# Liberating K8s from kube-proxy and iptables (and netfilter)

Martynas Pumputis, Cilium
(Daniel Borkmann, Thomas Graf, André Martins)

cilium

# Performance

```
# perf top -a -e cycles:k

PerfTop:  16326 irqs/sec  (all, 4 CPUs)

------------------------------------------------------------------------------------

      8.79% [kernel]         [k] native_sched_clock
      4.99% [ip_tables]      [k] ipt_do_table
      3.09% [e1000e]         [k] e1000_irq_enable
      2.51% [nf_conntrack]   [k] __nf_conntrack_find_get
      2.03% [kernel]         [k] fib_table_lookup
      1.98% [kernel]         [k] sched_clock_cpu
      1.75% [nf_conntrack]   [k] tcp_packet
      1.65% [nf_conntrack]   [k] nf_conntrack_tuple_taken
      [...]
```

# Reliability

DNS intermittent delays of 5s #56903

**Closed** **mikksoone** opened this issue on Dec 6, 2017 · 230 comments

**mikksoone** commented on Dec 6, 2017 · edited ▾

**Is this a BUG REPORT or FEATURE REQUEST?:**
/kind bug

**What happened:**
DNS lookup is sometimes taking 5 seconds.

**What you expected to happen:**
No delays in DNS.

**Assignees**
No one assigned

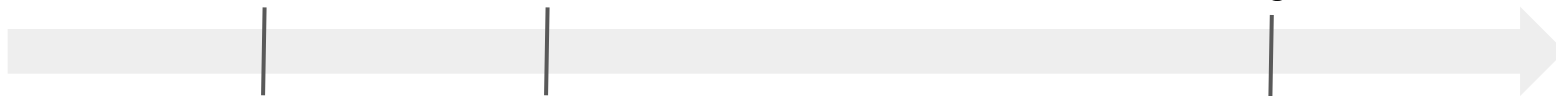**Labels**
area/dns
kind/bug
sig/network

Root cause

Patches submitted

Patches merged

May 27, 2018

Aug 5, 2018

Feb 11, 2019

# Reliability

DNS intermittent delays of 5s #56903

Closed  **mikksoone** opened this issue on Dec 6, 2017 · 230 comments

**mikksoone** commented on Dec 6, 2017 · edited ▾

**Is this a BUG REPORT or FEATURE REQUEST?**:
/kind bug

**What happened**:
DNS lookup is sometimes taking 5 seconds.

**What you expected to happen**:
No delays in DNS.

Assignees
No one assigned

Labels
area/dns
kind/bug
sig/network

First occurance
of bug

Patches
merged

Nov 11, 2010

Feb 11, 2019

# Debuggability

```
# iptables-save -c

*filter
:INPUT ACCEPT [0:0]
:FORWARD ACCEPT [0:0]
:OUTPUT ACCEPT [0:0]
[1:10] -A FORWARD -i eth0 -s 172.17.0.0/16 -j DROP
```

# Debuggability



https://www.reddit.com/r/networkingmemes/comments/8u7jyz/container_networking/
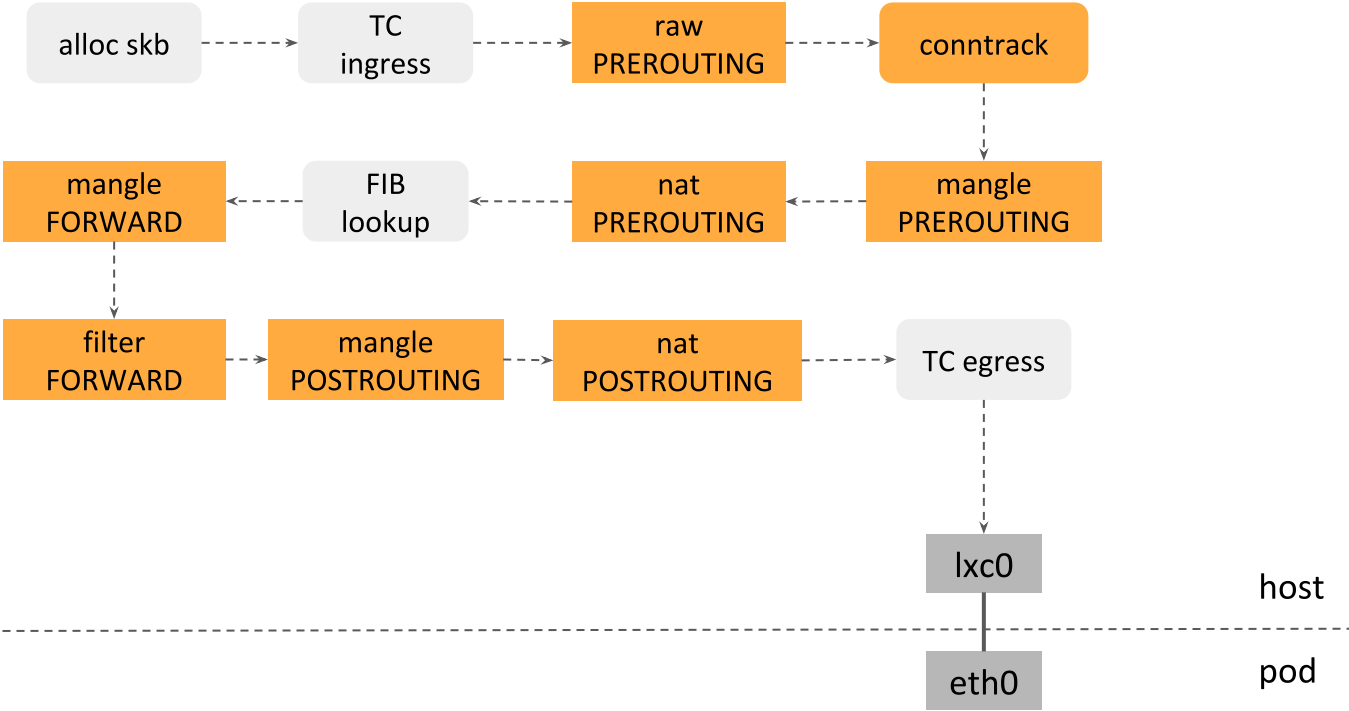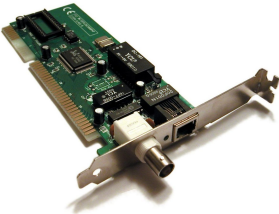
# Compatibility

kube-proxy currently incompatible with `iptables >= 1.8`
#71305

**⊘ Open** **drags** opened this issue on Nov 21, 2018 · 75 comments · May be fixed by #82966 or #84420

# Packet flow



alloc skb → TC ingress → raw PREROUTING → conntrack

mangle FORWARD ← FIB lookup ← nat PREROUTING ← mangle PREROUTING

filter FORWARD → mangle POSTROUTING → nat POSTROUTING → TC egress

lxc0

host

eth0

pod

# ClusterIP with iptables

```
$ kubectl get svc nginx
NAME   TYPE          CLUSTER-IP    EXTERNAL-IP   PORT(S)
nginx  ClusterIP     3.3.3.3       <none>        80/TCP

$ kubectl get endpoints nginx
NAME   ENDPOINTS
nginx  1.1.1.1:80, 1.1.2.2:80
```

```
-t nat -A PREROUTING -m conntrack --ctstate NEW -j KUBE-SERVICES

-A KUBE-SERVICES ! -s 1.1.0.0/16 -d 3.3.3.3/32 -p tcp -m tcp --dport 80 -j KUBE-MARK-MASQ
-A KUBE-SERVICES -d 3.3.3.3/32 -p tcp -m tcp --dport 80 -j KUBE-SVC-NGINX

-A KUBE-SVC-NGINX -m statistic --mode random --probability 0.50 -j KUBE-SEP-NGINX1
-A KUBE-SVC-NGINX -j KUBE-SEP-NGINX2

-A KUBE-SEP-NGINX1 -s 1.1.1.1/32 -j KUBE-MARK-MASQ
-A KUBE-SEP-NGINX1 -p tcp -m tcp -j DNAT --to-destination 1.1.1.1:80
-A KUBE-SEP-NGINX2 -s 1.1.2.2/32 -j KUBE-MARK-MASQ
-A KUBE-SEP-NGINX2 -p tcp -m tcp -j DNAT --to-destination 1.1.2.2:80
```
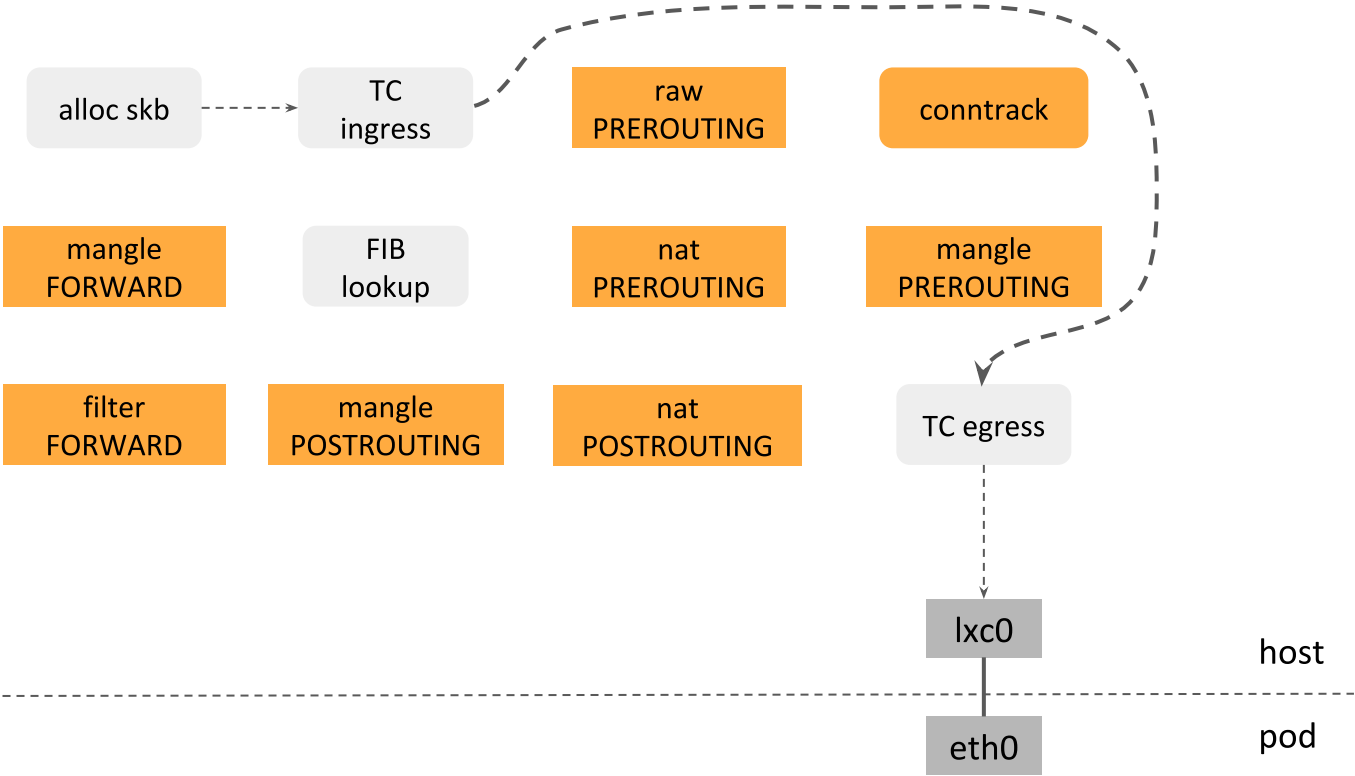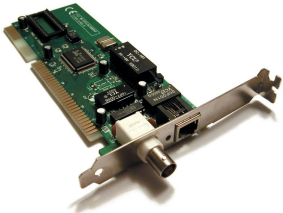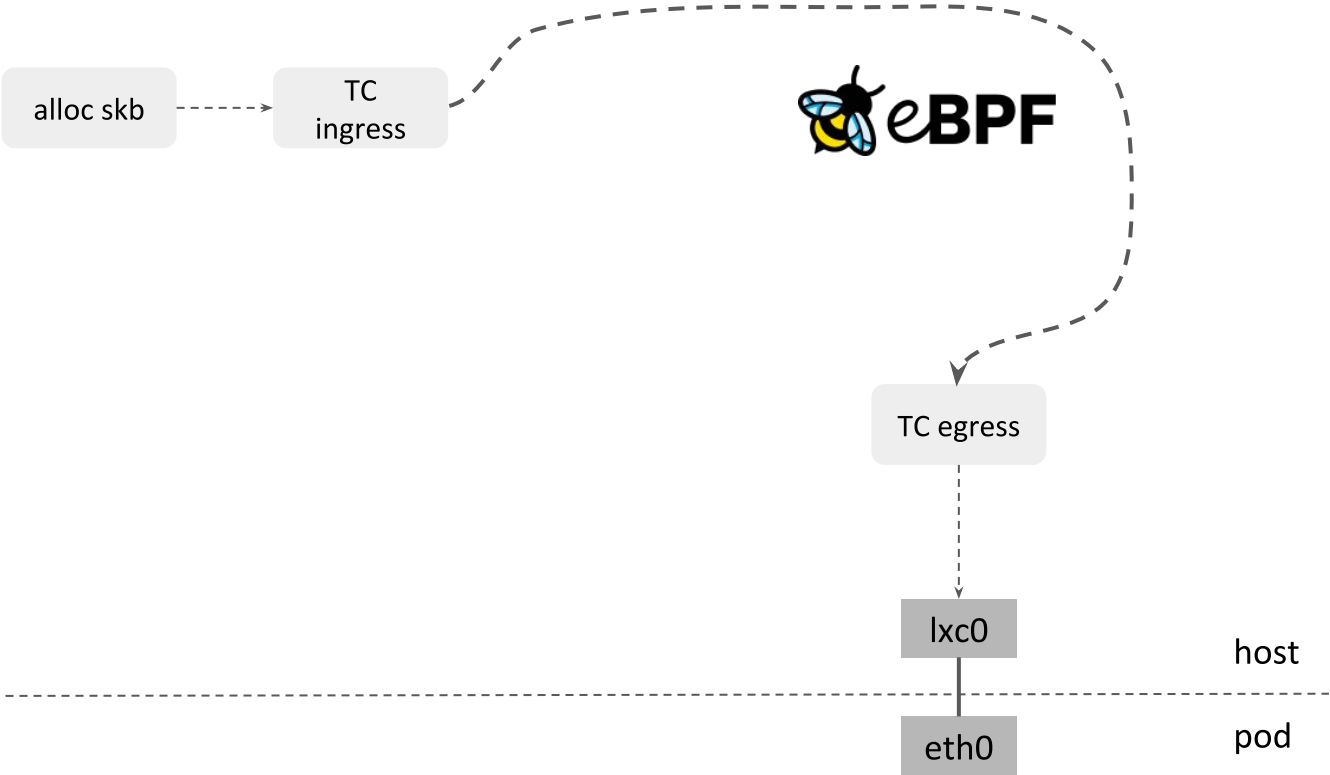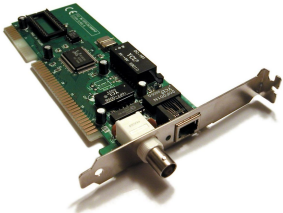
nat
PREROUTING

# Packet flow



alloc skb → TC ingress

raw PREROUTING

conntrack

mangle FORWARD

FIB lookup

nat PREROUTING

mangle PREROUTING

filter FORWARD

mangle POSTROUTING

nat POSTROUTING

TC egress

lxc0

host

eth0

pod

# Packet flow



alloc skb

TC ingress

eBPF

TC egress

lxc0

eth0

host

pod

eBPF

```
SEC("to_netdev")
int handle(struct sk_buff *skb) {
    ...
    if (tcp->dport == 80)
        redirect(nginx_pod);
    ...
}
```

agent

BPF maps

native code

eth0

clang -target bpf [...]

foo.o

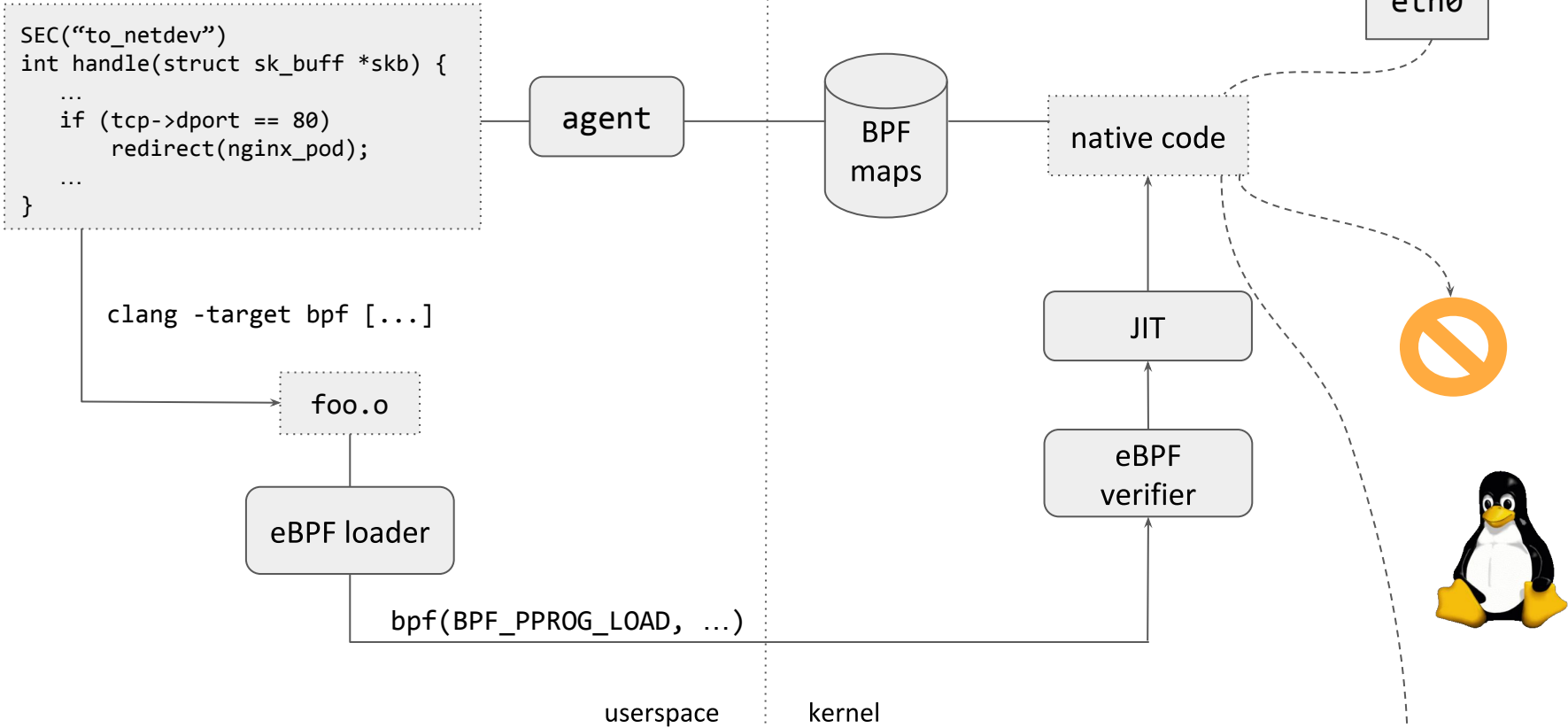eBPF loader

JIT

eBPF verifier

bpf(BPF_PPROG_LOAD, ...)

userspace     kernel

**eBPF**

**268 contributors (Jan 2016 to Nov 2019):**

➢   443  Daniel Borkmann (Cilium; maintainer)
➢   242  Alexei Starovoitov (Facebook; maintainer)
➢   210  Jakub Kicinski  (Netronome)
➢   195  Andrii Nakryiko (Facebook)
➢   161  Yonghong Song (Facebook)
➢   151  Stanislav Fomichev (Google)
➢   145  Quentin Monnet (Netronome)
➢   144  Martin KaFai Lau (Facebook)
➢   139  John Fastabend (Cilium)
➢   118 Jesper Dangaard Brouer (Red Hat)
➢   […]

**Users:**

**TheRustyTwit**
@rusty_twit

Replying to @LaF0rge

Well, iptables perf used to be "mostly good enough". Replacing it has taken so long because it requires a radically different approach; nice to see it finally happening!

12:46 AM · Apr 18, 2018 · Twitter for Android

```
$ kubectl -n kube-system delete ds kube-proxy
```

# kube-proxy

## 1. ClusterIP

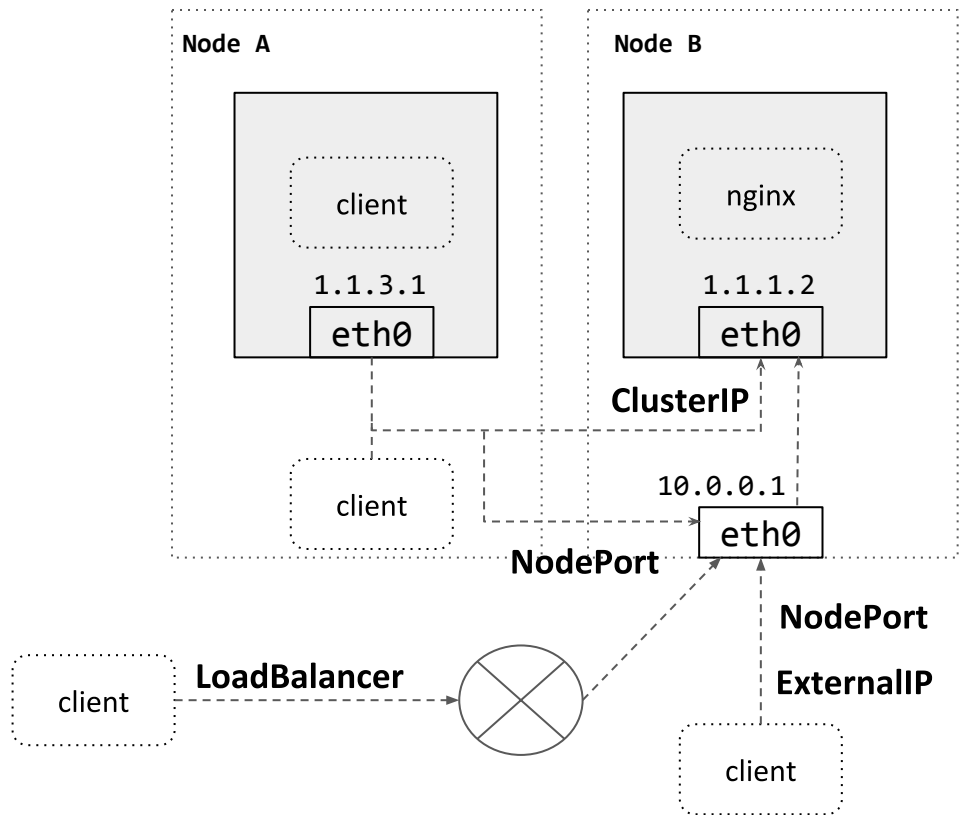- In-cluster access via virtual IP

## 2. NodePort

- Access from outside / inside via node IP + port

## 3. ExternalIP
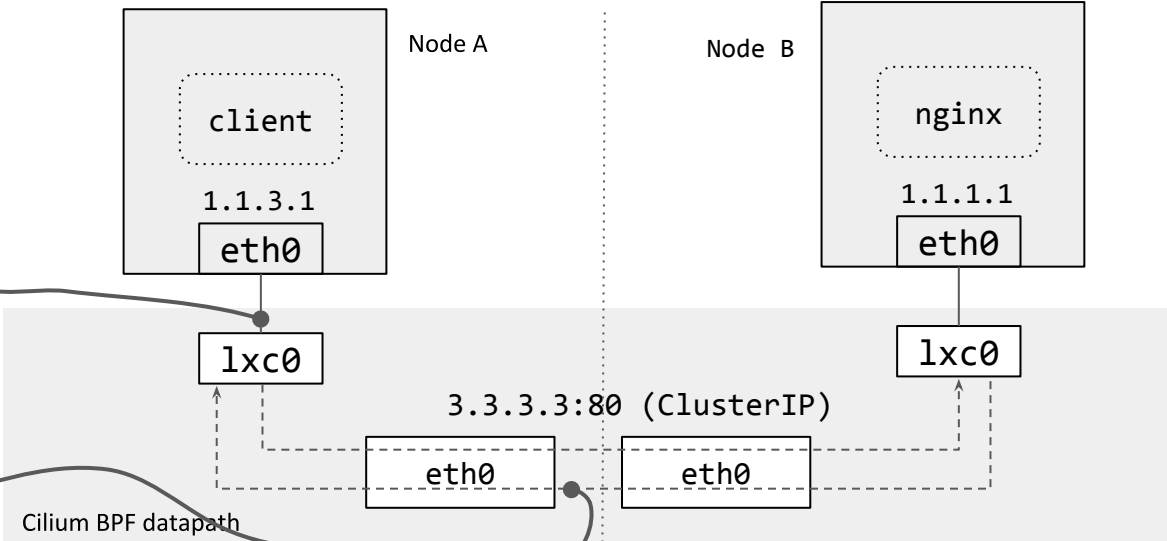
- Access from outside via external IP

## 4. LoadBalancer

- Access from outside via external LB

# ClusterIP (pod to pod) in Cilium

1. Lookup dst in SVC map
2. If found:
   a. Create SVC CT
   b. DNAT
3. Create Egress CT

1. Lookup Egress CT
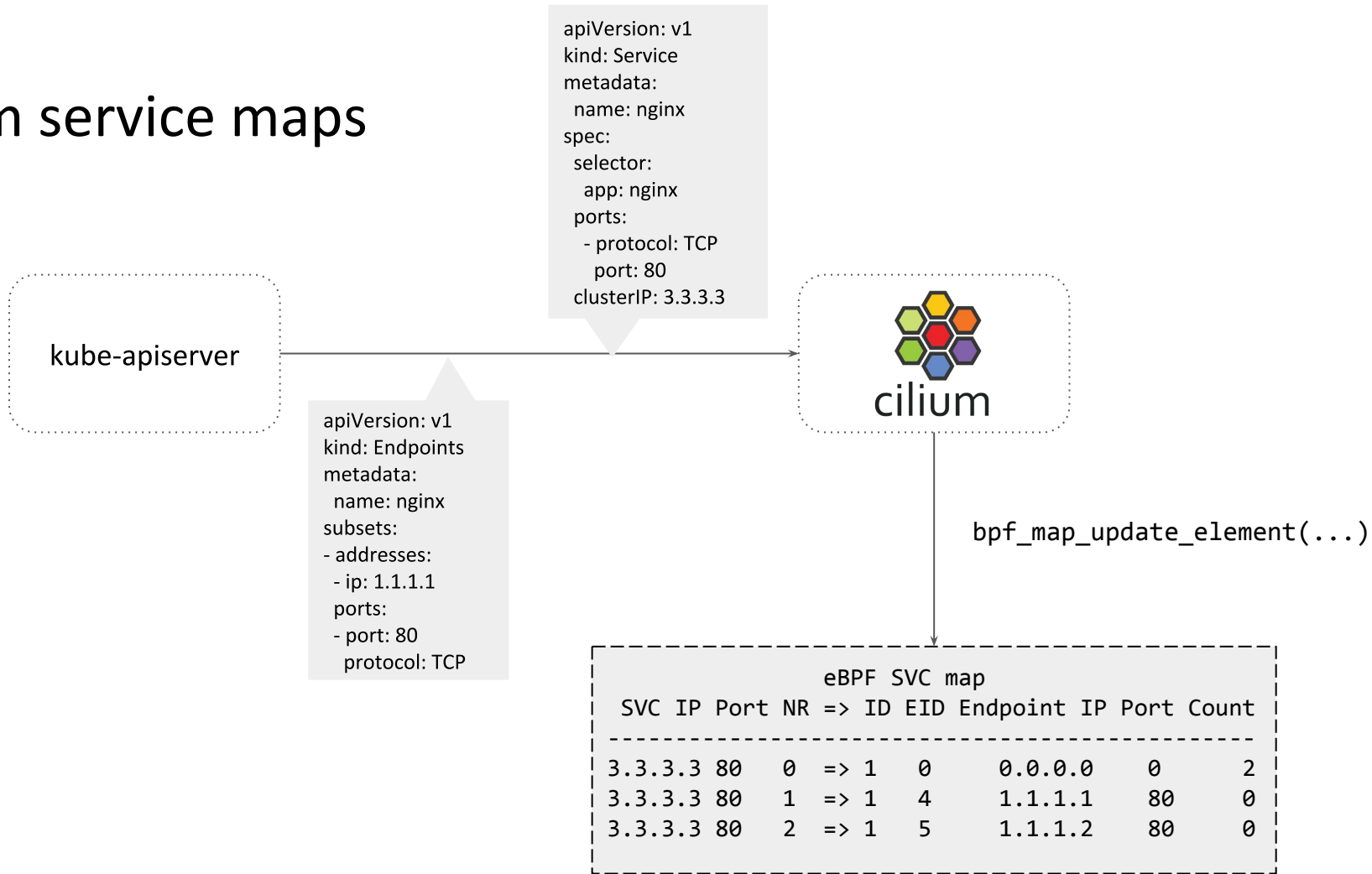2. If found:
   a. Rev-NAT xlation
3. Redirect to lxc0

Node A

Node B

```
client
```

```
nginx
```

1.1.3.1

1.1.1.1

```
eth0
```

```
eth0
```

```
lxc0
```

```
lxc0
```

3.3.3.3:80 (ClusterIP)

```
eth0
```

```
eth0
```

Cilium BPF datapath

```
                      eBPF SVC map
 SVC IP Port NR => ID EID Endpoint IP Port Count
 -----------------------------------------------
 3.3.3.3 80   0  => 1  0   0.0.0.0     0      2
 3.3.3.3 80   1  => 1  4   1.1.1.1     80     0
 3.3.3.3 80   2  => 1  5   1.1.1.2     80     0
```

```
                       eBPF conntrack LRU map
  srcIP    sPort  dstIP    dPort   Type      => EID|SVCID
  ------------------------------------------------------
  1.1.3.1  4321   3.3.3.3  80      SVC       =>  4
  1.1.3.1  4321   1.1.1.1  80      Egress    =>       1
  1.1.1.1  80     1.1.3.1  4321    Ingress =>
```

# Cilium service maps

```
apiVersion: v1
kind: Service
metadata:
 name: nginx
spec:
 selector:
  app: nginx
 ports:
  - protocol: TCP
    port: 80
 clusterIP: 3.3.3.3
```

kube-apiserver



```
apiVersion: v1
kind: Endpoints
metadata:
 name: nginx
subsets:
- addresses:
 - ip: 1.1.1.1
 ports:
 - port: 80
   protocol: TCP
```

`bpf_map_update_element(...)`

```
                      eBPF SVC map
  SVC IP Port NR => ID EID Endpoint IP Port Count
  --------------------------------------------------
  3.3.3.3 80   0  => 1   0   0.0.0.0     0     2
  3.3.3.3 80   1  => 1   4   1.1.1.1     80    0
  3.3.3.3 80   2  => 1   5   1.1.1.2     80    0
```
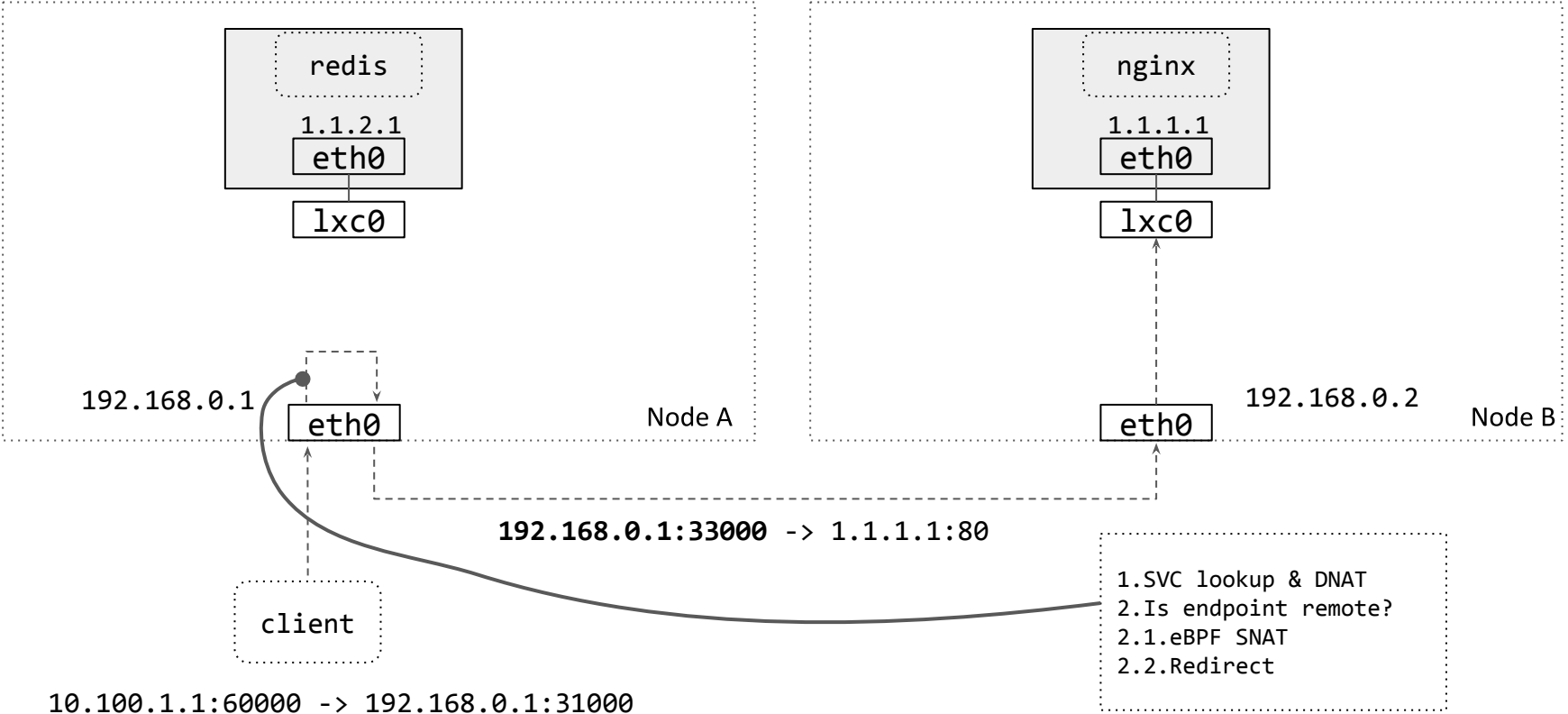
# ClusterIP (host or pod to pod) in Cilium

UDP

TCP

**connect()**

```
import "net/http"

func main() {
 r, err := http.Get("3.3.3.3")
 ...
}
```

client

nginx

1.1.1.1

eth0

lxc0

3.3.3.3:80 (ClusterIP)

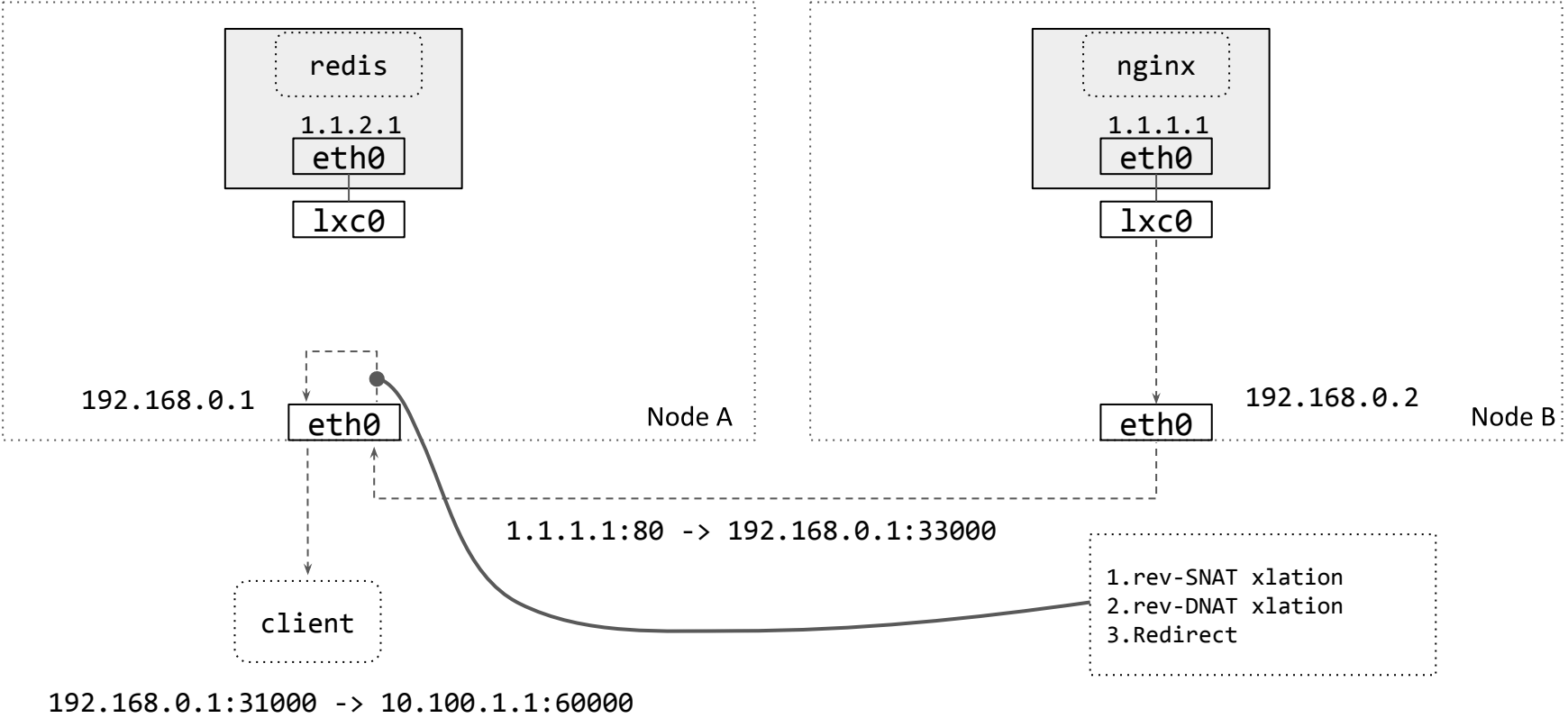1. Lookup dst in SVC map
2. If found:
   a. Change dst addr
      and port in socket

kernel

# NodePort with service endpoint on local node in Cilium

# NodePort with service endpoint on remote node in Cilium

redis

1.1.2.1

eth0

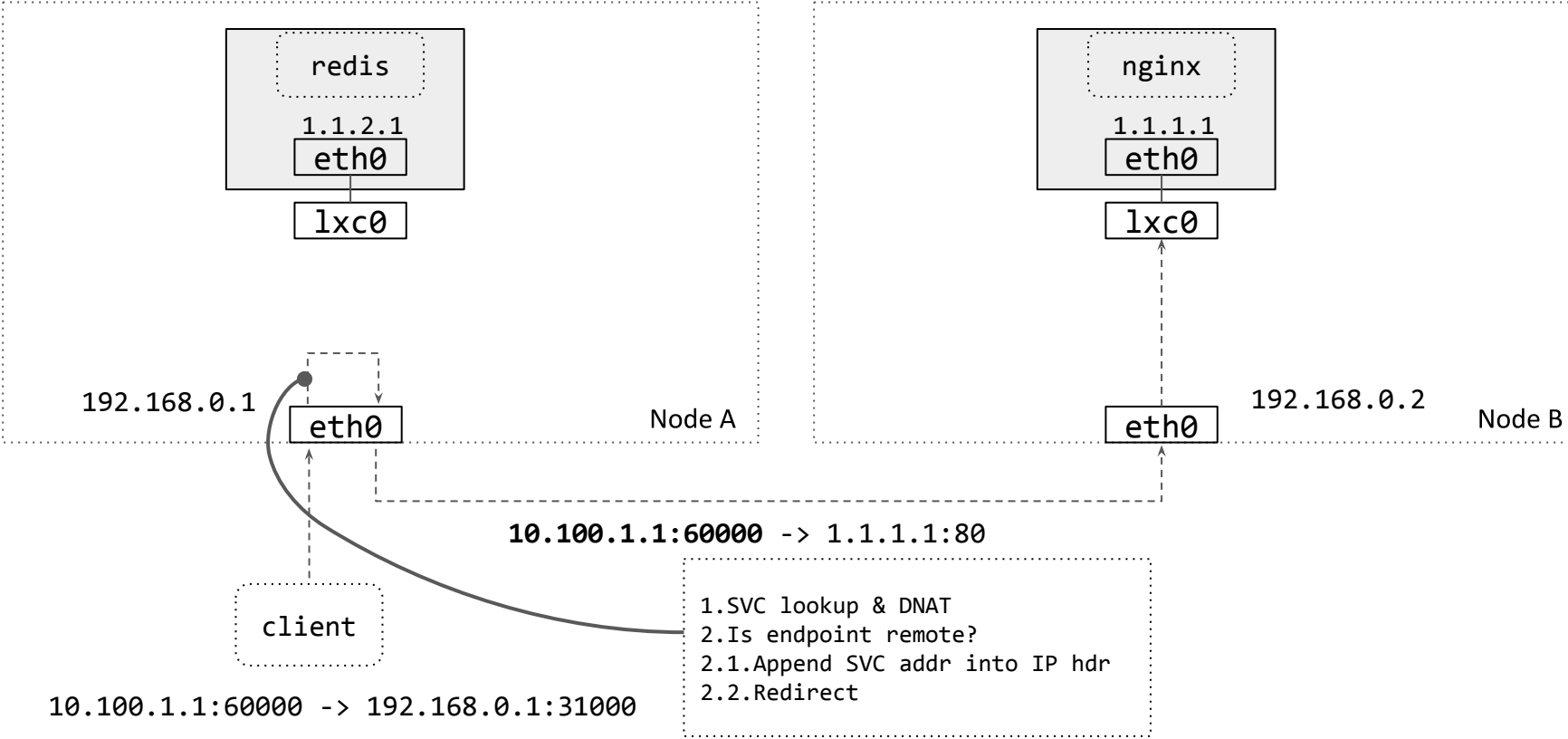lxc0

nginx

1.1.1.1

eth0

lxc0

192.168.0.1

eth0

Node A

192.168.0.2

eth0

Node B

**192.168.0.1:33000** -> 1.1.1.1:80

client

1.SVC lookup & DNAT
2.Is endpoint remote?
2.1.eBPF SNAT
2.2.Redirect

10.100.1.1:60000 -> 192.168.0.1:31000

# NodePort with service endpoint on remote node in Cilium

# NodePort externalTrafficPolicy=Local



10.100.1.1:60000 -> 192.168.0.1:31000

# NodePort (DSR) in Cilium

# NodePort (DSR) in Cilium

# Performance (lower is better)



TCP_CRR to direct backend via NodePort latency (µseq per tx)

# Performance (lower is better)



TCP_RR to remote backend via NodePort latency (µseq per tx)

# Summary

**Performance**
- Better performance and latency over kube-proxy (ipvs and iptables)

**Reliability**
- Less LOC in datapath
- No need to wait for a new kernel release to fix a bug

**Debuggability**
- Better tooling for introspection and troubleshooting

**Compatibility**
- No more exec iptables

**Customization**
- Ability to change LB behaviour

# ClusterIP (host to pod)

TCP

UDP

**sendmsg()**

**recvmsg()**

```
import "net/http"

func main() {
 r, err := http.Get("nginx")
 ...
}
```

client

nginx

1.1.1.1

eth0

lxc0

3.3.3.3:80 (ClusterIP)

kernel

1. Lookup dst in SVC map
2. If found:
   a. Change dst addr and port in socket
   b. Create rev NAT entry

1. Lookup src in rev NAT map
2. If found:
   a. Change src addr and port