

Tony Bai

一个程序员的心路历程

- 关于我
- 文章目录

理解Kubernetes网络之Flannel网络

- 一月 17, 2017
- 2条评论

第一次采用kubernetes脚本方式安装Kubernetes cluster目前运行良好，master node上的组件状态也始终是“没毛病”：

```
# kubectl get cs
NAME                STATUS    MESSAGE             ERROR
controller-manager  Healthy   ok
scheduler            Healthy   ok
etcd-0               Healthy   {"health": "true"}
```

不过在第二次尝试用kubernetes安装和初始化Kubernetes cluster时遇到的各种网络问题还是让我“心有余悸”。于是趁上个周末，对Kubernetes的网络原理进行了一些针对性的学习，这里把对Kubernetes网络的理解记录一下和大家一起分享。

Kubernetes支持Flannel、Calico、Weave network等多种网络Drivers，但由于学习过程使用的是第一个cluster的Flannel网络，这里的网络原理只针对k8s+Flannel网络。

一、环境+提示

凡涉及到Docker、Kubernetes这类正在Active dev的开源项目的文章，我都不得不提一嘴，那就是随着K8s以及flannel的演化，本文中的一些说法可能不再正确，提醒大家：阅读此类技术文章务必结合“环境”。

这里我们使用的环境就是我第一次建立k8s cluster的环境：

```
# kube-apiserver --version
Kubernetes v1.3.7

# /opt/bin/flanneld -version
0.5.5

# /opt/bin/etcd -version
etcd Version: 2.0.12
Git SHA: 2af268f
Go Version: go1.6.3
Go OS/Arch: linux/amd64
```

另外整个集群搭建在[百度云](#)上，每个ECS上的OS及kernel版本：Ubuntu 14.04.4 LTS, 3.19.0-70-generic。

在我的测试环境，有两个node：master node和一个minion node，master node参与workload的调度，所以你可以认为有两个minion node即可。

二、Kubernetes Cluster中的几个“网络”

之前的k8s cluster采用的是默认安装，即直接使用了配置脚本中(kubernetes.cluster/ubuntu/config.default.sh)自带的一些参数，比如：

```
//摘自kubernetes/cluster/ubuntu/config-default.sh
export nodes="{nodes:'root@master_node_ip root@minion_node_ip'}"
export SERVICE_CIDR_IP_RANGE={SERVICE_CIDR_IP_RANGE:-192.168.3.0/24}
export FLANNEL_NET={FLANNEL_NET:-172.16.0.0/16}
```

从这里我们能够识别出三个“网络”：

- node network：承载Kubernetes集群中各个“物理”Node(master和minion)通信的网络；
- service network：由Kubernetes集群的Service所组成的“网络”；
- flannel network：即Pod网络，集群中承载各个Pod相互通信的网络。

node network自不必多说，node间通过你的本地局域网（无论是物理的还是虚拟的）通信。

service network比较特殊，每个新创建的service会被分配一个service IP，在当前集群中，这个IP的分配范围是192.168.3.0/24，不过这个IP并不“真实”，更像一个“占位符”并且只有入口流量，所谓的“network”也是“名不副实”的，后续我们会详尽说明。

flannel network是我们理解的重点，cluster中各个Pod要实现相互通信，必须走这个网络，无论是在同一node上的Pod还是跨node的Pod，我们的cluster中，flannel net的分配范围是：172.16.0.0/16。

在进一步挖掘“原理”之前，我们先来直观认识一下service network和flannel network：

```
Service network(看cluster ip一列)：
# kubectl get services
NAME          CLUSTER-IP      EXTERNAL-IP      PORT(S)          AGE
kube-api      192.168.3.148    <none>            3000/TCP         1d
kubernetes    192.168.3.148    <none>            443/TCP          1d
my-nginx      192.168.3.179    <none>            80/TCP           90s
nginx-tilt    192.168.3.196    <none>            80/TCP           12s
rbd-rwtest-api 192.168.3.22     <none>            8080/TCP         60s

Flannel network (看IP那一列)：
# kubectl get pod -o wide
NAME          READY    STATUS    RESTARTS   AGE      IP              NODE
my-nginx-239515548-gpijv 1/1      Running   0          91d      172.16.99.3     (master node ip)
nginx-tilt-3872865736-cd8hr 2/2      Running   0          12d      172.16.57.7     (minion node ip)
...
```

三、平坦的Flannel网络

1、Kubernetes安装后的网络状态

首先让我们来看看：[kube-up.sh在安装k8s集群](#)时对各个K8s Node都动了什么手脚！

a) 修改etcd的默认配置

在ubuntu 14.04下，[docker的配置](#)都在/etc/defaults/docker文件中，如果你曾经修改过该文件，那么k8s脚本方式安装完Kubernetes后，你会发现/etc/defaults/docker已经变样了，只剩下一行：

```
master node:
DOCKER_OPTS="-m tcp://172.0.0.1:4243 -m unix:///var/run/docker.sock --hipp=172.16.99.1/24 --mtu=1450"

minion node:
DOCKER_OPTS="-m tcp://172.0.0.1:4243 -m unix:///var/run/docker.sock --hipp=172.16.57.1/24 --mtu=1450"
```

可以看出k8s脚本修改了Docker daemon的--hip选项，使得该node上的docker daemon在该node的flannel subnet范围以内为自动的Docker container分配IP地址。

b) 在etcd中初始化flannel网络数据

多个node上的Flannel依赖一个[etcd cluster](#)来集中配置服务，etcd保证了所有node上flannel所看到的配置是一致的，同时每个node上的flannel监听etcd上的数据变化，实时感知集群中node的变化。

我们可以通过etcdctl查询到这些配置数据：

```
master node:
//flannel network配置
# etcdctl --endpoints http://172.0.0.1:(etcd listen port) get /coreos.com/network/config
{"Network": "172.16.0.0/16", "Backend": {"Type": "vxlan"}, "Value": ""}
# etcdctl --endpoints http://172.0.0.1:(etcd listen port) ls /coreos.com/network/subnets
/coreos.com/network/subnets/172.16.99.0-24
/coreos.com/network/subnets/172.16.57.0-24
//某node上的flanne subnet的etcd配置
# etcdctl --endpoints http://172.0.0.1:(etcd listen port) get /coreos.com/network/subnets/172.16.99.0-24
{"PublicIP": "(master node ip)", "BackendType": "vxlan", "BackendData": {"VxlanPort": "b61d1ec1c1cf13b7"}, "Value": ""}
minion node:
# etcdctl --endpoints http://172.0.0.1:(etcd listen port) get /coreos.com/network/subnets/172.16.57.0-24
{"PublicIP": "(minion node ip)", "BackendType": "vxlan", "BackendData": {"VxlanPort": "d61512a803c1c1d9"}, "Value": ""}
```

或用etcd提供的rest api：

```
# curl -L http://172.0.0.1:(etcd listen port)/v2/keysp/coreos.com/network/config
{"action": "get", "node": {"key": "/coreos.com/network/config", "value": {"\"Network\": \"172.16.0.0/16\", \"Backend\": {\"Type\": \"vxlan\"}}\", \"modifiedIndex\": 5, \"createdIndex\": 5}}
```

c) 启动flanneld

kube-up.sh在每个Kubernetes node上启动了一个flanneld的程序：

```
# ps -ef|grep flanneld
master node:
root      1151      1  0 2016 ?        00:02:34 /opt/bin/flanneld --etcd-endpoints=http://172.0.0.1:(etcd listen port) --ip-masq --iface=(master node ip)
minion node:
root      11940     1  0 2016 ?        00:07:05 /opt/bin/flanneld --etcd-endpoints=http://(master node ip):(etcd listen port) --ip-masq --iface=(minion node ip)
```

一旦flanneld启动，它得从etcd中读取配置，并请求获取一个subnet lease(租约)，有效期目前是24hrs，并且监视etcd的数据更新，flanneld一旦获取subnet租约，配置完backend，它会将一些信息写入/run/flannel/subnet.env文件。

```
master node:
# cat /run/flannel/subnet.env
FLANNEL_NETWORK=172.16.0.0/16
FLANNEL_SUBNET=172.16.99.1/24
FLANNEL_MTU=1450
FLANNEL_IPMASQ=true

minion node:
# cat /run/flannel/subnet.env
FLANNEL_NETWORK=172.16.0.0/16
FLANNEL_SUBNET=172.16.57.0/16
FLANNEL_MTU=1450
FLANNEL_IPMASQ=true
```

当然flanneld的最大意义在于根据etcd中存储的全cluster的subnet信息，跨node传输flannel network中的数据包，这个后面会详细说明。

d) 创建flannel1网络设备，更新路由信息

各个node上的网络设备列表新增一个名为flannel.1的类型为vxlan的网络设备：

```
master node:
# ip -d link show
4: flannel.1: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1450 qdisc noqueue state UNKNOWN mode DEFAULT group default
    link/ether d6:b6:4c:81:c1cf:3b brd ff:ff:ff:ff:ff:ff promiscuity 0
    vxlan id 1 local (master node local ip) dev eth0 port 0 0 nolearning ageing 300

minion node:
349: flannel.1: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1450 qdisc noqueue state UNKNOWN mode DEFAULT group default
    link/ether d6:51:2a:80:3c1d:49 brd ff:ff:ff:ff:ff:ff promiscuity 0
    vxlan id 1 local (minion node local ip) dev eth0 port 0 0 nolearning ageing 300
```

从flannel.1的设备信息来看，它似乎与eth0存在着某种bind关系，这是在其他bridge、veth设备描述信息中所没有的。

flanneld设备的top：

```
master node:
flannel.1: link encap:ethernet S-Maddr b6:b6:4c:81:c1cf:3b
    link addr:172.16.99.0 S-Maddr:0.0.0.0 S-Mask:255.255.0.0
    UP BROADCAST RUNNING MULTICAST MTU:1450 Metric:1
    RX packets:193724 received:0 dropped:0 overruns:0 frame:0
    TX packets:15829444 sent:0 dropped:292 overruns:0 carrier:0
    collisions:0 txqueuelen:0
    RX bytes:168989045 (1.6 GB) TX bytes:1144725704 (1.1 GB)

minion node:
flannel.1: link encap:ethernet S-Maddr d6:51:2a:80:3c1d:49
    link addr:172.16.57.0 S-Maddr:0.0.0.0 S-Mask:255.255.0.0
    UP BROADCAST RUNNING MULTICAST MTU:1450 Metric:1
    RX packets:129440 received:0 dropped:0 overruns:0 frame:0
    TX packets:1755599 sent:0 dropped:25 overruns:0 carrier:0
    collisions:0 txqueuelen:0
    RX bytes:198936257 (198.3 MB) TX bytes:1861492847 (1.8 GB)
```

可以看到两个node上的flannel.1的ip与k8s cluster为两个node上分配subnet的ip范围是对应的。

下面是两个node上的当前路由表：

```
master node:
```

```
# ip route
172.16.0.0/16 dev flannel.1 proto kernel scope link src 172.16.99.0
172.16.99.0/24 dev dockero proto kernel scope link src 172.16.99.1
...

minion node:
# ip route
...
172.16.0.0/16 dev flannel.1
172.16.57.0/24 dev dockero proto kernel scope link src 172.16.57.1
...
```

以上信息得为后续数据包传输分析打下基础。

c) 平铺的flannel network

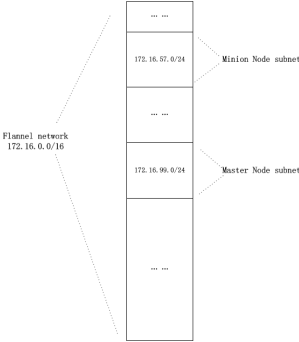
从以上kubernetes和flannel network安装之后获得的网络信息。我们能看到flannel network是一个flat network。在flannel: 172.16.0.0/16这个大网下，每个kubernetes node从中分配一个子网片段(24)：

```
master node:
--bip=172.16.99.1/24

minion node:
--bip=172.16.57.1/24

root@node1:~# etcdctl --endpoints http://127.0.0.1:2379 listen port} ls /coreos.com/network/subnets
/coreos.com/network/subnets/172.16.99.0-24
/coreos.com/network/subnets/172.16.57.0-24
```

用一张图来这样可能更为直观：

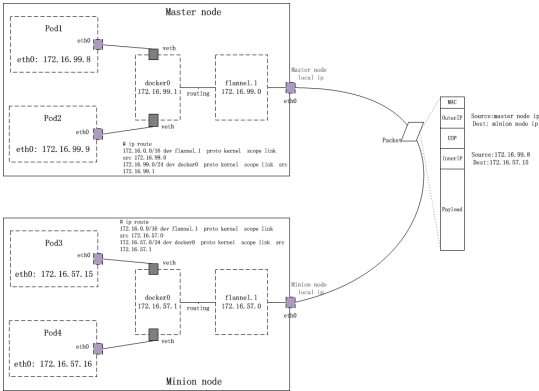


这个是不是有些像x86-64的虚拟内存寻址空间啊（同样是平铺内存地址访问模型）！

在平铺的flannel network中，每个pod都会被分配唯一的ip地址，且每个k8s node的子net各不重叠，没有交集。不过这样的子net分配模型也有一定弊端，那就是可能存在ip浪费：一个node上有200多个flannel ip地址(172.16.99.0-24)，如果仅仅启动了几个Pod，那么其余ip就处于闲置状态。

2、Flannel网络通信原理

这里我们模仿flannel官方的部署原理图，画了一幅与我们的实验环境匹配的图，作为后续讨论flannel网络通信流程的基础：



如上图所示，我们来看看从pod1: 172.16.99.8发出的数据包是如何到达pod3: 172.16.57.15的（比如：在pod1的某个container中ping -c 3 172.16.57.15）。

a) 从Pod出发

由于k8s更改了docker的DOCKER_OPTS，显式指定了-bip，这个值与分配给该node上的子net的范围是一致的。这样一来，docker引擎每次创建一个Docker container，该container被分配到的ip都在flannel subnet范围内。

当我们在Pod1下的某个容器内执行ping -c 3 172.16.57.15，数据包便开始了它在flannel network中的旅程。

Pod是Kubernetes调度的基本unit，Pod内的多个container共享一个network namespace。Kubernetes在创建Pod时，首先先创建pause容器，然后再以pause的network namespace为基础，创建pod内的其他容器（-net=container:xxx）。这样Pod内的所有容器共享一个network namespace，这些容器间的访问直接通过localhost即可。比如Pod下A容器启动了一个服务，监听8080端口，那么同一个Pod下面的另外一个B容器通过访问localhost:8080即可访问到A容器下面的那个服务。

在之前的《[详解Docker自画网络 2.1 Linux Network Namespace](#)》一文中，我相信我已经讲清楚了单机下Docker容器数据传的路径。在这个环节中，数据包的传输路径也并无不同。

我们看一下Pod1中某Container内的路由信息：

```
# docker exec ba7f81455c7 ip route
default via 172.16.99.1 dev eth0
172.16.99.0/24 dev eth0 proto kernel scope link src 172.16.99.8
```

目的地址172.16.57.15并不在直连网络中，因此数据包通过default路由出去，default路由的由路由地址是172.16.99.1，也就是上面的dockero bridge的IP地址。相当于dockero bridge以“二层的工件模式”直接接收来自容器的数据包而非从bridge的二层端口接收。

b) dockero与flannel.1之间的转发

数据包到达dockero后，dockero的内核线程处理程序发现这个数据包的目的地址是172.16.57.15，并不是真的要送给自己，于是开始为该数据包找下一hop，根据master node上的路由表：

```
master node:
# ip route
172.16.0.0/16 dev flannel.1 proto kernel scope link src 172.16.99.0
172.16.99.0/24 dev dockero proto kernel scope link src 172.16.99.1
...
```

我们匹配到“172.16.0/16”这条路由！这是一条直连路由，数据包被直接送到flannel.1设备上。

c) flannel.1设备以及flanneld的功用

flannel.1是否会重复dockero的套路呢？包不是发给自己，转发数据包？会，也不会。

“会”是指flannel.1肯定要转发数据包出去，因为毕竟包不是给自己的（包目的ip是172.16.57.15，vlan设备ip是172.16.99.0）。

“不会”是指flannel.1不会走寻常套路去转发包，因为它是一个vxlun类型的设备，也称为vtep，virtual tunnel end point。

那么它到底是怎么处理数据包的呢？这里涉及一些Linux内核对vxlan处理的内容，详细内容可参见本文末尾的参考资料。

flannel.1收到数据包后，由于自己不是目的地，也要尝试得数据包重新发送出去，数据包沿着网络协议栈向下流动，在二层时需要封二层以太包，填写目的mac地址，这时一般应该发出arp：“who is 172.16.57.15”，但vxlan设备的特殊性就在于它并没有真正在二层发出这个arp包，因为下面的这个内核参数设置：

```
# cat /proc/sys/net/ipv4/neigh/flannel.1/app_solicit
3
```

而是由linux kernel引发一个“L3 MISS”事件并得arp请求发到用户空间的flanneld程序。

flanneld程序收到“L3 MISS”内核事件以及arp请求(who is 172.16.57.15)后，并不会向外网发送arp request，而是尝试从etc/查找该地址匹配的子网的vtep信息。在前面章节我们曾经展示过etc/中flannel network的配置信息：

```
master node:
# etcdctl --endpoints http://127.0.0.1:2379 listen port} ls /coreos.com/network/subnets
/coreos.com/network/subnets/172.16.99.0-24
/coreos.com/network/subnets/172.16.57.0-24
# curl -s http://127.0.0.1:2379 listen port}/v4/range/coreos.com/network/subnets/172.16.57.0-24
{"action": "get", "node": {"key": "/coreos.com/network/subnets/172.16.57.0-24", "value": {"\"PublicIP\": \"(minion node local ip)\", \"BackendType\": \"(vxlan)\", \"BackendData\": {\"(\"vtepMac\": \"(d6:51:2e:80:5c:69)\", \"expiration\": \"2017-01-17T09:46:20.607339725s\", \"ttl\": 21466, \"modifiedIndex\": 2275460, \"createdIndex\": 2275460)}}}}
```

flanneld从etc/中找到了答案：

```
subnet: 172.16.57.0/24
public ip: (minion node local ip)
vtepMac: d6:51:2e:80:5c:69
```

我们查看minion node上的信息，发现minion node上的flannel.1设备mac就是d6:51:2e:80:5c:69：

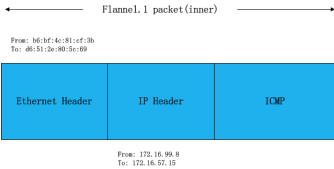
```
minion node:
# ip -d link show
34v flannel.1: 8000CAPT,MASTER,UP,LOWER_UP mru 1450 qlen response state UNKNOWN group default
link/ether d6:51:2e:80:5c:69 bond ffff:ffff:ffff:promiscuity 0
vxlan id 1 local 10.46.181.146 dev eth0 port 0 nolearning ageing 300
```

接下来，flanneld将查询到的信息放入master node head/arp cache表中：

```
master node:
# ip n [group 172.16.57.15
172.16.57.15 dev flannel.1 lladdr d6:51:2e:80:5c:69 REACHABLE
```

flanneld完成这项工作后，linux kernel就可以在arp table中找到172.16.57.15对应的mac地址并封装二层以太包了。

到目前为止，已经呈现在大家眼前的封装如下图：



不过这个封装也不能在物理网络上传输，因为它实际上只是vlan tunnel上的packet。

4) kernel的vlan封装

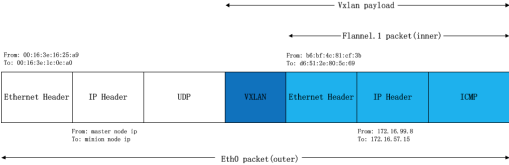
我们需要将上述的packet从master node传输到minion node，需要将上述packet再次封装。这个任务在backend为vlan的flannel network中由linux kernel来完成。

flannel.1为vlan设备，linux kernel可以自动识别，并将上面的packet进行vlan封装处理。在这个封装过程中，kernel需要知道该数据包究竟发到哪个node上去，kernel需要查看node上的fdb(forwarding database)以获得上面对等vtep设备（已经从arp table中查到其mac地址：dc:51:2c:80:5c:60）所在的node地址。如果db中没有这个信息，那么kernel会向用户空间flanneld程序发发L2 MISS事件，flanneld收到该事件后，会查询etcd，获取该vtep设备对应的node的Public IP，并将信息注册到db中。

这样Kernel就可以顺利查询到该信息并封装了：

```
master node:
# bridge fdb show dev flannel.1 | grep dc:51:2c:80:5c:60
dc:51:2c:80:5c:60 dev vtep (relation mode local ip) state permanent
```

由于目标ip是minion node，查找路由表，包应该从master node的eth0发出，这样src ip和src mac地址也就确定了。封装的包示意图如下：



5) kernel的vlan拆包

minion node上的eth0接收到上述vlan包，kernel将识别出这是一个vlan包，于是拆包后将flannel.1 packet转给minion node上的vtep (flannel.1)，minion node上的flannel.1再把这个数据包转到minion node上的docker0，随后由docker0传输到Pod3的某个容器里。

3、Pod内到外部网络

我们在Pod中除了可以与pod network中的其他pod通信外，还可以访问外部网络，比如：

```
master node:
# docker exec ba75f81455e7 ping -c 3 baifu.com
PING baifu.com (180.149.132.47): 16 data bytes
64 bytes from 180.149.132.47: icmp_seq=0 ttl=54 time=3.586 ms
64 bytes from 180.149.132.47: icmp_seq=1 ttl=54 time=3.752 ms
64 bytes from 180.149.132.47: icmp_seq=2 ttl=54 time=3.722 ms
--- baifu.com ping statistics ---
3 packets transmitted, 3 packets received, 0% packet loss
round-trip min/avg/med/stddev = 3.586/3.687/3.752/0.072 ms
```

这个通信与vlan就没什么关系了，主要是通过docker引擎在iptables的POSTROUTING chain中设置的MASQUERADE规则：

```
master node:
# iptables -t nat -nL
Chain POSTROUTING (policy ACCEPT)
target prot opt source destination
MASQUERADE all -- 172.16.99.0/24 0.0.0.0/0
```

docker将容器的pod network地址伪装为node ip出去，包回来时再nat回容器的pod network地址，这样网络就通了。

四、“不真实”的Service网络

每当我们在k8s cluster中创建一个service，k8s cluster就会在--service.cluster-ip-range的范围为service分配一个cluster-ip，比如本文开始时提到的：

```
# kubect get services
NAME CLUSTER-IP EXTERNAL-IP PORT(S) AGE
index-api 192.168.3.168 <none> 30080/TCP 18d
kubernetes 192.168.3.1 <none> 443/TCP 94d
my-nginx 192.168.3.179 <none> 80/TCP 90d
nginx-tls 192.168.3.196 <none> 80/TCP 12d
rbd-rest-api 192.168.3.22 <none> 8080/TCP 60d
```

这个cluster-ip只是一个虚拟的ip，并不真实绑定某个物理网络设备或虚拟网络设备，仅仅存在于iptables的规则中：

```
Chain PREROUTING (policy ACCEPT)
target prot opt source destination
KUBE-SERVICES all -- 0.0.0.0/0 0.0.0.0/0 /* kubernetes service portals */

# iptables -t nat -nL | grep 192.168.3
Chain KUBE-SERVICES (1 references)
target prot opt source destination
KUBE-SVC-MZGZMVGQ9VZ2RS tcp -- 0.0.0.0/0 192.168.3.182 /* kube-system/kubernetes-dashboard: cluster IP */ tcp dpt:80
KUBE-SVC-SE44K4K9P7Y7XKX tcp -- 0.0.0.0/0 192.168.3.1 /* default/kubernetes:https: cluster IP */ tcp dpt:443
KUBE-SVC-AU252PRZZQGOERSG tcp -- 0.0.0.0/0 192.168.3.22 /* default/rbd-rest-api: cluster IP */ tcp dpt:8080
KUBE-SVC-TCW7VQZGQERSGM tcp -- 0.0.0.0/0 192.168.3.10 /* kube-system/kube-dashboard: cluster IP */ udp dpt:53
KUBE-SVC-MEP9X80H9PFC1C tcp -- 0.0.0.0/0 192.168.3.179 /* default/my-nginx: cluster IP */ tcp dpt:80
KUBE-SVC-UG6H1872J2R278 tcp -- 0.0.0.0/0 192.168.3.196 /* default/nginx-tls: cluster IP */ tcp dpt:80
KUBE-SVC-DTRFSG8P77F0F4 tcp -- 0.0.0.0/0 192.168.3.10 /* kube-system/kube-dashboard: cluster IP */ tcp dpt:553
```

可以看到在PREROUTING环节，k8s设置了一个target: KUBE-SERVICES，而KUBE-SERVICES下面又设置了许多target，一旦destination和dport匹配，就会沿着chain进行处理。

比如：当我们在pod网络curl 192.168.3.22 8080时，匹配到下面的KUBE-SVC-AU252PRZZQGOERSG target：

```
KUBE-SVC-AU252PRZZQGOERSG tcp -- 0.0.0.0/0 192.168.3.22 /* default/rbd-rest-api: cluster IP */ tcp dpt:8080
```

沿着target，我们得到KUBE-SVC-AU252PRZZQGOERSG对应的内容如下：

```
Chain KUBE-SVC-AU252PRZZQGOERSG (1 references)
target prot opt source destination
KUBE-SEP-T6L4LR3J077P0RX all -- 0.0.0.0/0 0.0.0.0/0 /* default/rbd-rest-api */
KUBE-SEP-LBWOKUH4CUTN7XKH all -- 0.0.0.0/0 0.0.0.0/0 /* default/rbd-rest-api */

Chain KUBE-SEP-T6L4LR3J077P0RX (1 references)
target prot opt source destination
KUBE-MARK-MARK all -- 172.16.99.6 0.0.0.0/0 /* default/rbd-rest-api */
DNAT tcp -- 0.0.0.0/0 0.0.0.0/0 /* default/rbd-rest-api */ tcp to:172.16.99.618080

Chain KUBE-SEP-LBWOKUH4CUTN7XKH (1 references)
target prot opt source destination
KUBE-MARK-MARK all -- 172.16.99.7 0.0.0.0/0 /* default/rbd-rest-api */
DNAT tcp -- 0.0.0.0/0 0.0.0.0/0 /* default/rbd-rest-api */ tcp to:172.16.99.718080

Chain KUBE-MARK-MARK (17 references)
target prot opt source destination
MARK all -- 0.0.0.0/0 0.0.0.0/0 MARK or 0x000
```

请求被按5：5开的比例分发（起到负载均衡的作用）到KUBE-SEP-6L4LR3J077P0RX和KUBE-SEP-LBWOKUH4CUTN7XKH，而这两个chain的处理方式都是一样的，那就是先做mark，然后做dnat，将service ip改为pod network中的Pod IP，进而请求被实际传输到某个service下面的pod中处理了。

五、参考资料

- [How VXLAN works on Linux&NTP implementation with Flannel](#)
- [Virtual switching technologies and Linux bridge](#)
- [How Flannel's VXLAN backend works](#) 建议用google翻译得网页从日文翻译成英文再看“0”。
- [Software Defined Networking using VXLAN](#)

© 2017, [beyond](#)，版权所有。

0

Related posts:

1. [使用Kubernetes安装Kubernetes Part2](#)

2. [使用Kubernetes安装Kubernetes](#)

3. [一篇文章带你了解Kubernetes安装](#)

4. [理解Docker多主机高可用](#)

5. [理解Docker高可用原理](#)

5条评论

lucyGao

真的很好帮助。

1月24日

回复

顶

转发

王博

从头到尾，点到每个细节，写得真好！

3月4日

回复

顶

转发

社交账号登录:

微信

微博

QQ

人人

更多...

说点什么吧...

发布

加入关键字搜索

搜索



Tony Bai正在使用多链

这里是Tony Bai的个人Blog，欢迎访问、订阅和留言！[订阅Feed](#)请点击[上面图片](#)。

如果您觉得这里的文章对您有帮助，请扫描上方二维码进行捐赠，加油后的Tony Bai将会为您呈现更多精彩的文章，谢谢！

如果您喜欢通过微信App浏览本站内容，可以扫描下方二维码，订阅本站官方微信订阅号“iamtonybai”；点击二维码，可直达本人官方微博主页^_^：



本站Powered by Digital Ocean VPS。

[选择Digital Ocean VPS主机](#)，即可获得10美元现金充值，可免费使用两个月哟！

著名主机提供商Linux 10\$优惠码：linode10，在[这里注册](#)即可免费获得。

阿里云优惠券码：IWFZ0V，[立享9折](#)！

[bigdata.cn @ CSDN.com](#)



文章

- [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#)
- [在Kubernetes Pod中使用Service Account访问API Server](#)
- [在Kubernetes集群Pod中使用Host的本地磁盘](#)
- [Kubernetes Pod无法挂载Host的RBD存储的临时解决方法](#)
- [Kubernetes集群中Service的滚动更新](#)
- [TensorFlow入门：零基础建立一个神经网络](#)
- [Go 1.8中值得关注的几个变化](#)
- [在Kubernetes方式安装Kubernetes集群的坑](#)
- [Kubernetes Dashboard管理界面](#)
- [Kubernetes集群Dashboard组件安装](#)

评论

-  [Hux](#) 在 [Go 1.8中值得关注的几个变化](#) 多谢一赞的详细分享。
-  阿星 在 [使用Kubernetes安装Kubernetes](#) 网络。这个文件的内容是run/flannel/subnet.env，master节点...
-  王博 在 [理解Kubernetes网络之Flume网络](#) 从头到尾，点到每个细节，写得真好！
-  [Eridofei](#) 在 [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#) 其实我觉得，脚本码shana要能打开页面，能连接e获取数据是第二步...
-  [bigshike](#) 在 [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#) 似乎还是网络问题，这个name似乎和你的假想：https://discuss.elastic.co...
-  [Eridofei](#) 在 [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#) 我重新刷了kibana的pc，起来还是这样...
-  [Eridofei](#) 在 [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#) 你觉得还有什么可能性造成timeout
-  [bigshike](#) 在 [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#) 看这里的os的日志也都很正常，暂没看出问题在哪。
-  [Eridofei](#) 在 [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#) 而且我能够访问到\$SF curl -L http://localhost:8080/api/v1/pr...
-  [Eridofei](#) 在 [使用Flume和Elastic-Search Stack实现Kubernetes的集群Logging](#) # kubect logs elasticsearch-logging-v1-2akln --ta...
- 下一页 »

分类

- [杂影汇](#) (7)
- [影音区](#) (36)
- [资源区](#) (64)
- [技术区](#) (506)
- [教育区](#) (1)
- [杂闻区](#) (75)
- [生活区](#) (153)
- [测试区](#) (14)
- [读书区](#) (14)
- [运动区](#) (76)
- [交友区](#) (40)

标签

[Blog](#) [Blogger](#) [C](#) [Can](#) [docker](#) [traefik](#) [GCC](#) [c](#) [with](#) [GNU](#) [Go](#) [Golang](#) [Google](#) [Java](#) [Keras](#) [Kubernetes](#) [Linux](#) [mysql](#) [Opensource](#) [Programmer](#) [Python](#) [Tensorflow](#) [Ubuntu](#) [Unix](#) [Windows](#) [虚拟机](#) [博客](#) [工作](#) [生活](#) [开源](#) [思考](#) [感悟](#) [摄影](#) [音乐](#) [生活](#) [程序](#) [程序员](#) [测试](#) [网络](#) [足球](#) [篮球](#)

归档

- [2017年三月](#) (2)
- [2017年二月](#) (6)
- [2017年一月](#) (7)
- [2016年十二月](#) (7)
- [2016年十一月](#) (7)
- [2016年十月](#) (3)
- [2016年九月](#) (2)
- [2016年八月](#) (1)
- [2016年七月](#) (2)
- [2016年六月](#) (2)
- [2016年四月](#) (2)
- [2016年三月](#) (2)
- [2016年二月](#) (2)
- [2016年一月](#) (2)
- [2015年十二月](#) (1)
- [2015年十一月](#) (1)
- [2015年十月](#) (1)
- [2015年九月](#) (3)
- [2015年八月](#) (5)
- [2015年七月](#) (6)
- [2015年六月](#) (4)
- [2015年五月](#) (1)
- [2015年四月](#) (2)
- [2015年三月](#) (2)
- [2015年二月](#) (2)
- [2014年十二月](#) (5)
- [2014年十一月](#) (8)
- [2014年十月](#) (9)
- [2014年九月](#) (2)
- [2014年八月](#) (1)
- [2014年七月](#) (1)
- [2014年六月](#) (2)
- [2014年五月](#) (2)
- [2014年四月](#) (5)
- [2014年三月](#) (4)
- [2014年二月](#) (1)
- [2014年一月](#) (1)
- [2013年十二月](#) (3)
- [2013年十一月](#) (5)
- [2013年十月](#) (6)
- [2013年九月](#) (4)
- [2013年八月](#) (5)
- [2013年七月](#) (6)
- [2013年六月](#) (2)
- [2013年五月](#) (6)
- [2013年四月](#) (2)
- [2013年三月](#) (7)
- [2013年二月](#) (4)
- [2013年一月](#) (6)
- [2012年十二月](#) (8)
- [2012年十一月](#) (10)
- [2012年十月](#) (5)
- [2012年九月](#) (3)
- [2012年八月](#) (10)
- [2012年七月](#) (4)
- [2012年六月](#) (2)
- [2012年五月](#) (4)
- [2012年四月](#) (10)
- [2012年三月](#) (6)
- [2012年二月](#) (6)
- [2012年一月](#) (6)
- [2011年十二月](#) (4)
- [2011年十一月](#) (4)
- [2011年十月](#) (5)
- [2011年九月](#) (8)

