# Elastic Scheduling with TiKV

*Song Gao, PingCAP*
*Yutong Liang, PingCAP*

# Speaker

**Yutong Liang**
Engineer at PingCAP

Database engineer
Technical lead of TiKV SIG Scheduling

Github: @rleungx



**Song Gao**
Engineer at PingCAP

Database engineer
Maintainer of Chaos Mesh®
Committer of TiKV SIG scheduling

Github: @Yisaer

# Agenda

- Introduction to TiKV
- Elastic Scheduling background
- Implementation in TiKV
- Future work
- Q&A

# Introduction

What is TiKV?

**TiKV** is an open source **distributed transactional** key-value database.

**CLOUD NATIVE**
**COMPUTING FOUNDATION**

CNCF Graduated

8.2K

GitHub Stars

264

Contributors

# Introduction
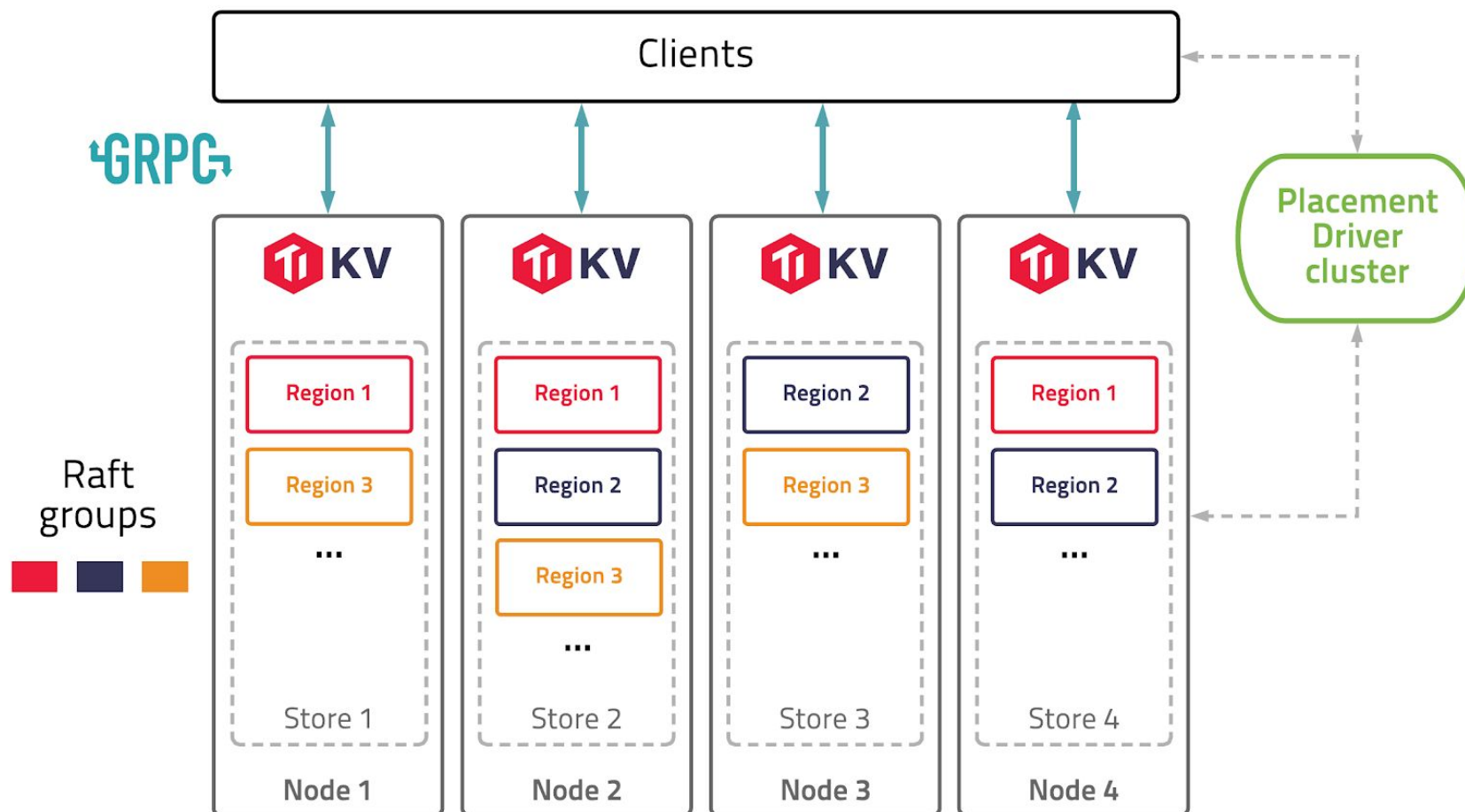
## TiKV Architecture

# Introduction

- Add a new Node D

**Add new node**

| Node A | Node B | Node C | Node D |
|--------|--------|--------|--------|
| **Region 1** | Region 1 | Region 1 | |
| Region 2 | **Region 2** | Region 2 | |
| Region 3 | Region 3 | **Region 3** | |

6

# Introduction

- Add a replica of Region 1 in Node D

**Add replica**

# Introduction
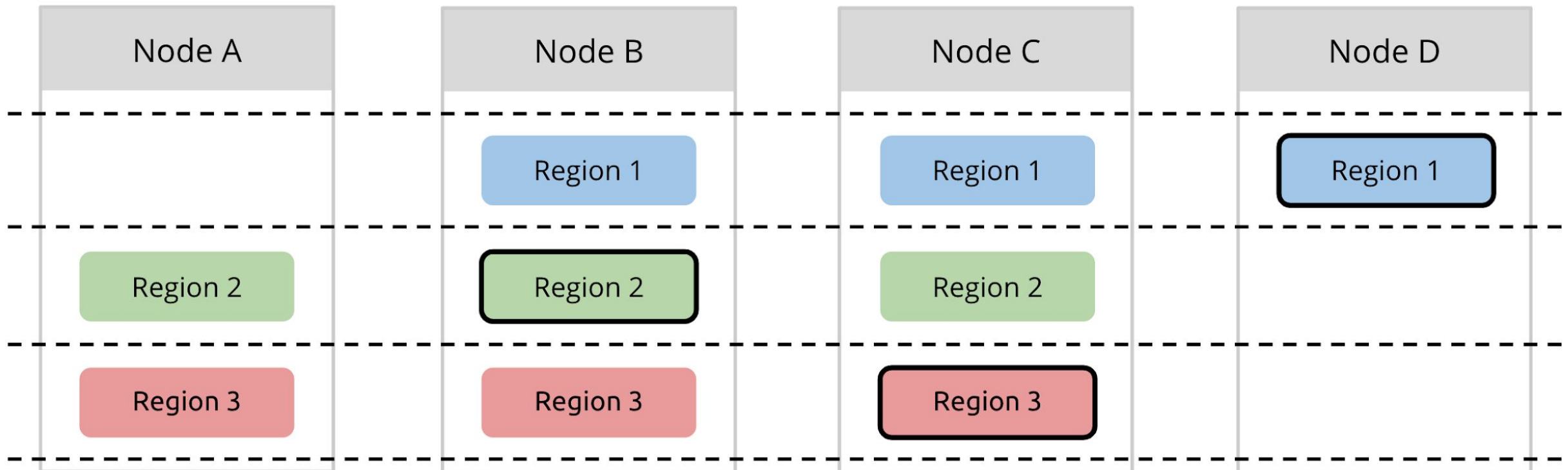
- Transfer leader of Region 1 from Node A to Node D

# Introduction

- Remove the original replica of Region 1 from Node A

**Remove replica**

# Elastic scheduling

# What is Elastic Scheduling?

# Elastic scheduling

Auto scaling by workloads



*3     The load increases     *5     The load decreases     *3

# Why Elastic Scheduling?

# Elastic scheduling

The traffic is unexpected

# Elastic scheduling

Some resources are wasted



**We need to pay the cost for the traffic hour.**

# Elastic scheduling

The cloud infra becomes mature.

# Implementation in TiKV

# Implementation

Elastic scheduling architecture

Scaling Plan

Placement driver

calculation

Operator

Prometheus

Scaling

TiKV cluster

Scraping
Metrics

# Implementation

## Operator side

User actions
New object, reconfigure

Scaling Plan

Placement driver

Kubernetes events
Current state of cluster

# Implementation

## Scheduling side



Operator

query plan →

← add 1 TiKV

Placement Driver

fetch metrics →

Prometheus

# Implementation

## Scheduling side

# Implementation

How does PD recognize the hot region?

- PD will maintain caches to record the top N Region write/read  flow of each store. The hot Region must meet two conditions:

    - continue to hit the cache
    - write/read flow no less than the minimum threshold

# Implementation

## For other schedulers

```
Select Target Node
```

```
Node Filter
(filter the node with label
specialUse:hotRegion)
```

```
Final Target Node
```

```
...
"node":{
    "id":1,
    "address":"host:port",
    "labels":[{
            "key":"specialUse",
            "value":"hotRegion"
        }],
    ...
}
...
```

22

# Demo

## The API Overview

```
...
spec:
  cluster:
    name: auto-scaling-demo
    namespace: default
  tikv:
    maxReplicas: 4
    metrics:
      - type: "Resource"
        resource:
          name: "cpu"
          target:
            type: "Utilization"
              averageUtilization: 80
```

# Demo

## Initial State

- 3 TiKV

- sysbench: oltp_read_only

# Demo

## Add 2 TiKV

# Demo

## Transferring hot regions
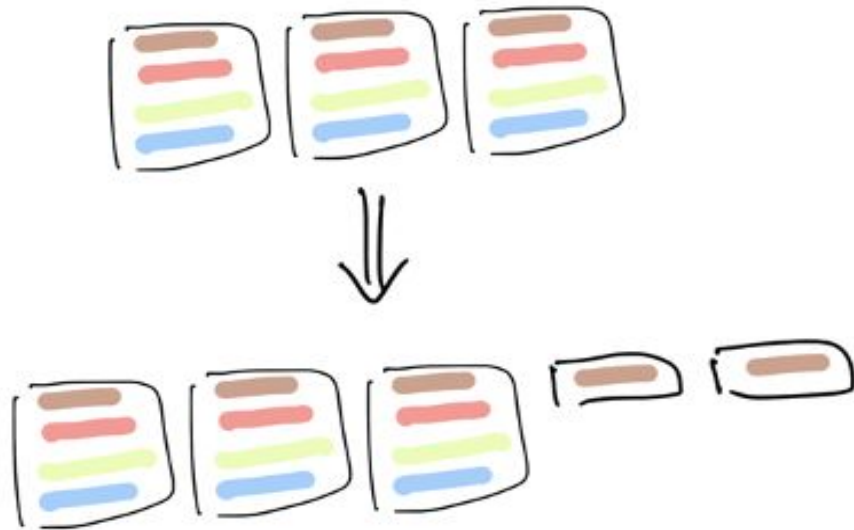
# Future work

# Future work



- **Replication by workloads**

Changing the replication for some regions according to different workloads.

# Future work

- **Separate hot and cold data**

Using cheaper storage media to store the cold data.

# Join us

- GitHub: https://github.com/tikv/tikv
- Website: https://tikv.org/
- Twitter: @tikvproject
- Slack: #sig-scheduling in Slack