

Go 夜读 & 云原生社区

TiDB Operator 架构与实现

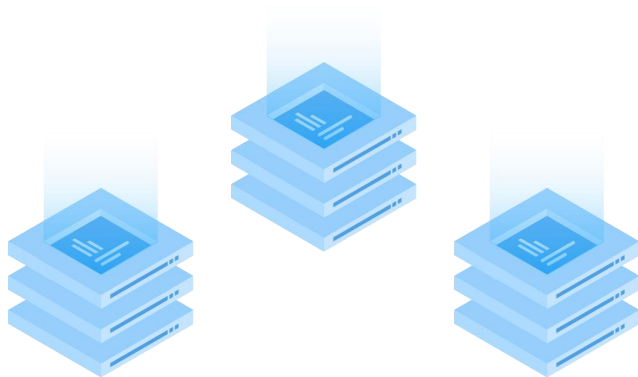
扫描下方二维码

填写你的问题



TiDB Operator: Design & Implementation

Presented by Yecheng Fu (@cofyc)



Agenda

- TiDB Operator 简介
- 扩展 Kubernetes 的几种方式
- TiDB Operator 实现

TiDB Operator 是什么



Cloud Native Era: Portable, Scalable, Automated

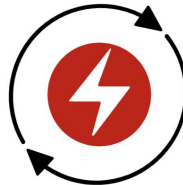
Full lifecycle management of TiDB cluster

- Deployment
- Upgrading
- Scaling
- Handle network, hardware failures, etc.
- Backup/Restore/Data migration
- ...

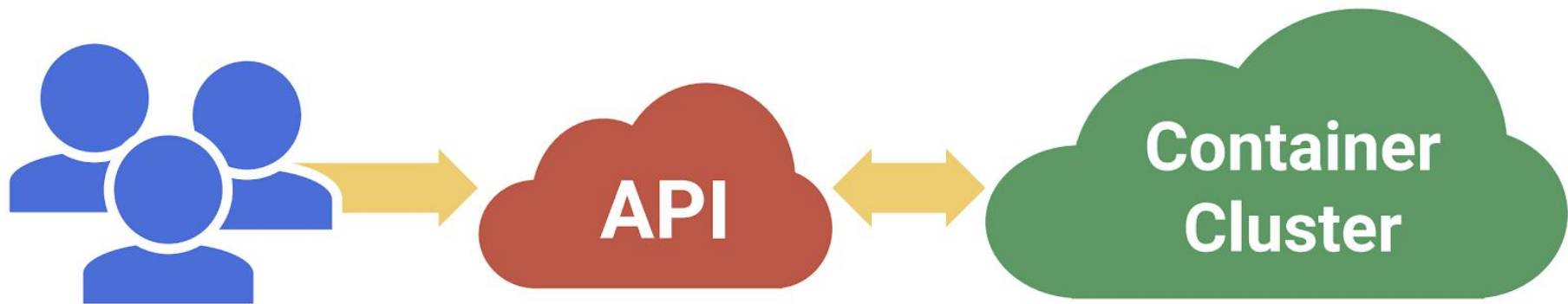
TiDB



TiDB Operator



Kubernetes Pattern - Declarative Model



- 用户描述自己的期望，提交给 Kubernetes API Server
- Kubernetes 根据用户的期望以及当前的状态，协调各方达成用户的期望

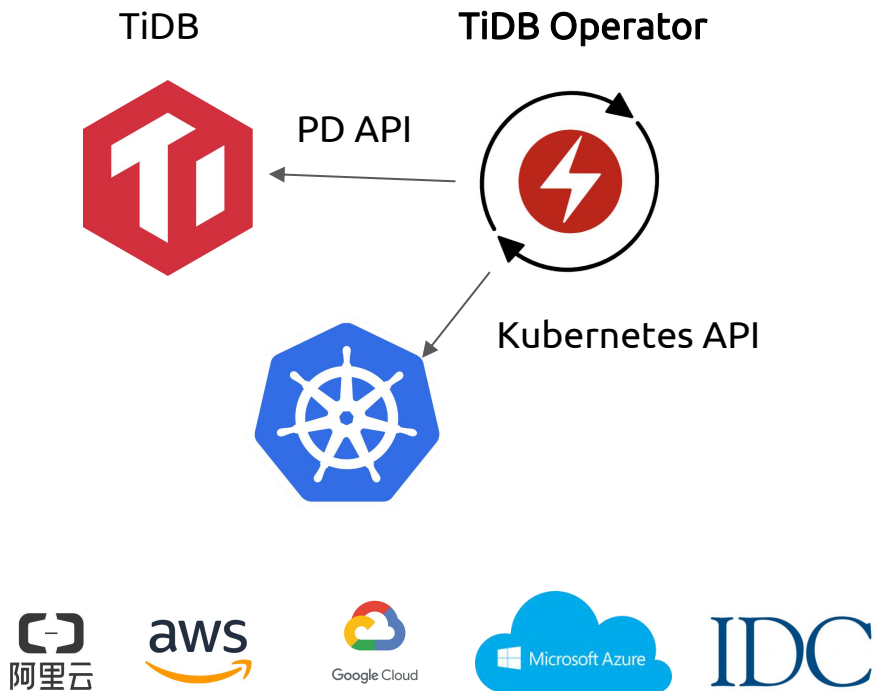
Kubernetes Pattern - Declarative Model

- Kubernetes has pod, deployment, statefulset, etc.
- But it does not know how to operate a TiDB cluster...

Kubernetes Pattern - Declarative Model

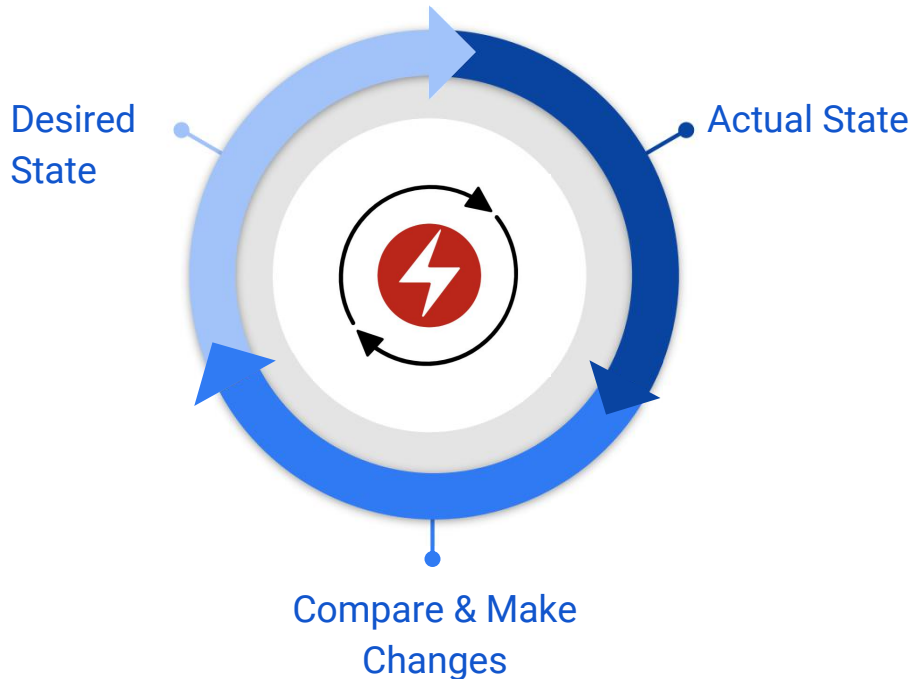
- Kubernetes has pod, deployment, statefulset, etc.
- But it does not know how to operate TiDB cluster...
- Until we implemented TiDB Operator

Operator Pattern - Extending Kubernetes



Operator Pattern - Custom Resource & Controller

```
1  apiVersion: pingcap.com/v1alpha1
2  kind: TidbCluster
3  spec:
4    version: v4.0.3
5    pd:
6      baseImage: pingcap/pd
7      replicas: 3
8      requests:
9        storage: "1Gi"
10     ...
11    tikv:
12      <desired state of TiKV>
13    tidb:
14      <desired state of TiDB>
15  status:
16    pd:
17      <observed state of PD>
18    tikv:
19      <observed state of TiKV>
20    tidb:
21      <observed state of TiDB>
```



扩展 Kubernetes 的几种方式



Common ways to extend Kubernetes

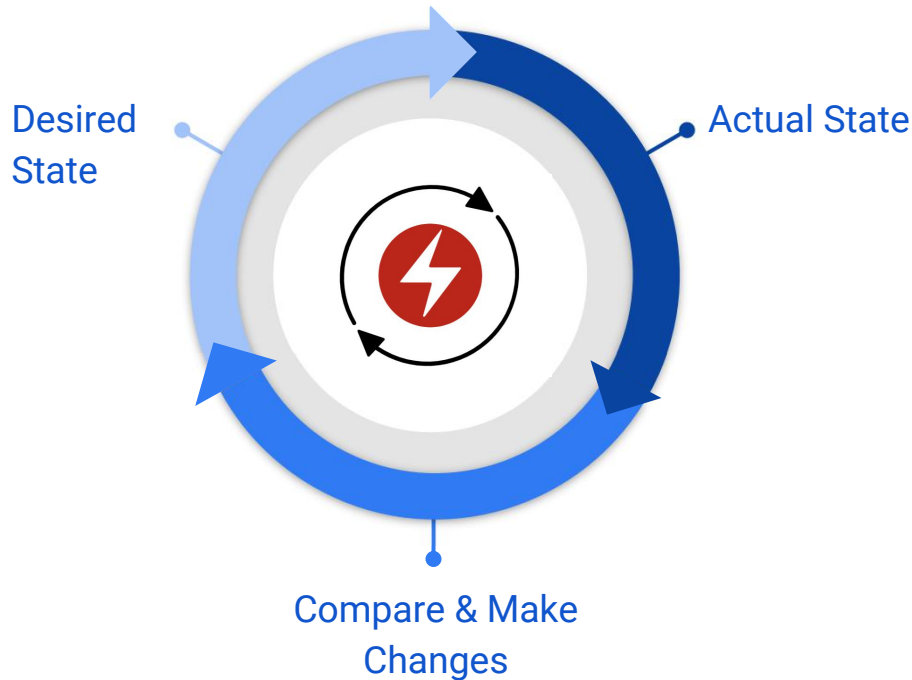
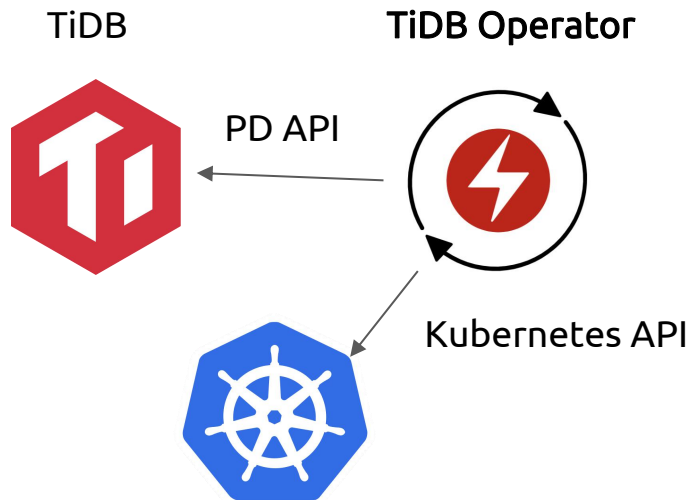
- Custom Resource
 - TidbCluster
 - TidbInitializer
 - TidbMonitor
 - Backup/Restore
- Custom Controller
- Scheduler Extender (optional)
- Admission Webhook (optional)
- ...Scheduler framework, Aggregated APIServer

Custom Resource

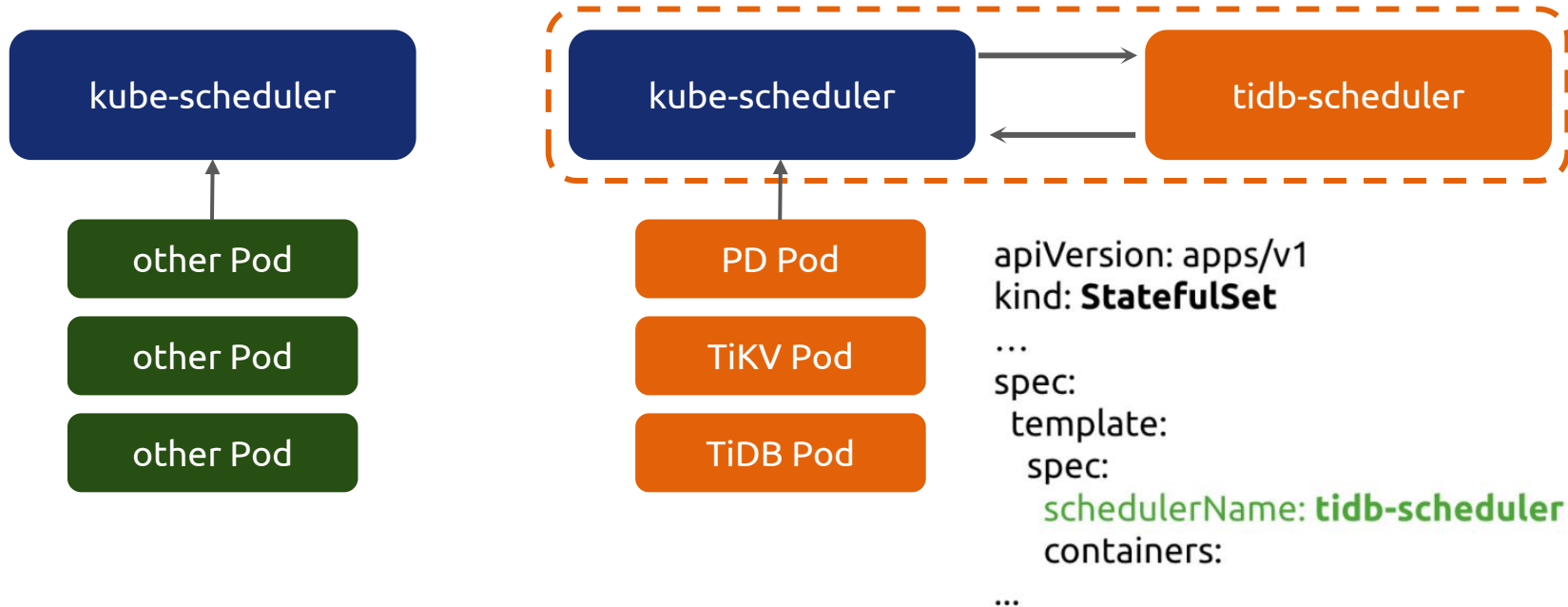
```
1  apiVersion: pingcap.com/v1alpha1
2  kind: TidbCluster
3  spec:
4    version: v4.0.3
5    pd:
6      baseImage: pingcap/pd
7      replicas: 3
8      requests:
9        storage: "1Gi"
10     ...
11    tikv:
12      <desired state of TiKV>
13    tidb:
14      <desired state of TiDB>
15  status:
16    pd:
17      <observed state of PD>
18    tikv:
19      <observed state of TiKV>
20    tidb:
21      <observed state of TiDB>
```

```
1  kind: CustomResourceDefinition
2  spec:
3    group: pingcap.com
4    names:
5      kind: TidbCluster
6      plural: tidbclusters
7      shortNames:
8        - tc
9    validation:
10      openAPIV3Schema:
11        type: object
12        properties:
13          spec:
14            type: object
15            properties:
16              pd:
17                type: object
18              tikv:
19                type: object
20              tidb:
21                type: object
```

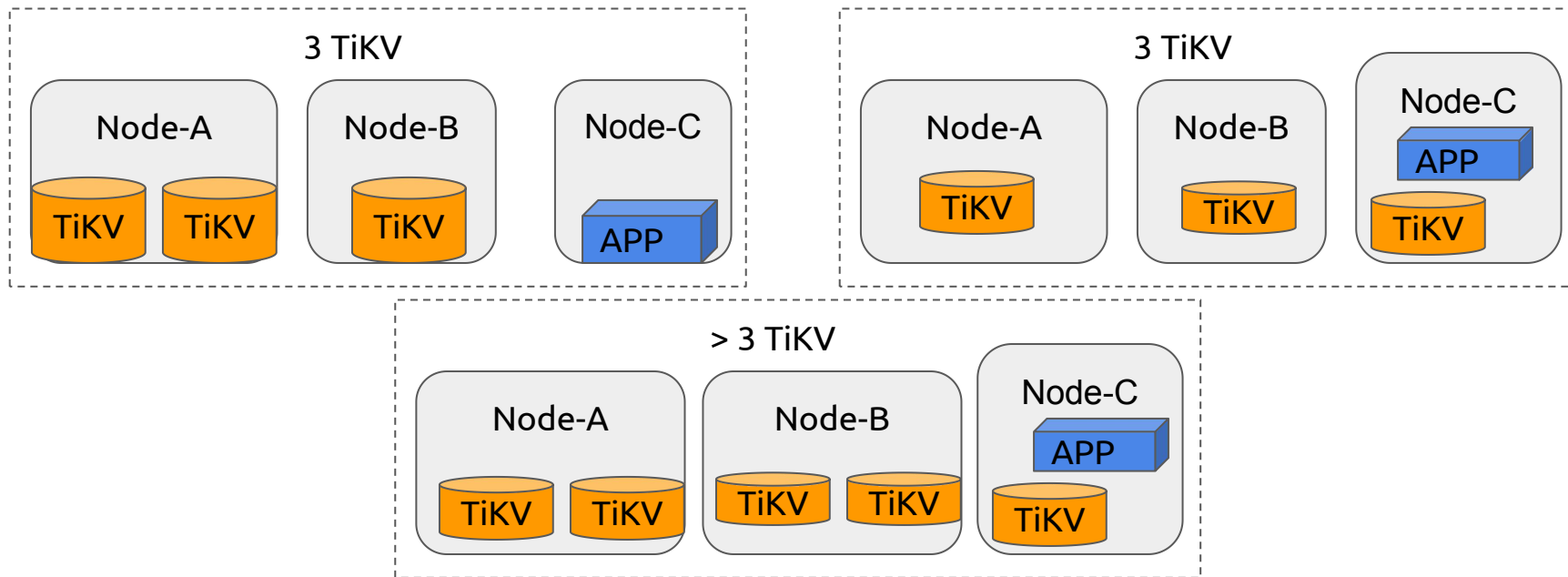
Custom Controller



Scheduler Extender

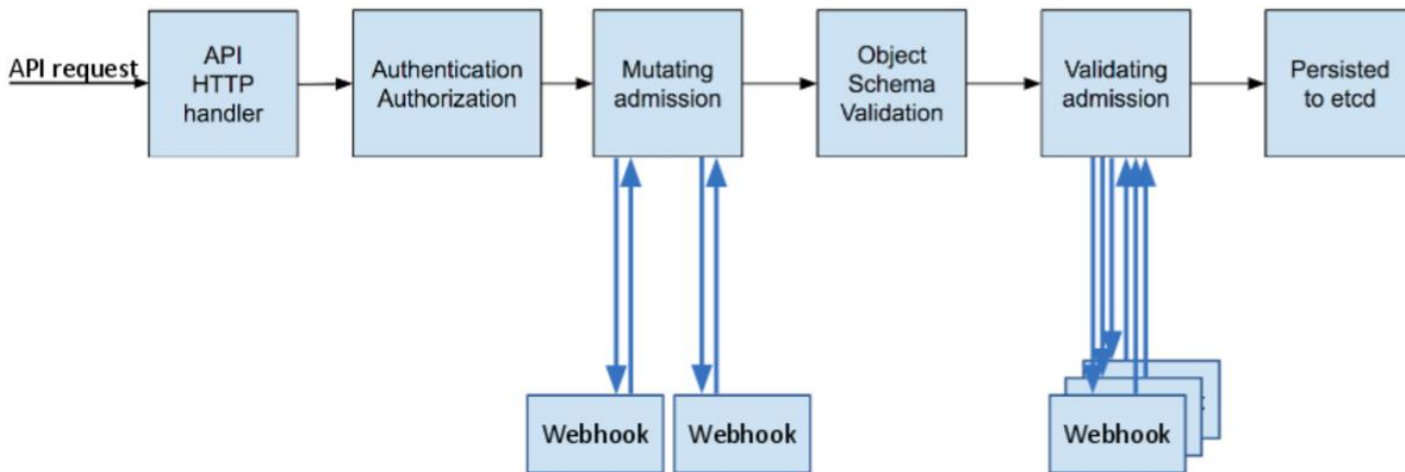


Scheduler Extender



* *EvenPodsSpread* feature introduced in Kubernetes 1.18 provides a more flexible solution

Admission Webhook



Able to intercept API requests and may change or deny it.

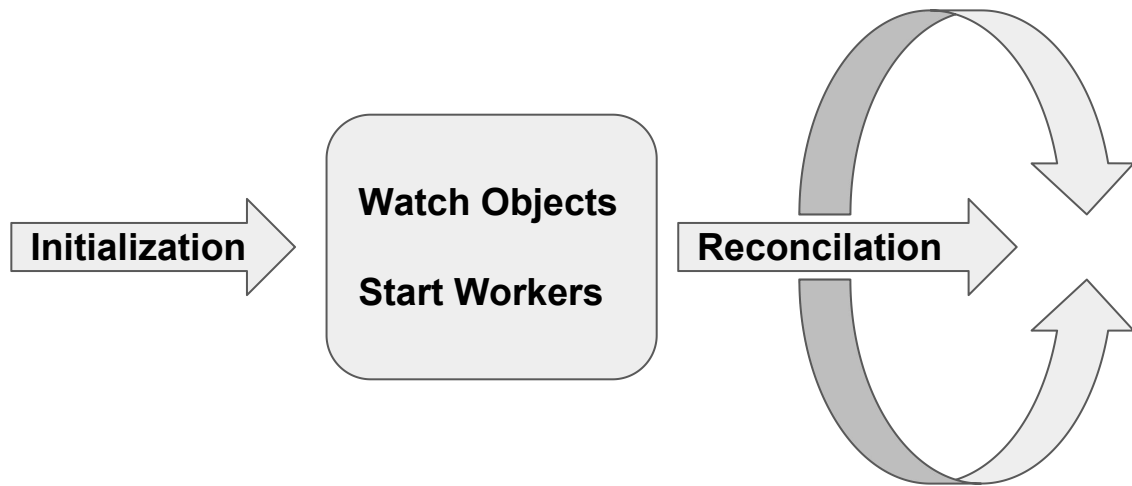
TiDB Operator 实现



TiDB Operator Implementation

- The tidb-controller-manager
 - Initialization
 - Watch objects (Tidb Cluster CR)
 - Reconciliation
 - on change of watched CR objects
 - periodically for TiDB clusters (30s by default)
- Full lifecycle management of a TiDB Cluster

tidb-controller-manager - Flowchart



tidb-controller-manager - Initialization

```
1 // Create informer factories for custom and Kubernetes
2 informerFactory := informers.NewSharedInformerFactoryWithOptions(cli, resyncDuration)
3 kubeInformerFactory = kubeinformers.NewSharedInformerFactoryWithOptions(kubeCli, resyncDuration)
4 // Create controllers
5 tcContrller := tidbcluster.NewController(cli, kubeCli, informerFactory, kubeInformerFactory)
6 ...
7 // Start informer factories after all controller are initialized
8 informerFactory.Start(ctx.Done())
9 kubeInformerFactory.Start(ctx.Done())
10 // Wait for all started informers' cache were synced
11 informerFactory.WaitForCacheSync(wait.NeverStop)
12 kubeInformerFactory.WaitForCacheSync(wait.NeverStop)
13 // Start controllers
14 go tcContrller.Run(5, ctx.Done())
15 ...
16 select {}
```

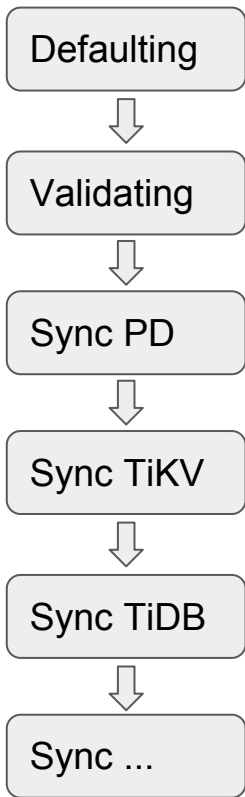
tidb-controller-manager - Watch Objects

```
1 func NewController(...) *Controller {
2     ...
3     // create a work queue
4     c.queue = workqueue.NewNamedRateLimitingQueue(workqueue.DefaultControllerRateLimiter(), "tidbcluster")
5     // add object to the queue on event
6     informerFactory.Pingcap().V1alpha1().TidbClusters().Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{
7         AddFunc: c.enqueue,
8         UpdateFunc: func(old, cur interface{}) {
9             c.enqueue(cur)
10        },
11        DeleteFunc: c.enqueue,
12    })
13    ...
14 }
15
16 func (c *Controller) enqueue(tc *v1alpha1.TidbCluster) {
17     key := getObjectKey(obj)
18     c.queue.Add(key)
19 }
```

tidb-controller-manager - Start workers

```
1 func (c *Controller) Run(workers int, stopCh <-chan struct{}) {
2     defer c.queue.ShutDown()
3     for i := 0; i < workers; i++ {
4         go wait.Until(c.worker, time.Second, stopCh)
5     }
6     <-stopCh
7 }
8
9 func (c *Controller) worker() {
10     for c.processNextWorkItem() {
11     }
12 }
13
14 func (c *Controller) processNextWorkItem() bool {
15     key, quit := c.queue.Get()
16     if quit {
17         return false
18     }
19     defer c.queue.Done(key)
20     if err := c.sync(key.(string)); err != nil {
21         c.queue.AddRateLimited(key)
22     } else {
23         c.queue.Forget(key)
24     }
25     return true
26 }
```

tidb-controller-manager - Reconciliation



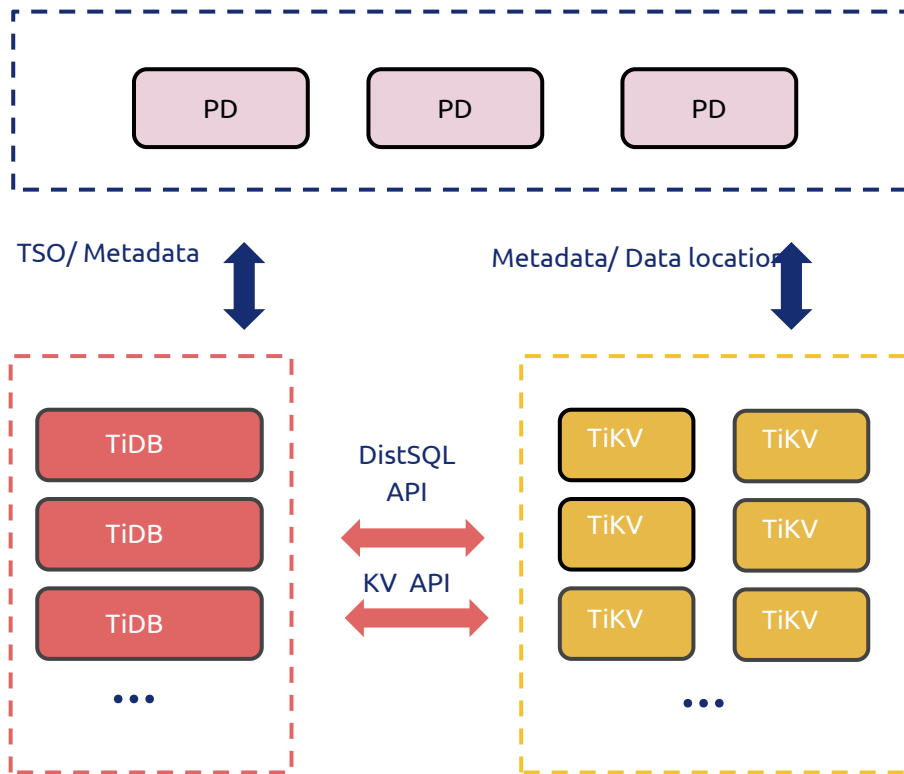
```
1 func (c *Controller) Sync(key string) error {
2     tc := getTidbClusterFromCache(key)
3     c.defaulting(tc)
4     if !c.validating(tc) {
5         return nil
6     }
7     defer func() {
8         c.updateStatus(tc)
9     }()
10    if err := c.syncPD(tc); err != nil {
11        return err
12    }
13    if err := c.syncTiKV(tc); err != nil {
14        return err
15    }
16    if err := c.syncTiDB(tc); err != nil {
17        return err
18    }
19    ...
20    return nil
21 }
```


Full Lifecycle management of a TiDB Cluster

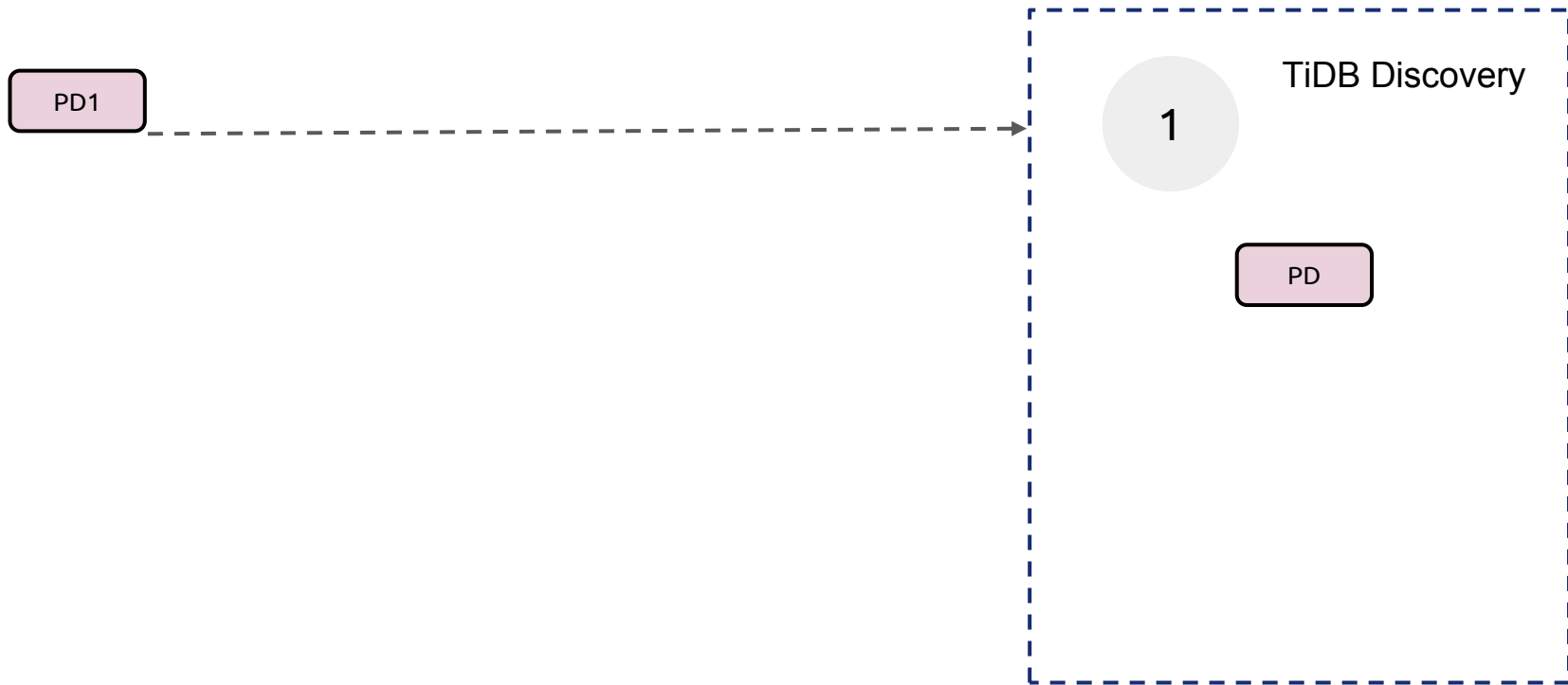
- Deploying
 - Bootstrapping
 - Configure services/configmaps, etc.
- Upgrading
 - Change version, config, etc.
- Scaling
- Automatic Failover
 - Create replacements for failed replicas

Deploying - Bootstrapping

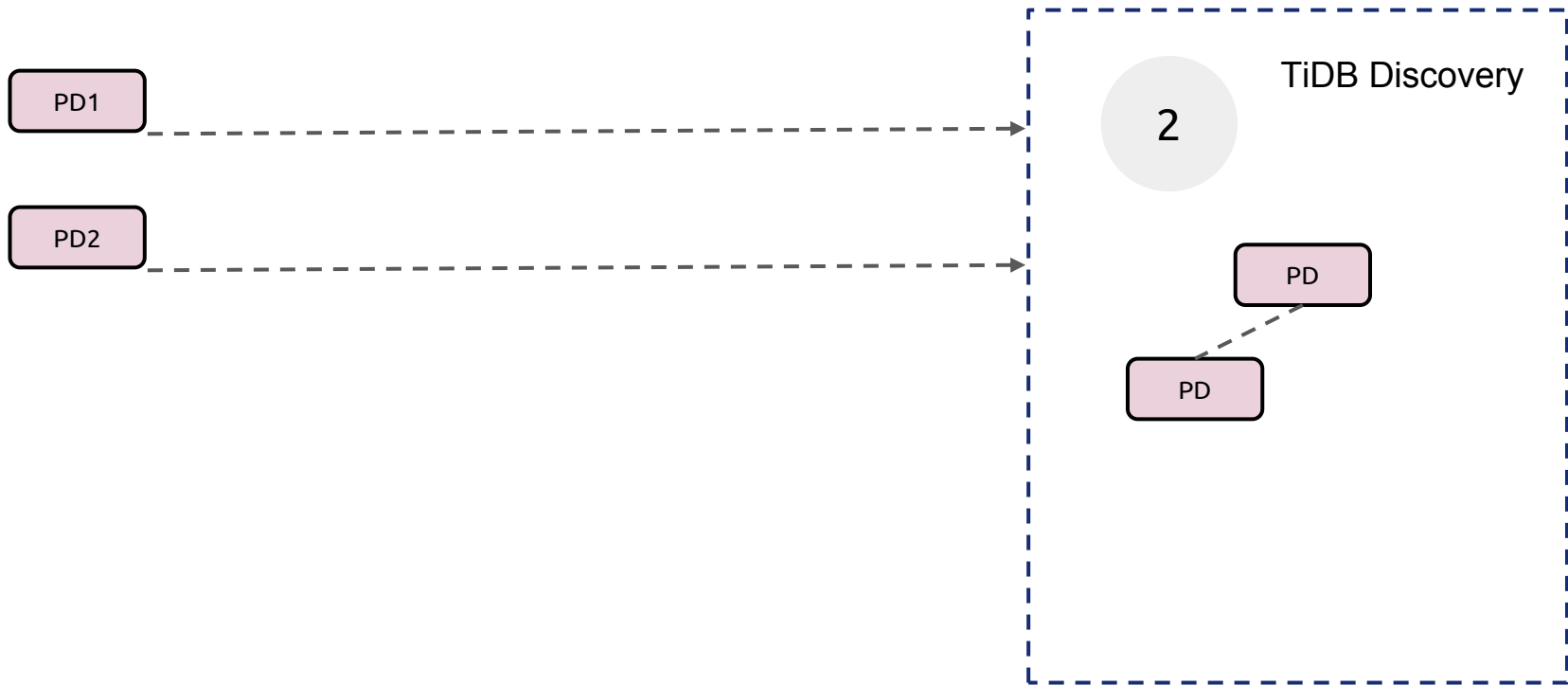
1. Bootstrap PD cluster with a discovery service
2. Start TiKV replicas and join the PD cluster
3. Start TiDB replicas and join the PD Cluster
4. Create TiDB Service...



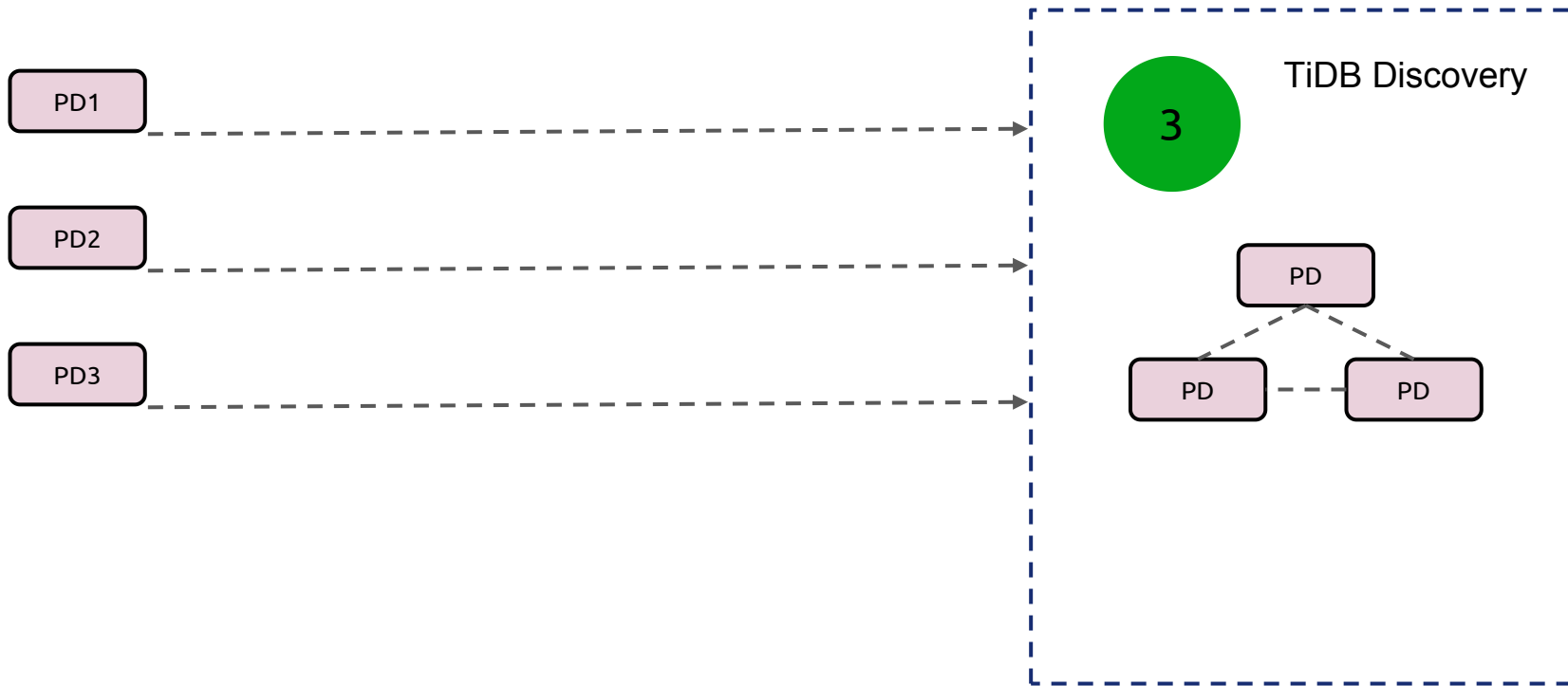
Deploying - Bootstrapping



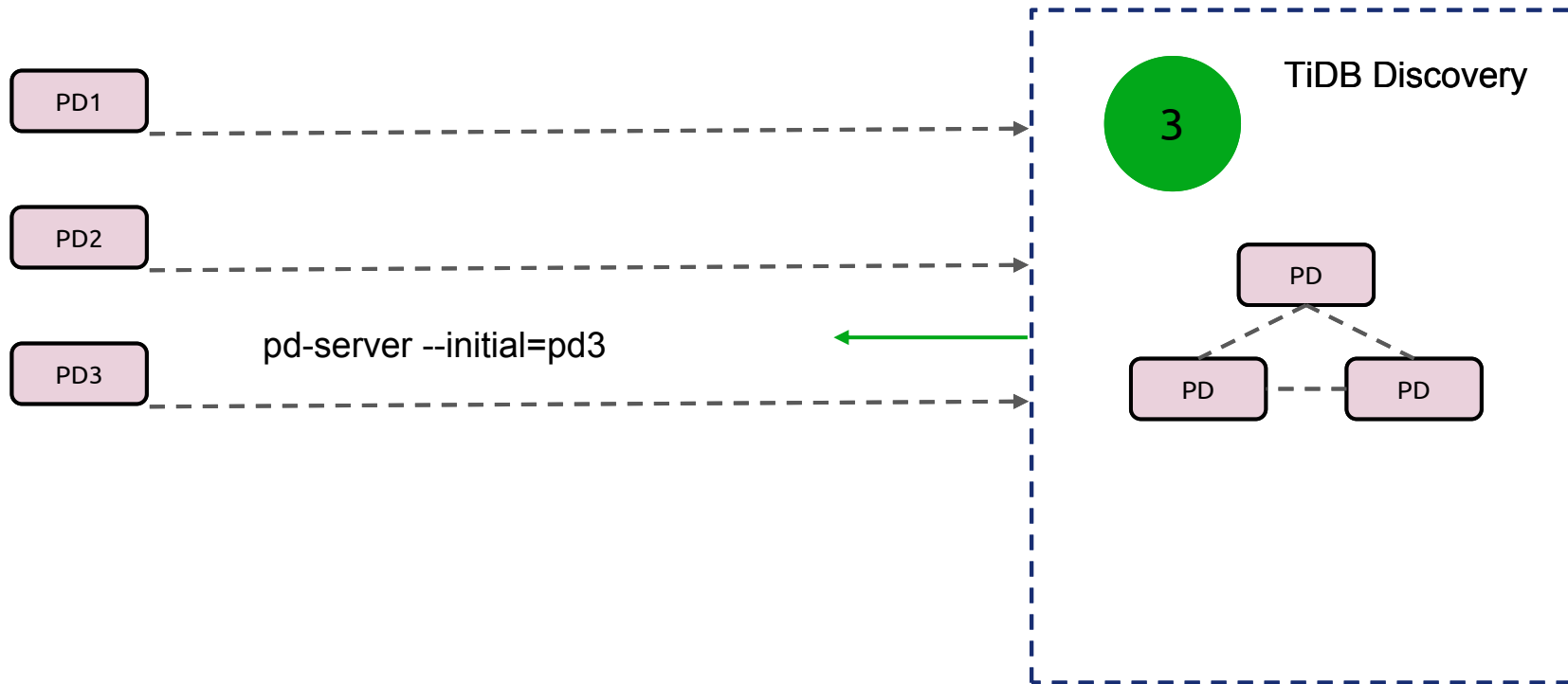
Deploying - Bootstrapping



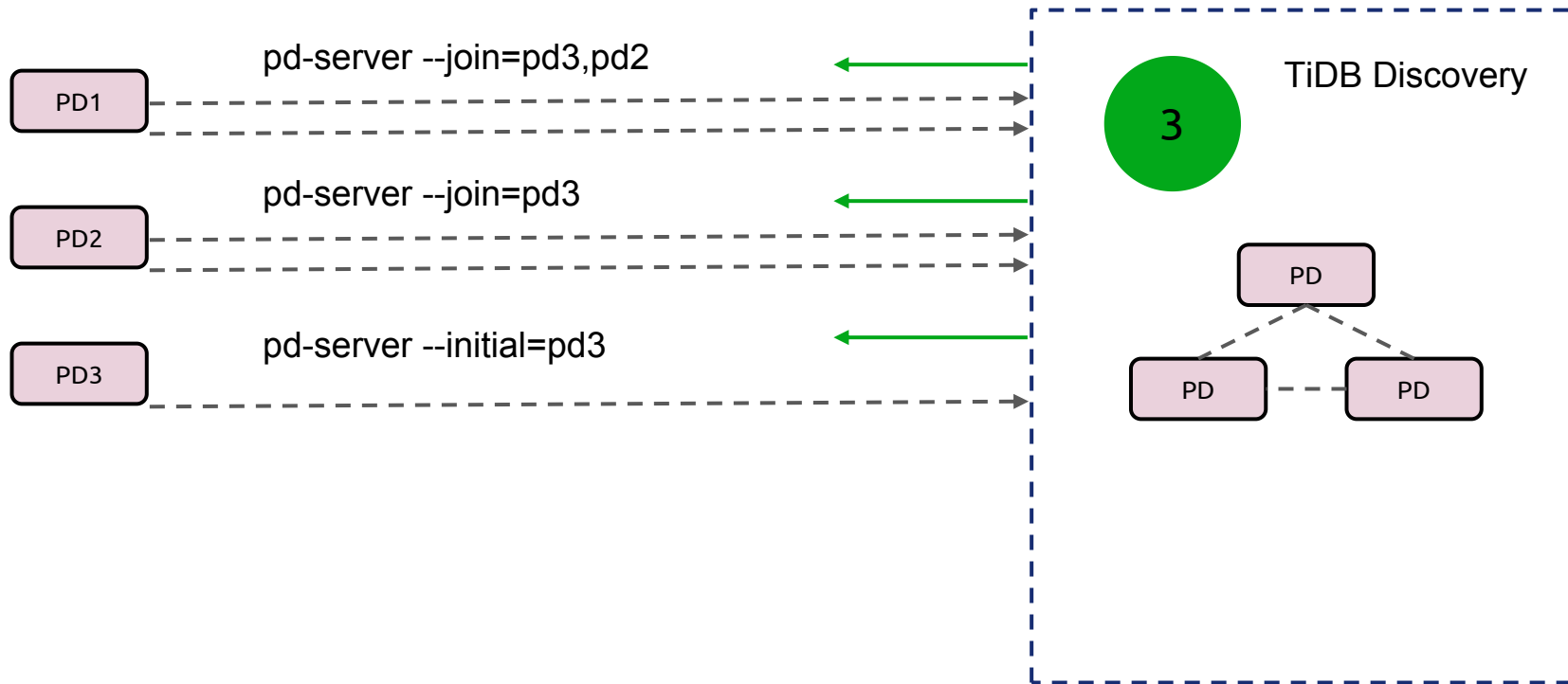
Deploying - Bootstrapping



Deploying - Bootstrapping



Deploying - Bootstrapping



Upgrading

It's ok to upgrade a High Availability system like TiDB directly.

But, to avoid spikes we use StatefulSet partition to perform a phased roll out.

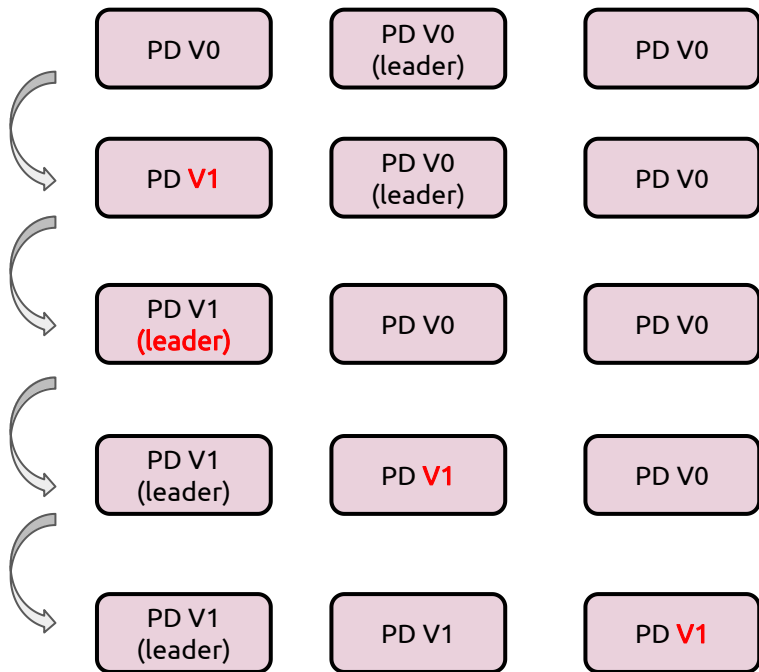
- Upgrade PD
- Upgrade TiKV
- Upgrade TiDB

```
1  status:
2      currentRevision: <revision-v0>
3      currentReplicas: 3
4      updateRevision: <revision-v1>
5      updateReplicas: 0
6
7  status:
8      currentRevision: <revision-v0>
9      currentReplicas: 2
10     updateRevision: <revision-v1>
11     updateReplicas: 1
12
13  status:
14     currentRevision: <revision-v0>
15     currentReplicas: 1
16     updateRevision: <revision-v1>
17     updateReplicas: 2
18
19  status:
20     currentRevision: <revision-v1>
21     currentReplicas: 3
22     updateRevision: <revision-v1>
23     updateReplicas: 3
```


Upgrading - PD

For each PD replica

- If it's the leader of PD cluster, transfer the leadership
- Upgrade



Upgrading - PD

For each PD replica

- If it's the leader of PD cluster, transfer the leadership
- Upgrade

```
1  for i in <pod ordinals from end to the start>:  
2      if isUpToDate(i):  
3          continue  
4      if not pd.IsLeader(i):  
5          setPartition(i)  
6      else  
7          pd.TransferLeader(i)
```

Upgrading - TiKV

For each TiKV replica

- Evict region leaders
- Wait for all region leaders are evicted or the timeout expired
- Upgrade

```
1  for i in <pod ordinals from end to the start>:  
2      if isUpToDate(i):  
3          continue  
4      if pd.RegionLeaderCount(i) == 0 or evictTimedOut(i):  
5          setPartition(i)  
6      else  
7          pd.evictRegionLeader(i)
```

Upgrading - TiDB

For each TiDB replica

- Upgrade

```
1  for i in <pod ordinals from end to the start>:  
2      if isUpToDate(i):  
3          continue  
4      setPartition(i)
```

Scaling In - PD

- Delete the member via PD API first
- After the member is deleted successfully, scale in the StatefulSet

```
1  for i in <pods need to delete>:  
2      if pd.Delete(i):  
3          delete(i)  
4      else  
5          return RequeueError
```

Scaling In - TiKV

- Delete the TiKV store via PD API first
- Wait for the store to be tombstoned
- Scale the statefulset

```
1  for i in <pods from the end to the start>:  
2      if pd.IsTombStoned(i):  
3          scaleIn(i)  
4      else  
5          pd.DeleteStore(i)
```

Scaling In - TiDB

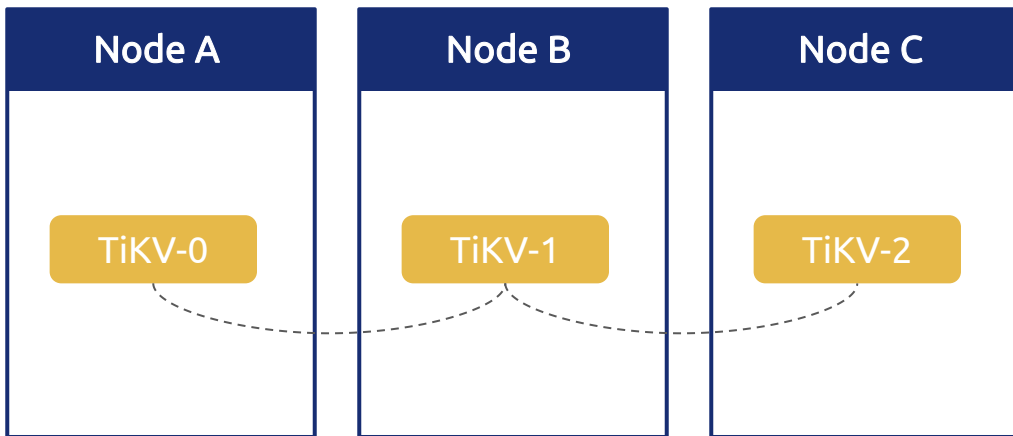
- Scale the statefulset directly

Scaling Out

- Clean retained PVCs if necessary
- Scale the statefulset
- Wait for new replicas to join the cluster

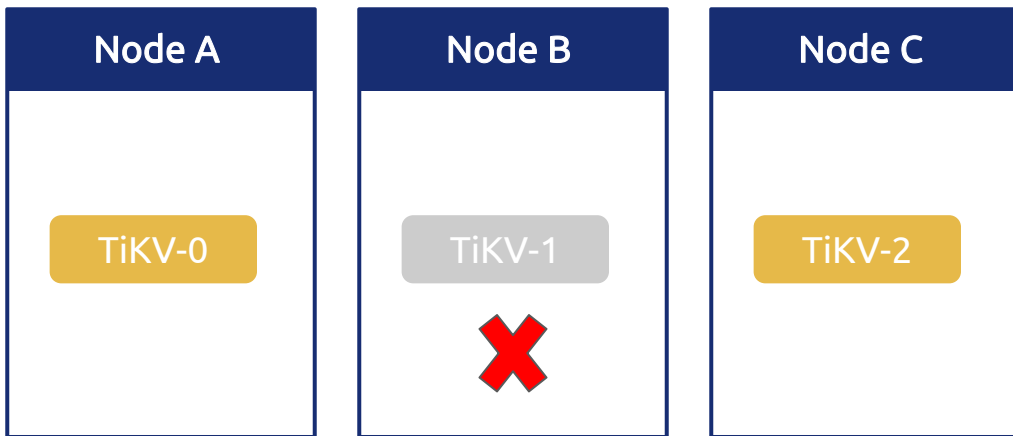
Automatic Failover

- 3 default replicas in each raft group can tolerate one member failure

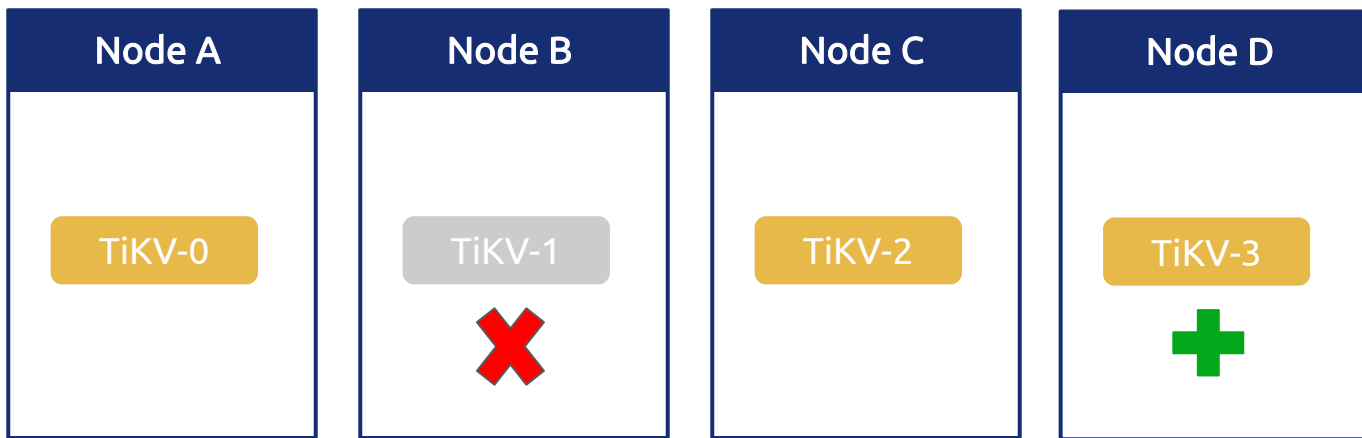


Automatic Failover

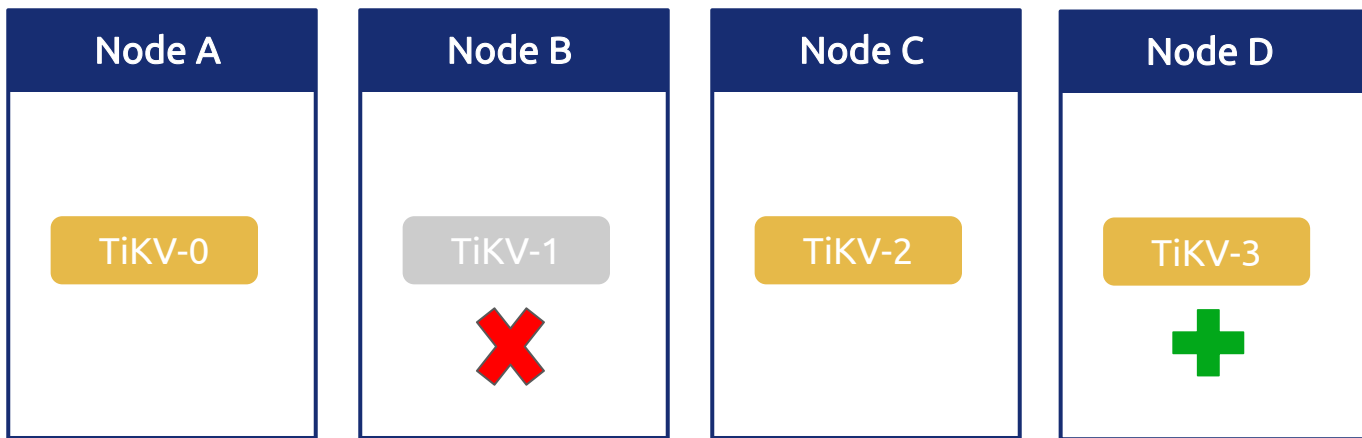
- What if a replicac is down for a long time
 - It's dangerous!
- What can we do?



Automatic Failover



Automatic Failover



Automatic Failover

```
1  for store in <all tikv stores>:  
2      if pd.IsDown(store) and sinceDownTime(store) > 5m:  
3          markStoreFailed(store)  
4  scaleStatefulSet(currentReplicas + failedStores)
```

- A store is marked **Down** by PD when it's unreachable for **30m** (by default)
- A store is marked **Failed** by TiDB Operator when it's **Down** for **5m** (by default)
- The controller will **increase the replicas** by the number of failed replicas

Summary

- TiDB Operator 简介
 - Cloud Native Era: Full Automation
 - Kubernetes Pattern: Declarative Model
 - Operator Pattern: Extending Kubernetes
- 扩展 Kubernetes 的几种方式
 - Custom Resource
 - Custom Controller
 - Scheduler Extender
 - Admission Webhook
- TiDB Operator 实现
 - the tidb-controller-manager
 - Full lifecycle management of a TiDB Cluster
 - Deploying, Upgrading, Scaling, Automatic Failover, etc.

Thank You!

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

Presented by Yecheng Fu (@cofyc)

