# Kubernetes二进制集群安装

本系列文档将介绍如何使用二进制部署 `Kubernetes v1.14` 集群的所有部署，而不是使用自动化部署 (kubeadm)集群。在部署过程中，将详细列出各个组件启动参数，以及相关配置说明。在学习完本文档后，将理解k8s各个组件的交互原理，并且可以快速解决实际问题。

> 本文档适用于 `Centos7.4` 及以上版本，随着各个组件的更新，本文档提供了相关镜像的包，及时版本更新也不会影响文档的使用。 如果有文档相关问题可以直接在网站下面注册回复，或者点击右下角加群，我将在12小时内回复您。 并且建议您使用的环境及配置和我相同！

## 组件版本

- Kubernetes 1.14.2
- Docker 18.09 (docker使用官方的脚本安装，后期可能升级为新的版本，但是不影响)
- Etcd 3.3.13
- Flanneld 0.11.0

## 组件说明

kube-apiserver
- 使用节点本地Nginx 4层透明代理实现高可用（也可以使用haproxy，只是起到代理apiserver的作用）
- 关闭非安全端口8080和匿名访问
- 使用安全端口6443接受https请求
- 严格的认知和授权策略 (x509、token、rbac)
- 开启bootstrap token认证，支持kubelet TLS bootstrapping；
- 使用https访问kubelet、etcd

kube-controller-manager
- 3节点高可用（在k8s中，有些组件需要选举，所以使用奇数为集群高可用方案）
- 关闭非安全端口，使用10252接受https请求

- 使用kubeconfig访问apiserver的安全扣
- 使用approve kubelet证书签名请求(CSR)，证书过期后自动轮转
- 各controller使用自己的ServiceAccount访问apiserver

kube-scheduler
- 3节点高可用；
- 使用kubeconfig访问apiserver安全端口

kubelet
- 使用kubeadm动态创建bootstrap token
- 使用TLS bootstrap机制自动生成client和server证书，过期后自动轮转
- 在kubeletConfiguration类型的JSON文件配置主要参数
- 关闭只读端口，在安全端口10250接受https请求，对请求进行认真和授权，拒绝匿名访问和非授权访问问
- 使用kubeconfig访问apiserver的安全端口

kube-proxy
- 使用kubeconfig访问apiserver的安全端口
- 在KubeProxyConfiguration类型JSON文件配置为主要参数
- 使用ipvs代理模式

集群插件
- DNS 使用功能、性能更好的coredns
- 网络 使用FlanneId 作为集群网络插件

---

# 一、初始化环境

**集群机器**

```
1  192.168.0.50 k8s-01
2  192.168.0.51 k8s-02
3  192.168.0.52 k8s-03
4  #node节点
5  192.168.0.53 k8s-04       #node节点只运行node，但是设置证书的时候要添加这个ip
```

本文档的所有etcd集群、master集群、worker节点均使用以上三台机器，并且初始化步骤需要在所有机器上执行命令。如果没有特殊命令，所有操作均在**192.168.0.50**上进行操作

> node节点后面会有操作，但是在初始化这步，是所有集群机器。包括node节点，我上面没有列出node节点

**修改主机名**

所有机器设置永久主机名

```
1 hostnamectl set-hostname k8s01   #所有机器按照要求修改
2 bash           #刷新主机名
```

接下来我们需要在所有机器上添加hosts解析

```
1 cat >> /etc/hosts <<EOF
2 192.168.0.50  k8s-01
3 192.168.0.51  k8s-02
4 192.168.0.52  k8s-03
5 192.168.0.53  k8s-04
6 EOF
```

**设置免密**
我们只在k8s-01上设置免密即可

```
1 wget -O /etc/yum.repos.d/epel.repo http://mirrors.aliyun.com/repo/ep
  el-7.repo
2 curl -o /etc/yum.repos.d/CentOS-Base.repo http://mirrors.aliyun.com/
  repo/Centos-7.repo
3 yum install -y expect
4 #分发公钥
5 ssh-keygen -t rsa -P "" -f /root/.ssh/id_rsa
6 for i in k8s-01 k8s-02 k8s-03 k8s-04;do
7 expect -c "
8 spawn ssh-copy-id -i /root/.ssh/id_rsa.pub root@$i
9         expect {
10                \"*yes/no*\" {send \"yes\r\"; exp_continue}
11                \"*password*\" {send \"123456\r\"; exp_continue}
12                \"*Password*\" {send \"123456\r\";}
13        } "
14 done
15 #我这里密码是123456   大家按照自己主机的密码进行修改就可以
```

**更新PATH变量**

本次的k8s软件包的目录全部存放在 `/opt` 下

```
1  [root@k8s01 ~]# echo 'PATH=/opt/k8s/bin:$PATH' >>/etc/profile
2  [root@k8s01 ~]# source  /etc/profile
3  [root@k8s01 ~]# env|grep PATH
4  PATH=/opt/k8s/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/
   root/bin
```

**安装依赖包**

在**每台**服务器上安装依赖包

```
1  yum install -y conntrack ntpdate ntp ipvsadm ipset jq iptables curl s
   ysstat libseccomp wget
```

**关闭防火墙 Linux 以及swap分区**

```
1  systemctl stop firewalld
2  systemctl disable firewalld
3  iptables -F && iptables -X && iptables -F -t nat && iptables -X -t na
   t
4  iptables -P FORWARD ACCEPT
5  swapoff -a
6  sed -i '/ swap / s/^\(.*\)$/#\1/g' /etc/fstab
7  setenforce 0
8  sed -i 's/^SELINUX=.*/SELINUX=disabled/' /etc/selinux/config
```

#如果开启了swap分区，kubelet会启动失败(可以通过设置参数——-fail-swap-on设置为false)

**升级内核**

Docker overlay2需要使用kernel 4.x版本，所以我们需要升级内核

我这里的内核使用4.18.9

CentOS 7.x 系统自带的 3.10.x 内核存在一些 Bugs，导致运行的 Docker、Kubernetes 不稳定，例如:

```
1  高版本的 docker(1.13 以后) 启用了 3.10 kernel 实验支持的 kernel memory acco
```

```
    unt 功能(无法关闭)，当节点压力大如频繁启动和停止容器时会导致 cgroup memory leak
    ;
  2 网络设备引用计数泄漏，会导致类似于报错: "kernel:unregister_netdevice: waiting
    for eth0 to become free. Usage count = 1";
```

解决方案如下：

```
  1 升级内核到 4.4.X 以上;
  2 或者，手动编译内核，disable CONFIG_MEMCG_KMEM 特性;
  3 或者，安装修复了该问题的 Docker 18.09.1 及以上的版本。但由于 kubelet 也会设置 k
    mem（它 vendor 了 runc），所以需要重新编译 kubelet 并指定 GOFLAGS="-tags=no
    kmem";
  4 export Kernel_Version=4.18.9-1
  5 wget  http://mirror.rc.usf.edu/compute_lock/elrepo/kernel/el7/x86_64/
    RPMS/kernel-ml{,-devel}-${Kernel_Version}.el7.elrepo.x86_64.rpm
  6 yum localinstall -y kernel-ml*
  7 #如果是手动下载内核rpm包，直接执行后面yum install -y kernel-ml*即可
```

修改内核启动顺序,默认启动的顺序应该为1,升级以后内核是往前面插入,为0（如果每次启动时需要手动选择哪个内核,该步骤可以省略）

```
  1 grub2-set-default  0 && grub2-mkconfig -o /etc/grub2.cfg
```

使用下面命令看看确认下是否启动默认内核指向上面安装的内核

```
  1 grubby --default-kernel
  2 #这里的输出结果应该为我们升级后的内核信息
```

重启加载新内核 （升级完内核顺便update一下）

```
  1 reboot
```

**加载内核模块**

首先我们要检查是否存在所需的内核模块

```
1 find  /lib/modules/`uname -r`/ -name "ip_vs_rr*"
2 find  /lib/modules/`uname -r`/ -name "br_netfilter*"
```

1.加载内核，加入开机启动（2选1即可）

```
1 cat > /etc/rc.local  << EOF
2 modprobe ip_vs_rr
3 modprobe br_netfilter
4 EOF
```

2.使用systemd-modules-load加载内核模块

```
1 cat > /etc/modules-load.d/ipvs.conf << EOF
2  ip_vs_rr
3  br_netfilter
4 EOF
5 systemctl enable --now systemd-modules-load.service
```

验证模块是否加载成功

```
1 lsmod |egrep " ip_vs_rr|br_netfilter"
2 为什么要使用IPVS,从k8s的1.8版本开始，kube-proxy引入了IPVS模式，IPVS模式与iptab
  les同样基于Netfilter，但是采用的hash表，因此当service数量达到一定规模时，hash查
  表的速度优势就会显现出来，从而提高service的服务性能。
3 ipvs依赖于nf_conntrack_ipv4内核模块,4.19包括之后内核里改名为nf_conntrack,1.1
  3.1之前的kube-proxy的代码里没有加判断一直用的nf_conntrack_ipv4,好像是1.13.1后
  的kube-proxy代码里增加了判断,我测试了是会去load nf_conntrack使用ipvs正常
```

**优化内核参数**

```
1  cat > kubernetes.conf <<EOF
2  net.bridge.bridge-nf-call-iptables=1
3  net.bridge.bridge-nf-call-ip6tables=1
4  net.ipv4.ip_forward=1
5  net.ipv4.tcp_tw_recycle=0
6  vm.swappiness=0 # 禁止使用 swap 空间，只有当系统 OOM 时才允许使用它
7  vm.overcommit_memory=1 # 不检查物理内存是否够用
8  vm.panic_on_oom=0 # 开启 OOM
9  fs.inotify.max_user_instances=8192
10 fs.inotify.max_user_watches=1048576
11 fs.file-max=52706963
12 fs.nr_open=52706963
13 net.ipv6.conf.all.disable_ipv6=1
14 net.netfilter.nf_conntrack_max=2310720
15 EOF
16 cp kubernetes.conf  /etc/sysctl.d/kubernetes.conf
17 sysctl -p /etc/sysctl.d/kubernetes.conf
```

需要关闭 `tcp_tw_recycle`，否则和NAT冲突，会导致服务不通

关闭IPV6，防止触发Docker BUG

**设置系统时区**

```
1  timedatectl set-timezone Asia/Shanghai
2   #将当前的 UTC 时间写入硬件时钟
3  timedatectl set-local-rtc 0
4   #重启依赖于系统时间的服务
5  systemctl restart rsyslog
6  systemctl restart crond
```

**创建相关目录**

```
1  mkdir -p  /opt/k8s/{bin,work} /etc/{kubernetes,etcd}/cert
```

#在所有节点上执行，因为flanneld是在所有节点运行的

**设置分发脚本参数**

后续所有的使用环境变量都定义在environment.sh中，需要根据个人机器及网络环境修改。并且需要拷贝到所有节点的/opt/k8s/bin目录下

```bash
#!/usr/bin/bash
# 生成 EncryptionConfig 所需的加密 key
export ENCRYPTION_KEY=$(head -c 32 /dev/urandom | base64)
# 集群各机器 IP 数组
export NODE_IPS=( 192.168.0.50 192.168.0.51 192.168.0.52 192.168.0.53 )
# 集群各 IP 对应的主机名数组
export NODE_NAMES=(k8s-01 k8s-02 k8s-03 k8s-04)
# 集群MASTER机器 IP 数组
export MASTER_IPS=(192.168.0.50 192.168.0.51 192.168.0.52 )
# 集群所有的master Ip对应的主机
export MASTER_NAMES=(k8s-01 k8s-02 k8s-03)
# etcd 集群服务地址列表
export ETCD_ENDPOINTS="https://192.168.0.50:2379,https://192.168.0.51:2379,https://192.168.0.52:2379"
# etcd 集群间通信的 IP 和端口
export ETCD_NODES="k8s-01=https://192.168.0.50:2380,k8s-02=https://192.168.0.51:2380,k8s-03=https://192.168.0.52:2380"
# etcd 集群所有node ip
export ETCD_IPS=(192.168.0.50 192.168.0.51 192.168.0.52 192.168.0.53 )
# kube-apiserver 的反向代理(kube-nginx)地址端口
export KUBE_APISERVER="https://192.168.0.54:8443"
# 节点间互联网络接口名称
export IFACE="eth0"
# etcd 数据目录
export ETCD_DATA_DIR="/data/k8s/etcd/data"
# etcd WAL 目录，建议是 SSD 磁盘分区，或者和 ETCD_DATA_DIR 不同的磁盘分区
export ETCD_WAL_DIR="/data/k8s/etcd/wal"
# k8s 各组件数据目录
export K8S_DIR="/data/k8s/k8s"
# docker 数据目录
#export DOCKER_DIR="/data/k8s/docker"
## 以下参数一般不需要修改
# TLS Bootstrapping 使用的 Token，可以使用命令 head -c 16 /dev/urandom | od -An -t x | tr -d ' ' 生成
#BOOTSTRAP_TOKEN="41f7e4ba8b7be874fcff18bf5cf41a7c"
# 最好使用 当前未用的网段 来定义服务网段和 Pod 网段
# 服务网段，部署前路由不可达，部署后集群内路由可达(kube-proxy 保证)
```

```
35  SERVICE_CIDR="10.254.0.0/16"
36  # Pod 网段，建议 /16 段地址，部署前路由不可达，部署后集群内路由可达(flanneld 保
    证)
37  CLUSTER_CIDR="172.30.0.0/16"
38  # 服务端口范围 (NodePort Range)
39  export NODE_PORT_RANGE="1024-32767"
40  # flanneld 网络配置前缀
41  export FLANNEL_ETCD_PREFIX="/kubernetes/network"
42  # kubernetes 服务 IP (一般是 SERVICE_CIDR 中第一个IP)
43  export CLUSTER_KUBERNETES_SVC_IP="10.254.0.1"
44  # 集群 DNS 服务 IP (从 SERVICE_CIDR 中预分配)
45  export CLUSTER_DNS_SVC_IP="10.254.0.2"
46  # 集群 DNS 域名（末尾不带点号）
47  export CLUSTER_DNS_DOMAIN="cluster.local"
48  # 将二进制目录 /opt/k8s/bin 加到 PATH 中
49  export PATH=/opt/k8s/bin:$PATH
```

**请根据IP进行修改**

分发环境变量脚本

```
1  source environment.sh
2  for node_ip in ${NODE_IPS[@]}
3    do
4      echo ">>> ${node_ip}"
5      scp environment.sh root@${node_ip}:/opt/k8s/bin/
6      ssh root@${node_ip} "chmod +x /opt/k8s/bin/* "
7  done
```

# 二、k8s集群部署

**创建CA证书和秘钥**

为确保安全，kubernetes各个组件需要使用x509证书对通信进行加密和认证
CA(Certificate Authority)是自签名的根证书，用来签名后续创建的其他证书。本文章使用CloudFlare的
PKI工具cfssl创建所有证书。

> **注意:** 如果没有特殊指明，本文档的所有操作均在k8s-01节点执行，远程分发到其他节点

**安装cfssl工具集**

```
1  mkdir -p /opt/k8s/cert && cd /opt/k8s
2  wget https://pkg.cfssl.org/R1.2/cfssl_linux-amd64
3  mv cfssl_linux-amd64 /opt/k8s/bin/cfssl
4  wget https://pkg.cfssl.org/R1.2/cfssljson_linux-amd64
5  mv cfssljson_linux-amd64 /opt/k8s/bin/cfssljson
6  wget https://pkg.cfssl.org/R1.2/cfssl-certinfo_linux-amd64
7  mv cfssl-certinfo_linux-amd64 /opt/k8s/bin/cfssl-certinfo
8  chmod +x /opt/k8s/bin/*
9  export PATH=/opt/k8s/bin:$PATH
```

### 创建根证书 (CA)

> CA证书是集群所有节点共享的，只需要创建一个CA证书，后续创建的所有证书都是由它签名

### 创建配置文件

> CA配置文件用于配置根证书的使用场景(profile)和具体参数
> (usage、过期时间、服务端认证、客户端认证、加密等)

```
1  cd /opt/k8s/work
2  cat > ca-config.json <<EOF
3  {
4    "signing": {
5      "default": {
6        "expiry": "87600h"
7      },
8      "profiles": {
9        "kubernetes": {
10         "usages": [
11             "signing",
12             "key encipherment",
13             "server auth",
14             "client auth"
15         ],
16         "expiry": "87600h"
17       }
18     }
19   }
20 }
21 EOF
22
```

```
23  #####################
24  signing 表示该证书可用于签名其它证书，生成的ca.pem证书找中CA=TRUE
25  server auth 表示client可以用该证书对server提供的证书进行验证
26  client auth 表示server可以用该证书对client提供的证书进行验证
```

创建证书签名请求文件

```
 1  cd /opt/k8s/work
 2  cat > ca-csr.json <<EOF
 3  {
 4    "CN": "kubernetes",
 5    "key": {
 6      "algo": "rsa",
 7      "size": 2048
 8    },
 9    "names": [
10      {
11        "C": "CN",
12        "ST": "BeiJing",
13        "L": "BeiJing",
14        "O": "k8s",
15        "OU": "xuyuntech"
16      }
17    ],
18    "ca": {
19      "expiry": "876000h"
20    }
21  }
22  EOF
23
24  #####################
25  CN CommonName,kube-apiserver从证书中提取该字段作为请求的用户名(User Name),
     浏览器使用该字段验证网站是否合法
26  O Organization,kube-apiserver 从证书中提取该字段作为请求用户和所属组(Group)
27  kube-apiserver将提取的User、Group作为RBAC授权的用户和标识
```

生成CA证书和私钥

```
1 cd /opt/k8s/work
2 cfssl gencert -initca ca-csr.json | cfssljson -bare ca
3 ls ca*
```

**分发证书**

#将生成的CA证书、秘钥文件、配置文件拷贝到所有节点的/etc/kubernetes/cert目录下

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${NODE_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     ssh root@${node_ip} "mkdir -p /etc/kubernetes/cert"
7     scp ca*.pem ca-config.json root@${node_ip}:/etc/kubernetes/cert
8   done
```

**部署kubectl命令行工具**

kubectl默认从 `~/.kube/config` 读取kube-apiserver地址和认证信息。kube/config只需要部署一次，生成的kubeconfig文件是通用的
下载和解压kubectl

```
1 cd /opt/k8s/work
2 wget http://down.i4t.com/k8s1.14/kubernetes-client-linux-amd64.tar.gz
3 tar -xzvf kubernetes-client-linux-amd64.tar.gz
```

分发所有使用kubectl节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${NODE_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kubernetes/client/bin/kubectl root@${node_ip}:/opt/k8s/bin/
7     ssh root@${node_ip} "chmod +x /opt/k8s/bin/*"
8   done
```

**创建admin证书和私钥**

kubectl与apiserver https通信，apiserver对提供的证书进行认证和授权。kubectl作为集群的管理工具，需要被授予最高权限，这里创建具有最高权限的admin证书

创建证书签名请求

```
1  cd /opt/k8s/work
2  cat > admin-csr.json <<EOF
3  {
4    "CN": "admin",
5    "hosts": [],
6    "key": {
7      "algo": "rsa",
8      "size": 2048
9    },
10   "names": [
11     {
12       "C": "CN",
13       "ST": "BeiJing",
14       "L": "BeiJing",
15       "O": "system:masters",
16       "OU": "xuyuntech"
17     }
18   ]
19 }
20 EOF
21
22 ##################
23 ● O 为system:masters, kube-apiserver收到该证书后将请求的Group设置为system:
   masters
24 ● 预定的ClusterRoleBinding cluster-admin将Group system:masters与Role c
   luster-admin绑定, 该Role授予API的权限
25 ● 该证书只有被kubectl当做client证书使用, 所以hosts字段为空
```

**生成证书和私钥**

```
1 cd /opt/k8s/work
2 cfssl gencert -ca=/opt/k8s/work/ca.pem \
3   -ca-key=/opt/k8s/work/ca-key.pem \
4   -config=/opt/k8s/work/ca-config.json \
5   -profile=kubernetes admin-csr.json | cfssljson -bare admin
6 ls admin*
```

**创建kubeconfig文件**

kubeconfig为kubectl的配置文件，包含访问apiserver的所有信息，如apiserver地址、CA证书和自身使用的证书

```
 1 cd /opt/k8s/work
 2 source /opt/k8s/bin/environment.sh
 3 # 设置集群参数
 4 kubectl config set-cluster kubernetes \
 5   --certificate-authority=/opt/k8s/work/ca.pem \
 6   --embed-certs=true \
 7   --server=${KUBE_APISERVER} \
 8   --kubeconfig=kubectl.kubeconfig
 9 #设置客户端认证参数
10 kubectl config set-credentials admin \
11   --client-certificate=/opt/k8s/work/admin.pem \
12   --client-key=/opt/k8s/work/admin-key.pem \
13   --embed-certs=true \
14   --kubeconfig=kubectl.kubeconfig
15 # 设置上下文参数
16 kubectl config set-context kubernetes \
17   --cluster=kubernetes \
18   --user=admin \
19   --kubeconfig=kubectl.kubeconfig
20 # 设置默认上下文
21 kubectl config use-context kubernetes --kubeconfig=kubectl.kubeconfi
   g
22
23 ###############
24 --certificate-authority 验证kube-apiserver证书的根证书
25 --client-certificate、--client-key 刚生成的admin证书和私钥，连接kube-apis
   erver时使用
26 --embed-certs=true 将ca.pem和admin.pem证书嵌入到生成的kubectl.kubeconfig
```

文件中（如果不加入，写入的是证书文件路径，后续拷贝kubeconfig到其它机器时，还需要单
独拷贝证书）

分发kubeconfig文件

分发到所有使用kubectl命令的节点

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  for node_ip in ${NODE_IPS[@]}
4    do
5      echo ">>> ${node_ip}"
6      ssh root@${node_ip} "mkdir -p ~/.kube"
7      scp kubectl.kubeconfig root@${node_ip}:~/.kube/config
8    done
9
10 #保存文件名为~/.kube/config
```

**部署ETCD集群**

这里使用的ETCD为三节点高可用集群，步骤如下

- 下载和分发etcd二进制文件
- 创建etcd集群各节点的x509证书，用于加密客户端(如kubectl)与etcd集群、etcd集群之间的数据流
- 创建etcd的system unit文件，配置服务参数
- 检查集群工作状态
- etcd集群各节点的名称和IP如下
- k8s-01 192.168.0.50
- k8s-02 192.168.0.51
- k8s-03 192.168.0.52
- 注意: 没有特殊说明都在k8s-01节点操作

**下载和分发etcd二进制文件**

```
1  cd /opt/k8s/work
2  wget http://down.i4t.com/k8s1.14/etcd-v3.3.13-linux-amd64.tar.gz
3  tar -xvf etcd-v3.3.13-linux-amd64.tar.gz
```

分发二进制文件到集群节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${ETCD_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp etcd-v3.3.13-linux-amd64/etcd* root@${node_ip}:/opt/k8s/bin
7     ssh root@${node_ip} "chmod +x /opt/k8s/bin/*"
8   done
```

**创建etcd证书和私钥**

```
1 cd /opt/k8s/work
2 cat > etcd-csr.json <<EOF
3 {
4   "CN": "etcd",
5   "hosts": [
6     "127.0.0.1",
7     "192.168.0.50",
8     "192.168.0.51",
9     "192.168.0.52"
10   ],
11   "key": {
12     "algo": "rsa",
13     "size": 2048
14   },
15   "names": [
16     {
17       "C": "CN",
18       "ST": "BeiJing",
19       "L": "BeiJing",
20       "O": "k8s",
21       "OU": "xuyuntech"
22     }
23   ]
24 }
25 EOF
26 #host字段指定授权使用该证书的etcd节点IP或域名列表，需要将etcd集群的3个节点都添加
   其中
```

生成证书和私钥

```
1  cd /opt/k8s/work
2  cfssl gencert -ca=/opt/k8s/work/ca.pem \
3      -ca-key=/opt/k8s/work/ca-key.pem \
4      -config=/opt/k8s/work/ca-config.json \
5      -profile=kubernetes etcd-csr.json | cfssljson -bare etcd
6  ls etcd*pem
```

分发证书和私钥到etcd各个节点

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  for node_ip in ${ETCD_IPS[@]}
4    do
5      echo ">>> ${node_ip}"
6      ssh root@${node_ip} "mkdir -p /etc/etcd/cert"
7      scp etcd*.pem root@${node_ip}:/etc/etcd/cert/
8    done
```

创建etcd的启动文件（这里将配置文件也存放在启动文件里）

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  cat > etcd.service.template <<EOF
4  [Unit]
5  Description=Etcd Server
6  After=network.target
7  After=network-online.target
8  Wants=network-online.target
9  Documentation=https://github.com/coreos
10 [Service]
11 Type=notify
12 WorkingDirectory=${ETCD_DATA_DIR}
13 ExecStart=/opt/k8s/bin/etcd \\
14     --data-dir=${ETCD_DATA_DIR} \\
15     --wal-dir=${ETCD_WAL_DIR} \\
16     --name=##NODE_NAME## \\
```

```
17    --cert-file=/etc/etcd/cert/etcd.pem \\
18    --key-file=/etc/etcd/cert/etcd-key.pem \\
19    --trusted-ca-file=/etc/kubernetes/cert/ca.pem \\
20    --peer-cert-file=/etc/etcd/cert/etcd.pem \\
21    --peer-key-file=/etc/etcd/cert/etcd-key.pem \\
22    --peer-trusted-ca-file=/etc/kubernetes/cert/ca.pem \\
23    --peer-client-cert-auth \\
24    --client-cert-auth \\
25    --listen-peer-urls=https://##NODE_IP##:2380 \\
26    --initial-advertise-peer-urls=https://##NODE_IP##:2380 \\
27    --listen-client-urls=https://##NODE_IP##:2379,http://127.0.0.1:237
   9 \\
28    --advertise-client-urls=https://##NODE_IP##:2379 \\
29    --initial-cluster-token=etcd-cluster-0 \\
30    --initial-cluster=${ETCD_NODES} \\
31    --initial-cluster-state=new \\
32    --auto-compaction-mode=periodic \\
33    --auto-compaction-retention=1 \\
34    --max-request-bytes=33554432 \\
35    --quota-backend-bytes=6442450944 \\
36    --heartbeat-interval=250 \\
37    --election-timeout=2000
38 Restart=on-failure
39 RestartSec=5
40 LimitNOFILE=65536
41 [Install]
42 WantedBy=multi-user.target
43 EOF
```

配置说明（此处不需要修改任何配置）

- WorkDirectory、—data-dir 指定etcd工作目录和数据存储为${ETCD_DATA_DIR},需要在启动前创建这个目录（后面跟着我操作就可以，会有创建步骤）
- —wal-dir 指定wal目录，为了提高性能，一般使用SSD和—data-dir不同的盘
- —name 指定节点名称，当—initial-cluster-state值为new时，—name的参数值必须位于—initial-cluster列表中
- —cert-file、—key-file ETCD server与client通信时使用的证书和私钥
- —trusted-ca-file 签名client证书的CA证书，用于验证client证书
- —peer-cert-file、—peer-key-file ETCD与peer通信使用的证书和私钥
- —peer-trusted-ca-file 签名peer证书的CA证书，用于验证peer证书

为各个节点分发启动文件

```
1 #分发会将配置文件中的#替换成ip
2 cd /opt/k8s/work
3 source /opt/k8s/bin/environment.sh
4 for (( i=0; i < 3; i++ ))
5   do
6     sed -e "s/##NODE_NAME##/${MASTER_NAMES[i]}/" -e "s/##NODE_IP##/
  ${ETCD_IPS[i]}/" etcd.service.template > etcd-${ETCD_IPS[i]}.service
7   done
8 ls *.service
9 #NODE_NAMES 和 NODE_IPS 为相同长度的 bash 数组，分别为节点名称和对应的 IP；
```

分发生成的etcd启动文件到对应的服务器

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp etcd-${node_ip}.service root@${node_ip}:/etc/systemd/system/e
  tcd.service
7   done
```

重命名etcd启动文件并启动etcd服务
etcd首次进程启动会等待其他节点加入etcd集群，执行启动命令会卡顿一会，为正常现象

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     ssh root@${node_ip} "mkdir -p ${ETCD_DATA_DIR} ${ETCD_WAL_DIR}"
7     ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
  etcd && systemctl restart etcd " &
8   done
9   #这里我们创建了etcd的工作目录
```

检查启动结果

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     ssh root@${node_ip} "systemctl status etcd|grep Active"
7   done
```

正常状态

```
[root@abcdocker-k8s01 work]# for node_ip in ${MASTER_IPS[@]}
>   do
>     echo ">>> ${node_ip}"
>     ssh root@${node_ip} "systemctl status etcd|grep Active"
>   done
>>> 192.168.0.50
   Active: active (running) since Sun 2019-08-11 16:54:58 CST; 1min 40s ago
>>> 192.168.0.51
   Active: active (running) since Sun 2019-08-11 16:54:48 CST; 1min 49s ago
>>> 192.168.0.52
   Active: active (running) since Sun 2019-08-11 16:54:48 CST; 1min 50s ago
[root@abcdocker-k8s01 work]#
```

如果etcd集群状态不是active (running)，请使用下面命令查看etcd日志

```
1 journalctl -fu etcd
```

验证ETCD集群状态
不是完etcd集群后，在任一etcd节点执行下命令

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     ETCDCTL_API=3 /opt/k8s/bin/etcdctl \
7     --endpoints=https://${node_ip}:2379 \
8     --cacert=/etc/kubernetes/cert/ca.pem \
```

```
  9       --cert=/etc/etcd/cert/etcd.pem \
 10       --key=/etc/etcd/cert/etcd-key.pem endpoint health
 11    done
```

正常状态如下

```
                                              1. root@abcdocker-k8s01:/opt/k8s/work (ssh)
[root@abcdocker-k8s01 work]# cd /opt/k8s/work
[root@abcdocker-k8s01 work]# source /opt/k8s/bin/environment.sh
[root@abcdocker-k8s01 work]# for node_ip in ${MASTER_IPS[@]}
>    do
>      echo ">>> ${node_ip}"
>      ETCDCTL_API=3 /opt/k8s/bin/etcdctl \
>      --endpoints=https://${node_ip}:2379 \
>      --cacert=/etc/kubernetes/cert/ca.pem \
>      --cert=/etc/etcd/cert/etcd.pem \
>      --key=/etc/etcd/cert/etcd-key.pem endpoint health
>    done
>>> 192.168.0.50
https://192.168.0.50:2379 is healthy: successfully committed proposal: took = 2.928499ms
>>> 192.168.0.51
https://192.168.0.51:2379 is healthy: successfully committed proposal: took = 14.197405ms
>>> 192.168.0.52
https://192.168.0.52:2379 is healthy: successfully committed proposal: took = 2.07568ms
[root@abcdocker-k8s01 work]#
```

我们还可以通过下面命令查看当前etcd集群leader

```
 1  source /opt/k8s/bin/environment.sh
 2  ETCDCTL_API=3 /opt/k8s/bin/etcdctl \
 3    -w table --cacert=/etc/kubernetes/cert/ca.pem \
 4    --cert=/etc/etcd/cert/etcd.pem \
 5    --key=/etc/etcd/cert/etcd-key.pem \
 6    --endpoints=${ETCD_ENDPOINTS} endpoint status
```

正常状态如下

```
[root@abcdocker-k8s01 work]#        source /opt/k8s/bin/environment.sh
[root@abcdocker-k8s01 work]#        ETCDCTL_API=3 /opt/k8s/bin/etcdctl \
>       -w table --cacert=/etc/kubernetes/cert/ca.pem \
>       --cert=/etc/etcd/cert/etcd.pem \
>       --key=/etc/etcd/cert/etcd-key.pem \
>       --endpoints=${ETCD_ENDPOINTS} endpoint status
+--------------------------+------------------+---------+---------+-----------+-----------+------------+
|         ENDPOINT         |        ID        | VERSION | DB SIZE | IS LEADER | RAFT TERM | RAFT INDEX |
+--------------------------+------------------+---------+---------+-----------+-----------+------------+
| https://192.168.0.50:2379 | 623a14f1769c93ee | 3.3.13  |  16 kB  |   false   |     5     |     14     |
| https://192.168.0.51:2379 | 1efbd37ff8237f14 | 3.3.13  |  16 kB  |   false   |     5     |     14     |
| https://192.168.0.52:2379 | a66681c44a794e52 | 3.3.13  |  16 kB  |   true    |     5     |     14     |
+--------------------------+------------------+---------+---------+-----------+-----------+------------+
[root@abcdocker-k8s01 work]#
```

部署Flannel网络

Kubernetes要求集群内各个节点(包括master)能通过Pod网段互联互通，Flannel使用vxlan技术为各个节点创建一个互通的Pod网络，使用的端口为8472.第一次启动时，从etcd获取配置的Pod网络，为本节点分配一个未使用的地址段，然后创建flannel.1网络接口(也可能是其它名称)flannel将分配给自己的Pod网段信息写入 `/run/flannel/docker` 文件，docker后续使用这个文件中的环境变量设置Docker0网桥，从而从这个地址段为本节点的所有Pod容器分配IP

**下载分发flanneld二进制文件** (本次flanneld不使用Pod运行)

```
1 cd /opt/k8s/work
2 mkdir flannel
3 wget http://down.i4t.com/k8s1.14/flannel-v0.11.0-linux-amd64.tar.gz
4 tar -xzvf flannel-v0.11.0-linux-amd64.tar.gz -C flannel
```

分发二进制文件到所有集群的节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${NODE_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp flannel/{flanneld,mk-docker-opts.sh} root@${node_ip}:/opt/k8s/bin/
7     ssh root@${node_ip} "chmod +x /opt/k8s/bin/*"
8   done
```

**创建Flannel证书和私钥**

flanneld从etcd集群存取网段分配信息，而etcd集群开启了双向x509证书认证，所以需要为flannel生成证书和私钥

创建证书签名请求

```
1 cd /opt/k8s/work
2 cat > flanneld-csr.json <<EOF
3 {
4   "CN": "flanneld",
5   "hosts": [],
6   "key": {
7     "algo": "rsa",
8     "size": 2048
9   },
```

```
10    "names": [
11      {
12        "C": "CN",
13        "ST": "BeiJing",
14        "L": "BeiJing",
15        "O": "k8s",
16        "OU": "xuyuntech"
17      }
18    ]
19  }
20  EOF
```

生成证书和私钥

```
1  cfssl gencert -ca=/opt/k8s/work/ca.pem \
2    -ca-key=/opt/k8s/work/ca-key.pem \
3    -config=/opt/k8s/work/ca-config.json \
4    -profile=kubernetes flanneld-csr.json | cfssljson -bare flanneld
5  ls flanneld*pem
```

将生成的证书和私钥分发到所有节点

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  for node_ip in ${NODE_IPS[@]}
4    do
5      echo ">>> ${node_ip}"
6      ssh root@${node_ip} "mkdir -p /etc/flanneld/cert"
7      scp flanneld*.pem root@${node_ip}:/etc/flanneld/cert
8    done
```

向etcd写入Pod网段信息

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  etcdctl \
4    --endpoints=${ETCD_ENDPOINTS} \
```

```
5    --ca-file=/opt/k8s/work/ca.pem \
6    --cert-file=/opt/k8s/work/flanneld.pem \
7    --key-file=/opt/k8s/work/flanneld-key.pem \
8    mk ${FLANNEL_ETCD_PREFIX}/config '{"Network":"'${CLUSTER_CIDR}'",
  "SubnetLen": 21, "Backend": {"Type": "vxlan"}}'
```

注意:

flanneld当前版本 `v0.11.0` 不支持etcd v3，故使用etcd v2 API写入配置Key和网段数据；

写入的Pod网段${CLUSTER_CIDR}地址段(如/16)必须小于SubnetLen，必须与kube-controller-manager的—cluster-cidr参数一致

**创建flanneld的启动文件**

```
 1  cd /opt/k8s/work
 2  source /opt/k8s/bin/environment.sh
 3  cat > flanneld.service << EOF
 4  [Unit]
 5  Description=Flanneld overlay address etcd agent
 6  After=network.target
 7  After=network-online.target
 8  Wants=network-online.target
 9  After=etcd.service
10  Before=docker.service
11  [Service]
12  Type=notify
13  ExecStart=/opt/k8s/bin/flanneld \\
14    -etcd-cafile=/etc/kubernetes/cert/ca.pem \\
15    -etcd-certfile=/etc/flanneld/cert/flanneld.pem \\
16    -etcd-keyfile=/etc/flanneld/cert/flanneld-key.pem \\
17    -etcd-endpoints=${ETCD_ENDPOINTS} \\
18    -etcd-prefix=${FLANNEL_ETCD_PREFIX} \\
19    -iface=${IFACE} \\
20    -ip-masq
21  ExecStartPost=/opt/k8s/bin/mk-docker-opts.sh -k DOCKER_NETWORK_OPTIO
  NS -d /run/flannel/docker
22  Restart=always
23  RestartSec=5
24  StartLimitInterval=0
25  [Install]
26  WantedBy=multi-user.target
```

```
27  RequiredBy=docker.service
28  EOF
```

- mk-docker-opts.sh 脚本将分配给 flanneld 的 Pod 子网段信息写入 /run/flannel/docker 文件，后续 docker 启动时使用这个文件中的环境变量配置 docker0 网桥；
- flanneld 使用系统缺省路由所在的接口与其它节点通信，对于有多个网络接口（如内网和公网）的节点，可以用 –iface 参数指定通信接口；
- flanneld 运行时需要 root 权限；
- –ip-masq: flanneld 为访问 Pod 网络外的流量设置 SNAT 规则，同时将传递给 Docker 的变量 –ip-masq（/run/flannel/docker 文件中）设置为 false，这样 Docker 将不再创建 SNAT 规则；Docker 的 –ip-masq 为 true 时，创建的 SNAT 规则比较"暴力"：将所有本节点 Pod 发起的、访问非 docker0 接口的请求做 SNAT，这样访问其他节点 Pod 的请求来源 IP 会被设置为 flannel.1 接口的 IP，导致目的 Pod 看不到真实的来源 Pod IP。 flanneld 创建的 SNAT 规则比较温和，只对访问非 Pod 网段的请求做 SNAT。

**分发启动文件到所有节点**

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  for node_ip in ${NODE_IPS[@]}
4    do
5      echo ">>> ${node_ip}"
6      scp flanneld.service root@${node_ip}:/etc/systemd/system/
7    done
```

启动flanneld服务

```
1  source /opt/k8s/bin/environment.sh
2  for node_ip in ${NODE_IPS[@]}
3    do
4      echo ">>> ${node_ip}"
5      ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
   flanneld && systemctl restart flanneld"
6    done
```

检查启动结果

```
1  source /opt/k8s/bin/environment.sh
```

```
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "systemctl status flanneld|grep Active"
6   done
```

检查分配给flanneld的Pod网段信息

```
1 source /opt/k8s/bin/environment.sh
2 etcdctl \
3   --endpoints=${ETCD_ENDPOINTS} \
4   --ca-file=/etc/kubernetes/cert/ca.pem \
5   --cert-file=/etc/flanneld/cert/flanneld.pem \
6   --key-file=/etc/flanneld/cert/flanneld-key.pem \
7   get ${FLANNEL_ETCD_PREFIX}/config
```

查看已分配的Pod子网网段列表

```
1 source /opt/k8s/bin/environment.sh
2 etcdctl \
3   --endpoints=${ETCD_ENDPOINTS} \
4   --ca-file=/etc/kubernetes/cert/ca.pem \
5   --cert-file=/etc/flanneld/cert/flanneld.pem \
6   --key-file=/etc/flanneld/cert/flanneld-key.pem \
7   ls ${FLANNEL_ETCD_PREFIX}/subnets
```

查看某Pod网段对应节点IP和flannel接口地址

```
1 source /opt/k8s/bin/environment.sh
2 etcdctl \
3   --endpoints=${ETCD_ENDPOINTS} \
4   --ca-file=/etc/kubernetes/cert/ca.pem \
5   --cert-file=/etc/flanneld/cert/flanneld.pem \
6   --key-file=/etc/flanneld/cert/flanneld-key.pem \
7   get ${FLANNEL_ETCD_PREFIX}/subnets/172.30.16.0-21
8
9   #后面节点IP需要根据我们查出来的地址进行修改
```

查看节点flannel网络信息

```
1 ip addr show
```

flannel.1网卡的地址为分配的pod自网段的第一个个IP (.0)，且是/32的地址

```
1 ip addr show|grep flannel.1
```

到其它节点 Pod 网段请求都被转发到 flannel.1 网卡；
flanneld 根据 etcd 中子网段的信息，如 `${FLANNEL_ETCD_PREFIX}/subnets/172.30.80.0-21`，来决定进请求发送给哪个节点的互联 IP；
验证各节点能通过 Pod 网段互通
在各节点上部署 flannel 后，检查是否创建了 flannel 接口(名称可能为 flannel0、flannel.0、flannel.1 等):

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh ${node_ip} "/usr/sbin/ip addr show flannel.1|grep -w inet"
6   done
```

**kube-apiserver 高可用**
- 使用Nginx 4层透明代理功能实现k8s节点(master节点和worker节点)高可用访问kube-apiserver的步骤
- 控制节点的kube-controller-manager、kube-scheduler是多实例部署，所以只要一个实例正常，就可以保证集群高可用
- 集群内的Pod使用k8s服务域名kubernetes访问kube-apiserver，kube-dns会自动解析多个kube-apiserver节点的IP，所以也是高可用的
- 在每个Nginx进程，后端对接多个apiserver实例，Nginx对他们做健康检查和负载均衡
- kubelet、kube-proxy、controller-manager、schedule通过本地nginx (监听我们vip 192.158.0.54)访问kube-apiserver，从而实现kube-apiserver高可用

**下载编译nginx** (k8s-01安装就可以，后面有拷贝步骤)

```
1 cd /opt/k8s/work
2 wget http://down.i4t.com/k8s1.14/nginx-1.15.3.tar.gz
```

```
 3  tar -xzvf nginx-1.15.3.tar.gz
 4  #编译
 5  cd /opt/k8s/work/nginx-1.15.3
 6  mkdir nginx-prefix
 7  ./configure --with-stream --without-http --prefix=$(pwd)/nginx-prefi
    x --without-http_uwsgi_module
 8  make && make install
 9  #############
10  --without-http_scgi_module --without-http_fastcgi_module
11  --with-stream: 开启 4 层透明转发(TCP Proxy)功能;
12  --without-xxx: 关闭所有其他功能，这样生成的动态链接二进制程序依赖最小;
```

查看 nginx 动态链接的库：

```
 1  ldd ./nginx-prefix/sbin/nginx
```

由于只开启了 4 层透明转发功能，所以除了依赖 libc 等操作系统核心 lib 库外，没有对其它 lib 的依赖
(如 libz、libssl 等)，这样可以方便部署到各版本操作系统中
创建目录结构

```
 1  cd /opt/k8s/work
 2  source /opt/k8s/bin/environment.sh
 3  for node_ip in ${NODE_IPS[@]}
 4    do
 5      echo ">>> ${node_ip}"
 6      mkdir -p /opt/k8s/kube-nginx/{conf,logs,sbin}
 7    done
```

拷贝二进制程序到其他主机 (有报错执行2遍就可以)

```
 1  cd /opt/k8s/work
 2  source /opt/k8s/bin/environment.sh
 3  for node_ip in ${NODE_IPS[@]}
 4    do
 5      echo ">>> ${node_ip}"
 6      scp /opt/k8s/work/nginx-1.15.3/nginx-prefix/sbin/nginx  root@${n
    ode_ip}:/opt/k8s/kube-nginx/sbin/kube-nginx
```

```
7       ssh root@${node_ip} "chmod a+x /opt/k8s/kube-nginx/sbin/*"
8       ssh root@${node_ip} "mkdir -p /opt/k8s/kube-nginx/{conf,logs,sbi
  n}"
9       sleep 3
10    done
```

配置Nginx文件，开启4层透明转发

```
1 cd /opt/k8s/work
2 cat > kube-nginx.conf <<EOF
3 worker_processes 1;
4 events {
5     worker_connections  1024;
6 }
7 stream {
8     upstream backend {
9         hash $remote_addr consistent;
10        server 192.168.0.50:6443      max_fails=3 fail_timeout=30s
  ;
11        server 192.168.0.51:6443      max_fails=3 fail_timeout=30s
  ;
12        server 192.168.0.52:6443      max_fails=3 fail_timeout=30s
  ;
13    }
14    server {
15        listen *:8443;
16        proxy_connect_timeout 1s;
17        proxy_pass backend;
18    }
19 }
20 EOF
21 #这里需要将server替换我们自己的地址
```

分发配置文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
```

```
4    do
5      echo ">>> ${node_ip}"
6      scp kube-nginx.conf  root@${node_ip}:/opt/k8s/kube-nginx/conf/kub
  e-nginx.conf
7    done
```

配置Nginx启动文件

```
1  cd /opt/k8s/work
2  cat > kube-nginx.service <<EOF
3  [Unit]
4  Description=kube-apiserver nginx proxy
5  After=network.target
6  After=network-online.target
7  Wants=network-online.target
8  [Service]
9  Type=forking
10 ExecStartPre=/opt/k8s/kube-nginx/sbin/kube-nginx -c /opt/k8s/kube-ng
   inx/conf/kube-nginx.conf -p /opt/k8s/kube-nginx -t
11 ExecStart=/opt/k8s/kube-nginx/sbin/kube-nginx -c /opt/k8s/kube-nginx
   /conf/kube-nginx.conf -p /opt/k8s/kube-nginx
12 ExecReload=/opt/k8s/kube-nginx/sbin/kube-nginx -c /opt/k8s/kube-ngin
   x/conf/kube-nginx.conf -p /opt/k8s/kube-nginx -s reload
13 PrivateTmp=true
14 Restart=always
15 RestartSec=5
16 StartLimitInterval=0
17 LimitNOFILE=65536
18 [Install]
19 WantedBy=multi-user.target
20 EOF
```

分发nginx启动文件

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  for node_ip in ${MASTER_IPS[@]}
4    do
```

```
5      echo ">>> ${node_ip}"
6      scp kube-nginx.service  root@${node_ip}:/etc/systemd/system/
7    done
```

启动 kube-nginx 服务

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
   kube-nginx && systemctl start kube-nginx"
7   done
```

检查 kube-nginx 服务运行状态

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     ssh root@${node_ip} "systemctl status kube-nginx |grep 'Active:'"
7   done
```

**KeepLived 部署**
前面我们也说了，高可用方案需要一个VIP，供集群内部访问
在所有**master**节点安装keeplived

```
1 yum  install -y keepalived
```

接下来我们要配置keeplive服务
`192.168.0.50`配置

```
1 cat > /etc/keepalived/keepalived.conf <<EOF
2 ! Configuration File for keepalived
```

```
 3 global_defs {
 4     router_id 192.168.0.50
 5 }
 6 vrrp_script chk_nginx {
 7     script "/etc/keepalived/check_port.sh 8443"
 8     interval 2
 9     weight -20
10 }
11 vrrp_instance VI_1 {
12     state MASTER
13     interface eth0
14     virtual_router_id 251
15     priority 100
16     advert_int 1
17     mcast_src_ip 192.168.0.50
18     nopreempt
19     authentication {
20         auth_type PASS
21         auth_pass 11111111
22     }
23     track_script {
24         chk_nginx
25     }
26     virtual_ipaddress {
27         192.168.0.54
28     }
29 }
30 EOF
31
32 ## 192.168.0.50 为节点IP, 192.168.0.54位VIP
```

将配置拷贝到其他节点，并替换相关IP

```
1 for node_ip in 192.168.0.50 192.168.0.51 192.168.0.52
2   do
3     echo ">>> ${node_ip}"
4     scp  /etc/keepalived/keepalived.conf $node_ip:/etc/keepalived/kee
  palived.conf
5   done
```

```
6  #替换IP
7  ssh root@192.168.0.51 sed -i 's#192.168.0.50#192.168.0.51#g'  /etc/ke
   epalived/keepalived.conf
8  ssh root@192.168.0.52 sed -i 's#192.168.0.50#192.168.0.52#g'  /etc/ke
   epalived/keepalived.conf
9  #192.168.0.50不替换是因为已经修改好了
```

创建健康检查脚本

```
1  vim  /opt/check_port.sh
2  CHK_PORT=$1
3   if [ -n "$CHK_PORT" ];then
4          PORT_PROCESS=`ss -lt|grep $CHK_PORT|wc -l`
5          if [ $PORT_PROCESS -eq 0 ];then
6                  echo "Port $CHK_PORT Is Not Used,End."
7                  exit 1
8          fi
9   else
10         echo "Check Port Cant Be Empty!"
11  fi
```

启动keeplived

```
1  for NODE in k8s-01 k8s-02 k8s-03; do
2      echo "--- $NODE ---"
3      scp -r /opt/check_port.sh  $NODE:/etc/keepalived/
4      ssh $NODE 'systemctl enable --now keepalived'
5  done
```

启动完毕后ping 192.168.0.54 (VIP)

```
1  [root@k8s03 ~]# ping 192.168.0.54
2  PING 192.168.0.54 (192.168.0.54) 56(84) bytes of data.
3  64 bytes from 192.168.0.54: icmp_seq=1 ttl=64 time=0.055 ms
4  ^C
5  --- 192.168.0.54 ping statistics ---
6  1 packets transmitted, 1 received, 0% packet loss, time 0ms
```

```
 7 rtt min/avg/max/mdev = 0.055/0.055/0.055/0.000 ms
 8 #如果没有启动，请检查原因。 ps -ef|grep keep 检查是否启动成功
 9 #没有启动成功，请执行下面命令，从新启动。启动成功vip肯定就通了
10 systemctl start keepalived
```

**部署master节点**

kubernetes master节点运行组件如下:kube-apiserver、kube-scheduler、kube-controller-manager、kube-nginx

- kube-apiserver、kube-scheduler、kube-controller-manager均以多实例模式运行
- kube-scheduler和kube-controller-manager会自动选举一个leader实例，其他实例处于阻塞模式，当leader挂了后，重新选举产生的leader，从而保证服务可用性
- kube-apiserver是无状态的，需要通过kube-nginx进行代理访问，从而保证服务可用性

> 以下操作都在K8s-01操作

下载kubernetes二进制包，并分发到所有master节点

```
1 cd /opt/k8s/work
2 wget http://down.i4t.com/k8s1.14/kubernetes-server-linux-amd64.tar.gz
3 tar -xzvf kubernetes-server-linux-amd64.tar.gz
4 cd kubernetes
5 tar -xzvf  kubernetes-src.tar.gz
```

cd /opt/k8s/work
将压缩包的文件拷贝到所有master节点上

```
 1 cd /opt/k8s/work
 2 source /opt/k8s/bin/environment.sh
 3 for node_ip in ${MASTER_IPS[@]}
 4   do
 5     echo ">>> ${node_ip}"
 6     scp kubernetes/server/bin/kube-apiserver root@${node_ip}:/opt/k8
   s/bin/
 7     scp kubernetes/server/bin/{apiextensions-apiserver,cloud-control
   ler-manager,kube-controller-manager,kube-proxy,kube-scheduler,kubead
   m,kubectl,kubelet,mounter} root@${node_ip}:/opt/k8s/bin/
 8     ssh root@${node_ip} "chmod +x /opt/k8s/bin/*"
 9   done
10 #同时将kubelet kube-proxy拷贝到所有节点
```

```
11  cd /opt/k8s/work
12  source /opt/k8s/bin/environment.sh
13  for node_ip in ${NODE_IPS[@]}
14    do
15      echo ">>> ${node_ip}"
16      scp kubernetes/server/bin/{kubelet,kube-proxy} root@${node_ip}:/
    opt/k8s/bin/
17      ssh root@${node_ip} "chmod +x /opt/k8s/bin/*"
18    done
```

**创建Kubernetes 证书和私钥**

创建签证签名请求

```
 1  cd /opt/k8s/work
 2  source /opt/k8s/bin/environment.sh
 3  cat > kubernetes-csr.json <<EOF
 4  {
 5    "CN": "kubernetes",
 6    "hosts": [
 7      "127.0.0.1",
 8      "192.168.0.50",
 9      "192.168.0.51",
10      "192.168.0.52",
11      "192.168.0.54",
12      "10.254.0.1",
13      "kubernetes",
14      "kubernetes.default",
15      "kubernetes.default.svc",
16      "kubernetes.default.svc.cluster",
17      "kubernetes.default.svc.cluster.local."
18    ],
19    "key": {
20      "algo": "rsa",
21      "size": 2048
22    },
23    "names": [
24      {
25        "C": "CN",
```

```
26          "ST": "BeiJing",
27          "L": "BeiJing",
28          "O": "k8s",
29          "OU": "xuyuntech"
30      }
31    ]
32 }
33 EOF
34 #需要将集群的所有IP及VIP添加进去
35 #如果要添加注意最后的逗号，不要忘记添加，否则下一步报错
```

hosts 字段指定授权使用该证书的IP和域名列表，这里列出了master节点IP、kubernetes服务的IP和域名 kubernetes serviceIP是apiserver自动创建的，一般是－service－cluster－ip－range参数指定的网段的第一个IP
$ kubectl get svc kubernetes

```
1 $ kubectl get svc kubernetes
2 NAME          TYPE        CLUSTER-IP     EXTERNAL-IP    PORT(S)    AGE
3 kubernetes    ClusterIP   10.254.0.1                   443/TCP    31d
4
5 #目前我们是看不到
```

生成证书和私钥

```
1 cfssl gencert -ca=/opt/k8s/work/ca.pem \
2       -ca-key=/opt/k8s/work/ca-key.pem \
3       -config=/opt/k8s/work/ca-config.json \
4       -profile=kubernetes kubernetes-csr.json | cfssljson -bare kuber
  netes
5     ls kubernetes*pem
```

将生成的证书和私钥文件拷贝到所有master节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
```

```
4   do
5     echo ">>> ${node_ip}"
6     ssh root@${node_ip} "mkdir -p /etc/kubernetes/cert"
7     scp kubernetes*.pem root@${node_ip}:/etc/kubernetes/cert/
8   done
```

创建加密配置文件

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  cat > encryption-config.yaml <<EOF
4  kind: EncryptionConfig
5  apiVersion: v1
6  resources:
7    - resources:
8        - secrets
9      providers:
10       - aescbc:
11           keys:
12             - name: key1
13               secret: ${ENCRYPTION_KEY}
14       - identity: {}
15 EOF
```

将加密配置文件拷贝到master节点的 `/etc/kubernetes` 目录下

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  for node_ip in ${MASTER_IPS[@]}
4    do
5      echo ">>> ${node_ip}"
6      scp encryption-config.yaml root@${node_ip}:/etc/kubernetes/
7    done
```

创建审计策略文件

```
1  cd /opt/k8s/work
```

```
 2  source /opt/k8s/bin/environment.sh
 3  cat > audit-policy.yaml <<EOF
 4  apiVersion: audit.k8s.io/v1beta1
 5  kind: Policy
 6  rules:
 7    # The following requests were manually identified as high-volume
      and low-risk, so drop them.
 8    - level: None
 9      resources:
10        - group: ""
11          resources:
12            - endpoints
13            - services
14            - services/status
15      users:
16        - 'system:kube-proxy'
17      verbs:
18        - watch
19    - level: None
20      resources:
21        - group: ""
22          resources:
23            - nodes
24            - nodes/status
25      userGroups:
26        - 'system:nodes'
27      verbs:
28        - get
29    - level: None
30      namespaces:
31        - kube-system
32      resources:
33        - group: ""
34          resources:
35            - endpoints
36      users:
37        - 'system:kube-controller-manager'
38        - 'system:kube-scheduler'
39        - 'system:serviceaccount:kube-system:endpoint-controller'
40      verbs:
```

```yaml
41          - get
42          - update
43      - level: None
44        resources:
45          - group: ""
46            resources:
47              - namespaces
48              - namespaces/status
49              - namespaces/finalize
50        users:
51          - 'system:apiserver'
52        verbs:
53          - get
54      # Don't log HPA fetching metrics.
55      - level: None
56        resources:
57          - group: metrics.k8s.io
58        users:
59          - 'system:kube-controller-manager'
60        verbs:
61          - get
62          - list
63      # Don't log these read-only URLs.
64      - level: None
65        nonResourceURLs:
66          - '/healthz*'
67          - /version
68          - '/swagger*'
69      # Don't log events requests.
70      - level: None
71        resources:
72          - group: ""
73            resources:
74              - events
75      # node and pod status calls from nodes are high-volume and can be
    large, don't log responses for expected updates from nodes
76      - level: Request
77        omitStages:
78          - RequestReceived
79        resources:
```

```
80        - group: ""
81          resources:
82            - nodes/status
83            - pods/status
84      users:
85        - kubelet
86        - 'system:node-problem-detector'
87        - 'system:serviceaccount:kube-system:node-problem-detector'
88      verbs:
89        - update
90        - patch
91    - level: Request
92      omitStages:
93        - RequestReceived
94      resources:
95        - group: ""
96          resources:
97            - nodes/status
98            - pods/status
99      userGroups:
100       - 'system:nodes'
101     verbs:
102       - update
103       - patch
104   # deletecollection calls can be large, don't log responses for ex
    pected namespace deletions
105   - level: Request
106     omitStages:
107       - RequestReceived
108     users:
109       - 'system:serviceaccount:kube-system:namespace-controller'
110     verbs:
111       - deletecollection
112   # Secrets, ConfigMaps, and TokenReviews can contain sensitive & b
    inary data,
113   # so only log at the Metadata level.
114   - level: Metadata
115     omitStages:
116       - RequestReceived
117     resources:
```

```yaml
118        - group: ""
119          resources:
120            - secrets
121            - configmaps
122        - group: authentication.k8s.io
123          resources:
124            - tokenreviews
125    # Get repsonses can be large; skip them.
126    - level: Request
127      omitStages:
128        - RequestReceived
129      resources:
130        - group: ""
131        - group: admissionregistration.k8s.io
132        - group: apiextensions.k8s.io
133        - group: apiregistration.k8s.io
134        - group: apps
135        - group: authentication.k8s.io
136        - group: authorization.k8s.io
137        - group: autoscaling
138        - group: batch
139        - group: certificates.k8s.io
140        - group: extensions
141        - group: metrics.k8s.io
142        - group: networking.k8s.io
143        - group: policy
144        - group: rbac.authorization.k8s.io
145        - group: scheduling.k8s.io
146        - group: settings.k8s.io
147        - group: storage.k8s.io
148      verbs:
149        - get
150        - list
151        - watch
152    # Default level for known APIs
153    - level: RequestResponse
154      omitStages:
155        - RequestReceived
156      resources:
157        - group: ""
```

```
158        - group: admissionregistration.k8s.io
159        - group: apiextensions.k8s.io
160        - group: apiregistration.k8s.io
161        - group: apps
162        - group: authentication.k8s.io
163        - group: authorization.k8s.io
164        - group: autoscaling
165        - group: batch
166        - group: certificates.k8s.io
167        - group: extensions
168        - group: metrics.k8s.io
169        - group: networking.k8s.io
170        - group: policy
171        - group: rbac.authorization.k8s.io
172        - group: scheduling.k8s.io
173        - group: settings.k8s.io
174        - group: storage.k8s.io
175    # Default level for all other requests.
176    - level: Metadata
177      omitStages:
178        - RequestReceived
179 EOF
```

分发审计策略文件：

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp audit-policy.yaml root@${node_ip}:/etc/kubernetes/audit-polic
  y.yaml
7   done
```

创建证书签名请求

```
1 cat > proxy-client-csr.json <<EOF
2 {
```

```
 3     "CN": "aggregator",
 4     "hosts": [],
 5     "key": {
 6       "algo": "rsa",
 7       "size": 2048
 8     },
 9     "names": [
10       {
11         "C": "CN",
12         "ST": "BeiJing",
13         "L": "BeiJing",
14         "O": "k8s",
15         "OU": "xuyuntech"
16       }
17     ]
18 }
19 EOF
```

- CN名称需要位于kube-apiserver的—requestherader-allowed-names参数中，否则后续访问
  metrics时会提示权限不足

生成证书和私钥

```
1 cfssl gencert -ca=/etc/kubernetes/cert/ca.pem \
2   -ca-key=/etc/kubernetes/cert/ca-key.pem  \
3   -config=/etc/kubernetes/cert/ca-config.json  \
4   -profile=kubernetes proxy-client-csr.json | cfssljson -bare proxy-c
  lient
5 ls proxy-client*.pem
```

将生成的证书和私钥文件拷贝到master节点

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${MASTER_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     scp proxy-client*.pem root@${node_ip}:/etc/kubernetes/cert/
6   done
```

创建kube-apiserver启动文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 cat > kube-apiserver.service.template <<EOF
4 [Unit]
5 Description=Kubernetes API Server
6 Documentation=https://github.com/GoogleCloudPlatform/kubernetes
7 After=network.target
8 [Service]
9 WorkingDirectory=${K8S_DIR}/kube-apiserver
10 ExecStart=/opt/k8s/bin/kube-apiserver \\
11   --advertise-address=##NODE_IP## \\
12   --default-not-ready-toleration-seconds=360 \\
13   --default-unreachable-toleration-seconds=360 \\
14   --feature-gates=DynamicAuditing=true \\
15   --max-mutating-requests-inflight=2000 \\
16   --max-requests-inflight=4000 \\
17   --default-watch-cache-size=200 \\
18   --delete-collection-workers=2 \\
19   --encryption-provider-config=/etc/kubernetes/encryption-config.yam
  l \\
20   --etcd-cafile=/etc/kubernetes/cert/ca.pem \\
21   --etcd-certfile=/etc/kubernetes/cert/kubernetes.pem \\
22   --etcd-keyfile=/etc/kubernetes/cert/kubernetes-key.pem \\
23   --etcd-servers=${ETCD_ENDPOINTS} \\
24   --bind-address=##NODE_IP## \\
25   --secure-port=6443 \\
26   --tls-cert-file=/etc/kubernetes/cert/kubernetes.pem \\
27   --tls-private-key-file=/etc/kubernetes/cert/kubernetes-key.pem \\
28   --insecure-port=0 \\
29   --audit-dynamic-configuration \\
30   --audit-log-maxage=15 \\
31   --audit-log-maxbackup=3 \\
32   --audit-log-maxsize=100 \\
33   --audit-log-truncate-enabled \\
34   --audit-log-path=${K8S_DIR}/kube-apiserver/audit.log \\
35   --audit-policy-file=/etc/kubernetes/audit-policy.yaml \\
36   --profiling \\
```

```
37    --anonymous-auth=false \\
38    --client-ca-file=/etc/kubernetes/cert/ca.pem \\
39    --enable-bootstrap-token-auth \\
40    --requestheader-allowed-names="aggregator" \\
41    --requestheader-client-ca-file=/etc/kubernetes/cert/ca.pem \\
42    --requestheader-extra-headers-prefix="X-Remote-Extra-" \\
43    --requestheader-group-headers=X-Remote-Group \\
44    --requestheader-username-headers=X-Remote-User \\
45    --service-account-key-file=/etc/kubernetes/cert/ca.pem \\
46    --authorization-mode=Node,RBAC \\
47    --runtime-config=api/all=true \\
48    --enable-admission-plugins=NodeRestriction \\
49    --allow-privileged=true \\
50    --apiserver-count=3 \\
51    --event-ttl=168h \\
52    --kubelet-certificate-authority=/etc/kubernetes/cert/ca.pem \\
53    --kubelet-client-certificate=/etc/kubernetes/cert/kubernetes.pem
   \\
54    --kubelet-client-key=/etc/kubernetes/cert/kubernetes-key.pem \\
55    --kubelet-https=true \\
56    --kubelet-timeout=10s \\
57    --proxy-client-cert-file=/etc/kubernetes/cert/proxy-client.pem \\
58    --proxy-client-key-file=/etc/kubernetes/cert/proxy-client-key.pem
   \\
59    --service-cluster-ip-range=${SERVICE_CIDR} \\
60    --service-node-port-range=${NODE_PORT_RANGE} \\
61    --logtostderr=true \\
62    --v=2
63 Restart=on-failure
64 RestartSec=10
65 Type=notify
66 LimitNOFILE=65536
67 [Install]
68 WantedBy=multi-user.target
69 EOF
```

参数配置说明

```
 1 --advertise-address：apiserver 对外通告的 IP（kubernetes 服务后端节点 IP）
```

```
   ;
 2 --default-*-toleration-seconds: 设置节点异常相关的阈值;
 3 --max-*-requests-inflight: 请求相关的最大阈值;
 4 --etcd-*: 访问 etcd 的证书和 etcd 服务器地址;
 5 --experimental-encryption-provider-config: 指定用于加密 etcd 中 secret
   的配置;
 6 --bind-address: https 监听的 IP, 不能为 127.0.0.1, 否则外界不能访问它的安全
   端口 6443;
 7 --secret-port: https 监听端口;
 8 --insecure-port=0: 关闭监听 http 非安全端口(8080);
 9 --tls-*-file: 指定 apiserver 使用的证书、私钥和 CA 文件;
10 --audit-*: 配置审计策略和审计日志文件相关的参数;
11 --client-ca-file: 验证 client (kue-controller-manager、kube-scheduler
   、kubelet、kube-proxy 等)请求所带的证书;
12 --enable-bootstrap-token-auth: 启用 kubelet bootstrap 的 token 认证;
13 --requestheader-*: kube-apiserver 的 aggregator layer 相关的配置参数, pr
   oxy-client & HPA 需要使用;
14 --requestheader-client-ca-file: 用于签名 --proxy-client-cert-file 和 --
   proxy-client-key-file 指定的证书;在启用了 metric aggregator 时使用;
15 --requestheader-allowed-names: 不能为空, 值为逗号分割的 --proxy-client-ce
   rt-file 证书的 CN 名称, 这里设置为 "aggregator";
16 --service-account-key-file: 签名 ServiceAccount Token 的公钥文件, kube-c
   ontroller-manager 的 --service-account-private-key-file 定私钥文件, 两者
   配对使用;
17 --runtime-config=api/all=true: 启用所有版本的 APIs, 如 autoscaling/v2al
   pha1;
18 --authorization-mode=Node,RBAC、--anonymous-auth=false: 开启 Node 和
   RBAC 授权模式, 拒绝未授权的请求;
19 --enable-admission-plugins: 启用一些默认关闭的 plugins;
20 --allow-privileged: 运行执行 privileged 权限的容器;
21 --apiserver-count=3: 指定 apiserver 实例的数量;
22 --event-ttl: 指定 events 的保存时间;
23 --kubelet-: 如果指定, 则使用 https 访问 kubelet APIs;需要为证书对应的用户(上
   面 kubernetes.pem 证书的用户为 kubernetes) 用户定义 RBAC 规则, 否则访问 kub
   elet API 时提示未授权;
24 --proxy-client-*: apiserver 访问 metrics-server 使用的证书;
25 --service-cluster-ip-range: 指定 Service Cluster IP 地址段;
26 --service-node-port-range: 指定 NodePort 的端口范围;
27 如果 kube-apiserver 机器没有运行 kube-proxy, 则还需要添加 --enable-aggrega
   tor-routing=true 参数;
```

```
28   关于 --requestheader-XXX 相关参数，参考：
29 https://github.com/kubernetes-incubator/apiserver-builder/blob/maste
   r/docs/concepts/auth.md
30 https://docs.bitnami.com/kubernetes/how-to/configure-autoscaling-cus
   tom-metrics/
```

注意: requestheader-client-ca-file指定的CA证书，必须具有client auth and server auth
如果—requestheader-allowed-names为空，或者—proxy-client-cert-file证书的CN名称不在
allowed-names中，则后续查看node或者Pods的metrics失败
**为各个节点创建和分发kube-apiserver启动文件**
替换模板文件的变量，为各个节点生成启动文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for (( i=0; i < 3; i++ ))   #这里是三个节点所以为3,请根据实际情况修改,后边不在提
  示,同理
4   do
5     sed -e "s/##NODE_NAME##/${MASTER_NAMES[i]}/" -e "s/##NODE_IP##/
  ${MASTER_IPS[i]}/" kube-apiserver.service.template > kube-apiserver-
  ${MASTER_IPS[i]}.service
6   done
7 ls kube-apiserver*.service
```

cd /opt/k8s/work
分发apiserver启动文件

```
 1 cd /opt/k8s/work
 2 source /opt/k8s/bin/environment.sh
 3 for node_ip in ${MASTER_IPS[@]}
 4   do
 5     echo ">>> ${node_ip}"
 6     scp kube-apiserver-${node_ip}.service root@${node_ip}:/etc/syste
   md/system/kube-apiserver.service
 7   done
 8 启动apiserver
 9 source /opt/k8s/bin/environment.sh
10 for node_ip in ${MASTER_IPS[@]}
11   do
```

```
12      echo ">>> ${node_ip}"
13      ssh root@${node_ip} "mkdir -p ${K8S_DIR}/kube-apiserver"
14      ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
    kube-apiserver && systemctl restart kube-apiserver"
15    done
```

检查服务是否正常

```
1  source /opt/k8s/bin/environment.sh
2  for node_ip in ${MASTER_IPS[@]}
3    do
4      echo ">>> ${node_ip}"
5      ssh root@${node_ip} "systemctl status kube-apiserver |grep 'Activ
    e:'"
6    done
```

确保状态为active (running)，否则查看日志，确认原因
journalctl -u kube-apiserver
打印kube-apiserver写入etcd数据

```
1  source /opt/k8s/bin/environment.sh
2  ETCDCTL_API=3 etcdctl \
3      --endpoints=${ETCD_ENDPOINTS} \
4      --cacert=/opt/k8s/work/ca.pem \
5      --cert=/opt/k8s/work/etcd.pem \
6      --key=/opt/k8s/work/etcd-key.pem \
7      get /registry/ --prefix --keys-only
```

检查kube-apiserver监听的端口

```
1  netstat -lntup|grep kube
2  tcp        0      0 192.168.0.50:6443        0.0.0.0:*                 L
   ISTEN      11739/kube-apiserve
```

检查集群信息

```
1 $ kubectl cluster-info
```

```
2  Kubernetes master is running at https://192.168.0.54:8443
3  To further debug and diagnose cluster problems, use 'kubectl cluster
   -info dump'.
4  $ kubectl get all --all-namespaces
5  NAMESPACE   NAME                      TYPE        CLUSTER-IP    EXTERNAL-I
   P   PORT(S)   AGE
6  default     service/kubernetes   ClusterIP   10.254.0.1           44
   3/TCP    3m5s
7  $ kubectl get componentstatuses
8  NAME                    STATUS      MESSAGE
   ERROR
9  scheduler               Unhealthy   Get http://127.0.0.1:10251/healthz:
   dial tcp 127.0.0.1:10251: connect: connection refused
10 controller-manager      Unhealthy   Get http://127.0.0.1:10252/healthz:
   dial tcp 127.0.0.1:10252: connect: connection refused
11 etcd-2                  Healthy     {"health":"true"}
12 etcd-0                  Healthy     {"health":"true"}
13 etcd-1                  Healthy     {"health":"true"}
```

如果提示有报错，请检查 `~/.kube/config` 以及配置证书是否有问题

**授权kube-apiserver访问kubelet API的权限**

在执行kubectl命令时，apiserver会将请求转发到kubelet的https端口。这里定义的RBAC规则，授权apiserver使用的证书(kubernetes.pem)用户名(CN:kubernetes)访问kubelet API的权限

```
1  kubectl create clusterrolebinding kube-apiserver:kubelet-apis --clust
   errole=system:kubelet-api-admin --user kubernetes
```

**部署高可用kube-controller-manager集群**

该集群包含三个节点，启动后通过竞争选举机制产生一个leader节点，其他节点为阻塞状态。当leader节点不可用时，阻塞节点将会在此选举产生新的leader，从而保证服务的高可用。为保证通信安全，这里采用x509证书和私钥，kube-controller-manager在与apiserver的安全端口（http 10252）通信使用；
创建kube-controller-manager证书和私钥
创建证书签名请求

```
1  cd /opt/k8s/work
2  cat > kube-controller-manager-csr.json <<EOF
3  {
```

```
 4        "CN": "system:kube-controller-manager",
 5        "key": {
 6            "algo": "rsa",
 7            "size": 2048
 8        },
 9        "hosts": [
10          "127.0.0.1",
11          "192.168.0.50",
12          "192.168.0.51",
13          "192.168.0.52"
14        ],
15        "names": [
16          {
17            "C": "CN",
18            "ST": "BeiJing",
19            "L": "BeiJing",
20            "O": "system:kube-controller-manager",
21            "OU": "xuyuntech"
22          }
23        ]
24 }
25 EOF
```

#这里的IP地址为master ip
- host列表包含所有的kube-controller-manager节点IP(VIP不需要输入)
- CN和O均为system:kube-controller-manager，kubernetes内置的ClusterRoleBindings
  system:kube-controller-manager赋予kube-controller-manager工作所需权限

生成证书和私钥

```
1 cd /opt/k8s/work
2 cfssl gencert -ca=/opt/k8s/work/ca.pem \
3   -ca-key=/opt/k8s/work/ca-key.pem \
4   -config=/opt/k8s/work/ca-config.json \
5   -profile=kubernetes kube-controller-manager-csr.json | cfssljson -b
  are kube-controller-manager
6 ls kube-controller-manager*pem
```

将生成的证书和私钥分发到所有master节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kube-controller-manager*.pem root@${node_ip}:/etc/kubernetes/
  cert/
7   done
```

创建和分发kubeconfig文件
#kube-controller-manager使用kubeconfig文件访问apiserver
#该文件提供了apiserver地址、嵌入的CA证书和kube-controller-manager证书

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 kubectl config set-cluster kubernetes \
4   --certificate-authority=/opt/k8s/work/ca.pem \
5   --embed-certs=true \
6   --server=${KUBE_APISERVER} \
7   --kubeconfig=kube-controller-manager.kubeconfig
8 kubectl config set-credentials system:kube-controller-manager \
9   --client-certificate=kube-controller-manager.pem \
10   --client-key=kube-controller-manager-key.pem \
11   --embed-certs=true \
12   --kubeconfig=kube-controller-manager.kubeconfig
13 kubectl config set-context system:kube-controller-manager \
14   --cluster=kubernetes \
15   --user=system:kube-controller-manager \
16   --kubeconfig=kube-controller-manager.kubeconfig
17 kubectl config use-context system:kube-controller-manager --kubeconf
   ig=kube-controller-manager.kubeconfig
```

分发kubeconfig到所有master节点

```
1 cd /opt/k8s/work
```

```
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kube-controller-manager.kubeconfig root@${node_ip}:/etc/kuber
  netes/
7   done
```

cd /opt/k8s/work

创建kube-controller-manager启动文件

```
 1 cd /opt/k8s/work
 2 source /opt/k8s/bin/environment.sh
 3 cat > kube-controller-manager.service.template <<EOF
 4 [Unit]
 5 Description=Kubernetes Controller Manager
 6 Documentation=https://github.com/GoogleCloudPlatform/kubernetes
 7 [Service]
 8 WorkingDirectory=${K8S_DIR}/kube-controller-manager
 9 ExecStart=/opt/k8s/bin/kube-controller-manager \\
10   --profiling \\
11   --cluster-name=kubernetes \\
12   --controllers=*,bootstrapsigner,tokencleaner \\
13   --kube-api-qps=1000 \\
14   --kube-api-burst=2000 \\
15   --leader-elect \\
16   --use-service-account-credentials\\
17   --concurrent-service-syncs=2 \\
18   --bind-address=0.0.0.0 \\
19   #--secure-port=10252 \\
20   --tls-cert-file=/etc/kubernetes/cert/kube-controller-manager.pem
  \\
21   --tls-private-key-file=/etc/kubernetes/cert/kube-controller-manage
  r-key.pem \\
22   #--port=0 \\
23   --authentication-kubeconfig=/etc/kubernetes/kube-controller-manage
  r.kubeconfig \\
24   --client-ca-file=/etc/kubernetes/cert/ca.pem \\
25   --requestheader-allowed-names="" \\
```

```
26    --requestheader-client-ca-file=/etc/kubernetes/cert/ca.pem \\
27    --requestheader-extra-headers-prefix="X-Remote-Extra-" \\
28    --requestheader-group-headers=X-Remote-Group \\
29    --requestheader-username-headers=X-Remote-User \\
30    --authorization-kubeconfig=/etc/kubernetes/kube-controller-manager
   .kubeconfig \\
31    --cluster-signing-cert-file=/etc/kubernetes/cert/ca.pem \\
32    --cluster-signing-key-file=/etc/kubernetes/cert/ca-key.pem \\
33    --experimental-cluster-signing-duration=876000h \\
34    --horizontal-pod-autoscaler-sync-period=10s \\
35    --concurrent-deployment-syncs=10 \\
36    --concurrent-gc-syncs=30 \\
37    --node-cidr-mask-size=24 \\
38    --service-cluster-ip-range=${SERVICE_CIDR} \\
39    --pod-eviction-timeout=6m \\
40    --terminated-pod-gc-threshold=10000 \\
41    --root-ca-file=/etc/kubernetes/cert/ca.pem \\
42    --service-account-private-key-file=/etc/kubernetes/cert/ca-key.pem
   \\
43    --kubeconfig=/etc/kubernetes/kube-controller-manager.kubeconfig \\
44    --logtostderr=true \\
45    --v=2
46 Restart=on-failure
47 RestartSec=5
48 [Install]
49 WantedBy=multi-user.target
50 EOF
```

cd /opt/k8s/work

参数解释

- —port=0: 关闭监听非安全端口（http），同时 —address 参数无效，—bind-address 参数有效；
- —secure-port=10252、—bind-address=0.0.0.0: 在所有网络接口监听 10252 端口的 https /metrics 请求；
- —kubeconfig: 指定 kubeconfig 文件路径，kube-controller-manager 使用它连接和验证 kube-apiserver；
- —authentication-kubeconfig 和 —authorization-kubeconfig: kube-controller-manager 使用它连接 apiserver，对 client 的请求进行认证和授权。kube-controller-manager 不再使用 —tls-ca-file 对请求 https metrics 的 Client 证书进行校验。如果没有配置这两个 kubeconfig 参数，则 client 连接 kube-controller-manager https 端口的请求会被拒绝(提示权限不足)。

- —cluster-signing-*-file：签名 TLS Bootstrap 创建的证书；
- —experimental-cluster-signing-duration：指定 TLS Bootstrap 证书的有效期；
- —root-ca-file：放置到容器 ServiceAccount 中的 CA 证书，用来对 kube-apiserver 的证书进行校验；
- —service-account-private-key-file：签名 ServiceAccount 中 Token 的私钥文件，必须和 kube-apiserver 的 —service-account-key-file 指定的公钥文件配对使用；
- —service-cluster-ip-range ： 指定 Service Cluster IP 网段，必须和 kube-apiserver 中的同名参数一致；
- —leader-elect=true：集群运行模式，启用选举功能；被选为 leader 的节点负责处理工作，其它节点为阻塞状态；
- —controllers=*,bootstrapsigner,tokencleaner：启用的控制器列表，tokencleaner 用于自动清理过期的 Bootstrap token；
- —horizontal-pod-autoscaler-*：custom metrics 相关参数，支持 autoscaling/v2alpha1；
- —tls-cert-file、—tls-private-key-file：使用 https 输出 metrics 时使用的 Server 证书和秘钥；
- —use-service-account-credentials=true: kube-controller-manager 中各 controller 使用 serviceaccount 访问 kube-apiserver；

**替换启动文件，并分发脚本**

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for (( i=0; i < 3; i++ ))
4   do
5     sed -e "s/##NODE_NAME##/${MASTER_NAMES[i]}/" -e "s/##NODE_IP##/
  ${MASTER_IPS[i]}/" kube-controller-manager.service.template > kube-co
  ntroller-manager-${MASTER_IPS[i]}.service
6   done
7 ls kube-controller-manager*.service
```

分发到所有master节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kube-controller-manager-${node_ip}.service root@${node_ip}:/e
  tc/systemd/system/kube-controller-manager.service
7   done
```

启动服务

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${MASTER_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "mkdir -p ${K8S_DIR}/kube-controller-manager"
6     ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
  kube-controller-manager && systemctl restart kube-controller-manage
  r"
7   done
```

检查运行状态

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${MASTER_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "systemctl status kube-controller-manager|gre
  p Active"
6   done
```

检查服务状态

```
1 netstat -lnpt | grep kube-cont
2 tcp6       0       0 :::10252                      :::*                          L
  ISTEN      13279/kube-controll
3 tcp6       0       0 :::10257                      :::*                          L
  ISTEN      13279/kube-controll
```

**kube-controller-manager 创建权限**

ClusteRole system:kube-controller-manager的权限太小，只能创建secret、serviceaccount等资源，
将controller的权限分散到ClusterRole system:controller:xxx中

```
1 $ kubectl describe clusterrole system:kube-controller-manager
2 Name:         system:kube-controller-manager
3 Labels:       kubernetes.io/bootstrapping=rbac-defaults
```

```
4 Annotations:  rbac.authorization.kubernetes.io/autoupdate: true
5 PolicyRule:
6   Resources                                    Non-Resource URLs  Reso
  urce Names  Verbs
7   ---------                                    ----------------   ----
  ----------  -----
8   secrets                                      []                 []
  [create delete get update]
9   endpoints                                    []                 []
  [create get update]
10   serviceaccounts                             []                 []
  [create get update]
11   events                                      []                 []
  [create patch update]
12   tokenreviews.authentication.k8s.io          []                 []
  [create]
13   subjectaccessreviews.authorization.k8s.io   []                 []
  [create]
14   configmaps                                  []                 []
  [get]
15   namespaces                                  []                 []
  [get]
16   *.*                                         []                 []
  [list watch]
```

需要在 kube-controller-manager 的启动参数中添加 —use-service-account-credentials=true 参数，这样 main controller 会为各 controller 创建对应的 ServiceAccount XXX-controller。内置的 ClusterRoleBinding system:controller:XXX 将赋予各 XXX-controller ServiceAccount 对应的 ClusterRole system:controller:XXX 权限。

```
1 $ kubectl get clusterrole|grep controller
2 system:controller:attachdetach-controller
  22m
3 system:controller:certificate-controller
  22m
4 system:controller:clusterrole-aggregation-controller
  22m
5 system:controller:cronjob-controller
```

```
     22m
  6 system:controller:daemon-set-controller
     22m
  7 system:controller:deployment-controller
     22m
  8 system:controller:disruption-controller
     22m
  9 system:controller:endpoint-controller
     22m
 10 system:controller:expand-controller
     22m
 11 system:controller:generic-garbage-collector
     22m
 12 system:controller:horizontal-pod-autoscaler
     22m
 13 system:controller:job-controller
     22m
 14 system:controller:namespace-controller
     22m
 15 system:controller:node-controller
     22m
 16 system:controller:persistent-volume-binder
     22m
 17 system:controller:pod-garbage-collector
     22m
 18 system:controller:pv-protection-controller
     22m
 19 system:controller:pvc-protection-controller
     22m
 20 system:controller:replicaset-controller
     22m
 21 system:controller:replication-controller
     22m
 22 system:controller:resourcequota-controller
     22m
 23 system:controller:route-controller
     22m
 24 system:controller:service-account-controller
     22m
 25 system:controller:service-controller
```

```
22m
26 system:controller:statefulset-controller
    22m
27 system:controller:ttl-controller
    22m
28 system:kube-controller-manager
    22m
```

以 deployment controller 为例：

```
 1 $ kubectl describe clusterrole system:controller:deployment-controll
   er
 2 Name:          system:controller:deployment-controller
 3 Labels:        kubernetes.io/bootstrapping=rbac-defaults
 4 Annotations:   rbac.authorization.kubernetes.io/autoupdate: true
 5 PolicyRule:
 6   Resources                            Non-Resource URLs  Resource Nam
   es  Verbs
 7   ---------                            -----------------  -----------
   --  -----
 8   replicasets.apps                     []                 []
   [create delete get list patch update watch]
 9   replicasets.extensions               []                 []
   [create delete get list patch update watch]
10   events                               []                 []
   [create patch update]
11   pods                                 []                 []
   [get list update watch]
12   deployments.apps                     []                 []
   [get list update watch]
13   deployments.extensions               []                 []
   [get list update watch]
14   deployments.apps/finalizers          []                 []
   [update]
15   deployments.apps/status              []                 []
   [update]
16   deployments.extensions/finalizers    []                 []
   [update]
17   deployments.extensions/status        []                 []
```

```
[update]
```

查看当前的 leader

```
1  $ kubectl get endpoints kube-controller-manager --namespace=kube-sys
   tem  -o yaml
2  apiVersion: v1
3  kind: Endpoints
4  metadata:
5    annotations:
6      control-plane.alpha.kubernetes.io/leader: '{"holderIdentity":"k8
   s01_56e187ed-bc5b-11e9-b4a3-000c291b8bf5","leaseDurationSeconds":1
   5,"acquireTime":"2019-08-11T17:13:29Z","renewTime":"2019-08-11T17:1
   9:06Z","leaderTransitions":0}'
7    creationTimestamp: "2019-08-11T17:13:29Z"
8    name: kube-controller-manager
9    namespace: kube-system
10   resourceVersion: "848"
11   selfLink: /api/v1/namespaces/kube-system/endpoints/kube-controller
   -manager
12   uid: 56e64ea1-bc5b-11e9-b77e-000c291b8bf5
```

**部署高可用kube-scheduler**

创建 kube-scheduler 证书和私钥

创建证书签名请求：

```
1  cd /opt/k8s/work
2  cat > kube-scheduler-csr.json <<EOF
3  {
4      "CN": "system:kube-scheduler",
5      "hosts": [
6        "127.0.0.1",
7        "192.168.0.50",
8        "192.168.0.51",
9        "192.168.0.52"
10     ],
11     "key": {
```

```
12          "algo": "rsa",
13          "size": 2048
14        },
15      "names": [
16        {
17          "C": "CN",
18          "ST": "BeiJing",
19          "L": "BeiJing",
20          "O": "system:kube-scheduler",
21          "OU": "xuyuntech"
22        }
23      ]
24  }
25  EOF
```

- hosts 列表包含所有 kube-scheduler 节点 IP；
- CN 和 O 均为 system:kube-scheduler，kubernetes 内置的 ClusterRoleBindings system:kube-scheduler 将赋予 kube-scheduler 工作所需的权限；

生成证书和私钥：

```
1 cd /opt/k8s/work
2 cfssl gencert -ca=/opt/k8s/work/ca.pem \
3   -ca-key=/opt/k8s/work/ca-key.pem \
4   -config=/opt/k8s/work/ca-config.json \
5   -profile=kubernetes kube-scheduler-csr.json | cfssljson -bare kube-scheduler
6 ls kube-scheduler*pem
```

将生成的证书和私钥分发到所有 master 节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kube-scheduler*.pem root@${node_ip}:/etc/kubernetes/cert/
7   done
```

创建和分发 kubeconfig 文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 kubectl config set-cluster kubernetes \
4   --certificate-authority=/opt/k8s/work/ca.pem \
5   --embed-certs=true \
6   --server=${KUBE_APISERVER} \
7   --kubeconfig=kube-scheduler.kubeconfig
8 kubectl config set-credentials system:kube-scheduler \
9   --client-certificate=kube-scheduler.pem \
10   --client-key=kube-scheduler-key.pem \
11   --embed-certs=true \
12   --kubeconfig=kube-scheduler.kubeconfig
13 kubectl config set-context system:kube-scheduler \
14   --cluster=kubernetes \
15   --user=system:kube-scheduler \
16   --kubeconfig=kube-scheduler.kubeconfig
17 kubectl config use-context system:kube-scheduler --kubeconfig=kube-s
  cheduler.kubeconfig
```

分发 kubeconfig 到所有 master 节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kube-scheduler.kubeconfig root@${node_ip}:/etc/kubernetes/
7   done
```

创建 kube-scheduler 配置文件

```
1 cd /opt/k8s/work
2 cat >kube-scheduler.yaml.template <<EOF
3 apiVersion: kubescheduler.config.k8s.io/v1alpha1
4 kind: KubeSchedulerConfiguration
5 bindTimeoutSeconds: 600
```

```
 6 clientConnection:
 7   burst: 200
 8   kubeconfig: "/etc/kubernetes/kube-scheduler.kubeconfig"
 9   qps: 100
10 enableContentionProfiling: false
11 enableProfiling: true
12 hardPodAffinitySymmetricWeight: 1
13 healthzBindAddress: 127.0.0.1:10251
14 leaderElection:
15   leaderElect: true
16 metricsBindAddress: ##NODE_IP##:10251
17 EOF
```

- —kubeconfig：指定 kubeconfig 文件路径，kube-scheduler 使用它连接和验证 kube-apiserver；
- —leader-elect=true：集群运行模式，启用选举功能；被选为 leader 的节点负责处理工作，其它节点为阻塞状态；

替换模板文件中的变量：

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for (( i=0; i < 3; i++ ))
4   do
5     sed -e "s/##NODE_NAME##/${NODE_NAMES[i]}/" -e "s/##NODE_IP##/${NODE_IPS[i]}/" kube-scheduler.yaml.template > kube-scheduler-${NODE_IPS[i]}.yaml
6   done
7 ls kube-scheduler*.yaml
```

分发 kube-scheduler 配置文件到所有 master 节点：

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kube-scheduler-${node_ip}.yaml root@${node_ip}:/etc/kubernetes/kube-scheduler.yaml
```

```
7    done
```

创建kube-scheduler启动文件

```
 1 cd /opt/k8s/work
 2 cat > kube-scheduler.service.template <<EOF
 3 [Unit]
 4 Description=Kubernetes Scheduler
 5 Documentation=https://github.com/GoogleCloudPlatform/kubernetes
 6 [Service]
 7 WorkingDirectory=${K8S_DIR}/kube-scheduler
 8 ExecStart=/opt/k8s/bin/kube-scheduler \\
 9   --config=/etc/kubernetes/kube-scheduler.yaml \\
10   --bind-address=##NODE_IP## \\
11   --secure-port=10259 \\
12   --port=0 \\
13   --tls-cert-file=/etc/kubernetes/cert/kube-scheduler.pem \\
14   --tls-private-key-file=/etc/kubernetes/cert/kube-scheduler-key.pem
   \\
15   --authentication-kubeconfig=/etc/kubernetes/kube-scheduler.kubecon
   fig \\
16   --client-ca-file=/etc/kubernetes/cert/ca.pem \\
17   --requestheader-allowed-names="" \\
18   --requestheader-client-ca-file=/etc/kubernetes/cert/ca.pem \\
19   --requestheader-extra-headers-prefix="X-Remote-Extra-" \\
20   --requestheader-group-headers=X-Remote-Group \\
21   --requestheader-username-headers=X-Remote-User \\
22   --authorization-kubeconfig=/etc/kubernetes/kube-scheduler.kubeconf
   ig \\
23   --logtostderr=true \\
24   --v=2
25 Restart=always
26 RestartSec=5
27 StartLimitInterval=0
28 [Install]
29 WantedBy=multi-user.target
30 EOF
```

分发配置文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for (( i=0; i < 3; i++ ))
4   do
5     sed -e "s/##NODE_NAME##/${NODE_NAMES[i]}/" -e "s/##NODE_IP##/${NODE_IPS[i]}/" kube-scheduler.service.template > kube-scheduler-${NODE_IPS[i]}.service
6   done
7 ls kube-scheduler*.service
```

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${MASTER_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     scp kube-scheduler-${node_ip}.service root@${node_ip}:/etc/systemd/system/kube-scheduler.service
7   done
```

启动kube-scheduler

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${MASTER_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "mkdir -p ${K8S_DIR}/kube-scheduler"
6     ssh root@${node_ip} "systemctl daemon-reload && systemctl enable kube-scheduler && systemctl restart kube-scheduler"
7 done
```

检查服务运行状态

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${MASTER_IPS[@]}
```

```
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "systemctl status kube-scheduler|grep Active"
6   done
```

查看输出的 metrics

注意：以下命令在 kube-scheduler 节点上执行。

kube-scheduler 监听 10251 和 10251 端口：

10251：接收 http 请求，非安全端口，不需要认证授权；

10259：接收 https 请求，安全端口，需要认证授权；

两个接口都对外提供 /metrics 和 /healthz 的访问。

```
 1 curl -s http://192.168.0.50:10251/metrics|head
 2 # HELP apiserver_audit_event_total Counter of audit events generated
   and sent to the audit backend.
 3 # TYPE apiserver_audit_event_total counter
 4 apiserver_audit_event_total 0
 5 # HELP apiserver_audit_requests_rejected_total Counter of apiserver
    requests rejected due to an error in audit logging backend.
 6 # TYPE apiserver_audit_requests_rejected_total counter
 7 apiserver_audit_requests_rejected_total 0
 8 # HELP apiserver_client_certificate_expiration_seconds Distribution
    of the remaining lifetime on the certificate used to authenticate a
   request.
 9 # TYPE apiserver_client_certificate_expiration_seconds histogram
10 apiserver_client_certificate_expiration_seconds_bucket{le="0"} 0
11 apiserver_client_certificate_expiration_seconds_bucket{le="1800"} 0
```

查看当前leader

```
 1 $ kubectl get endpoints kube-scheduler --namespace=kube-system  -o y
   aml
 2 apiVersion: v1
 3 kind: Endpoints
 4 metadata:
 5   annotations:
 6     control-plane.alpha.kubernetes.io/leader: '{"holderIdentity":"k8
   s01_72210df0-bc5d-11e9-9ca8-000c291b8bf5","leaseDurationSeconds":1
```

```
  5,"acquireTime":"2019-08-11T17:28:35Z","renewTime":"2019-08-11T17:3
  1:06Z","leaderTransitions":0}'
7   creationTimestamp: "2019-08-11T17:28:35Z"
8   name: kube-scheduler
9   namespace: kube-system
10   resourceVersion: "1500"
11   selfLink: /api/v1/namespaces/kube-system/endpoints/kube-scheduler
12   uid: 72bcd72f-bc5d-11e9-b77e-000c291b8bf5
```

**work节点安装**

kubernetes work节点运行如下组件: >docker、kubelet、kube-proxy、flanneld、kube-nginx
前面已经安装flanneld这就不在安装了

**安装依赖包**

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "yum install -y epel-release"
6     ssh root@${node_ip} "yum install -y conntrack ipvsadm ntp ntpdate
  ipset jq iptables curl sysstat libseccomp && modprobe ip_vs "
7   done
```

**部署Docker组件**

我们在所有节点安装docker，这里使用阿里云的yum安装

> Docker步骤需要在所有节点安装

```
1 yum install -y yum-utils device-mapper-persistent-data lvm2
2 yum-config-manager --add-repo http://mirrors.aliyun.com/docker-ce/lin
  ux/centos/docker-ce.repo
3 yum makecache fast
4 yum -y install docker-ce
```

创建配置文件

```
 1  mkdir -p /etc/docker/
 2  cat > /etc/docker/daemon.json <<EOF
 3  {
 4    "exec-opts": ["native.cgroupdriver=systemd"],
 5    "registry-mirrors": ["https://hjvrgh7a.mirror.aliyuncs.com"],
 6    "log-driver": "json-file",
 7    "log-opts": {
 8      "max-size": "100m"
 9    },
10    "storage-driver": "devicemapper"
11  }
12  EOF
```

#这里配置当时镜像加速器,可以不进行配置，但是建议配置

要添加我们harbor仓库需要在添加下面一行

 "insecure-registries": ["harbor.i4t.com"],

默认docker hub需要https协议，使用上面配置不需要配置https

修改Docker启动参数

这里需要在所有的节点上修改docker配置！！

```
[Unit]
Description=Docker Application Container Engine
Documentation=https://docs.docker.com
BindsTo=containerd.service
After=network-online.target firewalld.service containerd.service
Wants=network-online.target
Requires=docker.socket

[Service]
Type=notify
# the default is not to use systemd for cgroups because the delegate issues still
# exists and systemd currently does not support the cgroup feature set required
# for containers run by docker
ExecStart=/usr/bin/dockerd $DOCKER_NETWORK_OPTIONS -H fd:// --containerd=/run/containerd/containerd.sock
EnvironmentFile=-/run/flannel/docker
ExecReload=/bin/kill -s HUP $MAINPID
TimeoutSec=0
RestartSec=2
Restart=always
```

```
1  EnvironmentFile=-/run/flannel/docker
2  ExecStart=/usr/bin/dockerd $DOCKER_NETWORK_OPTIONS -H fd:// --contain
   erd=/run/containerd/containerd.sock
```

完整配置如下

```
$ cat /usr/lib/systemd/system/docker.service
[Unit]
Description=Docker Application Container Engine
Documentation=https://docs.docker.com
BindsTo=containerd.service
After=network-online.target firewalld.service containerd.service
Wants=network-online.target
Requires=docker.socket
[Service]
Type=notify
# the default is not to use systemd for cgroups because the delegate
issues still
# exists and systemd currently does not support the cgroup feature s
et required
# for containers run by docker
ExecStart=/usr/bin/dockerd  $DOCKER_NETWORK_OPTIONS -H fd:// --conta
inerd=/run/containerd/containerd.sock
EnvironmentFile=-/run/flannel/docker
ExecReload=/bin/kill -s HUP $MAINPID
TimeoutSec=0
RestartSec=2
Restart=always
# Note that StartLimit* options were moved from "Service" to "Unit"
 in systemd 229.
# Both the old, and new location are accepted by systemd 229 and up,
so using the old location
# to make them work for either version of systemd.
StartLimitBurst=3
# Note that StartLimitInterval was renamed to StartLimitIntervalSec
 in systemd 230.
# Both the old, and new name are accepted by systemd 230 and up, so
 using the old name to make
# this option work for either version of systemd.
StartLimitInterval=60s
# Having non-zero Limit*s causes performance problems due to account
ing overhead
# in the kernel. We recommend using cgroups to do container-local ac
counting.
LimitNOFILE=infinity
```

```
31 LimitNPROC=infinity
32 LimitCORE=infinity
33 # Comment TasksMax if your systemd version does not support it.
34 # Only systemd 226 and above support this option.
35 TasksMax=infinity
36 # set delegate yes so that systemd does not reset the cgroups of doc
   ker containers
37 Delegate=yes
38 # kill only the docker process, not all processes in the cgroup
39 KillMode=process
40 [Install]
41 WantedBy=multi-user.target
```

启动 docker 服务

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
   docker && systemctl restart docker"
6   done
```

检查服务运行状态

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "systemctl status docker|grep Active"
6   done
```

检查 docker0 网桥

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
```

```
4    echo ">>> ${node_ip}"
5    ssh root@${node_ip} "/usr/sbin/ip addr show flannel.1 && /usr/sbi
  n/ip addr show docker0"
6  done
```

查看 docker 的状态信息

```
1 docker info
```

#查看docker版本以及存储引擎是否是overlay2

> 以上Docker步骤，有很多需要进入每台服务器进行修改配置文件！！

**部署kubelet组件**

kubelet运行在每个worker节点上，接收kube-apiserver发送的请求，管理Pod容器，执行交互命令
kubelet启动时自动向kube-apiserver注册节点信息，内置的cAdivsor统计和监控节点的资源使用资源情
况。为确保安全，部署时关闭了kubelet的非安全http端口，对请求进行认证和授权，拒绝未授权的访问
创建kubelet bootstrap kubeconfig文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_name in ${NODE_NAMES[@]}
4   do
5     echo ">>> ${node_name}"
6     # 创建 token
7     export BOOTSTRAP_TOKEN=$(kubeadm token create \
8       --description kubelet-bootstrap-token \
9       --groups system:bootstrappers:${node_name} \
10      --kubeconfig ~/.kube/config)
11    # 设置集群参数
12    kubectl config set-cluster kubernetes \
13      --certificate-authority=/etc/kubernetes/cert/ca.pem \
14      --embed-certs=true \
15      --server=${KUBE_APISERVER} \
16      --kubeconfig=kubelet-bootstrap-${node_name}.kubeconfig
17    # 设置客户端认证参数
18    kubectl config set-credentials kubelet-bootstrap \
19      --token=${BOOTSTRAP_TOKEN} \
```

```
20        --kubeconfig=kubelet-bootstrap-${node_name}.kubeconfig
21      # 设置上下文参数
22      kubectl config set-context default \
23        --cluster=kubernetes \
24        --user=kubelet-bootstrap \
25        --kubeconfig=kubelet-bootstrap-${node_name}.kubeconfig
26      # 设置默认上下文
27      kubectl config use-context default --kubeconfig=kubelet-bootstra
   p-${node_name}.kubeconfig
28    done
```

- 向kubeconfig写入的是token，bootstrap结束后kube-controller-manager为kubelet创建client和server证书

查看kubeadm为各个节点创建的token

```
1 $ kubeadm token list --kubeconfig ~/.kube/config
2 TOKEN                      TTL      EXPIRES                      USAGE
  S                DESCRIPTION                EXTRA GROUPS
3 ds9td8.wazmxhtaznrweknk   23h      2019-08-13T01:54:57+08:00   authe
  ntication,signing   kubelet-bootstrap-token   system:bootstrappers:k8
  s-01
4 hy5ssz.4zi4e079ovxba52x   23h      2019-08-13T01:54:58+08:00   authe
  ntication,signing   kubelet-bootstrap-token   system:bootstrappers:k8
  s-03
5 pkkcl0.l7syoup3jedt7c3l   23h      2019-08-13T01:54:57+08:00   authe
  ntication,signing   kubelet-bootstrap-token   system:bootstrappers:k8
  s-02
6 tubfqq.mja239hszl4rmron   23h      2019-08-13T01:54:58+08:00   authe
  ntication,signing   kubelet-bootstrap-token   system:bootstrappers:k8
  s-04
```

- token有效期为1天，超期后将不能被用来bootstrap kubelet，且会被kube-controller-manager的token cleaner清理
- kube-apiserver接收kubelet的bootstrap token后，将请求的user设置为system:bootstrap; group设置为system:bootstrappers，后续将为这个group设置ClusterRoleBinding

查看各token关联的Secret

```
1 $ kubectl get secrets  -n kube-system|grep bootstrap-token
```

```
2 bootstrap-token-ds9td8                          bootstrap.kubernetes
  .io/token           7       3m15s
3 bootstrap-token-hy5ssz                          bootstrap.kubernetes
  .io/token           7       3m14s
4 bootstrap-token-pkkcl0                          bootstrap.kubernetes
  .io/token           7       3m15s
5 bootstrap-token-tubfqq                          bootstrap.kubernetes
  .io/token           7       3m14s
```

分发 bootstrap kubeconfig 文件到所有 worker 节点

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_name in ${NODE_NAMES[@]}
4   do
5     echo ">>> ${node_name}"
6     scp kubelet-bootstrap-${node_name}.kubeconfig root@${node_name}:/
  etc/kubernetes/kubelet-bootstrap.kubeconfig
7   done
```

创建和分发kubelet参数配置

```
 1 cd /opt/k8s/work
 2 source /opt/k8s/bin/environment.sh
 3 cat > kubelet-config.yaml.template <<EOF
 4 kind: KubeletConfiguration
 5 apiVersion: kubelet.config.k8s.io/v1beta1
 6 address: "##NODE_IP##"
 7 staticPodPath: ""
 8 syncFrequency: 1m
 9 fileCheckFrequency: 20s
10 httpCheckFrequency: 20s
11 staticPodURL: ""
12 port: 10250
13 readOnlyPort: 0
14 rotateCertificates: true
15 serverTLSBootstrap: true
16 authentication:
```

```yaml
17    anonymous:
18      enabled: false
19    webhook:
20      enabled: true
21    x509:
22      clientCAFile: "/etc/kubernetes/cert/ca.pem"
23  authorization:
24    mode: Webhook
25  registryPullQPS: 0
26  registryBurst: 20
27  eventRecordQPS: 0
28  eventBurst: 20
29  enableDebuggingHandlers: true
30  enableContentionProfiling: true
31  healthzPort: 10248
32  healthzBindAddress: "##NODE_IP##"
33  clusterDomain: "${CLUSTER_DNS_DOMAIN}"
34  clusterDNS:
35    - "${CLUSTER_DNS_SVC_IP}"
36  nodeStatusUpdateFrequency: 10s
37  nodeStatusReportFrequency: 1m
38  imageMinimumGCAge: 2m
39  imageGCHighThresholdPercent: 85
40  imageGCLowThresholdPercent: 80
41  volumeStatsAggPeriod: 1m
42  kubeletCgroups: ""
43  systemCgroups: ""
44  cgroupRoot: ""
45  cgroupsPerQOS: true
46  cgroupDriver: systemd
47  runtimeRequestTimeout: 10m
48  hairpinMode: promiscuous-bridge
49  maxPods: 220
50  podCIDR: "${CLUSTER_CIDR}"
51  podPidsLimit: -1
52  resolvConf: /etc/resolv.conf
53  maxOpenFiles: 1000000
54  kubeAPIQPS: 1000
55  kubeAPIBurst: 2000
56  serializeImagePulls: false
```

```
57 evictionHard:
58   memory.available:  "100Mi"
59 nodefs.available:  "10%"
60 nodefs.inodesFree: "5%"
61 imagefs.available: "15%"
62 evictionSoft: {}
63 enableControllerAttachDetach: true
64 failSwapOn: true
65 containerLogMaxSize: 20Mi
66 containerLogMaxFiles: 10
67 systemReserved: {}
68 kubeReserved: {}
69 systemReservedCgroup: ""
70 kubeReservedCgroup: ""
71 enforceNodeAllocatable: ["pods"]
72 EOF
```

- address: kubelet 安全端口（https，10250）监听的地址，不能为 127.0.0.1，否则 kube-apiserver、heapster 等不能调用 kubelet 的 API；
- readOnlyPort=0：关闭只读端口(默认 10255)，等效为未指定；
- authentication.anonymous.enabled：设置为 false，不允许匿名 访问 10250 端口；
- authentication.x509.clientCAFile：指定签名客户端证书的 CA 证书，开启 HTTP 证书认证；
- authentication.webhook.enabled=true：开启 HTTPs bearer token 认证；
- 对于未通过 x509 证书和 webhook 认证的请求(kube-apiserver 或其他客户端)，将被拒绝，提示 Unauthorized；
- authroization.mode=Webhook：kubelet 使用 SubjectAccessReview API 查询 kube-apiserver 某 user、group 是否具有操作资源的权限(RBAC)；
- featureGates.RotateKubeletClientCertificate、featureGates.RotateKubeletServerCertificate：自动 rotate 证书，证书的有效期取决于 kube-controller-manager 的 —experimental-cluster-signing-duration 参数；
- 需要 root 账户运行；

为各个节点创建和分发kubelet配置文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_ip in ${NODE_IPS[@]}
4   do
5     echo ">>> ${node_ip}"
6     sed -e "s/##NODE_IP##/${node_ip}/" kubelet-config.yaml.template >
```

```
    kubelet-config-${node_ip}.yaml.template
7       scp kubelet-config-${node_ip}.yaml.template root@${node_ip}:/etc/
    kubernetes/kubelet-config.yaml
8    done
```

创建和分发kubelet启动文件

```
1  cd /opt/k8s/work
2  source /opt/k8s/bin/environment.sh
3  cat > kubelet.service.template <<EOF
4  [Unit]
5  Description=Kubernetes Kubelet
6  Documentation=https://github.com/GoogleCloudPlatform/kubernetes
7  After=docker.service
8  Requires=docker.service
9  [Service]
10 WorkingDirectory=${K8S_DIR}/kubelet
11 ExecStart=/opt/k8s/bin/kubelet \\
12   --allow-privileged=true \\
13   --bootstrap-kubeconfig=/etc/kubernetes/kubelet-bootstrap.kubeconfi
   g \\
14   --cert-dir=/etc/kubernetes/cert \\
15   --cni-conf-dir=/etc/cni/net.d \\
16   --container-runtime=docker \\
17   --container-runtime-endpoint=unix:///var/run/dockershim.sock \\
18   --root-dir=${K8S_DIR}/kubelet \\
19   --kubeconfig=/etc/kubernetes/kubelet.kubeconfig \\
20   --config=/etc/kubernetes/kubelet-config.yaml \\
21   --hostname-override=##NODE_NAME## \\
22   --pod-infra-container-image=registry.cn-beijing.aliyuncs.com/abcdo
   cker/pause-amd64:3.1 \\
23   --image-pull-progress-deadline=15m \\
24   --volume-plugin-dir=${K8S_DIR}/kubelet/kubelet-plugins/volume/exec
   / \\
25   --logtostderr=true \\
26   --v=2
27 Restart=always
28 RestartSec=5
29 StartLimitInterval=0
```

```
30 [Install]
31 WantedBy=multi-user.target
32 EOF
```

- 如果设置了 —hostname-override 选项，则 kube-proxy 也需要设置该选项，否则会出现找不到 Node 的情况；
- —bootstrap-kubeconfig：指向 bootstrap kubeconfig 文件，kubelet 使用该文件中的用户名和 token 向 kube-apiserver 发送 TLS Bootstrapping 请求；
- K8S approve kubelet 的 csr 请求后，在 —cert-dir 目录创建证书和私钥文件，然后写入 —kubeconfig 文件；
- —pod-infra-container-image 不使用 redhat 的 pod-infrastructure:latest 镜像，它不能回收容器的僵尸；

分发启动文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_name in ${NODE_NAMES[@]}
4   do
5     echo ">>> ${node_name}"
6     sed -e "s/##NODE_NAME##/${node_name}/" kubelet.service.template >
  kubelet-${node_name}.service
7     scp kubelet-${node_name}.service root@${node_name}:/etc/systemd/s
  ystem/kubelet.service
8   done
```

Bootstrap Token Auth 和授予权限 kubelet 启动时查找 —kubeletconfig 参数对应的文件是否存在，如果不存在则使用 —bootstrap-kubeconfig 指定的 kubeconfig 文件向 kube-apiserver 发送证书签名请求 (CSR)。 kube-apiserver 收到 CSR 请求后，对其中的 Token 进行认证，认证通过后将请求的 user 设置为 system:bootstrap:，group 设置为 system:bootstrappers，这一过程称为 Bootstrap Token Auth。
创建user和group的CSR权限，不创建kubelet会启动失败

```
1 $ kubectl create clusterrolebinding kubelet-bootstrap --clusterrole=s
  ystem:node-bootstrapper --group=system:bootstrappers
```

启动 kubelet 服务

```
1  source /opt/k8s/bin/environment.sh
2  for node_ip in ${NODE_IPS[@]}
3    do
4      echo ">>> ${node_ip}"
5      ssh root@${node_ip} "mkdir -p ${K8S_DIR}/kubelet/kubelet-plugins/
   volume/exec/"
6      ssh root@${node_ip} "/usr/sbin/swapoff -a"
7      ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
   kubelet && systemctl restart kubelet"
8    done
```

> 关闭 swap 分区，否则 kubelet 会启动失败；

kubelet 启动后使用 --bootstrap-kubeconfig 向 kube-apiserver 发送 CSR 请求，当这个
CSR 被 approve 后，kube-controller-manager 为 kubelet 创建 TLS 客户端证书、私钥和
--kubeletconfig 文件。 注意：kube-controller-manager 需要配置 --cluster-signing-cert-file 和
--cluster-signing-key-file 参数，才会为 TLS Bootstrap 创建证书和私钥。

```
1  $ kubectl get csr
2  NAME         AGE    REQUESTOR                 CONDITION
3  csr-22kt2    38s    system:bootstrap:pkkcl0   Pending
4  csr-f9trc    37s    system:bootstrap:tubfqq   Pending
5  csr-v7jt2    38s    system:bootstrap:ds9td8   Pending
6  csr-zrww2    37s    system:bootstrap:hy5ssz   Pending
```

这里4个节点均处于pending(等待)状态
**自动approve CSR请求**
创建三个ClusterRoleBinding，分别用于自动approve client、renew client、renew server证书

```
1  cd /opt/k8s/work
2  cat > csr-crb.yaml <<EOF
3   # Approve all CSRs for the group "system:bootstrappers"
4   kind: ClusterRoleBinding
5   apiVersion: rbac.authorization.k8s.io/v1
6   metadata:
7     name: auto-approve-csrs-for-group
8   subjects:
9   - kind: Group
```

```
10      name: system:bootstrappers
11      apiGroup: rbac.authorization.k8s.io
12    roleRef:
13      kind: ClusterRole
14      name: system:certificates.k8s.io:certificatesigningrequests:nodec
   lient
15      apiGroup: rbac.authorization.k8s.io
16  ---
17   # To let a node of the group "system:nodes" renew its own credentia
   ls
18   kind: ClusterRoleBinding
19   apiVersion: rbac.authorization.k8s.io/v1
20   metadata:
21      name: node-client-cert-renewal
22   subjects:
23   - kind: Group
24      name: system:nodes
25      apiGroup: rbac.authorization.k8s.io
26   roleRef:
27      kind: ClusterRole
28      name: system:certificates.k8s.io:certificatesigningrequests:selfn
   odeclient
29      apiGroup: rbac.authorization.k8s.io
30  ---
31  # A ClusterRole which instructs the CSR approver to approve a node r
   equesting a
32  # serving cert matching its client cert.
33  kind: ClusterRole
34  apiVersion: rbac.authorization.k8s.io/v1
35  metadata:
36    name: approve-node-server-renewal-csr
37  rules:
38  - apiGroups: ["certificates.k8s.io"]
39    resources: ["certificatesigningrequests/selfnodeserver"]
40    verbs: ["create"]
41  ---
42   # To let a node of the group "system:nodes" renew its own server cr
   edentials
43   kind: ClusterRoleBinding
44   apiVersion: rbac.authorization.k8s.io/v1
```

```
45  metadata:
46    name: node-server-cert-renewal
47  subjects:
48  - kind: Group
49    name: system:nodes
50    apiGroup: rbac.authorization.k8s.io
51  roleRef:
52    kind: ClusterRole
53    name: approve-node-server-renewal-csr
54    apiGroup: rbac.authorization.k8s.io
55 EOF
```

```
1 kubectl apply -f csr-crb.yaml
```

- auto-approve-csrs-for-group 自动approve node的第一次CSR，注意第一次CSR时，请求的Group为system:bootstrappers
- node-client-cert-renewal 自动approve node后续过期的client证书，自动生成的证书Group为system:nodes
- node-server-cert-renewal 自动approve node后续过期的server证书，自动生成的证书Group

**查看kubelet**

等待1-10分钟，3个节点的CSR都会自动approved

```
 1 $ kubectl get csr
 2 NAME        AGE      REQUESTOR                   CONDITION
 3 csr-22kt2   4m48s    system:bootstrap:pkkcl0     Approved,Issued
 4 csr-d8tvc   77s      system:node:k8s-01          Pending
 5 csr-f9trc   4m47s    system:bootstrap:tubfqq     Approved,Issued
 6 csr-kcdvx   76s      system:node:k8s-02          Pending
 7 csr-m8k8t   75s      system:node:k8s-04          Pending
 8 csr-v7jt2   4m48s    system:bootstrap:ds9td8     Approved,Issued
 9 csr-wwvwd   76s      system:node:k8s-03          Pending
10 csr-zrww2   4m47s    system:bootstrap:hy5ssz     Approved,Issued
```

Pending的CSR用于创建kubelet serve证书，需要手动approve（后面步骤）

目前所有节点均为ready状态

```
1 [root@k8s01 work]# kubectl get node
```

```
2 NAME      STATUS    ROLES    AGE       VERSION
3 k8s-01    Ready              2m29s     v1.14.2
4 k8s-02    Ready              2m28s     v1.14.2
5 k8s-03    Ready              2m28s     v1.14.2
6 k8s-04    Ready              2m27s     v1.14.2
```

kube-controller-manager为各node生成了kubeconfig文件和公钥

```
1 $ ls -l /etc/kubernetes/kubelet.kubeconfig
2 -rw------- 1 root root 2313 Aug 12 02:04 /etc/kubernetes/kubelet.kube
  config
3 $ ls -l /etc/kubernetes/cert/|grep kubelet
4 -rw------- 1 root root 1273 Aug 12 02:07 kubelet-client-2019-08-12-02
  -07-59.pem
5 lrwxrwxrwx 1 root root   59 Aug 12 02:07 kubelet-client-current.pem -
  > /etc/kubernetes/cert/kubelet-client-2019-08-12-02-07-59.pem
```

**手动approve server cert csr**
基于安全考虑，CSR approving controllers不会自动approve kubelet server证书签名请求，需要手动
approve

kubectl get csr | grep Pending | awk '{print $1}' | xargs kubectl certificate approve



kubelet API接口
kubelet启动后监听多个端口，用于接受kube-apiserver或其他客户端发送的请求

```
1  netstat -lntup|grep kubelet
2  tcp         0       0 192.168.0.50:10248       0.0.0.0:*                    L
   ISTEN       49491/kubelet
3  tcp         0       0 127.0.0.1:45737          0.0.0.0:*                    L
   ISTEN       49491/kubelet
4  tcp         0       0 192.168.0.50:10250       0.0.0.0:*                    L
   ISTEN       49491/kubelet
```

- 10248: healthz http 服务；
- 10250: https 服务，访问该端口时需要认证和授权（即使访问 /healthz 也需要）；
- 未开启只读端口 10255；
- 从 K8S v1.10 开始，去除了 —cadvisor-port 参数（默认 4194 端口），不支持访问 cAdvisor UI & API

**bear token认证和授权**

创建一个ServiceAccount，将它和ClusterRole system:kubelet-api-admin绑定，从而具有调用 kubelet API的权限

```
1  kubectl create sa kubelet-api-test
2  kubectl create clusterrolebinding kubelet-api-test --clusterrole=syst
   em:kubelet-api-admin --serviceaccount=default:kubelet-api-test
3  SECRET=$(kubectl get secrets | grep kubelet-api-test | awk '{print
    $1}')
4  TOKEN=$(kubectl describe secret ${SECRET} | grep -E '^token' | awk
   '{print $2}')
5  echo ${TOKEN}
```

**部署kube-proxy组件**

kube-proxy运行在所有worker节点上，它监听apiserver中service和endpoint的变化情况，创建路由规则提供服务IP和负载均衡功能。这里使用ipvs模式的kube-proxy进行部署

> 在各个节点需要安装ipvsadm和ipset命令，加载ip_vs内核模块

**创建kube-proxy证书签名请求**

```
1  cd /opt/k8s/work
2  cat > kube-proxy-csr.json <<EOF
3  {
4    "CN": "system:kube-proxy",
```

```
 5    "key": {
 6      "algo": "rsa",
 7      "size": 2048
 8    },
 9    "names": [
10      {
11        "C": "CN",
12        "ST": "BeiJing",
13        "L": "BeiJing",
14        "O": "k8s",
15        "OU": "xuyuntech"
16      }
17    ]
18 }
19 EOF
```

- CN: 指定该证书的 User 为 system:kube-proxy;
- 预定义的 RoleBinding system:node-proxier 将User system:kube-proxy 与 Role system:node-proxier 绑定, 该 Role 授予了调用 kube-apiserver Proxy 相关 API 的权限;
- 该证书只会被 kube-proxy 当做 client 证书使用, 所以 hosts 字段为空;

生成证书和私钥:

```
1 cd /opt/k8s/work
2 cfssl gencert -ca=/opt/k8s/work/ca.pem \
3   -ca-key=/opt/k8s/work/ca-key.pem \
4   -config=/opt/k8s/work/ca-config.json \
5   -profile=kubernetes  kube-proxy-csr.json | cfssljson -bare kube-pro
  xy
6 ls kube-proxy*
```

创建和分发 kubeconfig 文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 kubectl config set-cluster kubernetes \
4   --certificate-authority=/opt/k8s/work/ca.pem \
5   --embed-certs=true \
6   --server=${KUBE_APISERVER} \
```

```
 7    --kubeconfig=kube-proxy.kubeconfig
 8 kubectl config set-credentials kube-proxy \
 9    --client-certificate=kube-proxy.pem \
10    --client-key=kube-proxy-key.pem \
11    --embed-certs=true \
12    --kubeconfig=kube-proxy.kubeconfig
13 kubectl config set-context default \
14    --cluster=kubernetes \
15    --user=kube-proxy \
16    --kubeconfig=kube-proxy.kubeconfig
17 kubectl config use-context default --kubeconfig=kube-proxy.kubeconfi
   g
```

- 一embed-certs=true：将 ca.pem 和 admin.pem 证书内容嵌入到生成的kubectl-proxy.kubeconfig文件中(不加时，写入的是证书文件路径);
- 分发 kubeconfig 文件:

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for node_name in ${NODE_NAMES[@]}
4   do
5     echo ">>> ${node_name}"
6     scp kube-proxy.kubeconfig root@${node_name}:/etc/kubernetes/
7   done
```

**创建kube-proxy配置文件**
从v1.10开始，kube-proxy部分参数可以配置在文件中，可以使用一write-config-to选项生成该配置文件

```
 1 cd /opt/k8s/work
 2 cat > kube-proxy-config.yaml.template <<EOF
 3 kind: KubeProxyConfiguration
 4 apiVersion: kubeproxy.config.k8s.io/v1alpha1
 5 clientConnection:
 6   burst: 200
 7   kubeconfig: "/etc/kubernetes/kube-proxy.kubeconfig"
 8   qps: 100
 9 bindAddress: ##NODE_IP##
10 healthzBindAddress: ##NODE_IP##:10256
```

```
11 metricsBindAddress: ##NODE_IP##:10249
12 enableProfiling: true
13 clusterCIDR: ${CLUSTER_CIDR}
14 hostnameOverride: ##NODE_NAME##
15 mode: "ipvs"
16 portRange: ""
17 kubeProxyIPTablesConfiguration:
18   masqueradeAll: false
19 kubeProxyIPVSConfiguration:
20   scheduler: rr
21   excludeCIDRs: []
22 EOF
```

- bindAddress: 监听地址；
- clientConnection.kubeconfig: 连接 apiserver 的 kubeconfig 文件；
- –clusterCIDR: kube–proxy 根据 —cluster–cidr判断集群内部和外部流量，指定 —cluster–cidr 或 —masquerade–all 选项后 kube–proxy 才会对访问 Service IP 的请求做 SNAT；
- hostnameOverride: 参数值必须与 kubelet 的值一致，否则 kube–proxy 启动后会找不到该 Node，从而不会创建任何 ipvs 规则；
- mode: 使用 ipvs 模式；

分发和创建kube–proxy配置文件

```
1 cd /opt/k8s/work
2 source /opt/k8s/bin/environment.sh
3 for (( i=0; i < 4; i++ ))
4   do
5     echo ">>> ${NODE_NAMES[i]}"
6     sed -e "s/##NODE_NAME##/${NODE_NAMES[i]}/" -e "s/##NODE_IP##/${NODE_IPS[i]}/" kube-proxy-config.yaml.template > kube-proxy-config-${NODE_NAMES[i]}.yaml.template
7     scp kube-proxy-config-${NODE_NAMES[i]}.yaml.template root@${NODE_NAMES[i]}:/etc/kubernetes/kube-proxy-config.yaml
8   done
```

#我这里一共有4个节点要运行，所以这里写4，这是整个集群的node节点的数量！ 这里一定要注意修改！！
创建和分发 kube–proxy systemd unit 文件

```
 1  cd /opt/k8s/work
 2  source /opt/k8s/bin/environment.sh
 3  cat > kube-proxy.service <<EOF
 4  [Unit]
 5  Description=Kubernetes Kube-Proxy Server
 6  Documentation=https://github.com/GoogleCloudPlatform/kubernetes
 7  After=network.target
 8  [Service]
 9  WorkingDirectory=${K8S_DIR}/kube-proxy
10  ExecStart=/opt/k8s/bin/kube-proxy \\
11    --config=/etc/kubernetes/kube-proxy-config.yaml \\
12    --logtostderr=true \\
13    --v=2
14  Restart=on-failure
15  RestartSec=5
16  LimitNOFILE=65536
17  [Install]
18  WantedBy=multi-user.target
19  EOF
```

分发 kube-proxy systemd unit 文件：

```
 1  cd /opt/k8s/work
 2  source /opt/k8s/bin/environment.sh
 3  for node_name in ${NODE_NAMES[@]}
 4    do
 5      echo ">>> ${node_name}"
 6      scp kube-proxy.service root@${node_name}:/etc/systemd/system/
 7    done
```

启动 kube-proxy 服务

```
 1  cd /opt/k8s/work
 2  source /opt/k8s/bin/environment.sh
 3  for node_ip in ${NODE_IPS[@]}
 4    do
 5      echo ">>> ${node_ip}"
 6      ssh root@${node_ip} "mkdir -p ${K8S_DIR}/kube-proxy"
```

```
7    ssh root@${node_ip} "modprobe ip_vs_rr"
8    ssh root@${node_ip} "systemctl daemon-reload && systemctl enable
   kube-proxy && systemctl restart kube-proxy"
9  done
```

检查启动结果

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "systemctl status kube-proxy|grep Active"
6   done
```

检查监听端口

```
1 [root@k8s01 work]# netstat -lnpt|grep kube-prox
2 tcp        0      0 192.168.0.50:10249      0.0.0.0:*               L
  ISTEN      55015/kube-proxy
3 tcp        0      0 192.168.0.50:10256      0.0.0.0:*               L
  ISTEN      55015/kube-proxy
```

查看ipvs路由规则

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh root@${node_ip} "/usr/sbin/ipvsadm -ln"
6   done
```

正常输出如下

```
[root@abcdocker-k8s01 work]# source /opt/k8s/bin/environment.sh
[root@abcdocker-k8s01 work]# for node_ip in ${NODE_IPS[@]}
>   do
>     echo ">>> ${node_ip}"
>     ssh root@${node_ip} "/usr/sbin/ipvsadm -ln"
>   done
>>> 192.168.0.50
IP Virtual Server version 1.2.1 (size=4096)
Prot LocalAddress:Port Scheduler Flags
  -> RemoteAddress:Port           Forward Weight ActiveConn InActConn
TCP  10.254.0.1:443 rr
  -> 192.168.0.50:6443            Masq    1      0          0
  -> 192.168.0.51:6443            Masq    1      0          0
  -> 192.168.0.52:6443            Masq    1      0          0
>>> 192.168.0.51
IP Virtual Server version 1.2.1 (size=4096)
Prot LocalAddress:Port Scheduler Flags
  -> RemoteAddress:Port           Forward Weight ActiveConn InActConn
TCP  10.254.0.1:443 rr
  -> 192.168.0.50:6443            Masq    1      0          0
  -> 192.168.0.51:6443            Masq    1      0          0
  -> 192.168.0.52:6443            Masq    1      0          0
>>> 192.168.0.52
IP Virtual Server version 1.2.1 (size=4096)
Prot LocalAddress:Port Scheduler Flags
  -> RemoteAddress:Port           Forward Weight ActiveConn InActConn
TCP  10.254.0.1:443 rr
  -> 192.168.0.50:6443            Masq    1      0          0
  -> 192.168.0.51:6443            Masq    1      0          0
  -> 192.168.0.52:6443            Masq    1      0          0
>>> 192.168.0.53                                        • i4t.com
IP Virtual Server version 1.2.1 (size=4096)
Prot LocalAddress:Port Scheduler Flags
  -> RemoteAddress:Port           Forward Weight ActiveConn InActConn
[root@abcdocker-k8s01 work]# █
```

可见所有通过 https 访问 K8S SVC kubernetes 的请求都转发到 kube-apiserver 节点的 6443 端口；

**验证集群功能**

现在使用daemonset验证master和worker节点是否正常

```
1 [root@k8s01 work]# kubectl get node
2 NAME       STATUS    ROLES     AGE     VERSION
3 k8s-01     Ready     20m       v1.14.2
4 k8s-02     Ready     20m       v1.14.2
5 k8s-03     Ready     20m       v1.14.2
6 k8s-04     Ready     20m       v1.14.2
```

创建测试yaml文件

执行测试

```
1 kubectl create -f nginx-ds.yml
```

这里pod已经启动成功

```
1 [root@k8s01 work]# kubectl get pod  -o wide
2 NAME              READY   STATUS     RESTARTS    AGE    IP            NO
  DE      NOMINATED NODE   READINESS GATES
3 nginx-ds-29n8p    1/1     Running    0           116s   172.17.0.2    k8
  s-02
4 nginx-ds-7zhbb    1/1     Running    0           116s   172.30.96.2   k8
  s-01
5 nginx-ds-kvr7q    1/1     Running    0           116s   172.17.0.2    k8
  s-04
6 nginx-ds-lk9dv    1/1     Running    0           116s   172.17.0.2    k8
  s-03
```

检查各节点的Pod IP 连通性
这里看到pod的IP，我们将ip复制一下

```
1 source /opt/k8s/bin/environment.sh
2 for node_ip in ${NODE_IPS[@]}
3   do
4     echo ">>> ${node_ip}"
5     ssh ${node_ip} "ping -c 1 172.17.0.2"
6     ssh ${node_ip} "ping -c 1 172.30.96.2"
7     ssh ${node_ip} "ping -c 1 172.17.0.2"
8   done
```

检查服务IP和端口可达性

```
1 [root@k8s01 work]# kubectl get svc |grep nginx-ds
2 nginx-ds    NodePort    10.254.248.73          80:15402/TCP    4m11s
```
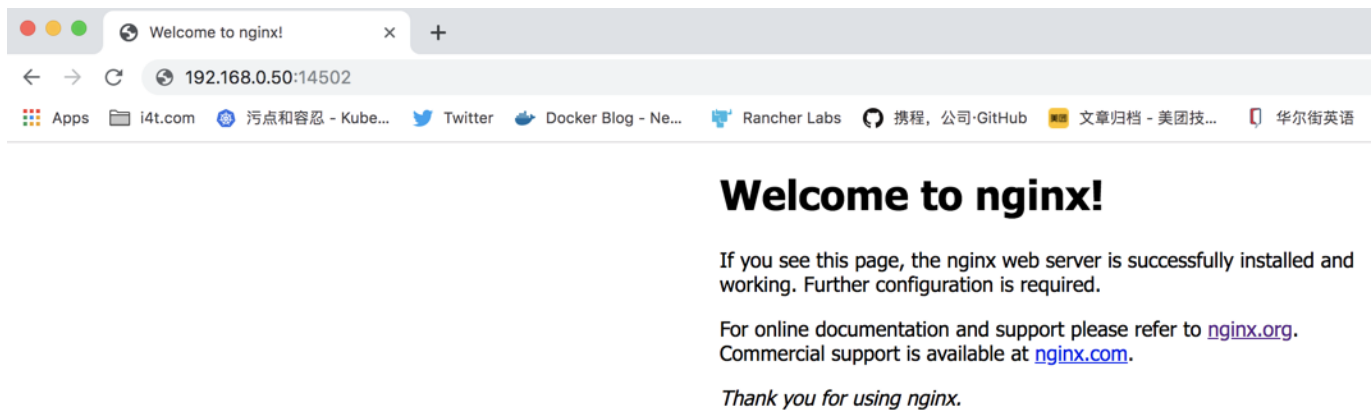
我们在任意节点访问server IP

```
1  source /opt/k8s/bin/environment.sh
2  for node_ip in ${NODE_IPS[@]}
3    do
4      echo ">>> ${node_ip}"
5      ssh ${node_ip} "curl -s 10.254.248.73"
6    done
```

#这里请根据上面查看的svc IP进行修改

此时我们已经可以使用任意节点IP+15402端口访问nginx（这个端口是通过kubectl get svc获取到的，每个人的端口可能不一样。请按照实际情况进行修改！）



## CoreDNS安装

上面我们验证的集群内部网络，已经没有问题。接下来进行安装DNS

这里的所有操作在k8s01上执行即可

```
1  source /opt/k8s/bin/environment.sh
2  for node_ip in ${NODE_IPS[@]}
3    do
4      echo "$node_ip"
5      ssh $node_ip "wget -P /opt/ http://down.i4t.com/coredns_v1.4.tar"
6      ssh $node_ip "docker load -i /opt/coredns_v1.4.tar"
7    done
```

#下载镜像并分发镜像

下载coredns yaml文件

```
1 wget -P /opt/ http://down.i4t.com/k8s1.14/coredns.yaml
```

创建coredns

```
1 kubectl create -f /opt/coredns.yaml
```

#这里已经镜像让你们手动下载了，没有下载请看docker步骤,最后一步
执行完毕后，pod启动成功（Running状态为正常）

```
1 kubectl get pod -n kube-system  -l k8s-app=kube-dns
2 NAME                         READY    STATUS    RESTARTS    AGE
3 coredns-d7964c8db-vgl5l      1/1      Running   0           21s
4 coredns-d7964c8db-wvz5k      1/1      Running   0           21s
```

coredns启动之后，我们需要测一下dns功能是否正常

> 温馨提示：busybox高版本有nslookup Bug，不建议使用高版本，请按照我的版本进行操作即可！

创建一个yaml文件测试是否正常

```
1 cat<<EOF | kubectl apply -f -
2 apiVersion: v1
3 kind: Pod
4 metadata:
5   name: busybox
6   namespace: default
7 spec:
8   containers:
9   - name: busybox
10    image: busybox:1.28.3
11    command:
12      - sleep
13      - "3600"
14    imagePullPolicy: IfNotPresent
15  restartPolicy: Always
```

```
16 EOF
```

创建后Pod我们进行检查

```
1 kubectl get pod
2 NAME       READY    STATUS     RESTARTS    AGE
3 busybox    1/1      Running    0           4s
```

使用nslookup查看是否能返回地址

```
1 kubectl exec -ti busybox -- nslookup kubernetes
2 Server:     10.254.0.2
3 Address 1: 10.254.0.2 kube-dns.kube-system.svc.cluster.local
4 Name:       kubernetes
5 Address 1: 10.254.0.1 kubernetes.default.svc.cluster.local
```

默认kubectl没有table补全命令，如果需要补全请参考下面文章

# 三、k8s集群验证（部署KubeStar）

如果遇到kubectl无法进入容器报错：

```
error: unable to upgrade connection: Forbidden (user=kubernetes,
verb=create, resource=nodes, subresource=proxy)
```

可以使用下面的命令解决：

```
1 kubectl exec -ti kubestar-75747f5bc9-77ptw -n kubestar-deploy /bin/sh
```

Service短域名调用超时：
在集群里kubestar的node-1上的Pod里访问kubestar自身的kubestar服务地址（轮询到本机），超时；
在node-1上访问localhost nodeport超时，访问node-2，node-3OK。