

20221010-Pod-Pod(Different Node)-[VxLAN Mode]

```

1 env 介绍: [kernel: 5.11.0-051100-generic]
2 root@bpf1:~# uname -a
3 Linux bpf1 5.11.0-051100-generic #202102142330 SMP Sun Feb 14 23:33:21 UTC 2021 x86_64 x86_64 x86_64 GNU/Linux
4 root@bpf1:~# kubectl get nodes -o wide
5 NAME STATUS ROLES AGE VERSION INTERNAL-IP EXTERNAL-IP OS-IMAGE KERNEL-VERSION CONTAINER-RUNTIME
6 bpf1 Ready control-plane,master 19h v1.23.2 192.168.2.61 <none> Ubuntu 20.04.3 LTS 5.11.0-051100-generic docker://20.10.12
7 bpf2 Ready <none> 19h v1.23.2 192.168.2.62 <none> Ubuntu 20.04.3 LTS 5.11.0-051100-generic docker://20.10.12
8 root@bpf1:~#
9 root@bpf1:~# kubectl -nkube-system exec -it cilium-dqnsk -- cilium status
10 Default container "cilium-agent" out of: cilium-agent, mount-cgroup (init), clean-cilium-state (init)
11 KVStore: Ok Disabled
12 Kubernetes: Ok 1.23 (v1.23.2) [linux/amd64]
13 Kubernetes APIs: ["cilium/v2::CiliumClusterwideNetworkPolicy", "cilium/v2::CiliumEndpoint", "cilium/v2::CiliumNetworkPolicy", "cilium/v2::CiliumNode", "core/v1::Namespace", "core/v1::Node", "core/v1::Pods", "core/v1::Service", "discovery/v1::EndpointSlice", "networking.k8s.io/v1::NetworkPolicy"]
14 KubeProxyReplacement: Strict [ens3 192.168.2.61 (Direct Routing)]
15 Host firewall: Disabled
16 Cilium: Ok 1.11.1 (v1.11.1-76d34db)
17 NodeMonitor: Disabled
18 Cilium health daemon: Ok
19 IPAM: IPv4: 6/254 allocated from 10.0.1.0/24,
20 BandwidthManager: Disabled
21 Host Routing: BPF
22 Masquerading: BPF [ens3] 10.0.1.0/24 [IPv4: Enabled, IPv6: Disabled]
23 Controller Status: 39/39 healthy
24 Proxy Status: OK, ip 10.0.1.159, 0 redirects active on ports 10000-20000
25 Hubble: Disabled
26 Encryption: Disabled
27 Cluster health: 2/2 reachable (2022-01-22T08:42:27Z)
28 root@bpf1:~#

```

1.Pod-Pod DIFF Node[轻量级隧道类型]

```

1 对于不同节点的通信，是我们通常需要关注的重点。此例以VxLAN Backend来说明此过程。对于VxLAN实际上我们已经相对来说较为熟悉了，但是在Cilium中VxLAN稍微有些不同：
2 root@bpf1:~# ip -d link show cilium_vxlan
3 6: cilium_vxlan: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc noqueue state UNKNOWN mode DEFAULT group default qlen 1000
4     link/ether de:41:40:10:e3:62 brd ff:ff:ff:ff:ff:ff promiscuity 0 minmtu 68 maxmtu 65535
5     vxlan external addrgenmode eui64 numtxqueues 1 numrxqueues 1 gso_max_size 65536 gso_max_segs 65535      # 1.这里有一个external的字段，需要做一个说明:
6 root@bpf1:~#
7 root@bpf1:~# man ip link -h #.查看帮助
8   VXLAN Type Support
9       For a link of type VXLAN the following additional arguments are supported:
10
11     ip link add DEVICE type vxlan id VNI [ dev PHYS_DEV ] [ { group | remote } IPADDR ] [ local { IPADDR | any } ] [ ttl TTL ] [ tos
12     TOS ] [ df DF ] [ flowlabel FLOWLABEL ] [ dstport PORT ] [ srcport MIN MAX ] [ [no]learning ] [ [no]proxy ] [ [no]rsc ] [
13     [no]l2miss ] [ [no]l3miss ] [ [no]udpcsum ] [ [no]udp6zerocsumtx ] [ [no]udp6zerocsumrx ] [ ageing SECONDS ] [ maxaddress NUMBER
14     ] [ [no]external ] [ gbp ] [ gpe ]
15
16         id VNI - specifies the VXLAN Network Identifier (or VXLAN Segment Identifier) to use.
17
18         dev PHYS_DEV - specifies the physical device to use for tunnel endpoint communication.
19
20         group IPADDR - specifies the multicast IP address to join. This parameter cannot be specified with the remote parameter.
21
22         remote IPADDR - specifies the unicast destination IP address to use in outgoing packets when the destination link layer
23         address is not known in the VXLAN device forwarding database. This parameter cannot be specified with the group parame-
24         ter.
25
26         local IPADDR - specifies the source IP address to use in outgoing packets.
27
28         ttl TTL - specifies the TTL value to use in outgoing packets.
29
30         tos TOS - specifies the TOS value to use in outgoing packets.
31
32         df DF - specifies the usage of the Don't Fragment flag (DF) in outgoing packets with IPv4 headers. The value inherit
33         causes the bit to be copied from the original IP header. The values unset and set cause the bit to be always unset or al-
34         ways set, respectively. By default, the bit is not set.
35
36         flowlabel FLOWLABEL - specifies the flow label to use in outgoing packets.
37
38         dstport PORT - specifies the UDP destination port to communicate to the remote
39             VXLAN tunnel endpoint.
40
41         srcport MIN MAX - specifies the range of port numbers to use as UDP source ports to communicate to the remote VXLAN tun-
42             nel endpoint.
43
44         [no]learning - specifies if unknown source link layer addresses and IP addresses are entered into the VXLAN device for-
45             warding database.
46
47         [no]rsc - specifies if route short circuit is turned on.
48
49         [no]proxy - specifies ARP proxy is turned on.
50
51         [no]l2miss - specifies if netlink LLADDR miss notifications are generated.
52
53         [no]l3miss - specifies if netlink IP ADDR miss notifications are generated.
54
55         [no]udpcsum - specifies if UDP checksum is calculated for transmitted packets over IPv4.
56
57         [no]udp6zerocsumtx - skip UDP checksum calculation for transmitted packets over IPv6.
58
59         [no]udp6zerocsumrx - allow incoming UDP packets over IPv6 with zero checksum field.
60
61         ageing SECONDS - specifies the lifetime in seconds of FDB entries learnt by the kernel.
62
63         maxaddress NUMBER - specifies the maximum number of FDB entries.
64
65         [no]external - specifies whether an external control plane (e.g. ip route encaps) or the internal FDB should be used.
66
67         gbp - enables the Group Policy extension (VXLAN-GBP).

```

```

68
69     Allows to transport group policy context across VXLAN network peers. If enabled, includes the mark of a packet in
70     the VXLAN header for outgoing packets and fills the packet mark based on the information found in the VXLAN header
71     for incoming packets.
72
73     Format of upper 16 bits of packet mark (flags):
74
75     +-----+
76     | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
77     +-----+
78
79     D := Don't Learn bit. When set, this bit indicates that the egress VTEP MUST NOT learn the source address of the
80     encapsulated frame.
81
82     A := Indicates that the group policy has already been applied to this packet. Policies MUST NOT be applied by de-
83     vices when the A bit is set.
84
85     Format of lower 16 bits of packet mark (policy ID):
86
87     +-----+
88     |           Group Policy ID           |
89     +-----+
90
91     Example:
92         iptables -A OUTPUT [...] -j MARK --set-mark 0x800FF
93
94     gpe - enables the Generic Protocol extension (VXLAN-GPE). Currently, this is only supported together with the external
95     keyword.
96 #
97     [no]external - specifies whether an external control plane (e.g. ip route encap) or the internal FDB should be used.
98 1.这里的解释：把隧道信息封装在路由中。
99 #
100 这里既然说到了轻量级Tunnel，我们就多说一点：
101
102 1.首先是常规情况下的Linux kernel中VxLAN配置：
103 ip link add vxlan0 type vxlan id 5 dstport 4789 remote 172.12.1.12 local 172.12.1.11 dev ens33
104 ip addr add 10.20.1.2/24 dev vxlan0
105 ip link set vxlan0 up
106 此时并没有路由相关的配置，那么当实现VxLAN的通信的时候，需要借助于FDB表的查询实现。
107
108 2.轻量级Tunnel：
109 ip link add vxlan1 type vxlan dstport 4789 external
110 ip link set dev vxlan1 up
111 ip addr add 20.1.1.1/24 dev vxlan1
112 ip route add 10.1.1.1 encapsip id 30001 dst 20.1.1.2 dev vxlan1
113 此时我们可以看到，隧道信息被封装在路由中。这时一种基础的轻量级实现方式。
114
115 3.还有一种是BPF介入：
116 有了 BPF 可编程性之后，能为出向流量（入向是只读的）做封装。
117 与 tc 类似，ip route 支持直接将 BPF 程序 attach 到网络设备：
118 $ ip route add 192.168.253.2/32 encapsip bpf out obj lwt_len_hist_kern.o section len_hist dev veth0
119 $ ip route add <route+prefix> encapsip bpf out obj <bpf object file.o> section <ELF section> dev <device>
120
121 且：通常情况下：我们的VxLAN在通信的时候，在同一VNI ID的我们认为是一个Tunnel。但是在Cilium的环境中，我们会有现在 (Request侧) 和 (Replay侧) 使用的VxLAN的 VNI ID是不同的。

```

✓ 2.Pod-Pod DIFF Nodes[cilium monitor -vv]分析

```

Properties ▾ ...
1 分析不同节点Pod之间通信，对应此前我们熟悉的CNI (Calico Flannel) 均是使用路由表，FDB表，ARP表等网络知识便可以分析的非常清楚，但是在Cilium中我们发现此种分析思路便“失效”了。究其原因，是由于Cilium的CNI实现结合eBPF技术实现了
2 datapath的“跳跃式”转发。
3 所以我们需要结合Cilium提供的Tools来辅助分析：
4 我们还是由诸多分析方式来进行分析：
5 1.cilium monitor -vv
6 2.pwru
7 3.iptables TRACE分析
8 4.源码分析(我们不做重点解析，有余力可结合pwru和bpf以及kernel来分析源码)
9 5.tcpdump抓包分析
10
11 首先使用cilium monitor -vv来辅助分析：
12 env介绍：
13 root@bpff1:~# kubectl get pods -o wide
14 NAME      READY   STATUS    RESTARTS   AGE   IP          NODE   NOMINATED NODE   READINESS GATES
15 cni-6vf5d  1/1     Running   0          41h   10.0.1.180  bpff1  <none>        <none>
16 cni-j22kt  1/1     Running   0          41h   10.0.0.23   bpff2  <none>        <none>
17
18 root@bpff1:~#
19
20 先给出地址信息，我们可以在cilium中查询：
21 root@bpff1:/home/cilium# cilium bpf endpoint list
22 IP ADDRESS      LOCAL ENDPOINT INFO
23 10.0.1.180:0   id=3663 flags=0x0000 ifindex=24 mac=FE:55:FA:02:EE:E5 nodemac=D2:F1:29:4B:BD:90 # 说明：这里实际上就是我们pod的eth0网卡和其对应的lxc网卡的信息。
24 root@bpff2:/home/cilium# cilium bpf endpoint list
25 IP ADDRESS      LOCAL ENDPOINT INFO
26 10.0.0.23:0   id=2720 flags=0x0000 ifindex=28 mac=42:4A:CA:D9:16:B5 nodemac=4E:51:12:28:98:B2 # 说明：这里实际上就是我们pod的eth0网卡和其对应的lxc网卡的信息。
27
28 我们在bpff1节点上monitor log：
29 root@bpff1:~# kubectl exec -it cni-6vf5d -- ping -c 1 10.0.0.23
30 PING 10.0.0.23 (10.0.0.23): 56 data bytes
31 64 bytes from 10.0.0.23: seq=0 ttl=63 time=1.067 ms
32
33 --- 10.0.0.23 ping statistics ---
34 1 packets transmitted, 1 packets received, 0% packet loss
35 round-trip min/avg/max = 1.067/1.067/1.067 ms
36
37
38 [1.bpf1节点上的log: ]
39 CPU 05: MARK 0x0 FROM 3663 DEBUG: Conntrack lookup 1/2: src=10.0.1.180:15872 dst=10.0.0.23:0
40 CPU 05: MARK 0x0 FROM 3663 DEBUG: Conntrack lookup 2/2: nexthdr=1 flags=1
41 CPU 05: MARK 0x0 FROM 3663 DEBUG: CT verdict: New, revnat=0
42 CPU 05: MARK 0x0 FROM 3663 DEBUG: Successfully mapped addr=10.0.0.23 to identity=37483
43 CPU 05: MARK 0x0 FROM 3663 DEBUG: Conntrack create: proxy-port=0 revnat=0 src-identity=37483 lb=0.0.0.0
44 CPU 05: MARK 0x0 FROM 3663 DEBUG: Encapsulating to node 3232236094 (0xc0a8023e) from seclabel 37483
45 ----- # 1.这里分为3个部分，第一个部分CT指的是Pod cni-6vf5d 内部数据出来的时候经过的iptables过程。
46 Ethernet  {Contents=[..14..] Payload=[..86..] SrcMAC=fe:02:ee:e5 DstMAC=d2:f1:29:4b:bd:9d EthernetType=IPv4 Length=0}
47 IPv4  {Contents=[..20..] Payload=[..64..] Version=4 IHL=5 TOS=0 Length=84 Id=55780 Flags=0F FragOffset=0 TTL=64 Protocol=ICMPv4 Checksum=19194 SrcIP=10.0.1.180 DstIP=10.0.0.23 Options=[] Padding=[]}
48 ICMPv4  {Contents=[..8..] Payload=[..56..] TypeCode=EchoRequest Checksum=47137 Id=15872 Seq=0}
49 Failed to decode layer: No decoder for layer type Payload
50 CPU 05: MARK 0x0 FROM 3663 toOverlay: 98 bytes (98 captured), state new, interface cilium_vxlan, , identity 37483->unknown, orig-ip 0.0.0.0
51 CPU 05: MARK 0x0 FROM 83 DEBUG: Conntrack lookup 1/2: src=192.168.2.61:36535 dst=192.168.2.62:8472
52 CPU 05: MARK 0x0 FROM 83 DEBUG: Conntrack lookup 2/2: nexthdr=17 flags=1
53 CPU 05: MARK 0x0 FROM 83 DEBUG: CT verdict: New, revnat=0
54

```

✓ 3.tcpdump抓包分析

```
1 root@bpf1:~# kubectl get pods -o wide
2   NAME      READY   STATUS    RESTARTS   AGE     IP           NODE   NOMINATED NODE   READINESS GATES
3   cni-6vf5d   1/1    Running   0          47h    10.0.1.180   bpf1   <none>        <none>
4   cni-j22kt   1/1    Running   0          47h    10.0.0.23    bpf2   <none>        <none>
5   root@bpf1:~#
6   先给出地址信息，我们可以在cilium中查询：
7   root@bpf1:/home/cilium# cilium bpf endpoint list
8   IP ADDRESS      LOCAL ENDPOINT INFO
9   10.0.1.180:0    id=3663  flags=0x0000 ifindex=24  mac=FE:55:FA:02:EE:E5 nodemac=D2:F1:29:4B:BD:9D  # 说明：这里实际上就是我们bpf1上pod的eth0网卡和其对应的lxc网卡的信息。
10  root@bpf2:/home/cilium# cilium bpf endpoint list
11  IP ADDRESS      LOCAL ENDPOINT INFO
12  10.0.0.23:0    id=2720  flags=0x0000 ifindex=28  mac=42:4A:CA:D9:16:B5 nodemac=4E:51:12:28:98:B2  # 说明：这里实际上就是我们bpf2上pod的eth0网卡和其对应的lxc网卡的信息。
13
14 具体抓包过程：[cni-6vf5d]pod数据egress的datapath:
15 [1.pod中的eth0网卡]
16 root@bpf1:~# kubectl exec -it cni-6vf5d -- tcpdump -pne -i eth0
17 tcpdump: verbose output suppressed, use -v[v]... for full protocol decode
18 listening on eth0, link-type EN10MB (Ethernet), snapshot length 262144 bytes
19 13:12:50.047213 fe:55:fa:02:ee:e5 > d2:f1:29:4b:bd:9d, ethertype IPv4 (0x0800), length 98: 10.0.1.180 > 10.0.0.23: ICMP echo request, id 19200, seq 0, length 64 # 1.此时在Pod中抓包，出去的ICMP Request。
20 13:12:50.047930 d2:f1:29:4b:bd:9d > fe:55:fa:02:ee:e5, ethertype IPv4 (0x0800), length 98: 10.0.0.23 > 10.0.1.180: ICMP echo reply, id 19200, seq 0, length 64
21 ^C
22 2 packets captured
```



```

137
138 1429 Disabled Disabled 4 reserved:host
139 1690 Disabled Disabled 44045 reserved:health
140 k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=kube-system
141 k8s:io.cilium.k8s.policy.cluster=default
142 k8s:io.cilium.k8s.policy.serviceaccount=coredns
143 k8s:io.kubernetes.pod.namespace=kube-system
144 2447 Disabled Disabled 44045 k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=kube-system 10.0.1.93 ready
145 k8s:io.cilium.k8s.policy.cluster=default
146 k8s:io.cilium.k8s.policy.serviceaccount=coredns
147 k8s:io.kubernetes.pod.namespace=kube-system
148 k8s:k8s-app=kube-dns
149 3663 Disabled Disabled 37483 k8s:app:cni
150 k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=default
151 k8s:io.cilium.k8s.policy.cluster=default
152 k8s:io.cilium.k8s.policy.serviceaccount=default
153 k8s:io.kubernetes.pod.namespace=default
154 root@bpf1:/home/cilium#
155
156 [bpf2 node# cilium endpoint list]
157 root@bpf2:/home/cilium# cilium endpoint list
158 ENDPOINT POLICY (ingress) POLICY (egress) IDENTITY LABELS (source:key[=value]) IPv6 IPv4 STATUS
159 ENFORCEMENT ENFORCEMENT
160 415 Disabled Disabled 1 reserved:host
161 2720 Disabled Disabled 37483 k8s:app:cni
162 k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=default
163 k8s:io.cilium.k8s.policy.cluster=default
164 k8s:io.cilium.k8s.policy.serviceaccount=default
165 k8s:io.kubernetes.pod.namespace=default
166 3590 Disabled Disabled 4 reserved:health
167 root@bpf2:/home/cilium#

```

train/clium datapath.png at main · BurlyLuo/train
Contribute to BurlyLuo/train development by creating an account on GitHub.
[GitHub](#)

BurlyLuo/train

ENDPOINT	POLICY (ingress)	POLICY (egress)	IDENTITY	LABELS (source:key[=value])	IPv6	IPv4	STATUS
83	Disabled	Disabled	1	k8s:node-role.kubernetes.io/control-plane k8s:node-role.kubernetes.io/master k8s:node.kubernetes.io/exclude-from-external-load-balancers reserved:host			ready
1429	Disabled	Disabled	4	reserved:health			
1690	Disabled	Disabled	44045	k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=kube-system k8s:io.cilium.k8s.policy.cluster=default k8s:io.cilium.k8s.policy.serviceaccount=coredns k8s:io.kubernetes.pod.namespace=kube-system	10.0.1.244	10.0.1.71	ready
2447	Disabled	Disabled	44045	k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=kube-system k8s:io.cilium.k8s.policy.cluster=default k8s:io.cilium.k8s.policy.serviceaccount=coredns k8s:io.kubernetes.pod.namespace=kube-system	10.0.1.93		ready
3663	Disabled	Disabled	37483	k8s:app:cni k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=default k8s:io.cilium.k8s.policy.cluster=default k8s:io.cilium.k8s.policy.serviceaccount=default k8s:io.kubernetes.pod.namespace=default	10.0.1.180		ready