



KubeCon

THE LINUX FOUNDATION



China 2024



CloudNativeCon





KubeCon



CloudNativeCon



China 2024

快手的 100% 资源利用率提升： 从裸机迁移 100K Redis 到 Kubernetes

刘裕惺 快手 高级软件工程师
吴学强 云猿生 研发总监 | KubeBlocks 维护者

 KUAISHOU - KubeBlocks ·

关于我

- 前阿里云 PolarDB-X 团队技术 TL
- 开源版 PolarDB-X Maintainer
- 阿里巴巴 2021 年度开源先锋人物
- 目前是云猿生的研发总监、KubeBlocks Maintainer
- 对操作系统、分布式系统、数据库等多个领域感兴趣



吴学强
GitHub ID: free6om

大纲



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



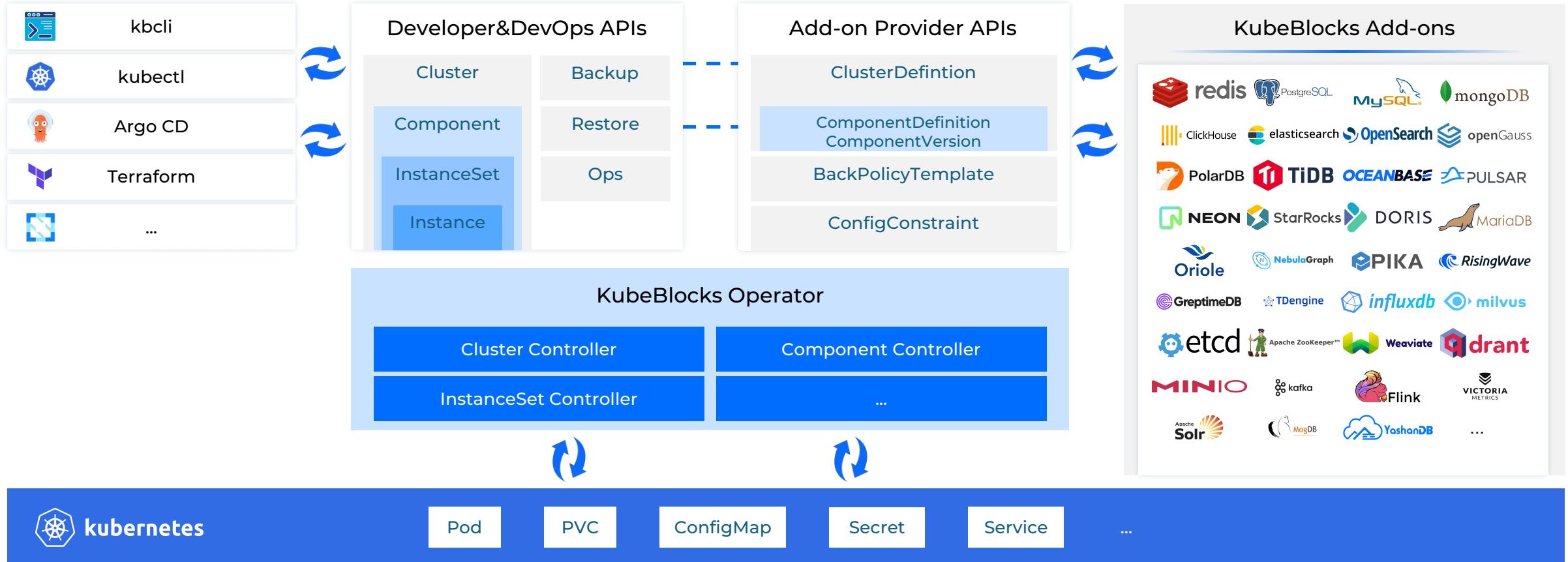
Open Source Dev & ML Summit

China 2024

- 什么是 KubeBlocks
- KubeBlocks 解决了什么问题 (单个 Redis 集群视角)
- 多 Redis 集群和多 K8s 集群 (大规模视角)
- Q&A

简介：什么是 KubeBlocks ?

什么是 KubeBlocks



一个专为将数据库运行在 K8s 设计的 Operator

- 支持 35 种数据库
- 可通过 add-on APIs 扩展更多
- 统一的 开发者&DevOps APIs
- 支持数据库生命周期管理
- 支持备份和恢复（包括PITR）
- 已与 Prometheus 和 Grafana 集成
- 更多特性

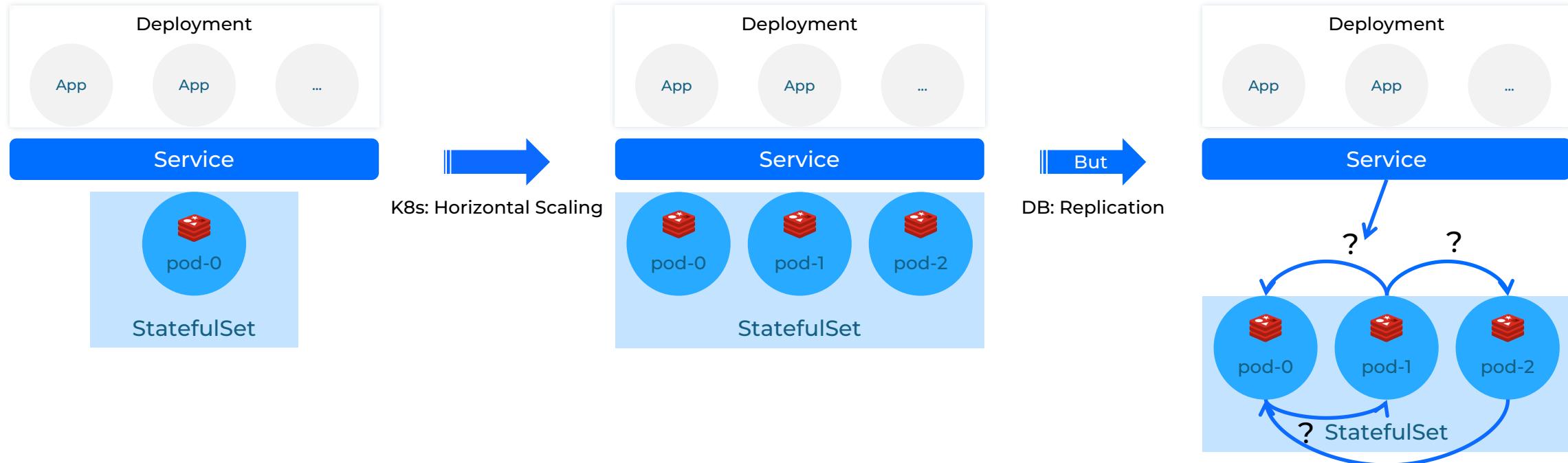
单个 Redis 集群视角：
KubeBlocks 解决了哪些问题从而让数据库更好的运行在 **K8s** 上？

如何处理数据复制 (Replication)



China 2024

在 K8s 上跑数据库是有非常挑战的



- 吞吐瓶颈
- 单点故障
- 数据丢失风险
- 如何找到具备读写能力的 Pod
- 如何知道并搭建正确的复制关系
- 更新时如何尽量降低对服务的影响
-

如何处理数据复制 (Replication)



China 2024

角色 (Role) : 在有状态之上增加的一层新的抽象

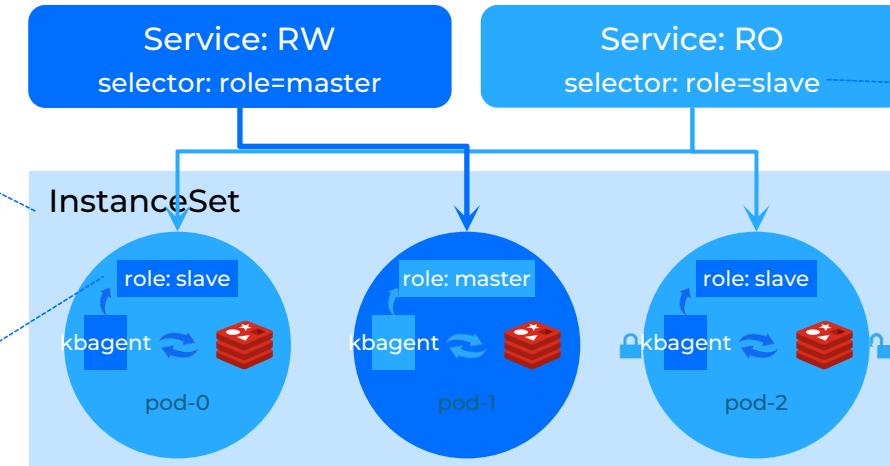
1. InstanceSet API

- 一个专为数据库设计的 Workload API
- 与 StatefulSet 类似，每个 Pod 都有固定的网络标识
- 在这之上，增加一层角色抽象

2. 每个 Pod 都有一个角色标签

- ① 角色周期性探测
- ② 触发角色事件
- ③ 观测角色事件
- ④ 更新 Pod 角色 Label

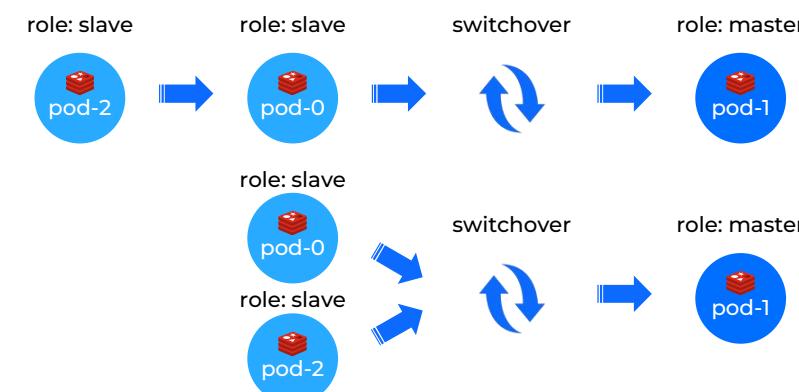
- 一个通用的角色探测框架
- 多种事件发送机制
- 如何处理网络分区
- 如何处理事件延迟



3. 基于角色的 Service Selector

4. 基于角色的成员管理

- 成员加入
- 成员离开
- switchover



5. 基于角色的更新策略

- Serial, Parallel, BestEffortParallel
- 计划内 switchover

RTO 从分钟级降到秒级

如何实现高可用 (HA)



KubeCon

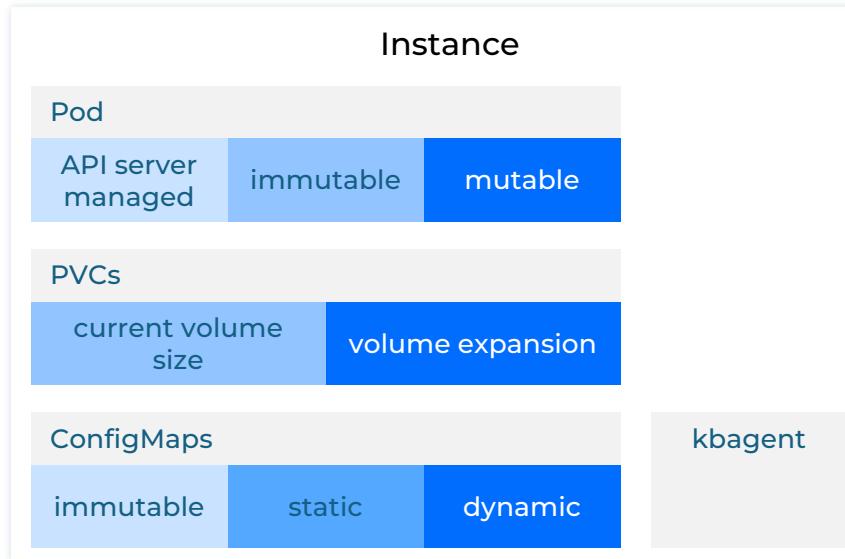


CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source DevOps & ML Summit

China 2024

实例原地更新 (Instance In-place Update)



- Pod 原地更新
- Volume 扩容
- 配置动态加载

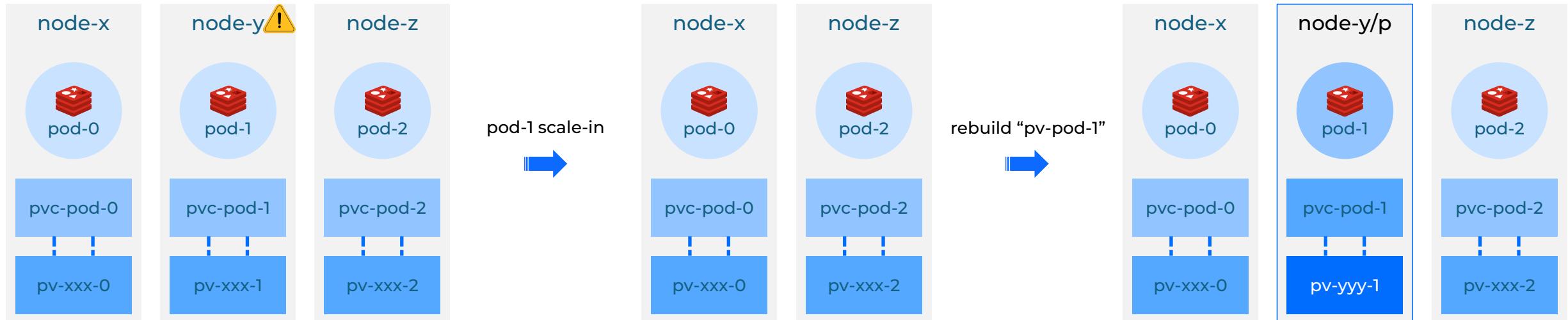
可用性 (Uptime) & MTBF 增加
RTO 从秒级降到接近于 0

如何实现高可用 (HA)



China 2024

实例重搭 (Instance Rebuilding)



Node node-y 磁盘故障。

node-y 下线，使用指定实例缩容 (specified instance scale-in) 功能将 pod-1 和 pvc-pod-1 删除。

在新的名字为 node-p 的 Node 上，使用实例重搭 (instance rebuilding) 功能创建跟老的 PV 包含相同数据的新 PV，即 pv-yyy-1。

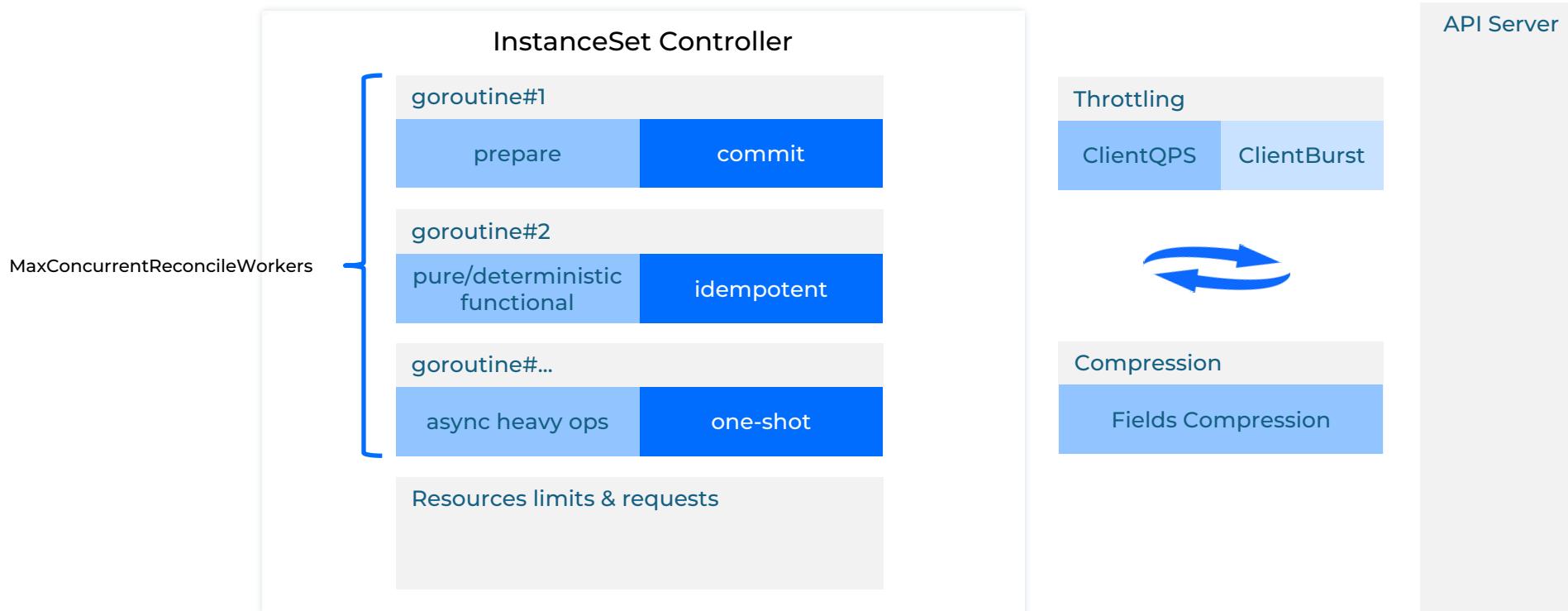
RTO 从几小时降到几分钟

如何管理单个大规模集群



China 2024

Operator P10K 问题



定义:

与 C10K 问题类似, P10K 问题是指通过优化 Operator 使其能够调谐包含大量 Pod 的单个 CR。

场景:

一个 Redis 集群接近 10000 个 Pods

小结



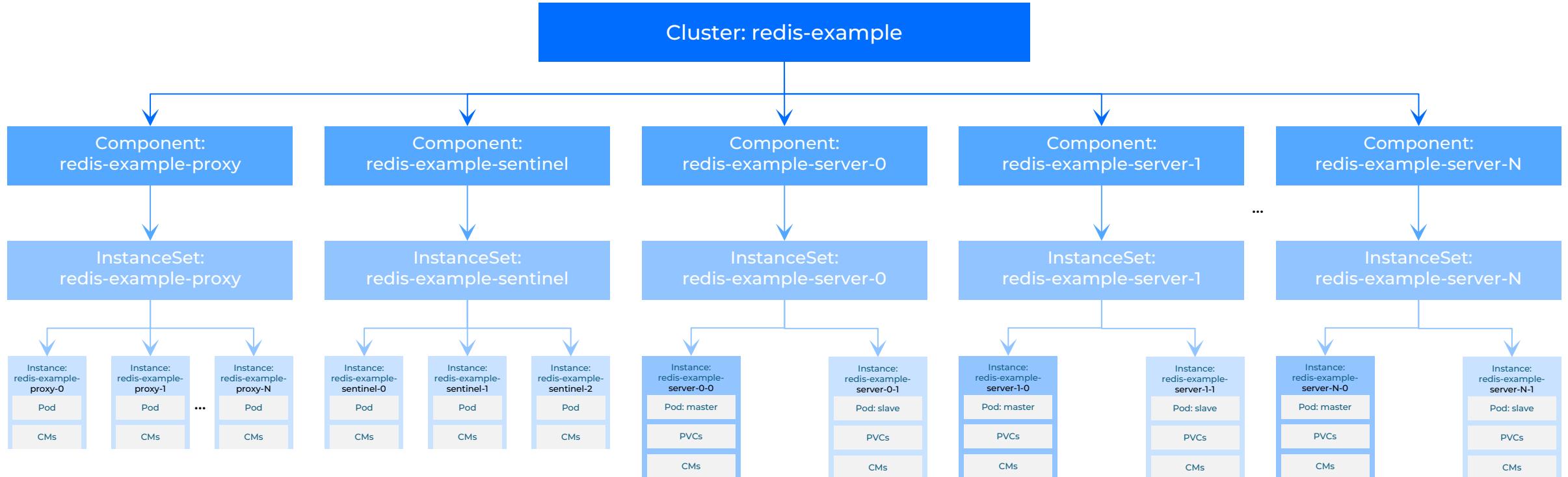
KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024



- 角色 (Role)
 - 基于角色的 Service Selector
 - 基于角色的成员管理
 - 基于角色的更新策略
- 高可用 (HA)
 - Pod 原地更新
 - Volume 扩缩容
 - 配置动态更新
 - 实例重搭 (Instance rebuilding)

- Operator P10K 问题
 - 流控
 - 对象压缩
 - 资源限制
 - **MaxConcurrentReconcileWorkers**
 - 两阶段调谐
 - 重 (heavy) 操作异步化

大规模视角：
快手如何借助 **KubeBlocks** 来将大规模 Redis
实例运行在多 **Kubernetes** 集群上？

关于我 - 刘裕惺



KubeCon



CloudNativeCon



China 2024



- 曾就职于阿里云，负责云原生应用分发产品商业化
- 开源&CNCF 项目 Dragonfly 和 Sealer 的 Maintainer
- 目前专注于快手数据库业务的云原生化
- Github ID: starnop

快手 - Redis 介绍



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



AI_dev
Open Source Dev & ML Summit

China 2024

有哪些特点? - 极端的大规模:

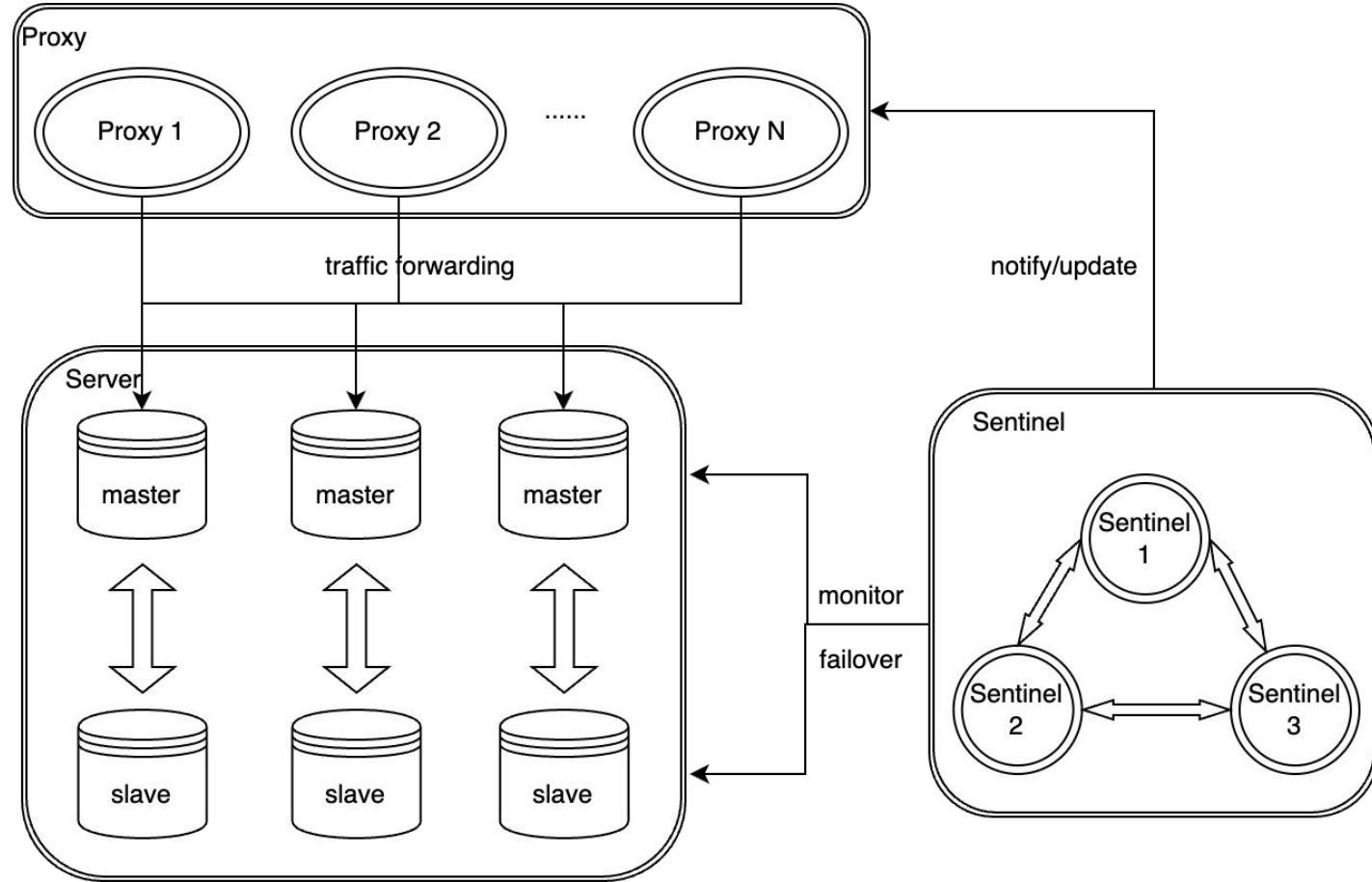
1. 超大规模总实例数
2. 单个 Redis 集群接近 10000 实例

为什么要云原生化?

1. 提高资源利用率并降低资源成本
2. 解耦基础设施与业务从而提高业务迭代敏捷性

组织结构?

1. Redis DBA 团队
2. 云原生团队



为什么选择 KubeBlocks



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



AI-dev
Open Source DevOps & ML Summit

China 2024

KubeBlocks - 面向有状态服务的 API

1. 强大的功能: 基于角色的管理能力等
2. 技术可复用: 支持多种数据库类型
3. 云原生化成本低:
 1. 实例生命周期与运维管理分离
 2. 通过 OpsRequest 提供面向过程的 API

Redis 集群编排定义



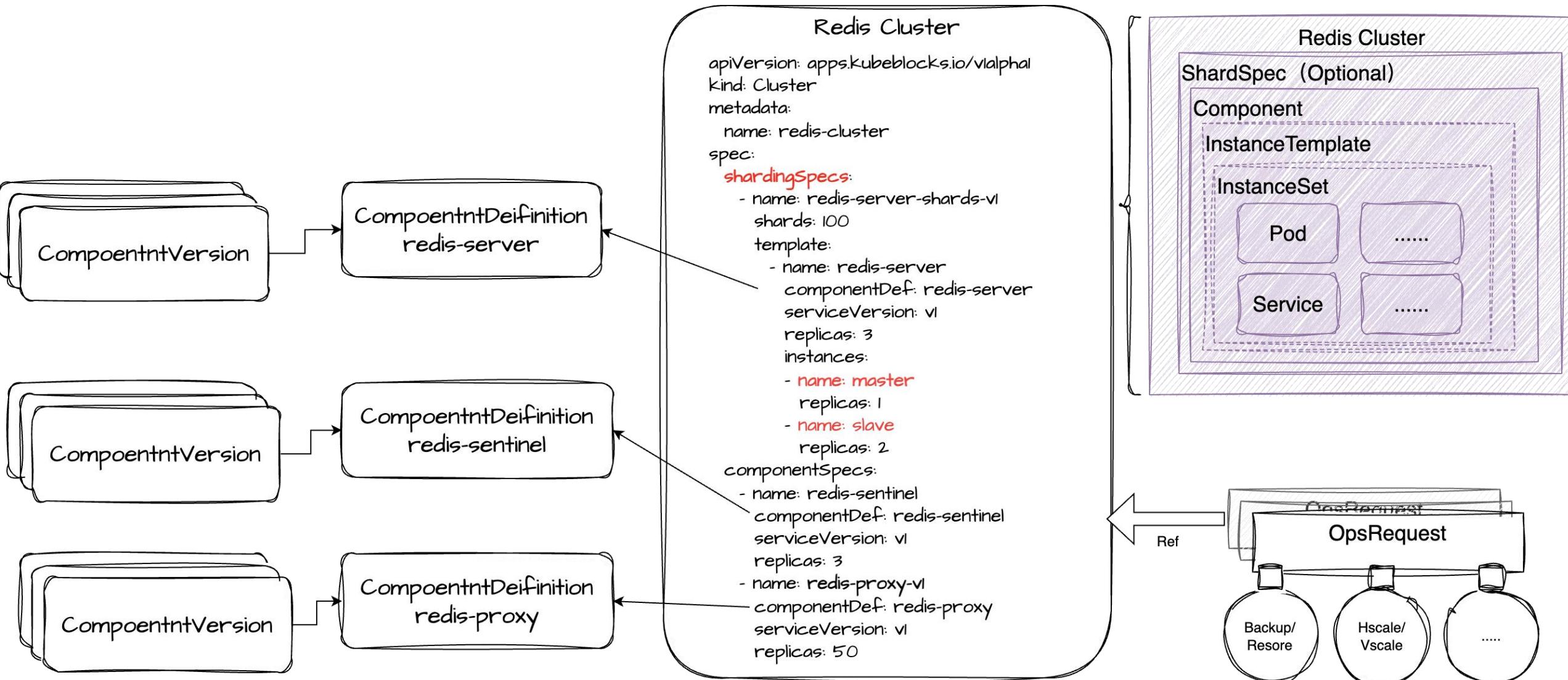
KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024



角色 (Role) 管理



KubeCon



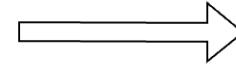
CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMIT

AI dev
Open Source Dev & ML Summit

China 2024

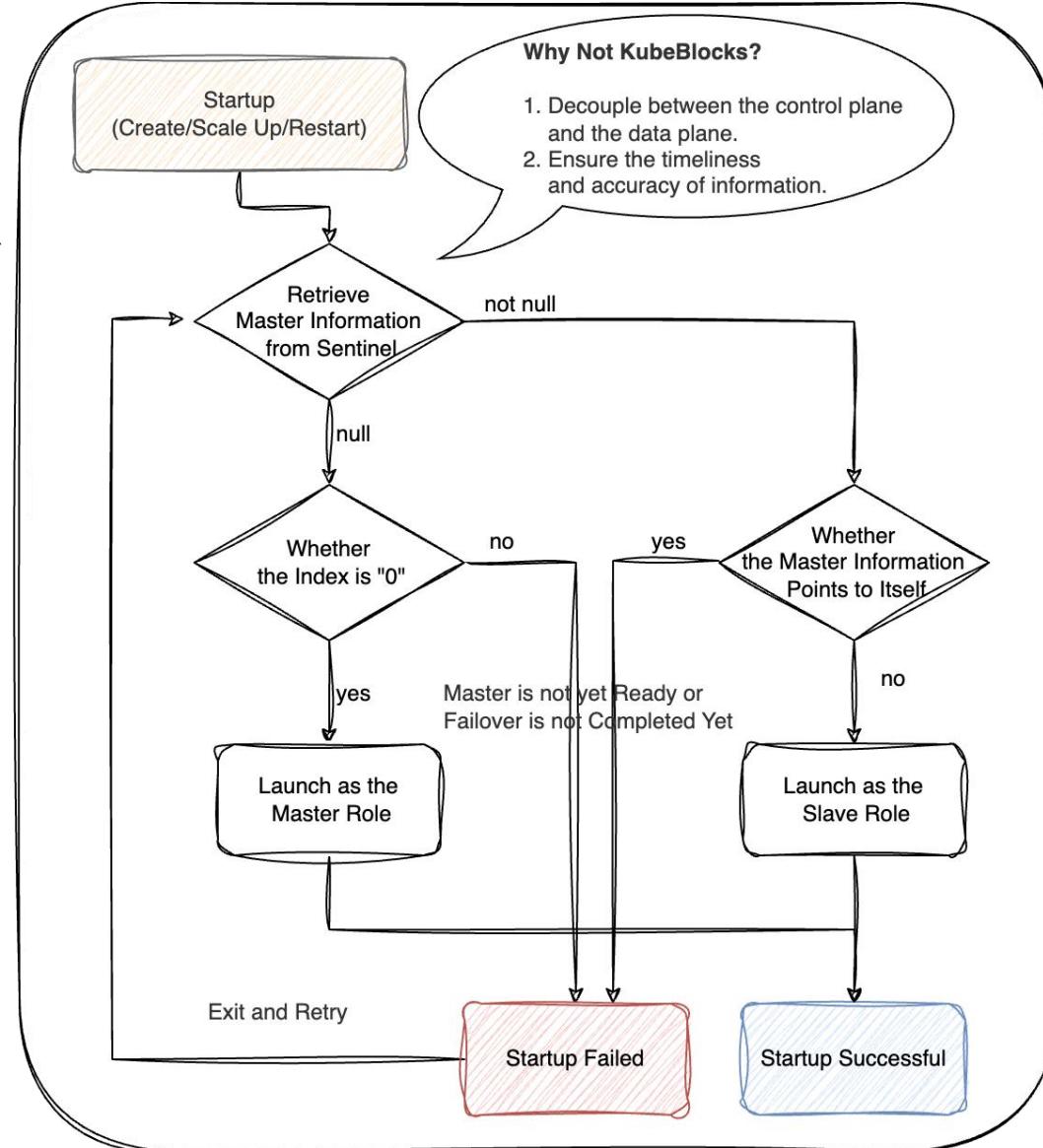
1. How to maintain the correct role relationships?



2. How to implement Role-Based O&M Management?



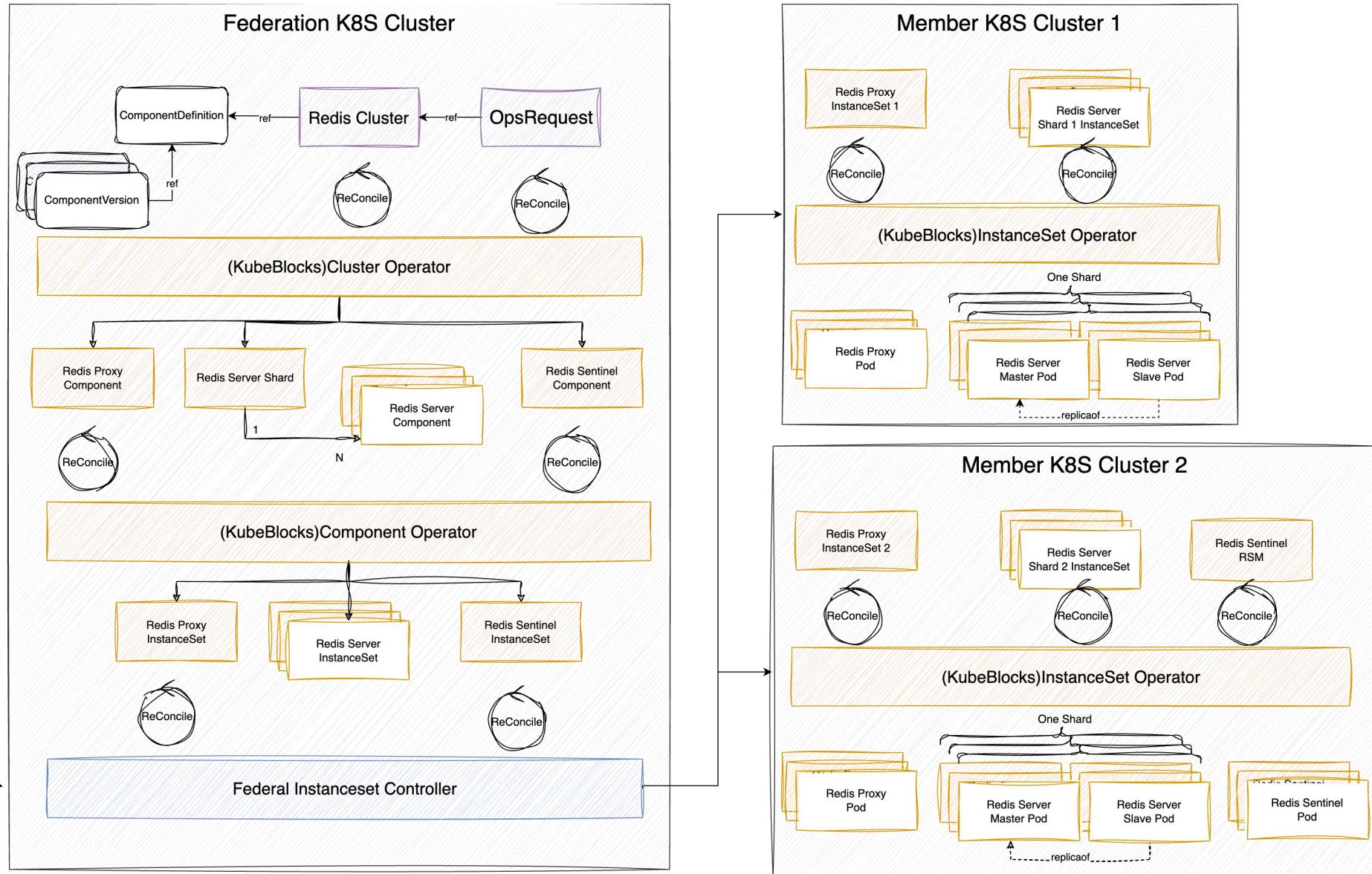
KubeBlocks has implemented it 👍



部署架构



China 2024



多集群分发管理



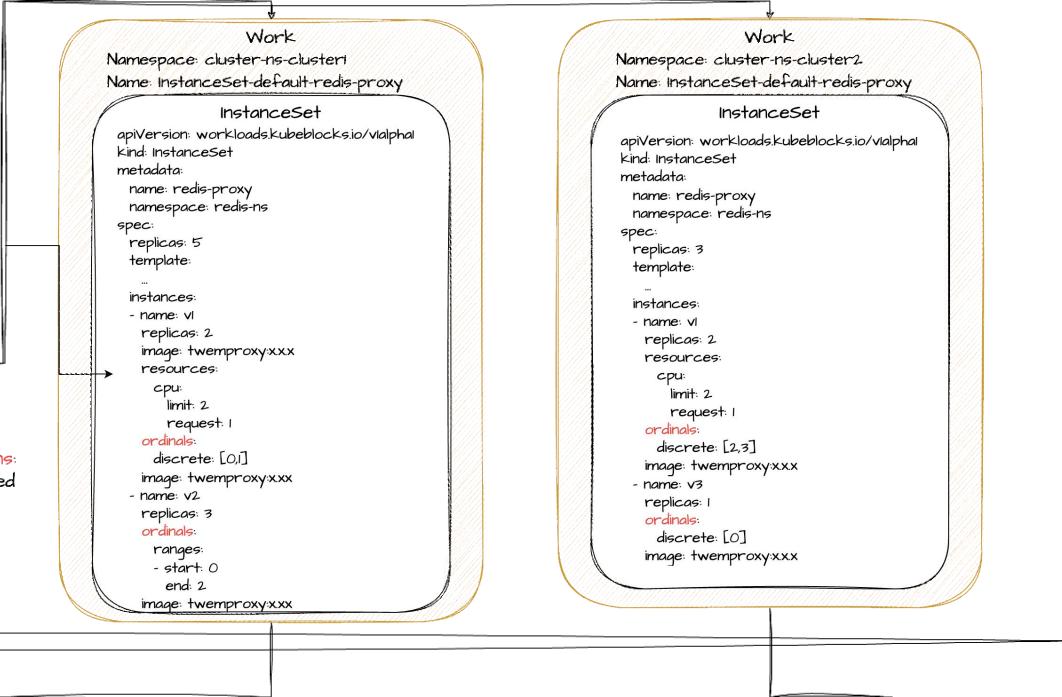
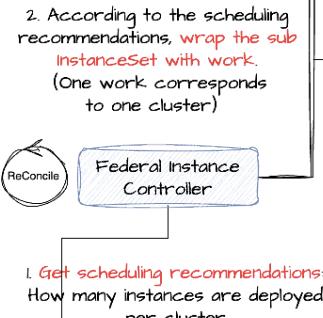
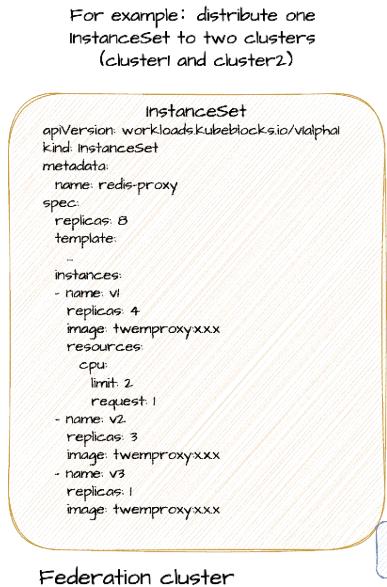
KubeCon



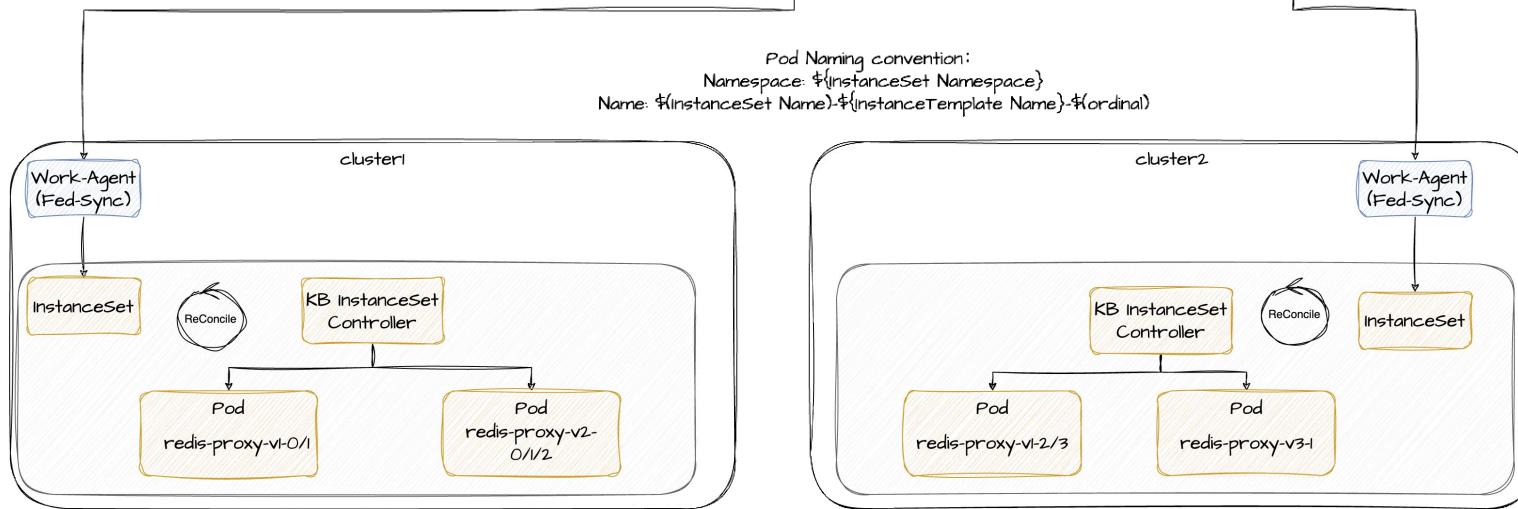
CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024



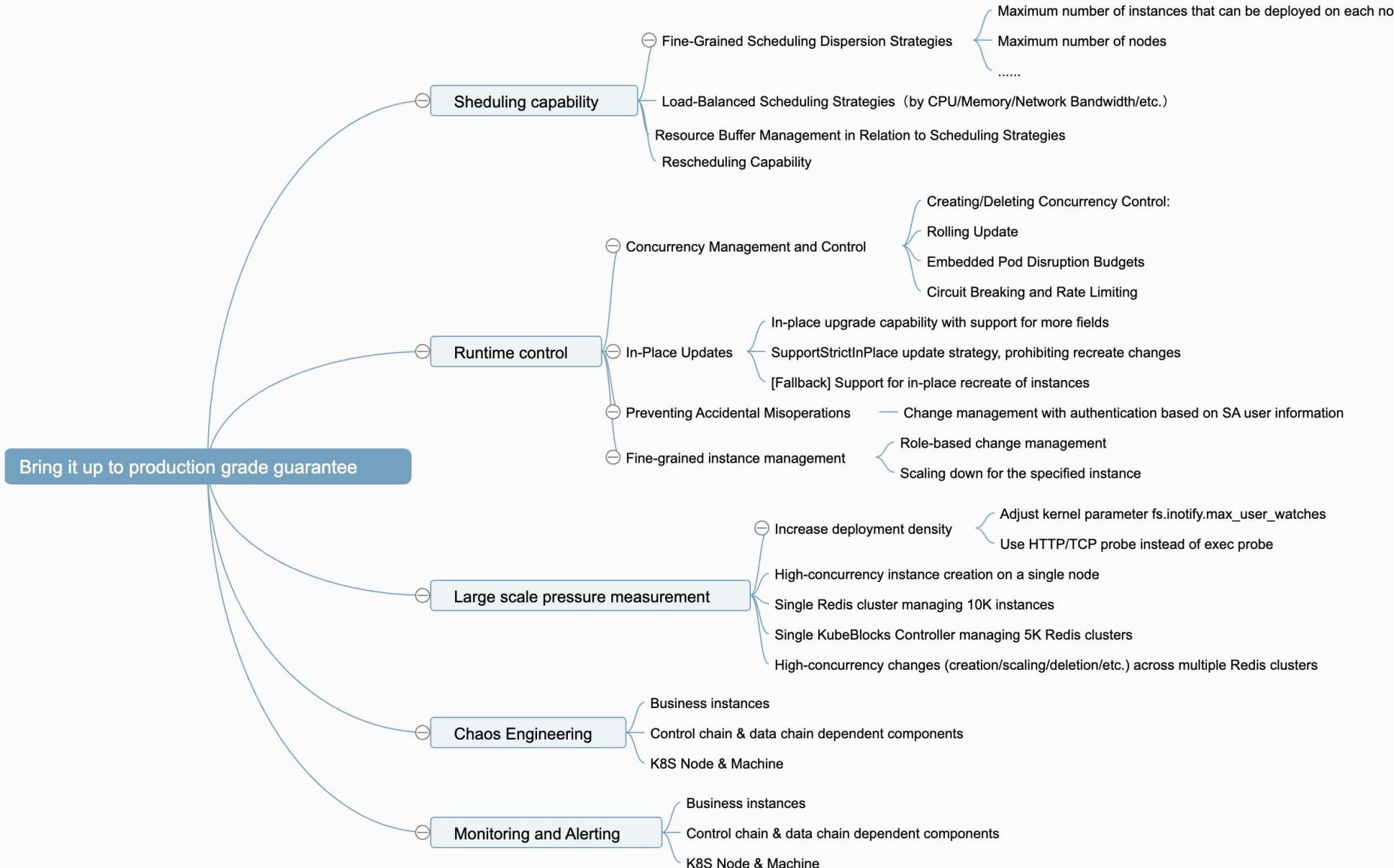
Member Clusters



稳定性保障



China 2024



关于 KubeBlocks 的一些思考

**The control plane for your
cloud-native data infrastructure**

Install, create, connect, and you have it all.



与 KubeBlocks 合作



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



AI Dev
Open Source Dev & ML Summit

China 2024

- InstanceSet 直管 Pod 和 PVC
- 实例模板 Instance Template(former heterogeneous pod)
- 与联邦集群集成
- 并发控制策略 (Parallel concurrency policy)
- 严格更新策略 (Restrict update policy)
- In-place update 和 in-place vertical scaling
- 指定实例下线
- 大规模/高并发场景性能优化
- ...



微信



Slack

加入 KubeBlocks 社区!

<https://github.com/apecloud/kubeblocks>



快手微信公众号

关注“快手技术”公众号
获取更多技术干货