리눅스 커널 BPF 소개와 적용사례

2- 리얼리눅스송태웅

http://reallinux.co.kr

"BPF in the kernel is similar with V8 in Chrome browser"

"Linux kernel code execution engine"

"Run code in the kernel"

"Run code in the kernel"

\$ readelf -h bpf-prog.o | grep Machine
Machine: Linux BPF

"Run code in the kernel"

\$ readelf -h bpf-prog.o | grep Machine
Machine: Linux BFF
Machine: Advanced Micro Devices X86-64

"Run code in the kernel"

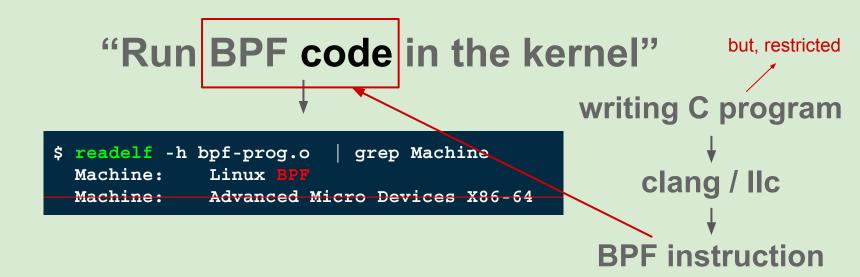
```
$ readelf -h bpf-prog.o | grep Machine
Machine: Linux BPF
Machine: Advanced Micro Devices X86-64
```

writing C program

to clang / Ilc

to BPF instruction





Is it safe?



Is it safe?

"BPF verifier(in-kernel) guarantees"

Is it safe?

"BPF verifier(in-kernel) guarantees"

Check BPF program & Prevent problems before the injection

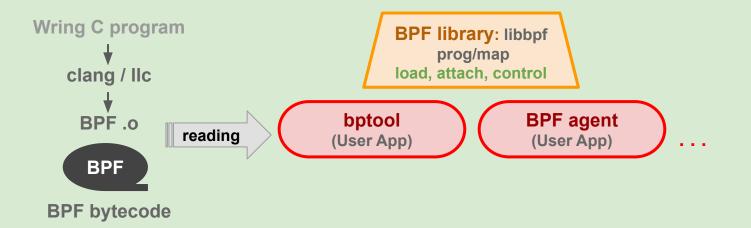
How to execute BPF program?

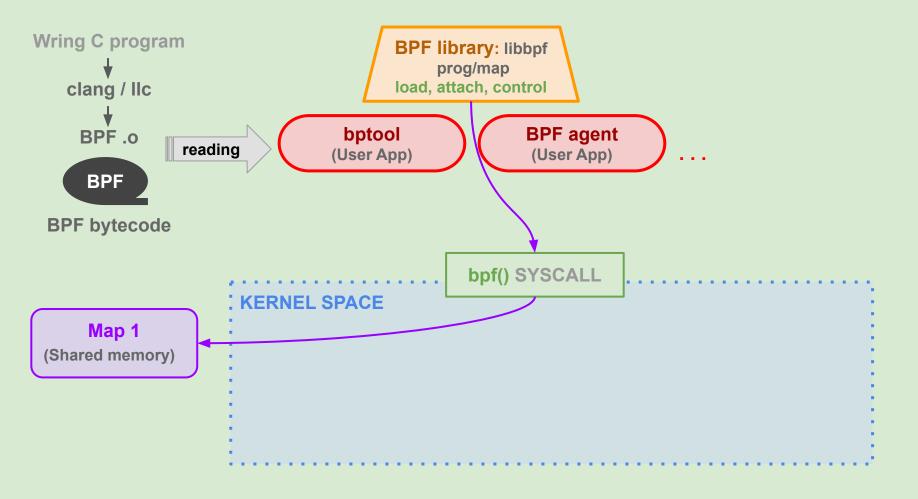
LOAD

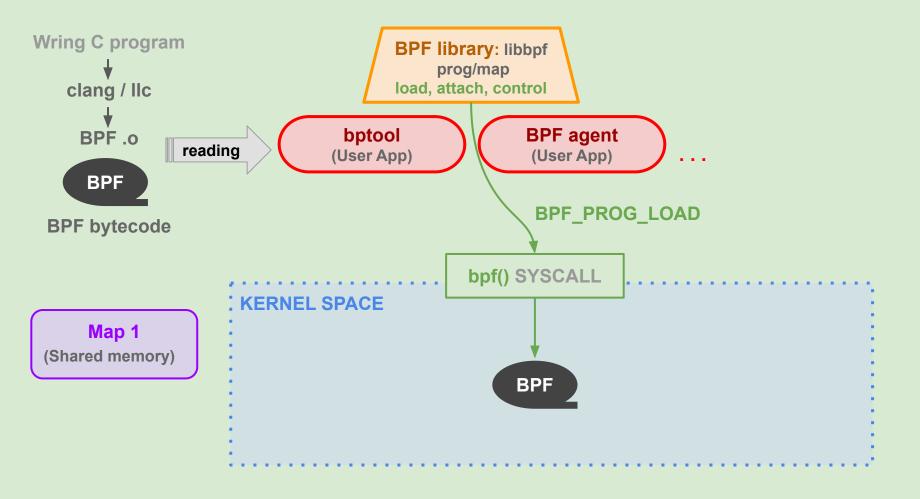
ATTACH

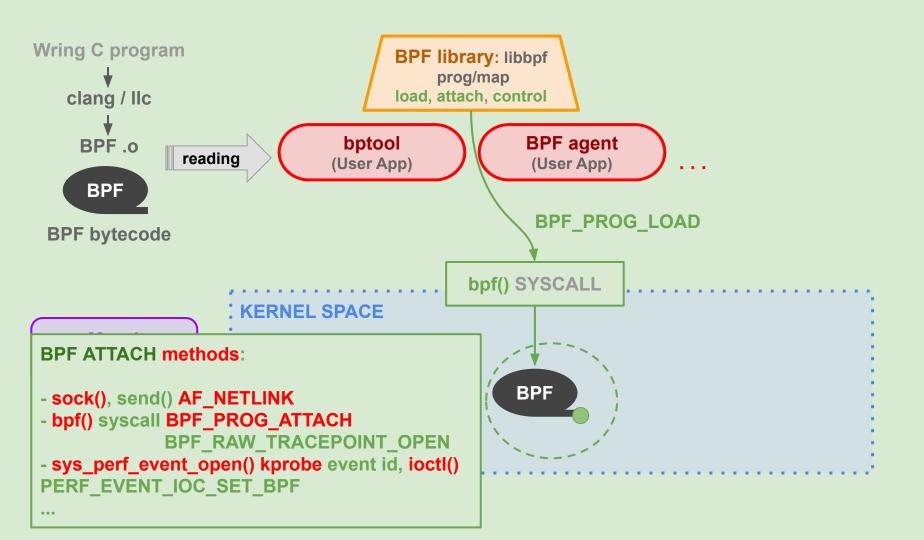
CALLBACK

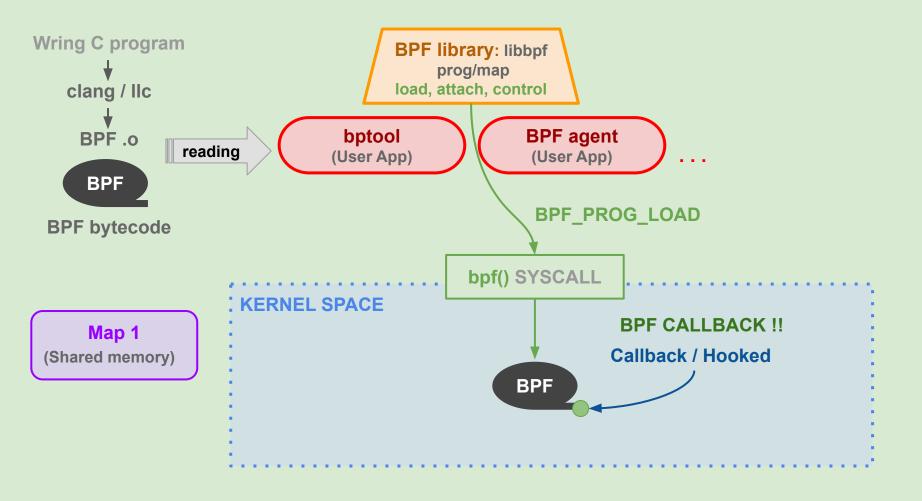


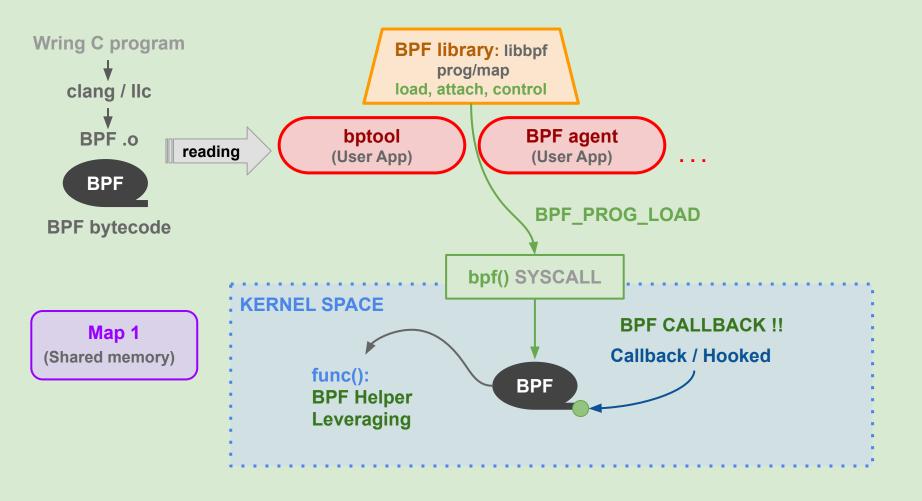


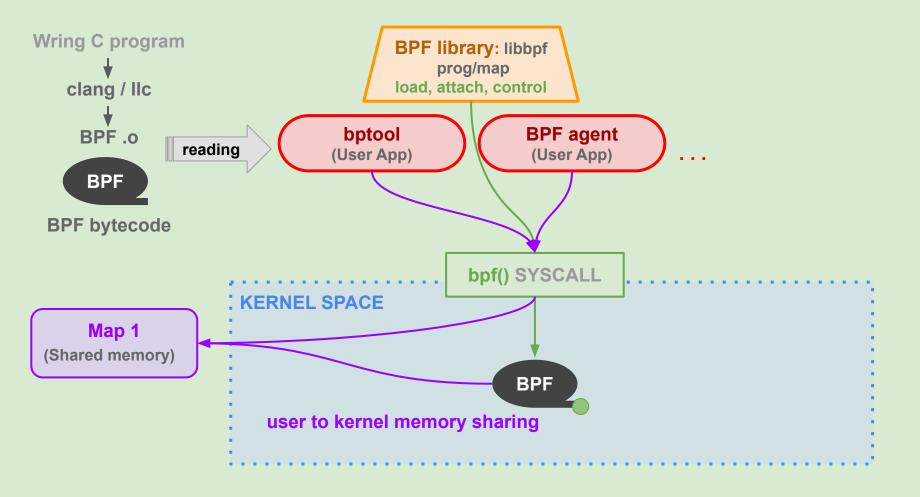


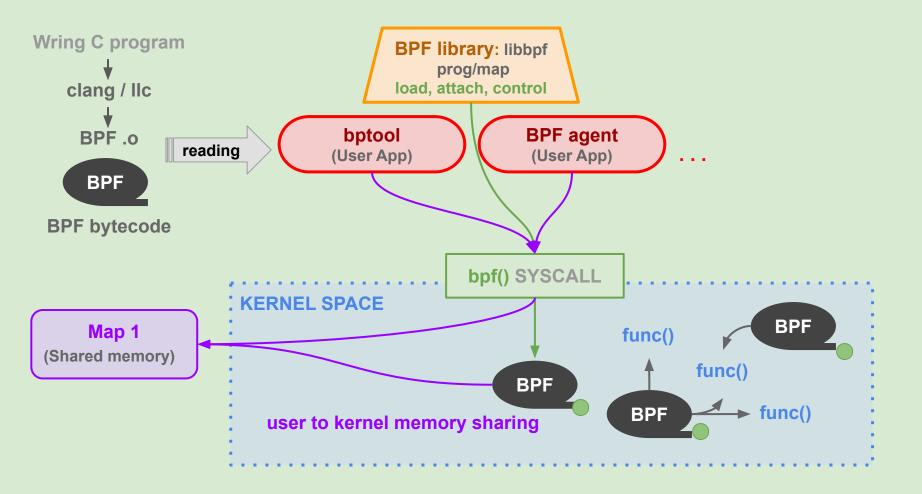














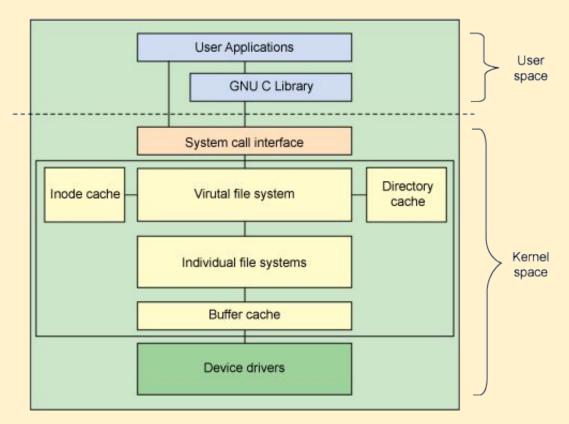
실습1

파일 write 과정

(추적용도 BPF 적용하기)

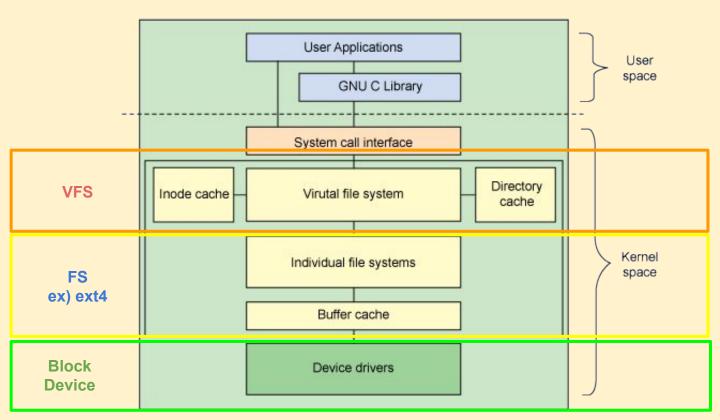
High-level architecture

(Linux Filesystem components)



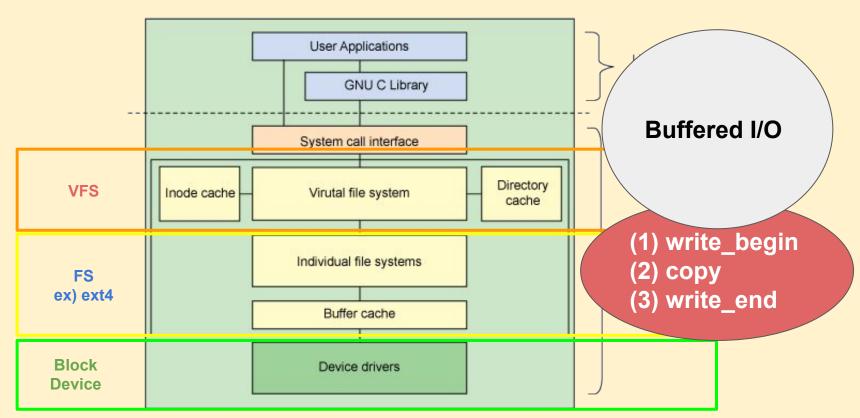
High-level architecture

(Linux Filesystem components)



High-level architecture

(Linux Filesystem components)



₩ 실습1: tracex1_kern.c 예제 파일을 만들고 컴파일 하자

```
. . .
$ cat tracex1 kern.c
#include "bpf helpers.h"
SEC("kprobe/ext4_da_write_begin")
int bpf_prog5(struct pt_regs *ctx)
    char fmt [] = "ext4: da_write_begin \n";
SEC("kprobe/ext4_da_write_end")
int bpf proa6(struct pt reas *ctx)
    char fmt [] = "ext4: da write end\n";
    bpf trace printk(fmt, sizeof(fmt));
 SEC("kprobe/mark_buffer_dirty")
 int bpf_proq7(struct pt_regs *ctx)
    char fmt1 [] = "buffer on pagecache(disk block):%s\n";
    struct buffer head *bh = (struct buffer head *) PT REGS PARM1(ctx):
    char *b_data = _(bh->b_data);
char _license[] SEC("license") = "GPL";
```

🏟실습1: hello.txt 파일 write 과정중 write_begin / write_end / copy 결과 확인

```
$ cd ~/git/linux/samples/bpf
$ make - j4
$ sudo ./tracex1
write-15764 [002] .... 7580.240905: 0: ext4: da_write_begin
write-15764 [002] .... 7580.241251: 0: ext4: da_write_end
write-15764 [002] .N.. 7580.241362: 0: buffer on pagecache(disk block):hello linux filesystem
$ ./write
$ cat hello.txt
hello linux filesystem
```

🖟 실습1: hello.txt 파일 write 과정중 write_begin / write_end / copy 결과 확인

```
$ cd ~/git/linux/samples/bpf
$ make - j4
$ sudo ./tracex1
write-15764 [002] .... 7580.240905: 0: ext4: da_write_begin
write-15764 [002] .... 7580.241251: 0: ext4: da_write_end
write-15764 [002] .N.. 7580.241362: 0: buffer on pagecache(disk block):hello linux filesystem
 ./write
 cat hello.txt
hello linux filesystem
```



실습2

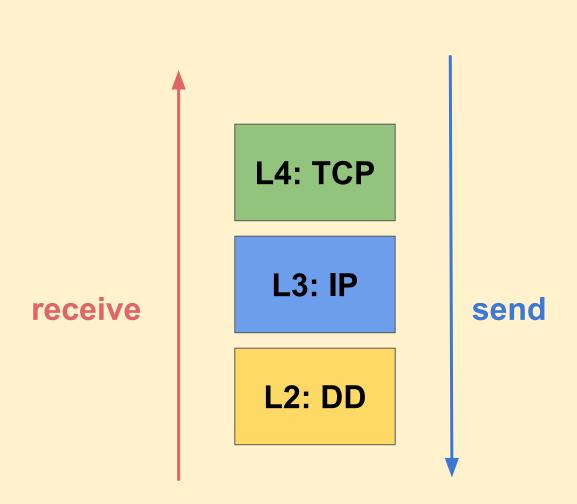
네트워크 패킷 receive 과정

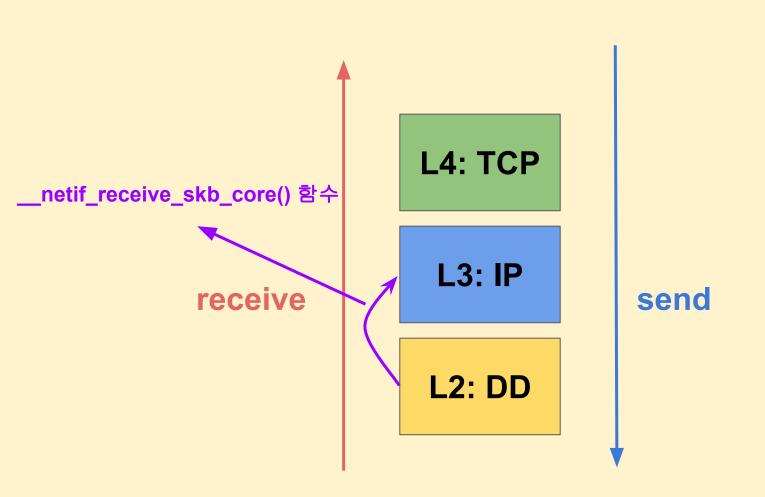
(추적용도 BPF 적용하기)

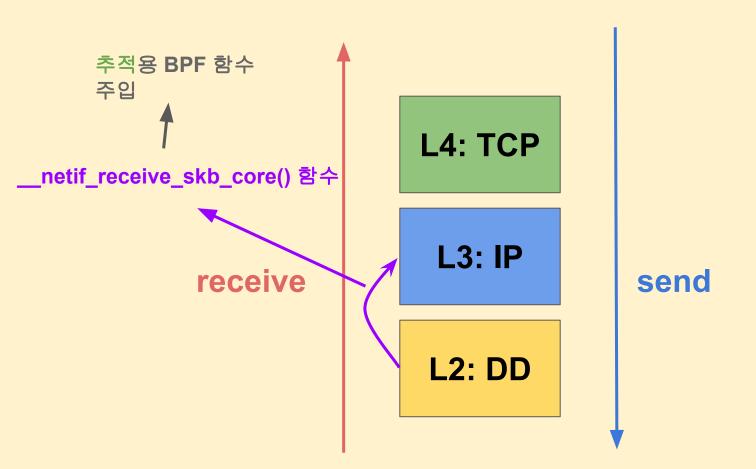
L4: TCP

L3: IP

L2: DD







🎢 실습2: 네트워크 패킷 receive 과정 중 loopback 에서 오는 패킷만 모니터링 하기

```
$ cd ~/git/linux/samples/bpf
$ git checkout -- tracex1_kern.c
$ make - 14
$ sudo ./tracex1
ping-16548 [000] ..s1 8007.417583: 0: skb 000000009218db0b len 84
ping-16548 [000] ..s1 8007.417618: 0: skb 0000000064bf60cb len 84
ping-16548 [000] ..s1 8008.441838: 0: skb 000000009218db0b len 84
ping-16548 [000] ..s1 8008.441982: 0: skb 0000000064bf60cb len 84
$ ping 127.0.0.1
PING 127.0.0.1 (127.0.0.1) 56(84) bytes of data.
64 bytes from 127.0.0.1: icmp_seq=1 ttl=64 time=0.322 ms
64 bytes from 127.0.0.1: icmp_seq=2 ttl=64 time=0.223 ms
64 bytes from 127.0.0.1: icmp_seq=3 ttl=64 time=0.192 ms
64 bytes from 127.0.0.1: icmp seg=4 ttl=64 time=0.191 ms
```

🎢 실습2: 네트워크 패킷 receive 과정 중 loopback 에서 오는 패킷만 모니터링 하기

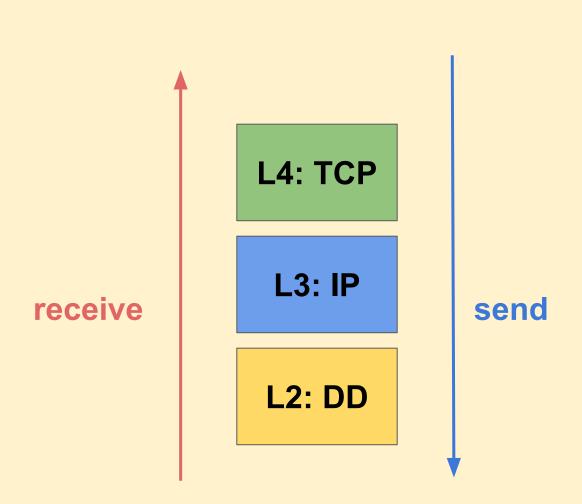
```
$ cd ~/git/linux/samples/bpf
$ git checkout -- tracex1_kern.c
$ make - 14
$ sudo ./tracex1
ping-16548 [000] ..s1 8007.417583: 0: skb 000000009218db0b len 84
ping-16548 [000] ..s1 8007.417618: 0: skb 0000000064bf60cb len 84
ping-16548 [000] ..s1 8008.441838: 0: skb 000000009218db0b len 84
ping-16548 [000] ..s1 8008.441982: 0: skb 0000000064bf60cb len 84
$ ping 127.0.0.1
PING 127.0.0.1 (127.0.0.1) 56(84) bytes of data.
64 bytes from 127.0.0.1: icmp_seq=1 ttl=64 time=0.322 ms
64 bytes from 127.0.0.1: icmp_seq=2 ttl=64 time=0.223 ms
64 bytes from 127.0.0.1: icmp_seq=3 ttl=64 time=0.192 ms
64 bytes from 127.0.0.1: icmp seg=4 ttl=64 time=0.191 ms
```

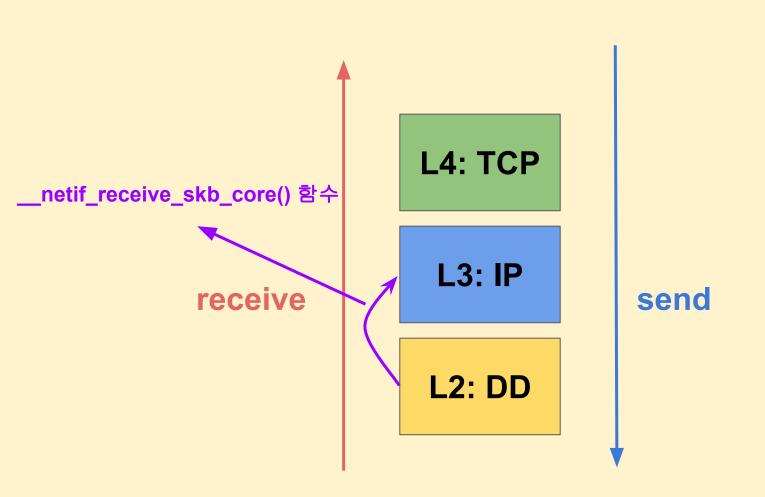


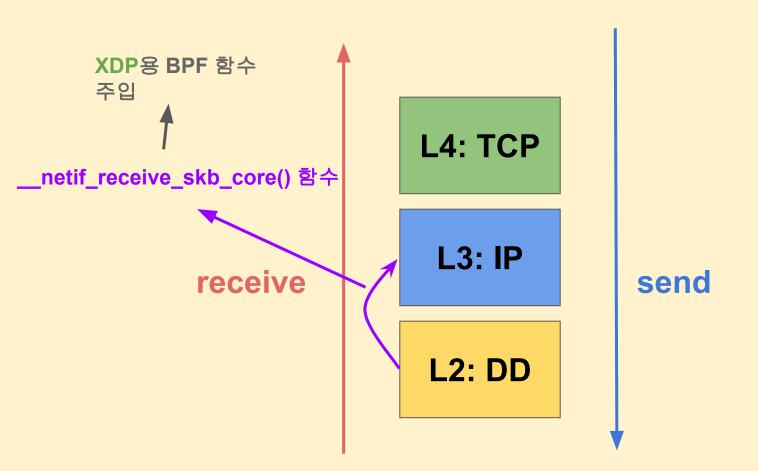
실습3

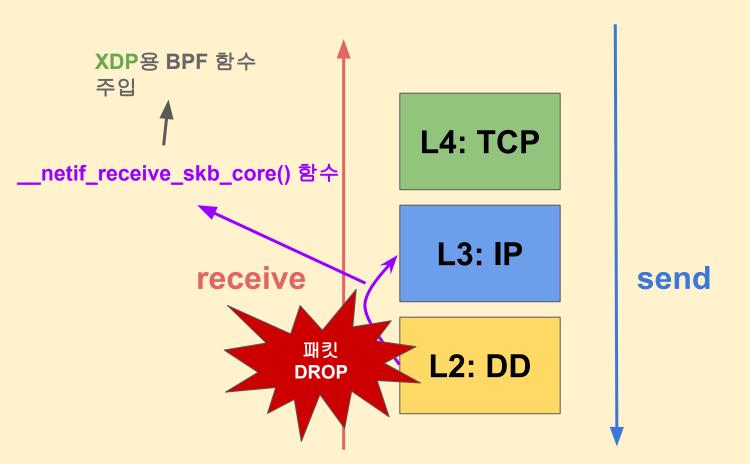
네트워크 패킷 receive 과정

(XDP용도 BPF 적용하기)









🎢 실습3: 예제소스 xdp-drop.c 를 만들자

```
$ vim xdp-drop.c
$ cat xdp-drop.c
#include <linux/bpf.h>
#ifndef __section
# define __section(NAME)
  __attribute__((section(NAME), used))
#endif
__section("prog")
int xdp_drop(struct xdp_md *ctx)
    return XDP_DROP;
char __license[] __section("license") = "GPL";
```



₩ 실습3: 예제소스 xdp-drop.c 를 clang으로 컴파일 하자

```
● ● ● ● ● # xdp-drop.c 예제소스 코드 clang으로 컴파일 하기
$ clang -02 -Wall -target bpf -I/usr/include/x86_64-linux-gnu/ -c xdp-drop.c -o xdp-drop.o
```



실습3: 예제 BPF프로그램 xdp-drop.o 를 ip link 명령으로 커널에 주입하자

* BPF로 짠 새로운 함수(xdp_drop)가 커널코드로 실시간 주입된다.



(주입된)xdp_drop 함수 call-back 테스트



₩ 실습3: 예제 BPF프로그램 xdp-drop.o 를 ip link 명령으로 커널에 제거하자

* BPF로 짠 새로운 함수(xdp_drop)가 커널코드에서 제거 된다.

```
# xdp 제거후 확인
$ sudo ip link set dev lo xdp off
$ ip a
```



🎢 참고: XDP가 지원이 Network Driver 들을 위해 L2 -> L3 로 올라오면서

호출되는 __netif_receive_skb_core() 함수에

do_xdp_generic() 함수가 있고

여기에 주입된 bpf함수를 호출해주는 hook 포인트가 있다.



추가 BPF 적용 사례 소개

1. 적외선 리모트 컨트롤러 decode 함수 주입용 (LIRC - Linux Infrared Remote Control)

2. 특정 커널 함수 return 값 변경하기

(error 테스트 목적: bpf_override_return 함수 활용)