

TiDB on Kubernetes

Huang Dongxu

CTO, PingCAP

极客邦企业培训与咨询

「帮助企业和技术人成长」

10 余年
经验技术专家

200+
国内外一线技术
专家团队

800+
企业研发团队
的选择

10000+
学员参与学习
交流

助力企业提升技术竞争壁垒，让技术驱动业务发展



扫码了解更多官方咨询

About Me

- Huang Dongxu (黄东旭)
- CTO, Co-founder of PingCAP
- Infrastructure Engineer / Open-source advocator
- Co-author of Codis / TiDB / TiKV
- MSRA => Netease => WandouLabs => PingCAP
- h@pingcap.com

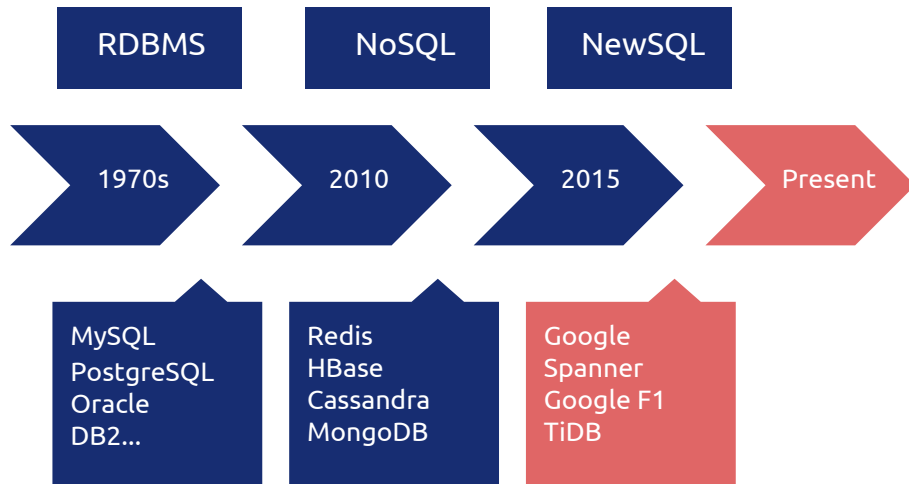
Agenda

- TiDB introduction
 - TiDB architecture
 - TiDB ecosystem
- Why combine TiDB & Kubernetes
 - Cloud vendor agnostic
 - Automation
- How we make it possible
 - TiDB Operator architecture & features
 - How we manage state
 - How we schedule stateful app

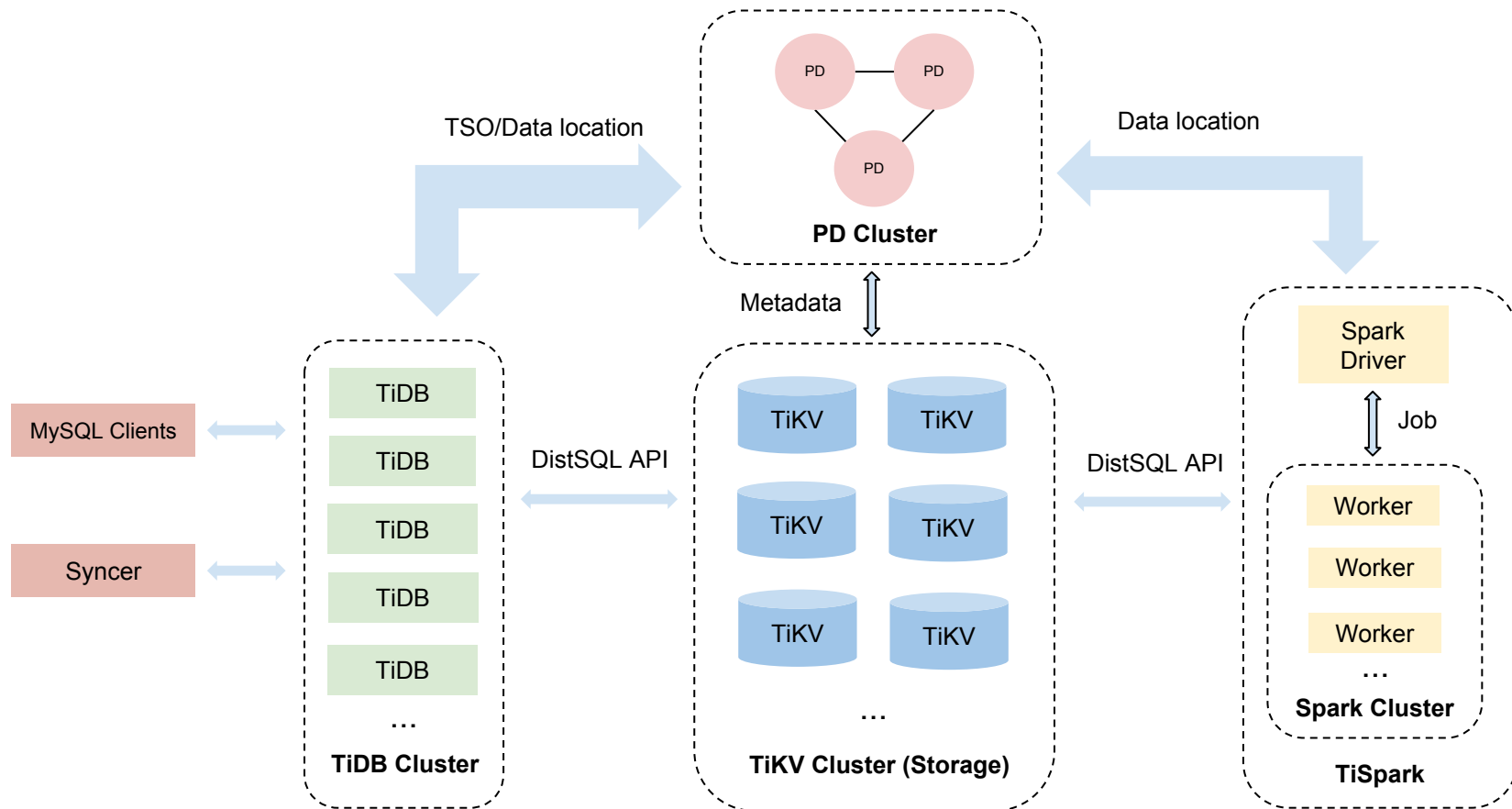
Part I - Intro to TiDB

Why we want to build a NewSQL Database

- From the beginning
- What's wrong with the existing DBs?
 - RDBMS
 - NoSQL & Middleware
- NewSQL: F1 & Spanner

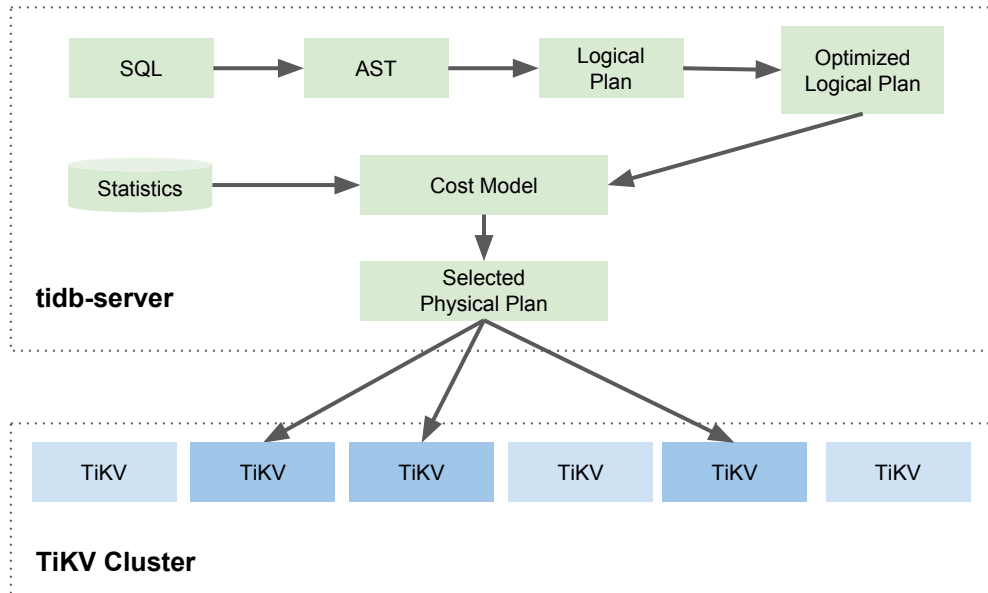


TiDB architecture



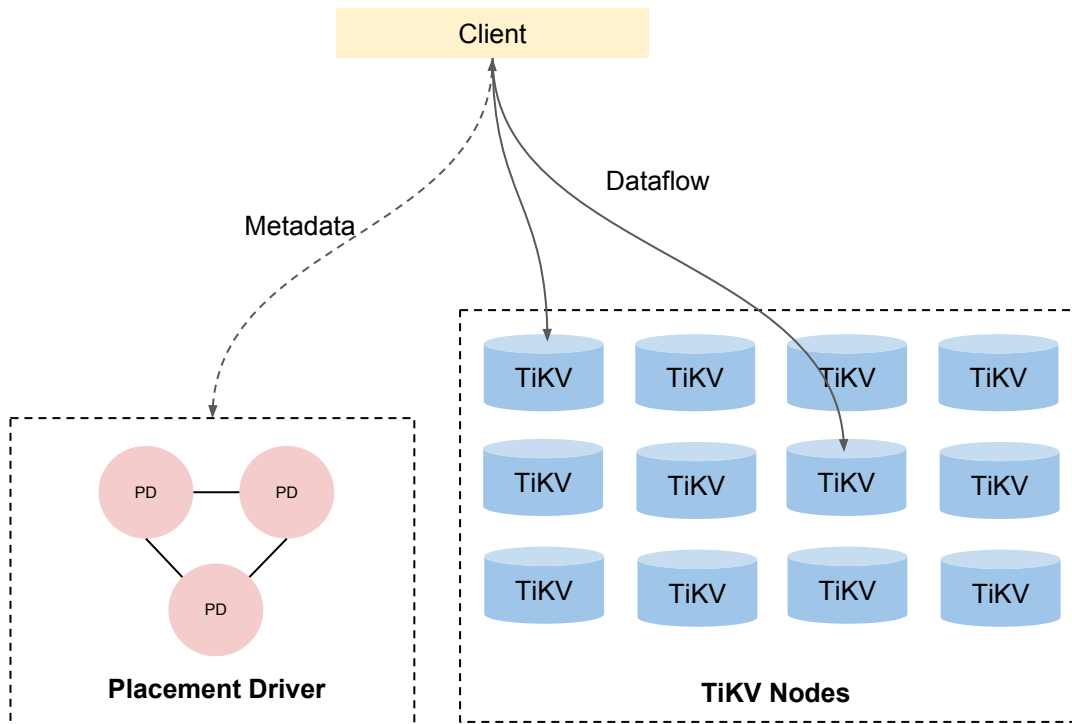
TiDB: Computing

- Stateless SQL layer
 - Client can connect to any existing tidb-server instance
 - TiDB ***will not*** re-shuffle the data across different tidb-servers
- Full-featured SQL Layer
 - Speak MySQL wire protocol
 - Why not reusing MySQL?
 - Homemade parser & lexer
 - RBO & CBO
 - Secondary index support
 - DML & DDL

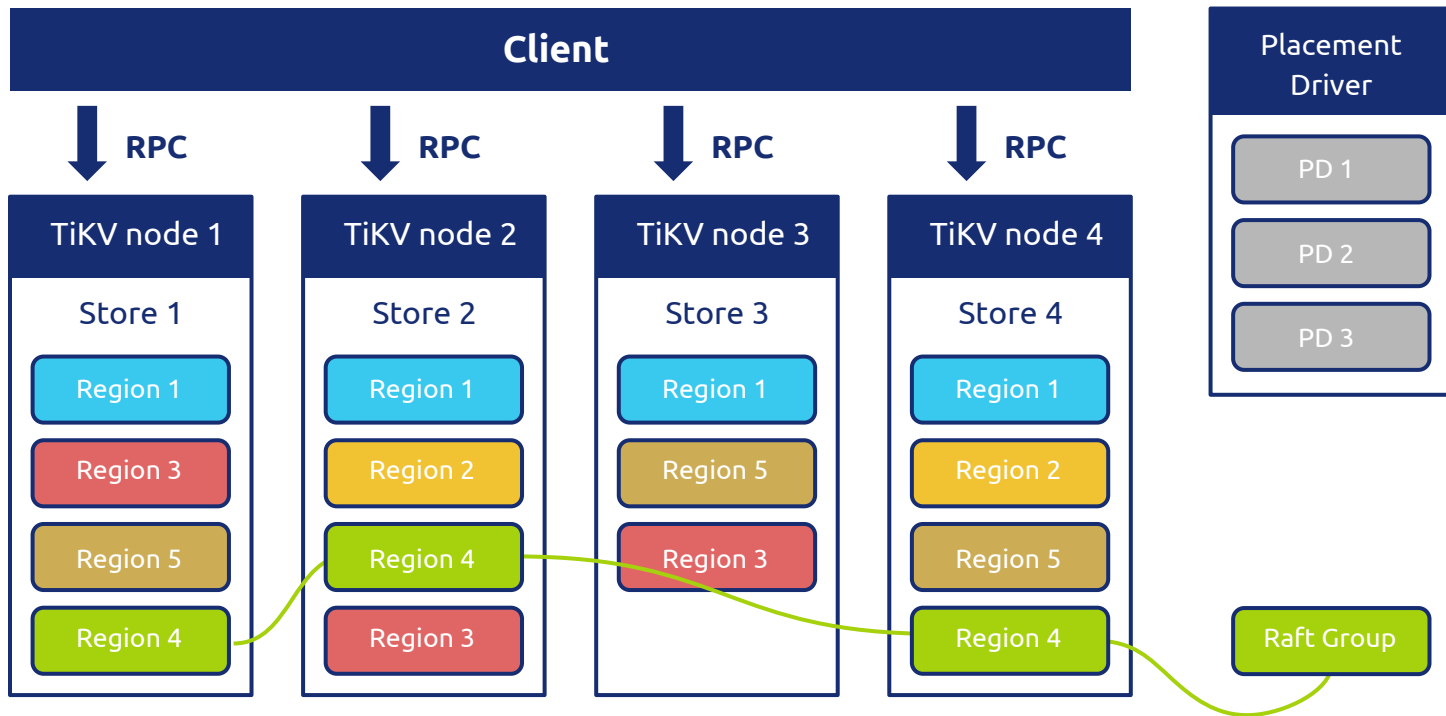


TiKV: The Storage

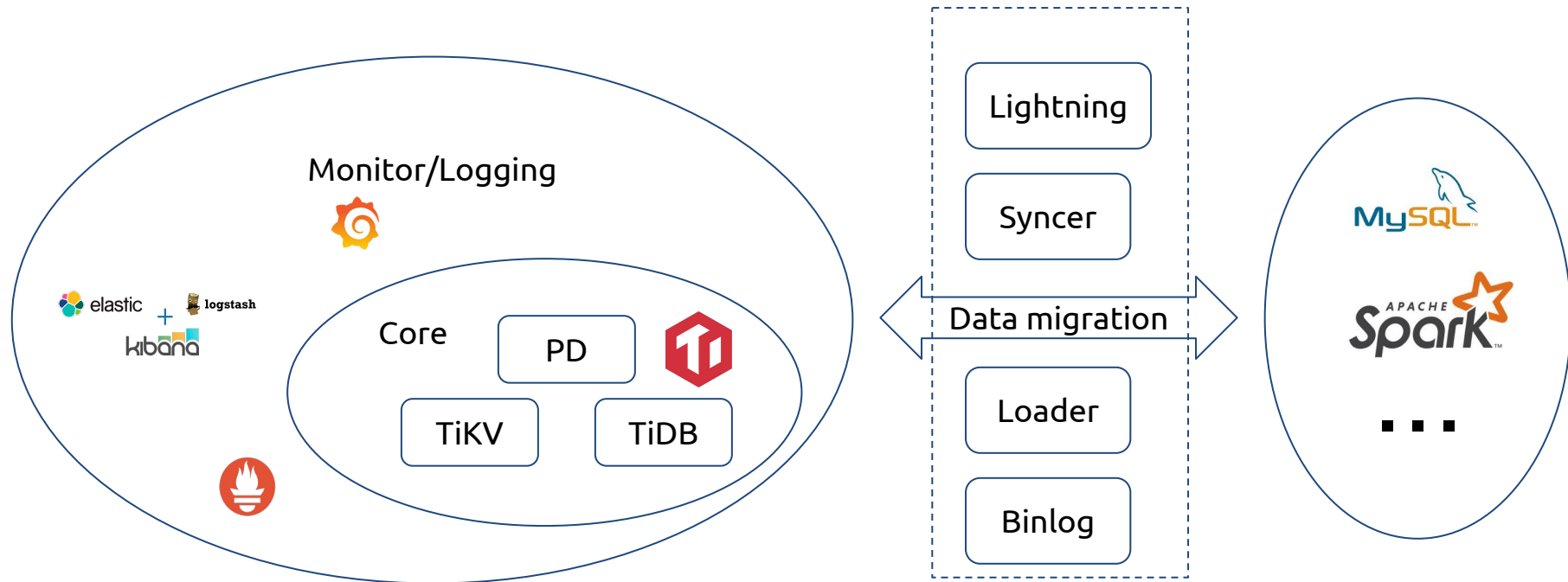
- The **storage layer** for TiDB
- Distributed **Key-Value** store
 - **Support ACID Transactions**
 - Replicate logs by **Raft**
 - **Range** partitioning
 - Split / merge **dynamically**
 - Support coprocessor for **SQL operators pushdown**



TiKV: The Storage



TiDB ecosystem



Part II - Why TiDB on Kubernetes

Cloud-Native applications

- Microservice architecture
- Easy deployment on any cloud
- Elastic scaling
- Highly available
- Automatic operation



Google Cloud



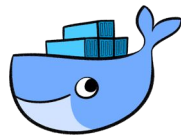
Microsoft Azure

Kubernetes: standard platform

De-facto Container Orchestration System (Google Sponsored)

Distributed, cloud provider agnostic OS

- CPU, Memory, Storage and other Devices across all nodes
- Container \Leftrightarrow Process
- Docker image \Leftrightarrow Executable artifacts
- Deployment, StatefulSet \Leftrightarrow Systemd/Supervisor ...
- Helm / Charts \Leftrightarrow apt yum / deb rpm

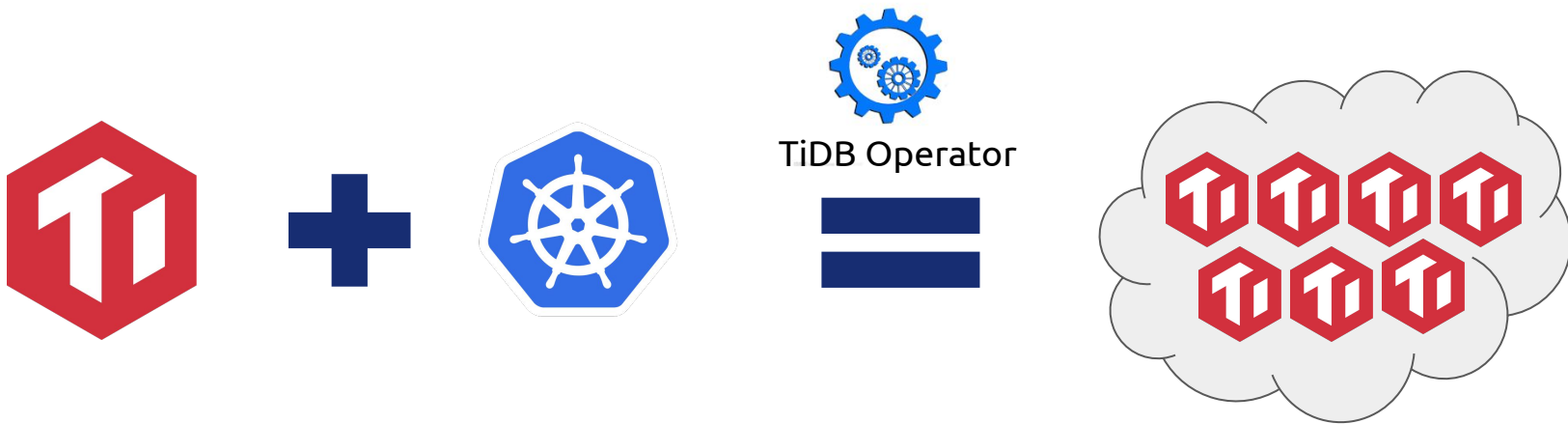


Kubernetes: powerful extensibility

- Standard interface: CNI, CRI, CSI
- Scheduler: scheduler extender
- Controller: CRD
- API Server: Aggregated API Server
- Kubelet: virtual kubelet
- Cloud Provider: LoadBalancer, PersistentVolume
- ...

Part III - How we make it possible

TiDB Operator

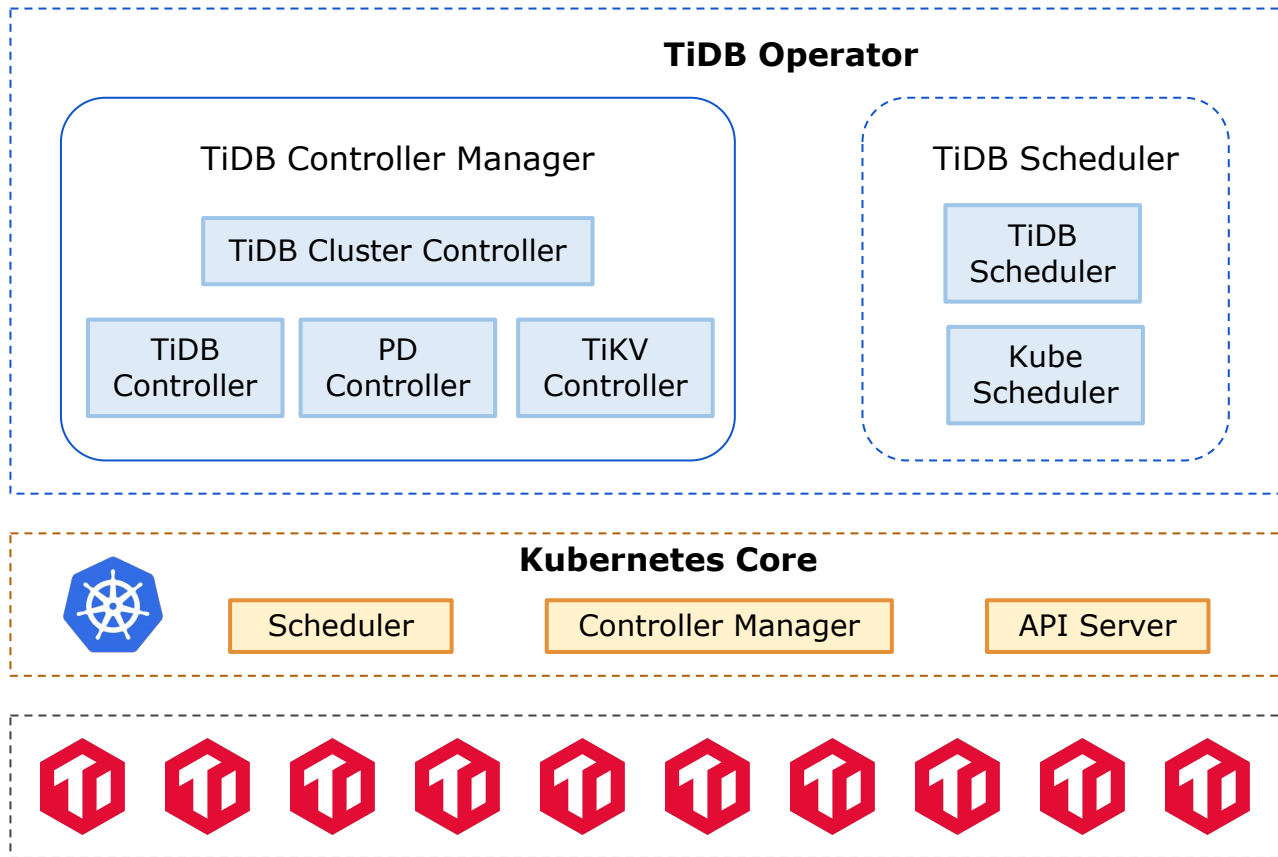


<https://github.com/pingcap/tidb-operator>

Features

- Manage multiple TiDB clusters
- Safely scale the TiDB cluster
- Easily installed with Helm charts
- Network/Local PV support
- Automatically monitoring the TiDB cluster
- Seamlessly perform rolling updates to the TiDB cluster
- Automatic failover
- TiDB related tools integration

Architecture



How we manage state

Kubernetes builtin controllers

Deployment:

- Start ✓
- Scale ✓
- Upgrade ✓
- Failover ✓

StatefulSet:

- Start ✓
- Scale ✓
- Upgrade ✓
- Failover ✗

TiDB Operation

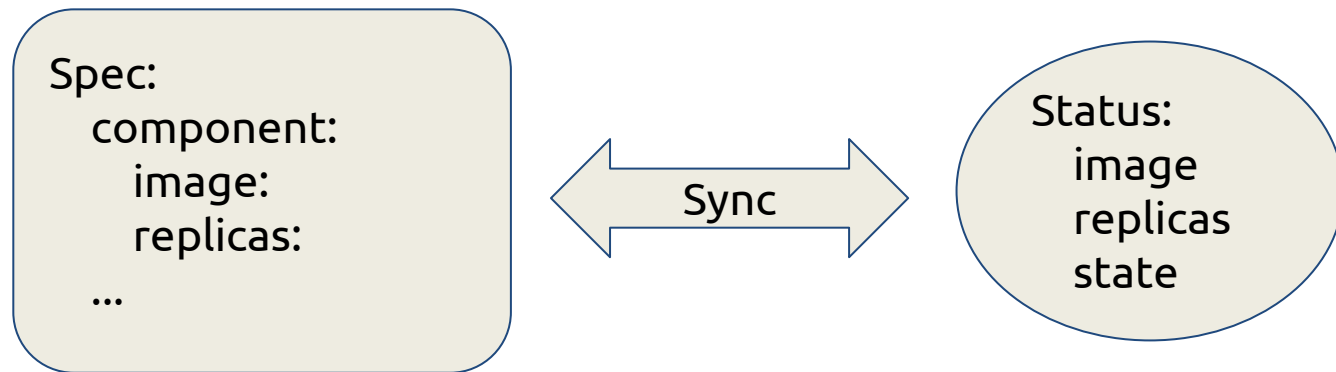
- Cluster bootstrap: Initial PD -> New PD joins existing cluster
- Safely delete PD
 - remove member using PD API
 - stop pd-server
- Safely delete TiKV
 - offline store using PD API
 - stop tikv-server
- Graceful upgrade
 - PD: transfer Raft leader
 - TiKV: evict Raft leaders
 - TiDB: evict DDL owner

Custom controller

Domain operation logic

- ThirdPartyResource (TPR), CustomResourceDefinition (CRD):
 - Simple & easy
 - Lack schema & versionning (added in newer version)
- Aggregated API Server (AA):
 - Powerful but complicated
 - Coupled with the built-in API Server, hard to deploy

Custom controller



Custom controller

```
type Manager interface {  
    Sync(*TidbCluster) error  
}
```

```
...  
status:  
  tikv:  
    stores:  
      "5":  
        podName: demo-tikv-2  
        state: Up  
...
```

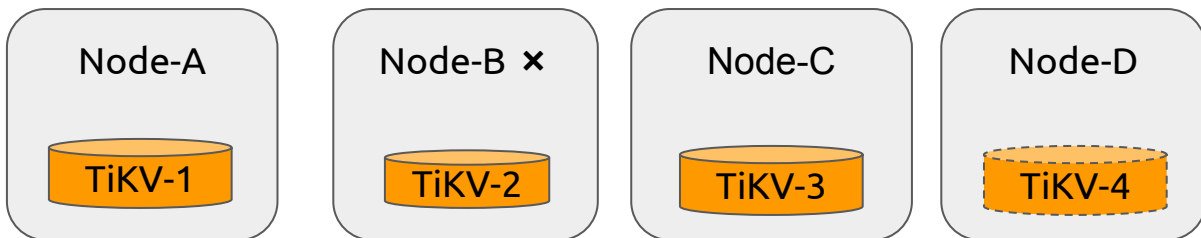
```
apiVersion: pingcap.com/v1alpha1  
kind: TidbCluster  
metadata:  
  name: demo  
spec:  
  pd:  
    image: pingcap/pd:v2.1.0  
    replicas: 3  
    requests:  
      cpu: "4"  
      memory: "8Gi"  
  ...  
  tikv:  
    image: pingcap/tikv:v2.1.0
```



Custom controller

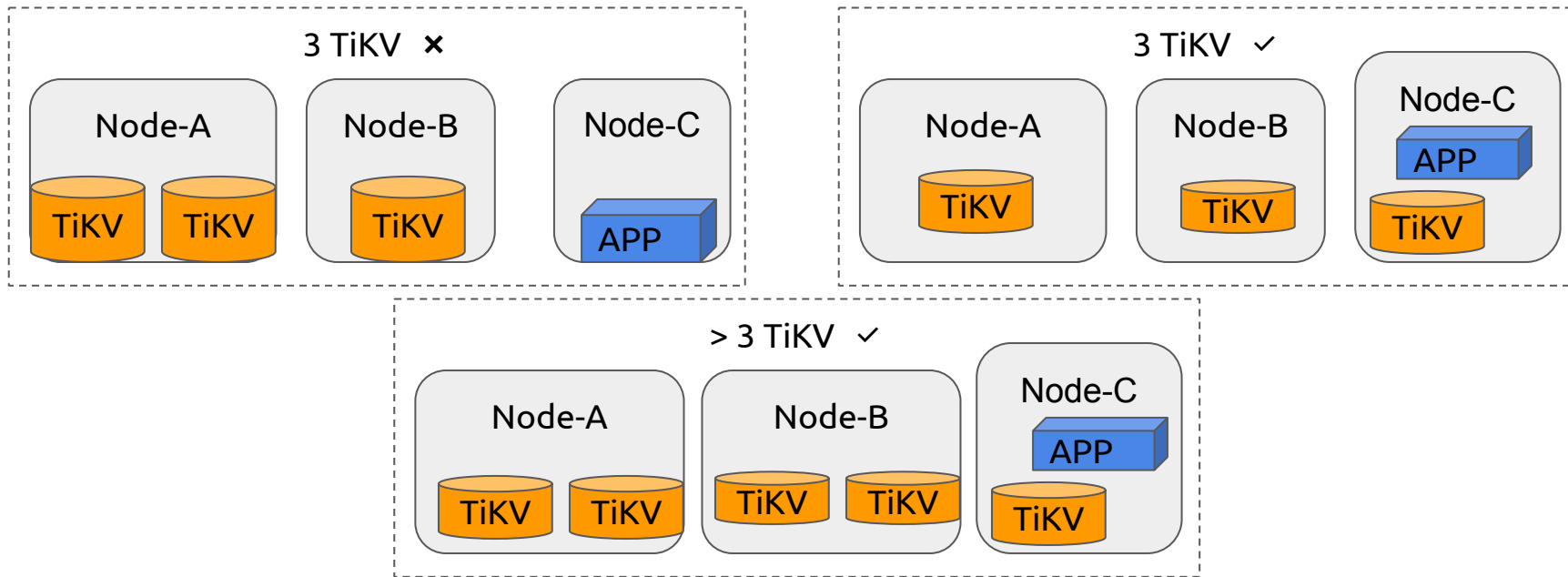
StatefulSet with Local PV failover:

1. Increase replicas when failure occurs
2. Decrease replicas when node come back (ordinal limitations with statefulset)



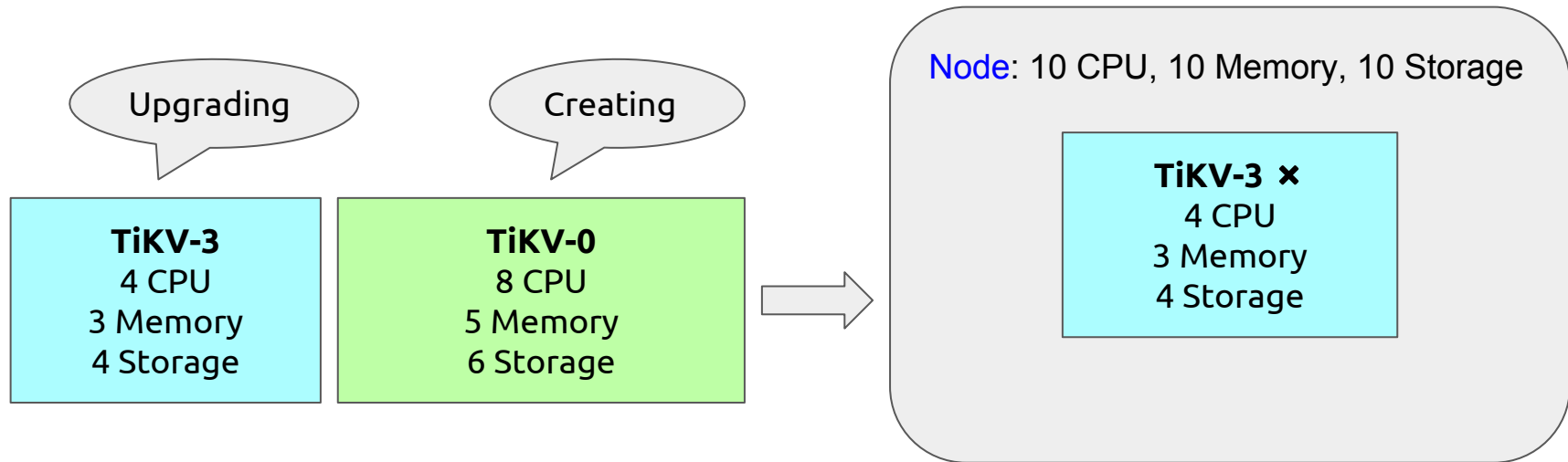
How we schedule stateful app

- Schedule consider existing pods topology



How we schedule stateful app

- Schedule consider *virtual* resource for local volume



TiDB Operator

Open-sourced! ヽ(=^▽^=)ﾉ 🎉

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

邀你一起探讨 人工智能商业化下的技术演进.



机器学习应用和实践

计算机视觉

NLP和语音技术

搜索推荐与算法

AI工具与框架

数据智能驱动业务



- Google、微软、亚马逊等国际巨头的AI产品技术干货
- 企业如何根据业务选择技术框架及搭建AI团队
- BAT、美团、京东如何利用AI技术赋能各业务线
- 人工智能的未来发展方向, 如何快人一步把握机会

AiCon
全球人工智能与机器学习技术大会

2018.12.20-23
北京·国际会议中心

入场券



TGO 鲲鹏会

汇聚全球科技领导者的高端社群

📍 全球9大城市

👤 700+ 高端科技领导者

使命
Mission

为社会输送更多
优秀的科技领导者

愿景
Vision

构建全球顶尖的有技术背景的
优秀人才成长平台



扫码了解更多内容

THANKS