

Selective Packet Capture at High Speed Rates

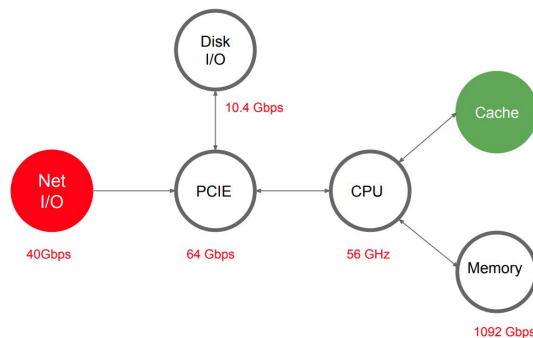
Reservoir Labs

Peter Cullen, James Ezick, Kelly Fox, Troy Hanson, Richard Lethin, Erik Mogus, Jordi Ros-Giralt, Alison Ryan,
{cullen, ezick, fox, hanson, lethin, mogus, giralt, ryan}@reservoir.com

Presenter: Jordi Ros-Giralt | giralt@reservoir.com

2nd European Zeek (Bro) Workshop

April 10, 2019



Reservoir Labs

632 Broadway
Suite 803
New York, NY 10012

- Selective Packet Capture: Problem definition
- Optimizations
 - Long queue emulation
 - Lockless bimodal queues
 - Tail early dropping
 - LFN tables
 - Multiresolution priority queues
- Zeek script

Packet Capturing at Very High Speed Rates

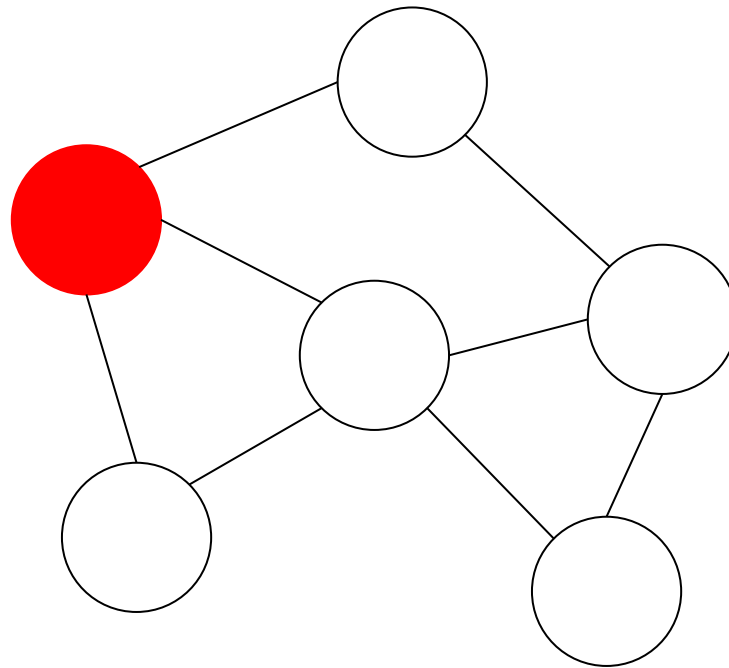
- Scalability issue: performing packet capture is either intractable or requires highly expensive hardware both in processing and storage.
- Liability issue: indiscriminate packet capture poses a liability issue.
- Selective Packet Capture (SPC) provides a sweet-spot solution to both of these problems.
- SPC gets a "free lunch" by leveraging all the heavy lifting work done by Zeek

Capturing packets at very high speed rate is an HPC problem...

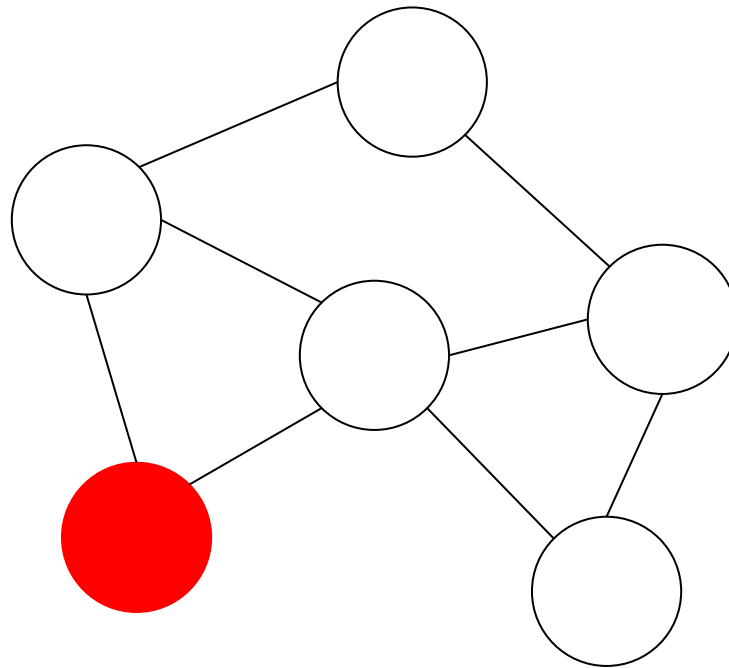
So let's talk first about performance optimization...

- System wide performance optimization of network components like routers, firewalls, or network analyzers such as a Zeek sensor is complex.
- Hundreds of different SW algorithms and data structures interrelated in subtle ways.
- Two interdependent problems:
 - Shifting micro-bottlenecks
 - Nonlinear performance collapse
- Special focus on the problem of packet capturing at very high speed rates

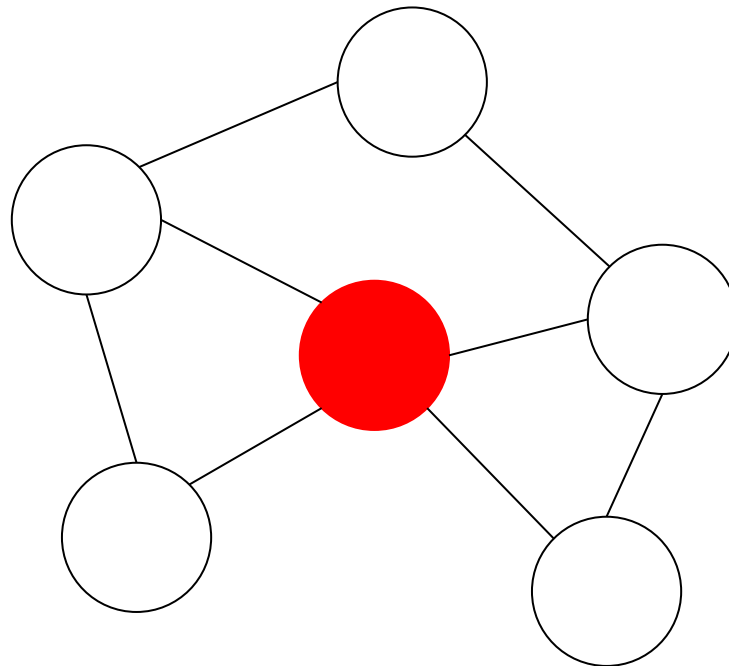
It's difficult...



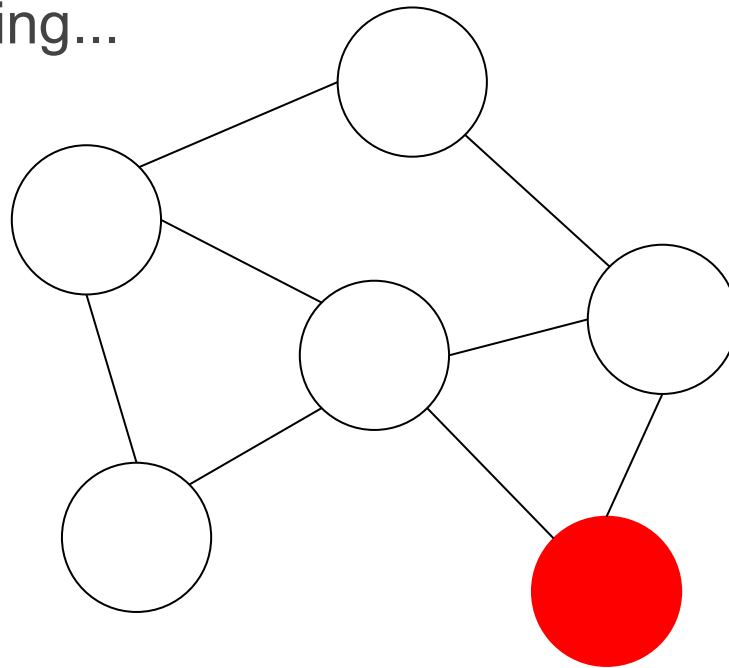
...to optimize...



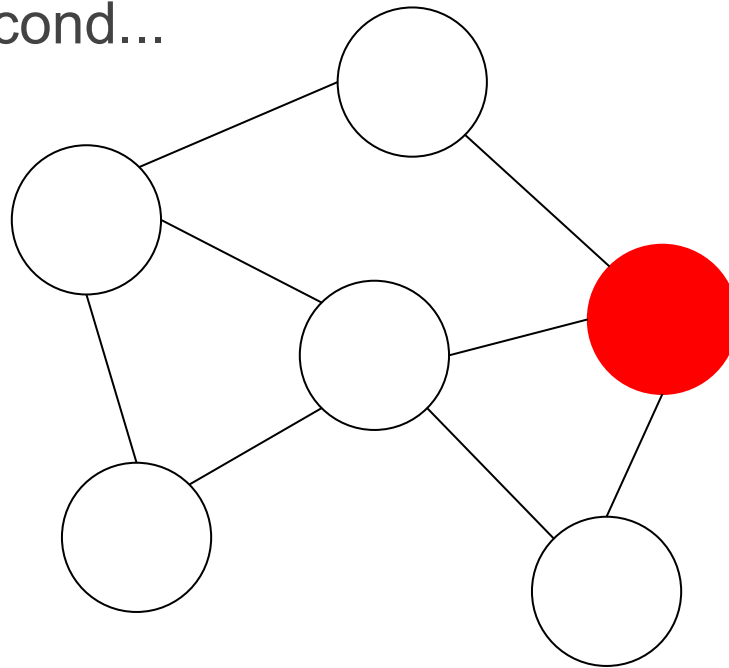
...bottlenecks...



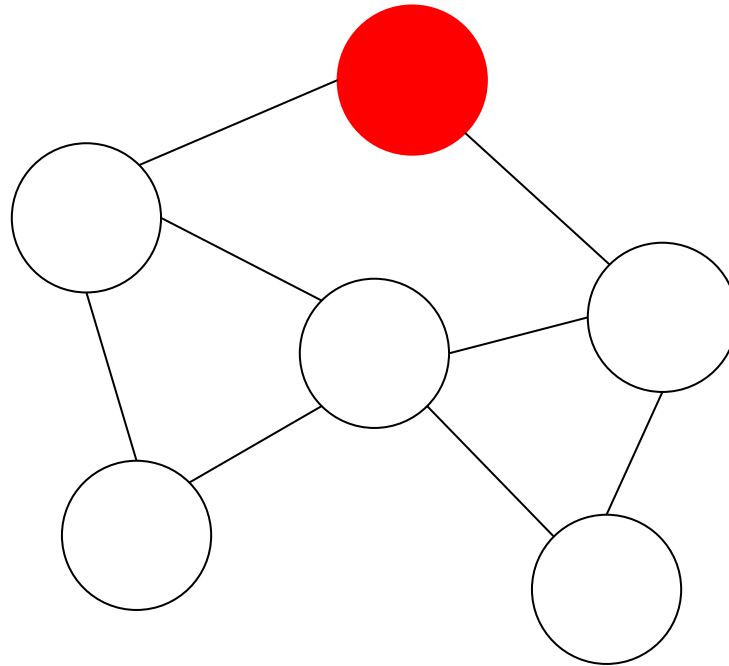
...that keep moving...



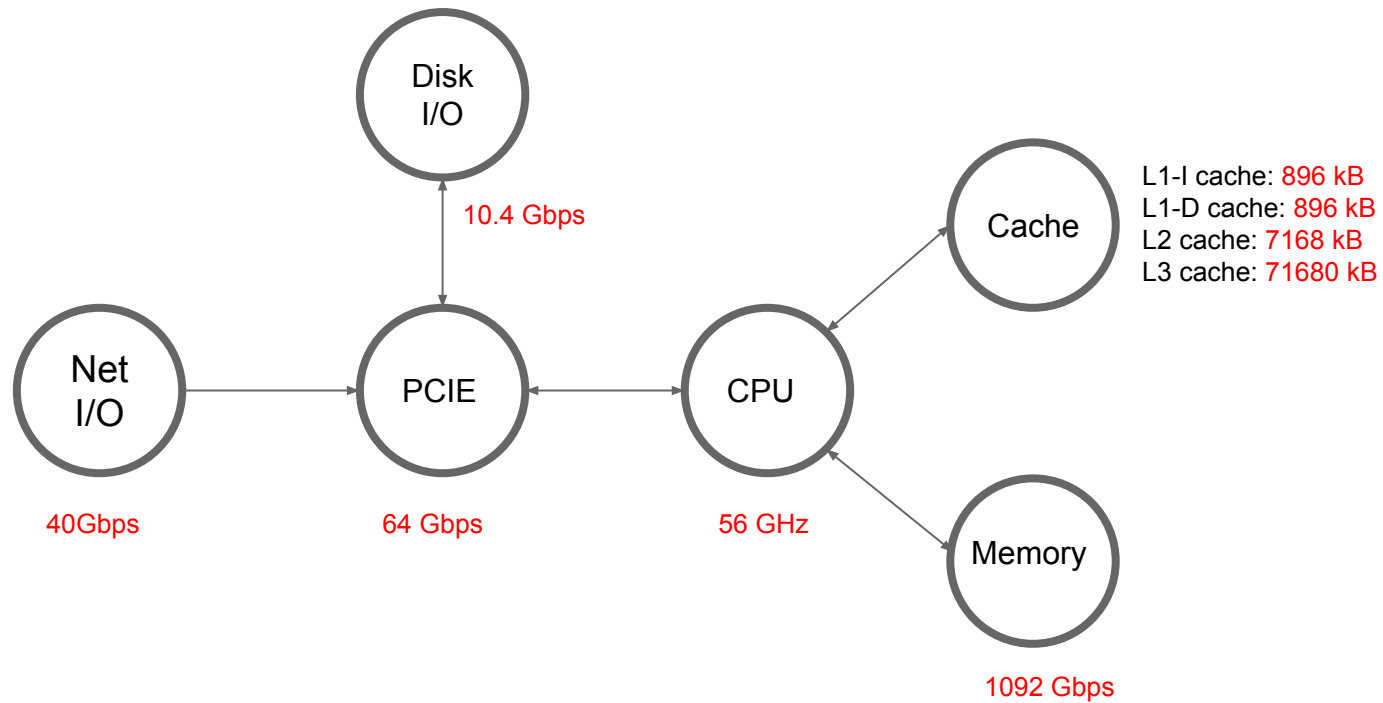
...every microsecond...



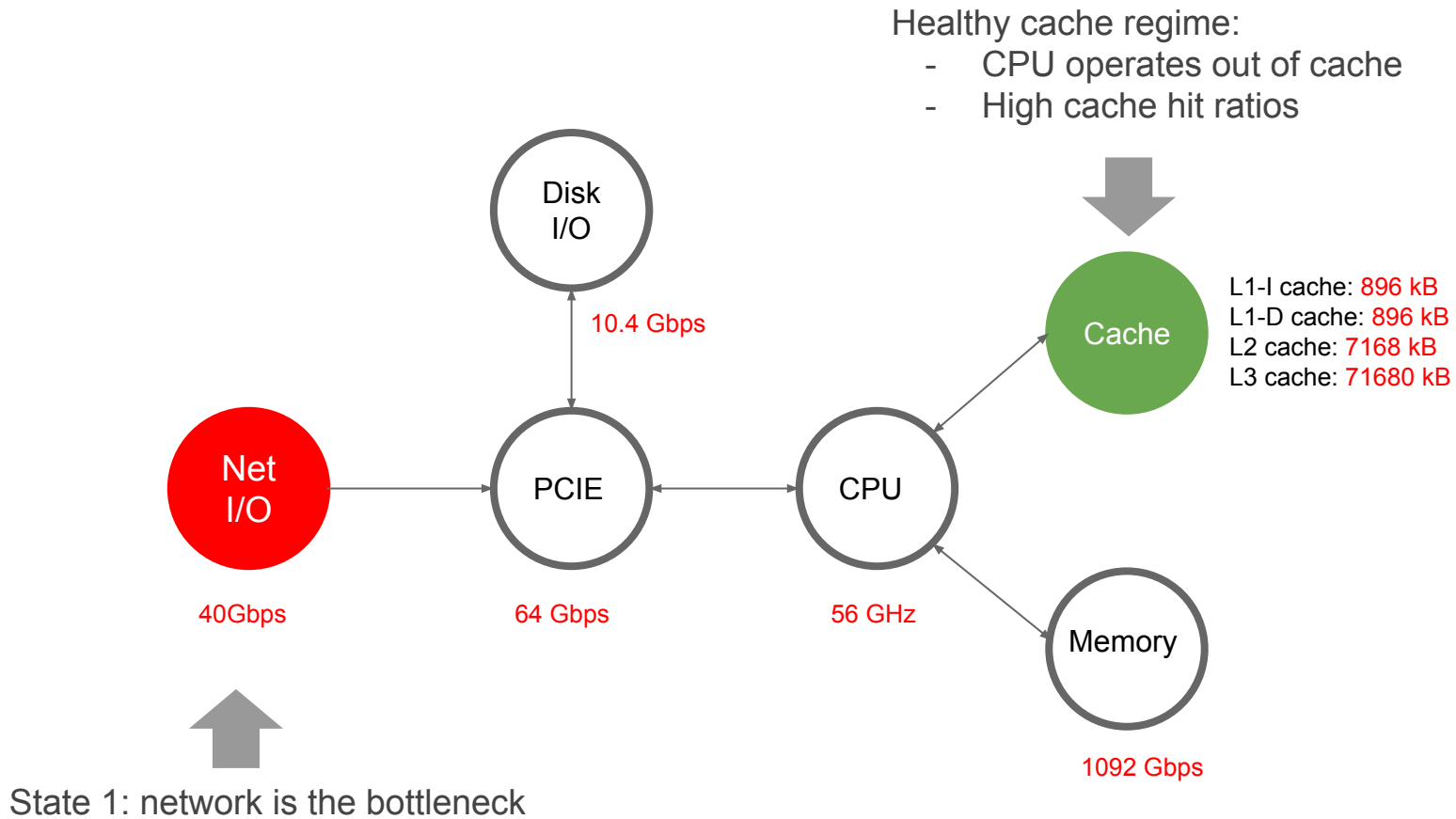
...or so.



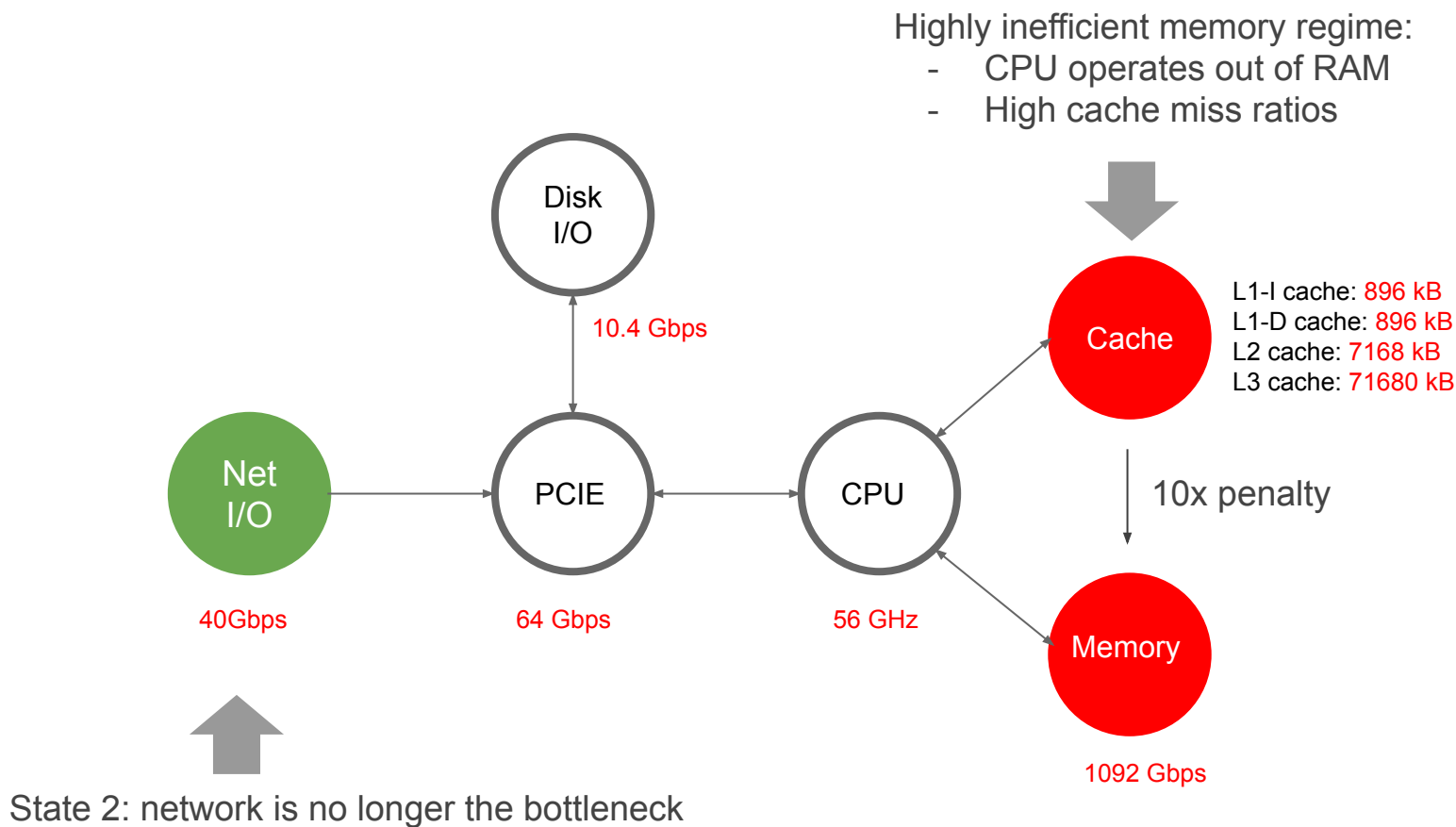
Non-linear Performance Collapse



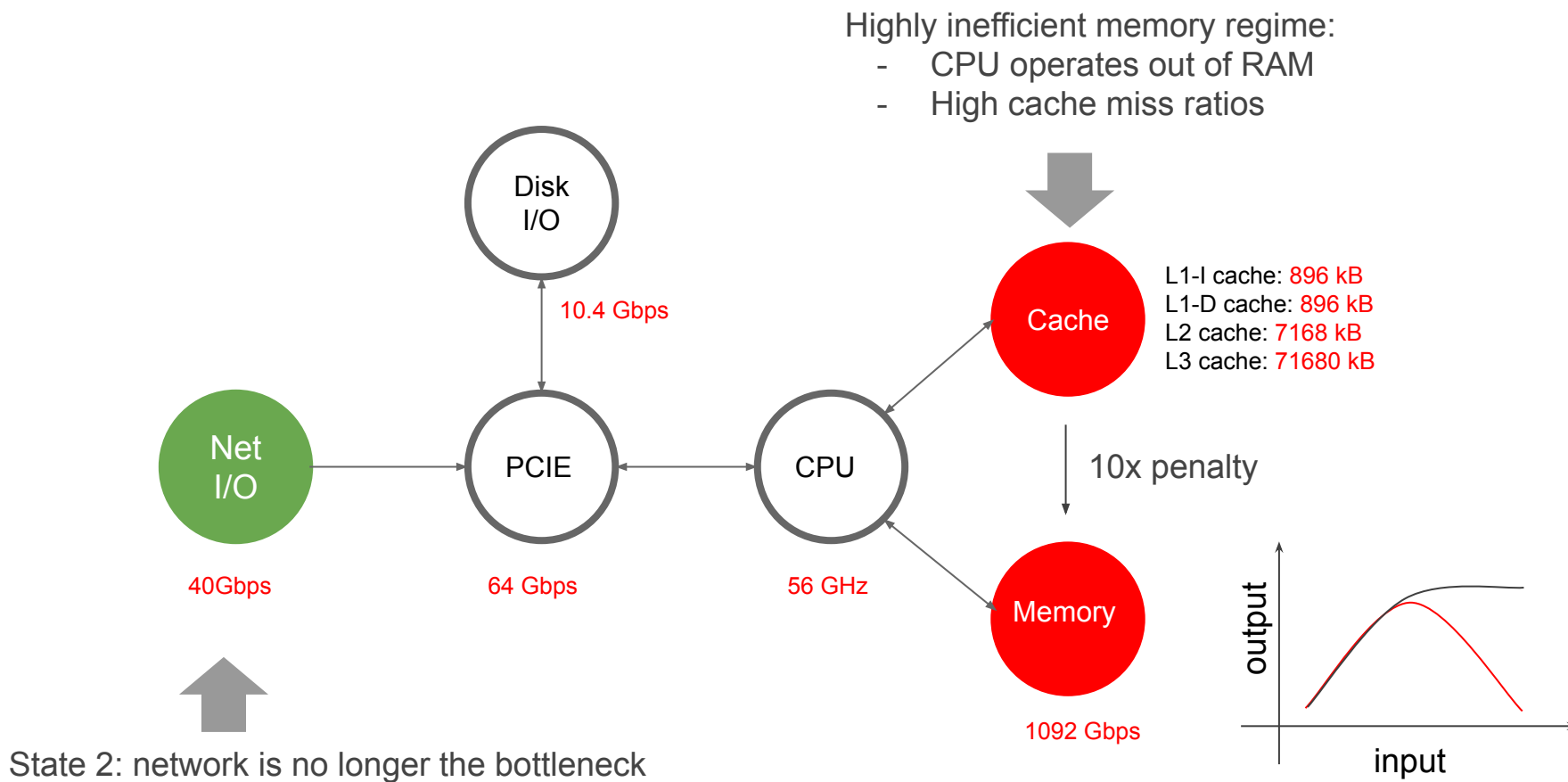
Non-linear Performance Collapse



Non-linear Performance Collapse



Non-linear Performance Collapse



By removing the network bottleneck, system spends more time processing packets that will need to be dropped anyway → net performance degradation (performance collapse)

Long queue emulation Reduces packet drops due to fixed-size hardware rings

Lockless bimodal queues Improves packet capturing performance

Tail early dropping (TED) Increases information entropy and extracted metadata

LFN tables Reduces state sharing overhead

Multiresolution priority queues Reduces cost of processing timers

Long queue emulation Reduces packet drops due to fixed-size hardware rings

Lockless bimodal queues Improves packet capturing performance

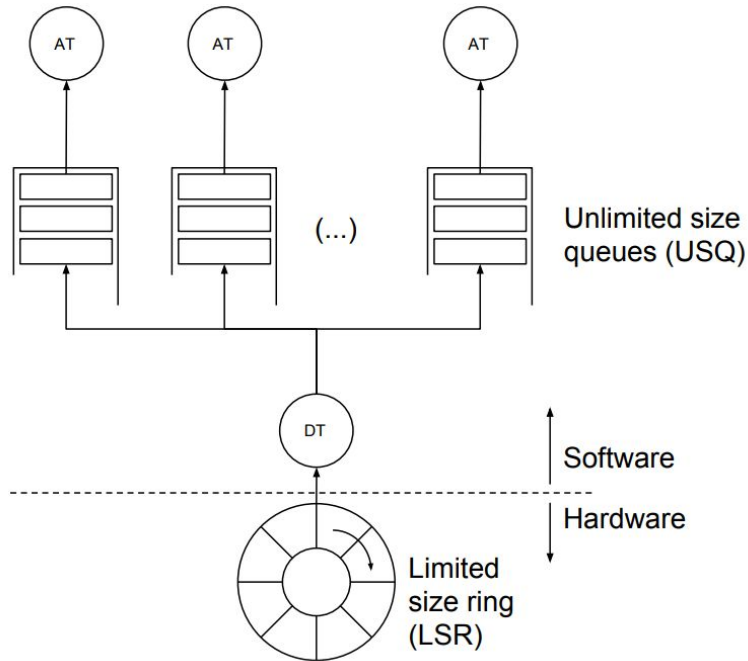
Tail early dropping (TED) Increases information entropy and extracted metadata

LFN tables Reduces state sharing overhead

Multiresolution priority queues Reduces cost of processing timers

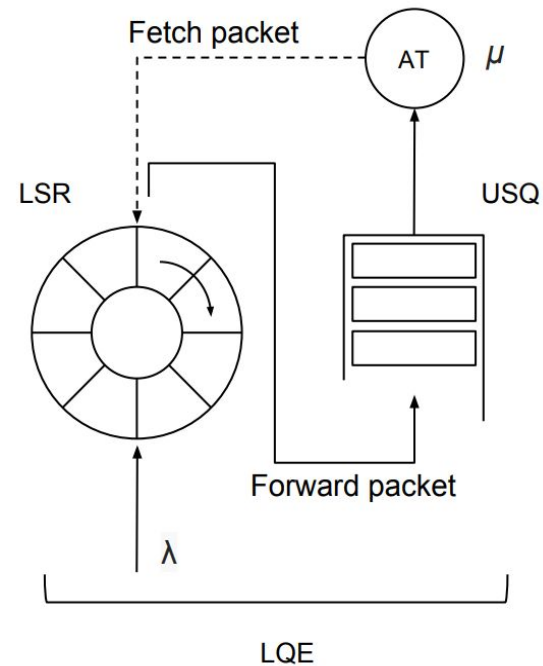
Long Queue Emulation

Dispatcher Model:



- Packet read cache penalty.
- Descriptor read cache penalty

Long queue emulation Model:



- Packet drop penalty under certain conditions

Long Queue Emulation: Operational Lemma

Lemma 1. Long queue emulation performance.

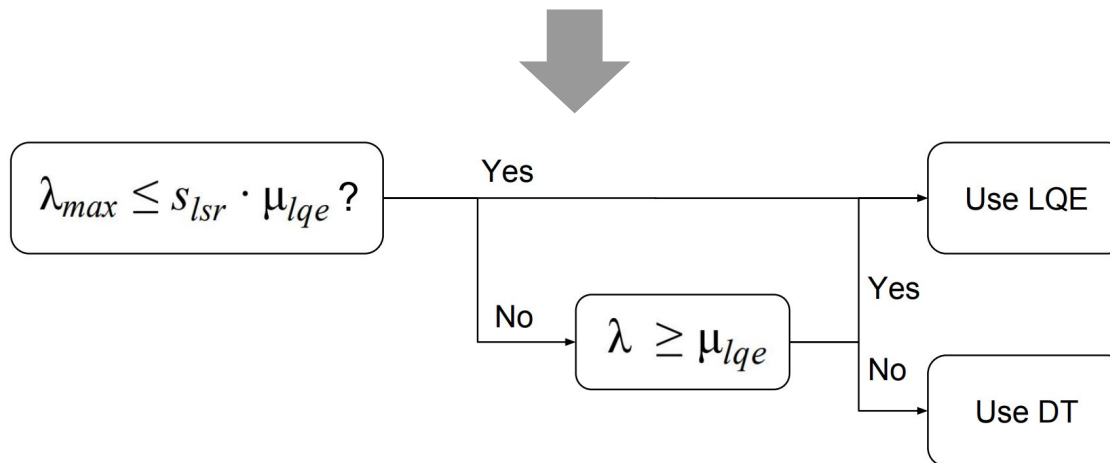
λ : average packet arrival rate

λ_{max} : maximum packet arrival rate

μ_{dt} : packet processing rate of the DT model

μ_{lqe} : packet processing rate of the LQE model

s_{lsr} : size of the LSR ring



Long Queue Emulation in Practice

Table 1. Maximum packet processing time for a Solarflare SFN7122F NIC

λ_{max} (Gbps)	1	2	4	6	8	10
s_{lsr}/λ_{max} (secs)	1.09	0.55	0.27	0.18	0.14	0.11

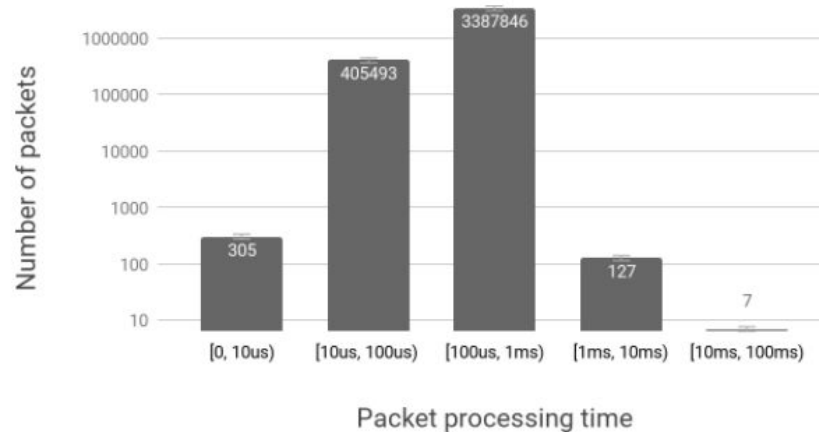
$$\lambda_{max} \leq s_{lsr} \cdot \mu_{lqe}?$$

Table 2. Packet processing time distribution.

[0, 10us)	[10us, 100us)	[100us, 1ms)	[1ms, 10ms)	[10ms, 100ms)
305	405493	3387846	127	7
Total packets:		3793778		

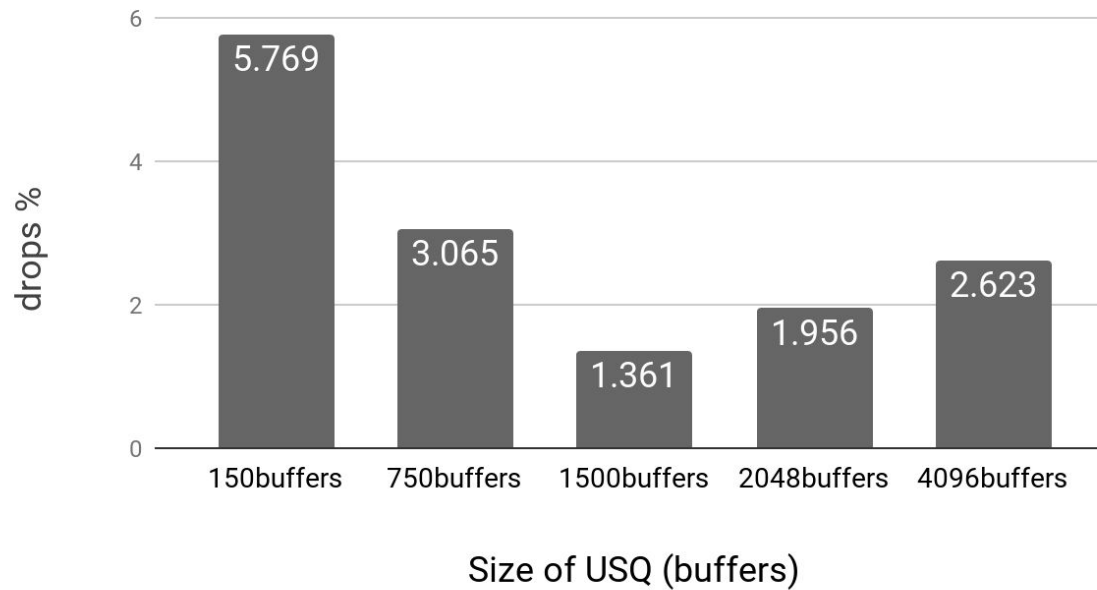
Use LQE

Packet processing time distribution



- Optimal LQE size

Packet drops at 10Gbps



Long queue emulation Reduces packet drops due to fixed-size hardware rings

Lockless bimodal queues Improves packet capturing performance

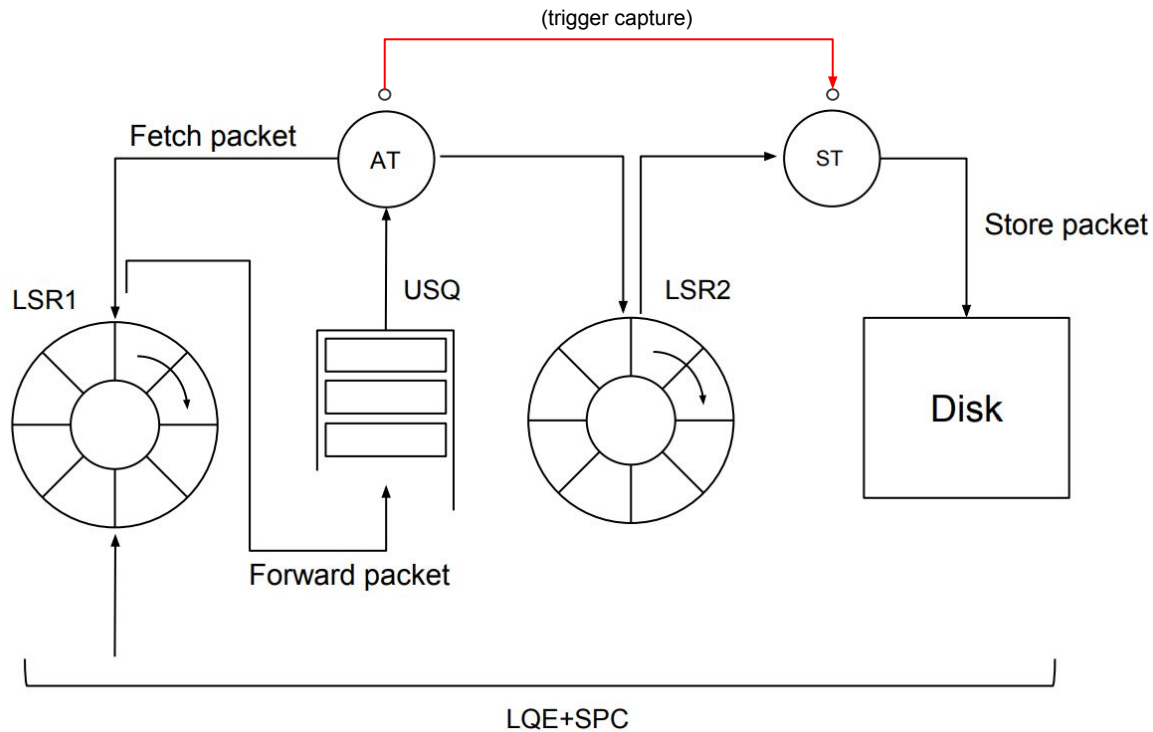
Tail early dropping Increases information entropy and extracted metadata

LFN tables Reduces state sharing overhead

Multiresolution priority queues Reduces cost of processing timers

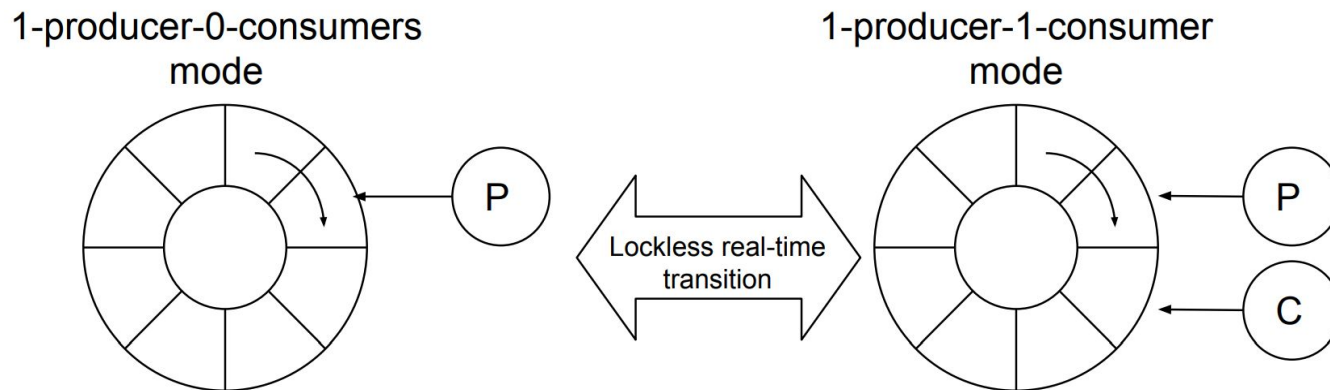
Lockless Bimodal Queues

- Goal: move packets from the memory ring to disk without using locks



Lockless Bimodal Queues

- Goal: move packets from the memory ring to disk without using locks



Lockless 1-producer-1-consumer queue

```
1  typedef struct {
2      volatile unsigned int offset_p;
3      volatile unsigned int offset_c;
4      packet_t* vector[RINGSIZE];
5  } ring_t;
6
7  void enqueue(ring_t* ring, packet_t* pkt) {
8      next_offset_p = ring->offset_p + 1 % RINGSIZE;
9      while(next_offset_p == ring->offset_c);
10     ring->vector[ring->offset_p] = pkt;
11     ring->offset_p = next_offset_p;
12 }
13
14 packet_t* dequeue(ring_t* ring) {
15     if(ring->offset_p == ring->offset_c)
16         return NULL;
17     current_offset_c = ring->offset_c;
18     ring->offset_c = ring->offset_c + 1 % RINGSIZE;
19     return(vector[current_offset_c]);
20 }
```

Lockless 1-producer-0-consumers queue

```
1  typedef struct {
2      unsigned int offset_p;
3      unsigned int offset_c;
4      packet_t* vector[RINGSIZE];
5  } ring_t;
6
7  void enqueue(ring_t* ring, packet_t* pkt) {
8      next_offset_p = ring->offset_p + 1 % RINGSIZE;
9      if(next_offset_p == ring->offset_c)
10         dequeue(ring);
11     ring->vector[ring->offset_p] = pkt;
12     ring->offset_p = next_offset_p;
13 }
14
15 packet_t* dequeue(ring_t* ring) {
16     if(ring->offset_p == ring->offset_c)
17         return NULL;
18     current_offset_c = ring->offset_c;
19     ring->offset_c = ring->offset_c + 1 % RINGSIZE;
20     return(vector[current_offset_c]);
21 }
```

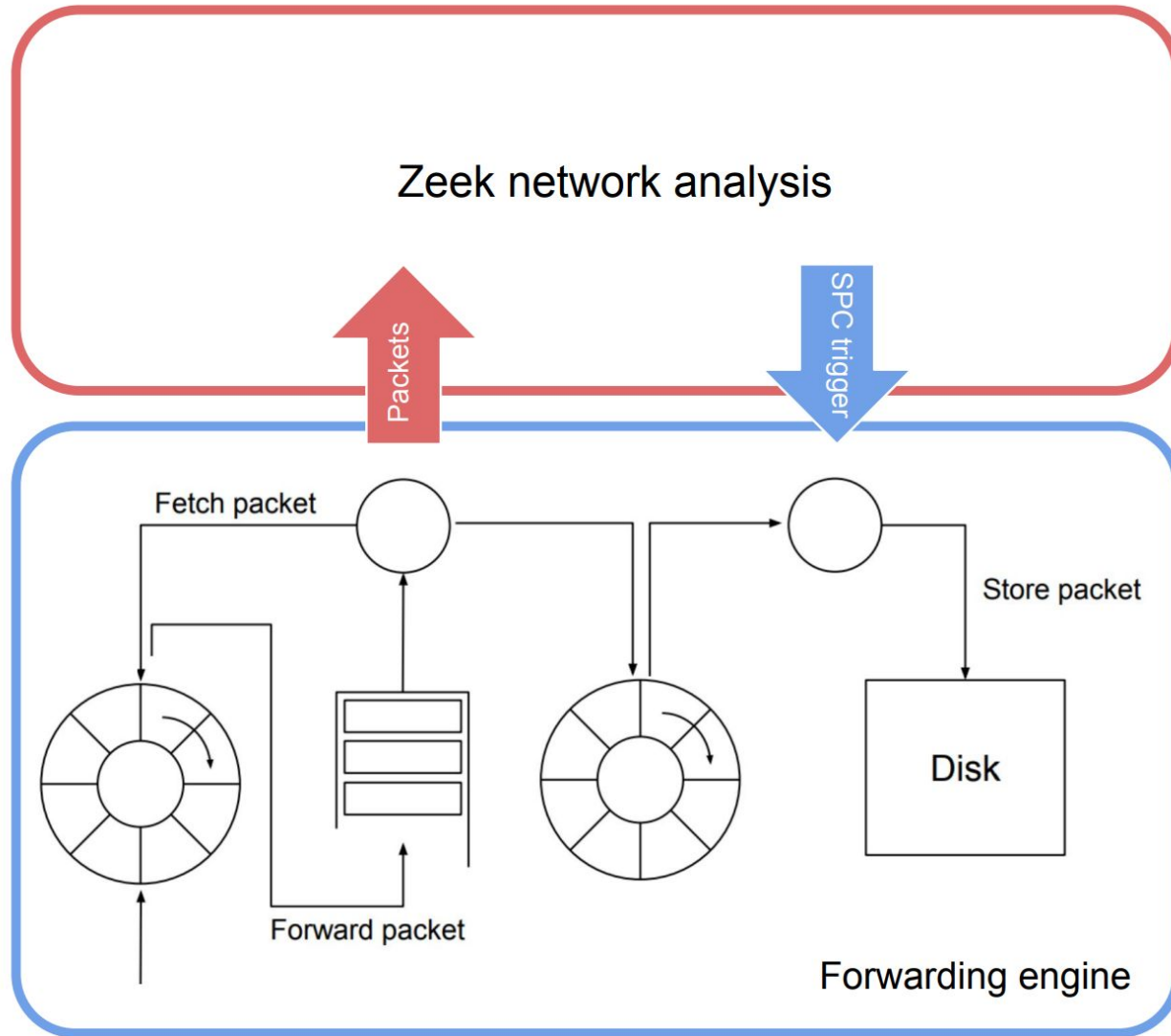
Lockless Bimodal Queues

Lockless bimodal queue using CAS
(producer does not need to be permanently active)

```
1  typedef struct {
2      volatile unsigned int offset_p;
3      volatile unsigned int offset_c;
4      volatile bool trans; // Used to transition modes
5      volatile bool state; // True, consumer is ON;
                           // False consumer is OFF
6      packet_t* vector[RINGSIZE];
7  } ring_t;
8
9  void enqueue(ring_t* ring, packet_t* pkt) {
10     next_offset_p = ring->offset_p + 1 % RINGSIZE;
11     while(!cas(&ring->trans, false, true));
12     if(!ring->state) {
13         if(next_offset_p == ring->offset_c)
14             dequeue(ring);
15     }
16     else
17         while(next_offset_p == ring->offset_c);
18     ring->trans = false;
19     ring->vector[ring->offset_p++] = pkt;
20     ring->offset_p = next_offset_p;
21 }
22
23 packet_t* dequeue(ring_t* ring) {
24     if(ring->offset_p == ring->offset_c)
25         return NULL;
26     current_offset_c = ring->offset_c;
27     ring->offset_c = ring->offset_c + 1 % RINGSIZE;
28     return(vector[current_offset_c]);
29 }
30
31 void start_c(ring_t* ring) {
32     while(!cas(&ring->trans, false, true));
33     ring->state = true;
34     ring->trans = false;
35 }
36
37 void stop_c(ring_t* ring) {
38     while(!cas(&ring->trans, false, true));
39     ring->state = false;
40     ring->trans = false;
41 }
```

Lockless bimodal queue without using CAS
(producer must be permanently active to avoid consumer starvation)

```
1  typedef struct {
2      volatile unsigned int offset_p;
3      volatile unsigned int offset_c;
4      volatile bool req; // owned by consumer
5      volatile bool ack; // owned by producer
6      packet_t* vector[RINGSIZE];
7  } ring_t;
8
9  void enqueue(ring_t* ring, packet_t* pkt) {
10     next_offset_p = ring->offset_p + 1 % RINGSIZE;
11     if(!ring->req) {
12         if(ring->ack)
13             ring->ack = false;
14         if(next_offset_p == ring->offset_c)
15             dequeue(ring);
16     }
17     else {
18         if(!ring->ack)
19             ring->ack = true;
20         while(next_offset_p == ring->offset_c);
21     }
22     ring->vector[ring->offset_p] = pkt;
23     ring->offset_p = next_offset_p;
24 }
25
26 packet_t* dequeue(ring_t* ring) {
27     if(ring->offset_p == ring->offset_c)
28         return NULL;
29     current_offset_c = ring->offset_c;
30     ring->offset_c = ring->offset_c + 1 % RINGSIZE;
31     return(vector[current_offset_c]);
32 }
33
34 void start_c(ring_t* ring) {
35     ring->req = true;
36     while(!ring->ack);
37 }
38
39 void stop_c(ring_t* ring) {
40     ring->req = false;
41     while(ring->ack);
42 }
```



- The function `spc_capture()` takes two arguments as shown by its function prototype:

```
## API for capturing a Pcap  
function spc_capture(prefix: string, filter: string);
```

- The `prefix` argument allows users to specify a prefix for the generated Pcap file name. The `filter` argument can be used to specify a BPF filter applied to the captured packets as they are written to the pcap file. See <https://www.tcpdump.org/manpages/pcap-filter.7.html> for the expression syntax of the BPF filter. If set to the empty string "", all packets (without any filtering) are written to the Pcap file.

Selective Packet Capture by Example

```
##! Selective Packet Capture trigger example

## This line contains the list of hosts that will cause a capture file to be generated
## when part of a DNS query
global spc_host_trigger_list: set[string] = { "www.purple.com" } &redef;

## This line contains the prefix to use when creating capture files.
global spc_prefix: string = "";

event dns_request(c: connection, msg: dns_msg, query: string, qtype: count, qclass: count)
{
    if (query in spc_host_trigger_list) {
        local spc_filter = fmt("src %s and port %s", c$id$orig_h, port_to_count(c$id$orig_p));
        spc_capture(spc_prefix, spc_filter);
    }
}
```

Thank You

Reservoir Labs

632 Broadway
Suite 803
New York, NY 10012

812 SW Washington St.
Suite 1200
Portland, OR 97205