



基于腾讯云的kubernetes实践

包健@腾讯云

为什么选择K8S

三大开源容器编排工具在 [Github](#) 上的活跃数据(2016.9月)

技术平台	社区关注数	社区喜爱数	社区参与数
Kubernetes	1461	17416	5647
Swarm	339	4033	814
Mesos	53	387	72
Marathon (Mesos 依赖框架)	353	2187	630

➤ 开源、社区活跃度最高

➤ 对公有云的支持完善，像volume, loadblance、router等都有可以通过plugin来提供

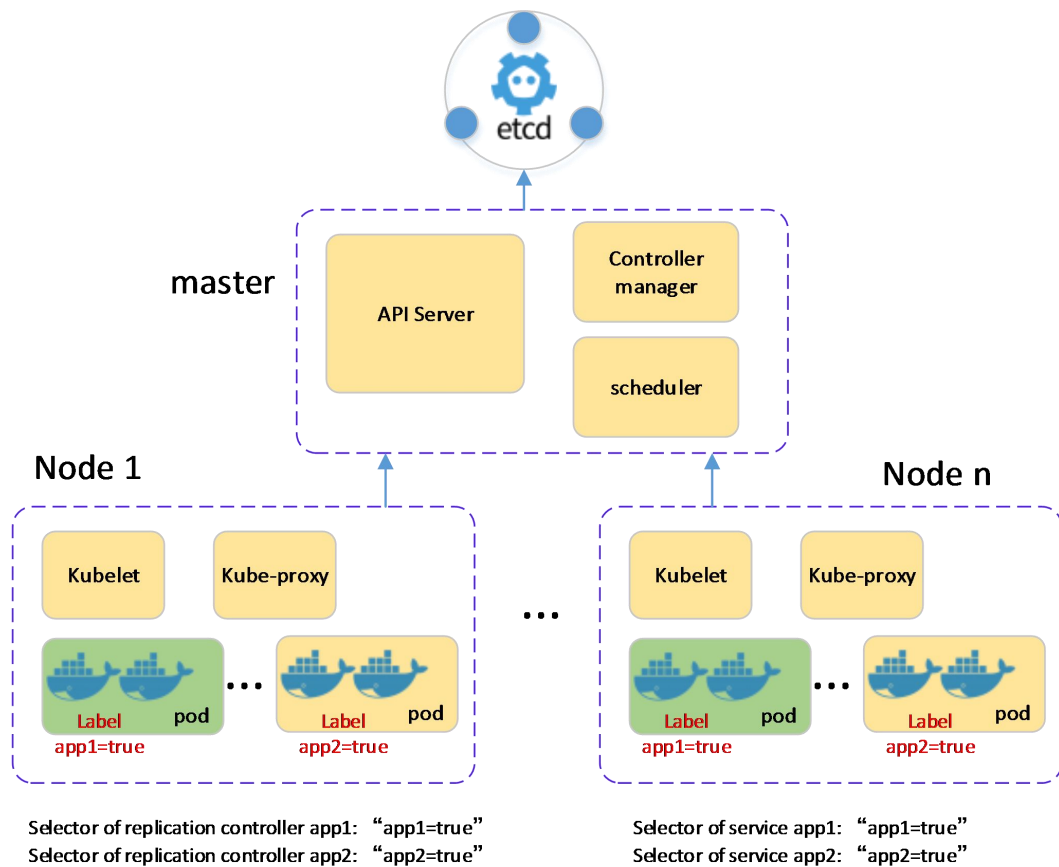
➤ 为分布式而生，相关概念完备，如服务发现、健康检查、ingress、命名空间等

➤ 使用人数多，70%以上的企业客户在使用 kubernetes 来搭建容器服务



相比之下，Docker swarm 的背后厂商以 Docker 公司自己为主，而 Mesos 则以 Apache 开源基金会为主。

k8s简介



Kubernetes常用名词	说明
pod	容器组, Kubernetes的基本资源
replication controller	Pod的副本控制器
service	Pod的流量访问控制器
label	标签
node	集群节点, 运行pod
master	集群master, 负责资源的管理、调度
etcd	集群资源持久化db



01

控制台

- a、集群管理：用户通过控制台管理集群，不需手动配置搭建
- b、封装服务、配置项、pod等概念封装到控制台
- c、提供服务事件、容器日志和webshell等工具

02

对接腾讯云IAAS层

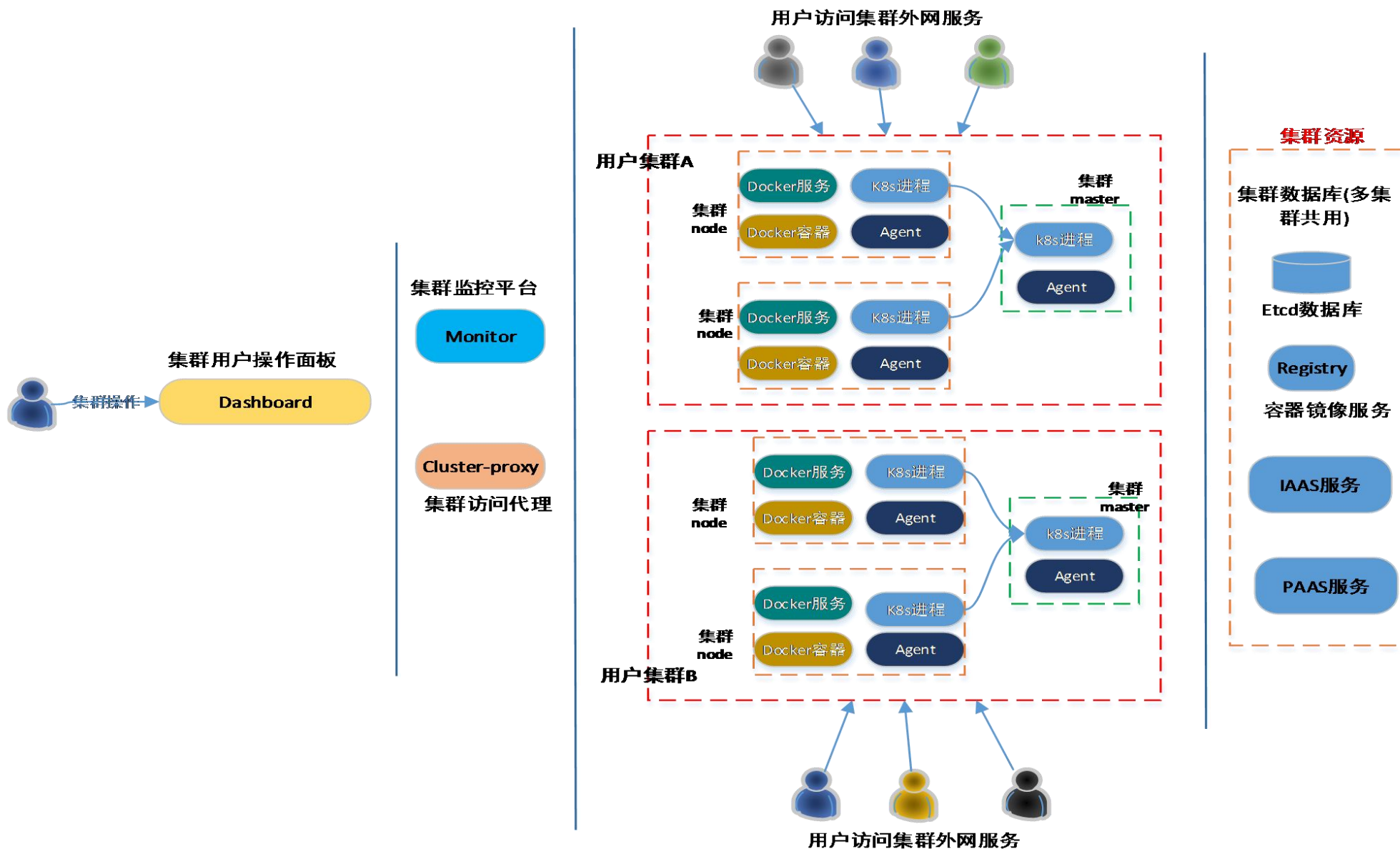
- a、基于VPC的容器网络，可和外部CDB等资源打通
- b、支持cbs盘的volume和storage class
- c、基于腾讯云LB的ingress controller

03

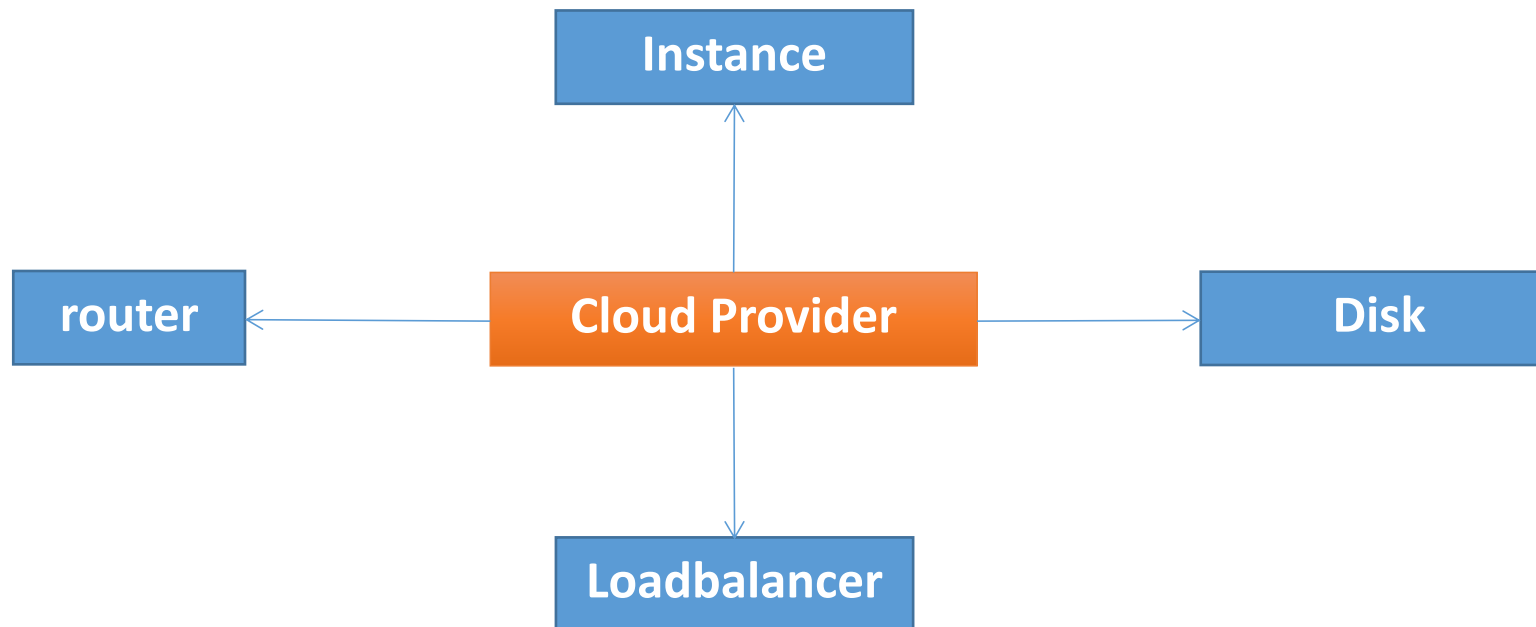
镜像仓库

- a、仓库优化，包括添加缓存层以及和cos存储的对接
- b、CI/CD：代码源对接github，通过镜像触发器可与企业自建CI流程打通
- c、提供mirror registry

总体架构



k8s cloud provider

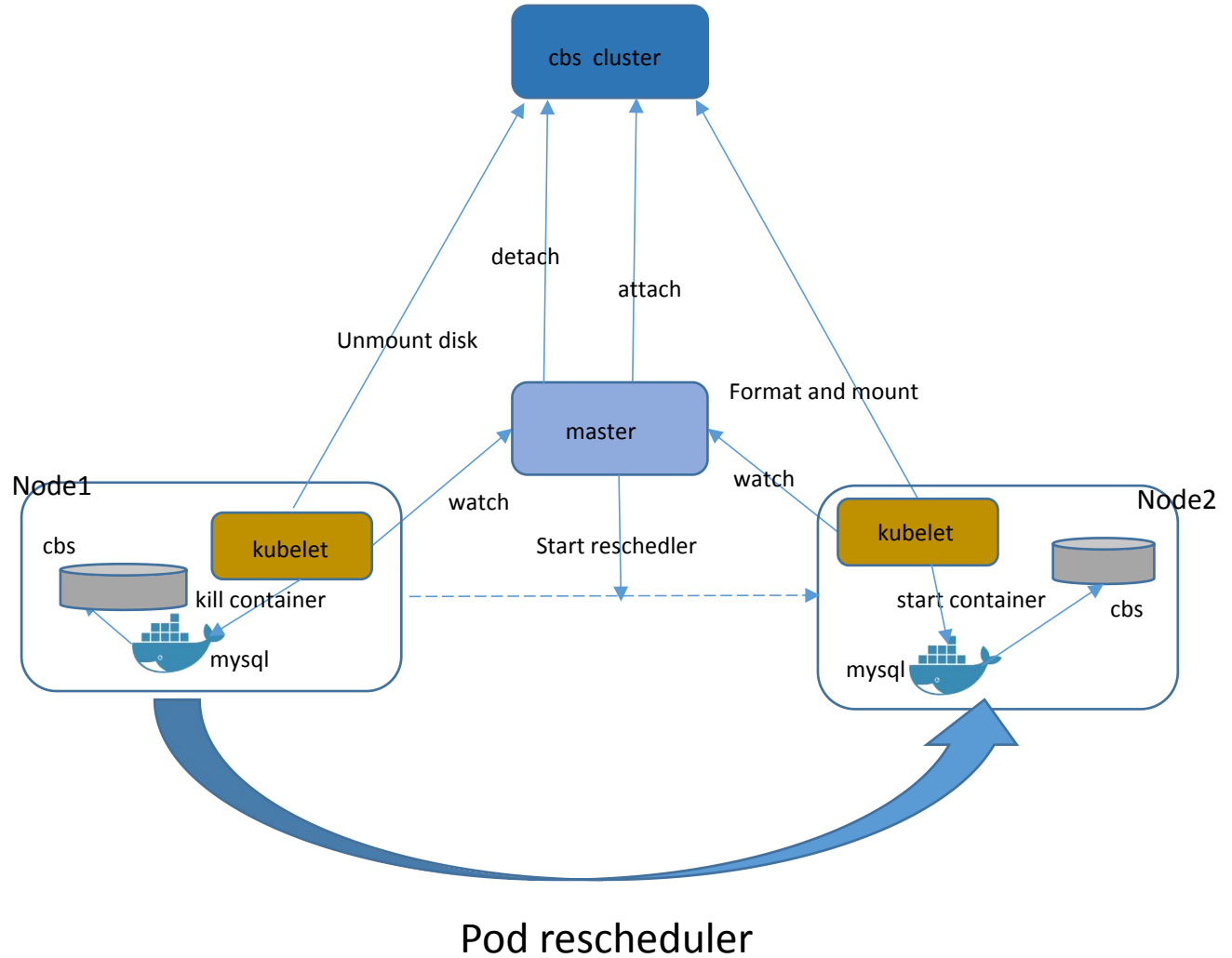


- router: 集群新增节点时添加容器子网在哪台VM上的路由信息
- instance: 云主机相关信息查询，在 `kubectl describe node` 中可看到
- loadbalancer: 当 `service` 的 `spec.type: LoadBalancer` 时使用
- Disk: cbs 盘 `attach`、`detach` 的接口实现，供 `volume` 和 `storage class` 使用
- 使用 `cloud provider` 时，`apiserver`，`controller-manager` 和 `kubelet` 的 `--cloud-config` 和 `--cloud-provider`，需配置成对应云厂家

支持cbs的volume

```
spec:
  containers:
  - image: nginx
    imagePullPolicy: Always
    name: nginx
    resources:
      requests:
        cpu: 200m
    securityContext:
      privileged: false
    terminationMessagePath: /dev/termination-log
    terminationMessagePolicy: File
    volumeMounts:
    - mountPath: /mnt
      name: test
  dnsPolicy: ClusterFirst
  imagePullSecrets:
  - name: qcloudregistrykey
  restartPolicy: Always
  schedulerName: default-scheduler
  securityContext: {}
  terminationGracePeriodSeconds: 30
  volumes:
  - name: test
    qcloudCbs:
      cbsDiskId: disk-echj2drf
      fsType: ext4
```

使用cbs的yaml文件



cbs storage class

```
apiVersion: storage.k8s.io/v1beta1
kind: StorageClass
metadata:
  annotations:
    storageclass.beta.kubernetes.io/is-default-class: "true"
  name: cbs
parameters:
  type: cbs
provisioner: cloud.tencent.com/qcloud-cbs
```

使用cbs 的storage class

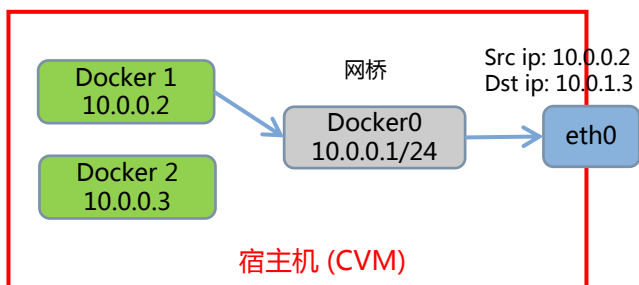
```
spec:
  terminationGracePeriodSeconds: 10
  serviceAccountName: mongo
  containers:
  - name: mongo
    image: mongo
    command:
      - mongod
      - "--replSet"
      - rs0
      - "--smallfiles"
      - "--noprealloc"
    ports:
      - containerPort: 27017
    volumeMounts:
      - name: mongo-persistent-storage
        mountPath: /data/db
  volumeClaimTemplates:
  - metadata:
      name: mongo-persistent-storage
      annotations:
        volume.beta.kubernetes.io/storage-class: "cbs"
    spec:
      accessModes: [ "ReadWriteOnce" ]
      resources:
        requests:
          storage: 10Gi
```

mongodb使用cbs storage class的yaml文件

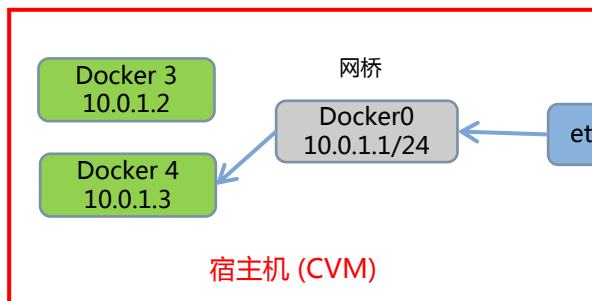
容器之间互访

宿主机 VPC网段 : 172.0.0.0/24
容器 VPC网段 : 10.0.0.0/14

CVM-1 :172.0.0.2 (母机IP 为 10.172.1.1)



CVM-2 172.0.0.3 (母机IP 为 10.172.1.2)



母机内核中 完成 路由查找, 封包解包过程

在 Node-1节点母机上 查询全局路由表 获取 10.0.1.1/24 网段容器 在 172.0.0.3 上

查询 vpc路由表 172.0.0.3 位于 母机 10.172.1.2 上

封包 并 发送请求到 10.172.1.2 机器

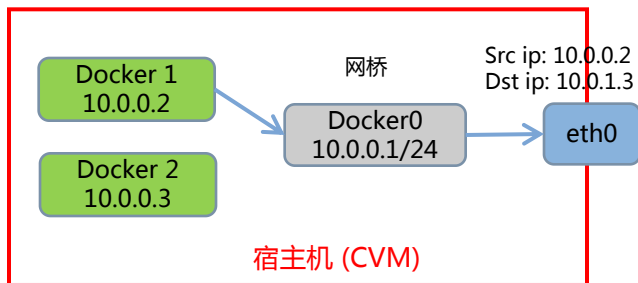
解包将请求发送 172.0.0.3 虚拟机

10.0.0.2
10.0.1.3
Vpc id
172.0.0.3
10.172.1.1
10.172.1.2

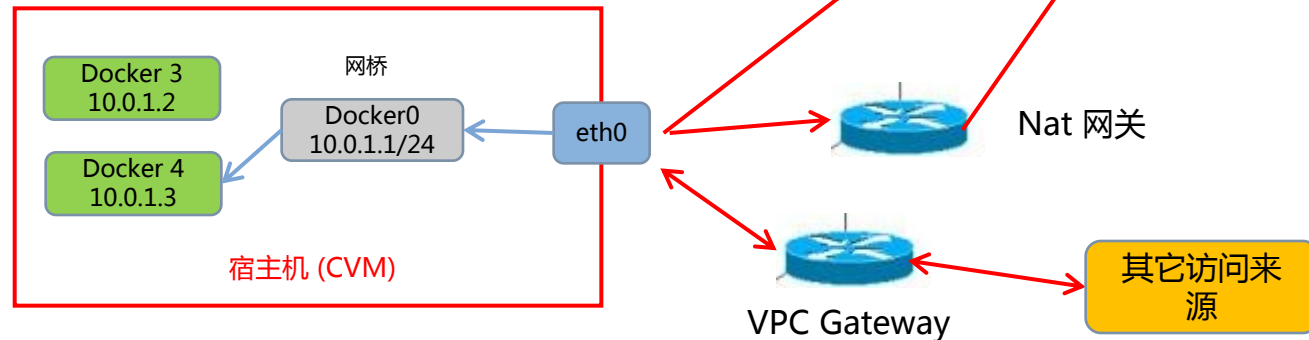
容器访问其他网络

宿主机 VPC网段：172.0.0.0/24
容器 VPC网段：10.0.0.0/14

CVM-1 :172.0.0.2 (母机IP 为 10.172.1.1)



CVM-2 172.0.0.3 (母机IP 为 10.172.1.2)

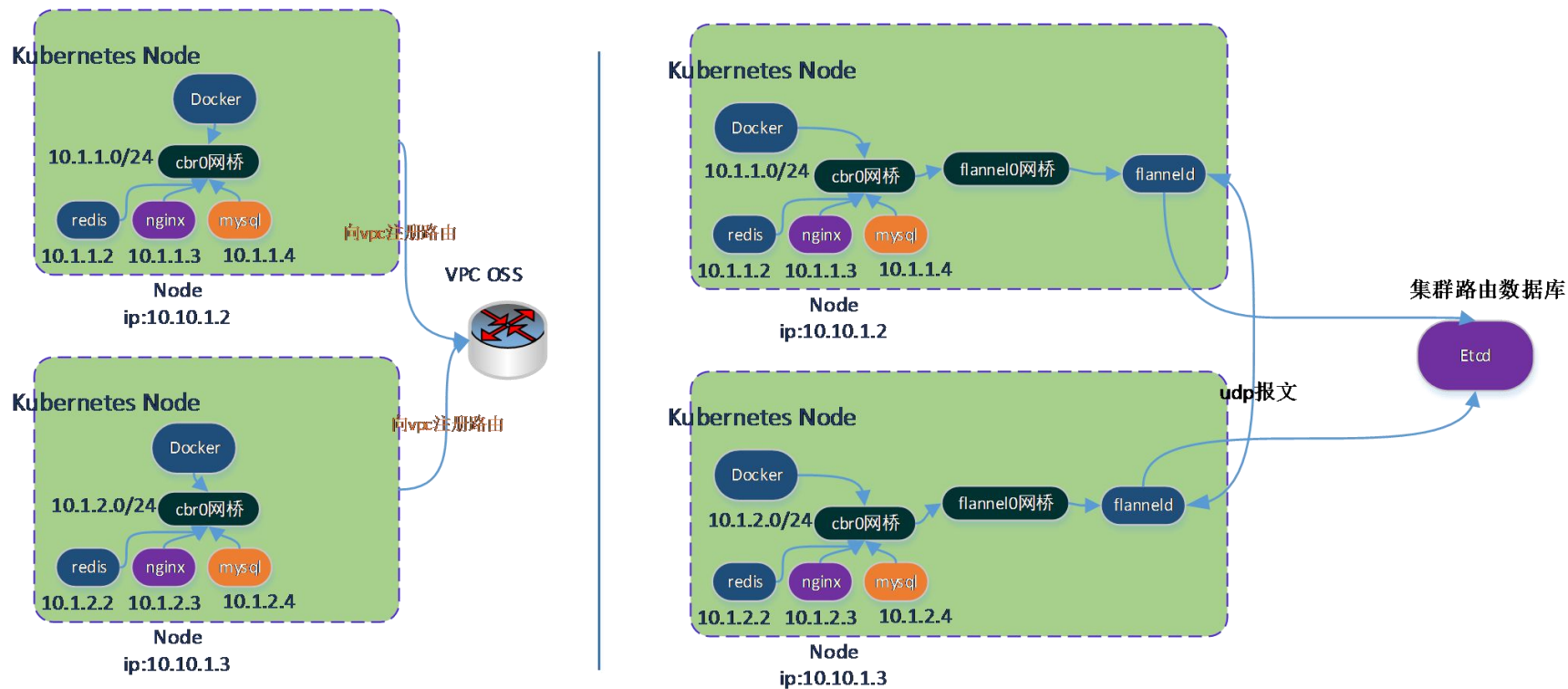


➤ 容器可以通过 VPC Gateway, 与其它VPC网段, 基础网络相互通信

➤ 容器可以通过宿主机的外网IP, 直接访问外网, 源地址做NAT转换

➤ 宿主机没有外网IP的话, 容器可以通过NAT网关直接访问外网, 源地址做NAT转换

容器网络对比云上自建flannel



容器路由由vpc软交换机负责

- 所有发往10.1.1.0/24下的流量，全路由到node1
- 无封包解包消耗

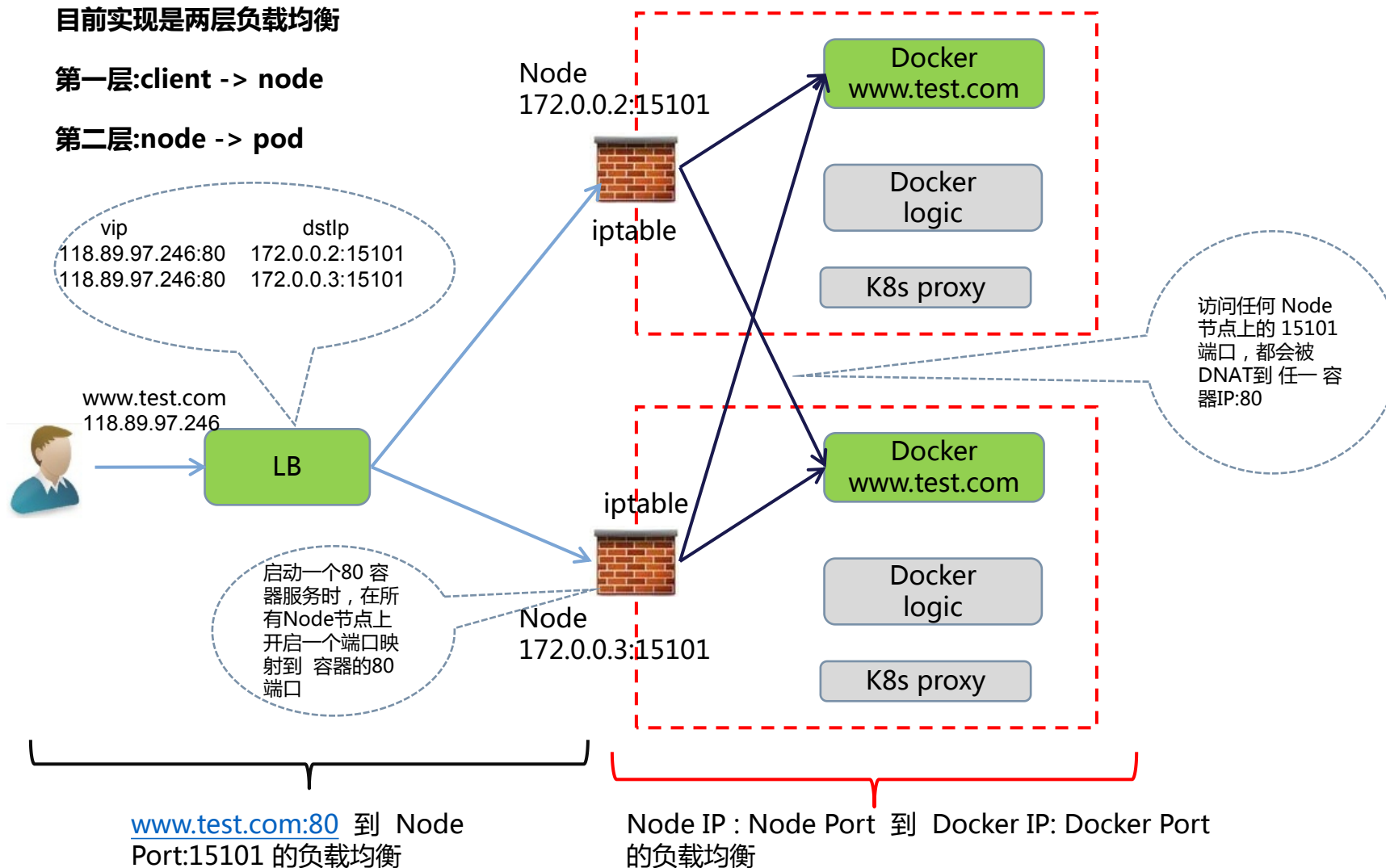
对接四层LB

LB 插件 将K8S 与 腾讯云负载均衡打通,

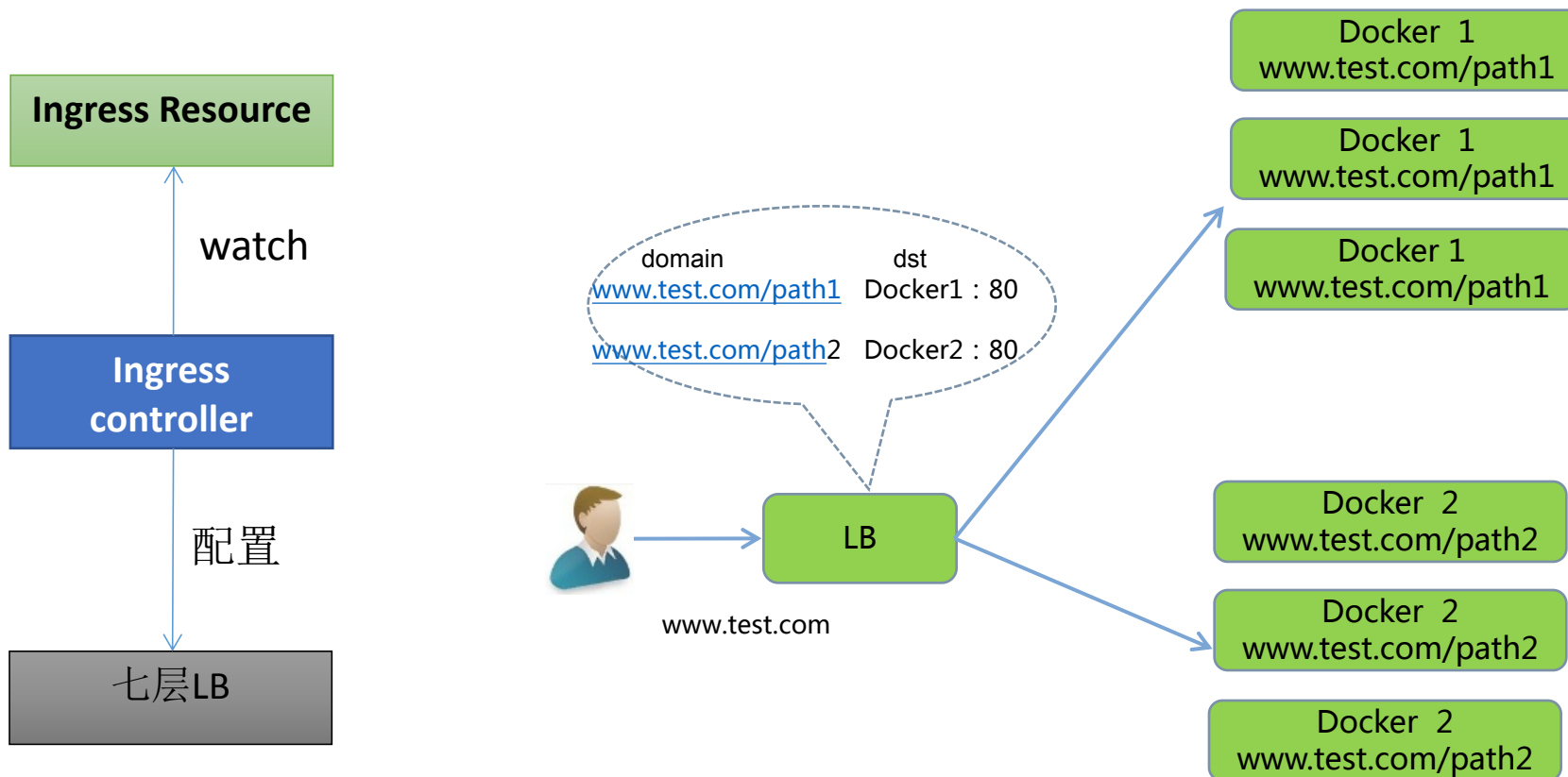
目前实现是两层负载均衡

第一层:client -> node

第二层:node -> pod

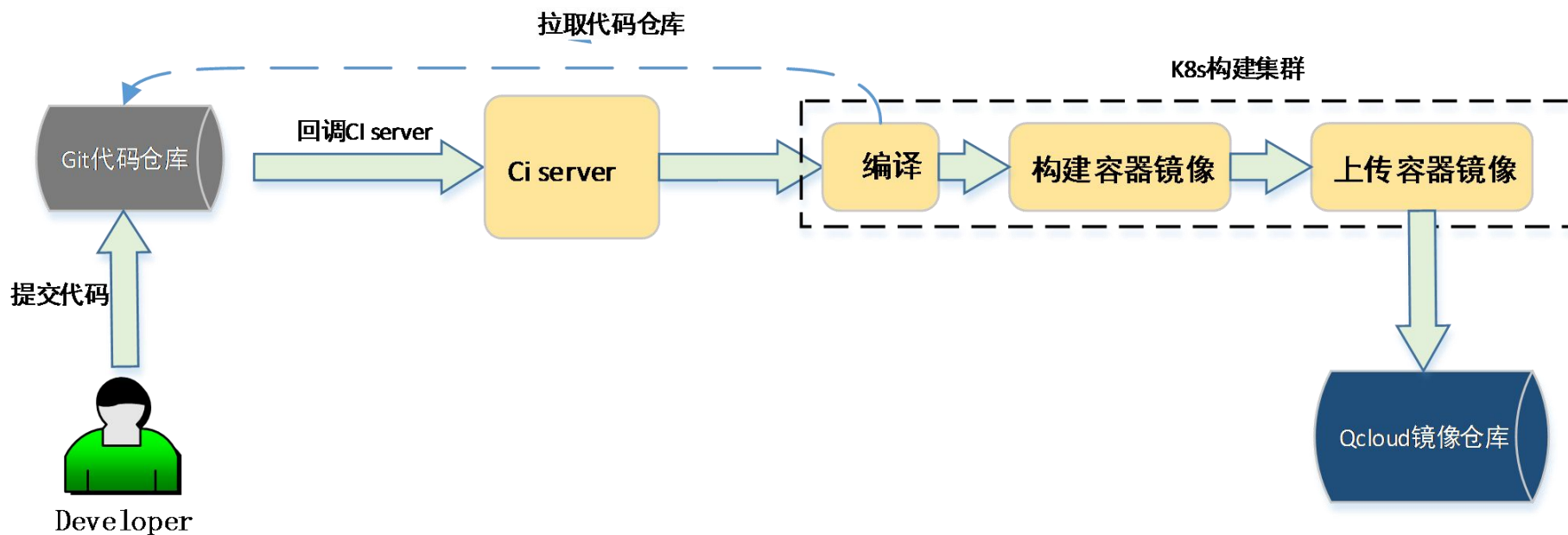


对接七层LB

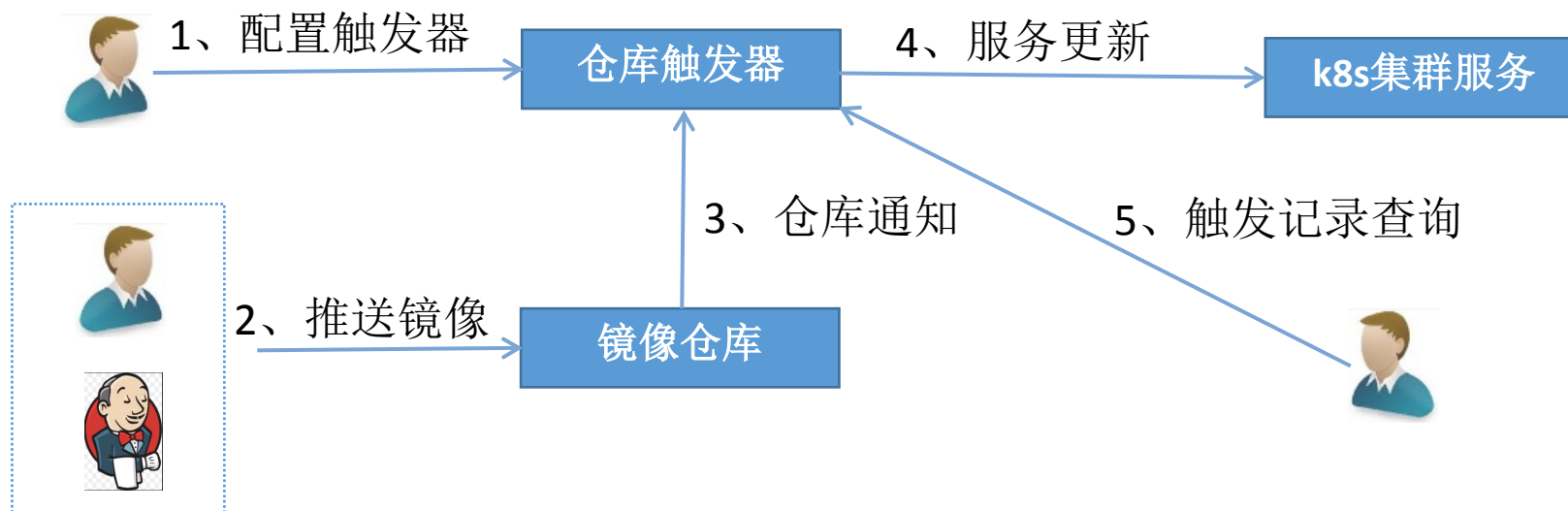


ingress 原理

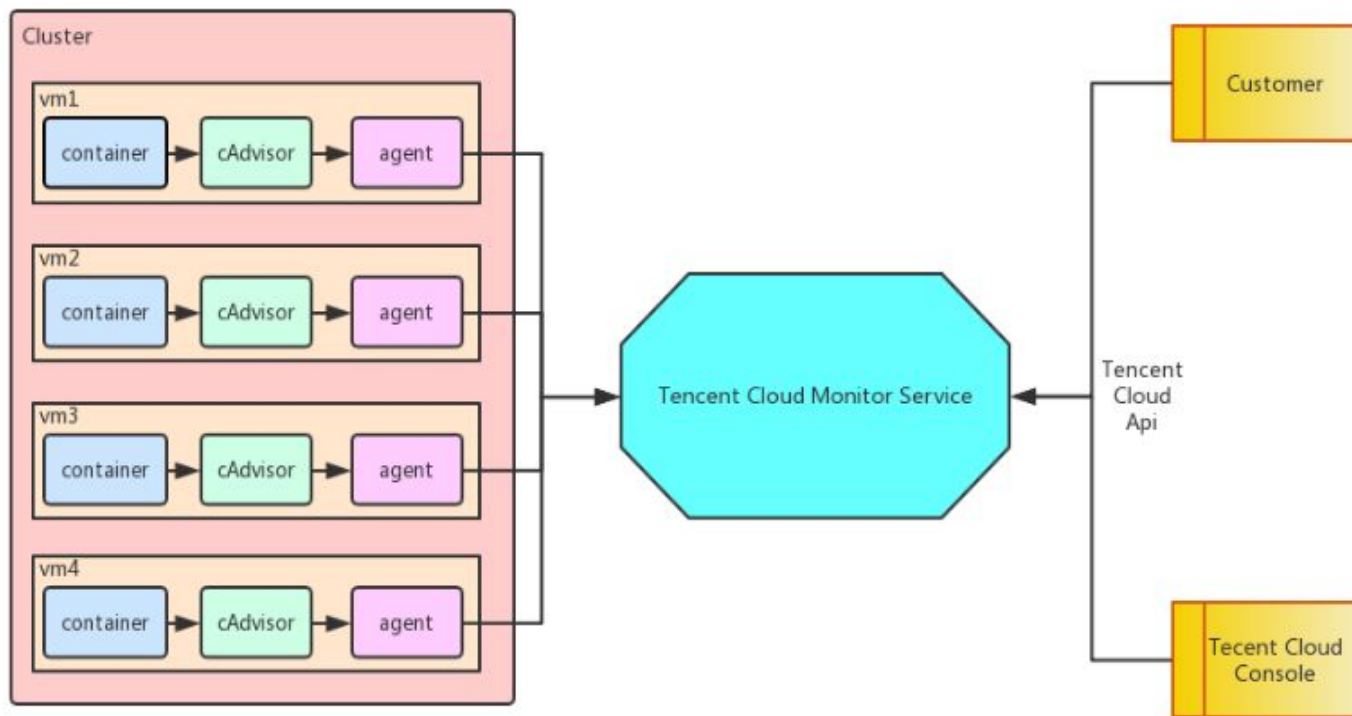
根据域名路径转发到 后端不同 容器服务



- 支持Github和gitlab仓库，在仓库提交代码后，在我们的CI平台上自动触发容器镜像构建。
- 支持在控制台上传dockerfile来构建容器镜像。
- 构建任务在海外执行，访问海外资源(基础镜像、安装包等)速度很快。
- 支持并发执行，用户默认2个并发构建任务，每个构建任务默认分配2C2G资源，构建时间限制为40Min。

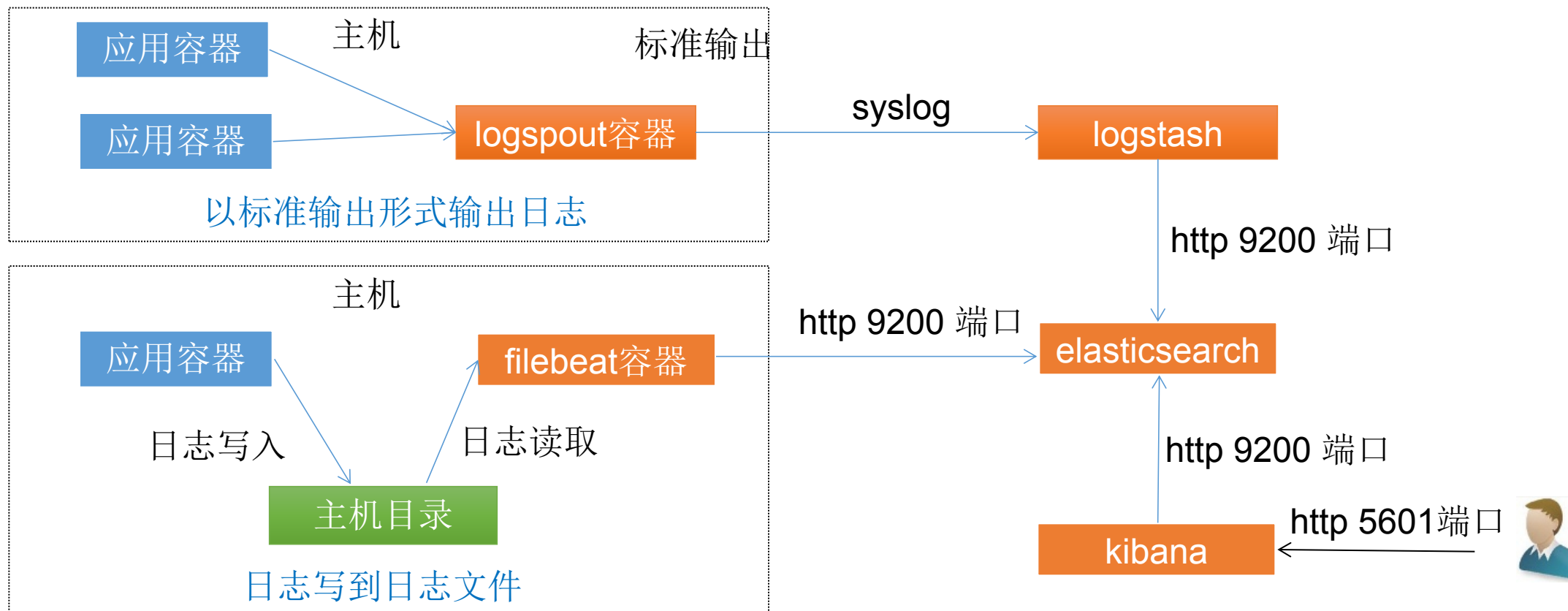


- 触发器配置包括触发条件，触发动作设置
- 支持手动和CI工具触发
- 支持CD结果查询，可查询触发动作是否成功及相关触发参数

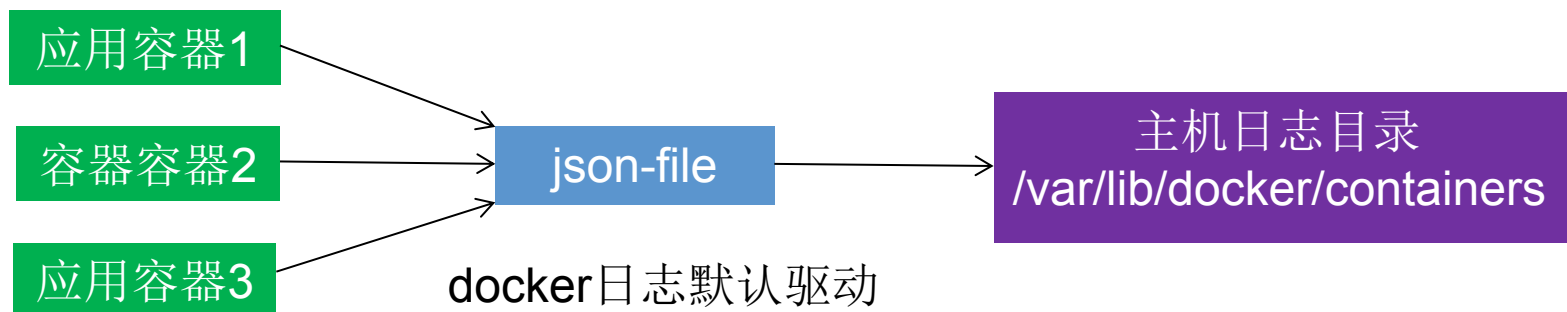


- agent每分钟会从cadvisor采集一次数据，报给监控平台
- 支持的指标，包括主机、容器和服务层面的监控数据，如cpu、内存、磁盘和网络

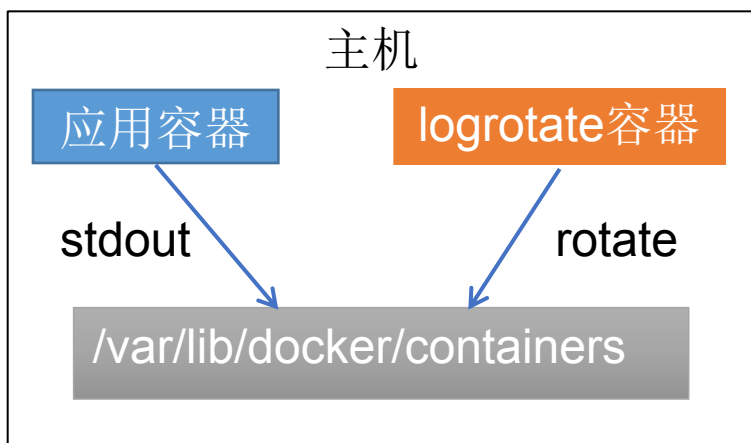
日志收集



本地日志清理



清理方式1:daemon set方式部署logrotate容器



清理方式2:修改dockerd 启动参数

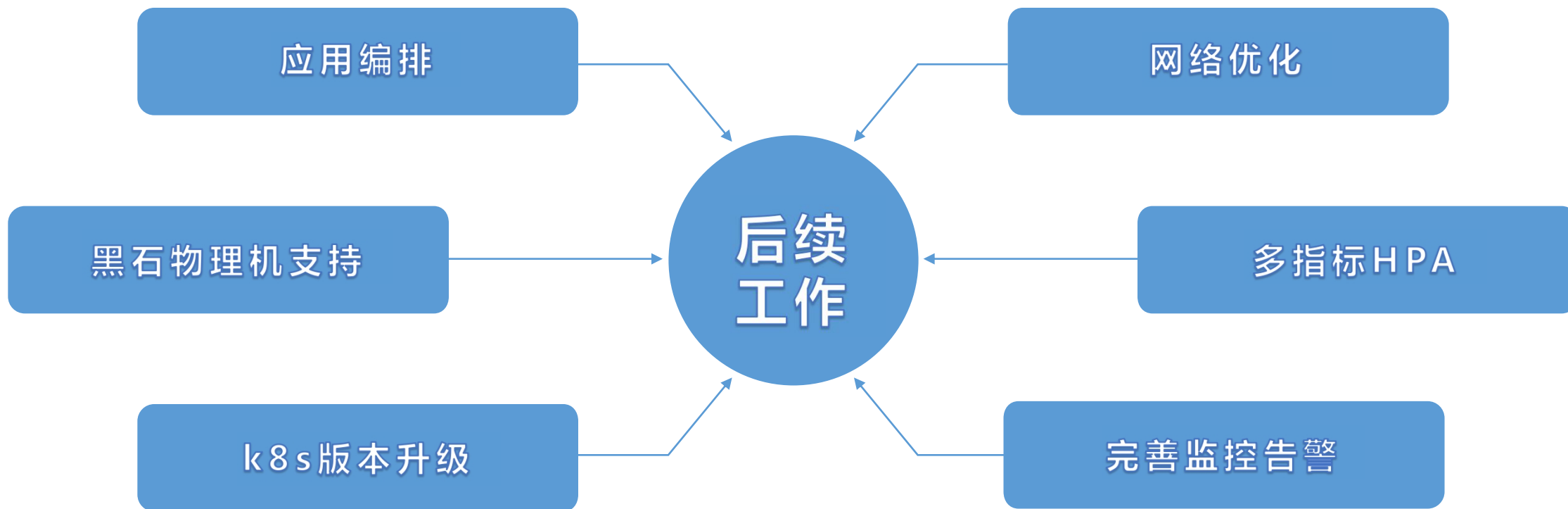
1、创建/etc/docker/daemon.json

```
{
  "log-driver": "json-file",
  "log-opts": {
    "max-size": "10m", "max-file": "3"
  }
}
```

2、修改dockerd 服务配置文件
/etc/systemd/system/multi-user.target.wants/dockerd.serviced

添加dockerd启动参数--config-file=/etc/docker/daemon.json

3、重启dockerd服务





Thanks

