



CSDN首页 > 云计算

订阅云计算RSS

# Kubernetes应用部署模型解析（原理篇）

发表于 2015-06-11 16:29 | 1277次阅读 | 来源 CSDN | 0条评论 | 作者 龚永生

Kubernetes 容器 Docker 云计算

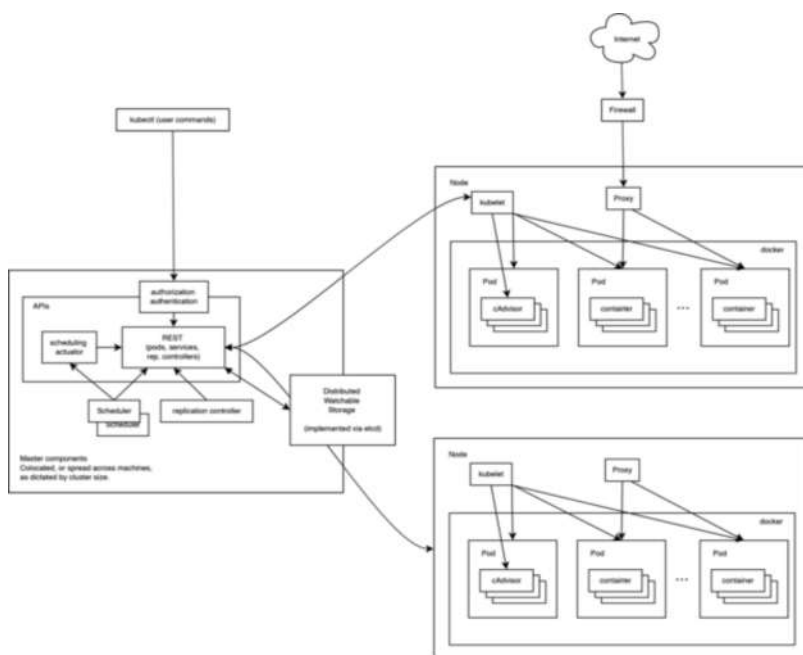
摘要：本系列文章着眼于实际部署，带您快速掌握Kubernetes。本文主要介绍部署之前需要了解的原理和概念，包括Kubernetes的组件结构，各个组件角色的功能，以及Kubernetes的应用模型等。

【编者按】Kubernetes可用来管理Linux容器集群，加速开发和简化运维（即DevOps）。但目前网络上关于Kubernetes的文章介绍性远多于实际使用。本系列文章着眼于实际部署，带您快速掌握Kubernetes。本文为上篇，主要介绍部署之前需要了解的原理和概念，包括Kubernetes的组件结构，各个组件角色的功能，以及Kubernetes的应用模型等。

十多年来Google一直在生产环境中使用容器运行业务，负责管理其容器集群的系统就是Kubernetes的前身Borg。其实现在很多工作在Kubernetes项目上的Google开发者先前就在Borg这个项目上工作。多数Kubernetes的应用部署模型的思想都起源于Borg，了解这些模型是掌握Kubernetes的关键。Kubernetes的API版本目前是v1，本文以代码0.18.2版为基础来介绍它的应用部署模型，最后我们用一个简单的用例来说明部署过程。在部署结束后，阐述了它是如何用iptables规则来实现各种类型Service的。

## Kubernetes架构

Kubernetes集群包括Kubernetes代理(agents)和Kubernetes服务(master node)两种角色，代理角色的组件包括Kube-proxy和Kubelet，它们同时部署在一个节点上，这个节点也就是代理节点。服务角色的组件包括kube-apiserver，kube-scheduler，kube-controller-manager，它们可以任意布属，它们可以部署在同一个节点上，也可以部署在不同的节点上（目前版本好像不行）。Kubernetes集群依赖的第三方组件目前有etcd和docker两个。前者提供状态存储，二者用来管理容器。集群还可以使用分布式存储给容器提供存储空间。下图显示了目前系统的组成部分：



## Kubernetes代理节点

Kubelet和Kube-proxy运行在代理节点上。他们监听服务节点的信息来启动容器和实现Kubernetes网络和其它业务模型，比如Service、Pod等。当然每个代理节点都运行Docker。Docker负责下载容器



CSDN官方微信

扫描二维码,向CSDN吐槽  
微信号: CSDNnews



程序员移动端订阅下载



## 每日资讯快速浏览

### 微博关注



CSDN云计算 北京 朝阳区

已关注

【2015中国SaaS生态“元素周期表”】去年以来，SaaS市场持续火爆，吸引了无数创业者或者投资机构的关注，为此我们特别策划了这期2015中国SaaS生态“元素周期表”的专题，希望从一个比较直观的角度勾勒出2015年中国SaaS大生态，共谱中国SaaS大势。http://t.cn/RLVvSY

今天 08:46

转发(2) | 评论



¥0.10/PCS

### 相关热门文章

新手福利：Apache Spark入门攻略

解密京东618技术：重构多中心交易平台 11000...

RebornDB：下一代分布式Key-Value数据库

镜像和运行容器。

## Kubelet

**Kubelet**组件管理**Pods**和它们的容器，镜像和卷等信息。

## Kube-Proxy

**Kube-proxy**是一个简单的网络代理和负载均衡器。它具体实现**Service**模型，每个**Service**都会所有的**Kube-proxy**节点上体现。根据**Service**的**selector**所覆盖的**Pods**，**Kube-proxy**会对这些**Pods**做负载均衡来服务于**Service**的访问者。

## Kubernetes服务节点

**Kubernetes**服务组件形成了**Kubernetes**的控制平面，目前他们运行在单一节点上，但是将来会分开来部署，以支持高可用性。

## etcd

所有的持久性状态都保存在**etcd**中。**Etcd**同时支持**watch**，这样组件很容易得到系统状态的变化，从而快速响应和协调工作。

## Kubernetes API Server

这个组件提供对**API**的支持，响应**REST**操作，验证**API**模型和更新**etcd**中的相应对象。

## Scheduler

通过访问**Kubernetes**中/binding API, **Scheduler**负责**Pods**在各个节点上的分配。**Scheduler**是插件式的，**Kubernetes**将来可以支持用户自定义的**scheduler**。

## Kubernetes Controller Manager Server

**Controller Manager Server**负责所有其它的功能，比如**endpoints**控制器负责**Endpoints**对象的创建，更新。**node**控制器负责节点的发现，管理和监控。将来可能会把这些控制器拆分并且提供插件式的实现。

## Kubernetes模型

**Kubernetes**的伟大之处就在于它的应用部署模型，主要包括**Pod**、**Replication controller**、**Label**和**Service**。

## Pod

**Kubernetes**的最小部署单元是**Pod**而不是容器。作为**First class API**公民，**Pods**能被创建，调度和管理。简单地来说，像一个豌豆荚中的豌豆一样，一个**Pod**中的应用容器同享同一个上下文：

1. **PID** 名字空间。但是在**docker**中不支持
2. 网络名字空间，在同一**Pod**中的多个容器访问同一个**IP**和端口空间。
3. **IPC**名字空间，同一个**Pod**中的应用能够使用**SystemV IPC**和**POSIX**消息队列进行通信。
4. **UTS**名字空间，同一个**Pod**中的应用共享一个主机名。
5. **Pod**中的各个容器应用还可以访问**Pod**级别定义的共享卷。

从生命周期来说，**Pod**应该是短暂的而不是长久的应用。**Pods**被调度到节点，保持在这个节点上直到被销毁。当节点死亡时，分配到这个节点的**Pods**将会被删掉。将来可能会实现**Pod**的迁移特性。在实际使用时，我们一般不直接创建**Pods**，我们通过**replication controller**来负责**Pods**的创建，复制，监控和销毁。一个**Pod**可以包括多个容器，他们直接往往相互协作完成一个应用功能。

## Replication controller

复制控制器确保**Pod**的一定数量的份数(**replica**)在运行。如果超过这个数量，控制器会杀死一些，如果少了，控制器会启动一些。控制器也会在节点失效、维护的时候来保证这个数量。所以强烈建议即使我们的份数是**1**，也要使用复制控制器，而不是直接创建**Pod**。

在生命周期上讲，复制控制器自己不会终止，但是跨度不会比**Service**强。**Service**能够横跨多个复制控制器管理的**Pods**。而且在一个**Service**的生命周期内，复制控制器能被删除和创建。**Service**和客户端程序是不知道复制控制器的存在的。

复制控制器创建的**Pods**应该是可以互相替换的和语义上相同的，这个对无状态服务特别合适。

**Pod**是临时性的对象，被创建和销毁，而且不会恢复。复制器动态地创建和销毁**Pod**。虽然**Pod**会分配到**IP**地址，但是这个**IP**地址都不是持久的。这样就产生了一个疑问：外部如何消费**Pod**提供的服务

开发者成功使用机器学习的十大诀窍

AWS Quick Start参考部署方案

2015中国SaaS生态“元素周期表”

Expedia如何对相互依赖的数据集进行准实时分析

【报名】DefCon黑客大会来了，不差钱，只差你！

MapReduce、Spark、Phoenix、Disco、Mars...

新的可视化帮助更好地了解Spark Streaming应...

## 热门标签

Hadoop	AWS	移动游戏
Java	Android	iOS
Swift	智能硬件	Docker
OpenStack	VPN	Spark
ERP	IE10	Eclipse
CRM	JavaScript	数据库
Ubuntu	NFC	WAP

## CSDN Share PPT下载



JAVA总结基础部分



指数级增长业务下的服务架构改造



3.张宁--移动大数据技术在互联网金融获客及经营中的应用

呢？

## Service

**Service**定义了一个Pod的逻辑集合和访问这个集合的策略。集合是通过定义**Service**时提供的**Label**选择器完成的。举个例子，我们假定有3个Pod的备份来完成一个图像处理的后端。这些后端备份逻辑上是相同的，前端不关心哪个后端在给它提供服务。虽然组成这个后端的实际Pod可能变化，前端客户端不会意识到这个变化，也不会跟踪后端。**Service**就是用来实现这种分离的抽象。

对于**Service**，我们还可以定义**Endpoint**，**Endpoint**把**Service**和Pod动态地连接起来。

### Service Cluster IP 和 kube proxy

每个代理节点都运行了一个**kube-proxy**进程。这个进程从服务进程那边拿到**Service**和**Endpoint**对象的变化。对每一个**Service**，它在本地打开一个端口。到这个端口的任意连接都会代理到后端Pod集合中的一个Pod IP和端口。在创建了服务后，服务**Endpoint**模型会体现后端Pod的IP和端口列表，**kube-proxy**就是从这个**endpoint**维护的列表中选择服务后端的。另外**Service**对象的**sessionAffinity**属性也会帮助**kube-proxy**来选择哪个具体的后端。缺省情况下，后端Pod的选择是随机的。可以设置**service.spec.sessionAffinity** 成"**ClientIP**"来指定同一个**ClientIP**的流量代理到同一个后端。在实现上，**kube-proxy**会用**IPtables**规则把访问**Service**的**Cluster IP**和端口的流量重定向到这个本地端口。下面的部分会讲什么是**service**的**Cluster IP**。

注意：在0.18以前的版本中**Cluster IP**叫**PortalNet IP**。

### 内部使用者的服务发现

Kubernetes在一个集群内创建的对象或者在代理集群节点上发出访问的客户端我们称之为内部使用者。要把服务暴露给内部使用者，Kubernetes支持两种方式：环境变量和DNS。

#### 环境变量

当**kubelet**在某个节点上启动一个Pod时，它会给这个Pod的容器为当前运行的**Service**设置一系列环境变量，这样Pod就可以访问这些**Service**了。一般地情况

是{**SVCNAME**}\_**SERVICE\_HOST**和{**SVCNAME**}\_**SERVICE\_PORT**变量, 其中{**SVCNAME**}是**Service**名字变成大写，中划线变成下划线。比如**Service "redis-master"**，它的端口是 **TCP 6379**，分配到的**Cluster IP**地址是 **10.0.0.11**，**kubelet**可能会产生下面的变量给新创建的Pod容器：

```
REDIS_MASTER_SERVICE_HOST= 10.0.0.11
REDIS_MASTER_SERVICE_PORT=6379
REDIS_MASTER_PORT=tcp://10.0.0.11:6379
REDIS_MASTER_PORT_6379_TCP=tcp://10.0.0.11:6379
REDIS_MASTER_PORT_6379_TCP_PROTO=tcp
REDIS_MASTER_PORT_6379_TCP_PORT=6379
REDIS_MASTER_PORT_6379_TCP_ADDR= 10.0.0.11
```

注意，只有在某个**Service**后创建的Pod才会有这个**Service**的环境变量。

#### DNS

一个可选的Kubernetes附件（强烈建议用户使用）是DNS服务。它跟踪集群中**Service**对象，为每个**Service**对象创建DNS记录。这样所有的Pod就可以通过DNS访问服务了。

比如说我们在Kubernetes 名字空间"**my-ns**"中有个叫**my-service**的服务，DNS服务会创建一条"**my-service.my-ns**"的DNS记录。同在这个命名空间的Pod就可以通过"**my-service**"来得到这个**Service**分配到的**Cluster IP**，在其它命名空间的Pod则可以用全限定名"**my-service.my-ns**"来获得这个**Service**的地址。

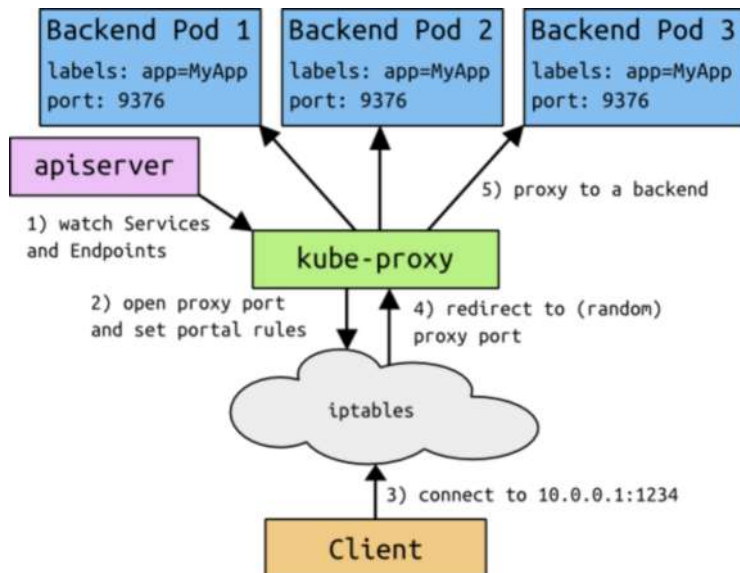
### Pod IP and Service Cluster IP

**Pod IP** 地址是实际存在于某个网卡(可以是虚拟设备)上的，但**Service Cluster IP**就不一样了，没有网络设备为这个地址负责。它是由**kube-proxy**使用**Iptables**规则重新定向到其本地端口，再均衡到后端Pod的。我们前面说的**Service**环境变量和DNS都使用**Service**的**Cluster IP**和端口。

就拿上面我们提到的图像处理程序为例。当我们的**Service**被创建时，Kubernetes给它分配一个地址**10.0.0.1**。这个地址从我们启动API的**service-cluster-ip-range**参数(旧版本为**portal\_net**参数)指定的地址池中分配，比如**--service-cluster-ip-range=10.0.0.0/16**。假设这个**Service**的端口是**1234**。集群内的所有**kube-proxy**都会注意到这个**Service**。当**proxy**发现一个新的**service**后，它会在本地节点打开一个任意端口，建相应的**iptables**规则，重定向服务的IP和port到这个新建的端口，开始接受到达这个服务的连接。

当一个客户端访问这个**service**时，这些**iptables**规则就开始起作用，客户端的流量被重定向到**kube-**

proxy为这个service打开的端口上，kube-proxy随机选择一个后端pod来服务客户。这个流程如下图所示：



根据Kubernetes的网络模型，使用Service Cluster IP和Port访问Service的客户端可以坐落在任意代理节点上。外部要访问Service，我们就需要给Service外部访问IP。

### 外部访问Service

Service对象在Cluster IP range池中分配到的IP只能在内部访问，如果服务作为一个应用程序内部的层次，还是很合适的。如果这个Service作为前端服务，准备为集群外的客户提供业务，我们就需要给这个服务提供公共IP了。

外部访问者是访问集群代理节点的访问者。为这些访问者提供服务，我们可以在定义Service时指定其spec.publicIPs，一般情况下publicIP是代理节点的物理IP地址。和先前的Cluster IP range上分配到的虚拟的IP一样，kube-proxy同样会为这些publicIP提供iptables重定向规则，把流量转发到后端的Pod上。有了publicIP，我们就可以使用load balancer等常用的互联网技术来组织外部对服务的访问了。

spec.publicIPs在新的版本中标记为过时了，代替它的是spec.type=NodePort，这个类型的service，系统会给它在集群的各个代理节点上分配一个节点级别的端口，能访问到代理节点的客户端都能访问这个端口，从而访问到服务。

## Label和Label selector

Label标签在Kubernetes模型中占着非常重要的作用。Label表现为key/value对，附加到Kubernetes管理的对象上，典型的就是Pods。它们定义了这些对象的识别属性，用来组织和选择这些对象。Label可以在对象创建时附加在对象上，也可以对象存在时通过API管理对象的Label。

在定义了对象的Label后，其它模型可以用Label选择器（selector）来定义其作用的对象。

Label选择器有两种，分别是Equality-based和Set-based。

比如如下Equality-based选择器样例：

```
environment = production
tier != frontend
environment = production, tier != frontend
```

对于上面的选择器，第一条匹配Label具有environment key且等于production的对象，第二条匹配具有tier key，但是值不等于frontend的对象。由于kubernetes使用AND逻辑，第三条匹配production但不是frontend的对象。

Set-based选择器样例：

```
environment in (production, qa)
tier notin (frontend, backend)
partition
```

第一条选择具有**environment key**，而且值是**production**或者**qa**的**label**附加的对象。第二条选择具有**tier key**，但是其值不是**frontend**和**backend**。第三条选则具有**partition key**的对象，不对**value**进行校验。

**replication controller**复制控制器和**Service**都用**label**和**label selector**来动态地配备作用对象。复制控制器在定义的时候就指定了其要创建**Pod**的**Label**和自己要匹配这个**Pod**的**selector**，**API**服务器应该校验这个定义。我们可以动态地修改**replication controller**创建的**Pod**的**Label**用于调式，数据恢复等。一旦某个**Pod**由于**Label**改变从**replication controller**移出来后，**replication controller**会马上启动一个新的**Pod**来确保复制池子中的份数。对于**Service**，**Label selector**可以用来选择一个**Service**的后端**Pods**。

下篇：[Kubernetes应用部署模型解析（部署篇）](#)（责编/周建丁）

作者简介：龚永生，九州云架构师。多年Linux系统开发，J2EE产品和云计算相关技术研发经验。目前活跃在OpenStack社区的各个项目上，主要技术方向是虚拟网络项目Neutron，是Neutron项目早期的主要贡献者之一。

本文为CSDN原创文章，未经允许不得转载，如需转载请联系[market#csdn.net](mailto:market#csdn.net)(#换成@)

顶  
1

踩  
0



推荐阅读相关主题：

相关文章

最新报道

谁是容器中的“战斗机”？ Docker与Chef、LXC等容器...  
思科和红帽拟正式推出Linux应用程序容器技术  
京东田琪：让Container发威，你应该了解这些核心...  
运用Kubernetes进行分布式负载测试  
基于容器的自动构建——Docker在美国的应用  
灵雀云：CaaS和微服务架构终结传统PaaS？

已有0条评论

还可以再输入500个字



有什么感想，你也来说说吧！

您还没有登录! 请 登录 或 注册

发表评论

最新评论

最热评论

还没有评论，赶快来抢沙发吧。

请您注意

- 自觉遵守：爱国、守法、自律、真实、文明的原则
- 尊重网上道德，遵守《全国人大常委会关于维护互联网安全的决定》及中华人民共和国其他各项有关法律法规
- 严禁发表危害国家安全，破坏民族团结、国家宗教政策和社会稳定，含侮辱、诽谤、教唆、淫秽等内容的作品
- 承担一切因您的行为而直接或间接导致的民事或刑事法律责任
- 您在CSDN新闻评论发表的作品，CSDN有权在网站内保留、转载、引用或者删除
- 参与本评论即表明您已经阅读并接受上述条款

热门专区



容联云通讯开发者技术专区



腾讯云技术社区



IBM新兴技术大学



高效能团队解决方案



高通开发者专区

---

[公司简介](#) | [招贤纳士](#) | [广告服务](#) | [银行汇款帐号](#) | [联系方式](#) | [版权声明](#) | [法律顾问](#) | [问题报告](#) | [合作伙伴](#) | [论坛反馈](#)

---

[网站客服](#)   [杂志客服](#)   [微博客服](#)   [webmaster@csdn.net](mailto:webmaster@csdn.net)   400-600-2320 | 北京创新乐知信息技术有限公司 版权所有 | 江苏乐知网络技术有限公司 提供商务支持  
京 ICP 证 070598 号 | Copyright © 1999-2014, CSDN.NET, All Rights Reserved 