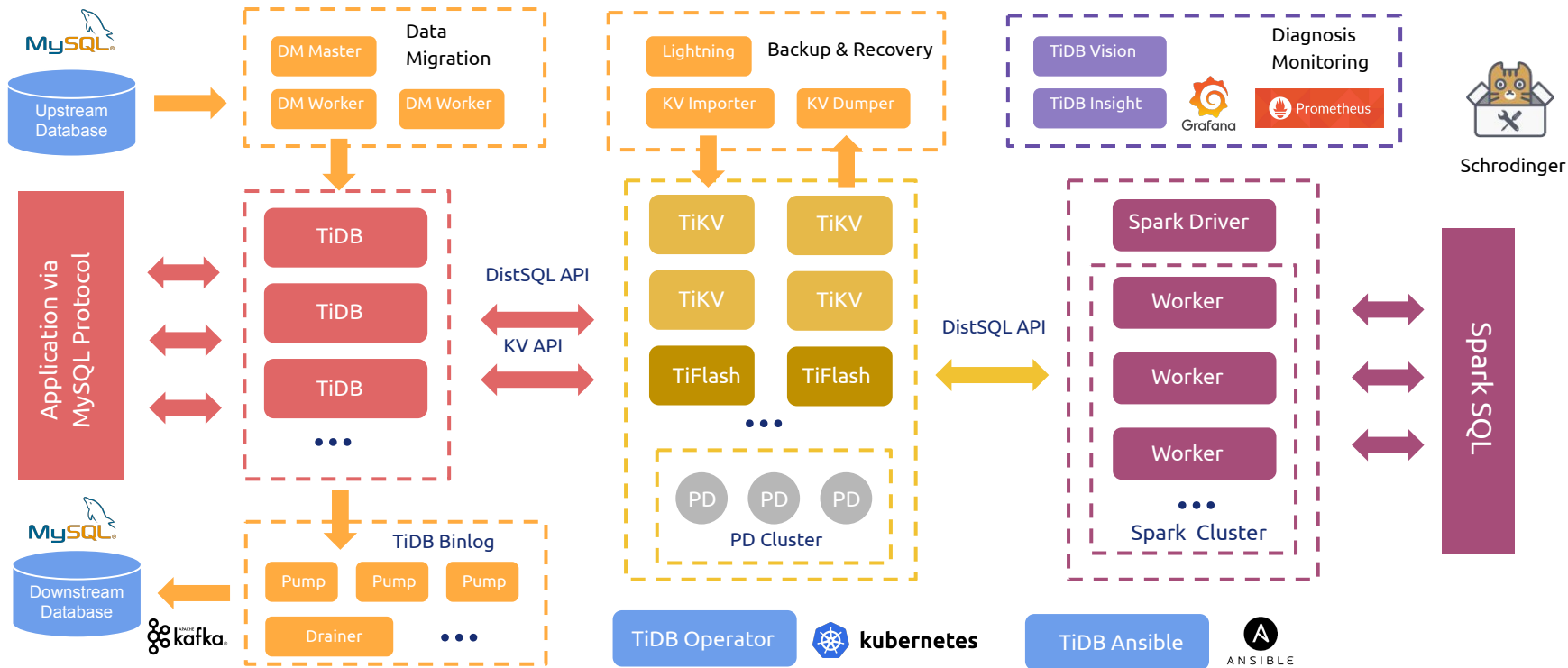


# What is TiDB

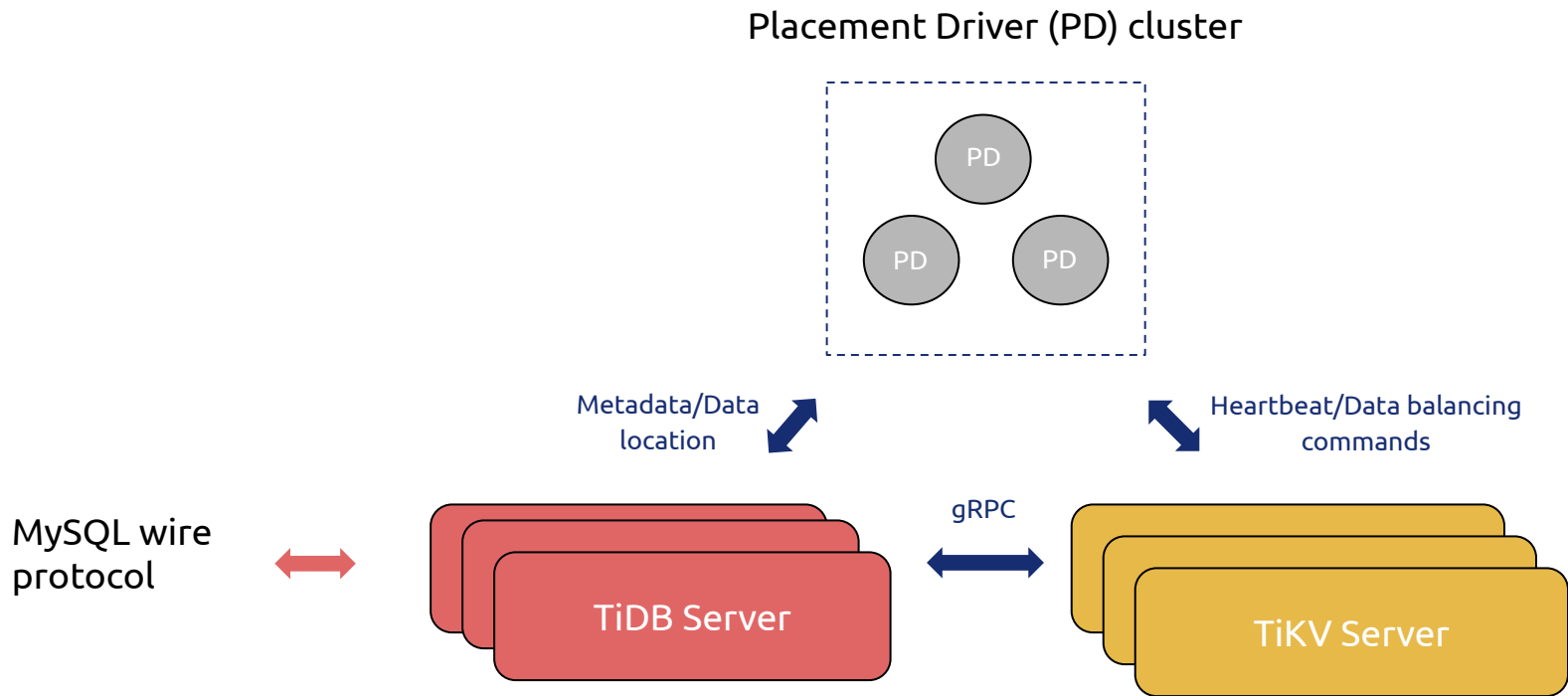
- Elastic scaling-out
  - Transparent to applications, no more manual sharding!
- Always-on
  - HA with strong consistency
- SQL
  - MySQL dialect
  - HTAP = OLTP + OLAP
- ACID semantics



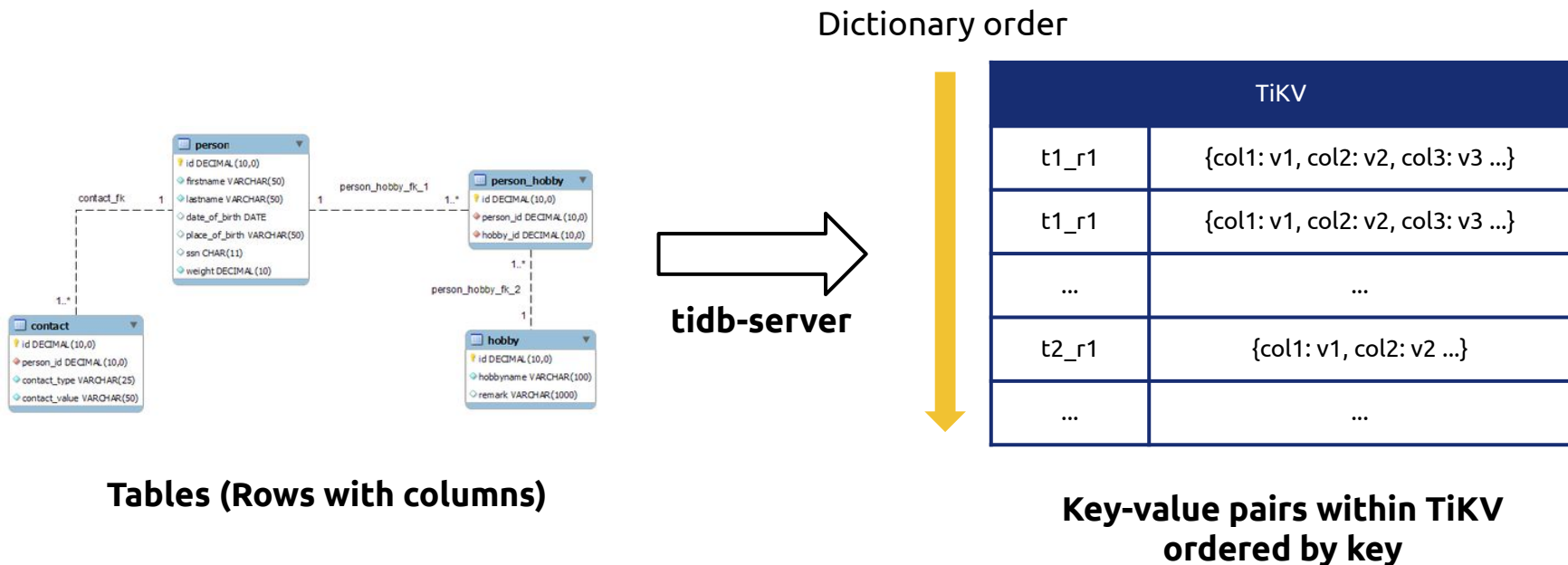
# The whole picture



# A closer review



# Data organization within TiDB





# Data organization within TiDB

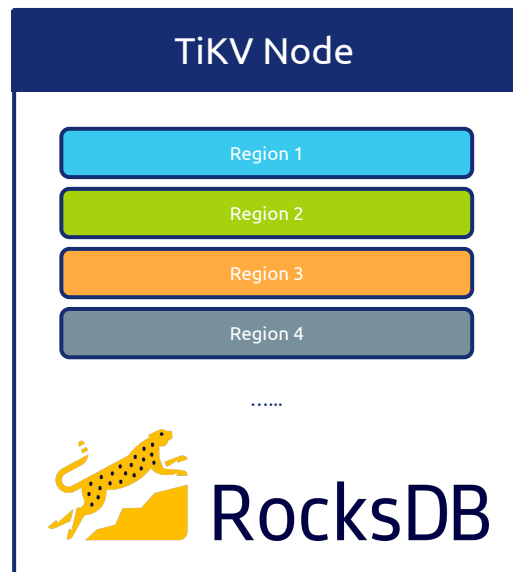


A diagram of a "Local RocksDB instance" showing a table with two columns. The table is divided into four regions, each indicated by a bracket on the right. Region 1 (light blue) contains rows (t1\_r1, v1), (t1\_r2, v2), and (...). Region 2 (light green) contains rows (t5\_r1, ...) and (t5\_r10, ...). Region 3 (light orange) contains rows (t1\_i1\_1\_1, ...), (t1\_i1\_2\_2, ...), and (...). Region 4 (light grey) contains rows (t1\_i6\_1\_3, ...) and (...).

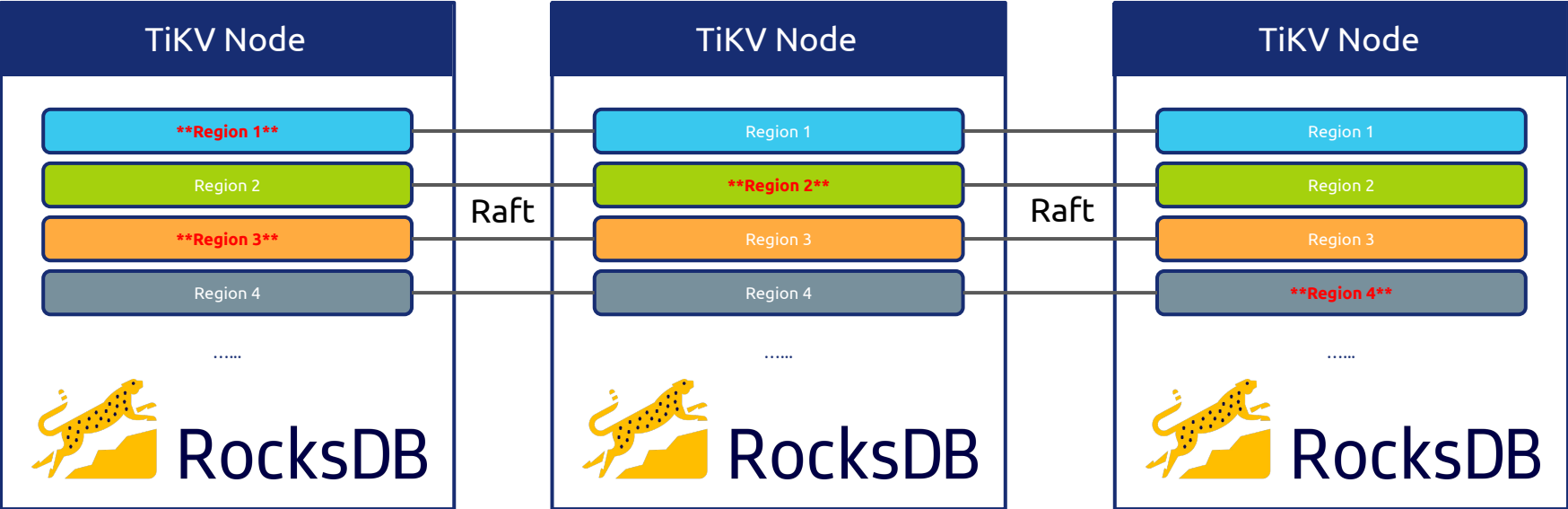
Local RocksDB instance	
t1_r1	v1
t1_r2	v2
...	...
t5_r1	...
t5_r10	...
t1_i1_1_1	...
t1_i1_2_2	...
...	...
t1_i6_1_3	...
...	...

# Data organization within TiDB

- Region (A bunch of key-value pairs, or Split)
  - Default total size of one Region: 96MB
    - You can change it in configuration
- Region is a logical concept
  - Region meta : [Start key, End key)
  - All regions within a TiKV node share a same **RocksDB** instance
- Each Region is a Raft group
  - Default: 3 replicas



# Region: multiple replicas across different nodes

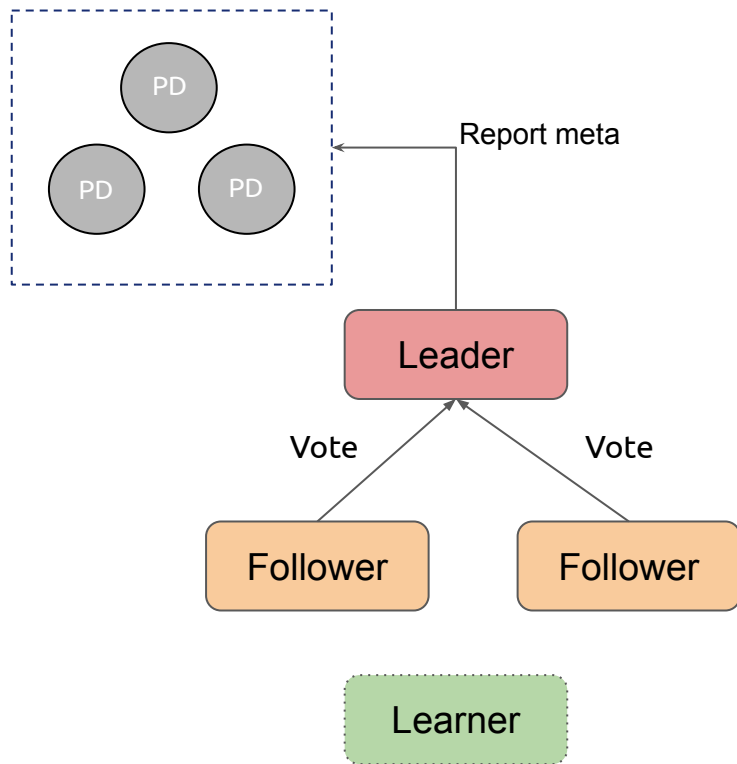


# Part II - How does replication work



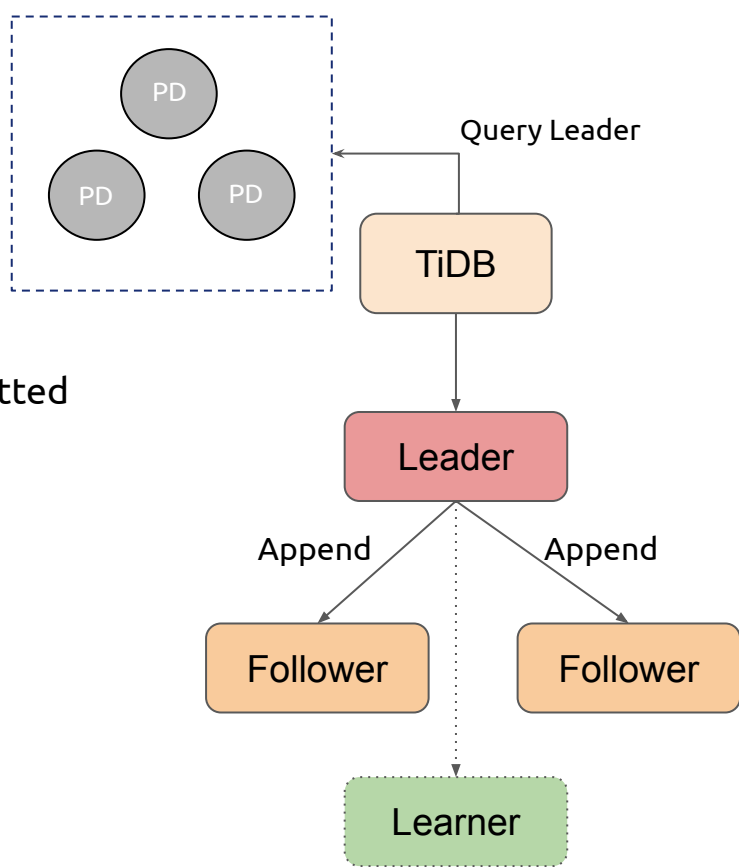
# Raft group

- A region is managed by a raft group
- Possible roles
  - Leader
  - Follower
  - Learner
- Leader receives votes from majority followers
- Leader manages group and report to PD
- On failure
  - Follower starts campaign and take over



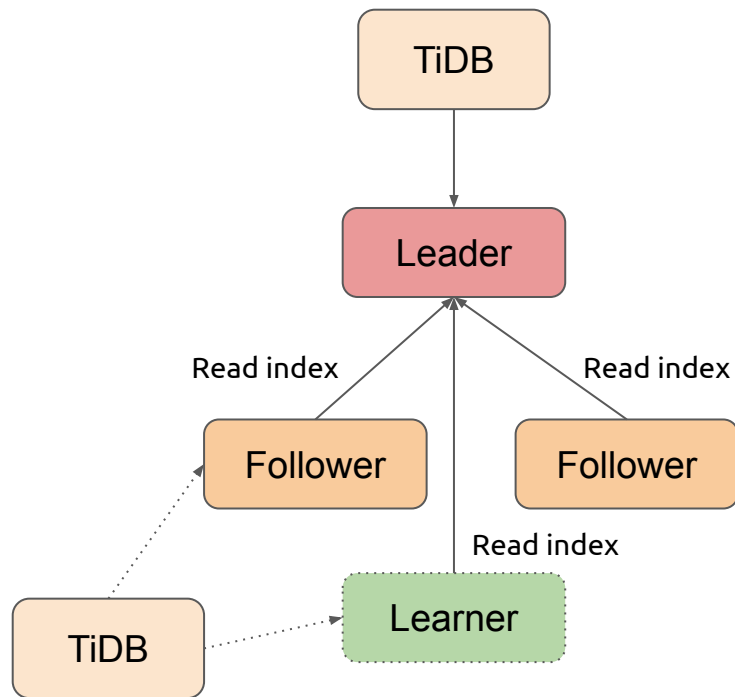
# Write

- Data is written to leader as logs
- Leader replicates logs to followers and learners
- Logs replicated to majority followers are committed



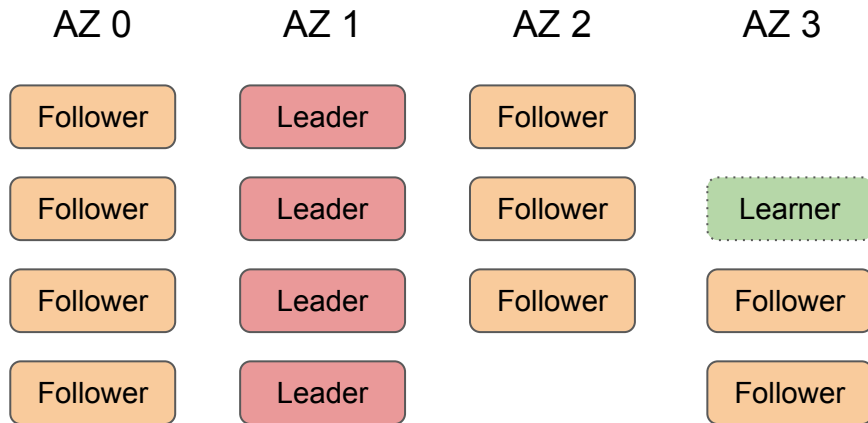
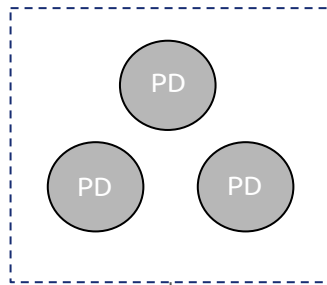
# Read

- All roles can read
- Read on leader
  - Read immediately if in lease
  - Renew lease otherwise
- Read on follower and learner
  - Read log index on leader
  - Read data on follower



# Configuration change

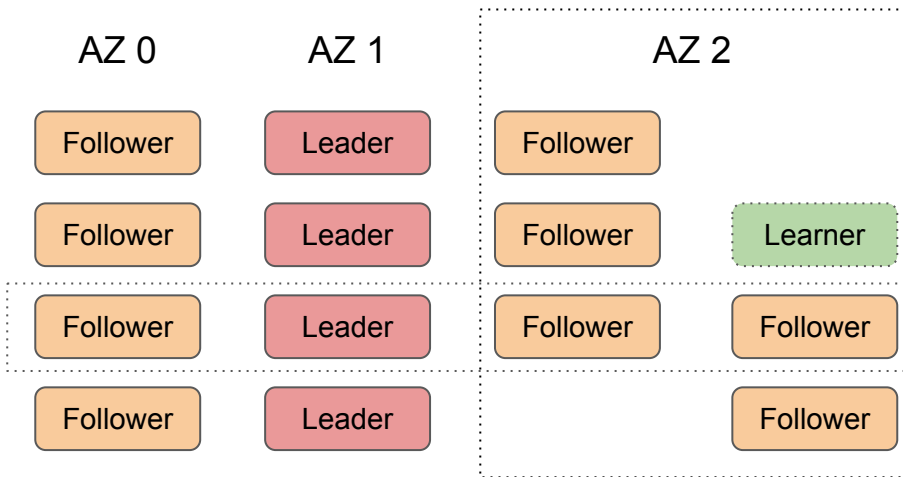
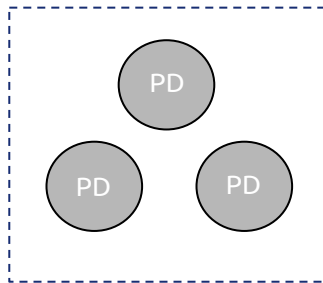
- Simple change
  - Can only handle one change at a time
  - Process as a special write
  - Change in one step
  - Quick and easy





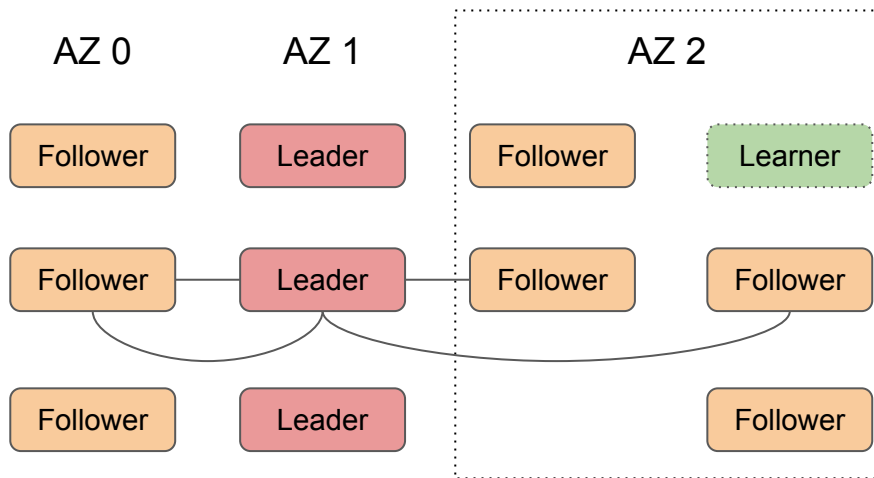
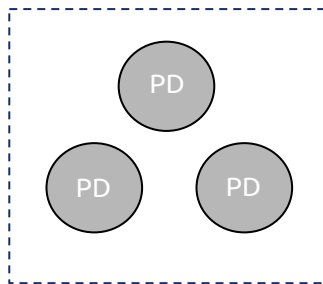
# Configuration change

- Simple changes can lead to unavailability



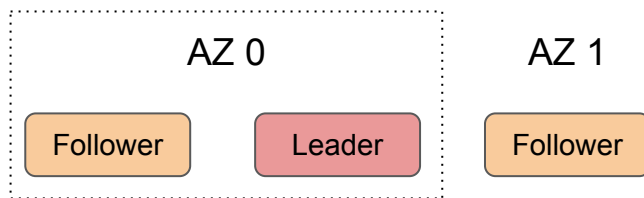
# Configuration change

- Simple changes can lead to unavailability
- Joint Consensus
  - Enter joint state first
    - Both new and old configuration takes effects
  - Complicated



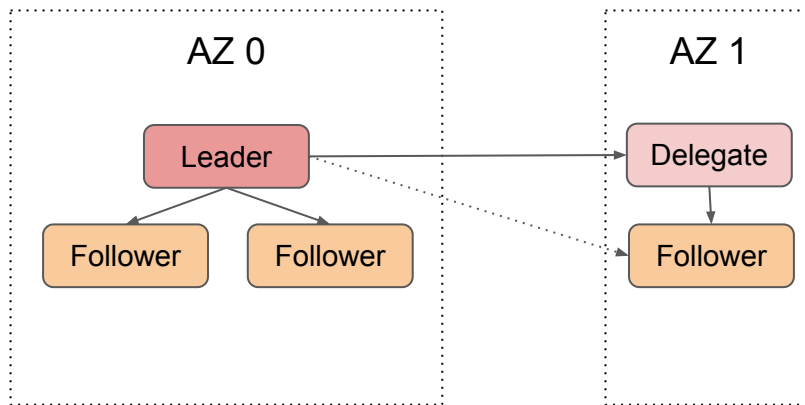
# Commit group

- Majority may not ensure data safety
  - Destroy of one AZ can lost writes
- Replicas are assigned to different groups
- Group is calculated according to AZs
- Logs commit
  - Majority from configuration
  - Replicated to at least 2 groups



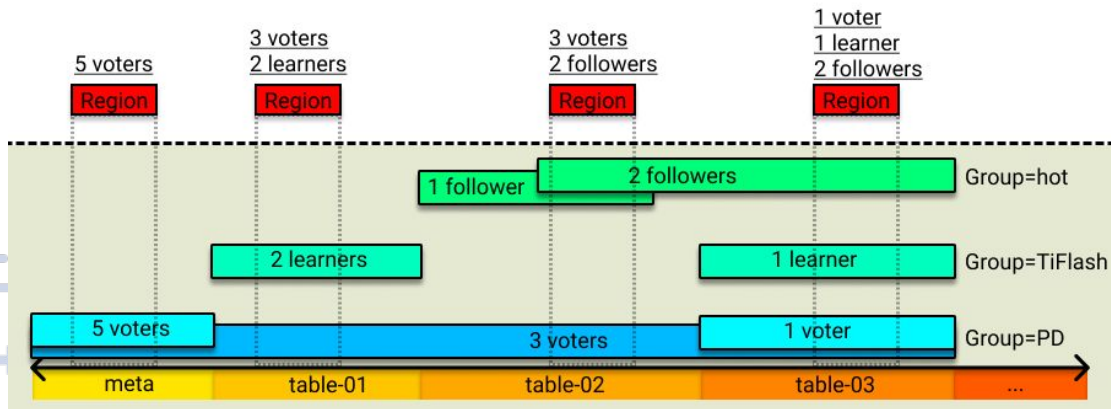
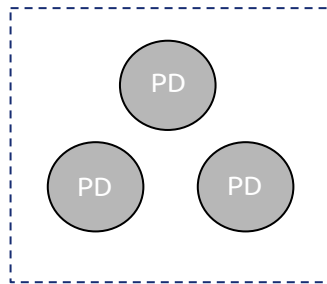
# Commit group

- A delegate is assigned to a group
  - Forwarding logs from leader
- Reduce bandwidth by half



# Placement rules

- Place replicas by rules
  - Constraints on replication numbers
  - Raft roles
  - Geolocation
- Work on ranges



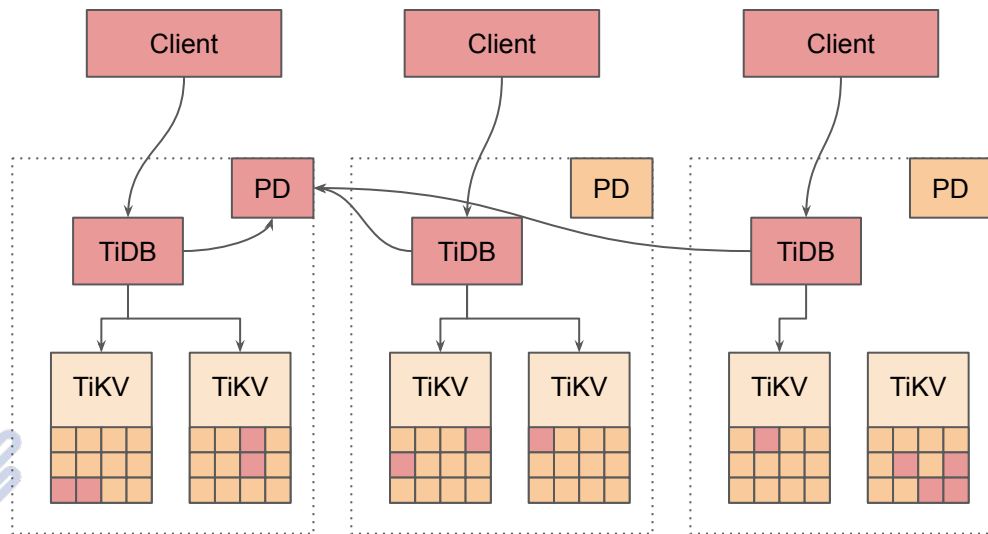
# Part III - Deployment





# Three AZs

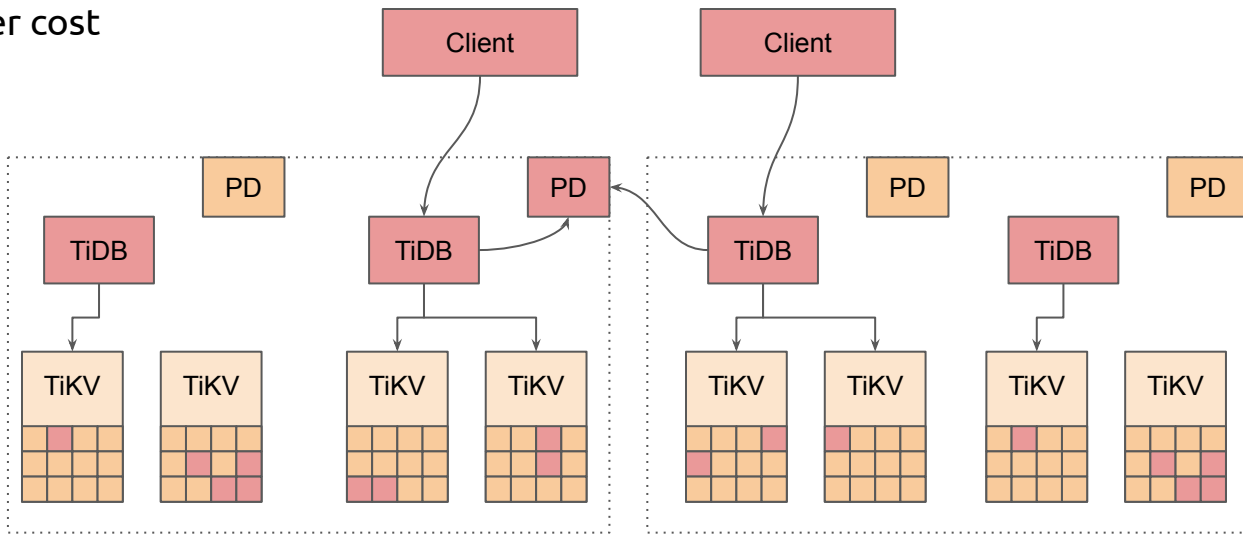
- Replicas are distributed among 3 AZs
- Placement rules in PD
  - controls leaders and replicas
- Leaders are scheduled relative to client
- Read can be optimized by follower read





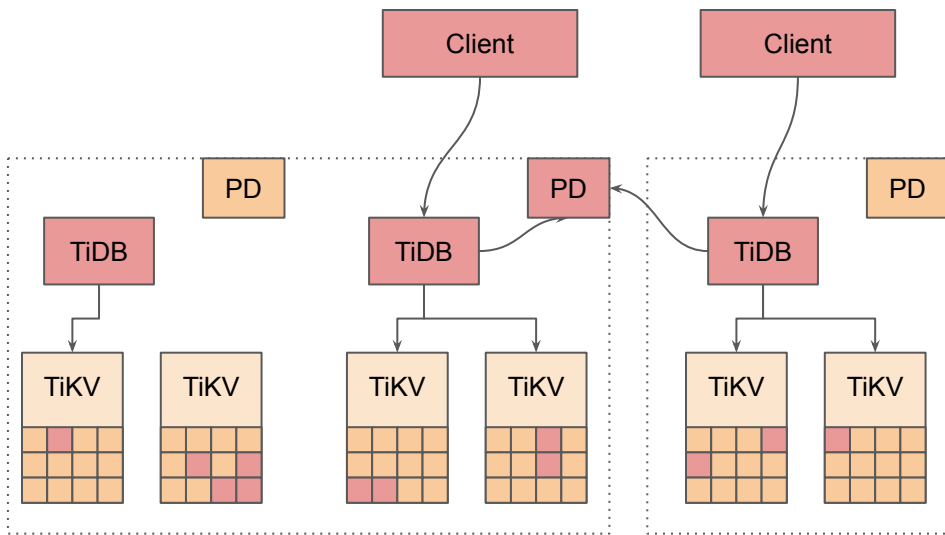
# Two AZs

- Use even number of replicas
- Safety
  - No data lost when either AZ fails
- Availability
  - Unavailable when either AZ fails
- Higher cost



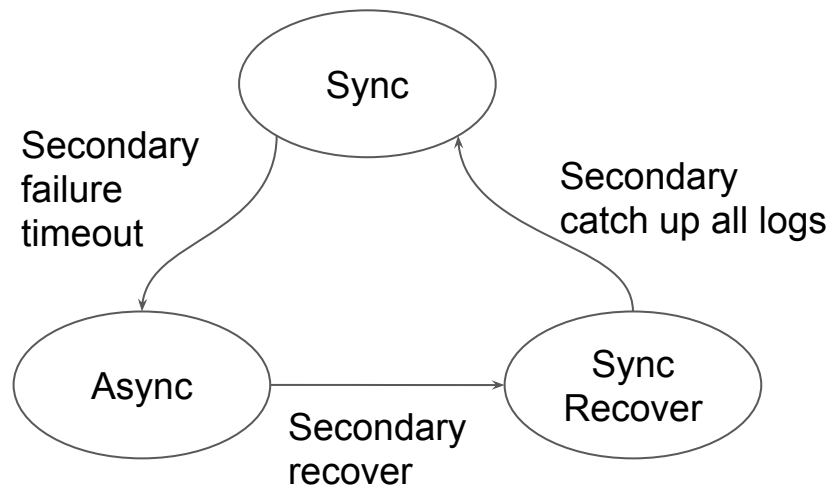
# Two AZs

- Odd number of replicas + group commit
- Safety is guaranteed by group commit
- Primary AZ has more replicas
- Failure of secondary AZ can be recover by removing group commit



# Two AZs

- PD manages replication states
- Only Sync state guarantees safety



# Summary

- Use raft algorithm to ensure atomicity and consistency
- Schedule Leaders to reduce latency using placement rules
- Introduce group commit to ensure safety across even AZs
- Follower read and replication to reduce bandwidth



# Thank You !

Twitter: @pingcap @busyjaylee

<https://www.pingcap.com>

Slack: #everyone in [Slack](#)

