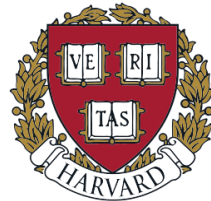


# eTran: Extensible Kernel Transport with eBPF

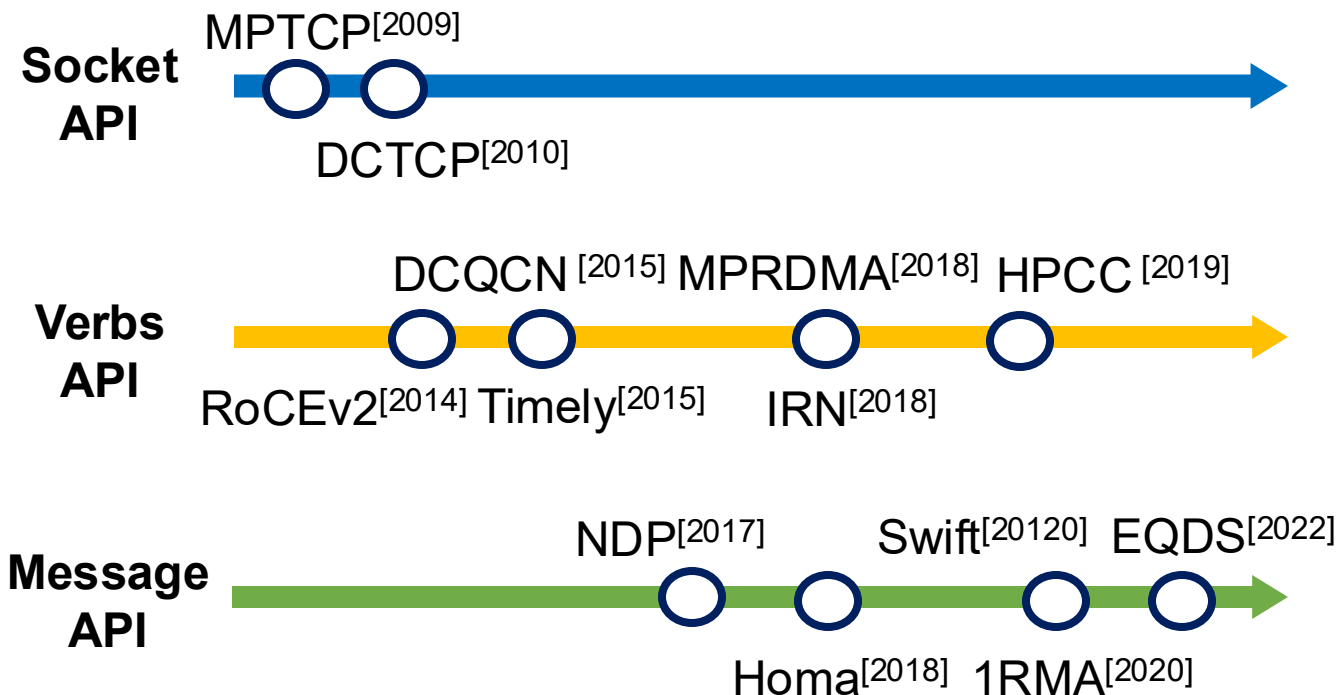
**Zhongjie Chen**<sup>1</sup> Qingkai Meng<sup>2</sup> ChonLam Lao<sup>3</sup> Yifan Liu<sup>1</sup>  
Fengyuan Ren<sup>1</sup> Minlan Yu<sup>3</sup> Yang Zhou<sup>4</sup>

<sup>1</sup>*Tsinghua University* <sup>2</sup>*Nanjing University*  
<sup>3</sup>*Harvard University* <sup>4</sup>*UC Berkeley & UC Davis*



# Evolving DC transports over the years

- Datacenter applications are evolving and increasingly diverse.
  - microservice, storage, etc.
- There is no *one-size-fits-all* transport for all workloads.



## Semantics

- stream, message, connection, connectionless

## Congestion control

- sender-driven, receiver-driven, switch-driven

## Loss recovery

- Go-Back-N, SACK, packet trimming

# However.....

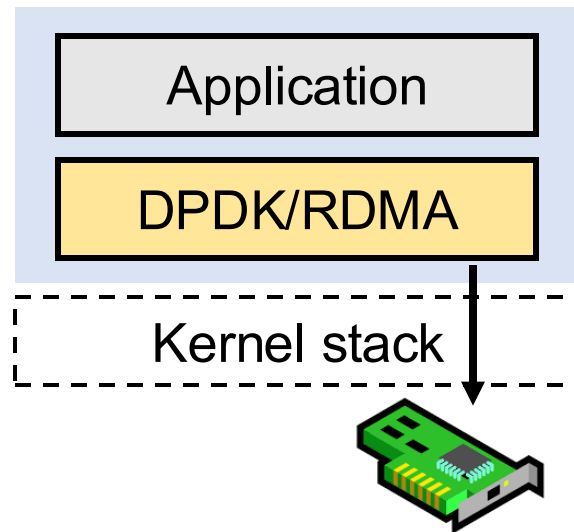
Transport	Proposal	→	Landing kernel
Softwarp	2009	~9 years	2019
SoftRoCE	2014	~2 years	2016
Homa	2021	—	—

And more transports stay at the kernel-bypass prototype stage.

It takes a long time to **extend**, **customize** and **evolve** kernel transport.

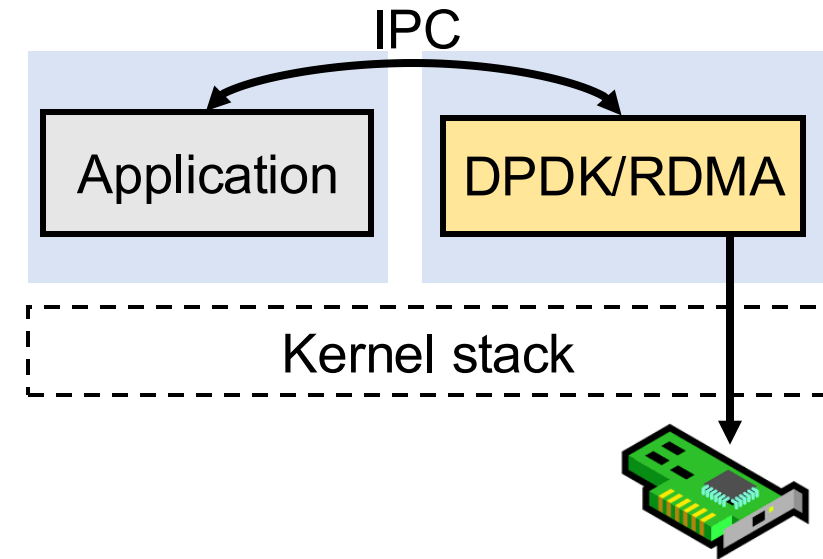
# Why don't we focus on kernel-bypass?

**(a) LibraryOS**



- Direct HW access, security concern
- Dedicated resources

**(b) Microkernel**



- Moderate performance
- Implementation from scratch

# Why is it hard to extend kernel transport?



*How to achieve **agile customization**, **kernel safety**, **strong protection**, and **high performance** for kernel transport ?*

Heavy development and  
maintenance effort

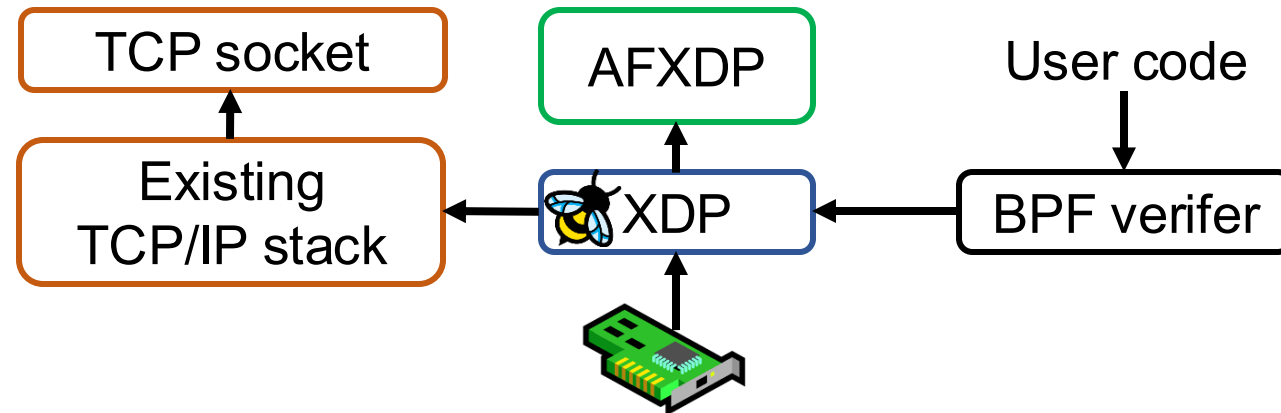
Unsatisfactory  
performance

Long time to land on  
kernel mainline

But it does provide **safety, multi-tenancy, out-of-the-box deployment.....**

# **eBPF** to the rescue

- eBPF (extended Berkeley Packet Filter)
  - ✓ **Safely** run programs in kernel at runtime guaranteed by BPF verifier
  - ✓ With an active community, it has impacted many OS subsystems
- XDP (eXpress Data Path) and AF\_XDP
  - ✓ Fast kernel packet processing
  - ✓ Fast userspace packet processing based on XDP



# Challenges in implementing kernel transport with eBPF/XDP

**C1:** Full transport logic (e.g., TCP) is too **complex** for the eBPF programming model, even though eBPF is increasingly more powerful

- e.g., memory allocator, kfuncs, dynptr, rb-tree, etc

**C2:** **Limited** operating points with eBPF/XDP

No queue/buffer

Existing TCP/IP stack

No packet sending

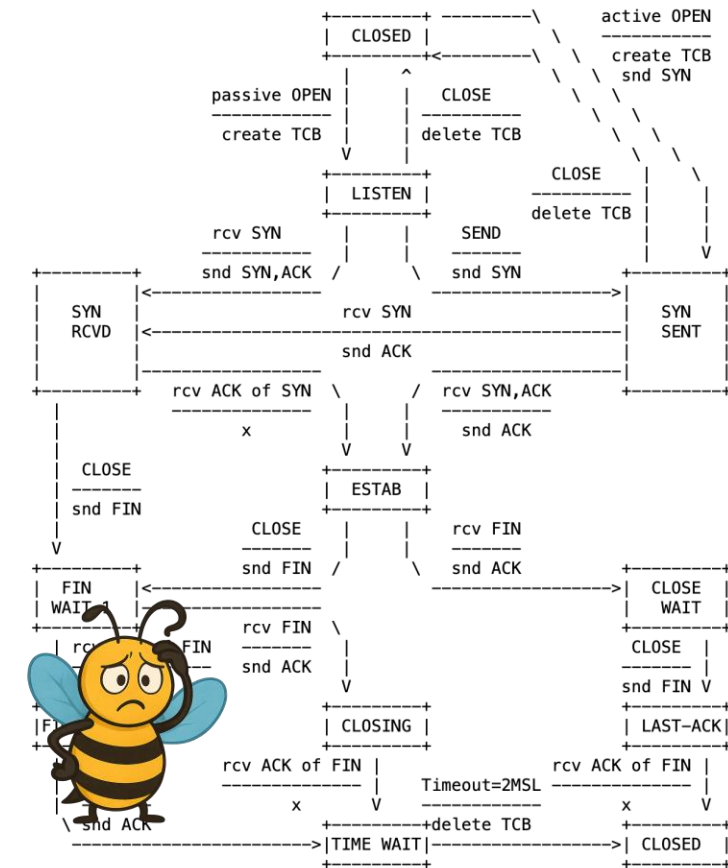
AF\_XDP

XDP

RX

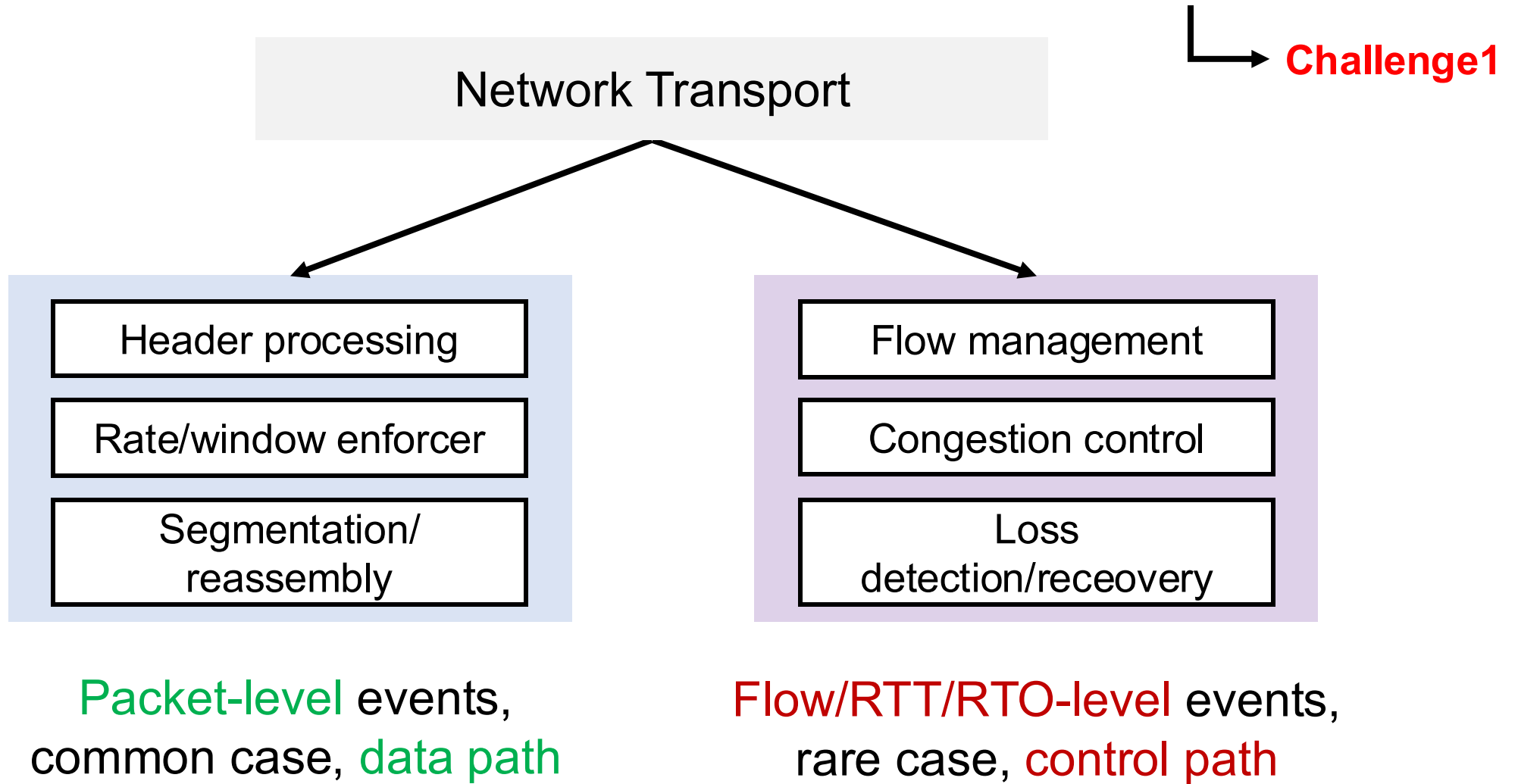
TX

Ingress traffic only



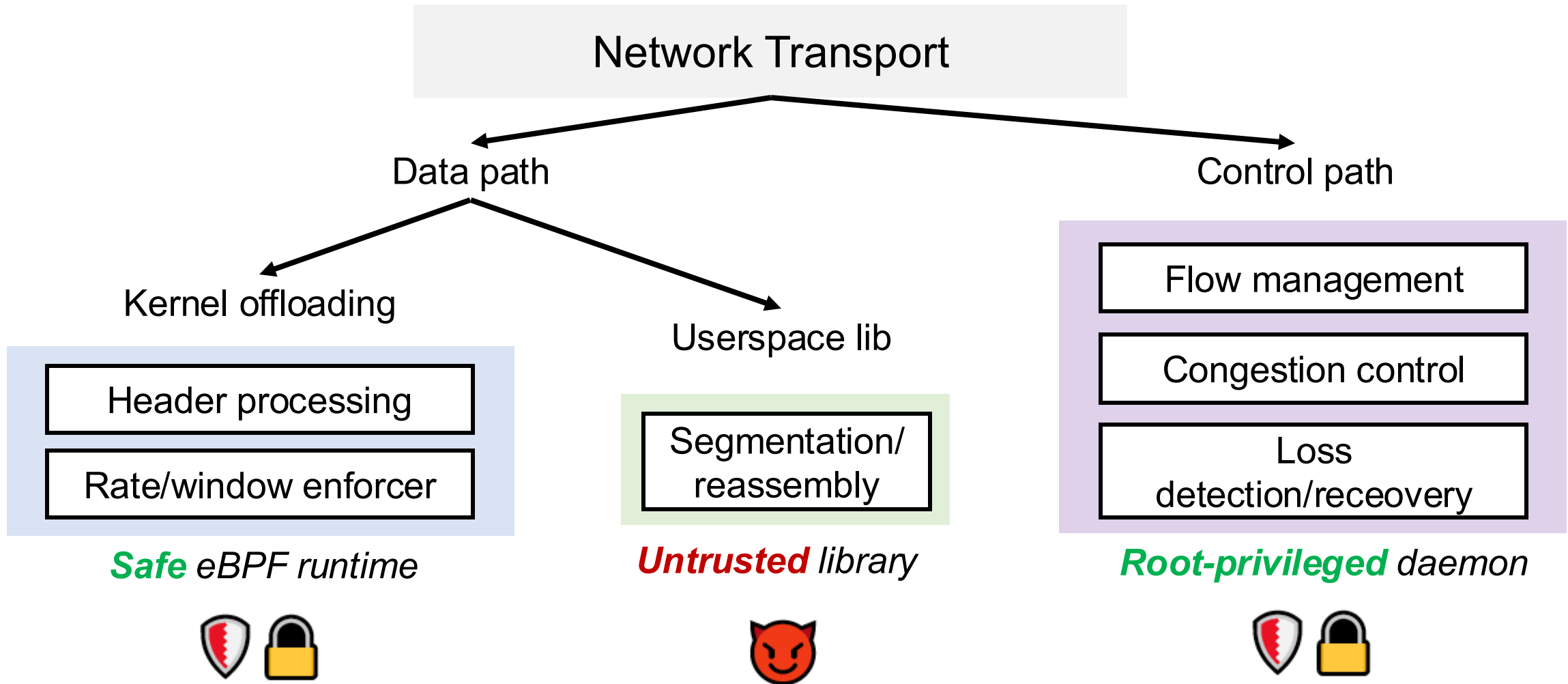
RFC793: TCP state machine

# eTran design I: split transport paths





# eTran design I: split transport paths



# eTran design II: extend eBPF subsystem

└─ Challenge2

## New XDP hooks

### XDP\_EGRESS

- Intercept egress traffic

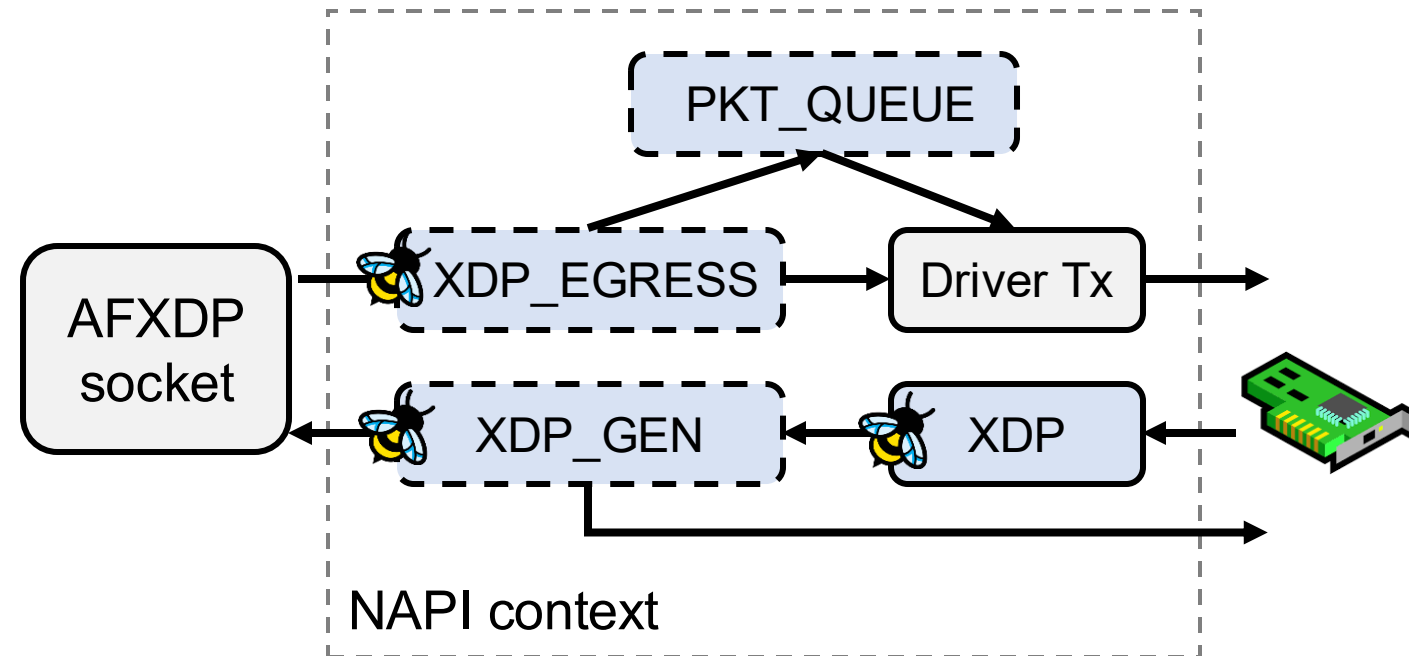
### XDP\_GEN

- Send packets proactively

## New eBPF map

### BPF\_MAP\_TYPE\_PKT\_QUEUE

- buffer/queue XDP frames

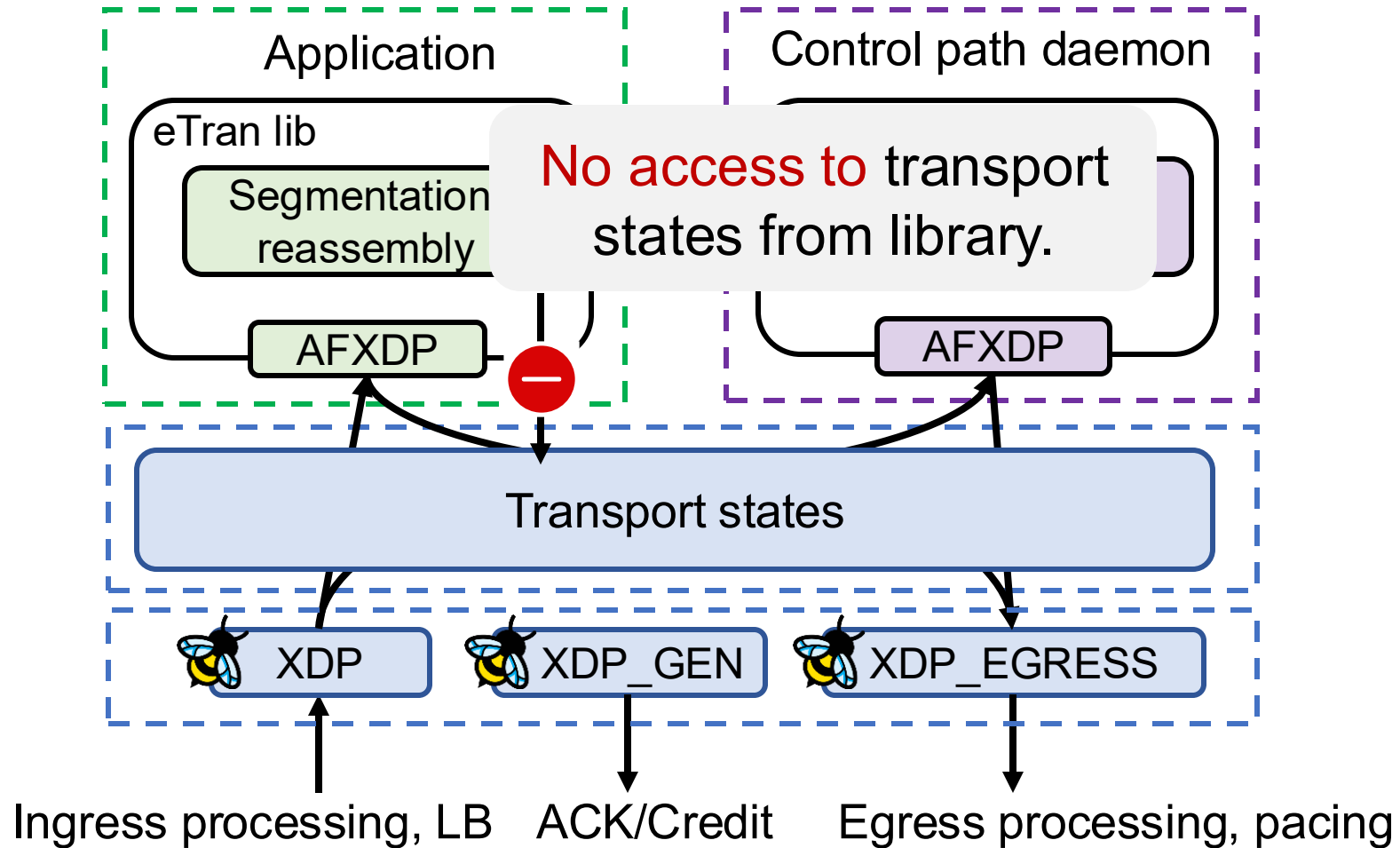


A bunch of optimizations for high performance.

- pre-allocated memory pool, batching, etc.

→ See the paper

# eTran overall architecture



# Extend eBPF subsystem: new hooks

```
int tcp(struct xdp_md *ctx) { // TCP logic...}
int homa(struct xdp_md *ctx){ // Homa logic...}
SEC('XDP_EGRESS')
int xdp_egress_prog(struct xdp_md *ctx)
{
    int proto;
    struct hdr_cursor nh;
    void *data=(void *) (long)ctx->data;
    void *data_end=(void *) (long)ctx->data_end;
    nh.pos=data;
    proto=parse_ethhdr(&nh, data_end, &eth);
    if (proto!=bpf_htons(ETH_P_IP))
        return XDP_DROP;
    proto=parse_iphdr(&nh, data_end, &iph);
    if (proto==IPPROTO_TCP || proto==IPPROTO_HOMA)
        bpf_tail_call(ctx, 0);
    return XDP_DROP;
}
SEC('XDP_GEN')
int xdp_gen_prog(struct xdp_md *ctx)
{
    prepare_ack_pkt(ctx);
    return XDP_TX;
}
```

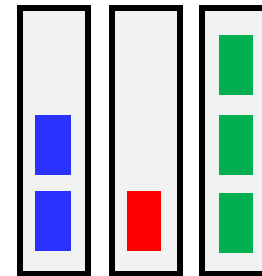
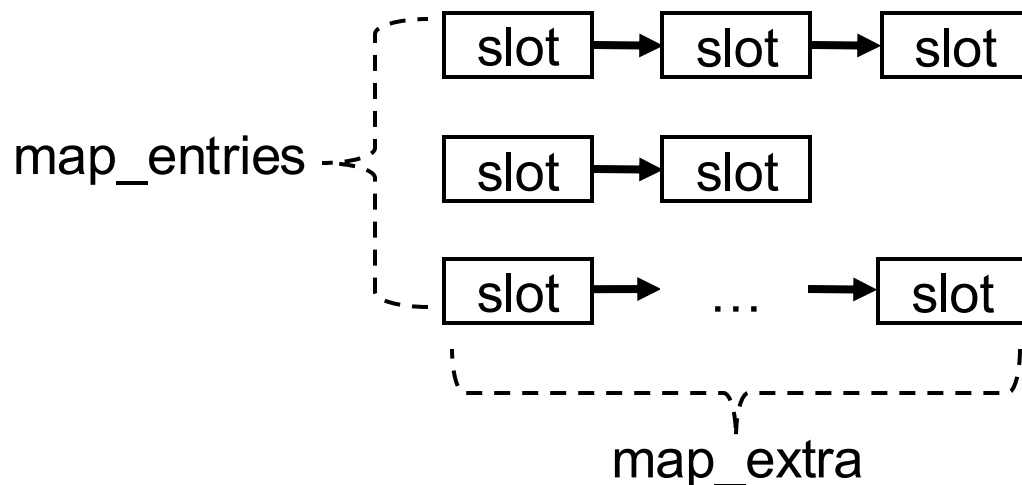
- Comply with all eBPF/XDP program constraints
  - # of instructions, bounded loop, etc
- Reuse most existing XDP infrastructures
  - Similar data structures and return codes
  - Most BPF helper functions for XDP

# Extend eBPF subsystem: new BPF map

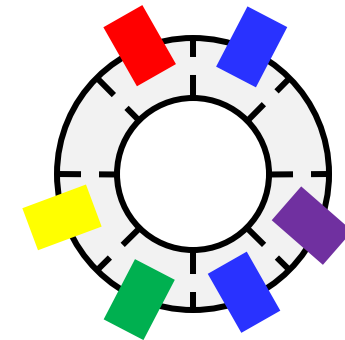
```
struct {  
  __uint(type, BPF_MAP_TYPE_PKT_QUEUE);  
  __type(key, __u32);  
  __uint(max_entries, MAX_BKT);  
  __uint(map_extra, MAX_PKT_PER_BKT);  
} xdp_pkt_queue SEC(''.maps'')
```

- Only store **frame pointers**
- Memory size is **static**
- **kfuncs** for operations
- Cooperate with **BPF timer**

See the paper



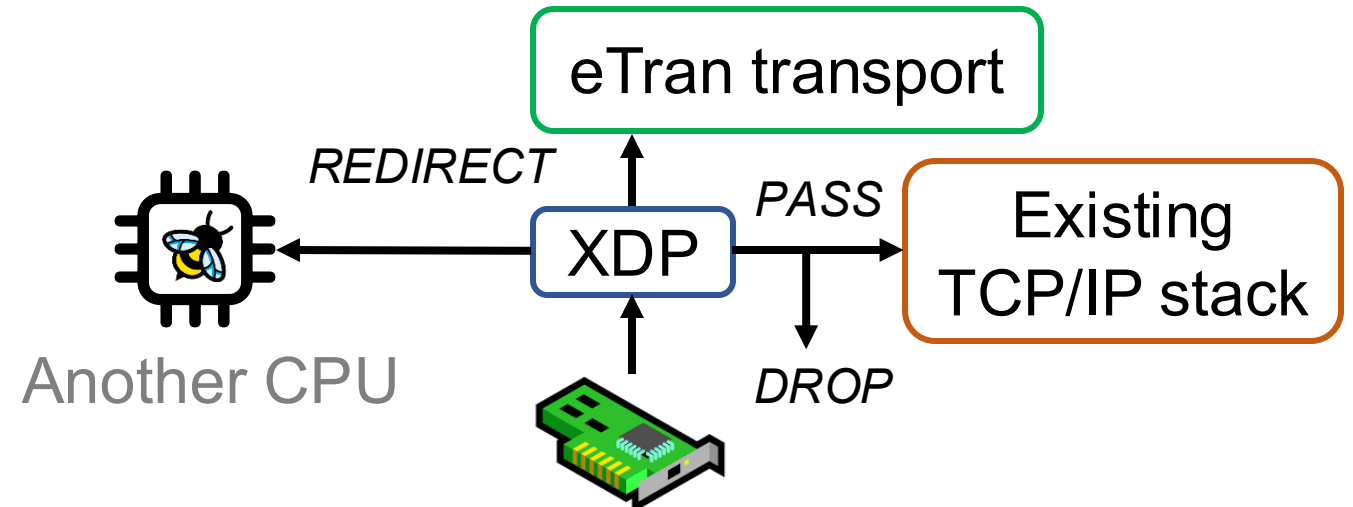
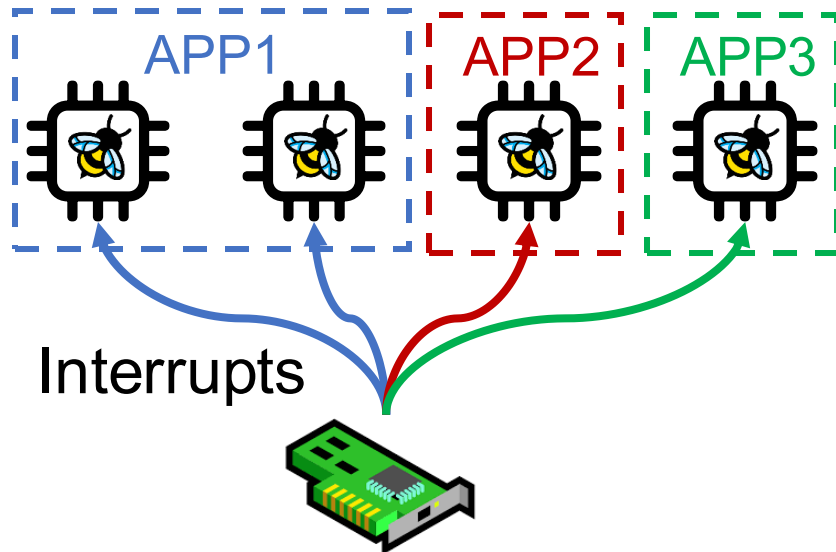
Per-flow queue



TimingWheel  
[SIGCOMM'17]

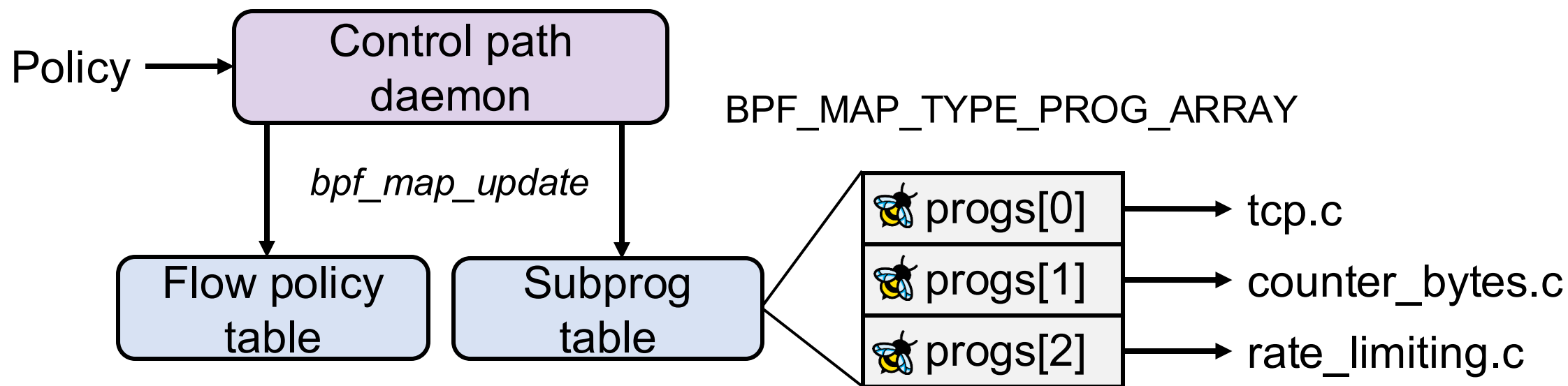
# More features: multi-tenancy

- Queue-level resource allocation by control path daemon
  - Isolate NAPI/IO buffers and avoid monopolizing the NIC (e.g., DPDK)
- Magic of XDP
  - Co-exist with other kernel transports flexibly



# More features: traffic management

- Call subprograms with *bpf\_tail\_call* feature.
  - Traffic monitoring, ACL, rate limiting, etc.



# More details of eTran

- Hook performance optimizations.
- Flow scheduling with BPF timer.
- AF\_XDP virtual sockets for multi-queue.
- Loss recovery under eTran.
- Receiver-driven CC under eTran.



Not covered in this talk



# Case studies & Implementation

## **DCTCP**<sup>[SIGCOMM'10]</sup>

- Connection-based
- Stream semantic
- Sender-driven CC
- POSIX API

## **Homa**<sup>[SIGCOMM'18]</sup>

- Connection-less
- Message semantic
- Receiver-driven CC
- RPC API

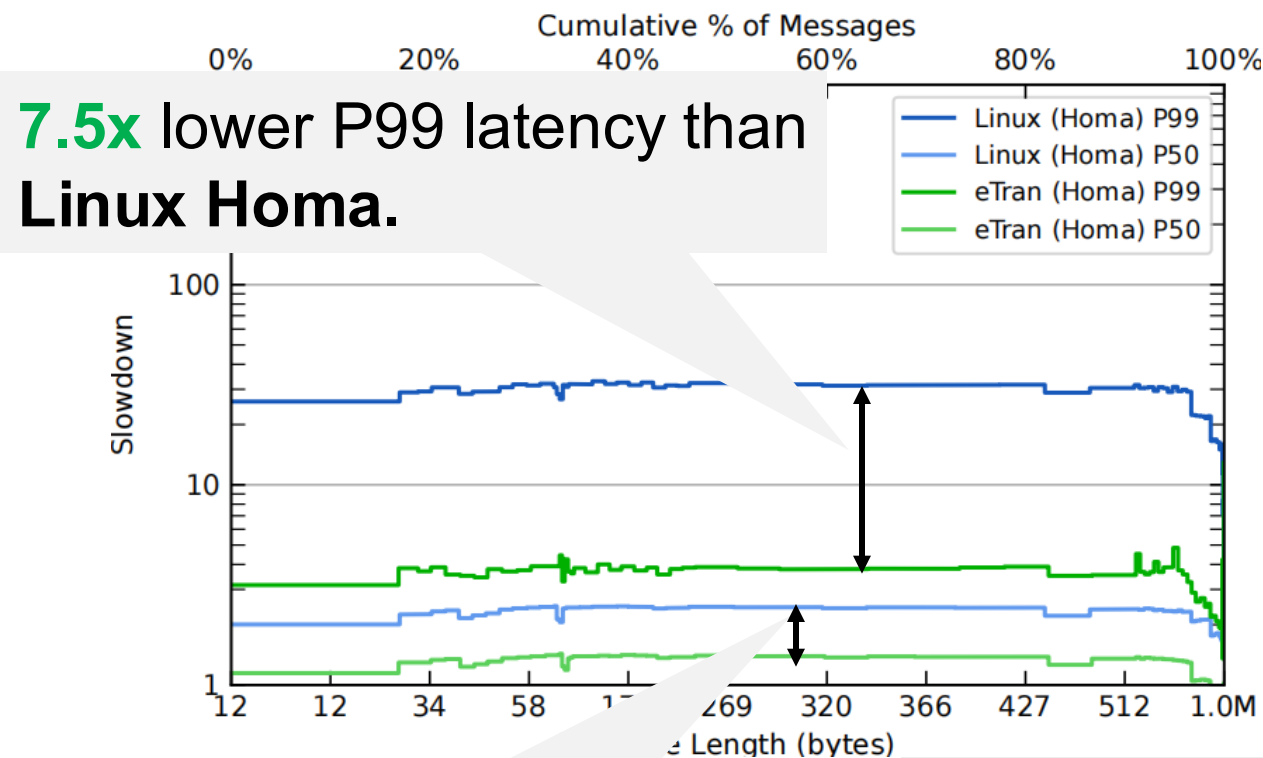
- Kernel modification
  - 2.6K LoC for extending eBPF subsystem
- eBPF code
  - 7.5K LoC eBPF for TCP and Homa
- Userspace code
  - 8K LoC for control path dameon and 5.6K for transport libs

# Evaluation

- Baseline
  - Linux Homa
  - Linux TCP, TAS<sup>[Eurosys'19]</sup>(kernel-bypass TCP/IP stack)
- Experiment setup
  - 10 Cloudfab xl170 machines running **customized** Linux kernel 6.6.0.
  - RPC workload for Homa
  - Key-value store application for TCP.
- Open source
  - <https://github.com/eTran-NSDI25/eTran>

# Linux Homa vs. eTran Homa

RPC workload



3.6x lower P50 latency than Linux Homa.

4.3/4.9x lower P50/P99 latency for large message workload .

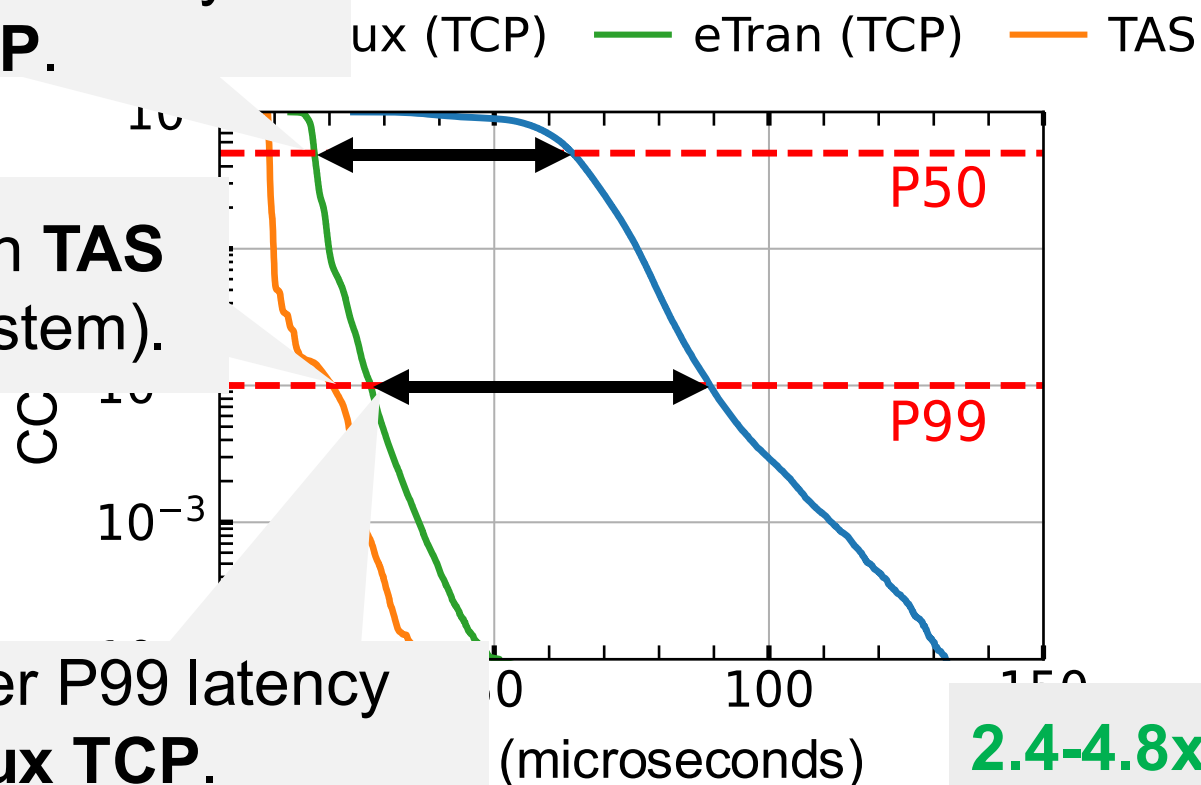
# Linux DCTCP, TAS vs. eTran DCTCP

Key-value store application

**3.7x** lower P50 latency than **Linux TCP**.

Slightly worse than **TAS** (kernel-bypass system).

**3.2x** lower P99 latency than **Linux TCP**.



**2.4-4.8x** higher throughput than **Linux TCP**.

# eTran summary

- eTran is an extensible kernel transport system achieving
  - ✓ Agile customization
  - ✓ Kernel safety
  - ✓ Strong protection
  - ✓ High performance
- eTran achieves these goals by
  - Extending the kernel-safe eBPF
  - Hiding transport states inside the kernel
  - Absorbing techniques from user-space transports



Thank you!