



# Speedrun through Splicing Sockets with SOCKMAP

Jakub Sitnicki  
Systems Engineer  
Cloudflare



Cilium +  
**eBPF Day**  
**EUROPE**

19 March 2024 | Paris, France

# \$ whoami

## Linux / OS Team @ Cloudflare

🌽 roll out fresh kernels

🐛 squash bugs





🧑🔧 troubleshoot stuff

✨ prototype features







# \$ whoami

## Linux / OS Team @ Cloudflare

-  roll out fresh kernels
-  squash bugs
-  troubleshoot stuff
-  prototype features

## SOCKMAP co-maintainer @ Linux upstream

-  small-time (= feature) maintainer
-  fix bugs
-  review patches
-  answer questions



# About this talk

## Good to know:

- ❑ network programming (`socket`, `connect`, `sendmsg`, `recvmsg`)
- ❑ basics of eBPF (what are BPF maps, programs, hooks, `bpftool`)
- ❑ building blocks of containers (cgroups, namespaces)

## Goals:

- ❑ know that SOCKMAP exists
- ❑ have idea how / when / what for use it
- ❑ feel ready to dive deeper

# Agenda

- 1 What can SOCKMAP do?
- 2 What is SOCKMAP?
- 3 How to set up SOCKMAP?
- 4 How to get sockets into a SOCKMAP?
- 5 Supported socket splicing setups
- 6 Real life use cases

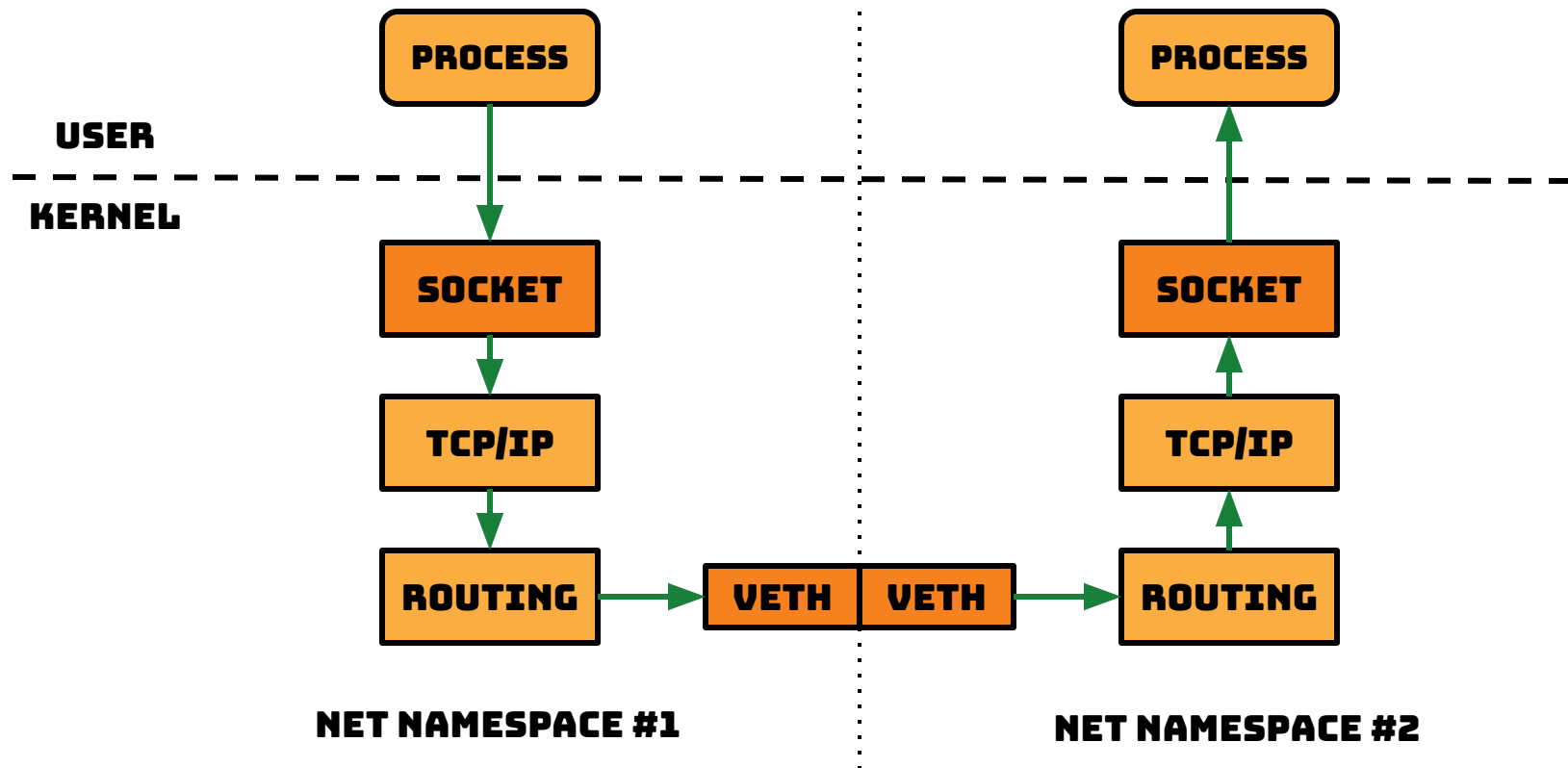


**REST AREA  
1 MILE**

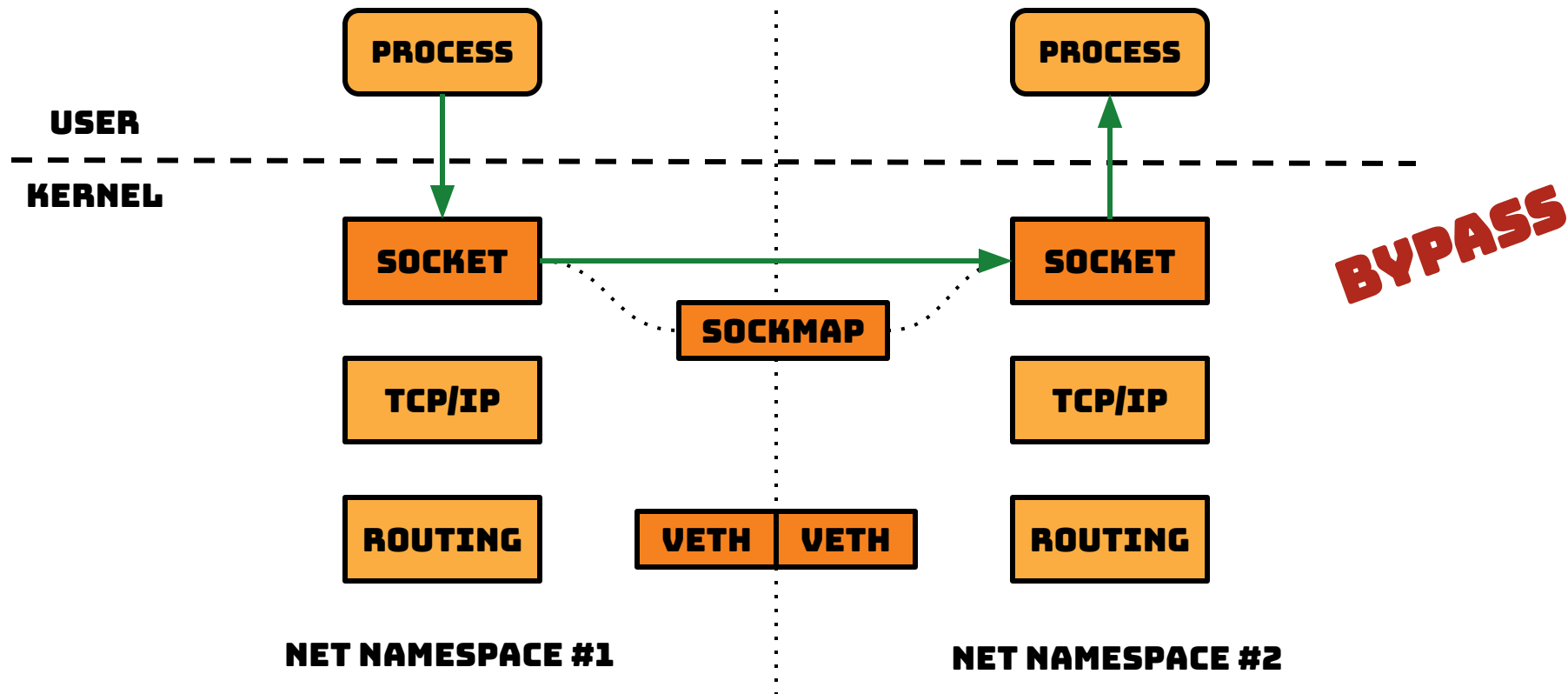
# What can SOCKMAP do for you?



# What can SOCKMAP do for... container networking



# What can SOCKMAP do for... container networking





# Let's set it up!

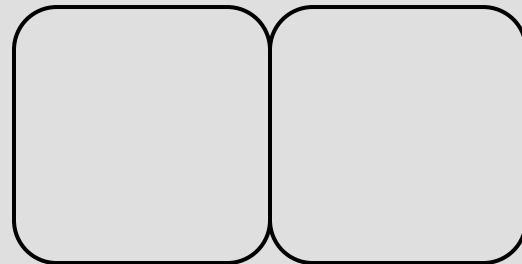
*Create two network namespaces*

```
# ip netns add A
```

```
# ip netns add B
```

**NETNS A**

**NETNS B**



# Let's set it up!

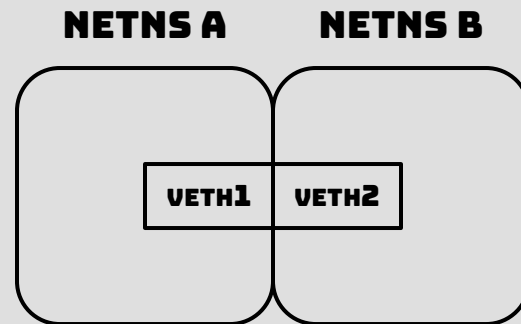
*Create two network namespaces*

```
# ip netns add A
```

```
# ip netns add B
```

*Link network namespaces with a veth pair*

```
# ip -n A link add name veth1 type veth \  
    peer name veth2 netns B
```



# Let's set it up!

*Create two network namespaces*

```
# ip netns add A
```

```
# ip netns add B
```

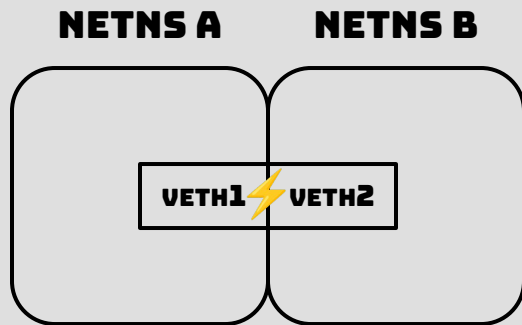
*Link network namespaces with a veth pair*

```
# ip -n A link add name veth1 type veth \  
    peer name veth2 netns B
```

*Bring up the links inside network namespaces*

```
# ip -n A link set dev veth1 up
```

```
# ip -n B link set dev veth2 up
```



# Let's set it up!

*Create two network namespaces*

```
# ip netns add A
```

```
# ip netns add B
```

*Link network namespaces with a veth pair*

```
# ip -n A link add name veth1 type veth \  
    peer name veth2 netns B
```

*Bring up the links inside network namespaces*

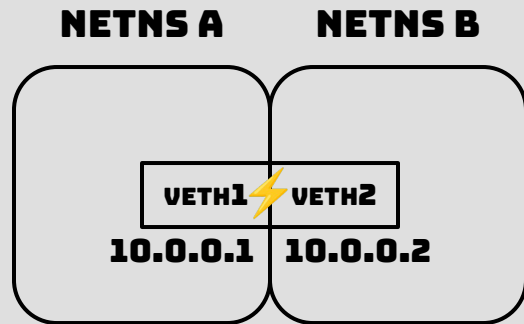
```
# ip -n A link set dev veth1 up
```

```
# ip -n B link set dev veth2 up
```

*Assign addresses to links inside network namespaces*

```
# ip -n A addr add 10.0.0.1/24 dev veth1
```

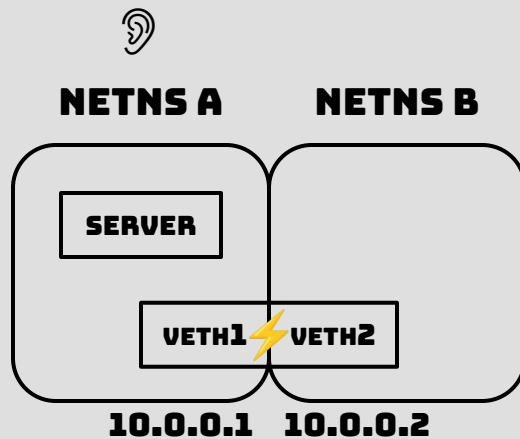
```
# ip -n B addr add 10.0.0.2/24 dev veth2
```



# Measure latency, no SOCKMAP first

*Run TCP server in netns A*

```
# ip netns exec A \  
  sockperf server -i 10.0.0.1 --tcp --daemonize
```



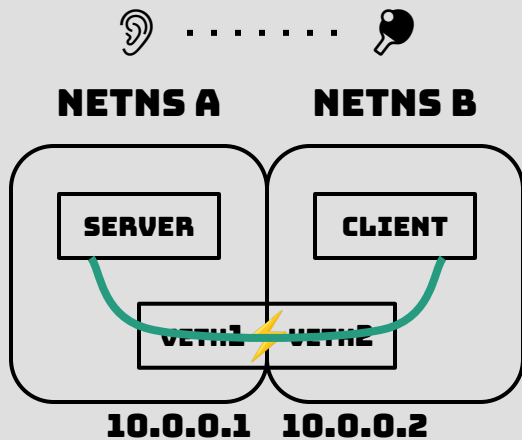
# Measure latency, no SOCKMAP first

*Run TCP server in netns A*

```
# ip netns exec A \  
  sockperf server -i 10.0.0.1 --tcp --daemonize
```

*Run TCP client in netns B*

```
# ip netns exec B \  
  sockperf ping-pong -i 10.0.0.1 --tcp --time 30
```



# Measure latency, no SOCKMAP first

*Run TCP server in netns A*

```
# ip netns exec A \  
  sockperf server -i 10.0.0.1 --tcp --daemonize
```

*Run TCP client in netns B*

```
# ip netns exec B \  
  sockperf ping-pong -i 10.0.0.1 --tcp --time 30
```

...

```
sockperf: [Total Run] RunTime=30.000 sec; Warm up time=400 msec; SentMessages=2599753;  
ReceivedMessages=2599752
```

...

```
sockperf: ==> avg-latency=5.748 (std-dev=2.010, mean-ad=0.322, median-ad=0.220,  
sigr=0.239, cv=0.350, std-error=0.001, 99.0% ci=[5.745, 5.751])  
sockperf: # dropped messages = 0; # duplicated messages = 0; # out-of-order messages = 0  
sockperf: Summary: Latency is 5.748 usec
```

$5.8 \pm 2.0 \mu\text{sec}$

# Set up SOCKMAP bypass

*Load BPF programs and create BPF maps*

```
# bpftool prog loadall redir_bypass.bpf.o /sys/fs/bpf \  
    pinmaps /sys/fs/bpf
```

**BPF  
PROGS**

**SK\_MSG**

**SK\_OPS**

**BPF MAPS**

**SOCKMAP**

**NETNS A**

**NETNS B**

**VETH1**

**VETH2**

**10.0.0.1**

**10.0.0.2**



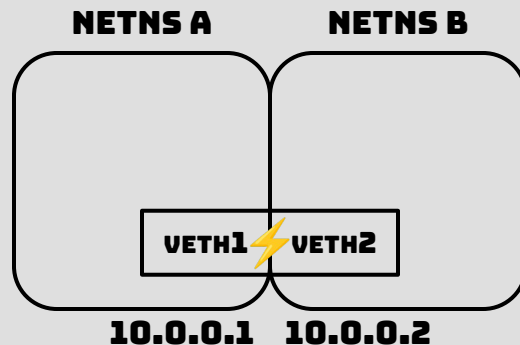
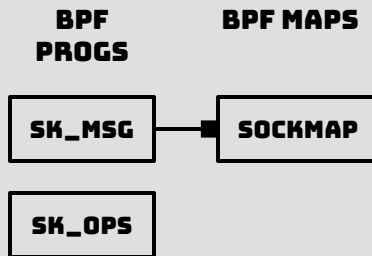
# Set up SOCKMAP bypass

*Load BPF programs and create BPF maps*

```
# bpftool prog loadall redir_bypass.bpf.o /sys/fs/bpf \  
    pinmaps /sys/fs/bpf
```

*Attach BPF program to BPF map*

```
# bpftool prog attach \  
    pinned /sys/fs/bpf/sk_msg_prog sk_msg_verdict \  
    pinned /sys/fs/bpf/sock_map
```



# Set up SOCKMAP bypass

*Load BPF programs and create BPF maps*

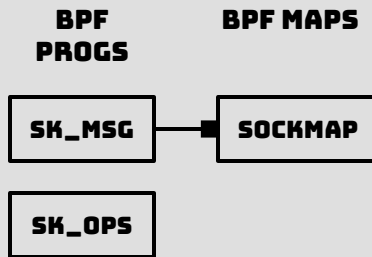
```
# bpftool prog loadall redir_bypass.bpf.o /sys/fs/bpf \  
    pinmaps /sys/fs/bpf
```

*Attach BPF program to BPF map*

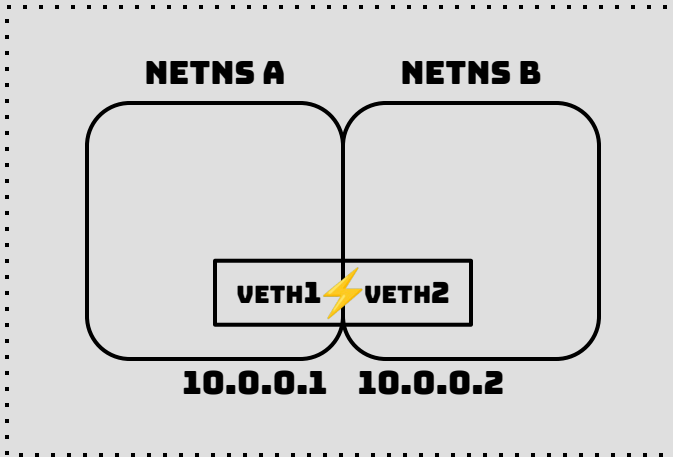
```
# bpftool prog attach \  
    pinned /sys/fs/bpf/sk_msg_prog sk_msg_verdict \  
    pinned /sys/fs/bpf/sock_map
```

*Create a test cgroup*

```
# mkdir /sys/fs/cgroup/unified/test.slice
```



## TEST.SLICE CGROUP



# Set up SOCKMAP bypass

*Load BPF programs and create BPF maps*

```
# bpftool prog loadall redir_bypass.bpf.o /sys/fs/bpf \  
    pinmaps /sys/fs/bpf
```

*Attach BPF program to BPF map*

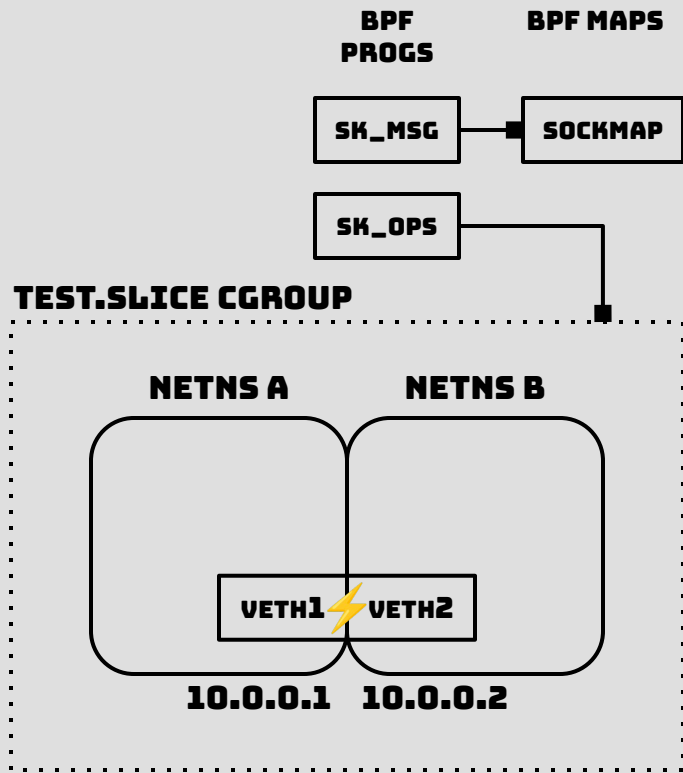
```
# bpftool prog attach \  
    pinned /sys/fs/bpf/sk_msg_prog sk_msg_verdict \  
    pinned /sys/fs/bpf/sock_map
```

*Create a test cgroup*

```
# mkdir /sys/fs/cgroup/unified/test.slice
```

*Attach BPF program to cgroup*

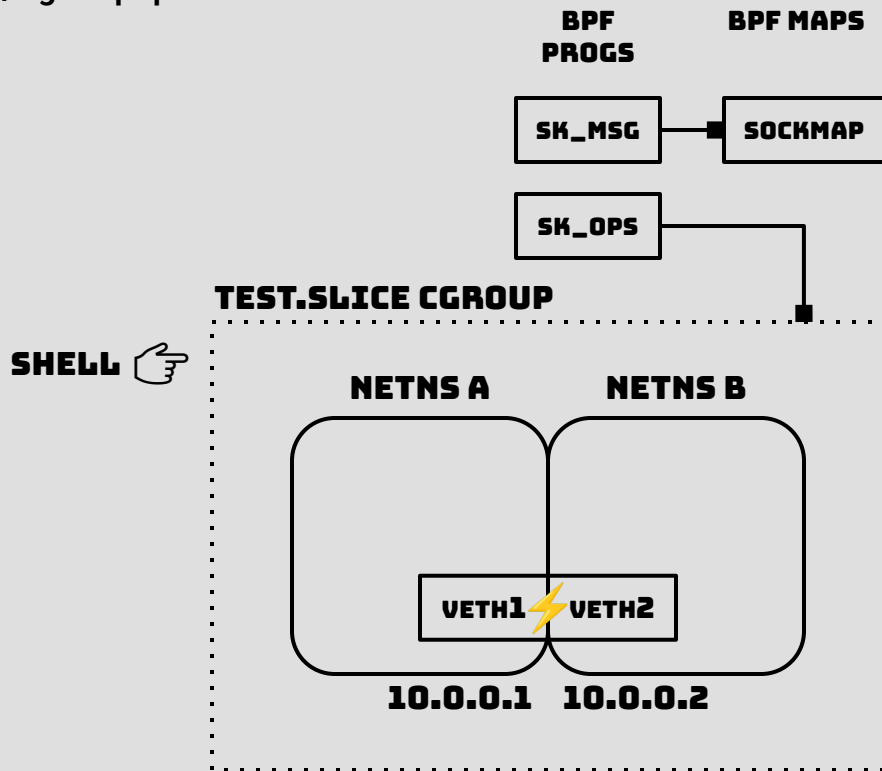
```
# bpftool cgroup attach \  
    /sys/fs/cgroup/unified/test.slice \  
    cgroup_sock_ops pinned /sys/fs/bpf/sockops_prog
```



# Repeat test with SOCKMAP bypass

*Spawn client and server inside the test cgroup*

```
# echo $$ > /sys/fs/cgroup/unified/test.slice/cgroup.procs
```



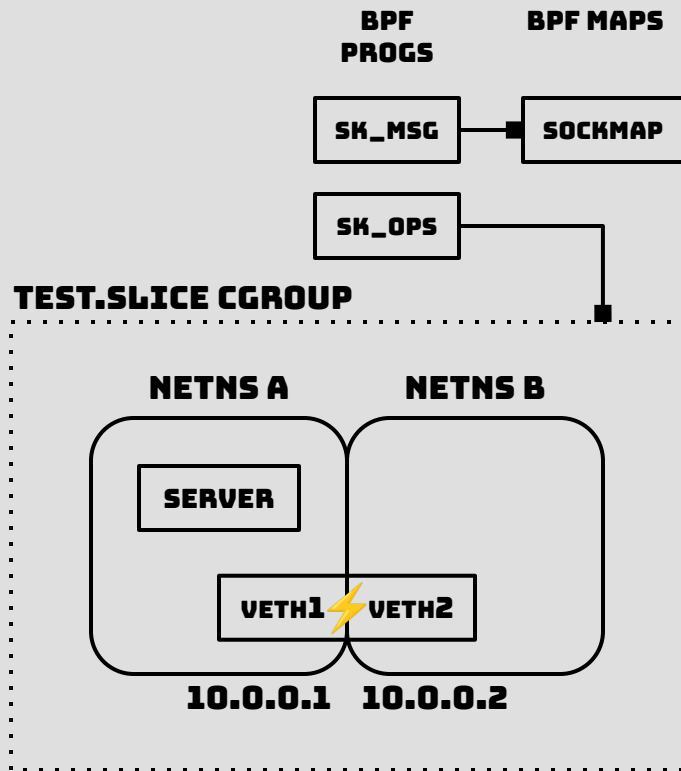
# Repeat test with SOCKMAP bypass

*Spawn client and server inside the test cgroup*

```
# echo $$ > /sys/fs/cgroup/unified/test.slice/cgroup.procs
```

*Run TCP server in netns A*

```
# ip netns exec A \  
  sockperf server -i 10.0.0.1 --tcp --daemonize
```



# Repeat test with SOCKMAP bypass

*Spawn client and server inside the test cgroup*

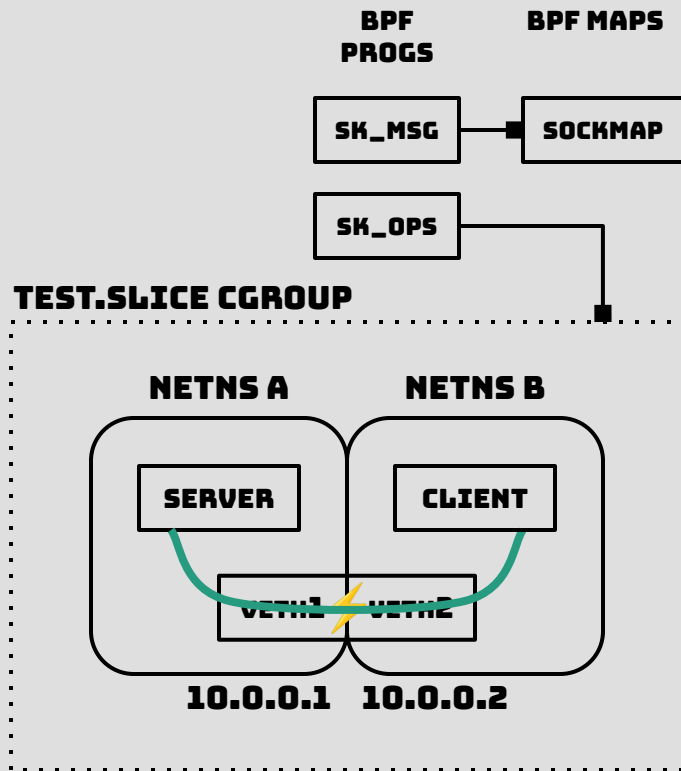
```
# echo $$ > /sys/fs/cgroup/unified/test.slice/cgroup.procs
```

*Run TCP server in netns A*

```
# ip netns exec A \  
    sockperf server -i 10.0.0.1 --tcp --daemonize
```

*Run TCP client*

```
# ip netns exec B \  
    sockperf ping-pong -i 10.0.0.1 --tcp --time 30
```



# Repeat test with SOCKMAP bypass

*Spawn client and server inside the test cgroup*

```
# echo $$ > /sys/fs/cgroup/unified/test.slice/cgroup.procs
```

*Run TCP server in netns A*

```
# ip netns exec A \  
    sockperf server -i 10.0.0.1 --tcp --daemonize
```

*Run TCP client*

```
# ip netns exec B \  
    sockperf ping-pong -i 10.0.0.1 --tcp --time 30
```

```
sockperf: [Total Run] RunTime=30.000 sec; Warm up time=400 msec; SentMessages=3189584;  
ReceivedMessages=3189583
```

...

```
sockperf: ==> avg-latency=4.686 (std-dev=2.862, mean-ad=0.250, median-ad=0.216,  
siqr=0.173, cv=0.611, std-error=0.002, 99.0% ci=[4.682, 4.690])
```

```
sockperf: # dropped messages = 0; # duplicated messages = 0; # out-of-order messages = 0
```

```
sockperf: Summary: Latency is 4.686 usec
```

$4.7 \pm 2.9 \mu\text{sec}$

# Without and with SOCKMAP bypass

before:  $5.8 \pm 2.0 \mu\text{sec}$

↓ - 18.5%

after:  $4.7 \pm 2.9 \mu\text{sec}$

Run the benchmark yourself:

<https://github.com/jsitnicki/kubecon-2024-sockmap/tree/main/examples/redirect-bypass>



# What is SOCKMAP?



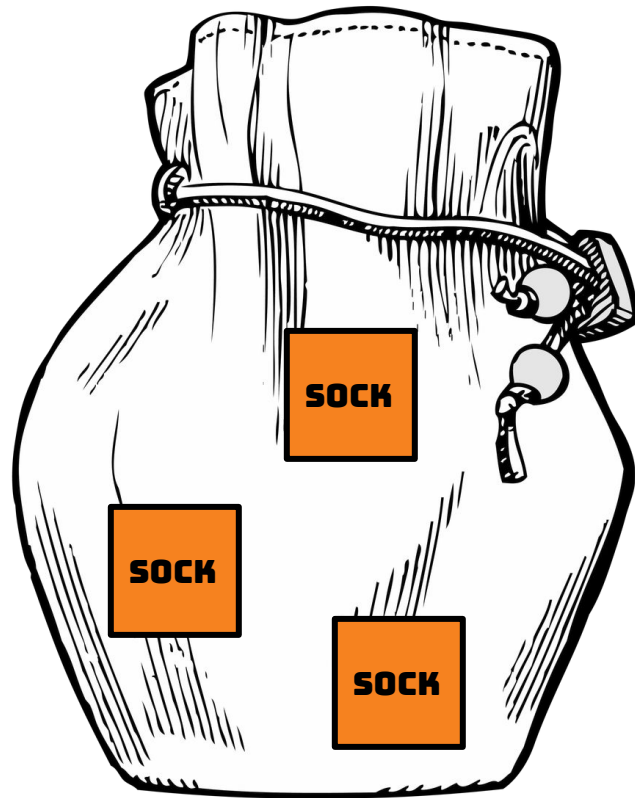
# What is SOCKMAP?



Two things

# What is SOCKMAP?

**Collection / container  
for socket references  
in Linux kernel**



# What is SOCKMAP?

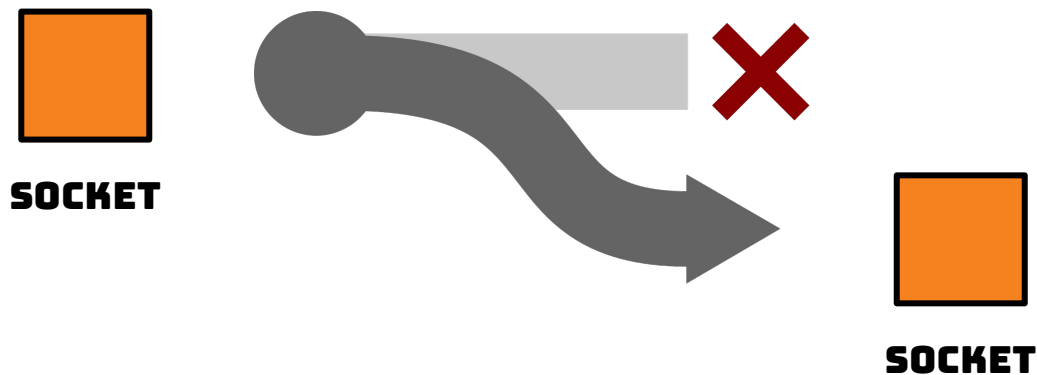
SOCKMAP  
API



## 1. container for sockets

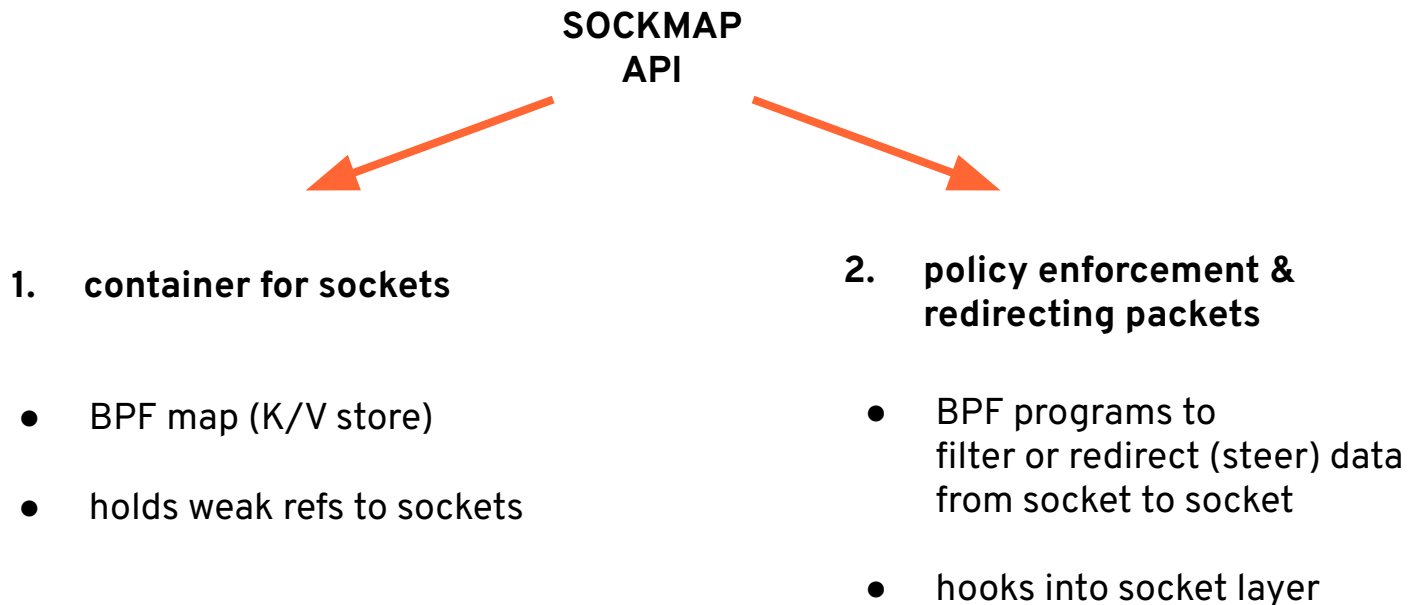
- BPF map (K/V store)
- holds weak refs to sockets

# What is SOCKMAP?



**API for enforcing policy  
and  
redirecting data between sockets**

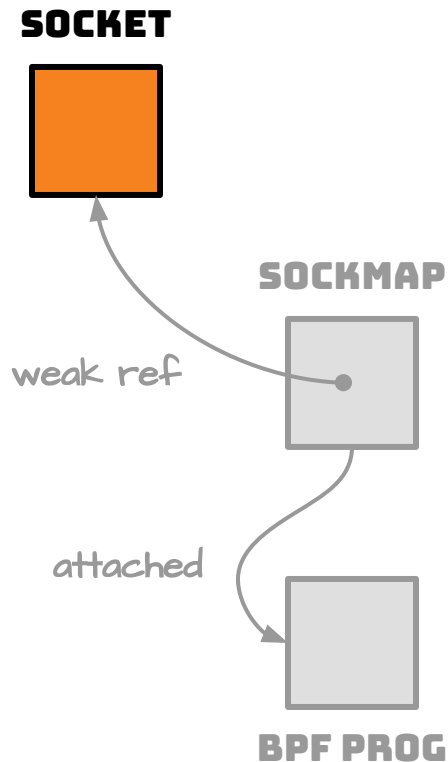
# What is SOCKMAP?



# How to set up SOCKMAP?



# ① Open a connected (established) socket



## active open

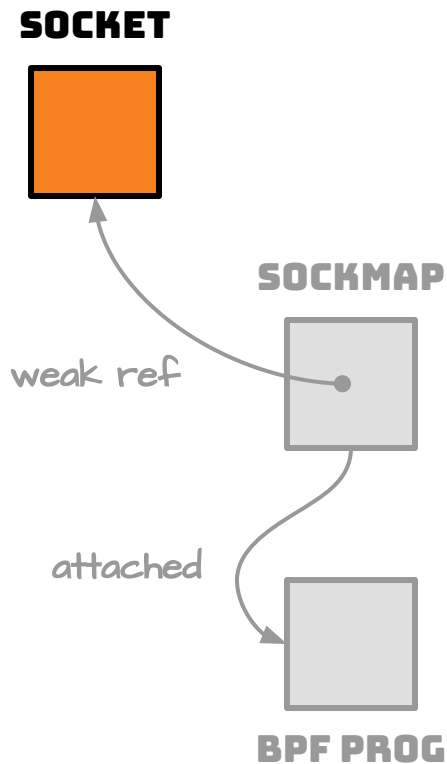
```
socket(AF_INET, SOCK_STREAM, IPPROTO_IP) = 8
connect(8, {sa_family=AF_INET,
          sin_port=htons(41895),
          sin_addr=inet_addr("127.0.0.1")}, 16) = 0
```

## passive open

```
socket(AF_INET, SOCK_STREAM, IPPROTO_IP) = 7
bind(7, {sa_family=AF_INET,
        sin_port=htons(41895),
        sin_addr=inet_addr("127.0.0.1")}, 16) = 0
listen(7, 4096) = 0
accept(7, NULL, NULL) = 9
```



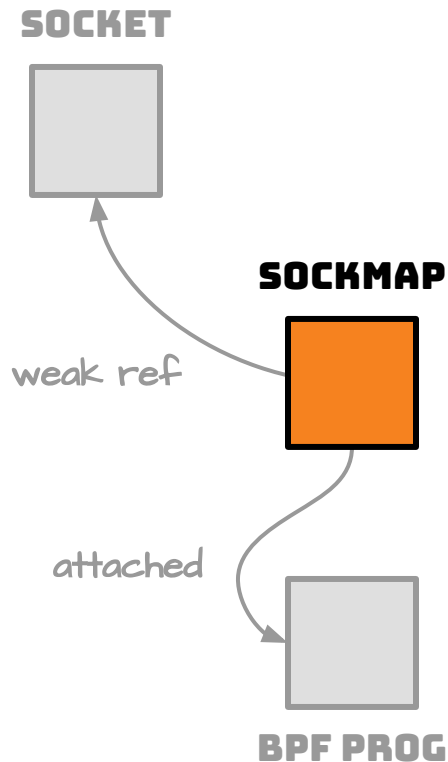
# What sockets can you use?



connected (established) socket:

- ❑ TCP
- ❑ UDP
- ❑ **UNIX** (STREAM, DGRAM)
- ❑ **VSOCK** (STREAM, SEQPACKET)

## ② Create a BPF map - SOCKMAP or SOCKHASH

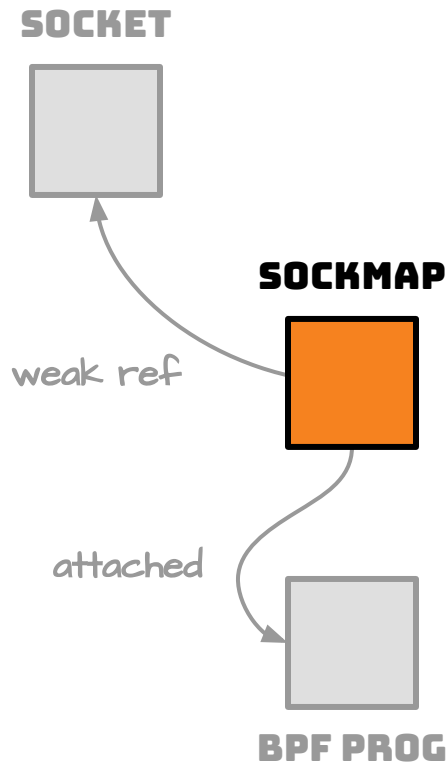


use the `bpf()` syscall

or a library wrapper (`ebpf-go`, `libbpf`)

```
bpf(BPF_MAP_CREATE, {map_type=BPF_MAP_TYPE_SOCKMAP,  
                      key_size=4,  
                      value_size=8,  
                      max_entries=1,  
                      map_flags=0,  
                      ...}, 72) = 5
```

## ② Create a BPF map

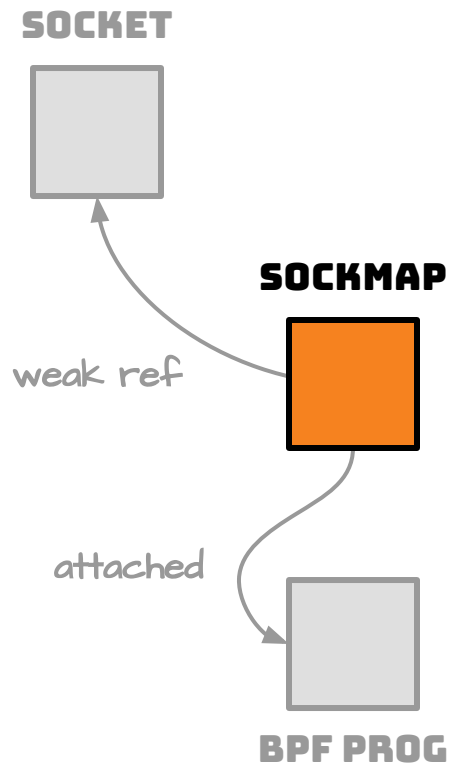


or use **bpftool map create** command

```
bpftool map create \
  /sys/fs/bpf/sockmap `# path on bpffs` \
  type sockmap `# sockmap or sockhash` \
  key 4 `# always 4 bytes for sockmap` \
  value 8 `# use 8 bytes for dump to work` \
  entries 1 \
  name sockmap
```

```
bpftool map show pinned /sys/fs/bpf/sockmap
3: sockmap name sockmap flags 0x0
    key 4B value 8B max_entries 1 memlock 328B
```

# What BPF maps can you use?



## Map types:

- ❑ `BPF_MAP_TYPE_SOCKMAP`
  - ❑ **32-bit integer key**
- ❑ `BPF_MAP_TYPE_SOCKHASH`
  - ❑ **binary blob key**

Not to be confused with  
`BPF_MAP_TYPE_REUSEPORT_SOCKARRAY`

### ③ Load a BPF program

**SOCKET**

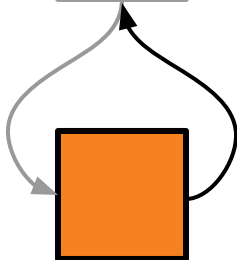


**SOCKMAP**



*weak ref*

*attached*



*uses*

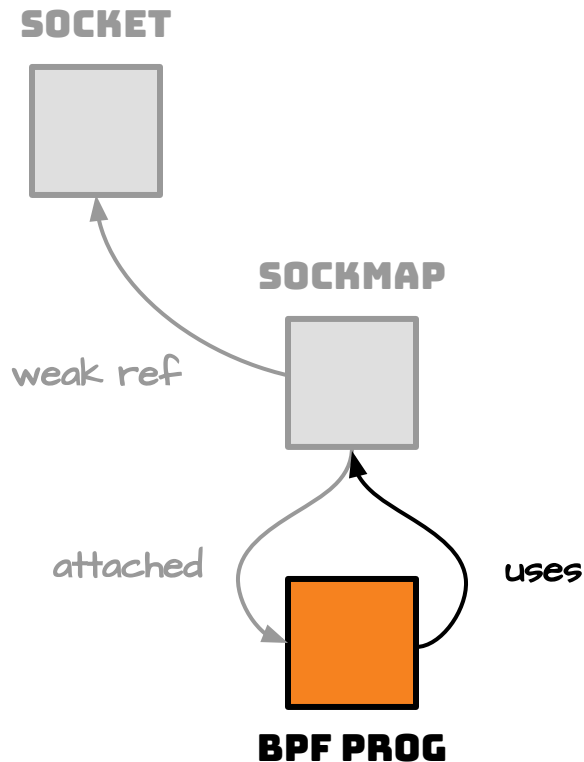
**BPF PROG**

```
bpf(BPF_PROG_LOAD, {prog_type=BPF_PROG_TYPE_SK_MSG,  
                    insn_cnt=6,  
                    insns=0xcf2a70,  
                    expected_attach_type=BPF_SK_MSG_VERDICT,  
                    ...}, 128) = 6
```

**Program types:**

- ❑ BPF\_PROG\_TYPE\_SK\_MSG
- ❑ BPF\_PROG\_TYPE\_SK\_SKB

### ③ Load a BPF program - it uses SOCKMAP



```
# bpftool prog dump xlated id 42
int prog_msg_redir_ingress(struct sk_msg_md * msg):
    0: (18) r2 = map[id:17]
    ...
    5: (95) exit
# bpftool map show id 17
17: sockmap name output flags 0x0
    key 4B value 8B max_entries 1 memlock 328B
    pids sockmap-redir-m(331)
```

## ④ Attach BPF program to SOCKMAP

**SOCKET**

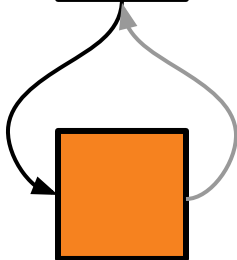


**SOCKMAP**



*weak ref*

*attached*



**BPF PROG**

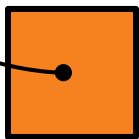
```
bpf(BPF_PROG_ATTACH, {target_fd=5,  
    attach_bpf_fd=6,  
    attach_type=BPF_SK_MSG_VERDICT,  
    attach_flags=0,  
    replace_bpf_fd=0}, 20) = 0
```

## ⑤ Insert socket into SOCKMAP

**SOCKET**

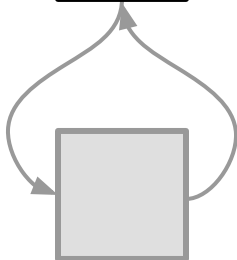


**SOCKMAP**



weak ref

attached



**BPF PROG**

uses

```
bpf(BPF_MAP_UPDATE_ELEM, {map_fd=5,  
    key=0x7ffeb3803870,  
    value=0x7ffeb3803868,  
    flags=BPF_NOEXIST}, 32) = 0
```

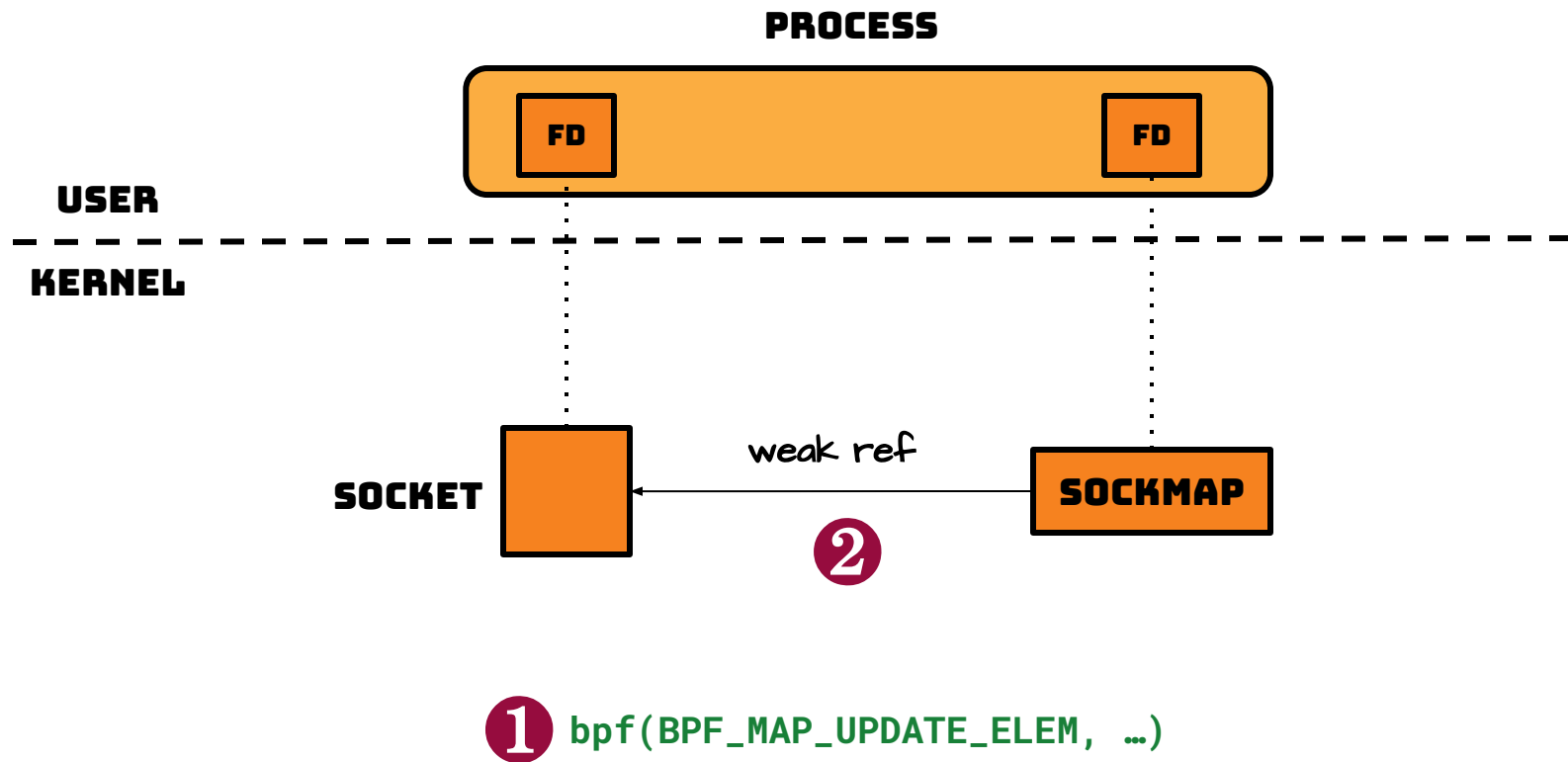
⚠ must be done after attaching the program



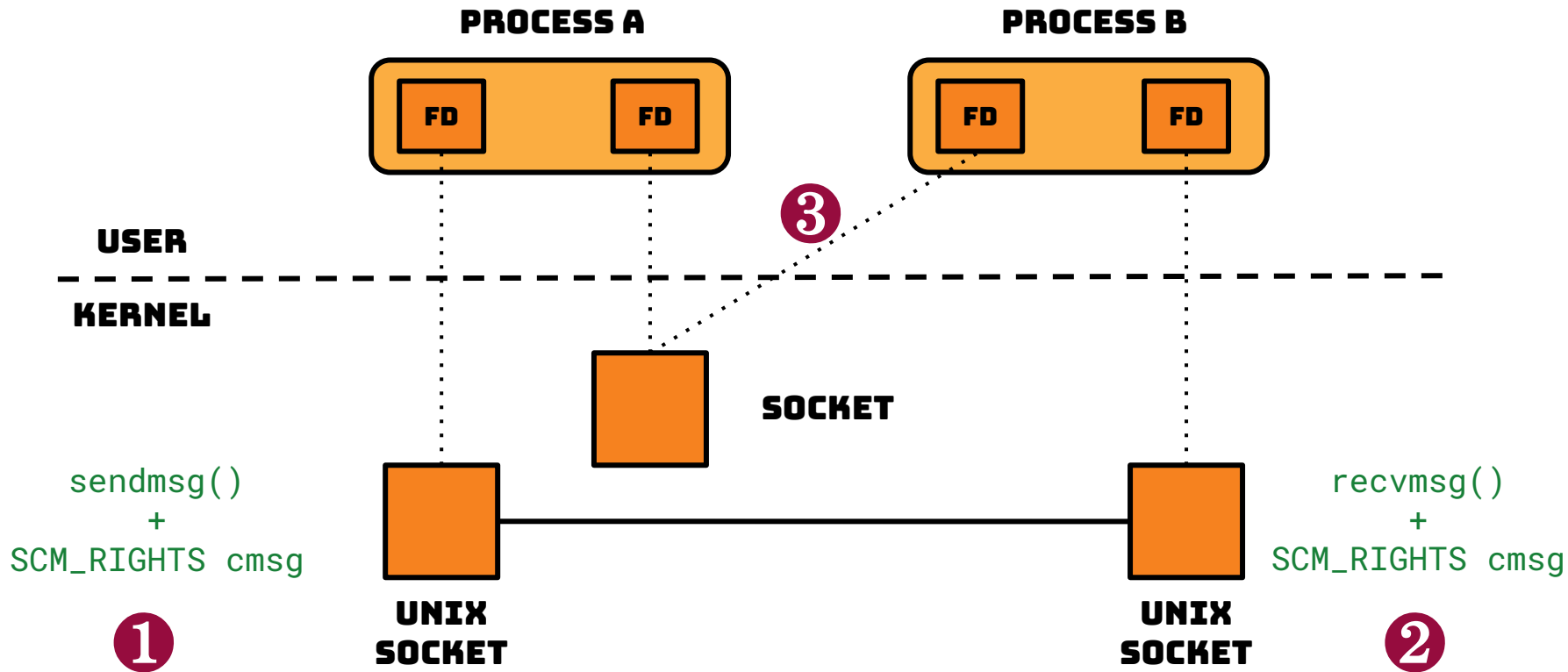
# How to get sockets into a SOCKMAP?



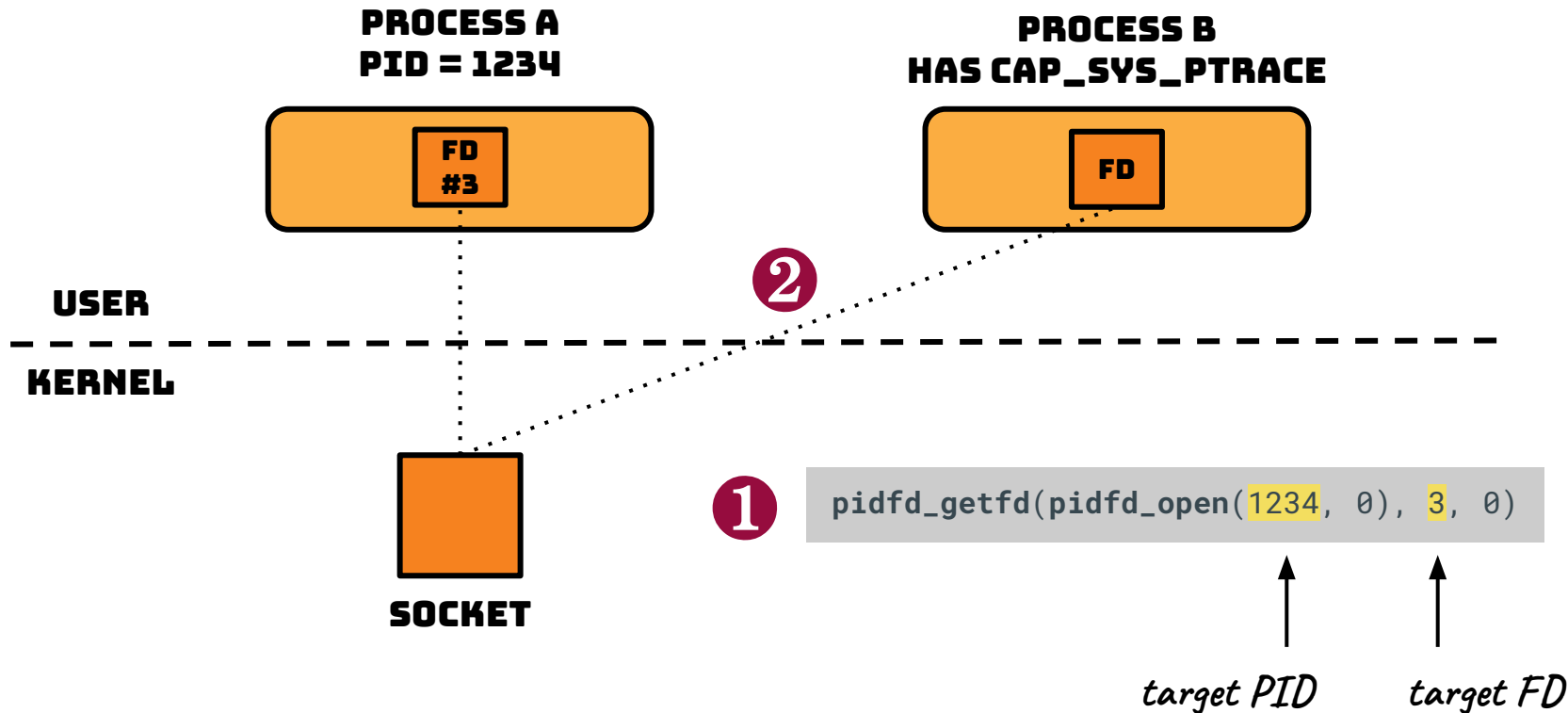
# Easy case - Single process



# Socket FD handover with SCM\_RIGHTS



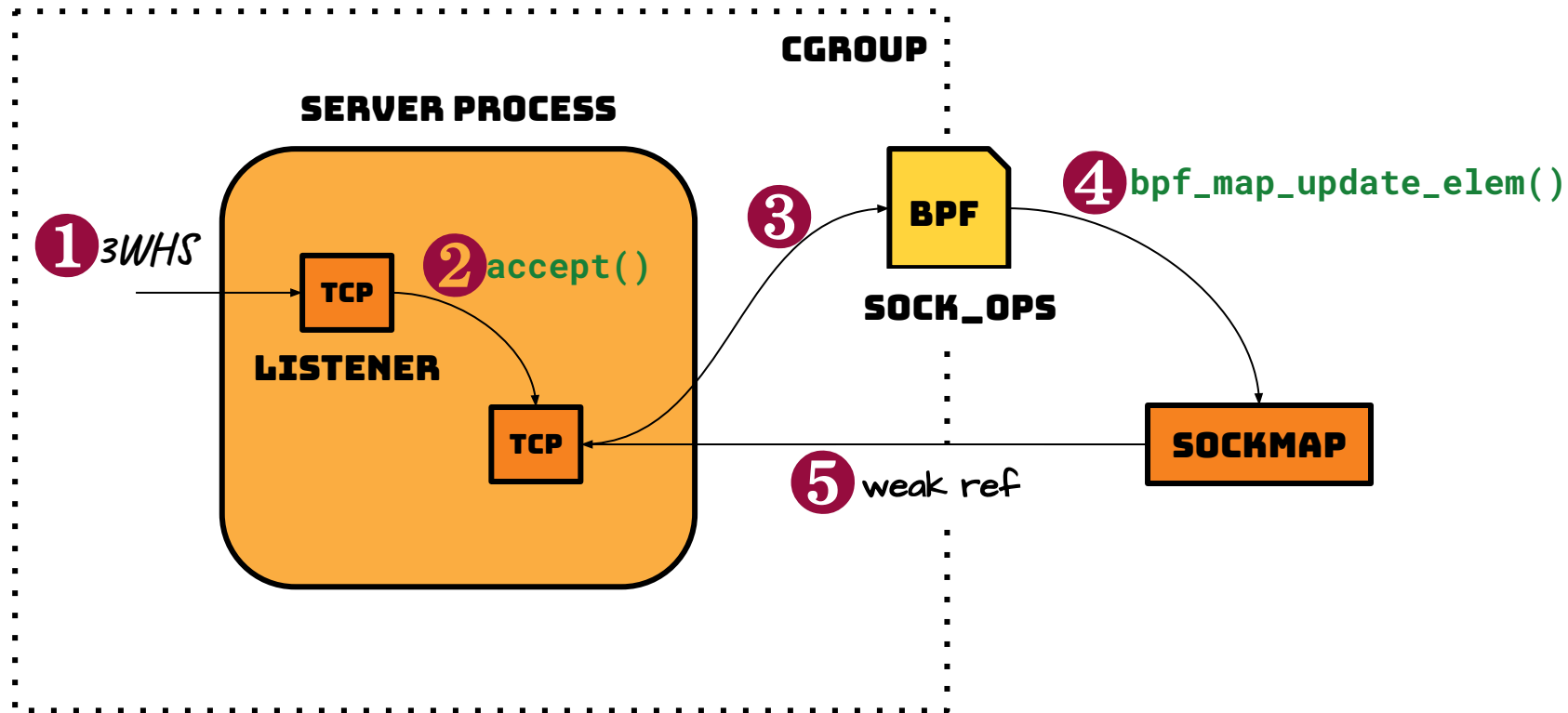
# “Steal” a socket FD



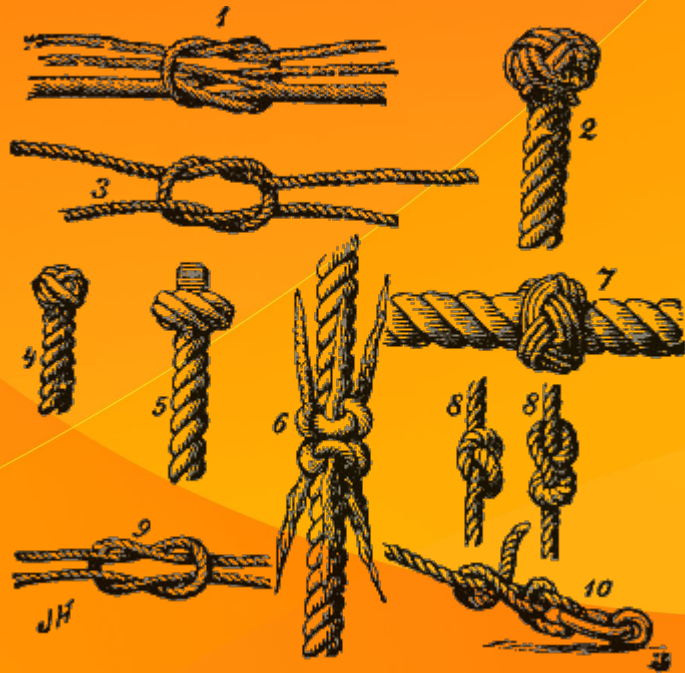
Linux 5.6+, requires `CAP_SYS_PTRACE`

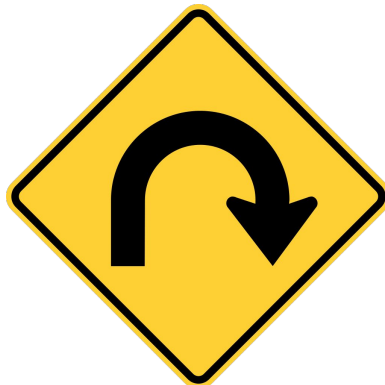
[https://manpages.debian.org/unstable/manpages-dev/pidfd\\_getfd.2.en.html](https://manpages.debian.org/unstable/manpages-dev/pidfd_getfd.2.en.html)

# BPF `sock_ops` program attached to cgroup (TCP only)

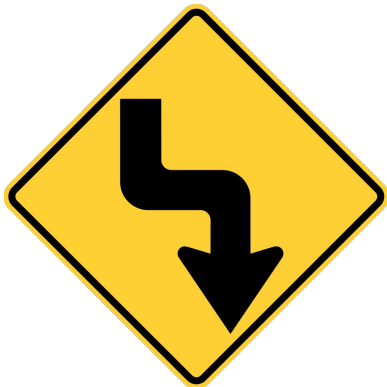


# Supported Socket Splicing Setups





**Redirect**



# Redirect

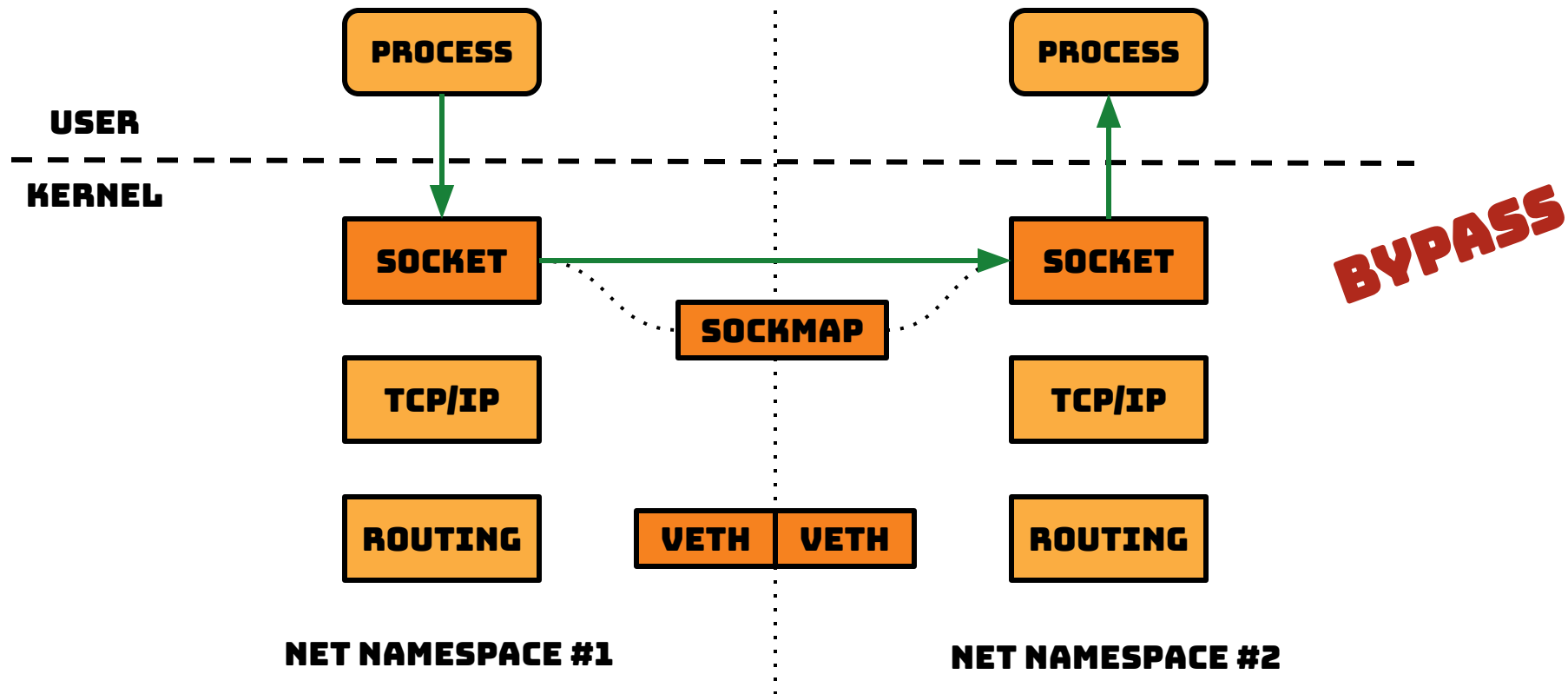


**send to local**



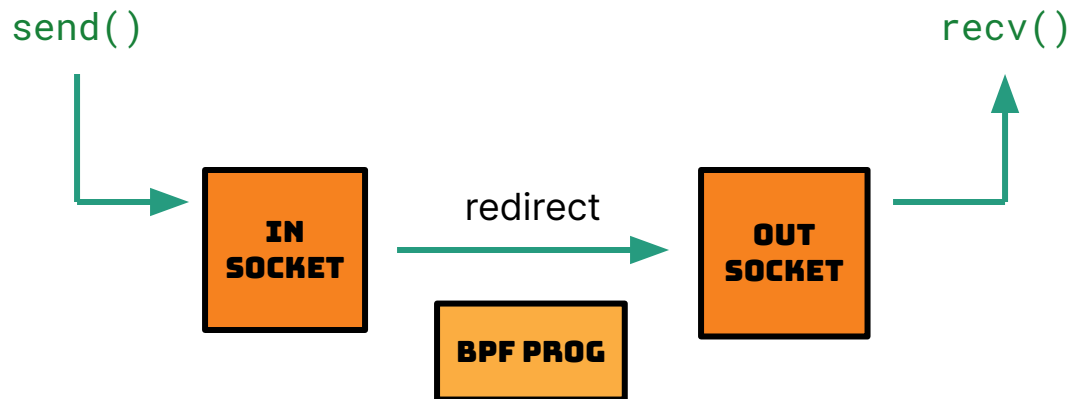


# Redirect use case → Bypass for containers





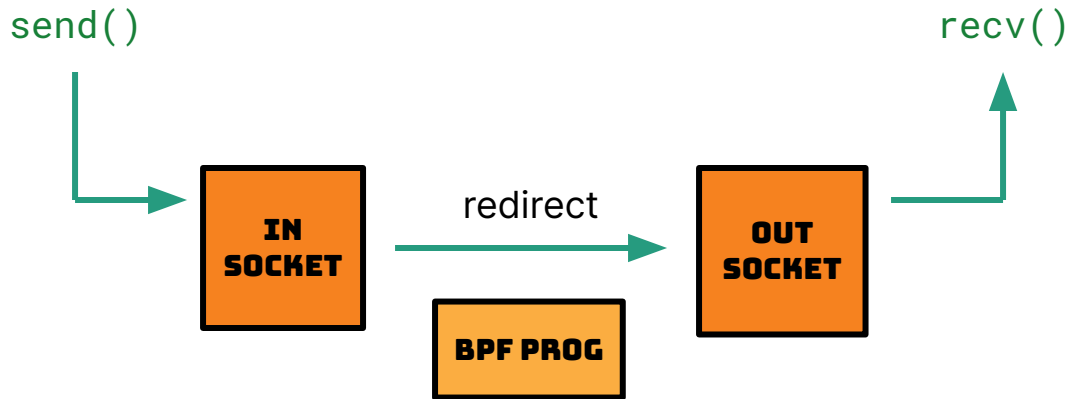
## Redirect → send to local



Like `socketpair()` or `pipe()`



# Redirect → send to local → How?



**BPF\_PROG\_TYPE\_SK\_MSG** program

→ attached to **BPF\_SK\_MSG\_VERDICT** hook

→ calls **bpf\_msg\_redirect\_hash/map()** with **BPF\_F\_INGRESS** flag

→ returns **SK\_PASS**

*selects target socket*



## Redirect → send to local → Example

```
SEC("sk_msg")
int sk_msg_redir_ingress(struct sk_msg_md *msg)
{
    __u32 key = 0;

    if (msg->remote_port == bpf_htonl(53))
        key = 1;

    return bpf_msg_redirect_map(msg, &sockmap, key, BPF_F_INGRESS);
}
```



## Redirect → send to local → What?

IN → OUT	TCP	UDP	UNIX STR	UNIX DGR	VSOCK STR	VSOCK SEQ
TCP	●	●	●	●	●	●
UDP	●	●	●	●	●	●
UNIX STR	●	●	●	●	●	●
UNIX DGR	●	●	●	●	●	●
VSOCK STR	●	●	●	●	●	●
VSOCK SEQ	●	●	●	●	●	●

*TCP to any but VSOCK*

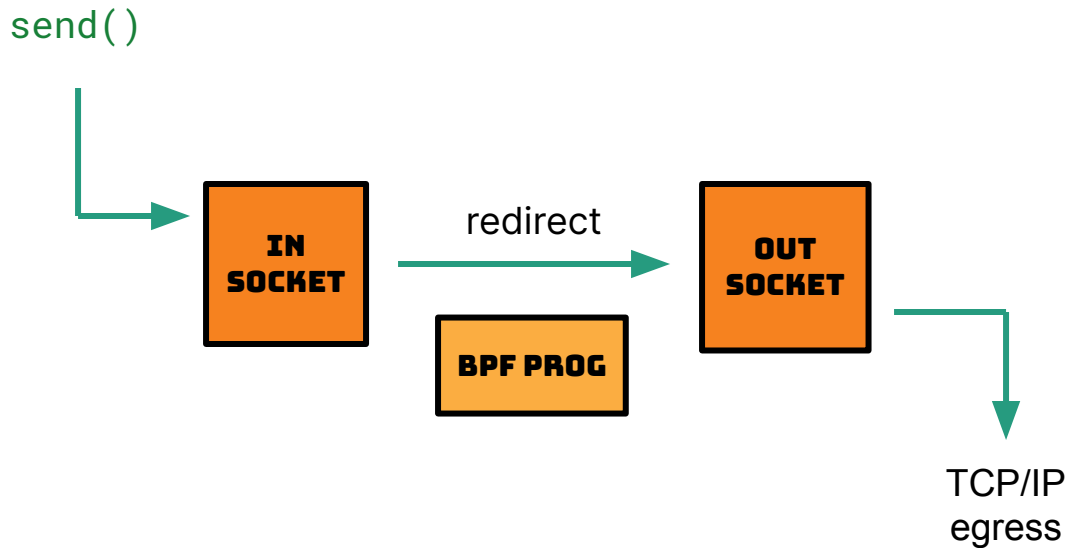
# Redirect



**send to egress**



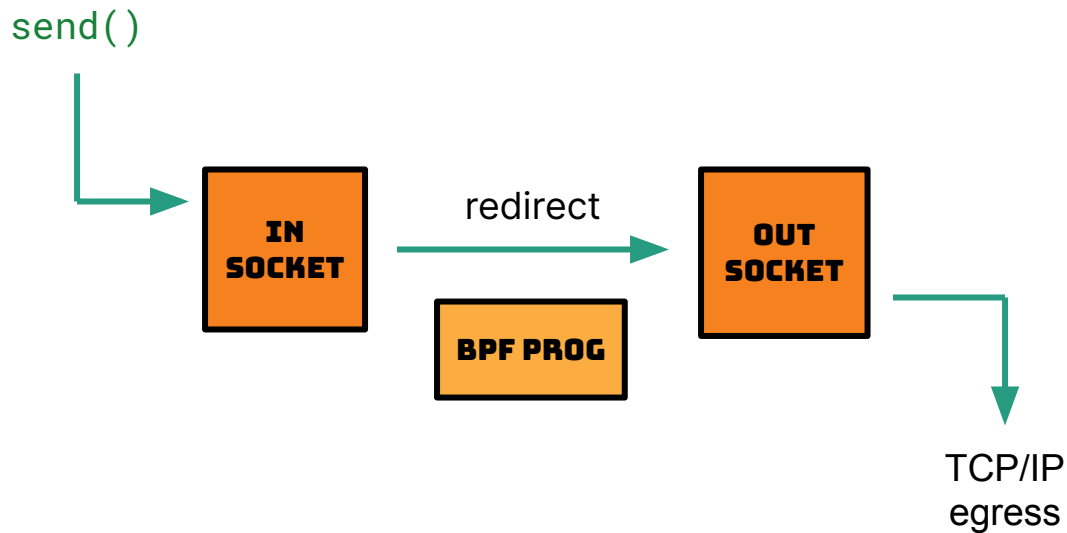
## Redirect → send to egress



Sort of like `splice()`  
(pipe → socket)



# Redirect → send to egress → How?



`BPF_PROG_TYPE_SK_MSG` prog

→ attached to `BPF_SK_MSG_VERDICT` hook

→ calls `bpf_msg_redirect_hash/map()` without any flags

→ returns `SK_PASS`



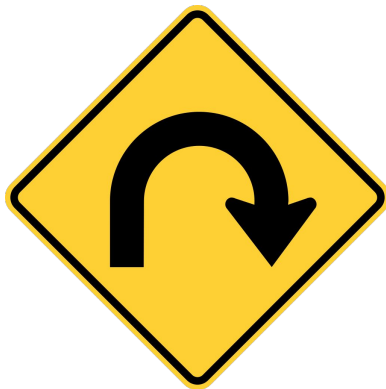


## Redirect → send to egress → What?

IN → OUT	TCP	UDP	UNIX STR	UNIX DGR	VSOCK STR	VSOCK SEQ
TCP	●	●	●	●	●	●
UDP	●	●	●	●	●	●
UNIX STR	●	●	●	●	●	●
UNIX DGR	●	●	●	●	●	●
VSOCK STR	●	●	●	●	●	●
VSOCK SEQ	●	●	●	●	●	●

*TCP to TCP only*

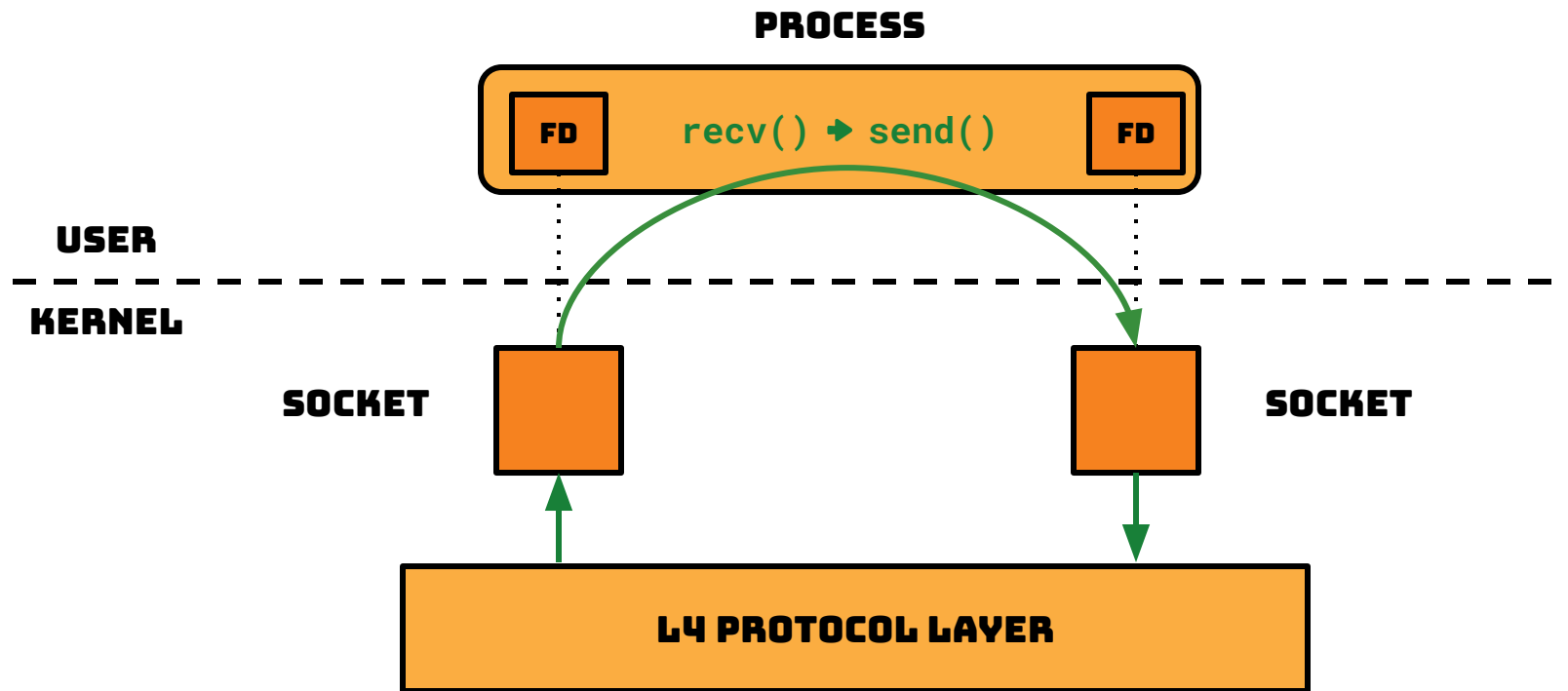
# Redirect



**ingress to egress**



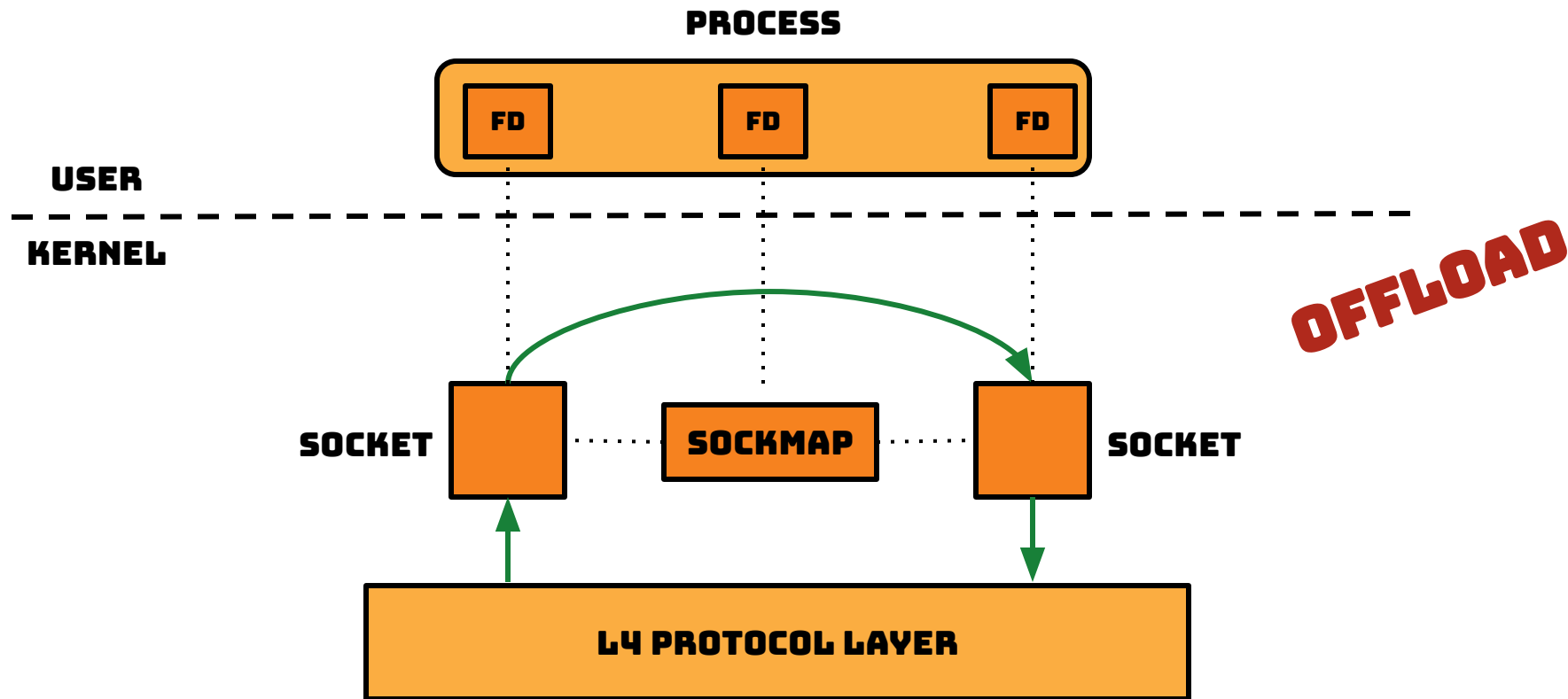
## Redirect use case → L7 network proxy



Examples: `socat`, `systemd-socket-proxyd`

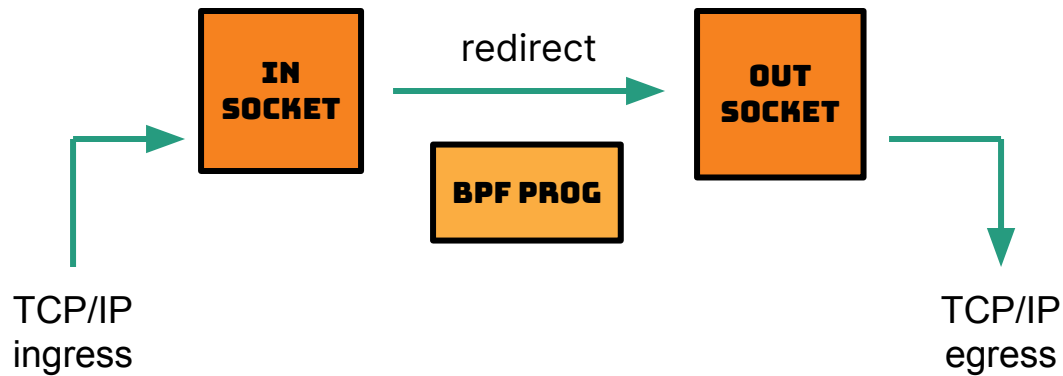


## Redirect use case → L7 network proxy





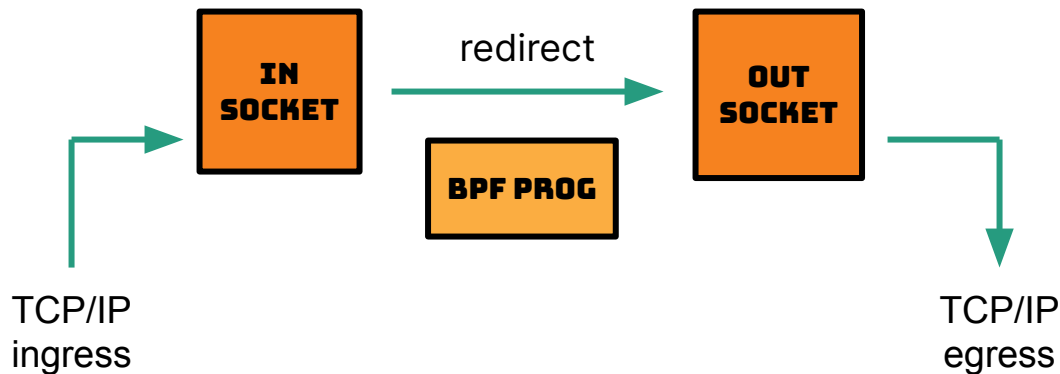
## Redirect → ingress to egress



Like double `splice()`  
(socket → pipe → socket)



# Redirect → ingress to egress → How?



`BPF_PROG_TYPE_SK_SKB` prog

→ attached to `BPF_SK_SKB_VERDICT` hook

→ calls `bpf_sk_redirect_hash/map()` without any flags

→ returns `SK_PASS`



# Redirect → ingress to egress → What?

IN → OUT	TCP	UDP	UNIX STR	UNIX DGR	VSOCK STR	VSOCK SEQ
TCP	●	●	●	●	●	●
UDP	●	●	●	●	●	●
UNIX STR	●	●	●	●	●	●
UNIX DGR	●	●	●	●	●	●
VSOCK STR	●	●	●	●	●	●
VSOCK SEQ	●	●	●	●	●	●

*any to any*

# Redirect

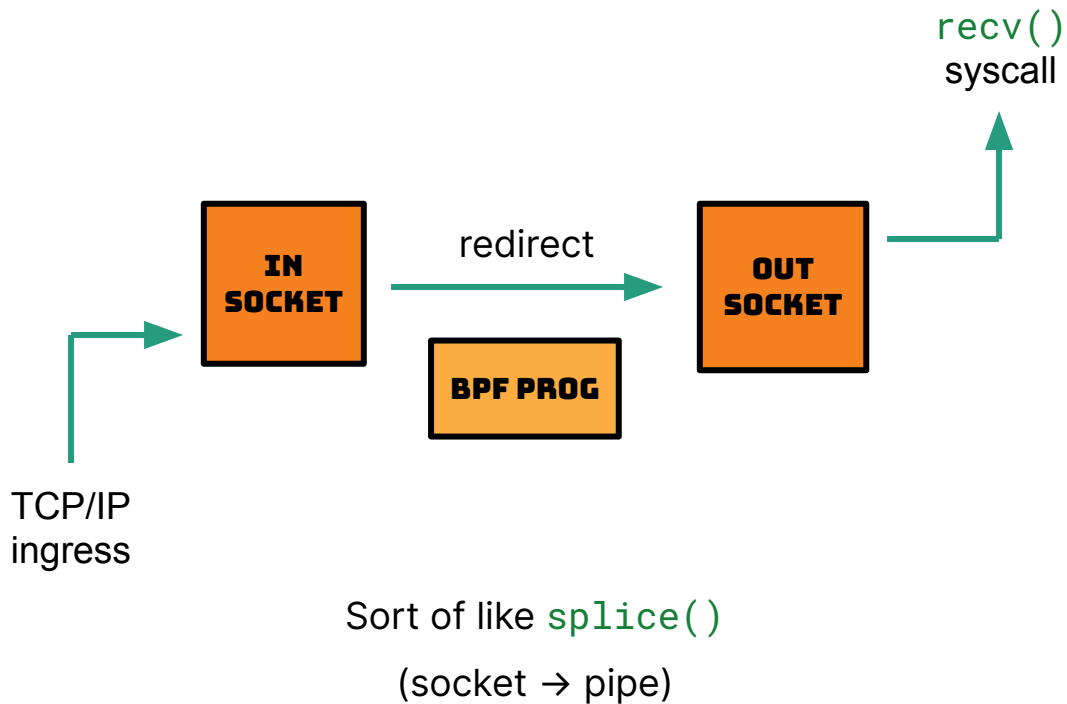


**ingress to local**



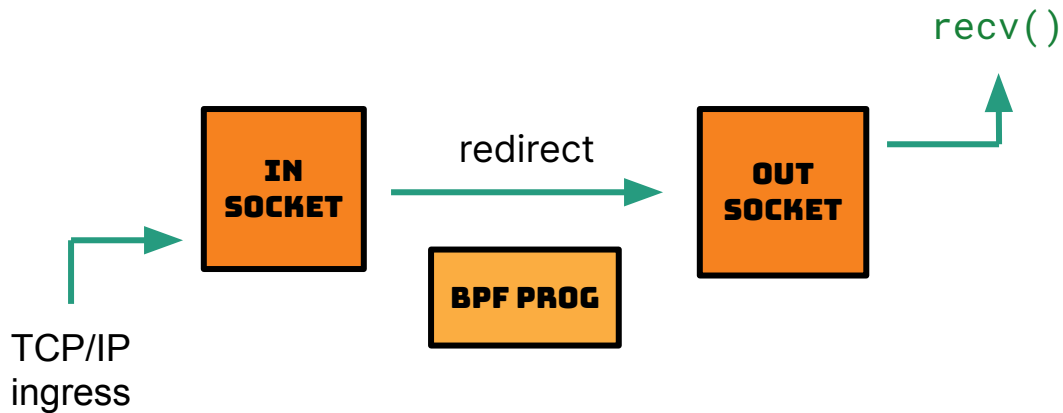


## Redirect → ingress to local





# Redirect → ingress to local → How?



`BPF_PROG_TYPE_SK_SKB` prog

→ attached to `BPF_SK_SKB_VERDICT` hook

→ calls `bpf_sk_redirect_hash/map()` with `BPF_F_INGRESS` flag

→ returns `SK_PASS`







## Redirect → ingress to local → What?

IN → OUT	TCP	UDP	UNIX STR	UNIX DGR	VSOCK STR	VSOCK SEQ
TCP	●	●	●	●	●	●
UDP	●	●	●	●	●	●
UNIX STR	●	●	●	●	●	●
UNIX DGR	●	●	●	●	●	●
VSOCK STR	●	●	●	●	●	●
VSOCK SEQ	●	●	●	●	●	●

*any to any but VSOCK*

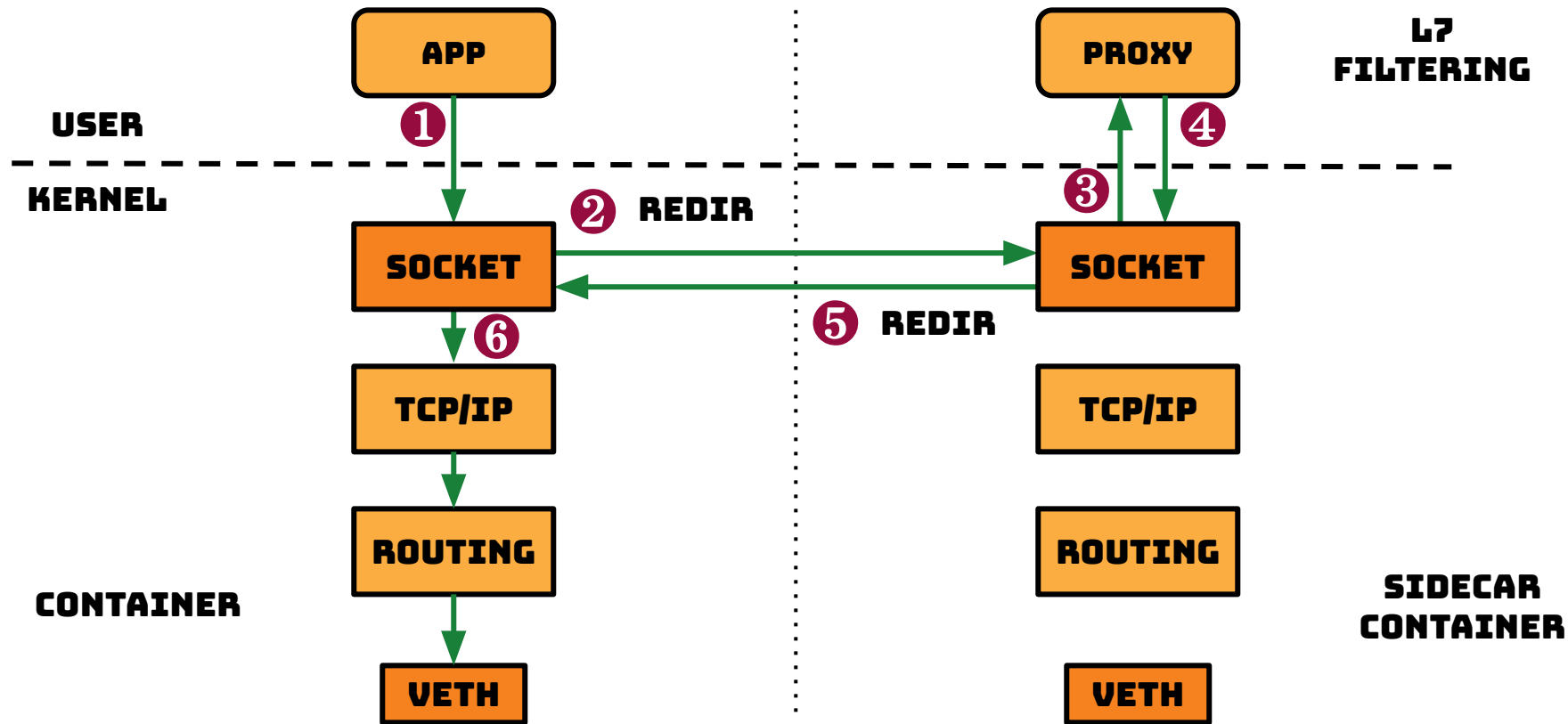
# Cheatsheet - Redirect with SOCKMAP

redirect scenario	program type BPF_PROG_TYPE_*	attach type BPF_*	redirect function	redirect flag	in socket type	out socket type
send to local 	SK_MSG	SK_MSG_VERDICT	bpf_msg_redirect_*	BPF_F_INGRESS	TCP	any but VSOCK
send to egress 	SK_MSG	SK_MSG_VERDICT	bpf_msg_redirect_*	none	TCP	TCP
ingress to egress 	SK_SKB	SK_SKB_VERDICT	bpf_sk_redirect_*	none	any	any
ingress to local 	SK_SKB	SK_SKB_VERDICT	bpf_sk_redirect_*	BPF_F_INGRESS	any	any but VSOCK

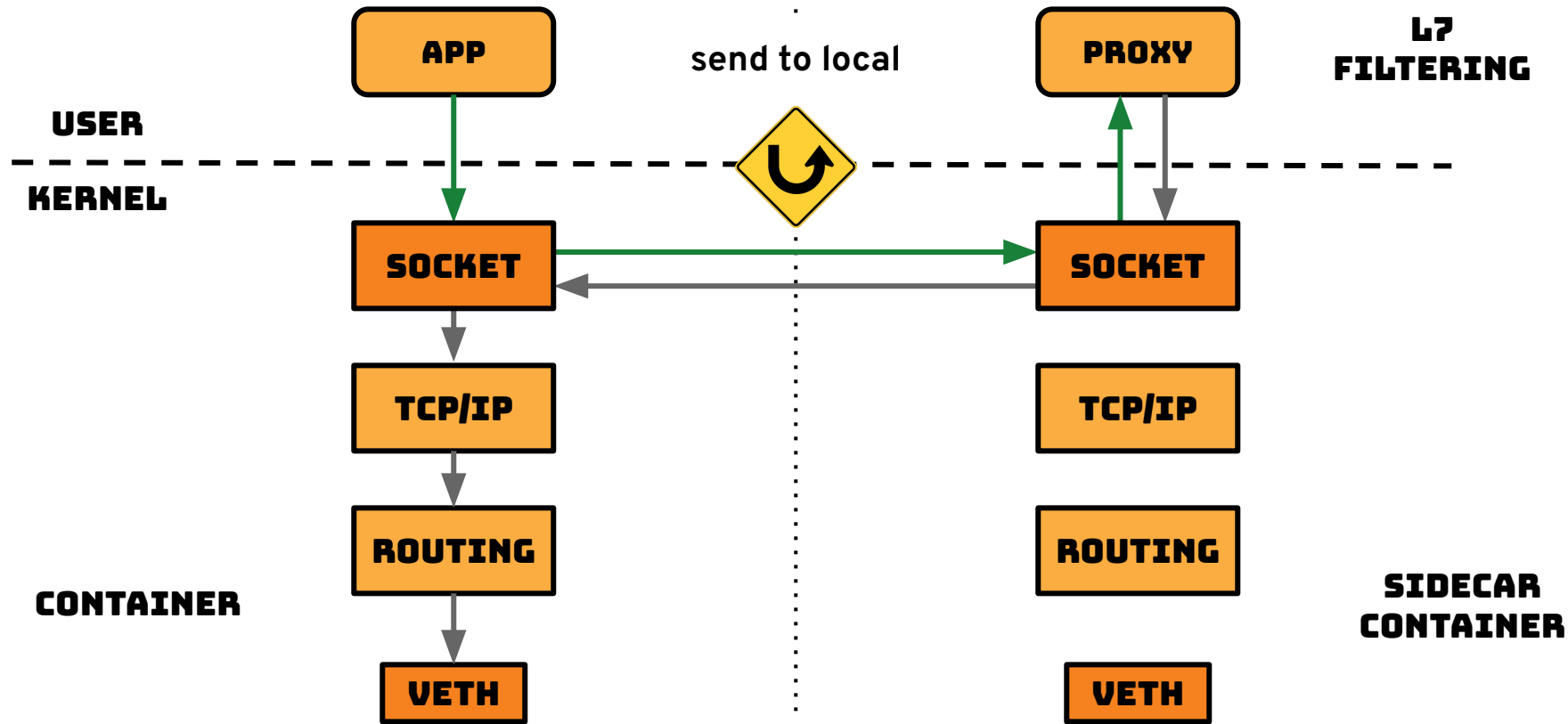
# Real life use-cases



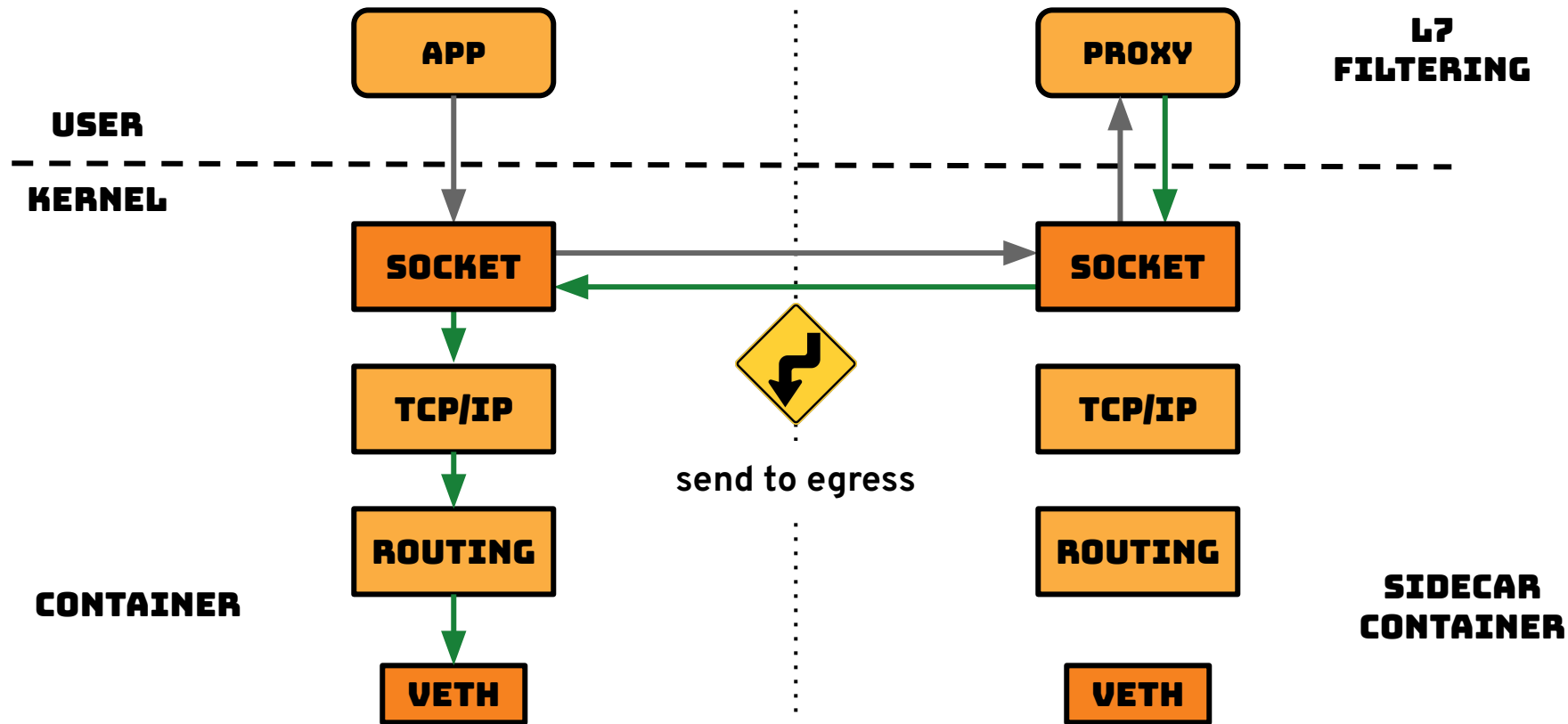
# Cilium project (CNI for K8S)



# Cilium project (CNI for K8S)

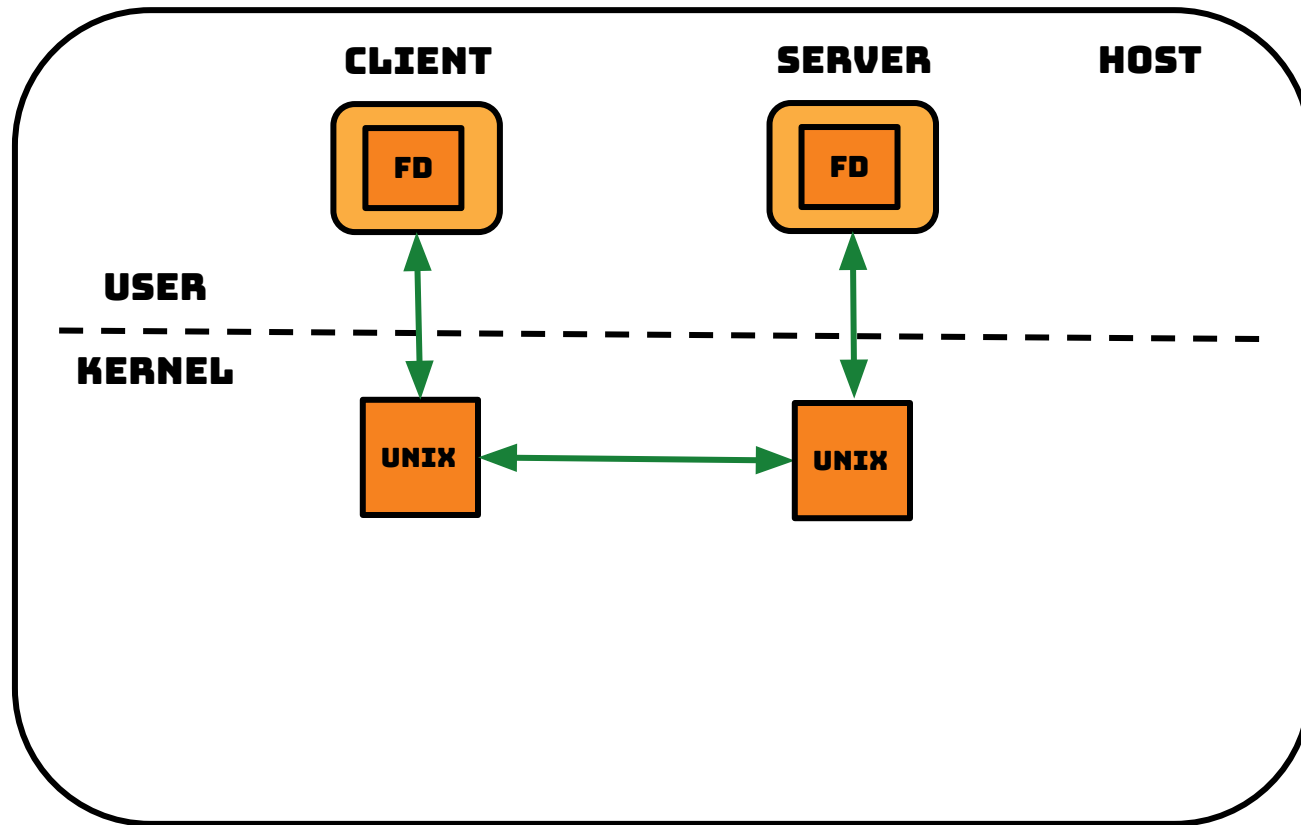


# Cilium project (CNI for K8S)

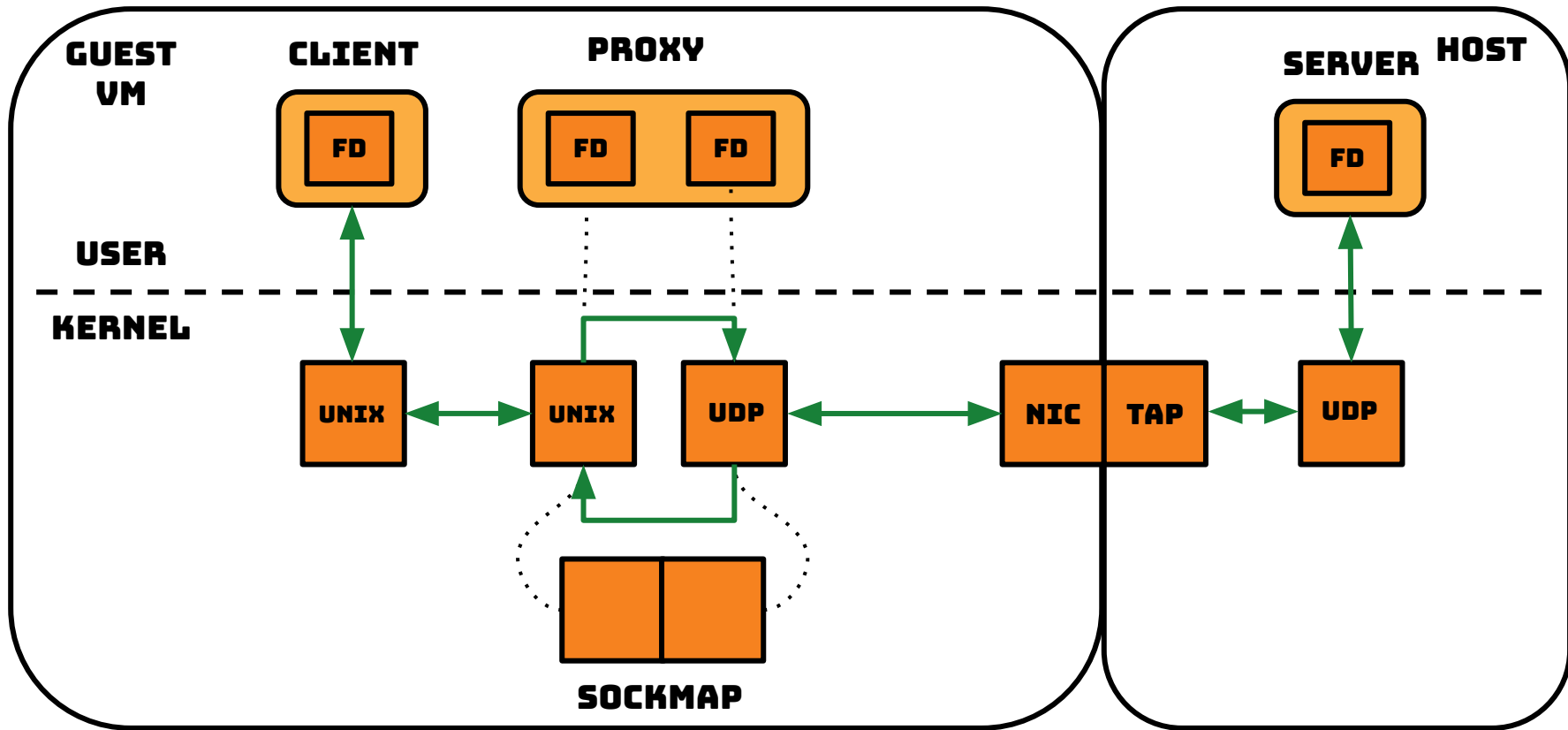




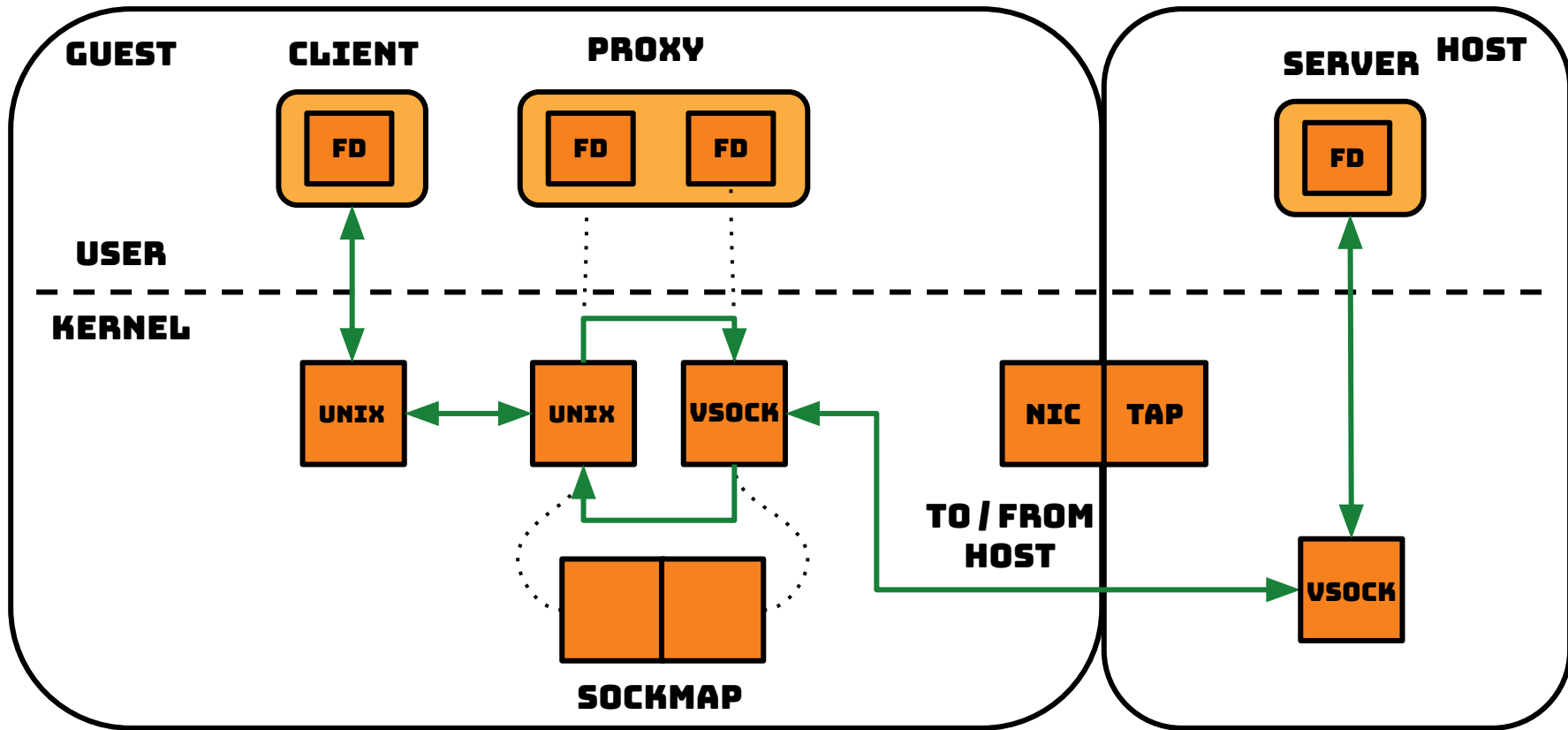
# Bytedance (TikTok)



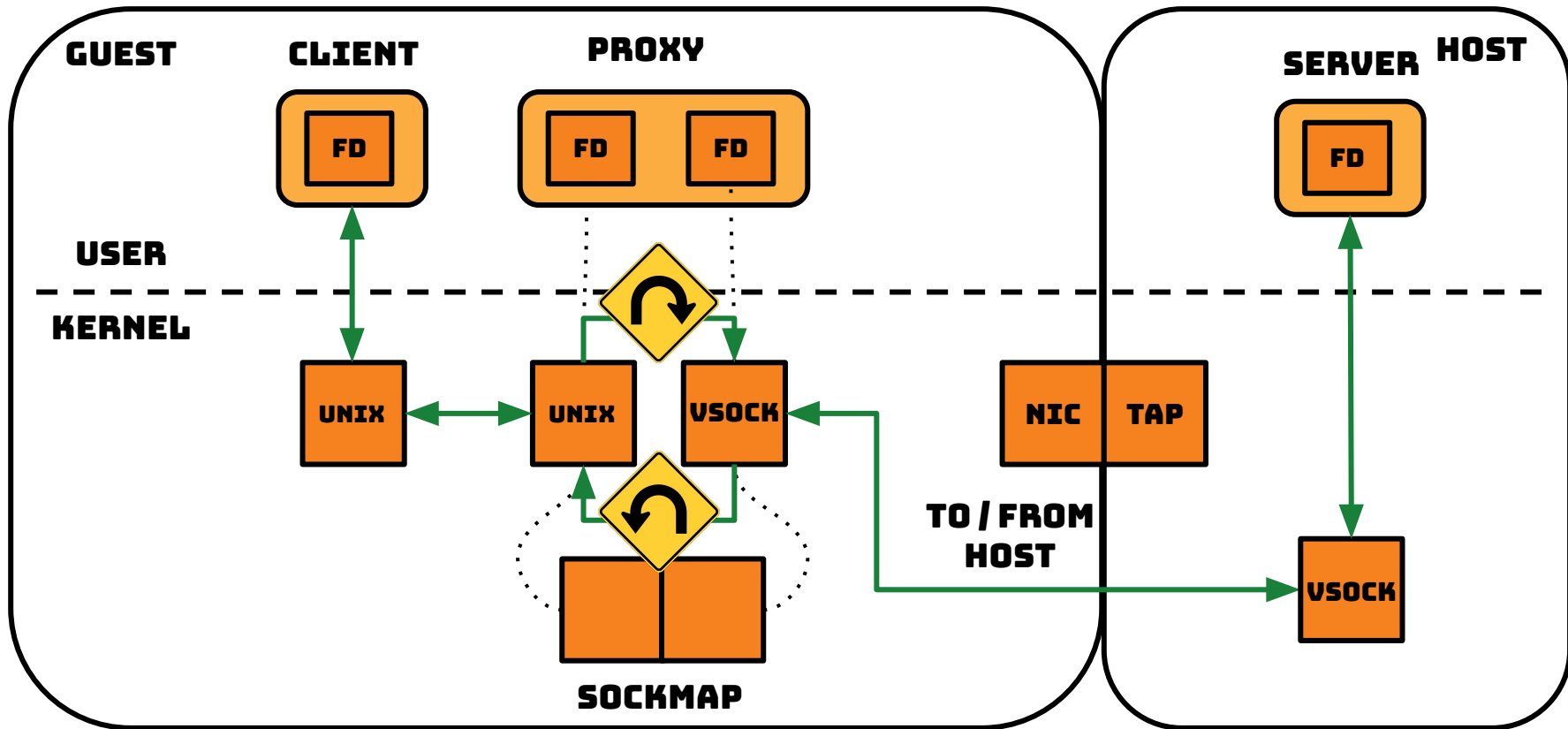
# Bytedance (TikTok)



# Bytedance (TikTok) → Improved



# Bytedance (TikTok) → Improved v2



# Where to learn more?





# RESOURCES

1) **Linux Kernel → BPF Documentation → SOCKMAP and SOCKHASH map**

Includes links to unit tests with API usage examples

2) **LPC 2018: Combining kTLS and BPF for Introspection and Policy Enforcement**

See Daniel & John talk about Cilium SOCKMAP + kTLS use case ([video](#), [slides](#), [paper](#))

3) **Cloudflare Blog: SOCKMAP - TCP splicing of the future**

Read Marek review SOCKMAP from L7 proxy perspective

4) **eBPF Summit 2020: Steering connections to sockets with BPF socket lookup hook**

Another use case for SOCKMAP as a container ([video](#), [slides](#), [code](#))

# THANK YOU



Reach out:

[jakub@cloudflare.com](mailto:jakub@cloudflare.com)

[@jkbs0 @ Twitter X](#)

[jsitnicki @ LinkedIn](#)

Mailing lists:

[bpf@vger.kernel.org](mailto:bpf@vger.kernel.org)

[netdev@vger.kernel.org](mailto:netdev@vger.kernel.org)

Slack:

[#ebpf-kernel-dev @ cilium.slack.com](#)



↑ code & slides repo ↑

<https://github.com/jsitnicki/kubecon-2024-sockmap>

# Attribution

- [Uncle Sam Wants You](#), Public Domain Pictures, CC0
- [A bunch of question marks](#), Wikimedia Commons, CC BY-SA 4.0 DEED
- [Nf knots](#), getarchive.net, CC0
- [Construction blocks](#), vectorportal.com, CC BY 4.0 DEED
- [Container crane image](#), vectorportal.com, CC BY 4.0 DEED
- [Golden Gate bridge silhouette](#), rawpixel.com, CC0
- [Evolution-des-wissens.jpg](#), Wikimedia Commons, CC BY-SA 3.0 DEED
- [Leather pouch](#), rawpixel.com, CC0
- [Redirection.svg](#), Wikimedia Commons, CC BY 3.0 DEED



# Overflow slides

# EVOLUTION OF SOCKMAP

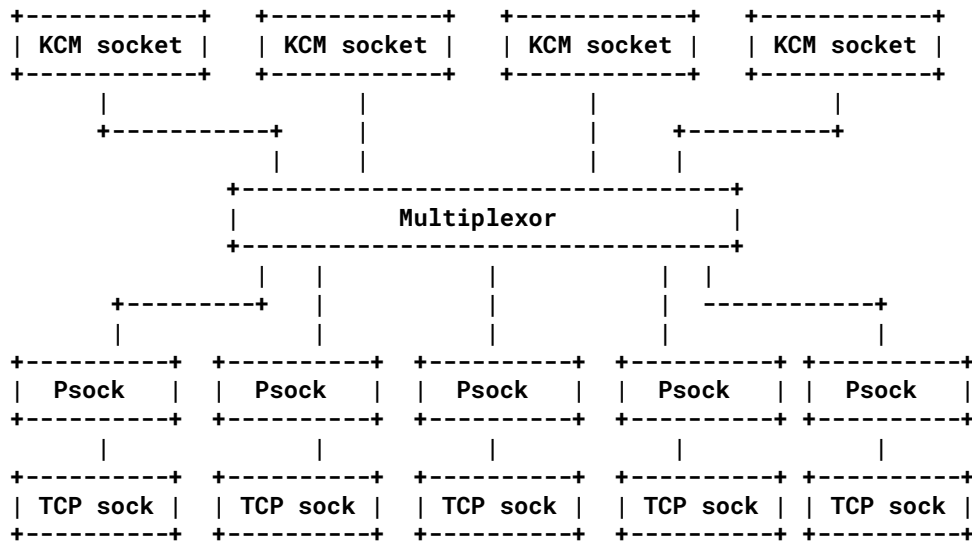


# How did we get here?

## Kernel Connection Multiplexor

Added infrastructure later reused by SOCKMAP - Psock and stream parser program.

2016  
v4.6



# How did we get here?



## Kernel Connection Multiplexor

Added infrastructure later reused by SOCKMAP - Psock and stream parser program.

2016  
v4.6

2017  
v4.14

## SOCKMAP initial version

Filtering and redirect on ingress to socket layer.

Counterpart of XDP DEVMAP.

```
author      John Fastabend <john.fastabend@gmail.com> 2017-08-15 22:32:47 -0700
committer   David S. Miller <davem@davemloft.net> 2017-08-16 11:27:53 -0700
commit      174a79ff9515f400b9a6115643dafd62a635b7e6 (patch)
tree        f48f1fc407adb9bce6fb0e5cddaabd7141acd071
parent      a6f6df69c48b86cd84f36c70593eb4968fceb34a (diff)
download    linux-174a79ff9515f400b9a6115643dafd62a635b7e6.tar.gz
```

## bpf: sockmap with sk redirect support

Recently we added a new map type called dev map used to forward XDP packets between ports (6093ec2dc313). This patches introduces a similar notion for sockets.

A sockmap allows users to add participating sockets to a map. When sockets are added to the map enough context is stored with the map entry to use the entry with a new helper

```
bpf_sk_redirect_map(map, key, flags)
```

This helper (analogous to `bpf_redirect_map` in XDP) is given the map and an entry in the map. When called from a sockmap program, discussed below, the skb will be sent on the socket using `skb_send_sock()`.

commit 174a79ff9515 ("bpf: sockmap with sk redirect support")

# How did we get here?



## Kernel Connection Multiplexor

Added infrastructure later reused by SOCKMAP - Psock and stream parser program.

2016  
v4.6

2017  
v4.14

2018  
v4.17

## SK\_MSG program

Filter and redirect at sendmsg() time (TCP only).

## SOCKMAP initial version

Filter and redirect on ingress to socket layer (TCP only).

Counterpart of XDP DEVMAP.

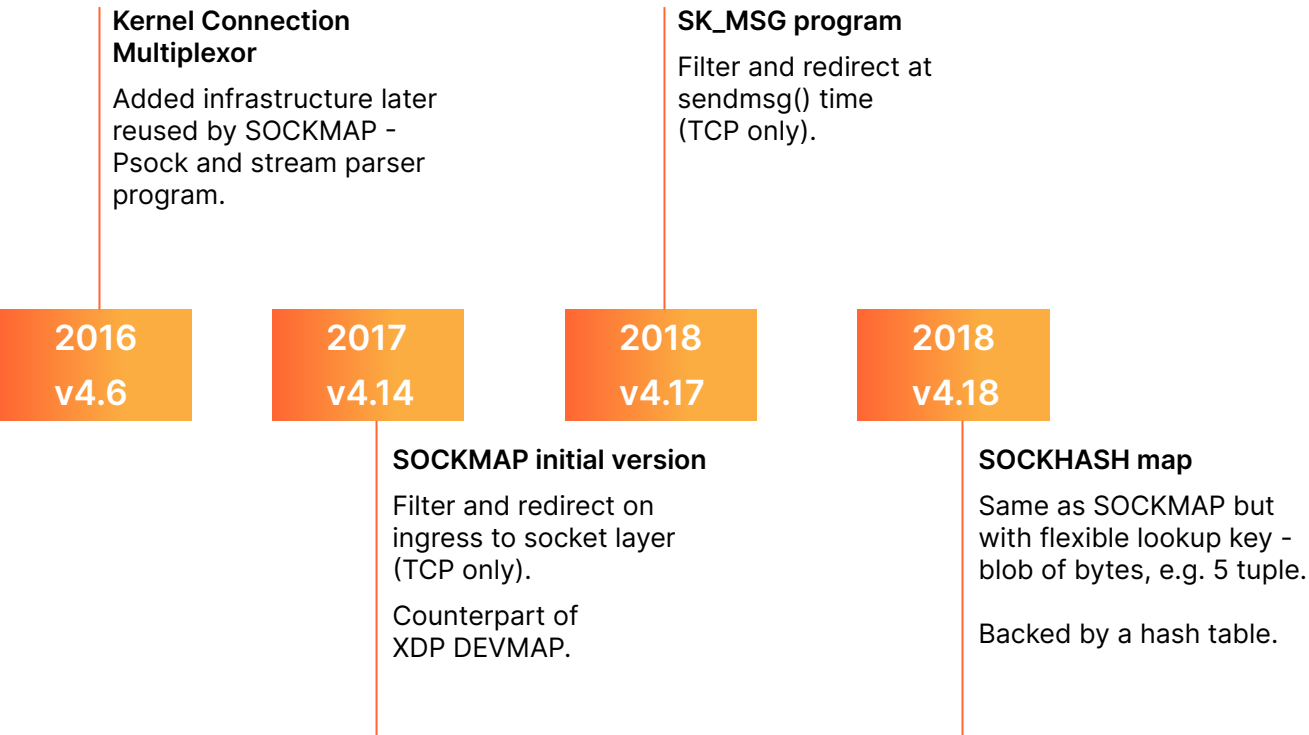
```
author      John Fastabend <john.fastabend@gmail.com> 2018-03-18 12:57:10 -0700
committer   Daniel Borkmann <daniel@iogearbox.net> 2018-03-19 21:14:38 +0100
commit      4f738adba30a7cfc006f605707e7aee847ffefa0 (patch)
tree        6603749a44356d3a44110c44f890a45b88d7e935
parent      8c05dbf04b2882c3c0bc43fe7668c720210877f3 (diff)
download    linux-4f738adba30a7cfc006f605707e7aee847ffefa0.tar.gz
```

## bpf: create tcp\_bpf\_ulp allowing BPF to monitor socket TX/RX data

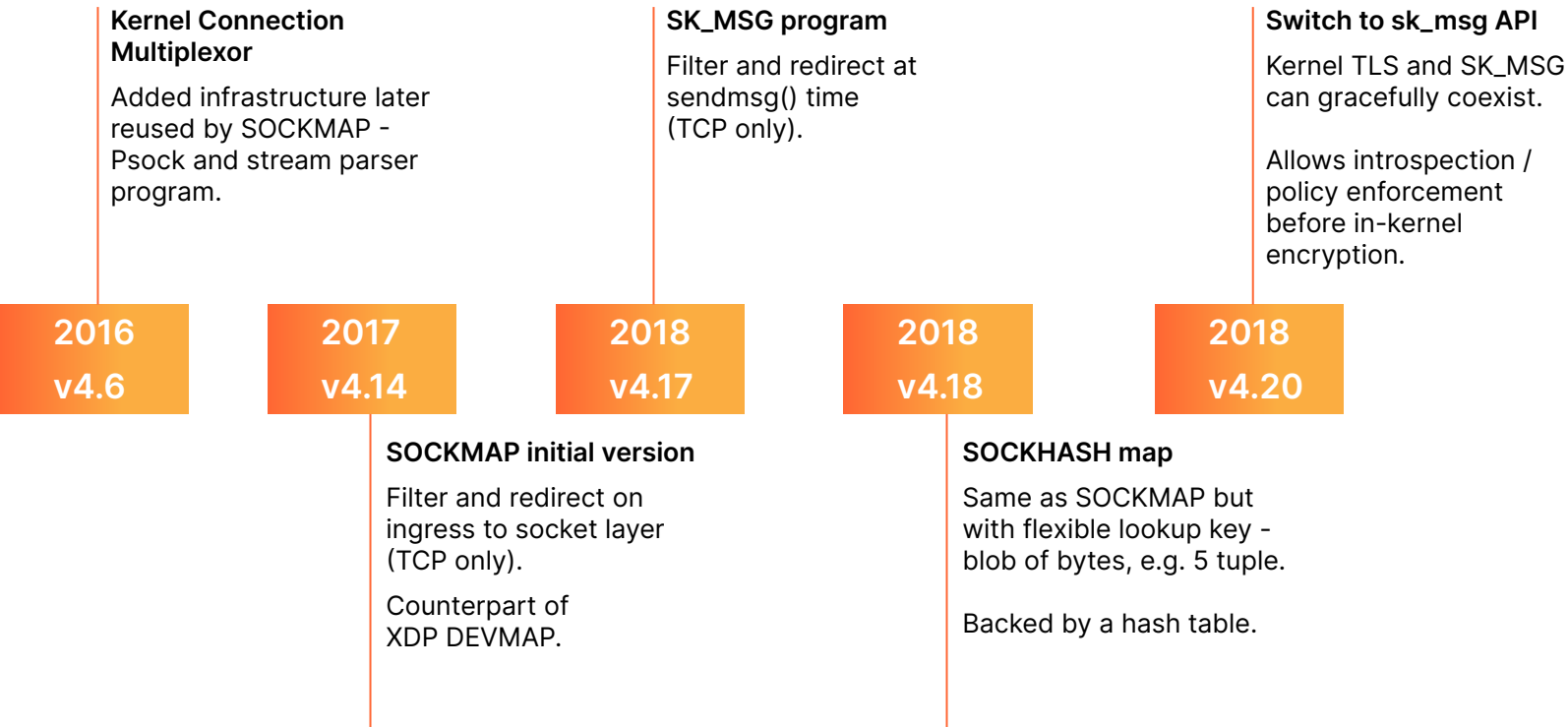
This implements a BPF ULP layer to allow policy enforcement and monitoring at the socket layer. In order to support this a new program type BPF\_PROG\_TYPE\_SK\_MSG is used to run the policy at the sendmsg/sendpage hook. To attach the policy to sockets a sockmap is used with a new program attach type BPF\_SK\_MSG\_VERDICT.

commit 4f738adba30a ("bpf: create tcp\_bpf\_ulp allowing BPF to monitor socket TX/RX data")

# How did we get here?



# How did we get here?



# How did we get here?

## Store TCP or UDP sockets

SOCKMAP becomes a generic BPF map for sockets.

It can hold both connected and listening TCP sockets, and any bound UDP socket.

2020

v5.7



# How did we get here?

## Store TCP or UDP sockets

SOCKMAP becomes a generic BPF map for sockets.

It can hold both connected and listening TCP sockets, and any bound UDP socket.

2020  
v5.7

2020  
v5.10

## More BPF programs can update SOCKMAP

Sockets can be inserted into SOCKMAP by a few selected types of BPF programs.

Initially only SOCK\_OPS programs could do it.

# How did we get here?

## Store TCP or UDP sockets

SOCKMAP becomes a generic BPF map for sockets.

It can hold both connected and listening TCP sockets, and any bound UDP socket.

2020  
v5.7

2020  
v5.10

## BPF iterators support

Iterate over SOCKMAP from BPF context.

Allows copying socket references from one SOCKMAP to another.

2020  
v5.10

## More BPF programs can update SOCKMAP

Sockets can be inserted into SOCKMAP by a few selected types of BPF programs.

Initially only SOCK\_OPS programs could do it.

# How did we get here?

## Store TCP or UDP sockets

SOCKMAP becomes a generic BPF map for sockets.

It can hold both connected and listening TCP sockets, and any bound UDP socket.

2020  
v5.7

## BPF iterators support

Iterate over SOCKMAP from BPF context.

Allows copying socket references from one SOCKMAP to another.

2020  
v5.10

2021  
v5.13

## More BPF programs can update SOCKMAP

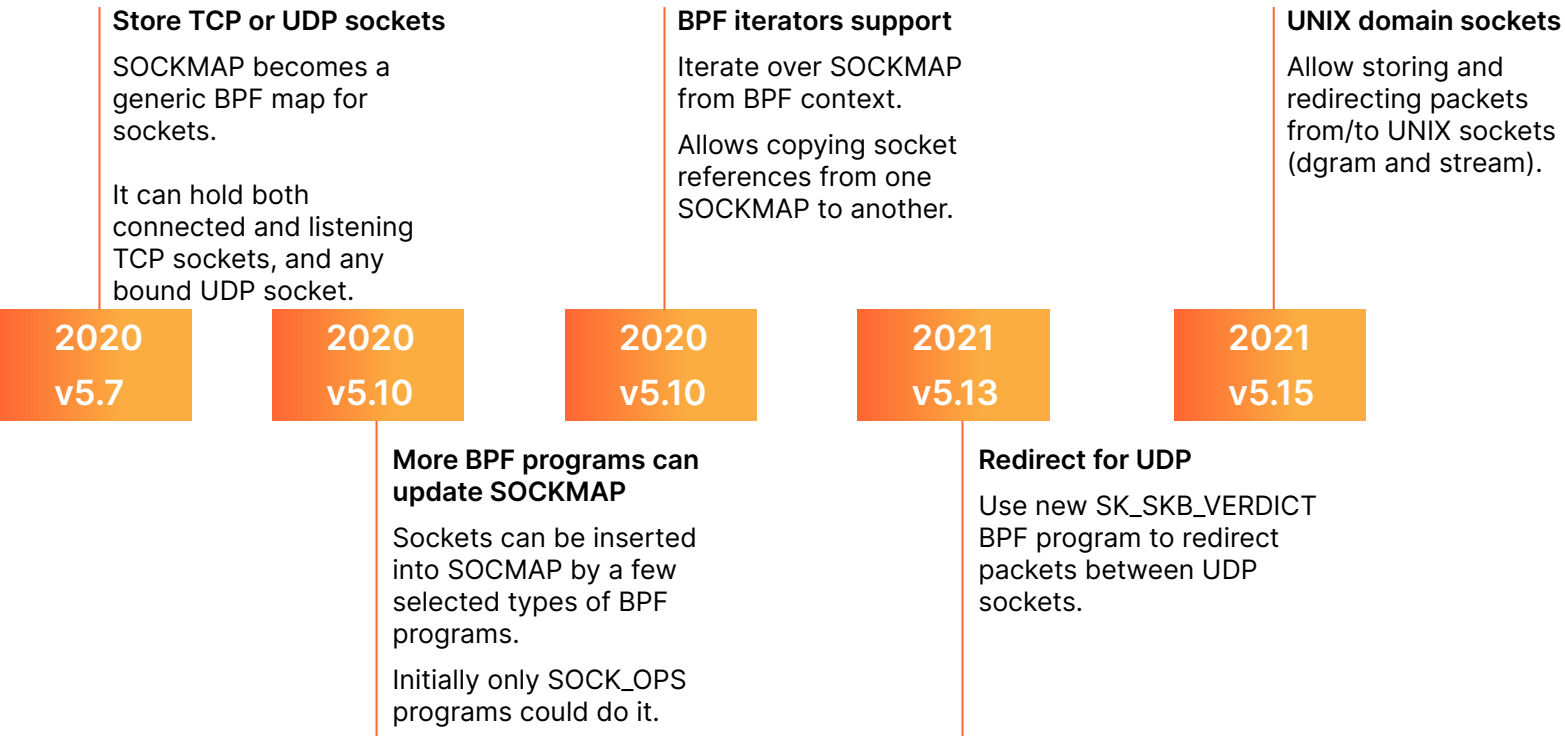
Sockets can be inserted into SOCKMAP by a few selected types of BPF programs.

Initially only SOCK\_OPS programs could do it.

## Redirect for UDP

Use new SK\_SKB\_VERDICT BPF program to redirect packets between UDP sockets.

# How did we get here?



# How did we get here?

## VSOCK domain sockets

Redirecting from/to  
VSOCK sockets (stream  
and seqpacket).

2023

v6.4

# How did we get here?

