

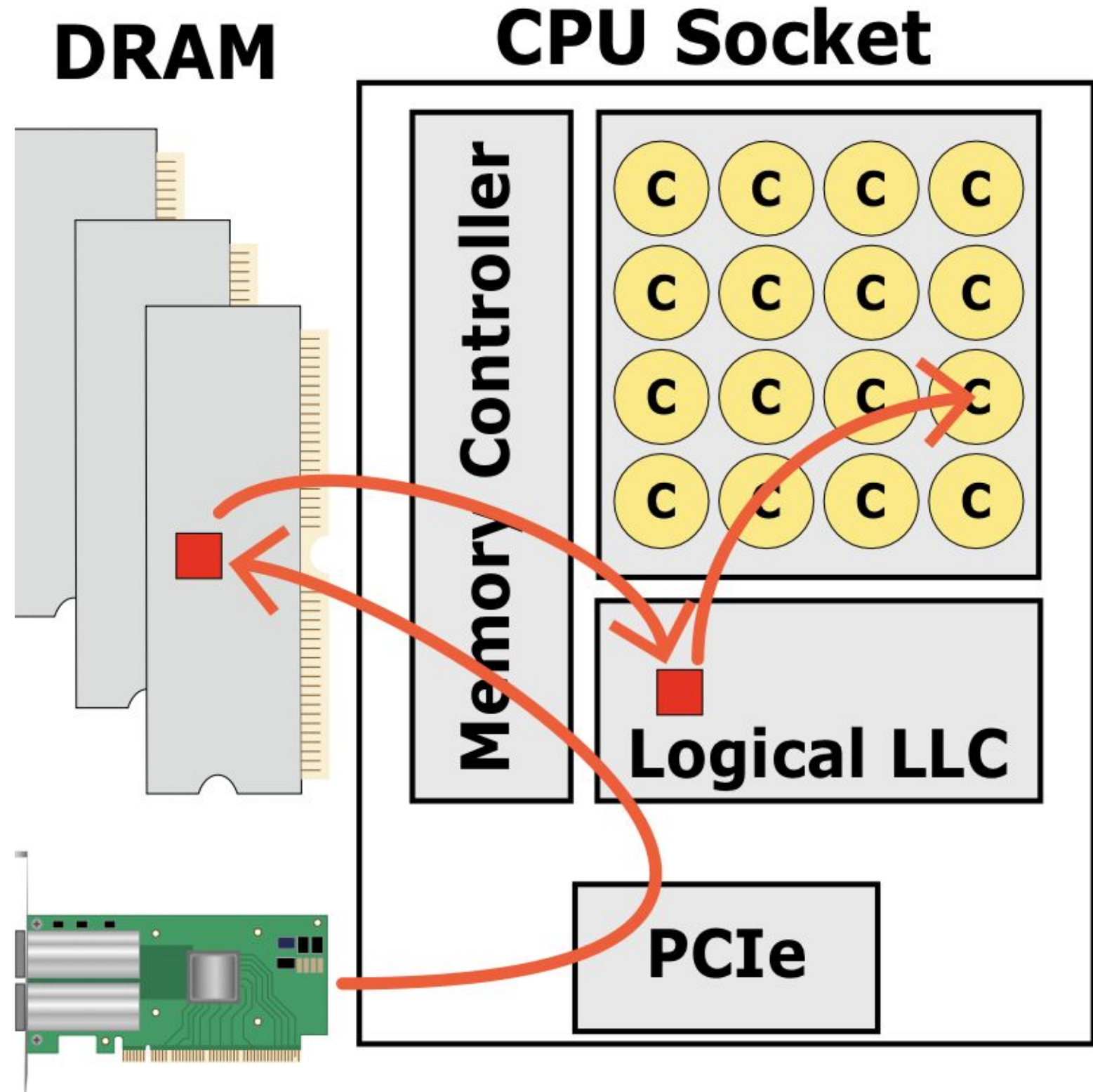
Zero Copy Rx with io_uring

Pavel Begunkov and David Wei

01 Problem Statement

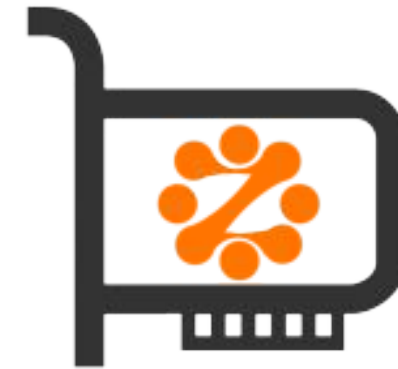
Linux networking Rx overheads

- Memory bandwidth bottlenecks
- Memcpy CPU overheads



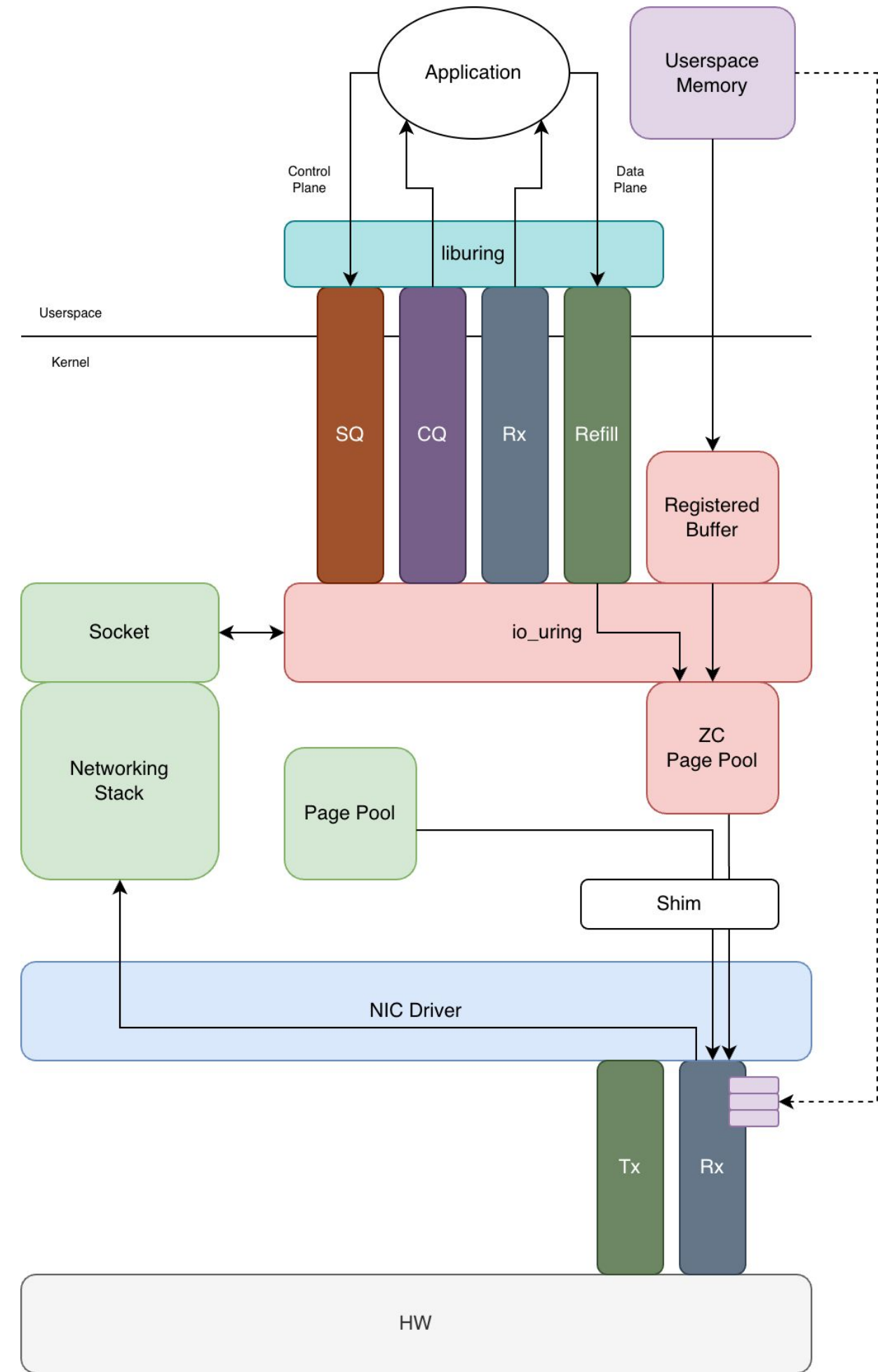
Kernel bypass

- High throughput! Low latency!
- Libraries and applications expect kernel TCP/IP stack
- Re-architecting an entire system around kernel bypass is expensive



Proposal

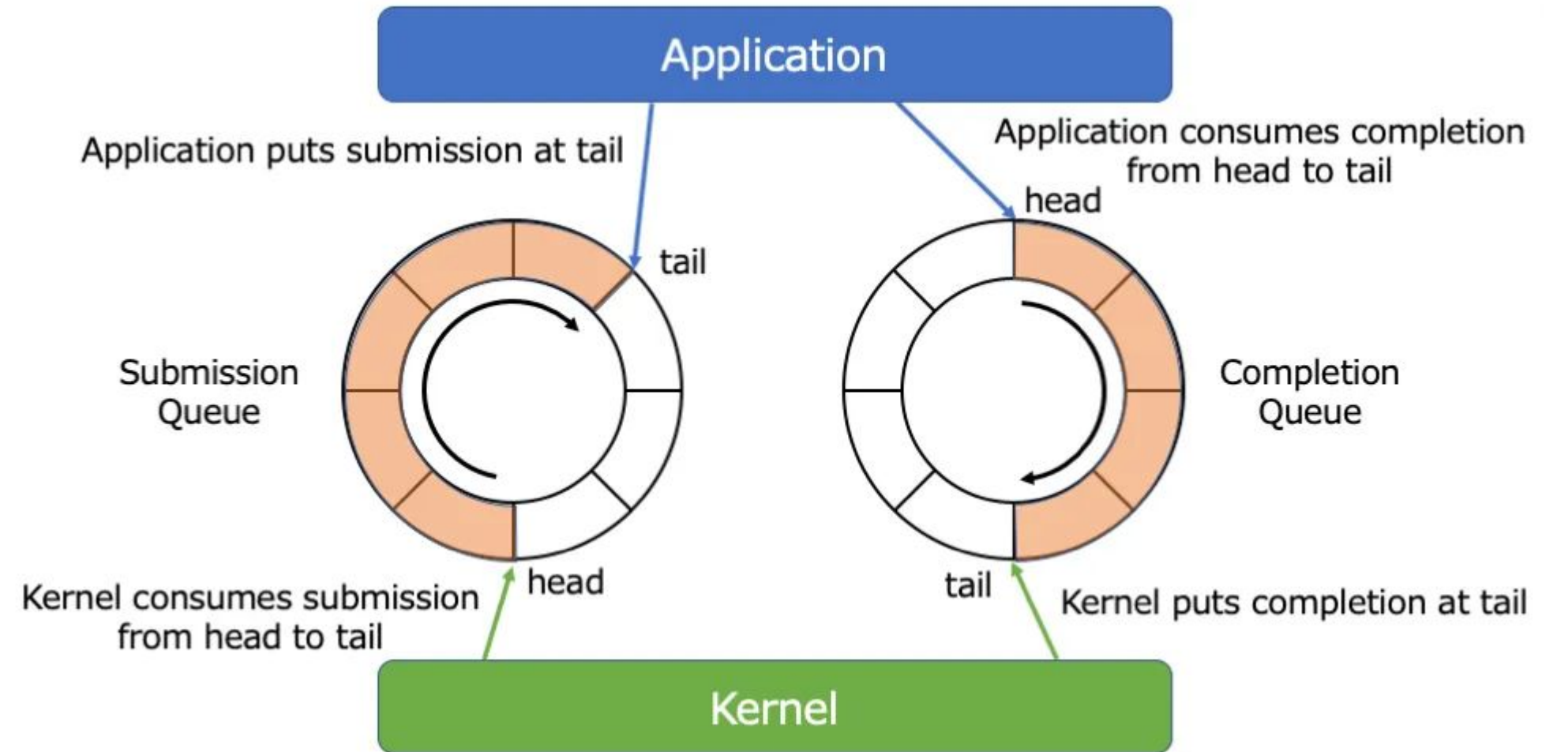
- Hybrid solution
 - Standard control plane using kernel networking stack
 - Fast ZC Rx data plane using io_uring



02 io_uring Primer

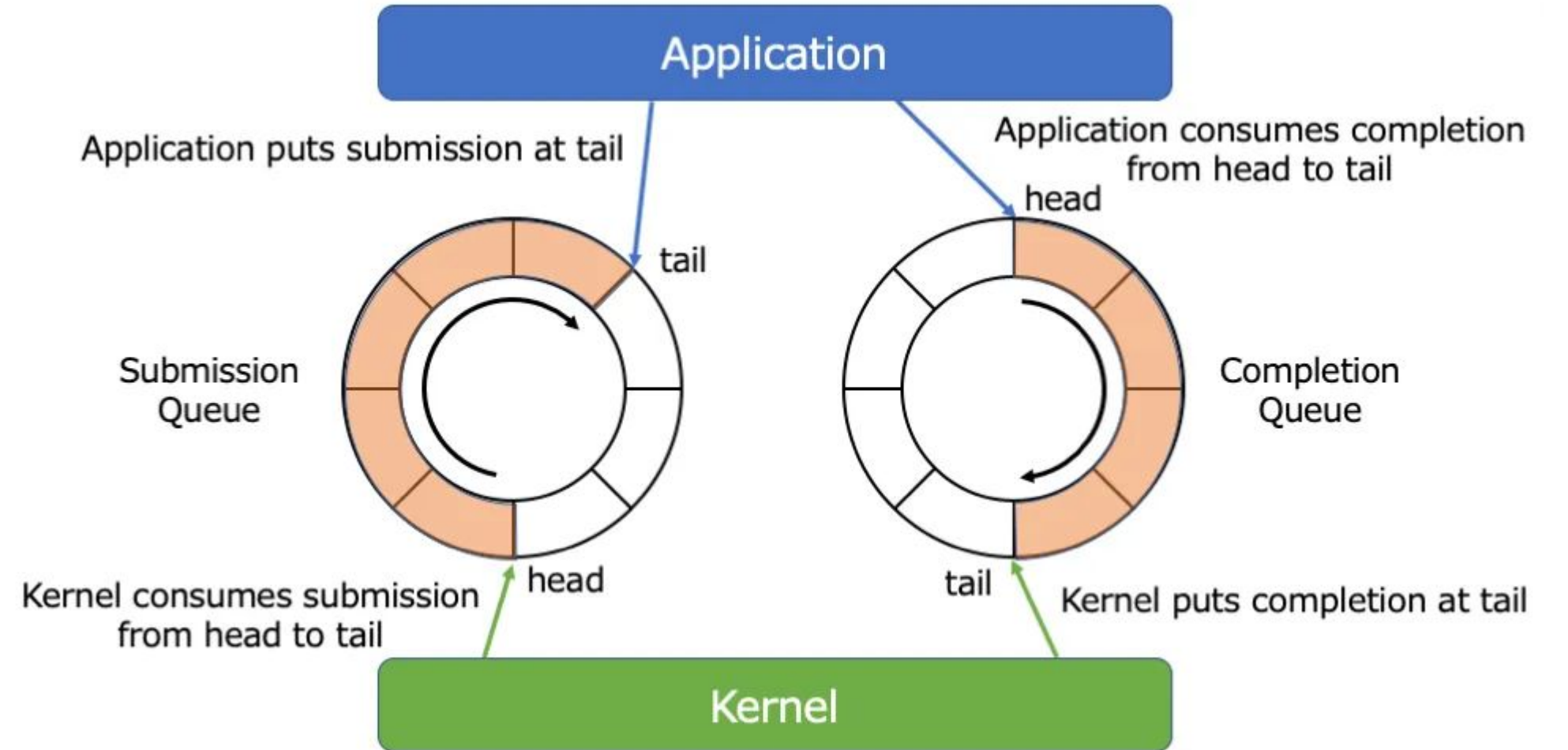
io_uring

- Ring buffers are not new... Similar to what we know and love!
- Userspace submit requests into Submission Queue (SQ)
- Kernel posts completions into Completion Queue (CQ)



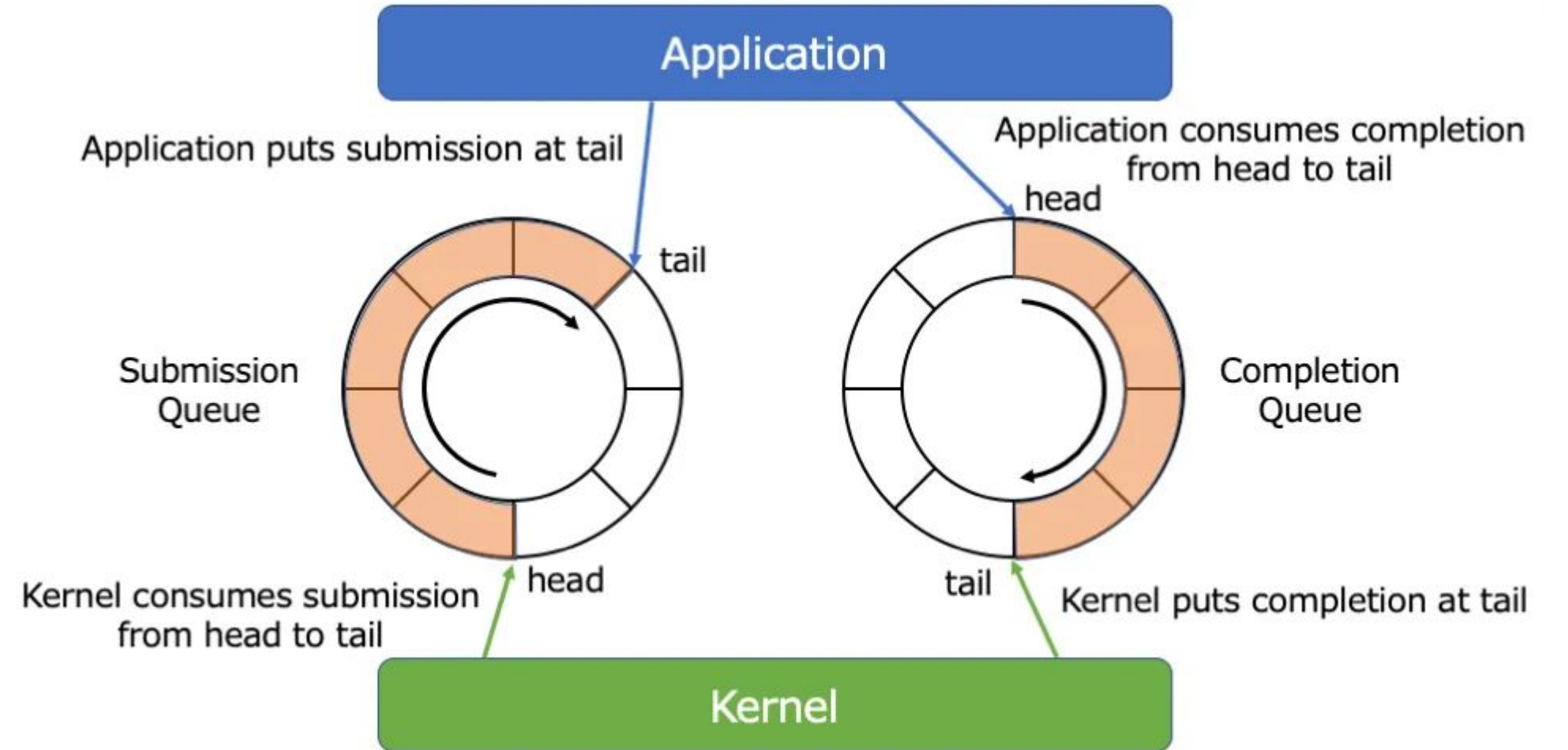
Prepare request

```
struct io_uring_sqe *sqe;  
sqe = io_uring_get_sqe(ring);  
io_uring_prep_recv(sqe, sockfd, buf, len, flags);  
Note this already moves the SQ tail
```



Submit

```
io_uring_submit_and_wait(ring, nr_completions);
```



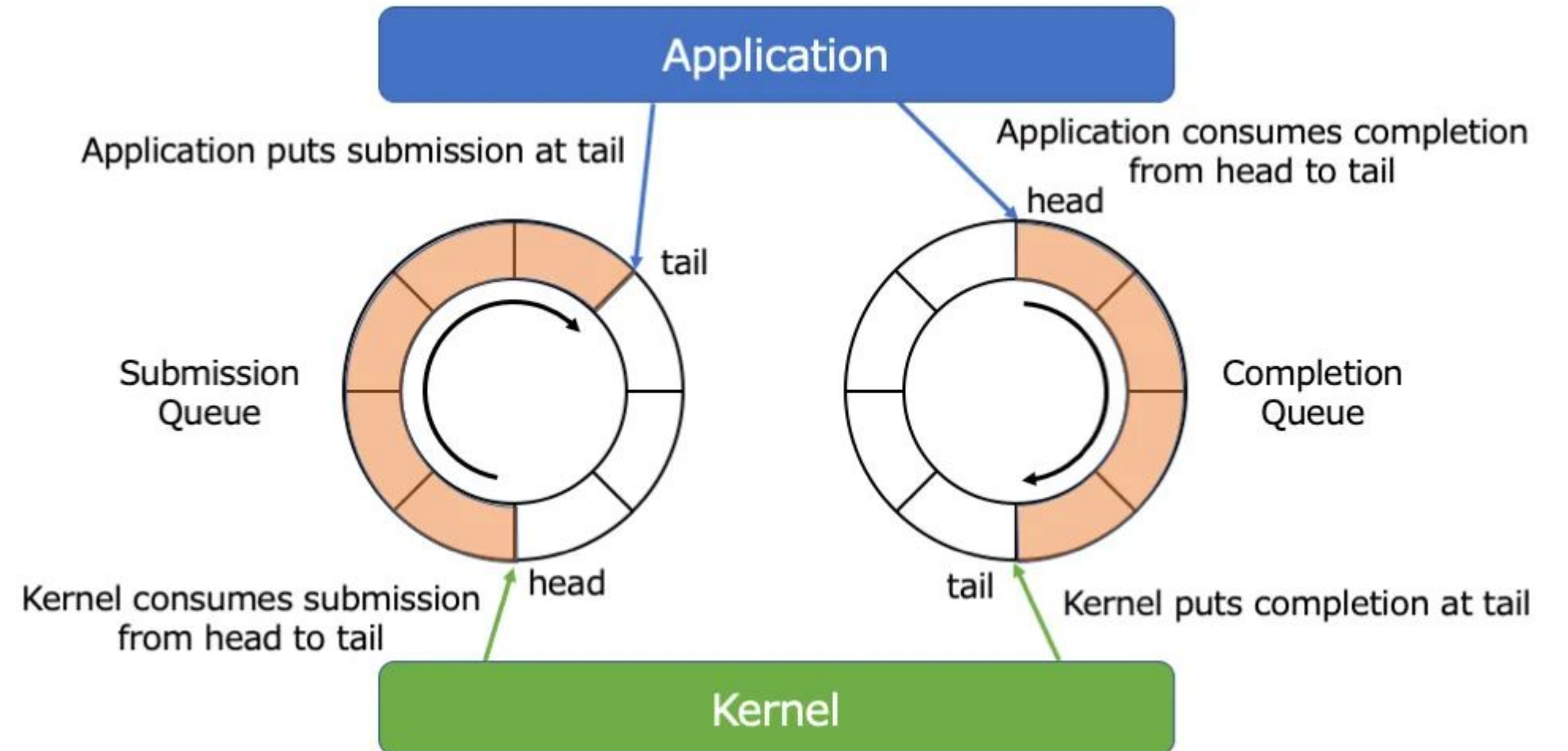
Process completions

```

unsigned head;
int count = 0;
io_uring_for_each_cqe(ring, head, cqe) {
    // do stuff
    count++;
}

io_uring_cq_advance(ring, count);

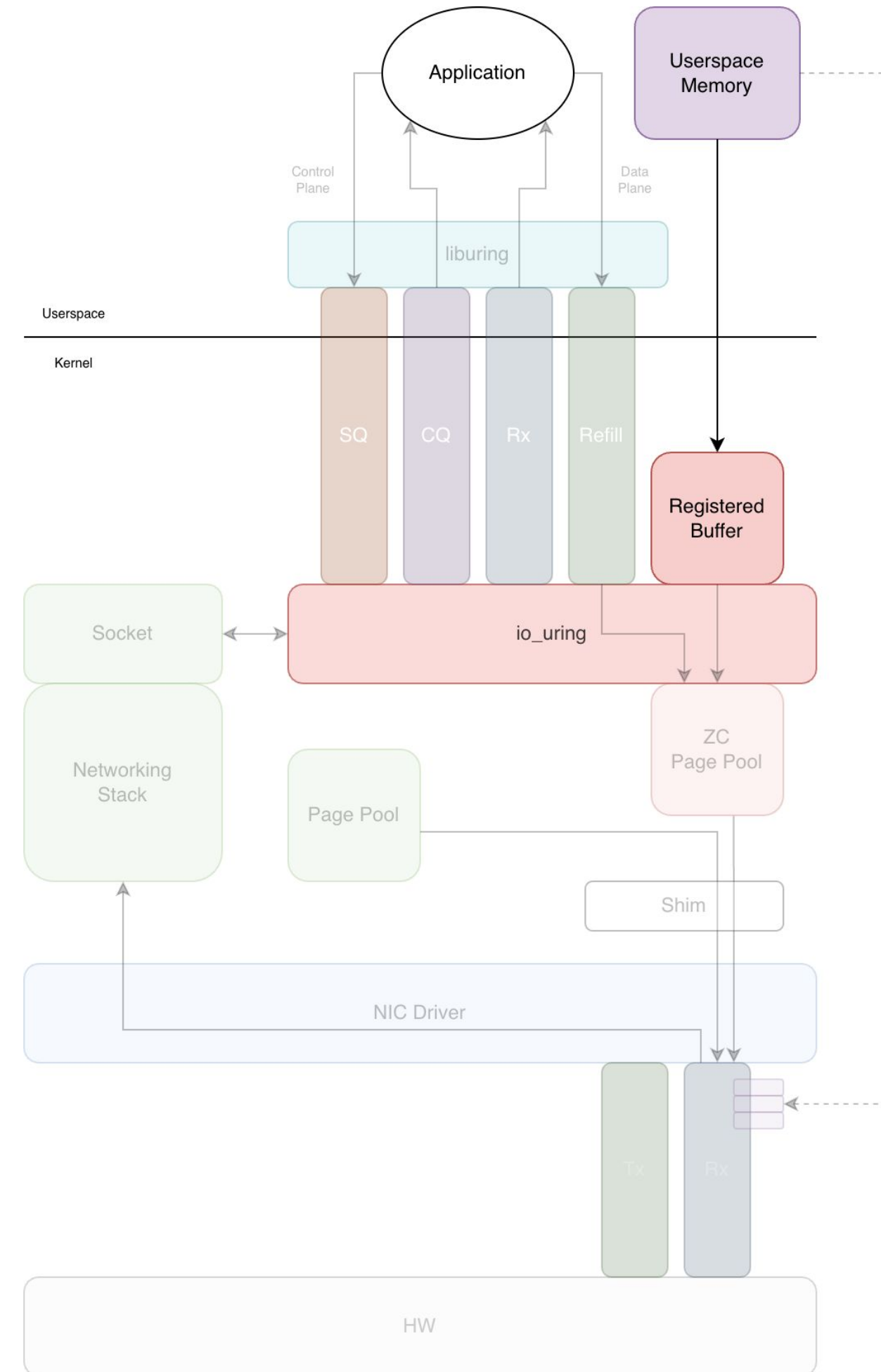
```



03 Design

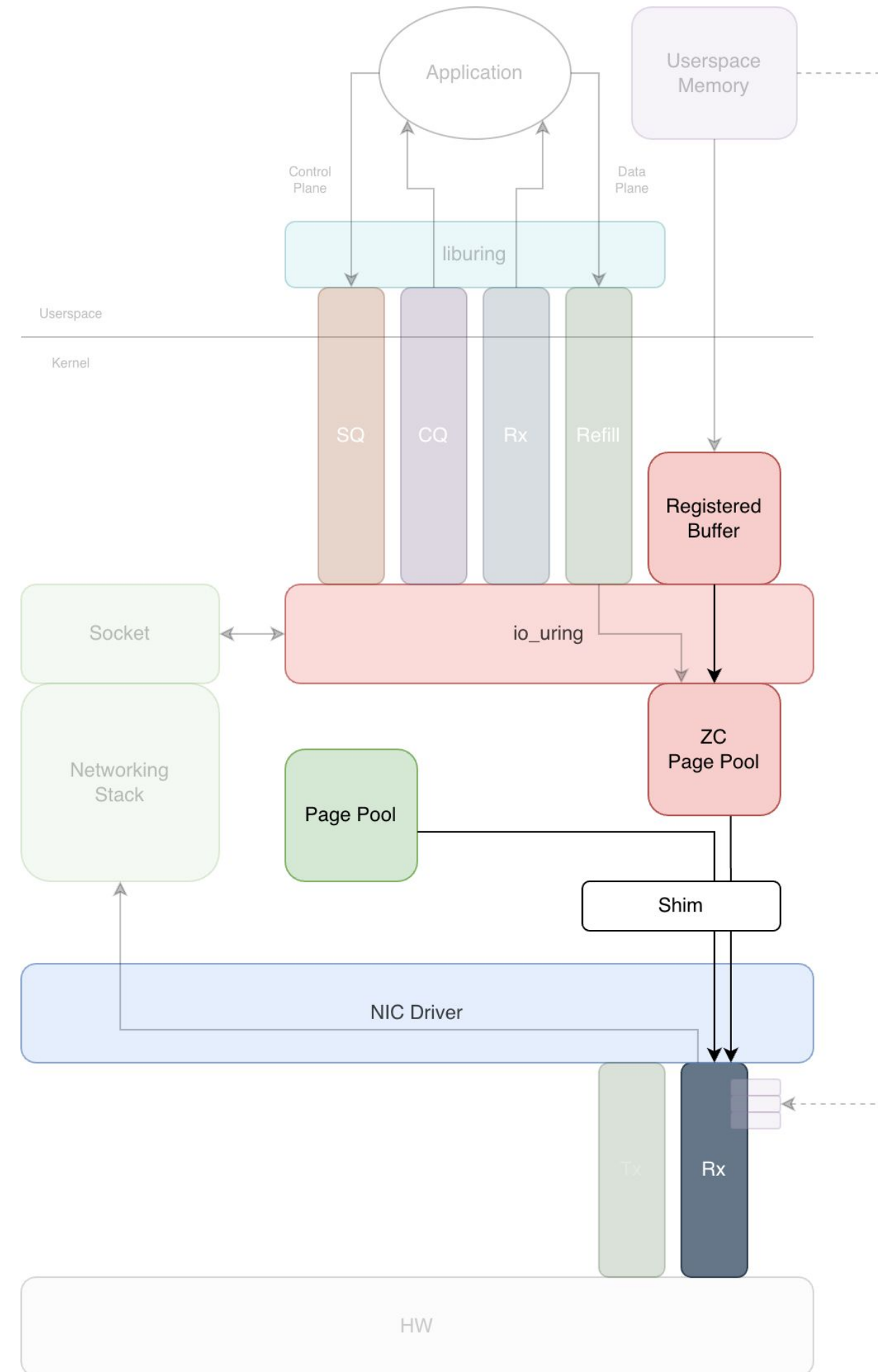
Buffer management

- Register userspace memory with io_uring
- Pin pages
- `struct bio_vec bvec[]`



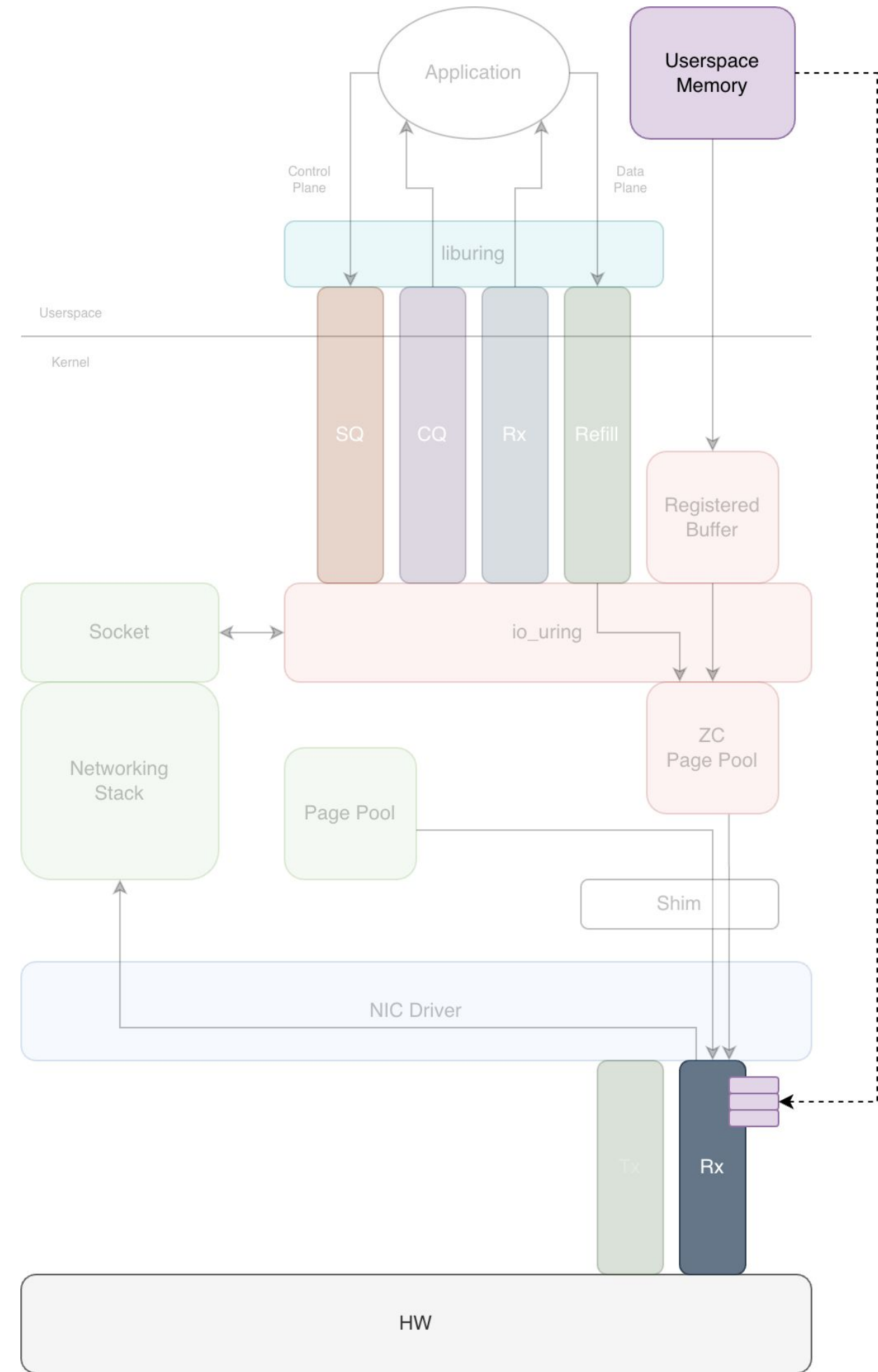
Buffer management

- ZC page pool “inspired” by page pool
- Thin shim layer + driver changes



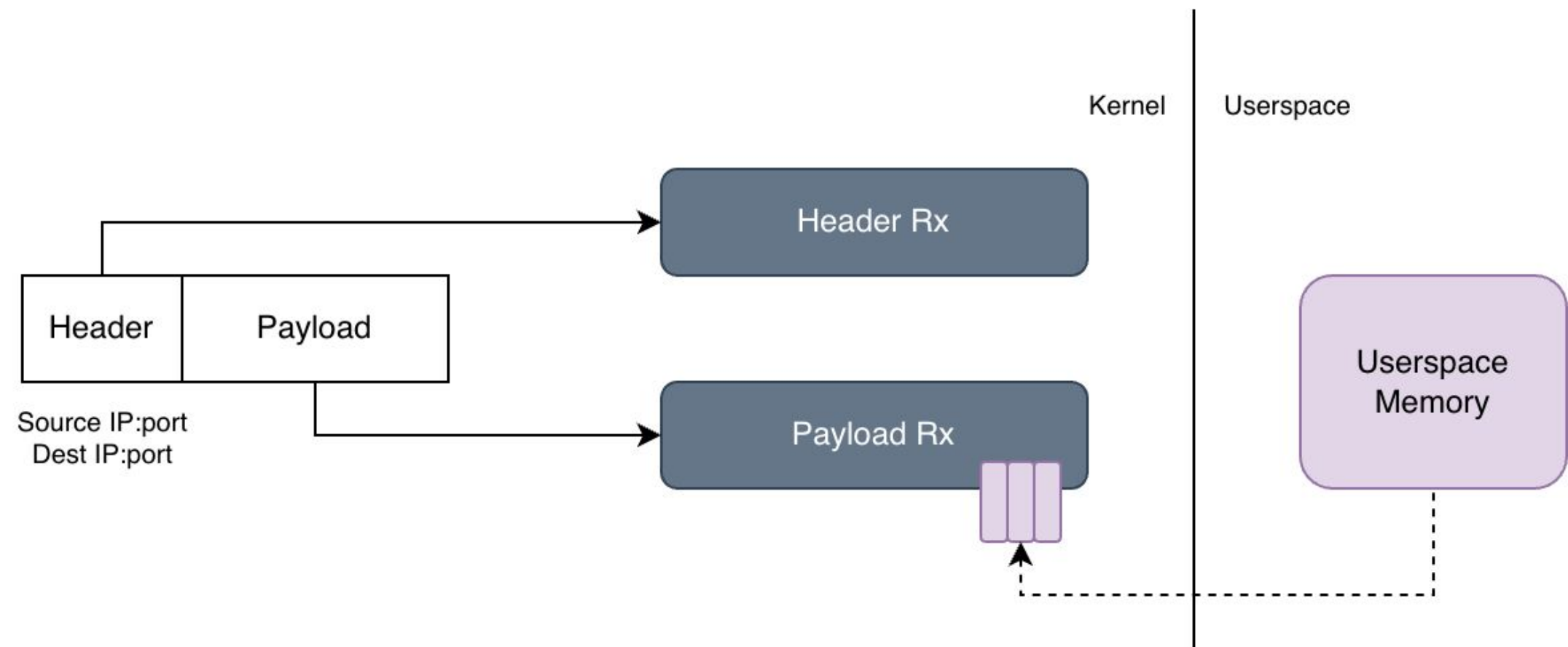
Buffer management

- End result: hardware Rx queue filled with userspace pages



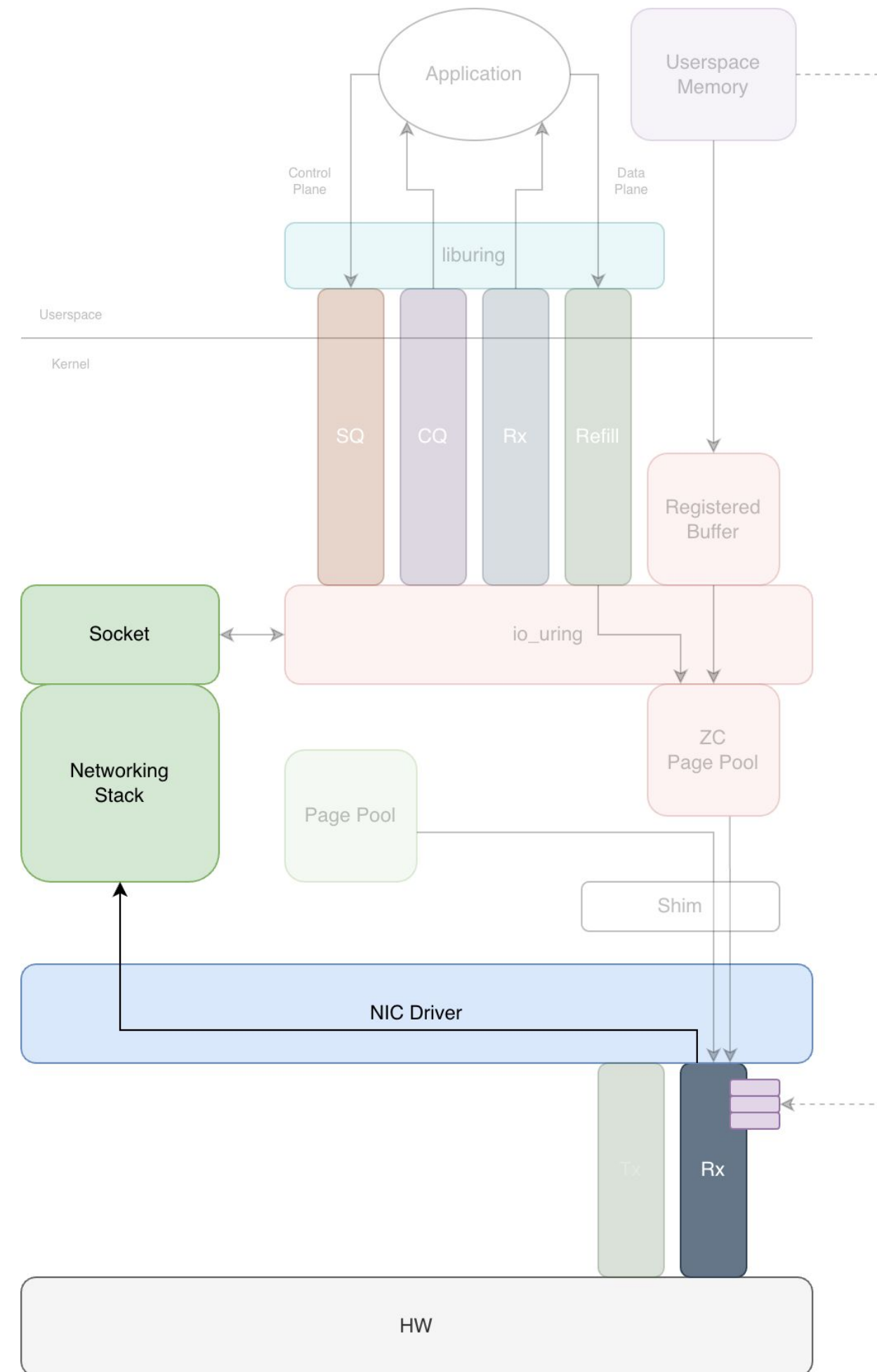
Header splitting + flow steering

- Only want payload
- Header splitting
- Only want our specific application flows to hit our ZC hardware Rx queues
- Flow steering
- RSS



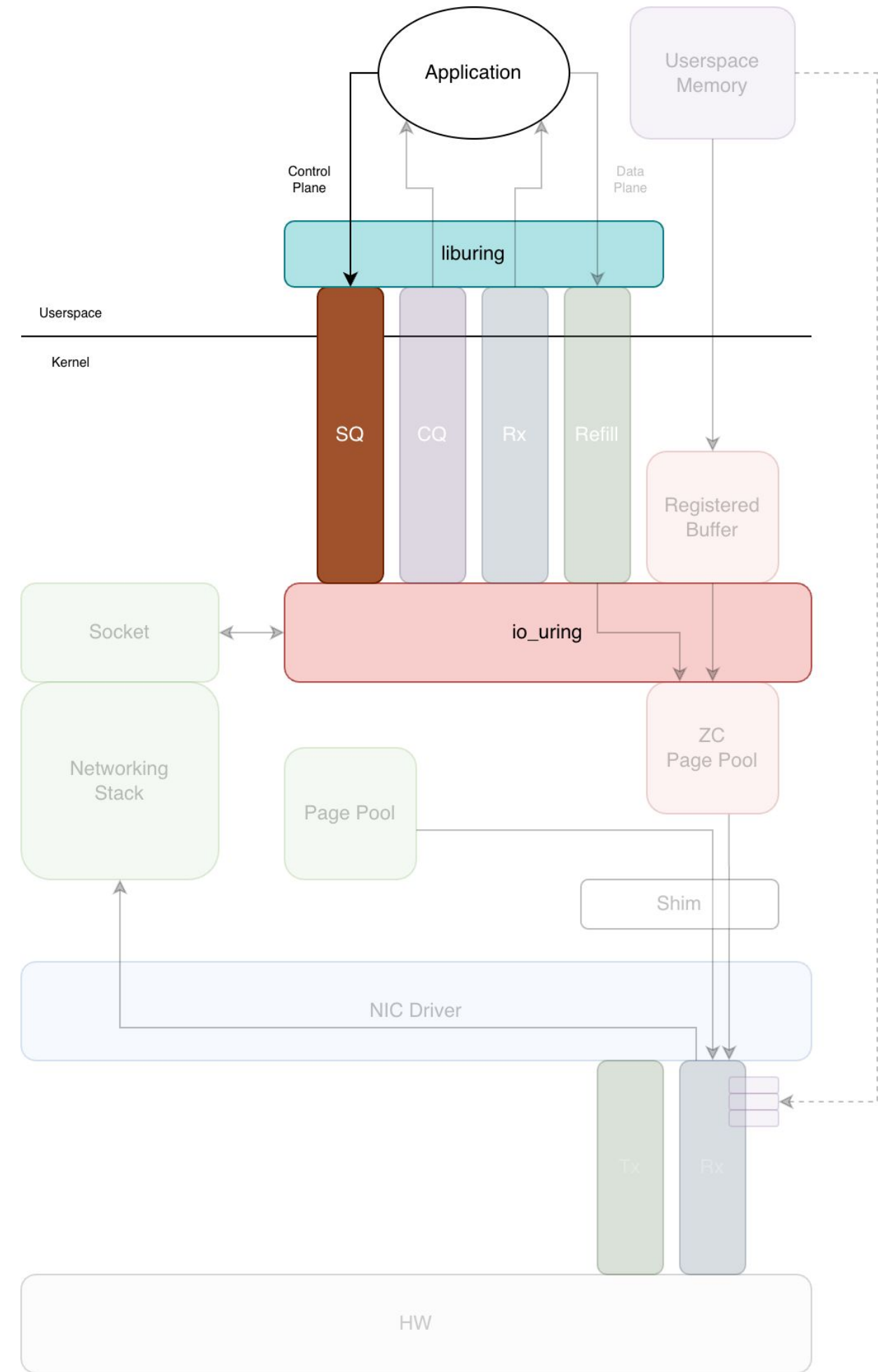
Kernel network stack

- Hardware side fully set up
- Hard IRQs
- NAPI poll
- Construct sk_buffs
 - Marked as ZC Rx
 - Page frags → userspace pages
- Goes through networking stack



Userspace: control plane

- Submit ZC receive request to io_uring

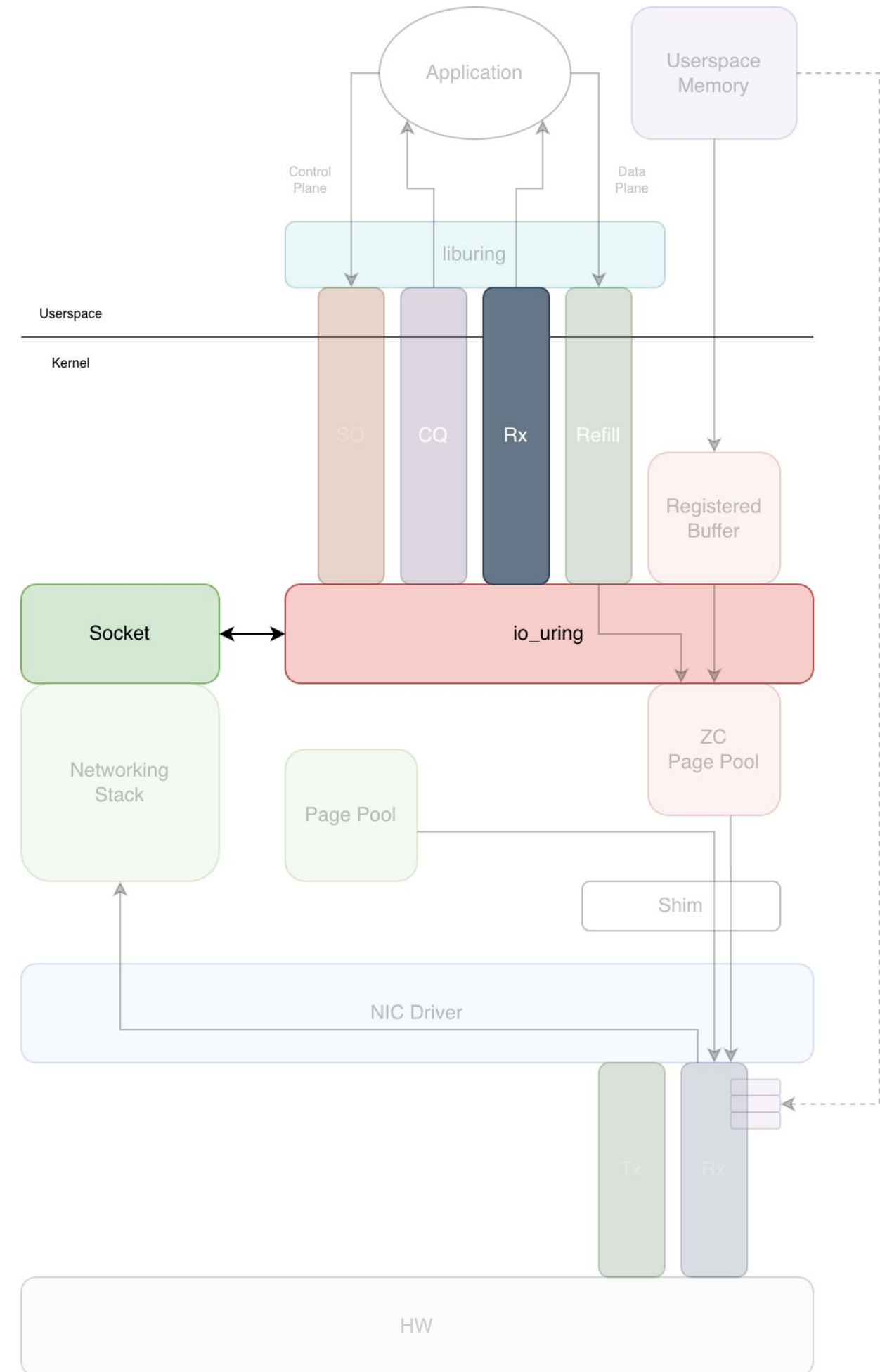


Userspace: control plane

- Handle ZC receive request
- Read sk_buffs from socket
- No copy - payload already in userspace
- Post one ZC Rx queue entry per skb page frag

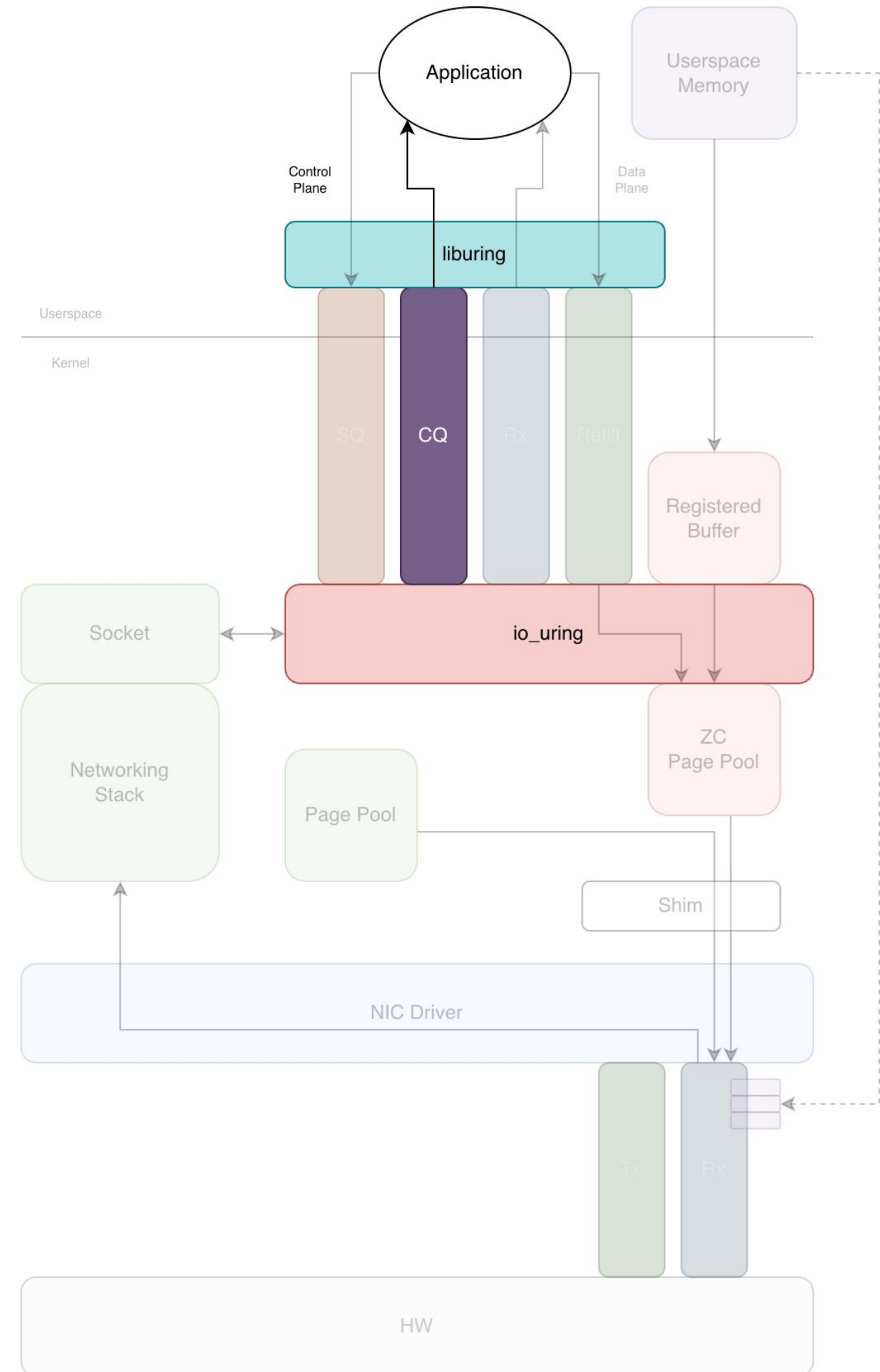
```

struct io_uring_rbuf_cqe {
    u32 off;
    u32 len;
    u16 region;
    u8 sock;
    u8 flags;
}
    
```



Userspace: control plane

- Post completion event into CQ
- Tells userspace to go look at which ZC Rx queue

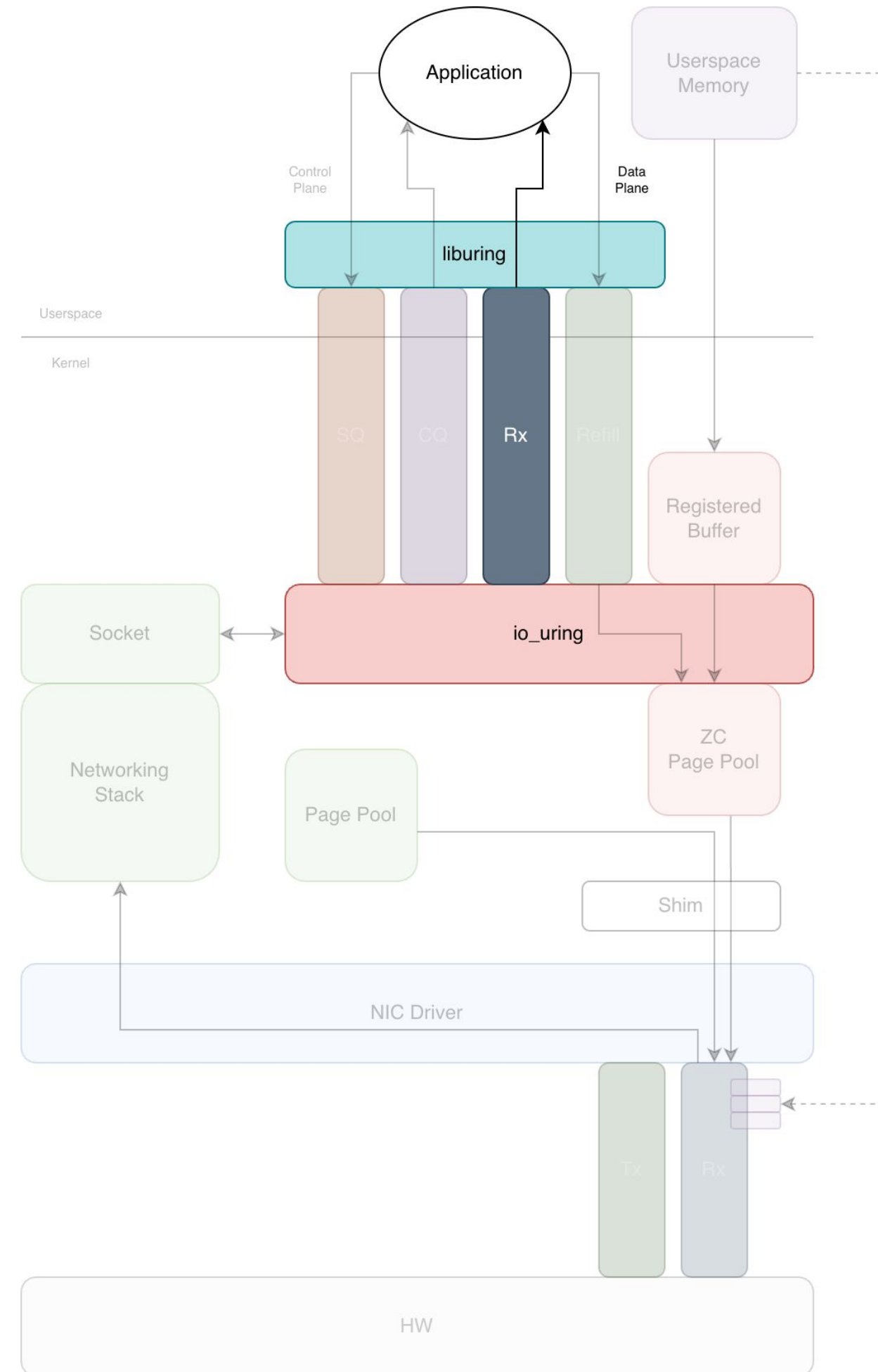


Userspace: data plane

- Look at a ZC Rx queue
- Each entry tells user where the payload is relative to the registered memory region

```

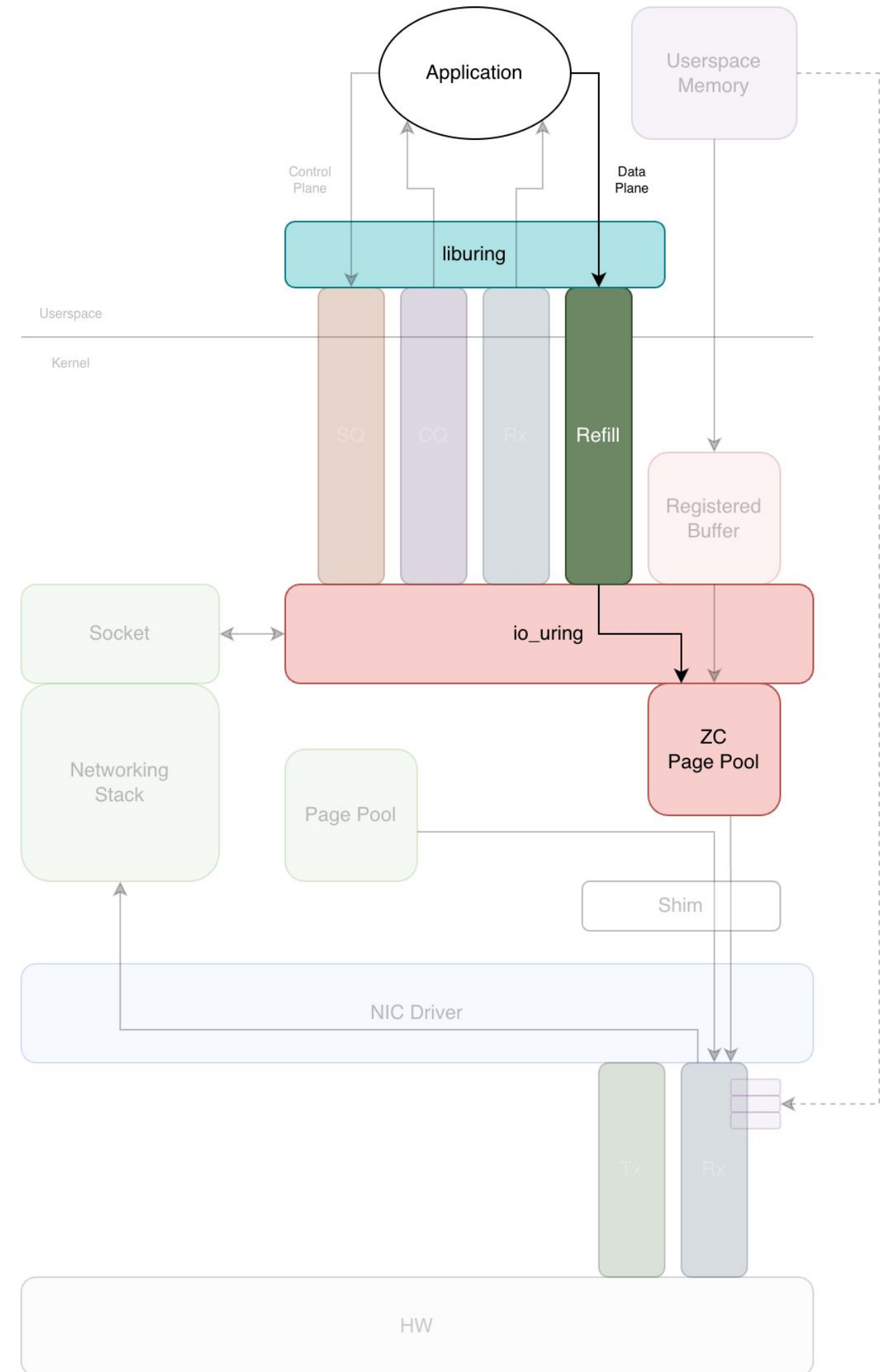
struct io_uring_rbuf_cqe {
    u32 off;
    u32 len;
    u16 region;
    u8 sock;
    u8 flags;
}
    
```



Userspace: data plane

- Return buffers to ZC page pool via refill queue
- Eventually used by NIC driver to refill hardware Rx queue

```
struct io_uring_rbuf_rqe {
    u32 off;
    u32 len;
    u16 region;
}
```



04 Preliminary Results

MemBW

Broadcom BCM57504 NIC @ 25 Gbps link

62 GB DRAM

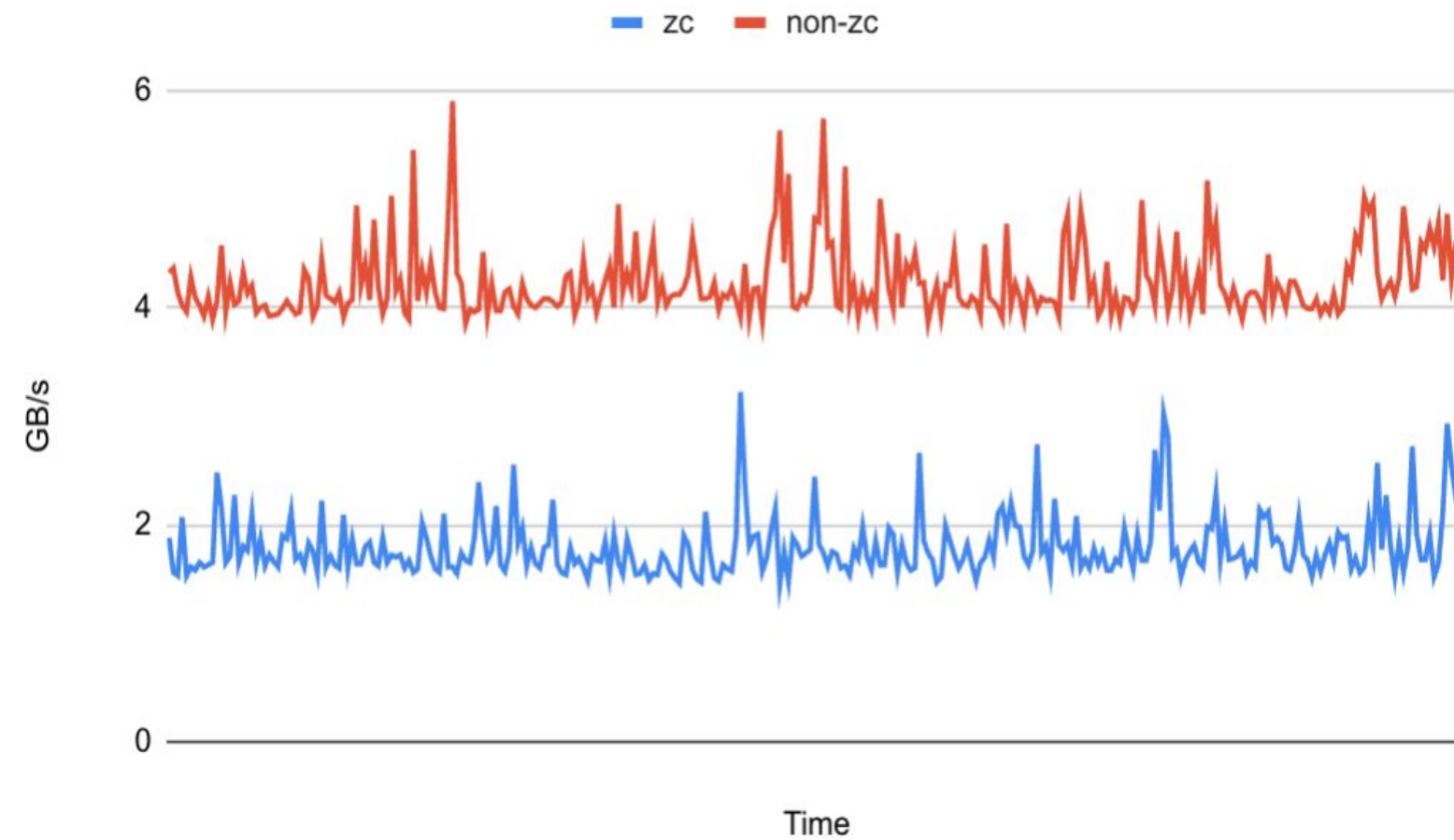
iperf3 + io_uring + ZC Rx

AMD EPYC 7D13

iperf3

uProf

System Memory Read Bandwidth



MemBW

Broadcom BCM57504 NIC @ 25 Gbps link

62 GB DRAM

iperf3 + io_uring + ZC Rx

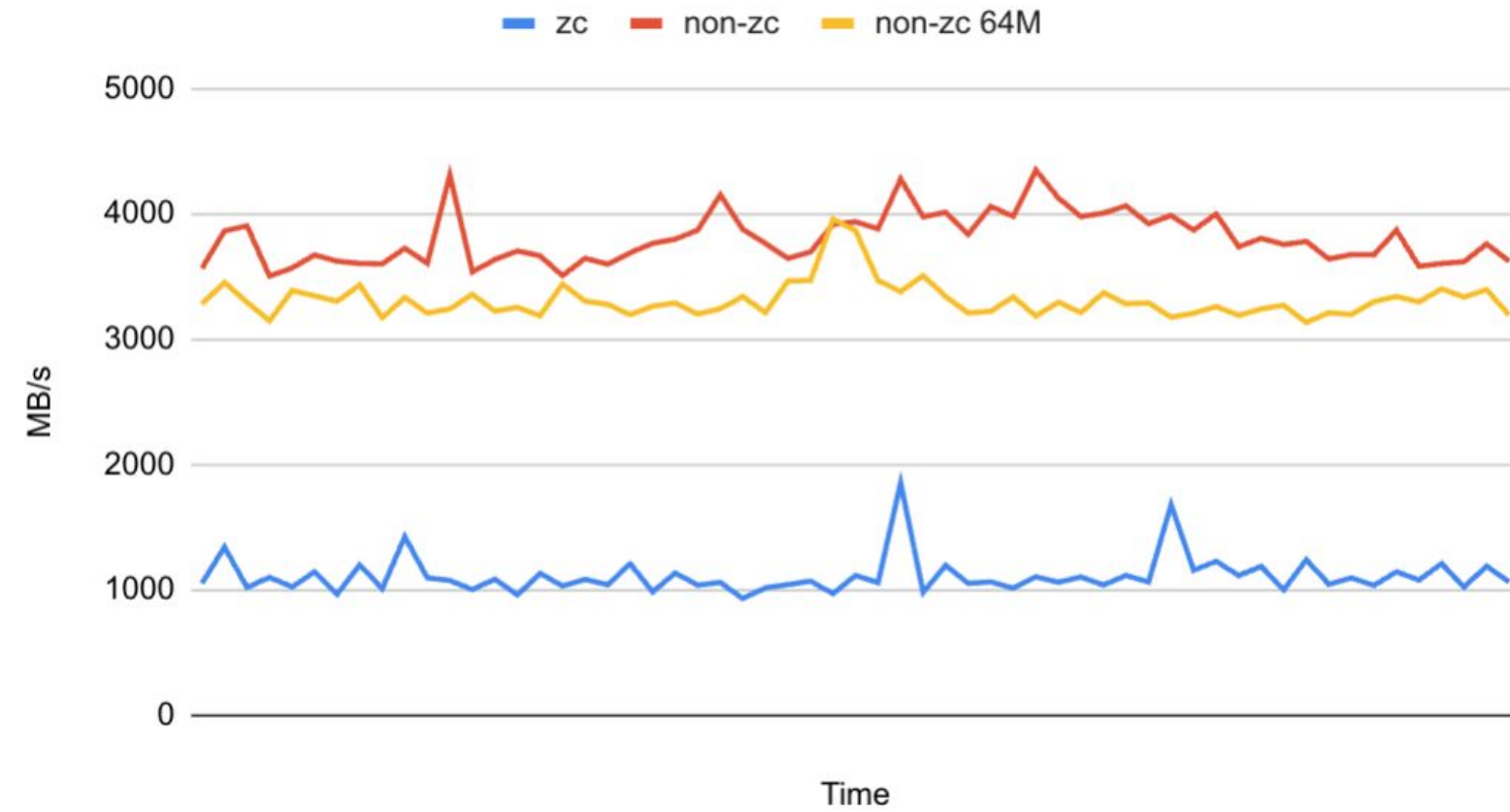
Intel Xeon Platinum 8321HC

iperf3

pcm-memory

DDIO is off

System Memory Read Bandwidth



05 Questions?

 Meta

06 Discussion

Handling errors

- How much to allocate ahead of time?
- What if it runs out?
- What if header splitting fails?
 - Split too little - header malformed
 - Split too much - payload included
- What if flow steering fails?
 - ZC Rx packet ends up in non-ZC Rx queue
 - Non-ZC Rx packet ends up in ZC Rx queue

Copy fallback

- What if we run out of userspace memory allocated for ZC Rx?
- Fill HW Rx queue with kernel pages - as before
- When io_uring ZC receive finds sk_buffs with page frags that are not ZC pages, copy into a page from refill queue
- Turn OFF ZC Rx! Then tell application
- Application must fix the problem then kick ZC Rx back on

Integrating ZC Rx well

- NIC → userspace memory is only one hop in a long end to end pipeline
- What if data needs to be modified after ZC Rx? Another copy...
- API need to expose fine control over the placement of data to satisfy constraints e.g. alignment
 - Hardware also needs to support this too
- TLS and kTLS?