

AF_XDP

Collins Huff

2021-06-15

Motivation

I'm interested in scanning the internet.

How to Scan 0.0.0.0/0 - TCP

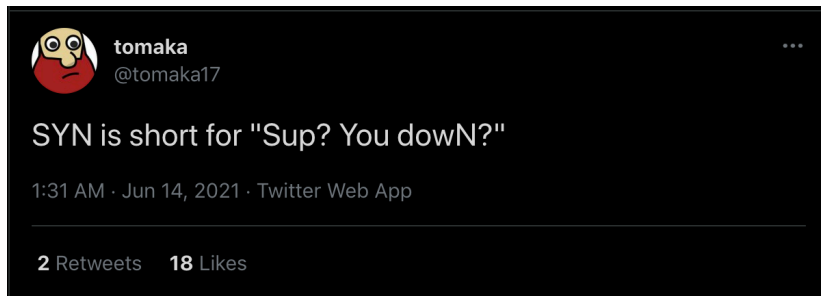
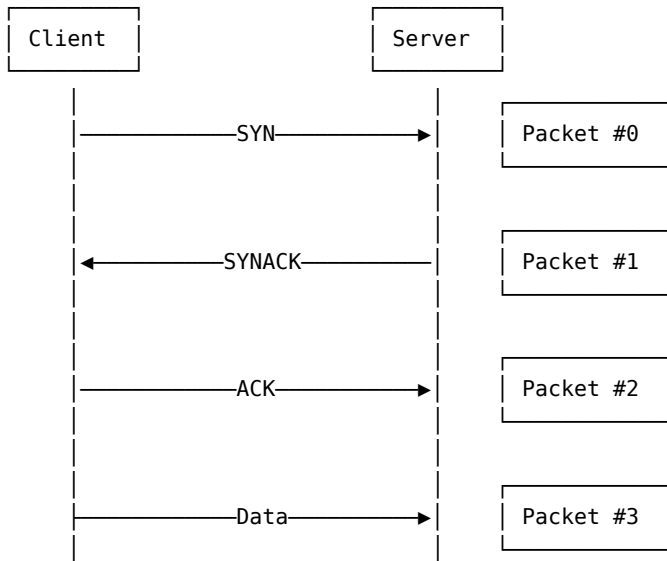


Figure 1: SYN

How to Scan 0.0.0.0/0 - TCP



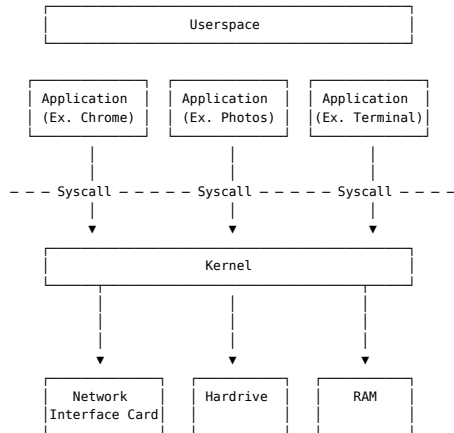
Zmap

Sends TCP SYN packets, listens for SYNACK to determine open ports.

Speed Matters

- ▶ There are 4,294,967,296 IPv4 Addresses
- ▶ Scanning all of IPv4 at 100,000 packets per second takes 12 hours
- ▶ Scanning all of IPv4 at 1,000,000 per second takes 71 minutes
- ▶ Scanning all of IPv4 at 10,000,000 packets per second takes 7 minutes

OS/Kernel Review



Fast Packet Processing

There are two main methods for fast packet processing:

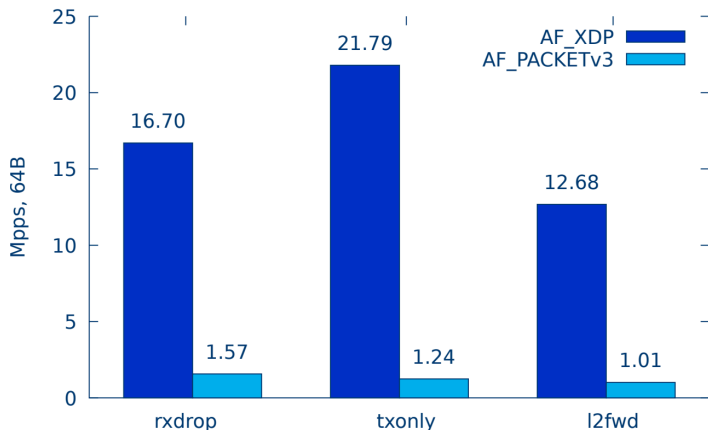
- ▶ In-Kernel: AF_PACKET, in kernel, slow but easy to use
- ▶ Kernel Bypass (DPDK, Netmap, PF_RING), fast but hard to use

Zmap

- ▶ AF_PACKET by default
- ▶ PF_RING if you buy a license

AF_XDP

AF_XDP is a third way: an in-kernel fast path. It is nearly as fast as kernel bypass, but it is built into the kernel.



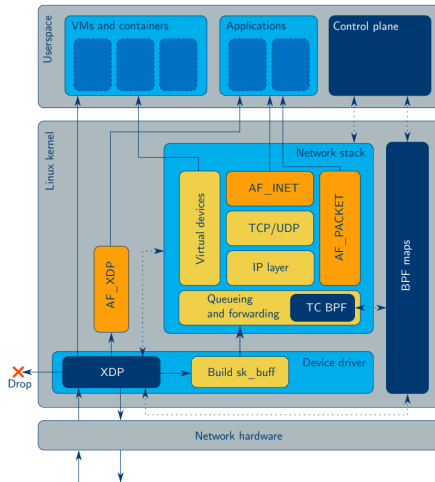
Analogy

Imagine going to the airport

- ▶ In-Kernel packet processing is like going through TSA
- ▶ Kernel bypass is like showing up to the airport and getting on a private jet
- ▶ AF_XDP is like TSA Precheck

AF_XDP

AF_XDP is an address family that is optimized for high performance packet processing. AF_XDP is built on top of two layers of abstraction - eBPF - XDP



AF_XDP components

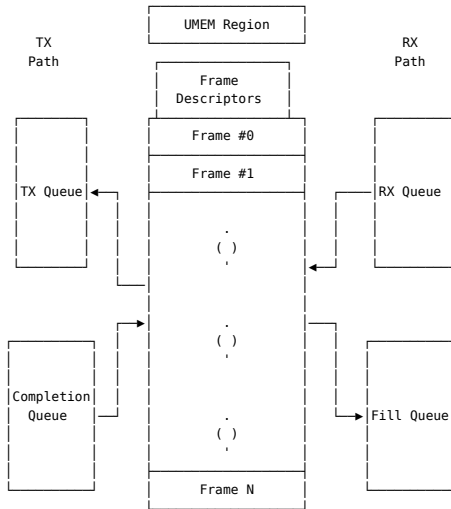
- ▶ 2 queues for TX: the TX Queue and Completion Queue
- ▶ 2 queues for RX: the RX Queue and Fill Queue
- ▶ 1 region of memory called the UMEM, shared between userspace and the kernel.

UMEM and MMAP

We need to allocate a big block of memory to use AF_XDP. The best way to do this is with mmap.

<https://man7.org/linux/man-pages/man2/mmap.2.html>

AF_XDP and xdpsock



Rewrite it in Rust

Starting point: <https://github.com/DouglasGray/xsk-rs>.

Similar to the `af_xdp` example in the kernel source tree.

Uses <https://github.com/alexforster/libbpf-sys>, which is used to set up the shared queues.

Issues

Two problems for my use case:

- ▶ Can't send and receive from multiple threads
- ▶ Complicated API

Design Issue

Original Design

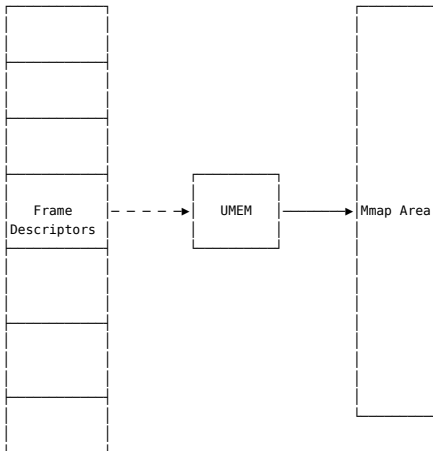
```
pub struct Umem<'a> {
    config: Config,
    frame_size: usize,
    umem_len: usize,
    mtu: usize,
    inner: Box<xsk_umem>,
    mmap_area: MmapArea,
    _marker: PhantomData<&'a ()>,
}

impl Umem<'a_> {
    pub unsafe fn read_from_umem(&self, addr: &usize, len: &usize) -
    > &[u8]

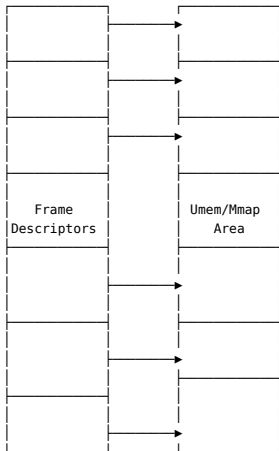
    pub unsafe fn write_to_umem(&mut self,
        frame_desc: &mut FrameDesc, data: &[u8])
}
```

Ownership Diagram

We can represent this with the following ownership diagram (Solid lines represent ownership, dashed lines represent references).



Revised Ownership Diagram



Unsafe

Each frame holds an Arc to the Umem region and constructs it's corresponding slice of bytes using a call to `slice::from_raw_parts_mut`.

Unsafe

```
pub struct Frame<'umem> {  
    addr: usize,  
    len: usize,  
    options: u32,  
    mtu: usize,  
    mmap_area: Arc<MmapArea>,  
    pub status: FrameStatus,  
}
```

Unsafe

```
impl Frame {  
    ...  
    pub unsafe fn read_from_umem(&self, len: usize) -> &[u8] {  
        self.mmap_area.mem_range(self.addr, len)  
    }  
}
```

Unsafe

```
...
```

```
pub unsafe fn write_to_umem(&mut self, data: &[u8]) {  
    let data_len = data.len();  
  
    if data_len > 0 {  
        let umem_region = self.mmap_area.mem_range_mut(&self.addr(), &data_len);  
        umem_region[..data_len].copy_from_slice(data);  
    }  
  
    self.set_len(data_len);  
}  
...  
}
```

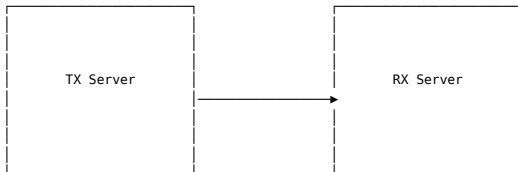

Simplifying the API

```
// Sending a packet
let pkt: Vec<u8> = vec![];
xsk.tx.send(&pkt);

// Receiving a packet
let mut pkt: Vec<u8> = vec![];
let len = xsk.recv(&mut pkt);
```

Performance Test Setup

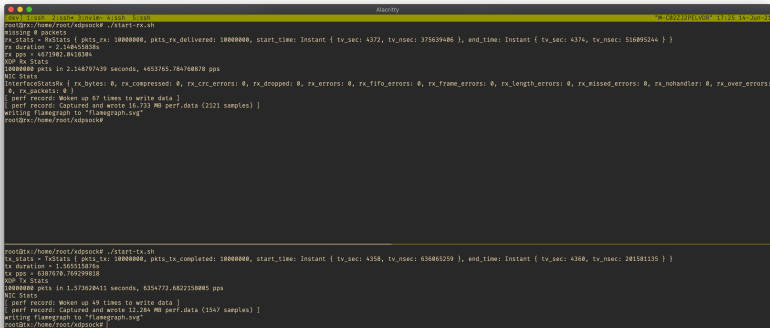
https://github.com/seeyarh/xdpsock/blob/master/examples/dev_to_dev.rs



Performance

Too slow

Should be able to get 14 million pps, only getting 5 million pps



```
Alacritty
root@rx:/home/root/xdpsoc# ./start-rx.sh
missing 0 packets
tx_stats = TStats { pkts_tx: 10000000, pkts_rx_delivered: 10000000, start_time: Instant { tv_sec: 4372, tv_nsec: 375639406 }, end_time: Instant { tv_sec: 4374, tv_nsec: 516095244 } }
rx_duration = 2.148455038s
rx_pps = 4671982.0418304
XDP Rx Stats
10000000 pkts in 2.148797439 seconds, 4653765.784760878 pps
NIC Stats
InterfaceStatsRx { rx_bytes: 0, rx_compressed: 0, rx_crc_errors: 0, rx_dropped: 0, rx_errors: 0, rx_fifo_errors: 0, rx_frame_errors: 0, rx_length_errors: 0, rx_missed_errors: 0, rx_nohandler: 0, rx_over_errors: 0, rx_packets: 0 }
perf record: Woken up 67 times to write data ]
perf record: Captured and wrote 16.733 MB perf.data (2121 samples) ]
writing flamegraph to "flamegraph.svg"
root@rx:/home/root/xdpsoc#

root@rx:/home/root/xdpsoc# ./start-tx.sh
tx_stats = TStats { pkts_tx: 10000000, pkts_tx_completed: 10000000, start_time: Instant { tv_sec: 4358, tv_nsec: 636065259 }, end_time: Instant { tv_sec: 4360, tv_nsec: 201581135 } }
tx_duration = 1.573615874s
tx_pps = 6307670.769299818
XDP Tx Stats
10000000 pkts in 1.573678411 seconds, 6304772.082350003 pps
NIC Stats
perf record: Woken up 49 times to write data ]
perf record: Captured and wrote 13.264 MB perf.data (1547 samples) ]
writing flamegraph to "flamegraph.svg"
root@rx:/home/root/xdpsoc#
```

Optimizing TX

Flamegraphs are a tool to visualize where your program is spending time. `cargo-flamegraph`

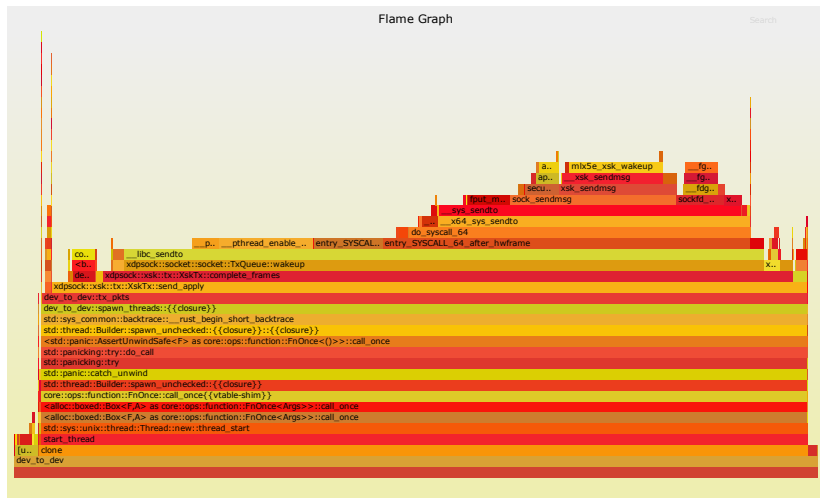


Figure 5: before

Send method unoptimized

The send method calls the complete frames method.

```
pub fn send(&mut self, data: &[u8])
    -> Result<(), XskSendError> {
    self.complete_frames();
    ...

    // Add consumed frames back to the tx queue
    if self.cur_batch_size == self.batch_size {
        self.put_batch_on_tx_queue();
    }

    Ok(())
}
```

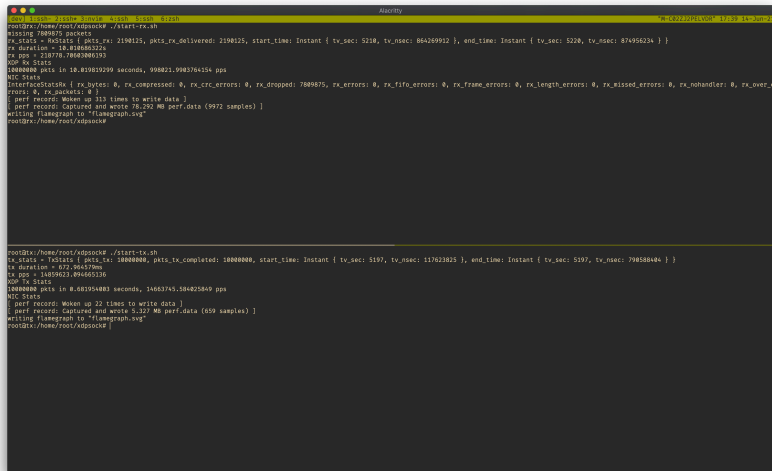
Send method unoptimized

```
fn put_batch_on_tx_queue(&mut self) {  
    ...  
    while unsafe {  
        self.tx_q  
            .produce_and_wakeup(&self.tx_frames[start..end])  
            .expect("failed to add frames to tx queue")  
    } != self.cur_batch_size  
    {  
        // Loop until frames added to the tx ring.  
    }  
    ...  
}
```

Send method unoptimized

```
/// Read frames from completion queue
fn complete_frames(&mut self) -> u64 {
    ...
    if n_free_frames == 0 {
        log::debug!("comp_q.consume() consumed 0 frames");
        if self.tx_q.needs_wakeup() {
            self.tx_q.wakeup()
                .expect("failed to wake up tx queue");
        }
    }
    ...
}
```

Optimizing TX



```
root@kali:~/root/xdpsock# ./start-tx.sh
missing 7889875 packets
rx_stats = RxStats { pkts_rx: 2198125, pkts_rx_delivered: 2198125, start_time: Instant { tv_sec: 5210, tv_nsec: 864269912 }, end_time: Instant { tv_sec: 5220, tv_nsec: 874956234 } }
rx_duration = 10.838086322s
rx_pos = 2187778.7880860193
XDP Tx Stats
10000000 pkts in 10.839819299 seconds, 908023.9983784154 pps
NIC Stats
InterfaceStatsRx { rx_bytes: 0, rx_compressed: 0, rx_crc_errors: 0, rx_dropped: 7889875, rx_errors: 0, rx_fifo_errors: 0, rx_frame_errors: 0, rx_length_errors: 0, rx_missed_errors: 0, rx_nohandler: 0, rx_over_e
rrors: 0, rx_packets: 0 }
[ perf record: Woken up 313 times to write data ]
[ perf record: Captured and wrote 78.292 MB perf.data (9972 samples) ]
writing flamegraph to "flamegraph.svg"
root@kali:~/root/xdpsock#

root@kali:~/root/xdpsock# ./start-tx.sh
tx_stats = TxStats { pkts_tx: 10000000, pkts_tx_completed: 10000000, start_time: Instant { tv_sec: 5197, tv_nsec: 117623825 }, end_time: Instant { tv_sec: 5197, tv_nsec: 798588404 } }
tx_duration = 0.727864579ms
tx_pos = 14889623.89465136
XDP Tx Stats
10000000 pkts in 0.683954083 seconds, 14633745.584825849 pps
NIC Stats
[ perf record: Woken up 22 times to write data ]
[ perf record: Captured and wrote 0.327 MB perf.data (659 samples) ]
writing flamegraph to "flamegraph.svg"
root@kali:~/root/xdpsock#
```

Figure 6: after

Optimizing TX

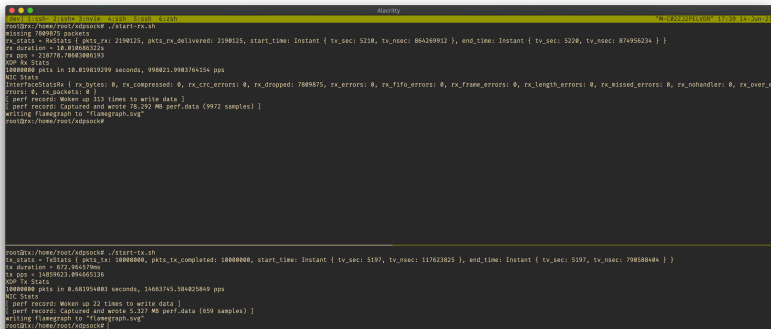


Figure 7: after

Optimizing RX

Now that we have optimized the TX path, we have a new problem: the RX path can't keep up.

We are missing 7,809,875 packets out of 10,000,000 packets, or 78%.

A terminal window titled 'Alacritty' showing the output of a script. The script reports that 7,809,875 packets were missed. It then shows RX statistics for a duration of 10.820863274 seconds, with a rate of 218776.7080806193 pps. The RX path is significantly slower than the TX path, which completed 10,000,000 packets. The terminal output is as follows:

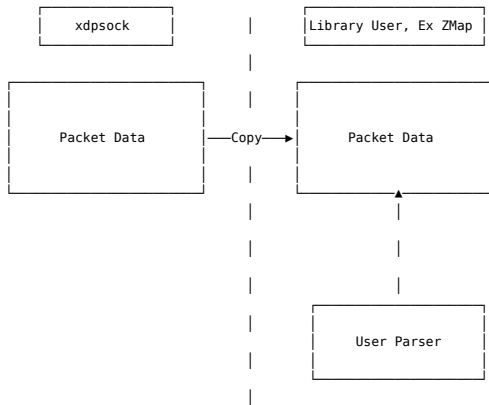
```
root@rx:/home/root/xdpsoc# ./start-rx.sh
missing 7809875 packets
rx_stats = RxStats { pkts_rx: 2198125, pkts_rx_delivered: 2198125, start_time: Instant { tv_sec: 5210, tv_nsec: 864269912 }, end_time: Instant { tv_sec: 5220, tv_nsec: 874956234 } }
rx_duration = 10.820863274
rx_pps = 218776.7080806193
UDP Rx Stats
10000000 pkts in 10.820863274 seconds, 908221.9903764154 pps
NIC Stats
InterfaceStatsRx { rx_bytes: 0, rx_compressed: 0, rx_crc_errors: 0, rx_dropped: 7809875, rx_errors: 0, rx_fifo_errors: 0, rx_frame_errors: 0, rx_length_errors: 0, rx_missed_errors: 0, rx_nohandler: 0, rx_over_e
rrors: 0, rx_packets: 0 }
perf record: Woken up 313 times to write data
perf record: Captured and wrote 78.292 MB perf.data (9972 samples)
writing flamegraph to "flamegraph.svg"
root@rx:/home/root/xdpsoc#

root@rx:/home/root/xdpsoc# ./start-tx.sh
tx_stats = TxStats { pkts_tx: 10000000, pkts_tx_completed: 10000000, start_time: Instant { tv_sec: 5197, tv_nsec: 117623825 }, end_time: Instant { tv_sec: 5197, tv_nsec: 798588404 } }
tx_duration = 0.72964579ms
tx_pps = 14059673.894665136
UDP Tx Stats
10000000 pkts in 0.687954083 seconds, 14663743.584825849 pps
NIC Stats
perf record: Woken up 22 times to write data
perf record: Captured and wrote 0.327 MB perf.data (659 samples)
writing flamegraph to "flamegraph.svg"
root@rx:/home/root/xdpsoc#
```

Optimizing RX

```
pub fn recv(&mut self, pkt_receiver: &mut [u8]) -> usize {
```

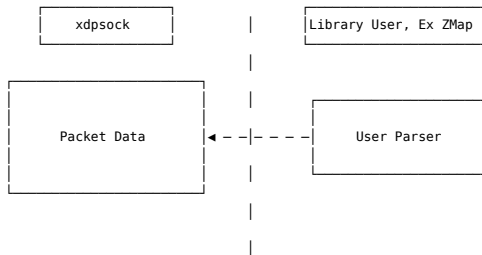
Optimizing RX - Copy



Optimizing RX - Zerocopy

Accept a function, use a closure

Optimizing RX - Zerocopy



Optimizing RX: avoiding copies

```
pub fn recv_apply<F>(&mut self, f: F)
where
    F: FnMut(&[u8]),
{
    ...
    if n_frames_recv > 0 {
        self.apply_batch(n_frames_recv, f);
    }
    ...
}
```

Optimizing RX: avoiding copies

```
fn apply_batch<F>(&mut self, n_frames_recv: usize, mut f: F)
where
    F: FnMut(&[u8]),
{
    ...

    for filled_frame in filled_frames {

        let data = unsafe { filled_frame.read_from_umem(frame.len()) };
        f(data);
    }
    ...
}
```


Optimizing RX: avoiding copies

Now we are only missing 403,862 packets out of 10,000,000 packets, or 4%.

```
root@rx:/home/root# ./start-rx.sh
missing 403862 packets
rx_stats = RxStats { pkts_rx: 9596138, pkts_rx_delivered: 9596138, start_time: Instant { tv_sec: 5453, tv_nsec: 254987839 }, end_time: Instant { tv_sec: 5463, tv_nsec: 267714482 } }
rx duration = 10.022087443s
rx pps = 959386.332162296
NOP Rx Stats
10000000 pkts in 10.025371474 seconds, 997469.2734263465 pps
NIC Stats
InterfaceStatsRx { rx_bytes: 0, rx_compressed: 0, rx_crc_errors: 0, rx_dropped: 403862, rx_errors: 0, rx_fifo_errors: 0, rx_frame_errors: 0, rx_length_errors: 0, rx_missed_errors: 0, rx_nohandler: 0, rx_over_er
rors: 0, rx_packets: 0 }
root@rx:/home/root# ./start-rx.sh

root@rx:/home/root# ./start-tx.sh
tx_stats = TxStats { pkts_tx: 10000000, pkts_tx_completed: 10000000, start_time: Instant { tv_sec: 5439, tv_nsec: 739795926 }, end_time: Instant { tv_sec: 5440, tv_nsec: 422387106 } }
tx duration = 0.7259178ss
tx pps = 14067872.63549827
NOP Tx Stats
10000000 pkts in 0.68853917 seconds, 14694231.340129916 pps
NIC Stats
[ perf record: Woken up 22 times to write data ]
[ perf record: Captured and wrote 0.343 MB perf.data (661 samples) ]
writing flamegraph to "flamegraph.svg"
root@rx:/home/root# ./start-rx.sh
```

C FFI

The Rust FFI Omnibus

<http://jakegoulding.com/rust-ffi-omnibus/>

C FFI

```
#[no_mangle]
pub unsafe extern "C" fn xsk_new(iframe: *const c_char) -> *mut Xsk2 {
    let iframe = {
        assert!(!iframe.is_null());
        CStr::from_ptr(iframe)
    };

    let iframe = iframe.to_str().unwrap();
    let umem_config = UmemConfigBuilder::new()
        ...
    let socket_config = SocketConfigBuilder::new()
        ...
    let n_tx_frames = umem_config.frame_count() / 2;
    let n_tx_batch_size = 1024;

    let xsk = Xsk2::new(
        &iframe,
        0,
        umem_config,
        socket_config,
        n_tx_frames as usize,
        n_tx_batch_size,
    )
    .expect("failed to build xsk");
    Box::into_raw(Box::new(xsk))
}
```

C FFI

```
#[no_mangle]
pub unsafe extern "C" fn xsk_send(xsk_ptr: *mut Xsk2,
    pkt: *const u8, len: size_t) {

    let xsk = {
        assert!(!xsk_ptr.is_null());
        &mut *xsk_ptr
    };

    let pkt = {
        assert!(!pkt.is_null());
        slice::from_raw_parts(pkt, len as usize)
    };

    xsk.tx.send(&pkt).expect("failed to send pkt");
}
```

C FFI

```
#[no_mangle]
pub unsafe extern "C" fn xsk_recv(xsk_ptr: *mut Xsk2,
    pkt: *mut u8, len: size_t) -> u16 {

    let xsk = {
        assert!(!xsk_ptr.is_null());
        &mut *xsk_ptr
    };

    let pkt = {
        assert!(!pkt.is_null());
        slice::from_raw_parts_mut(pkt, len as usize)
    };

    xsk.rx.recv(pkt) as u16
}
```

C FFI

```
int main() {
    char* ifname = "veth0";
    void* xsk = xsk_new(ifname);
    uint16_t len_rcvd;

    int i, j;
    int pkts_to_rcv = 10;
    size_t len = 1500;

    for(i = 0; i < pkts_to_rcv; i++) {
        char buf[MAX_PKT_SIZE] = {0};
        len_rcvd = xsk_rcv(xsk, &buf, len);
        for(j = 0; j < len_rcvd; j++) {
            printf("0x%hhx", buf[j]);
        }
        printf("\n");
    }

    char pkt_to_send[50] = {...};

    for(i = 0; i < pkts_to_rcv; i++) {
        xsk_send(xsk, &pkt_to_send, 50);
    }

    xsk_delete(xsk);
    return 0;
}
```

C FFI

https://github.com/seeyarh/zmap/tree/feature/af_xdp