# Dynamic Multi-Raft

## Dongxu Huang

PingCAP

# StuQ

ueue

斯 达 克 学 院

**实践驱动的IT教育**

**斯达克学院（StuQ），** 极客邦旗下实践驱动的IT教育平台。通过线下和线上多种形式的综合学习解决方案 ， 帮助IT从业者和研发团队提升技能水平。

| | | | | |
|---|---|---|---|---|
| 人工智能 | 大数据 | 前端开发 | 后端开发 | 架构设计 |
| 移动开发 | 运维设计 | 产品测试 | 产品经理 | 技术管理 |

**10大职业技术领域课程**
http://www.stuq.org

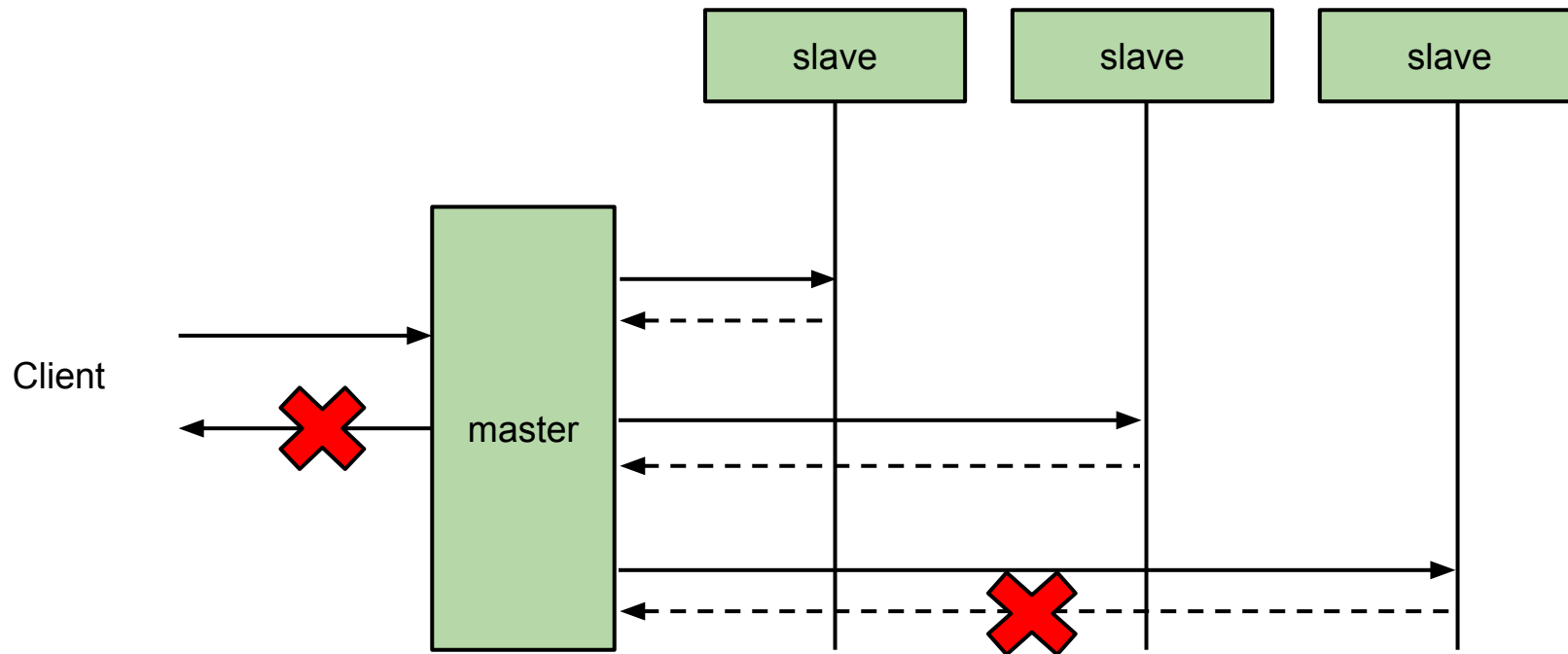SPEAKER
# INTRODUCE

## 黄东旭    CTO & Cofounder, PingCAP

- MSRA / Netease / WandouLabs / PingCAP

- Hacker / Infrastructure software engineer

- Distributed system / Database / PL / …

- Codis / TiDB / TiKV

- Golang / Rust / Python

# Consensus

is the only problem in distributed system...

# Modern HA

- Master-slave is not an option, why?



Client

# MySQL MHA + Semi-Sync?

**MySQL**

**Bug #80395**    **semi-sync: incorrect crash recovery handling**

| | | | |
|---|---|---|---|
| Submitted: | 16 Feb 2016 12:54 | Modified: | 16 Feb 2016 18:17 |
| Reporter: | Matt Lord | Email Updates: | Subscribe |
| Status: | Need Doc Info | Impact on me: | None  Affects Me |
| Category: | MySQL Server: Replication | Severity: | S2 (Serious) |
| Version: | 5.7.11 | OS: | Any |
| Assigned to: | David Moss | | |

| View | Add Comment | Files | Developer | Edit Submission | View Progress Log | Contributions |

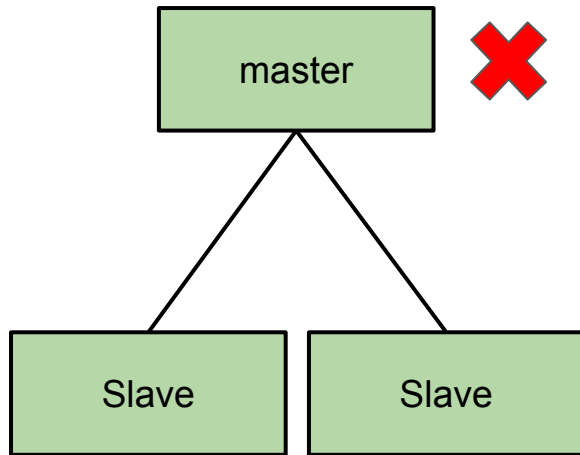**[16 Feb 2016 12:54] Matt Lord**

Description:

When mysqld is killed while an open semi-sync replication transaction is waiting for the master timeout, that prepared *but uncommitted* transaction is NOT property rolled back when the master performs its subsequent automated crash recovery.
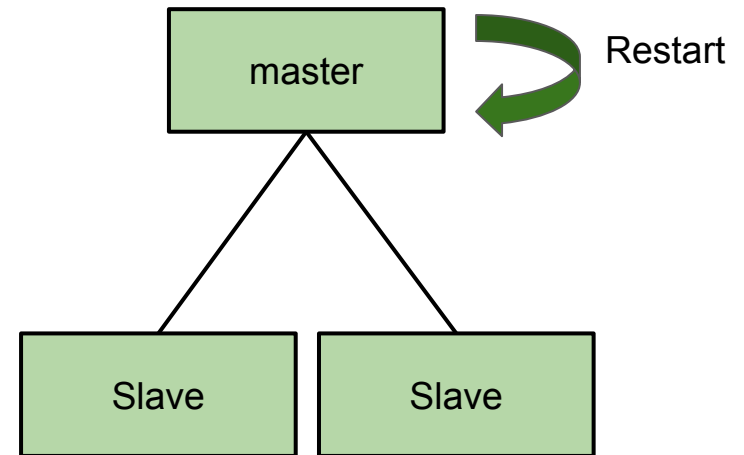
This was verified on OL 7.2 x86_64, using MySQL 5.7.11-community.

How to repeat:

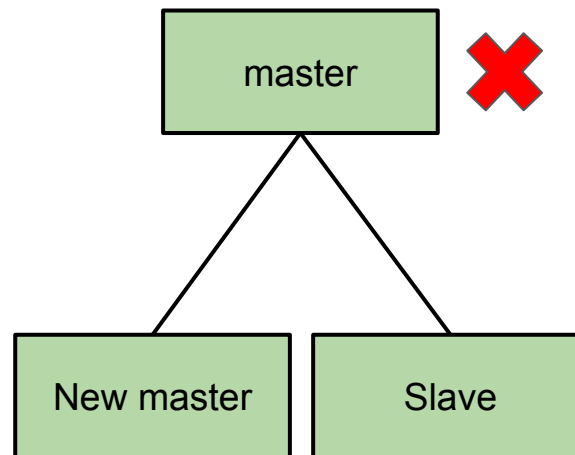(1) Insert A, but master crash when doing semi-sync
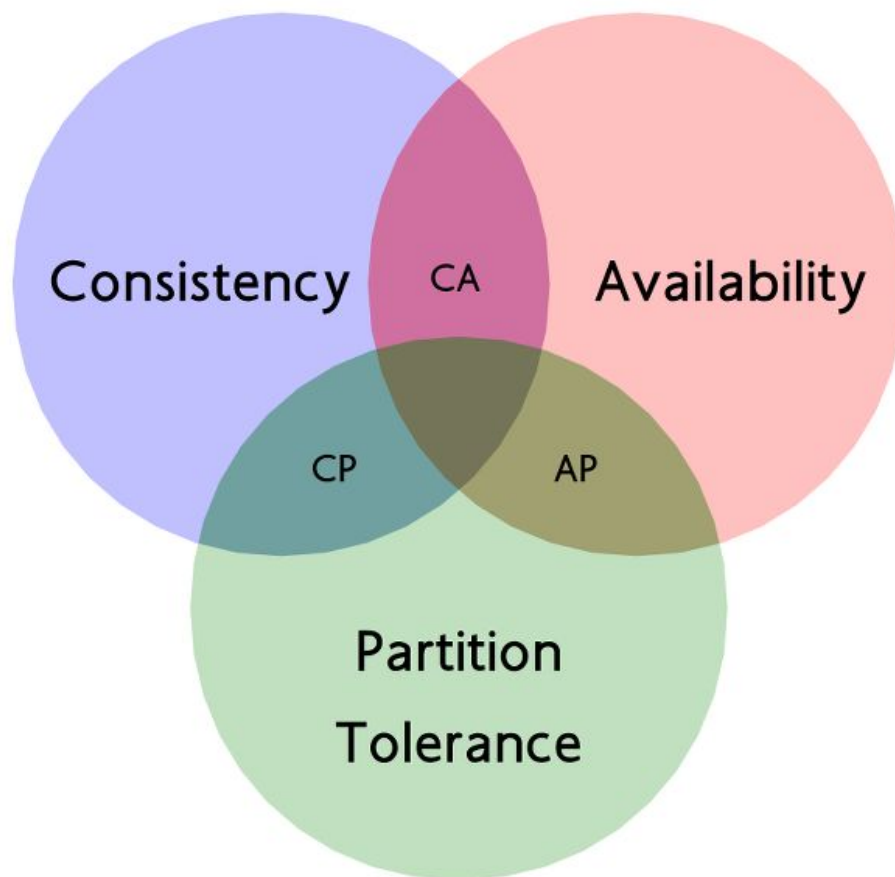


(2) Read A, OK; But degrade to async-replication



Restart

(3) Master crash again...then new master is elected
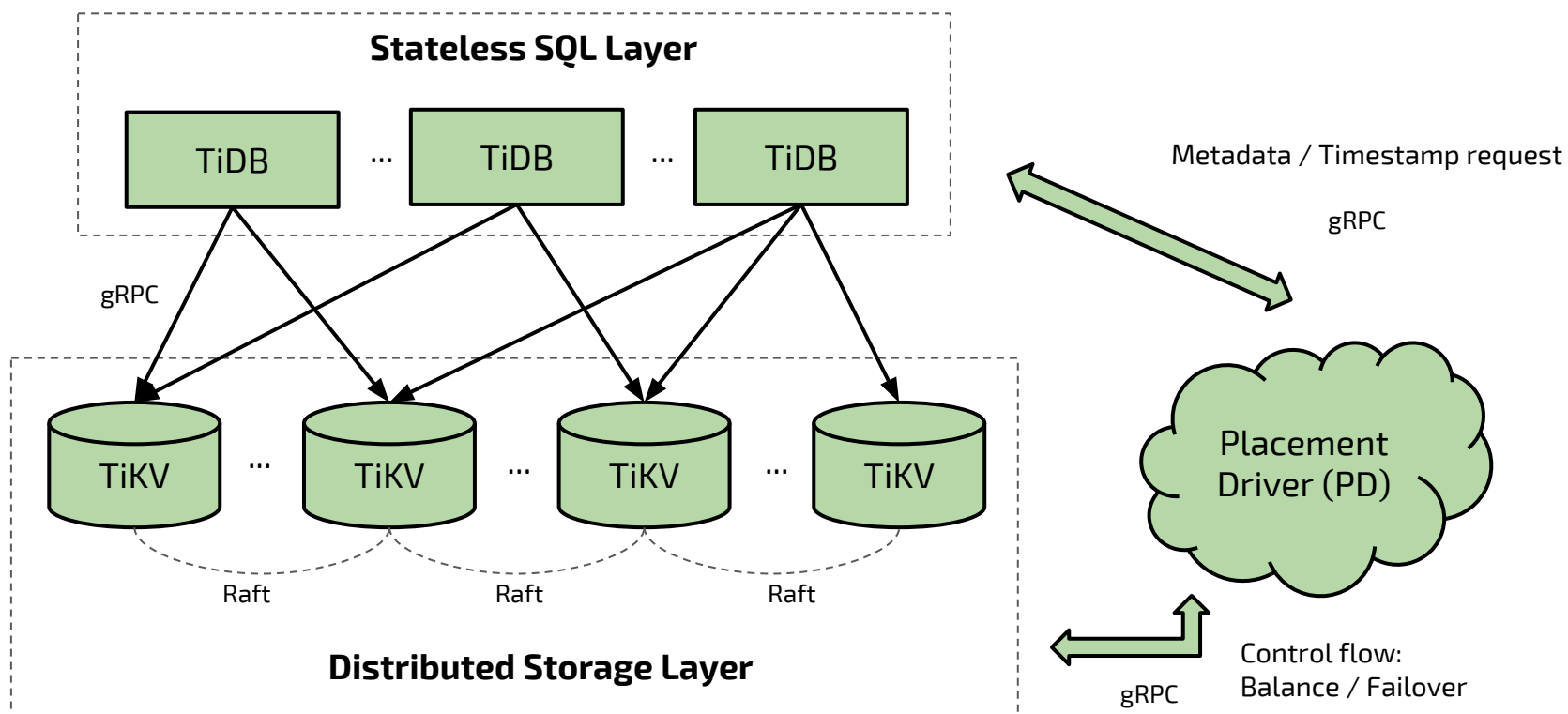


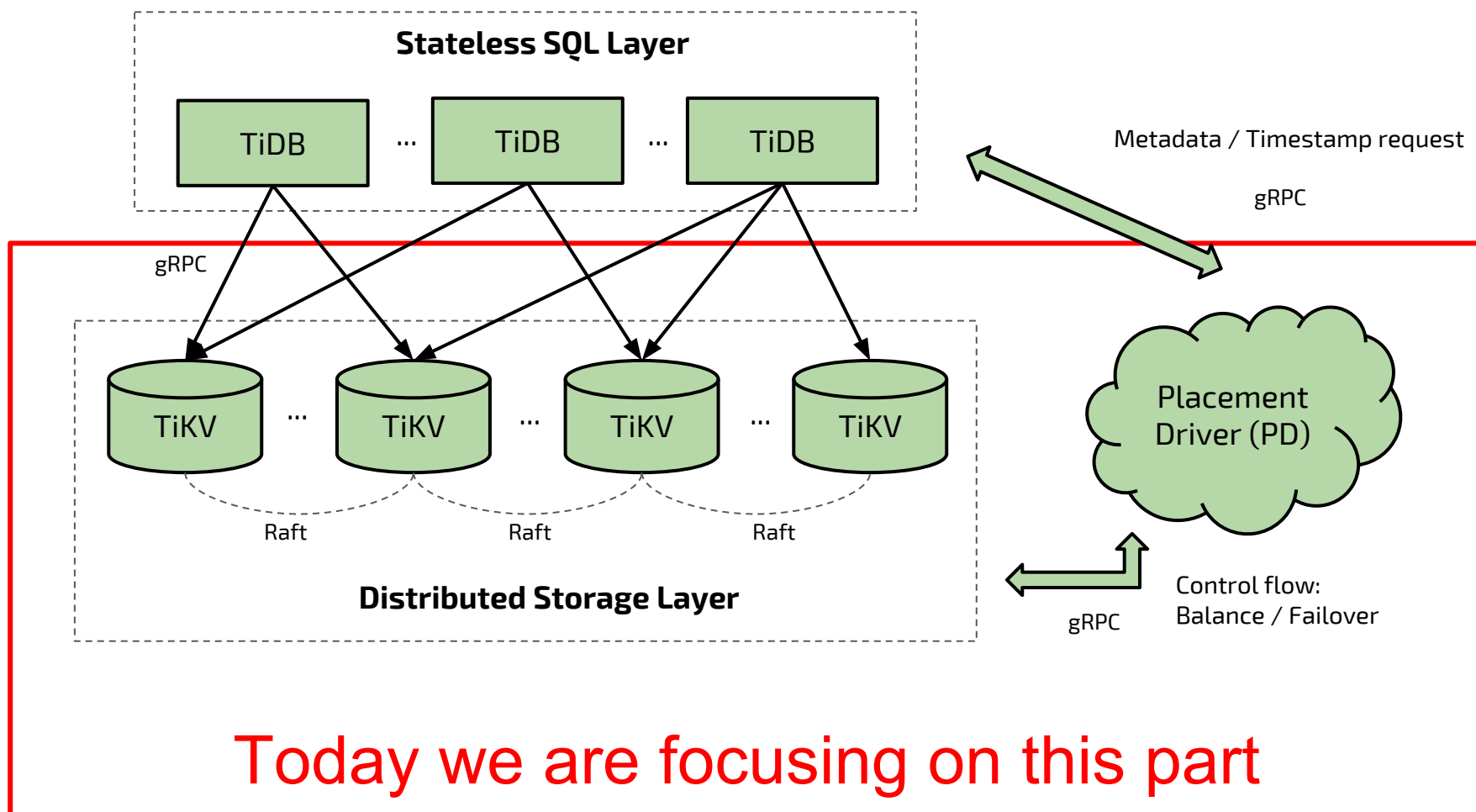Read A, fail, oops...

# TiDB Project (Requirement)

- Strong consistency
- Scalability
- High availability

# TiDB Project Overview

# TiDB Project Overview



Today we are focusing on this part

# Replicated State Machines

- All servers execute same commands in same order
- System makes progress as long as any majority of servers up
- Agreement on shared state (single system image)
- Recovers from server failures autonomously
  - Minority of servers fail: no problem
  - Majority fail: lose availability, retain strong consistency
- IMHO, there are only two RSM implementations:
  - Multi-Paxos / Raft

# The problem in Paxos (Multi Paxos)

*"The dirty little secret of the NSDI community is*

*that at most five people really, truly understand*

*every part of Paxos;-)."*

*—NSDI reviewer*

# Raft saves the day

- **Leader election**

  - Select one of the servers to act as cluster leader

  - Detect crashes, choose new leader

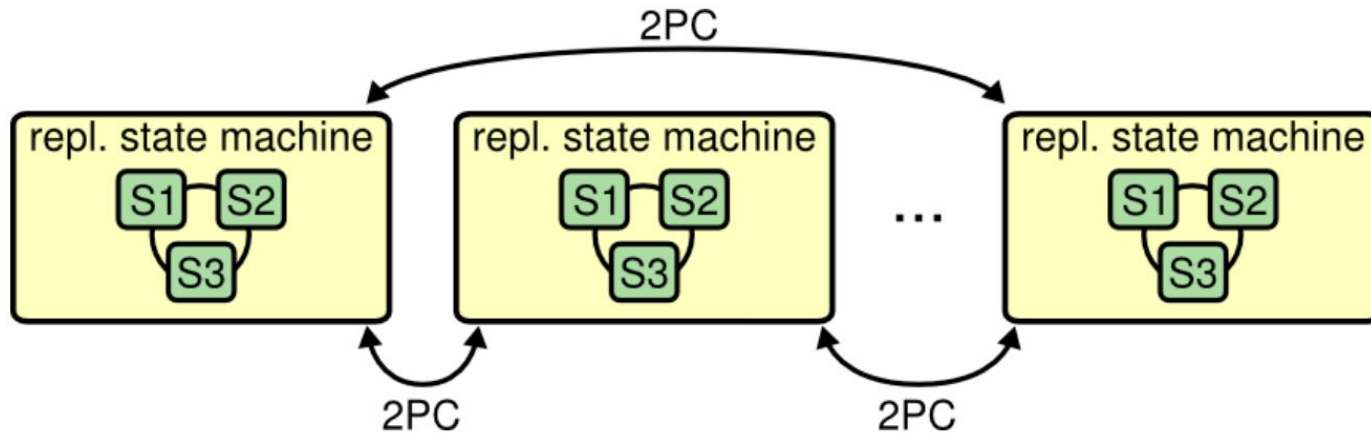- **Log replication**

  - Leader takes commands from clients, appends to its log

  - Leader replicates its log to other servers (overwriting inconsistencies)

- **Safety**

  - Only a server with an up-to-date log can become leader

# Use Raft in database

- Single RSM is **NOT** gonna work.
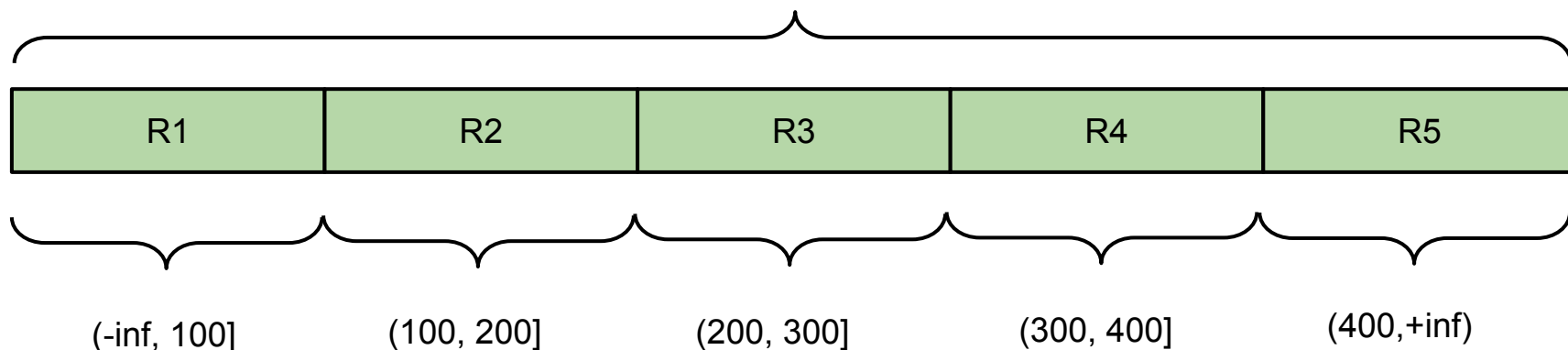- You need 2PC to retain strong consistency across different RSMs.

# Raft in database

- How to shard?
- How to split / merge dynamically?
- How to balance the workload?
- How to improve the throughput?

# Sharding Raft in TiKV

- Split key space into **Regions (normally in byte-order) logically**
- Each region is a raft group
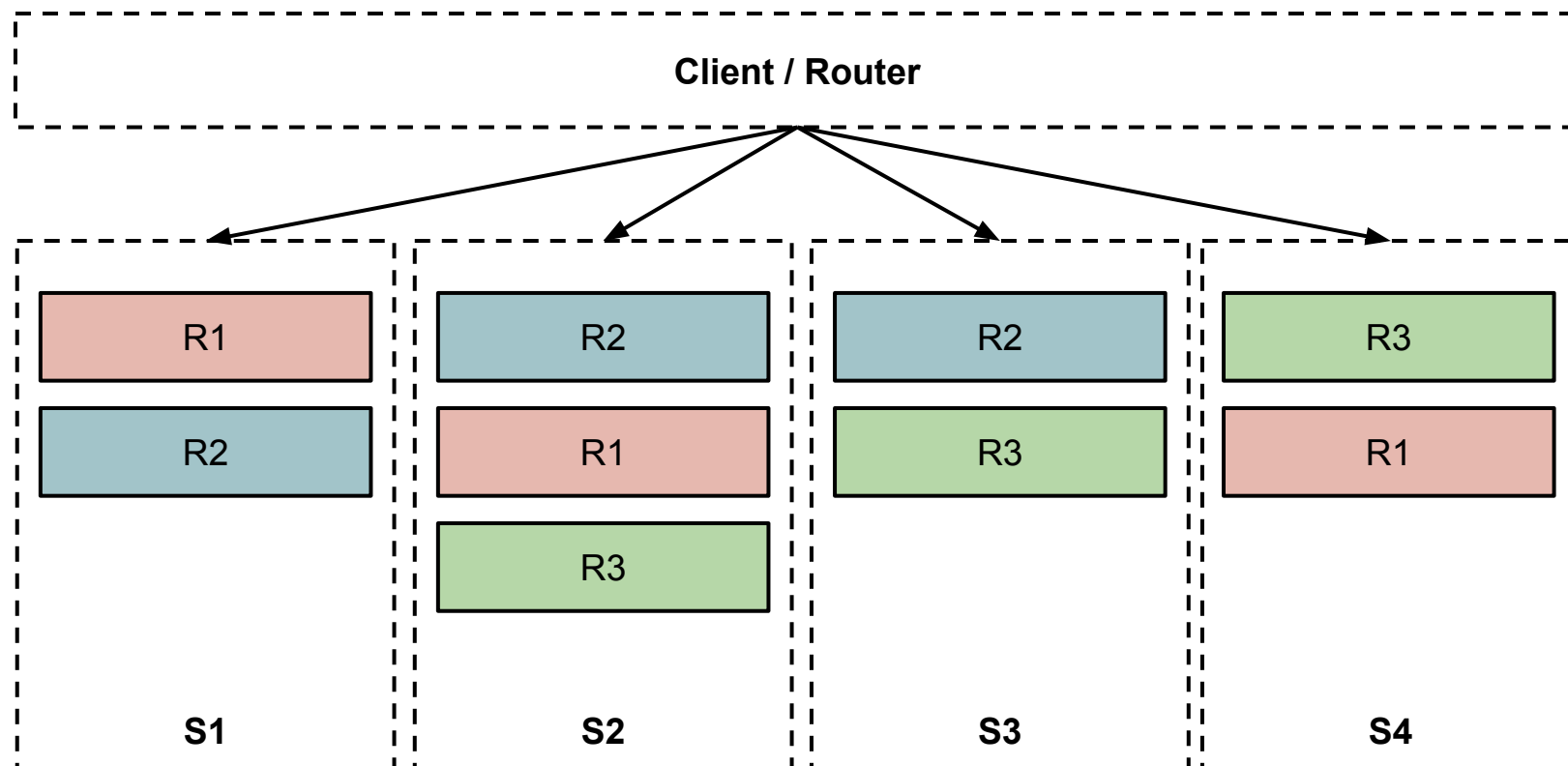  - Default size: 96 ~ 128 MB
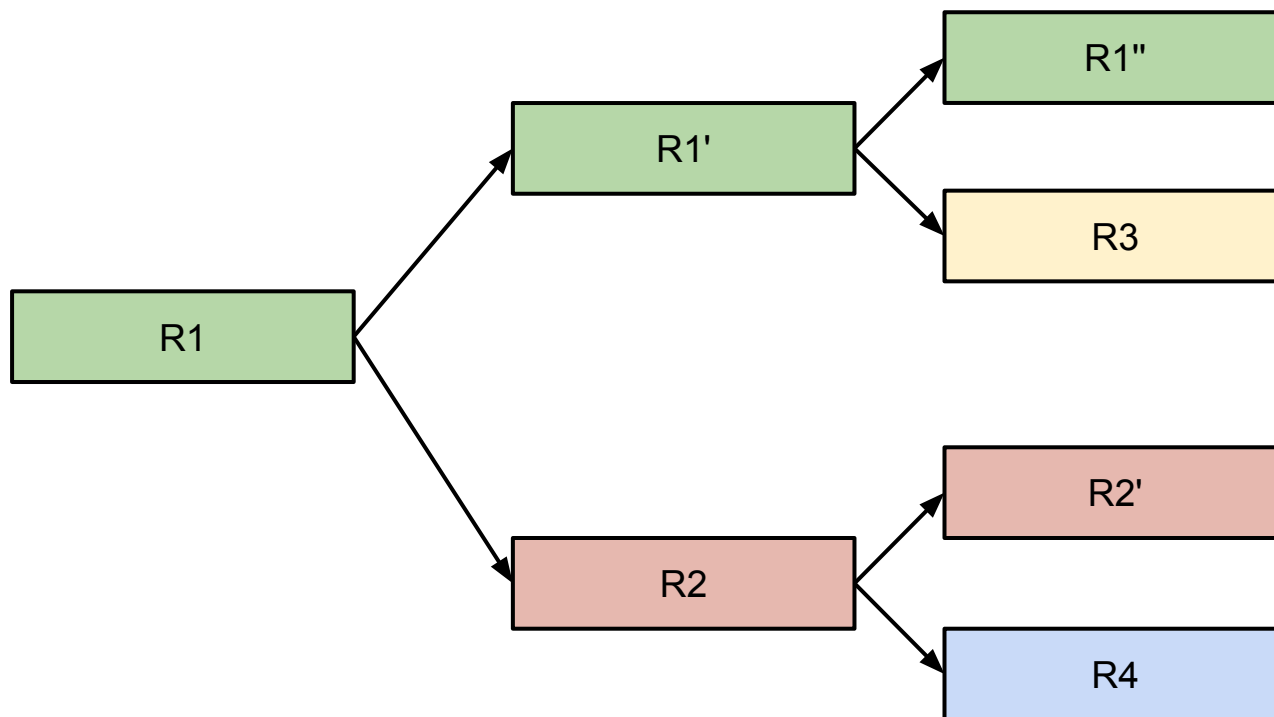  - Why?

Key space (-inf, +inf)

| R1 | R2 | R3 | R4 | R5 |
|---|---|---|---|---|

(-inf, 100]   (100, 200]   (200, 300]   (300, 400]   (400,+inf)

# Meta data storage

- We stores region meta in an in-memory B-Tree (in PD)
  - Sorted by the start key of region
  - We can find the right region which contains specific key in O(log N)
- PD is not '**the source of truth'**, data server is. Why?
  - Split is always happening
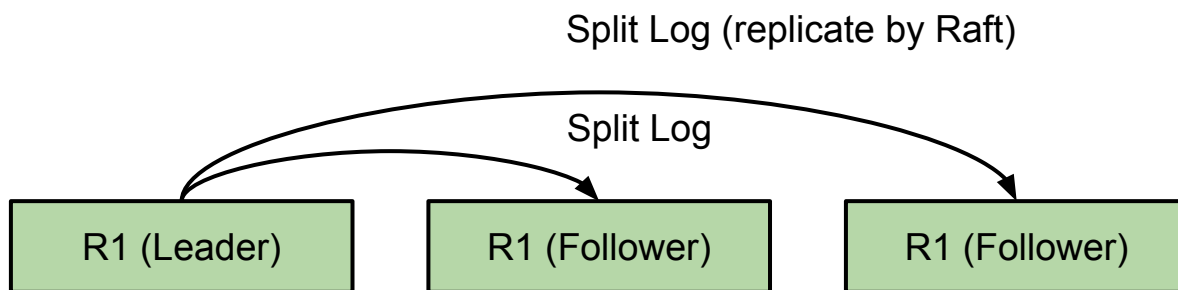  - The metadata stored in PD may be out-of-date
  - Retry is important

# Sharding Raft in TiKV

# Sharding Raft in TiKV

# Dynamic split / merge

Split Log (replicate by Raft)

Split Log

| R1 (Leader) | R1 (Follower) | R1 (Follower) |

# Simple...Huh?

# An abnormal situation...



(1)

R1 (Leader)  —  R1 (Follower)  —  R1 (Follower)

S1  —  S2  —  S3

(2)

R1 (Leader)  —  R1 (Follower)  —  R1 (Follower)

R2 (Leader)  —  R2 (Follower)

S1  —  S2  —  S3

# An abnormal situation...

(3) After N rounds of split or membership changes...



| R1 (Leader) | R1 (Follower) | R1 (Follower) | R1 (Follower) |
| R2 (Leader) | R2 (Follower) | R2 (Follower) | |
| **S6** | **S4** | **S5** | **S3** |

| Rx (Leader) | Rx (Follower) | |
| Ry (Leader) | Ry (Follower) | ... |
| **S1** | **S2** | |

# An abnormal situation...



(4)

| S3 | S4 | S5 | S3 |
|---|---|---|---|
| R1 (Leader) — R1 (Follower) — R1 (Follower) | | | R1 (Follower) |
| R2 (Leader) — R2 (Follower) — R2 (Follower) | | | |

| S1 | S2 |
|---|---|
| Rx (Leader) — Rx (Follower) | |
| Ry (Leader) — Ry (Follower) | ... |

Request votes for R1

# Another abnormal situation

| | | |
|---|---|---|
| R1 (Leader) | R1 (Follower) | R1 (When lease is not outdated) |
| R2 (Leader) | R2 (Follower) | |
| **S1** | **S2** | **S3** |

R1 is [A, C]                                    R1 is [A, D]

**PD**

**???**

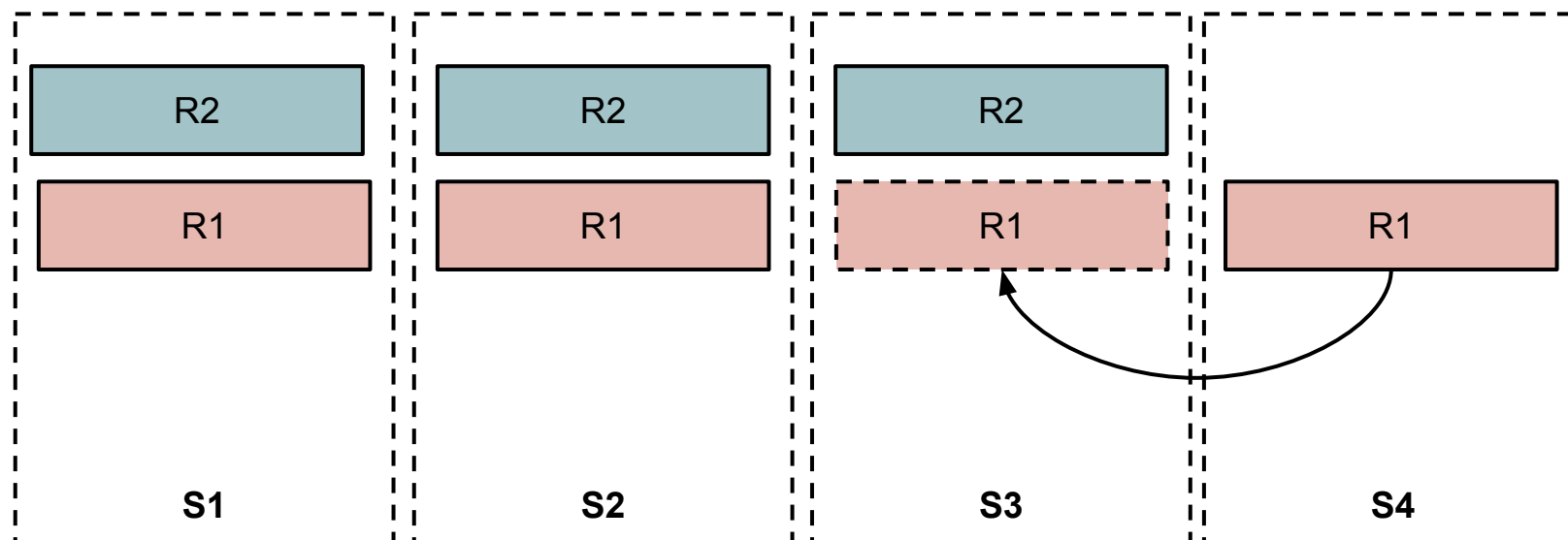# Introduce Region Epoch

- Epoch(Region X) := {ConfVer, SplitVer}
- Every configuration change in Region X will increase the ConfVer
- Every split occurs in Region X will increase the SplitVer
- Let's say Epoch(R1) >= Epoch(R2), if and only if:
  - ConfVer(R1) >= ConfVer(R2) and SplitVer(R1) >= SplitVer(R2)
- Larger epoch always win

# What about merge?

- Make sure all replica for these two regions are in same nodes
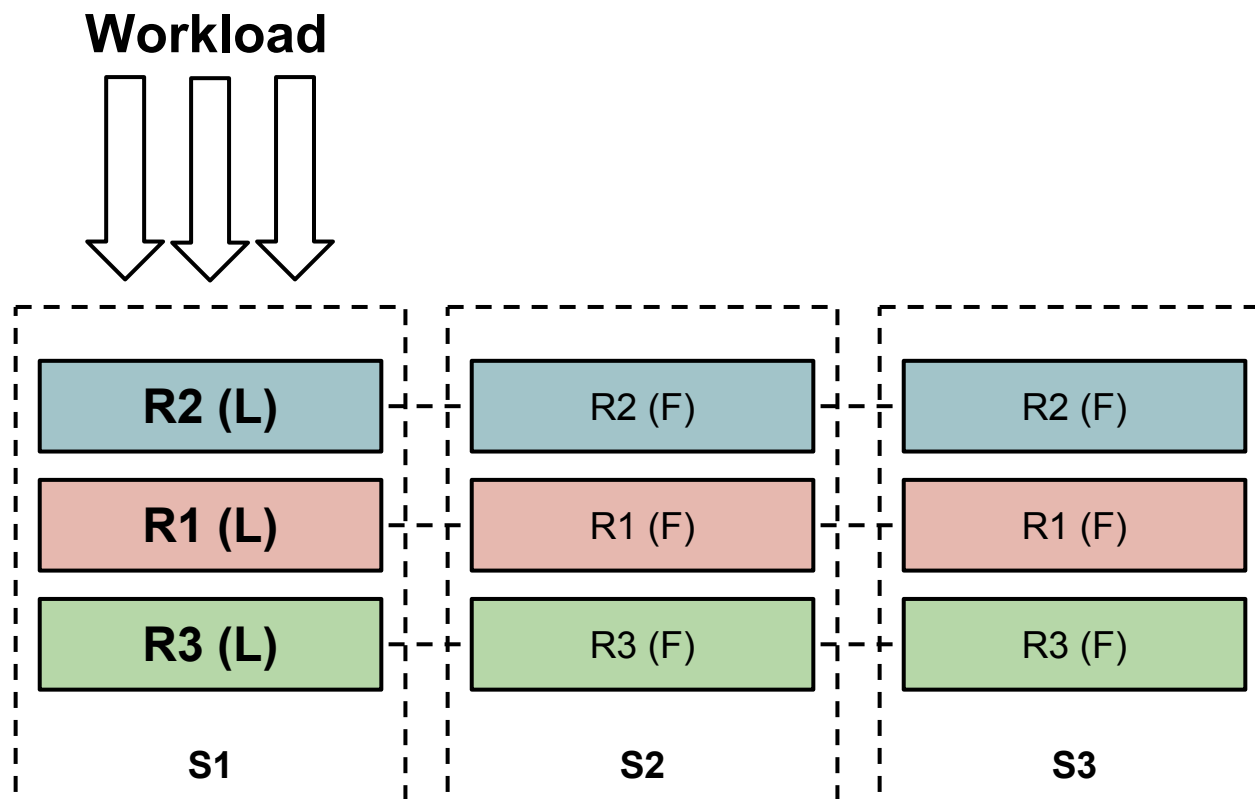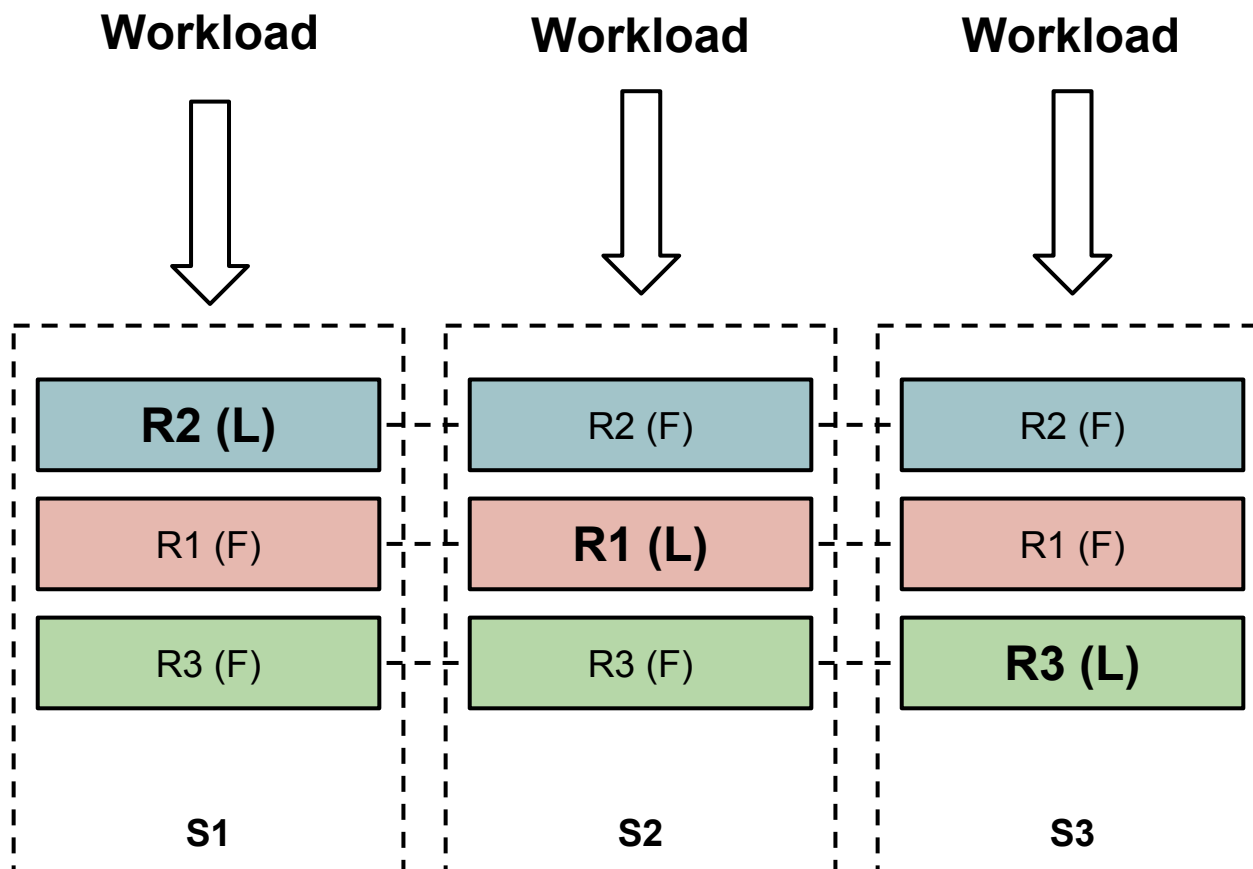- And no more rebalance for these two regions

# Storage

- RSM storage
  - All regions in the same physical node share one RocksDB instance
- Log storage
  - Journal-like storage
  - Share with region storage

# Leadership transfer

- For fast rebalance, since Raft is a randomized algorithm, there is a certain probability that one node has too many leaders.
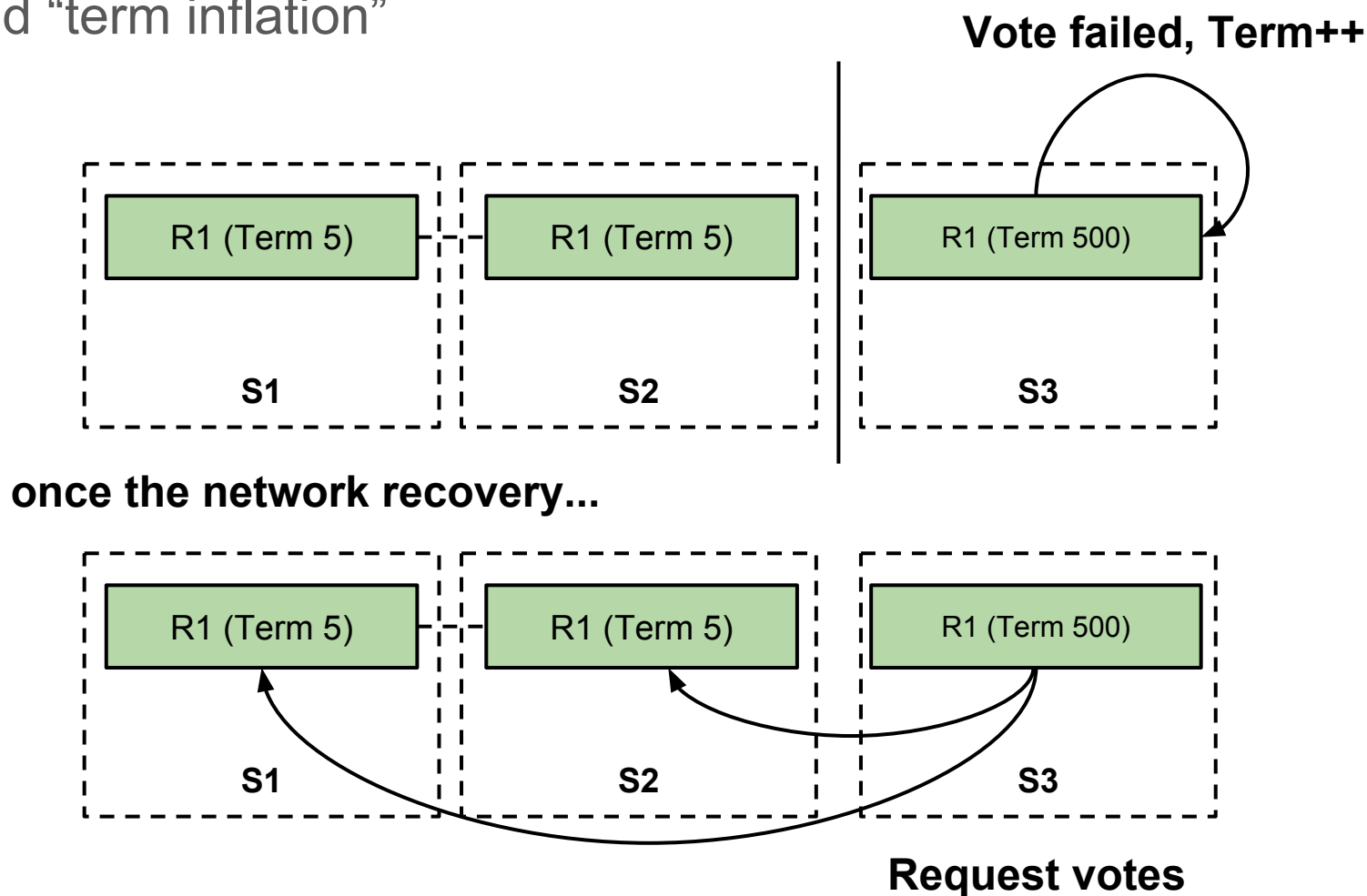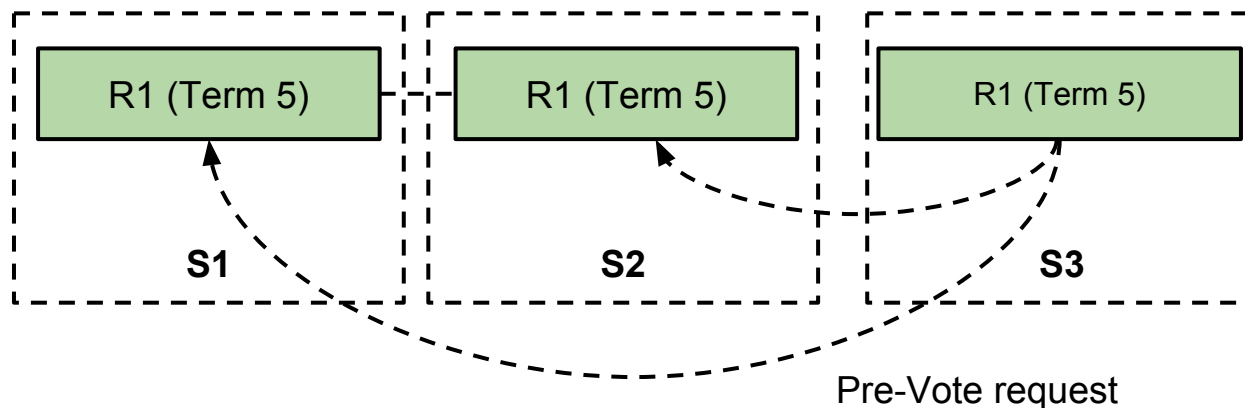
**Workload**

| S1 | S2 | S3 |
|----|----|----|
| **R2 (L)** | R2 (F) | R2 (F) |
| **R1 (L)** | R1 (F) | R1 (F) |
| **R3 (L)** | R3 (F) | R3 (F) |

# Leadership transfer

# Pre-vote algorithm

- Avoid "term inflation"



**Vote failed, Term++**

| R1 (Term 5) | R1 (Term 5) | R1 (Term 500) |
|:---:|:---:|:---:|
| **S1** | **S2** | **S3** |

**once the network recovery...**

| R1 (Term 5) | R1 (Term 5) | R1 (Term 500) |
|:---:|:---:|:---:|
| **S1** | **S2** | **S3** |

**Request votes**

# Pre-vote algorithm



S3 sends Pre-vote request to S1 and S2 to make sure S3's log is up-to-date, when S3 receives responses from a majority of the cluster, S3 will increase its term and start a normal election

# How to test

- Testing in distributed system is really hard
- Test-Driven Development
- Test cases from community
  - Lots of tests in MySQL drivers/connectors
  - Lots of ORMs
  - Lots of applications (Record---replay)
- Fault injection
  - Hardware: disk error, network card, cpu, clock
  - Software: file system, network and protocol
- Simulate everything：Network
- Distribute testing
  - Jepsen
  - Namazu

# Benchmark

- 46 Physical nodes
- 460 TiKV instances  (1 tikv instance for 1 HDD)
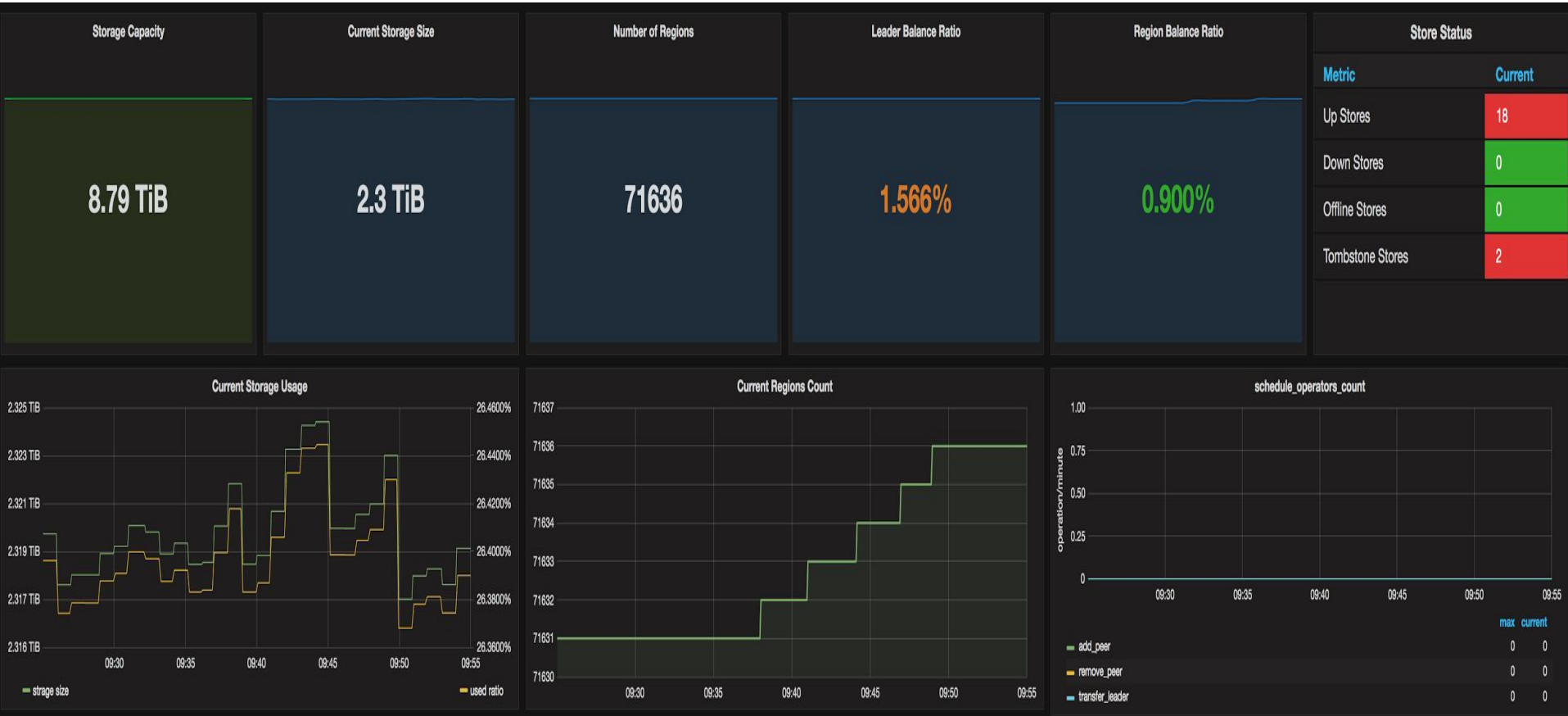- TiKV Raw API Put (Raft)

Put(key, value)

key size: 21 bytes

value size: random (1~100 bytes)

# Benchmark

# Benchmark

# THANKS!

让创新技术推动社会进步
HELP TO BUILD A BETTER SOCIETY WITH
INNOVATIVE TECHNOLOGIES

# Geekbang>.

## 极 客 邦 科 技

**InfoQ** ueue

专注中高端技术人员的技术媒体

**EGO** EXTRA GEEKS' ORGANIZATION NETWORKS

高端技术人员学习型社交平台

**StuQ** ueue
斯 达 克 学 院

实践驱动的 IT 教育平台