

网易云如何基于大规模Kubernetes集群 支撑高并发应用

刘 超 网易云解决方案总架构师

容器的使用场景

大规模云原生应用的支撑痛点

Kubernetes的性能问题

公有云的支撑痛点

Kubernetes的规模问题

主办：



容器化的本质？

容器

虚拟机

我秒级启动

你没算应用吧，我也秒级

我一台机器上启动1000个nginx

然并卵

我有镜像，原子升级回滚

说的谁好像没有似得

我秒级自修复

你确认这一秒没丢啥

我服务发现

Dubbo和springcloud笑而不语，熔断，降级

我弹性伸缩

你听说过autoscaling group么？

主办：



容器化的本质?

如果是个传统应用，启动慢，进程少，不更新，IaaS就够了

主办：



如果你遇到了以下问题

变化快

扛不住

拆

微服务

主办：



微服务

扛不住，进程多

变化快，常更新

容器：100个进程，每天一个镜像

虚拟机：我有点大

虚拟机：可不可以不用镜像，Ansible也能部署

微服务

扛不住，进程多

变化快，常更新

开发：运维 10:1.5

运维：开发写完代码就不管了，这么多环境都是我的

容器镜像的本质：

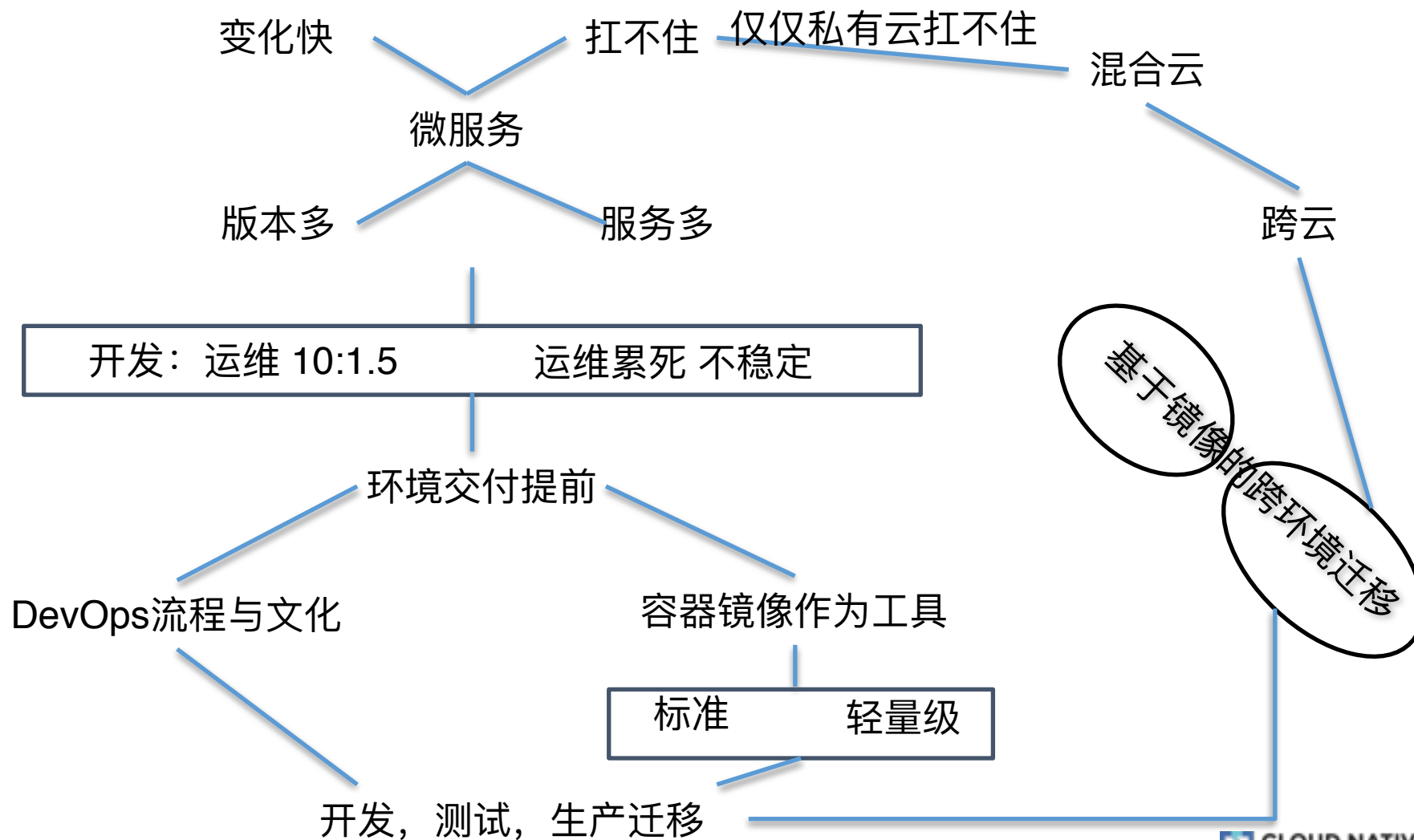
环境交付提前，每个研发5%的工作量，换取运维200%的工作量，增加稳定性

DevOps文化与流程

主办：



容器的本质?



主办：



容器的使用场景

大规模云原生应用的支撑痛点

Kubernetes的性能问题

公有云的支撑痛点

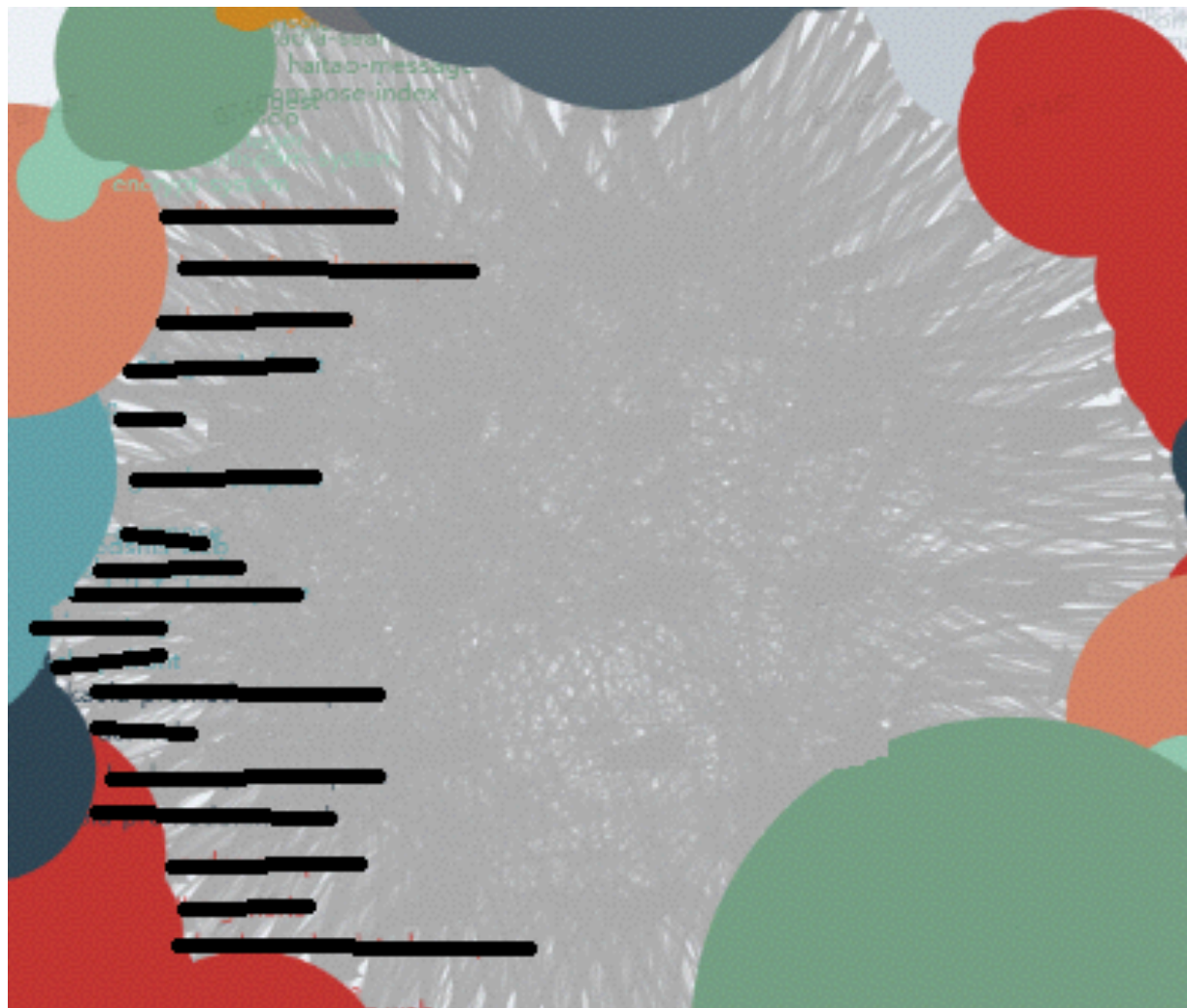
Kubernetes的规模问题

主办：



2017
China Kubernetes
End User Conference
kubernetes 中国用户大会 —2017—
2017.10.15 / 中国·杭州

主办：



- 入口：高性能负载均衡器
- 互访：高性能私有网络
- 访问PaaS平台：容器网络和虚拟机网络互通
- 高可用：高性能机房网络
- 服务发现：灵活策略，熔断，降级
- 配置复杂：统一配置中心

容器的使用场景

大规模云原生应用的支撑痛点

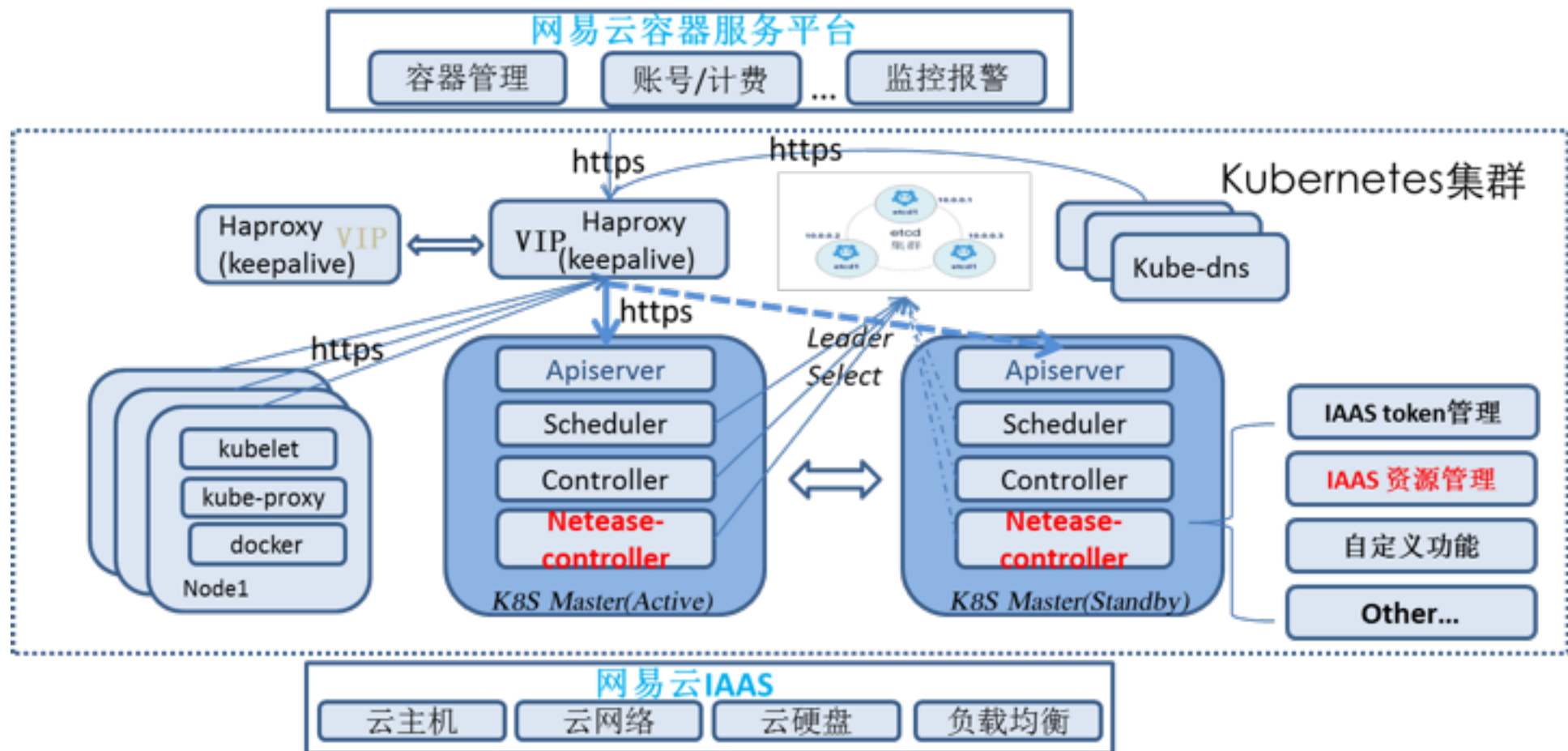
Kubernetes的性能问题

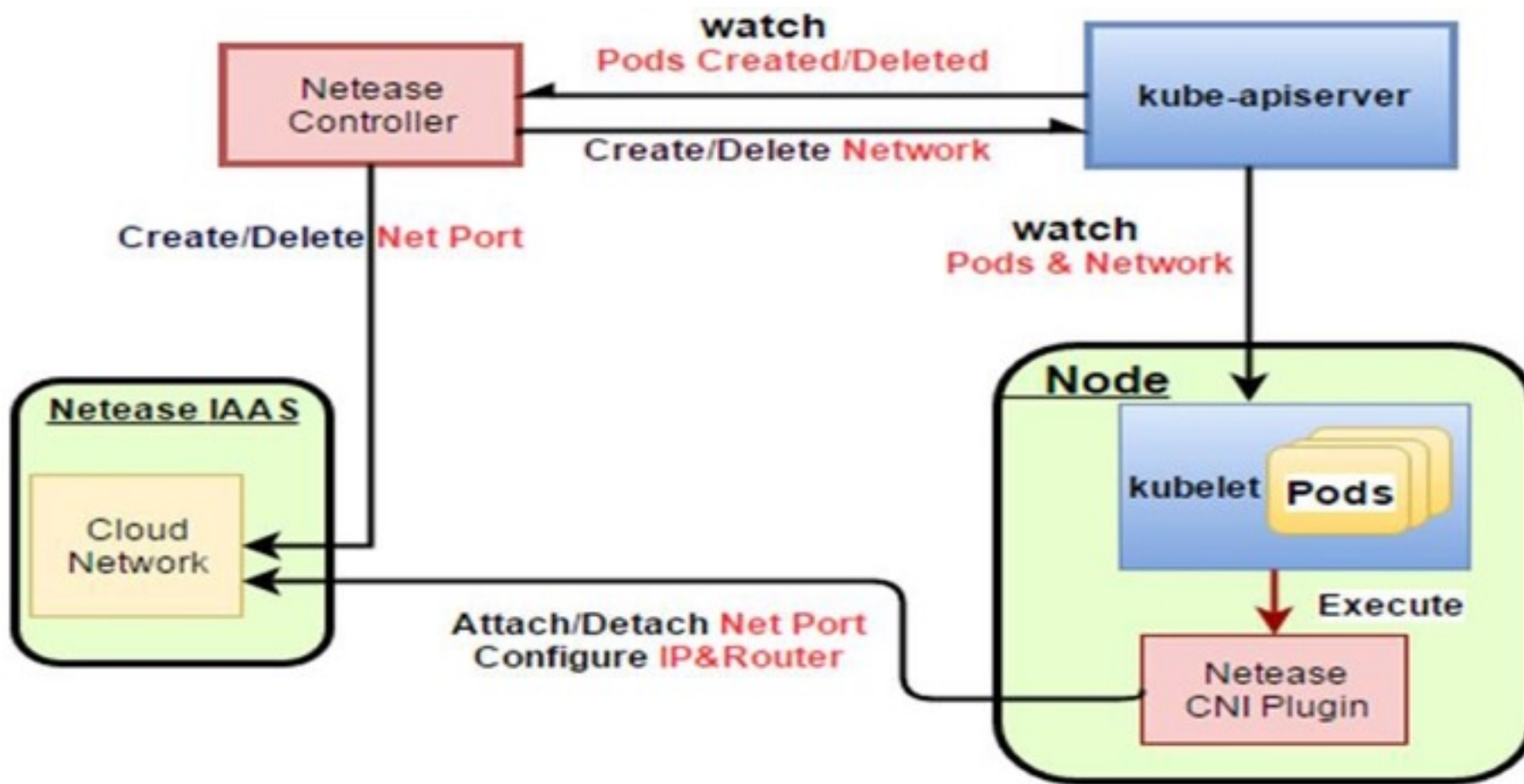
公有云的支撑痛点

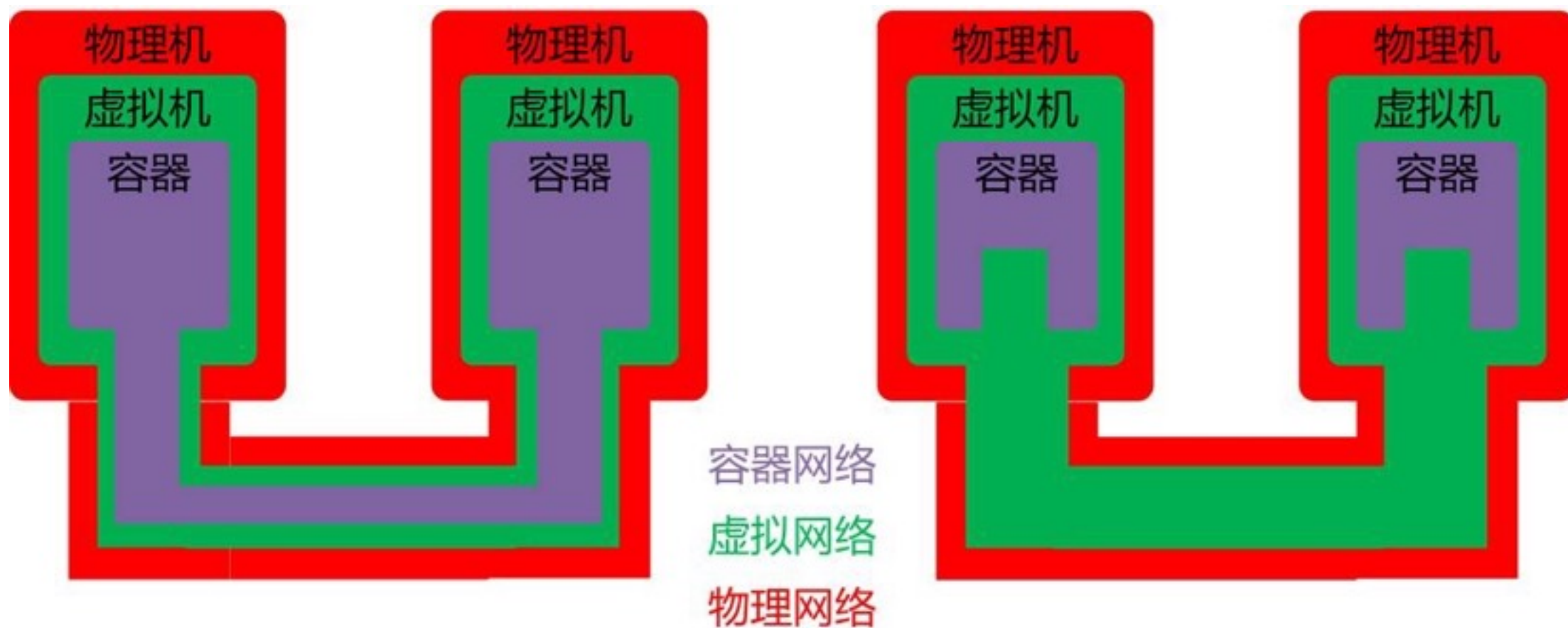
Kubernetes的规模问题

主办：



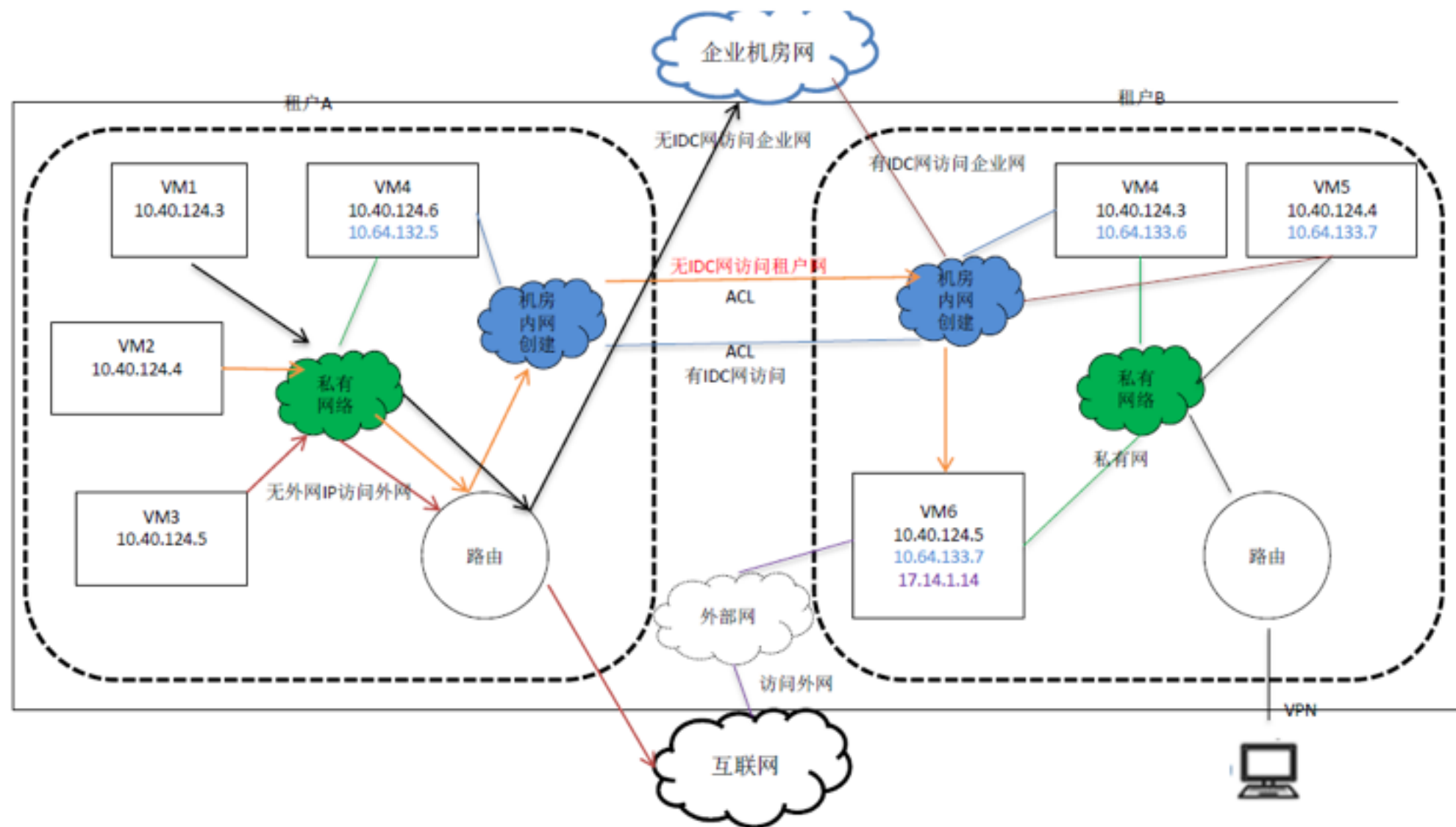


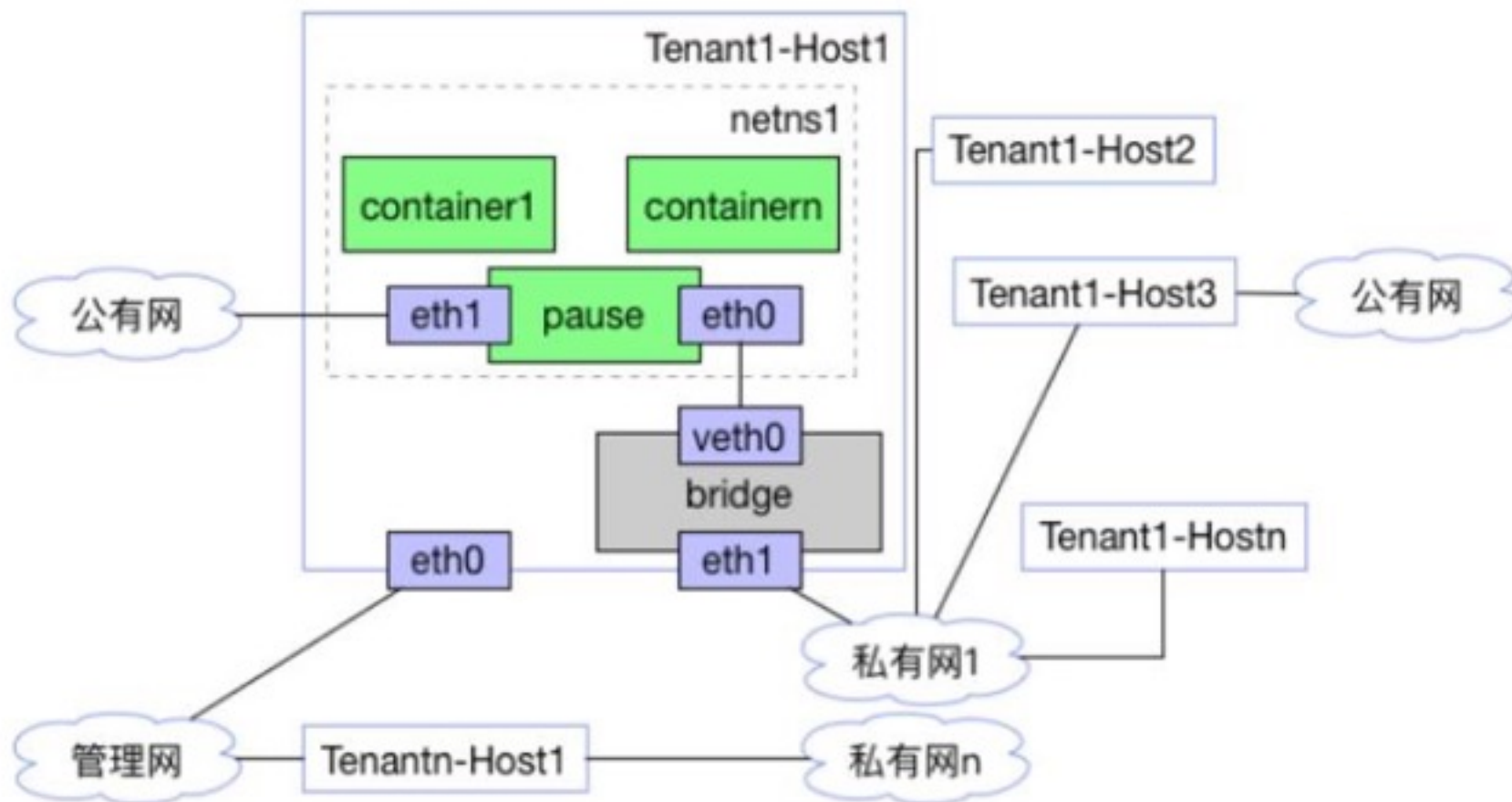


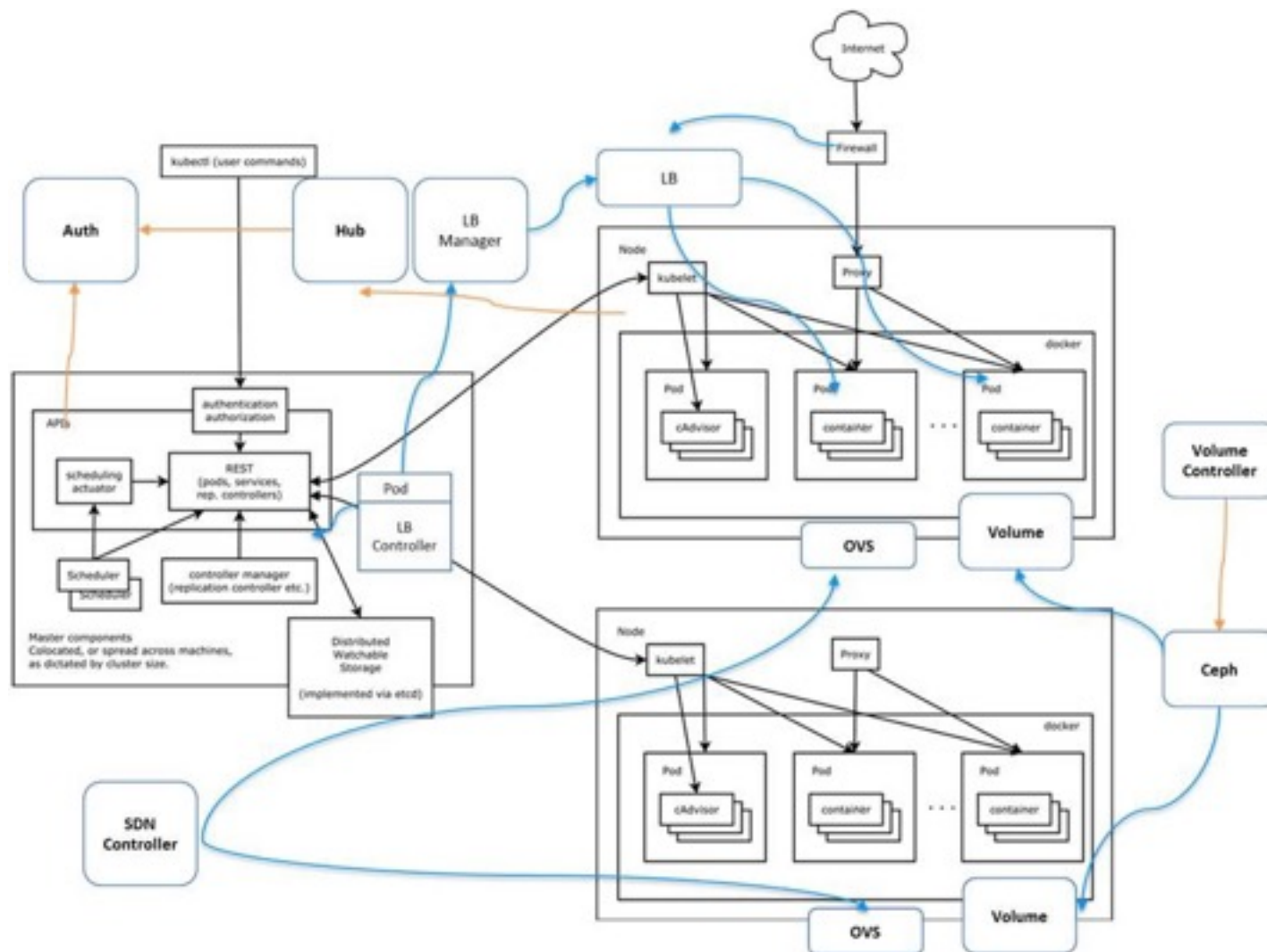


主办：



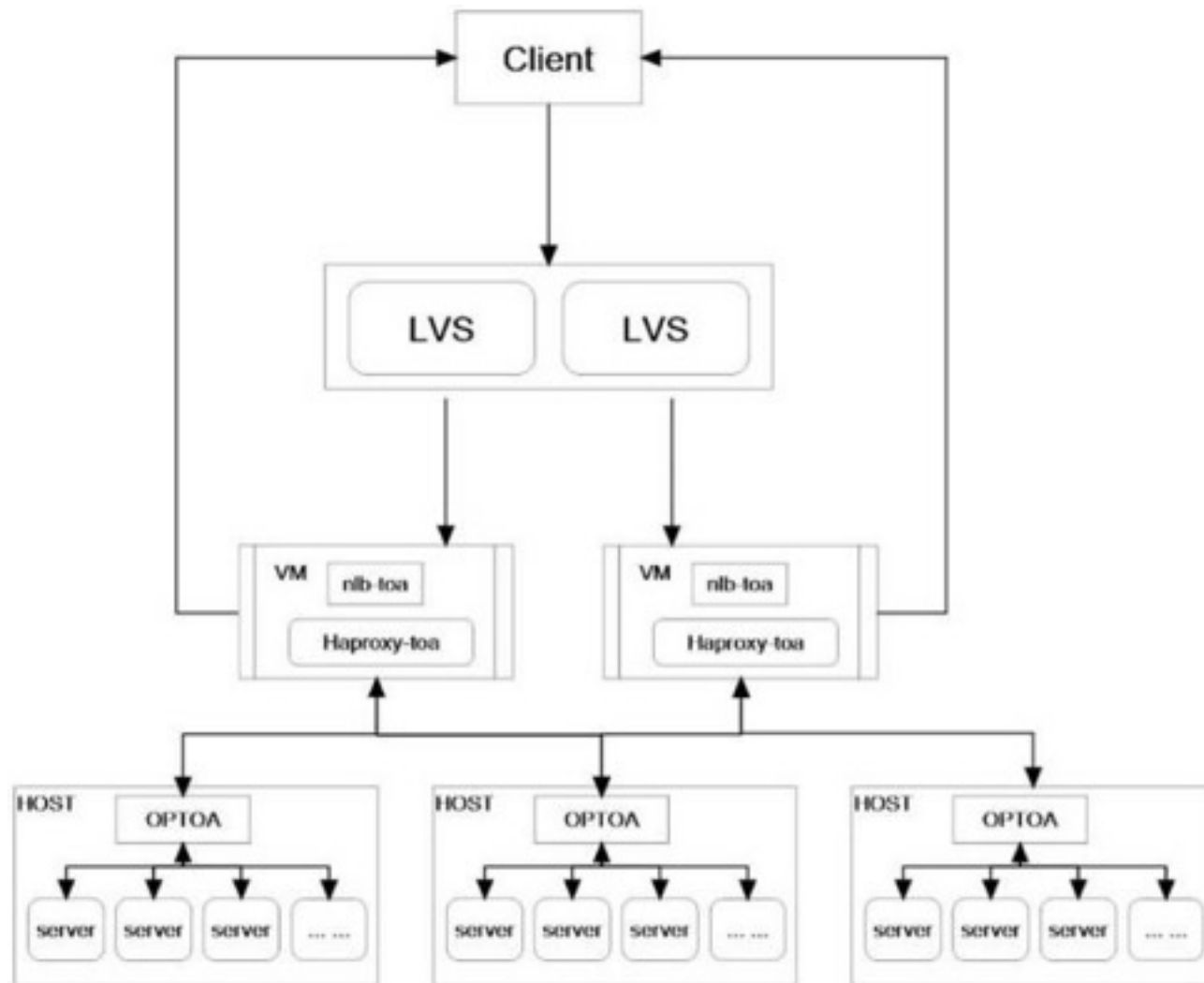




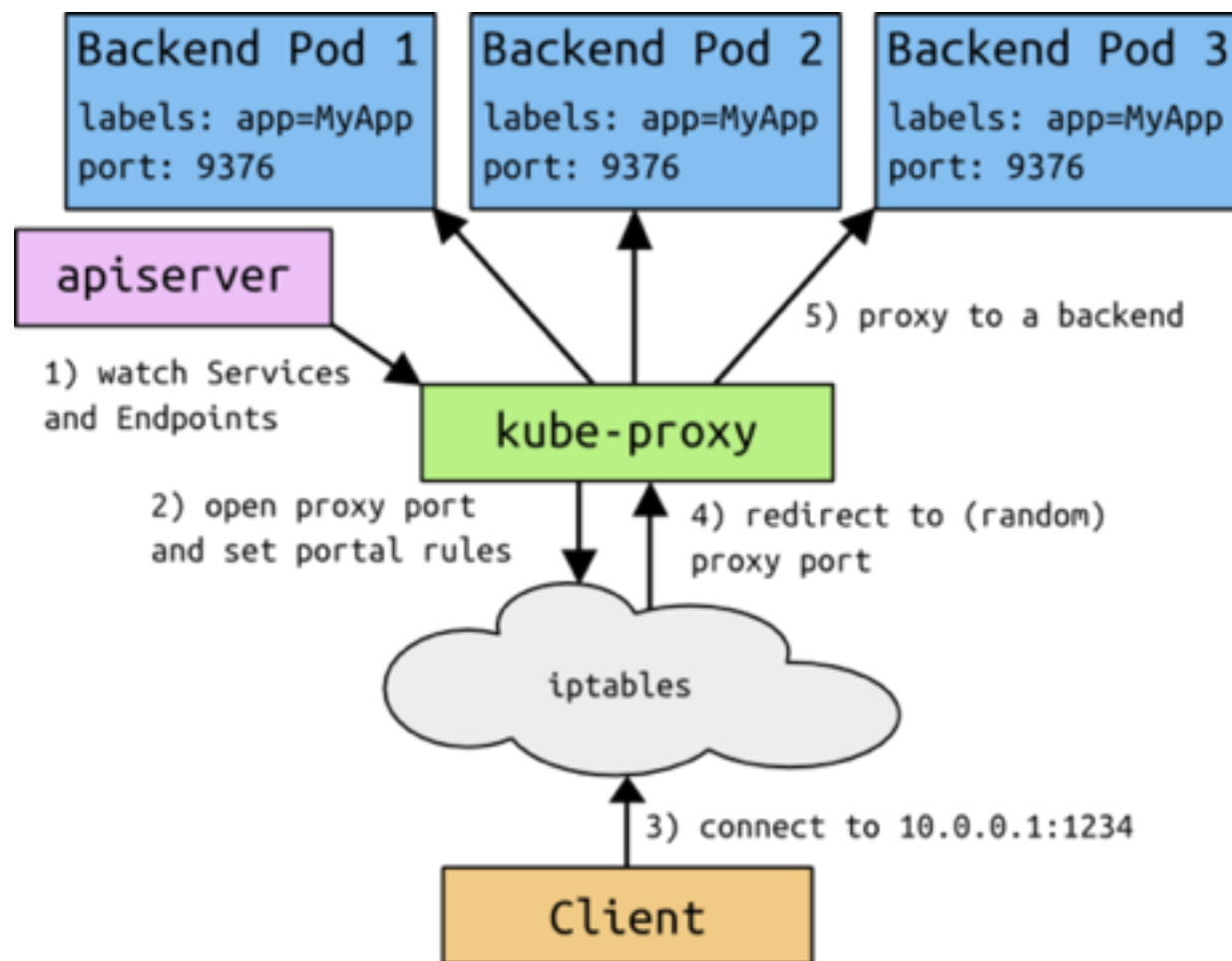


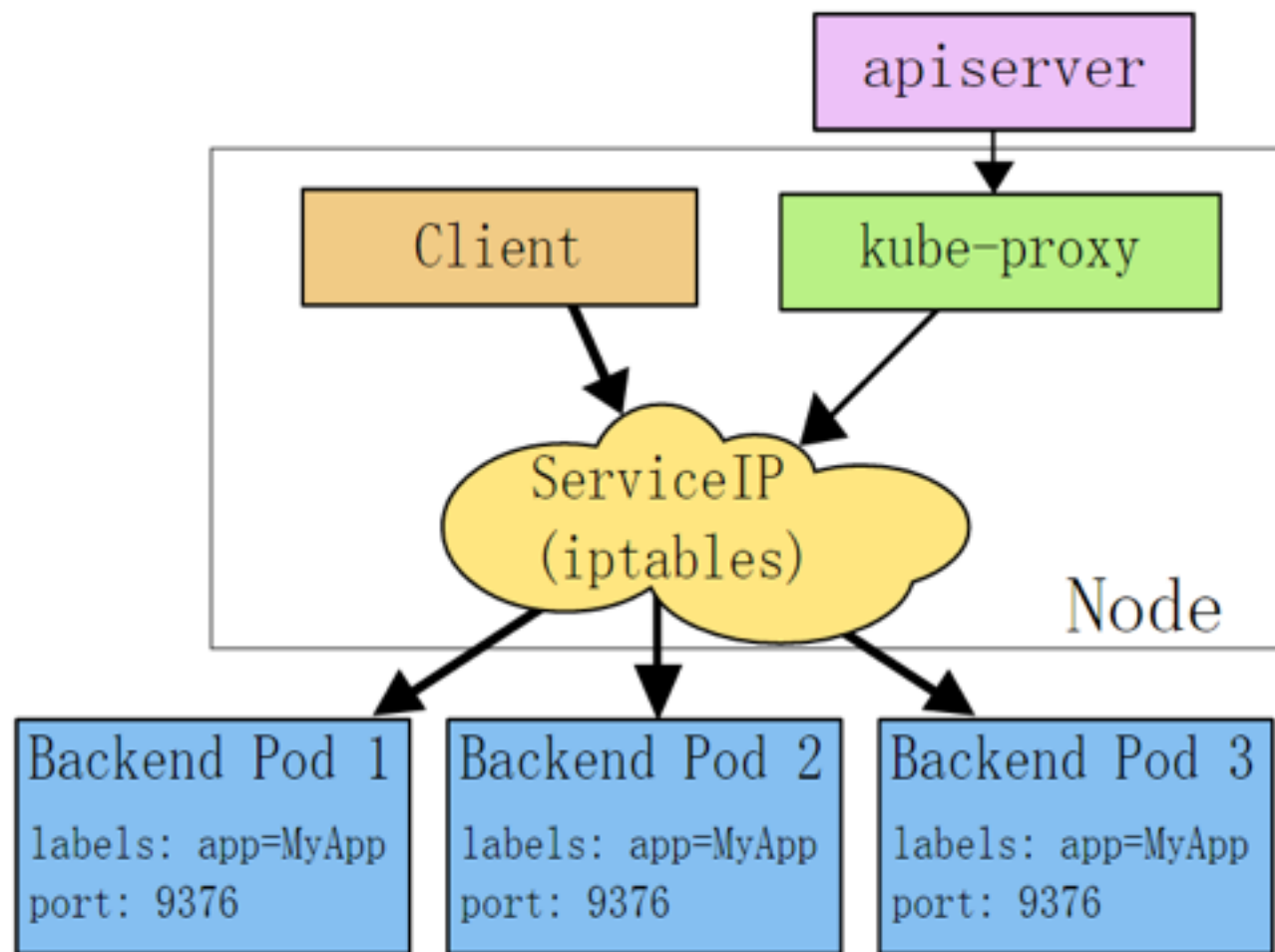
主办:





主办:





租户之间容器网络完全隔离，
无需配置多余的转发规则

只watch本租户的Service，生
成iptables 规则



摘要：微服务化是当前电商产品演化的必然趋势，网易考拉海购通过微服务化打破了业务爆发增长的架构瓶颈。本文结合网易考拉海购引用的开源Dubbo框架，分享支持考拉微服务工作的基本原理。文章分析了使用Dubbo过程中遇到的问题，讲解了团队所做的一些问题修复和功能集成工作，在此基础上最终形成了考拉内部持续维护升级的Dubbok框架。

本文背景还要从网易考拉海购（下文简称“考拉”）微服务化说起，现在任何大型的互联网应用，尤其是电商应用从Monolithic单体应用走向微服务化已经是必然趋势。微服务化是一个比较宽泛的概念，涉及到一个产品生命周期的多个方面，首先它作为一个指导原则指引业务划分、架构解耦等；技术层面实施微服务需要开发测试阶段、运行阶段、发布阶段、部署阶段等一系列基础框架的支撑。我们在享受服务化易扩展易部署等便利性的同时，也面临新的问题，如数据一致性、分布式调用链路追踪、异常定位、日志采集等。

Headless services

With selectors

For headless services that define selectors, the endpoints controller creates `Endpoints` records in the API, and modifies the DNS configuration to return A records (addresses) that point directly to the `Pods` backing the `Service`.

Without selectors

For headless services that do not define selectors, the endpoints controller does not create `Endpoints` records. However, the DNS system looks for and configures either:

- CNAME records for `ExternalName`-type services.
- A records for any `Endpoints` that share a name with the service, for all other types.

客户端服务发现

简化外部配置

主办：



容器的使用场景

大规模云原生应用的支撑痛点

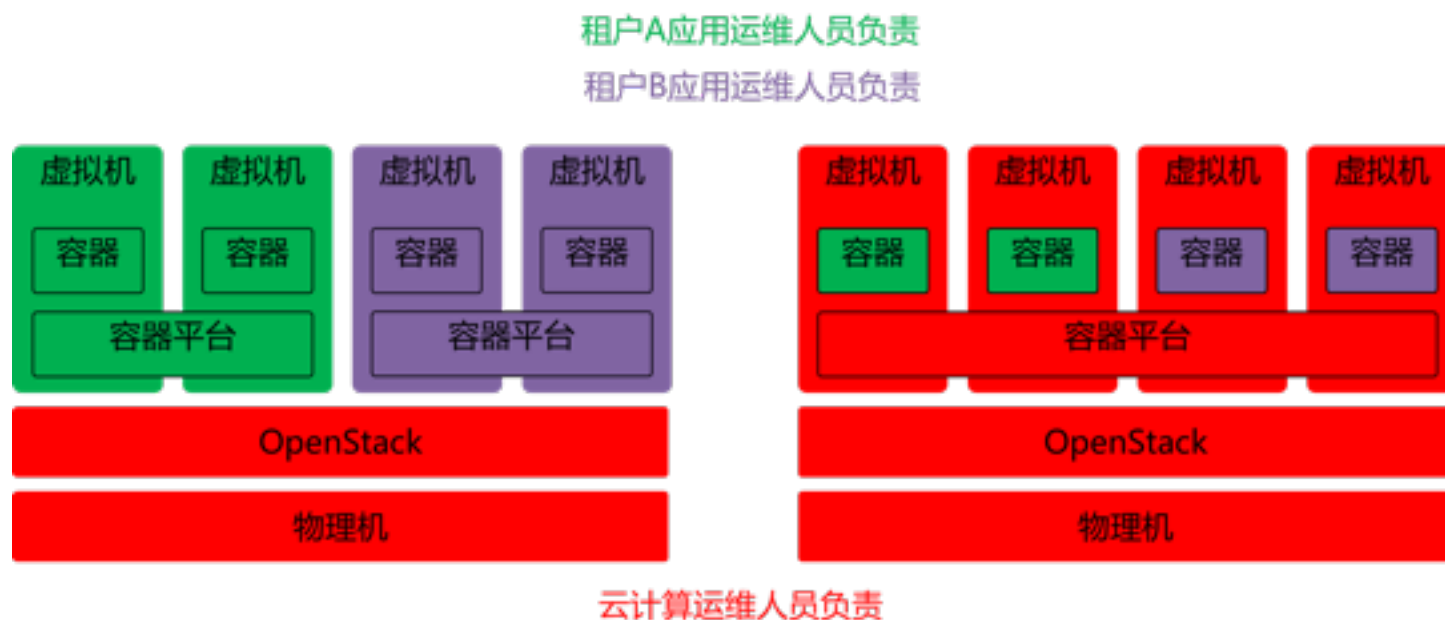
Kubernetes的性能问题

公有云的支撑痛点

Kubernetes的规模问题

主办：





- 一个集群 or 多个集群
- Docker内核隔离不好
- Node无租户隔离
- 虚拟机启动速度问题
- Kube-proxy转发规则太多
- 集群规模问题

容器的使用场景

大规模云原生应用的支撑痛点

Kubernetes的性能问题

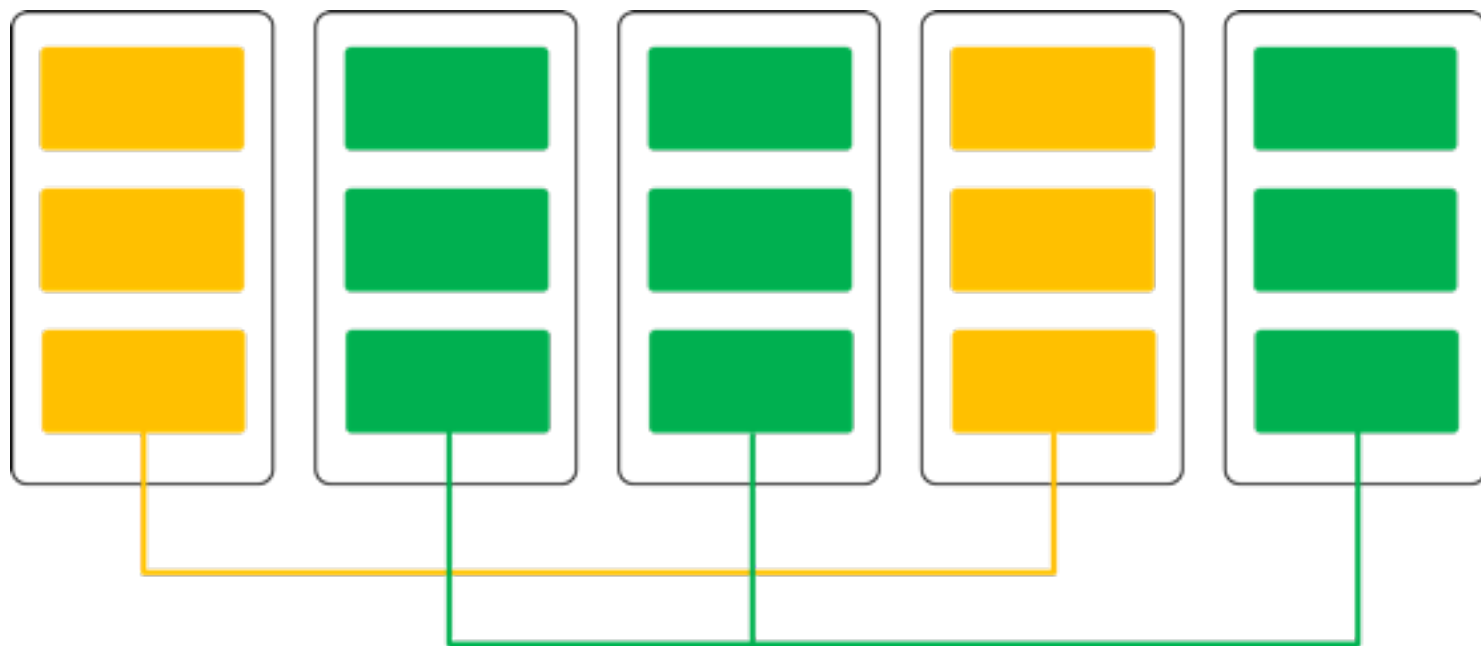
公有云的支撑痛点

Kubernetes的规模问题

主办：

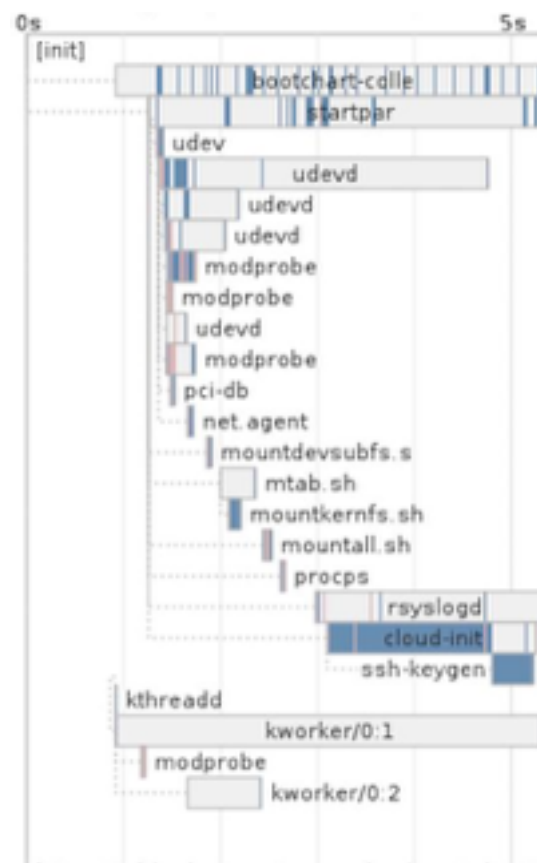
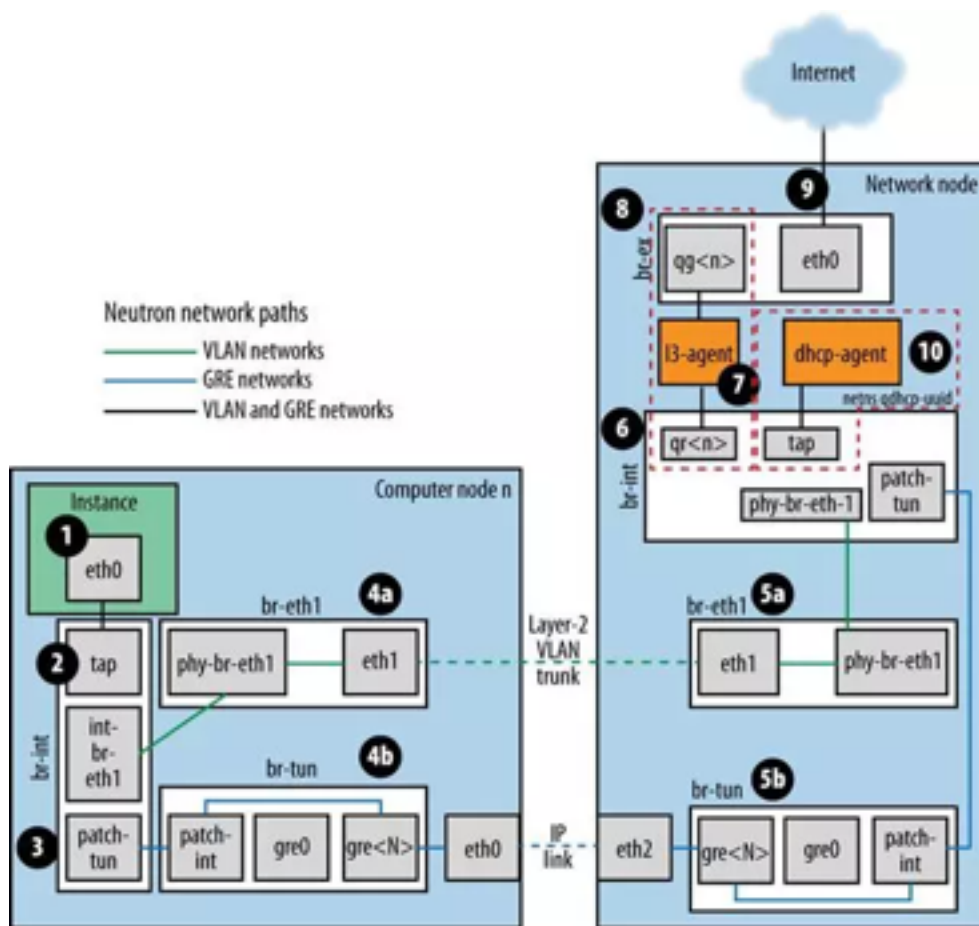


Kubernetes的kubelet运行容器



主办：





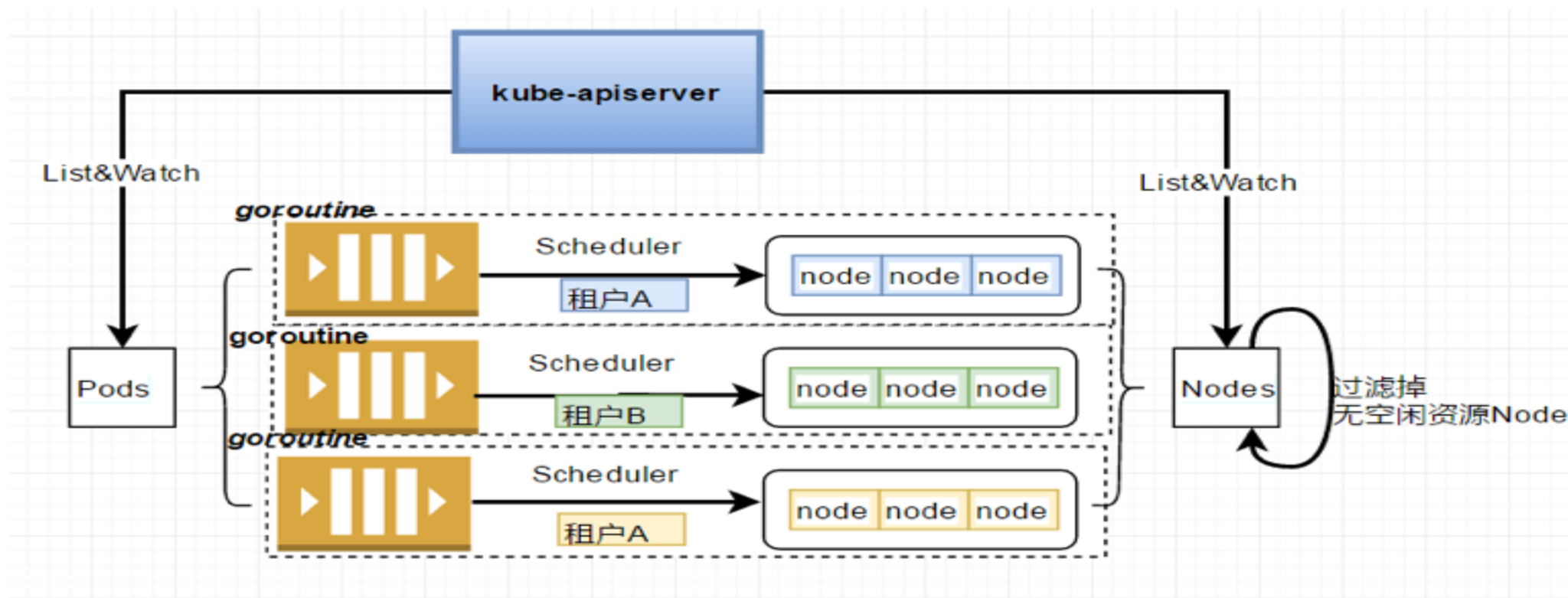
80秒

5秒

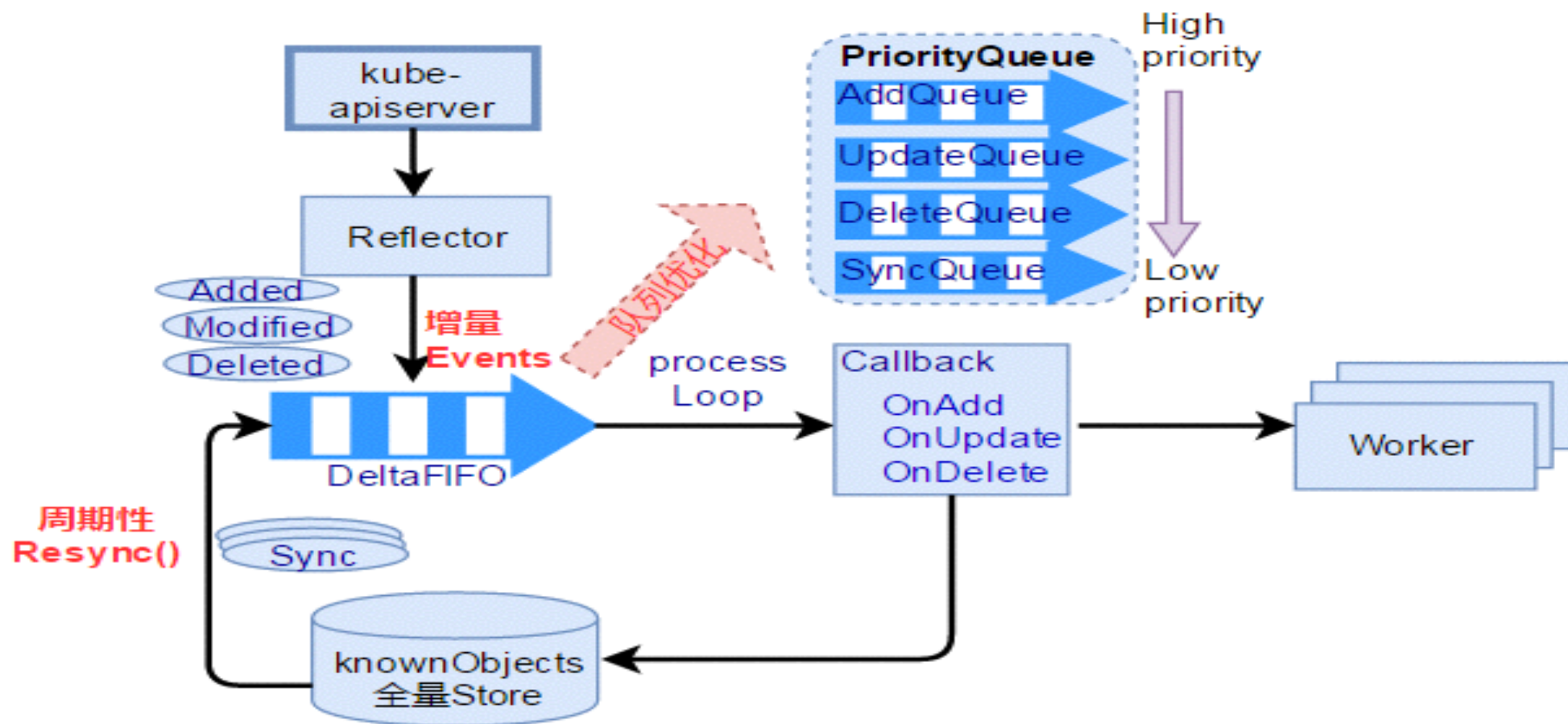
- 1.网卡IP初始化
- 2.网络路由注入
- 3.DNS服务IP配置
- 4.网卡udev规则

主办：



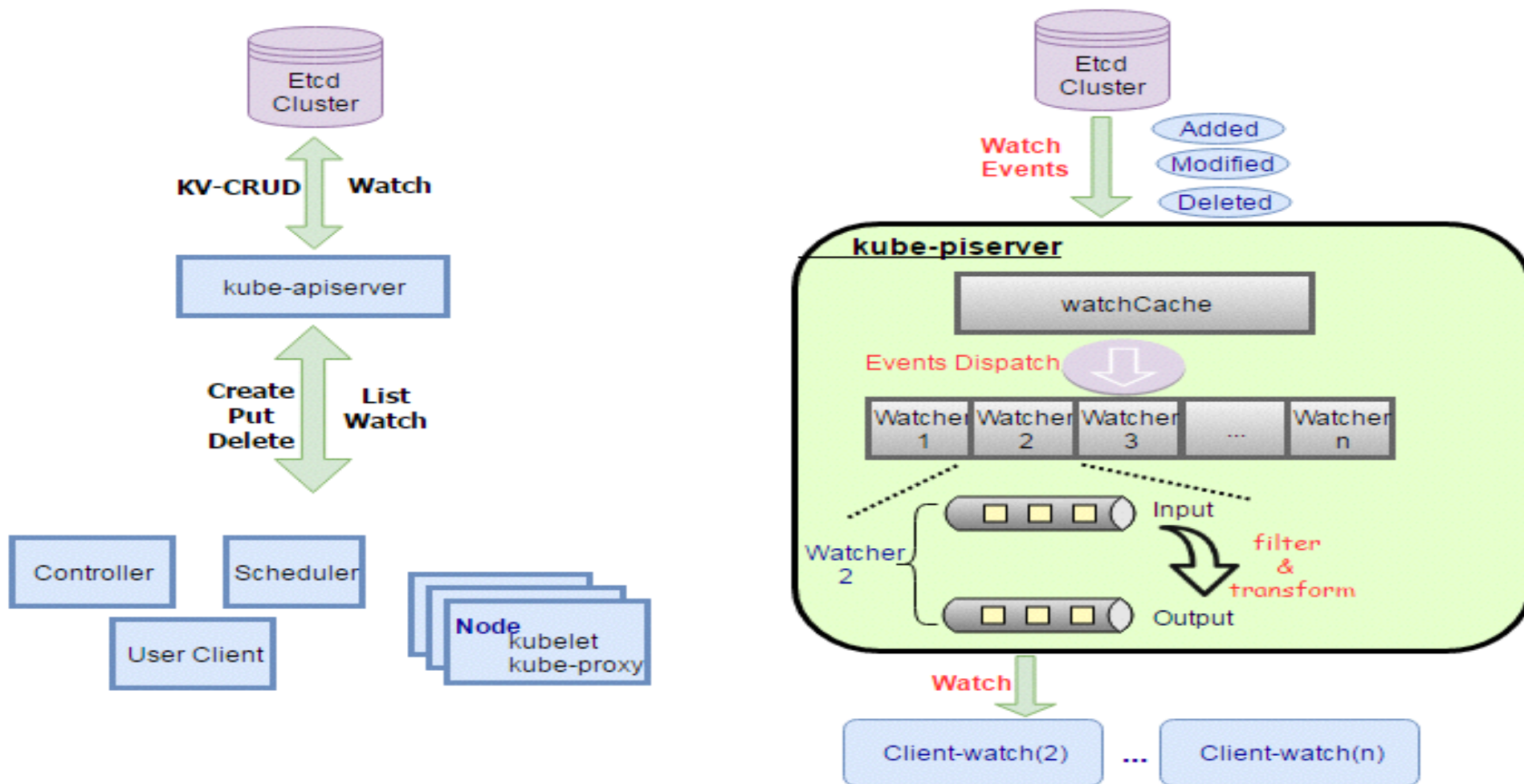


公有云租户间资源完全隔离，天然适合并行调度
预过滤无空闲资源的node
调整predicate调度算法快速过滤



按事件类型进入多优先级workqueue (Add > Update > Delete > Sync)

主办:



Apiserver 就是一个proxy代理，而且goroutine对web服务完美支持，最终性能瓶颈在对 etcd 的访问上

主办：



Q&A

Thank you for your time

主办：

