Advanced Computer Networks

Network Virtualization in Data Center Networks (II)

https://pages.cs.wisc.edu/~mgliu/CS740/F25/index.html

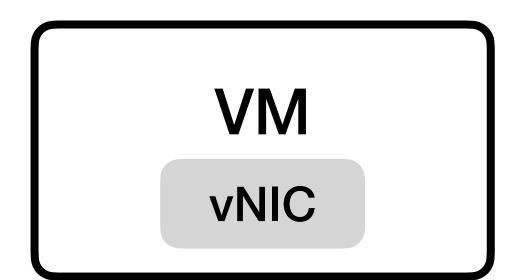
Ming Liu mgliu@cs.wisc.edu

Outline

- Last lecture
 - Network virtualization in data center networks (I)

- Today
 - Network virtualization in data center networks (II)

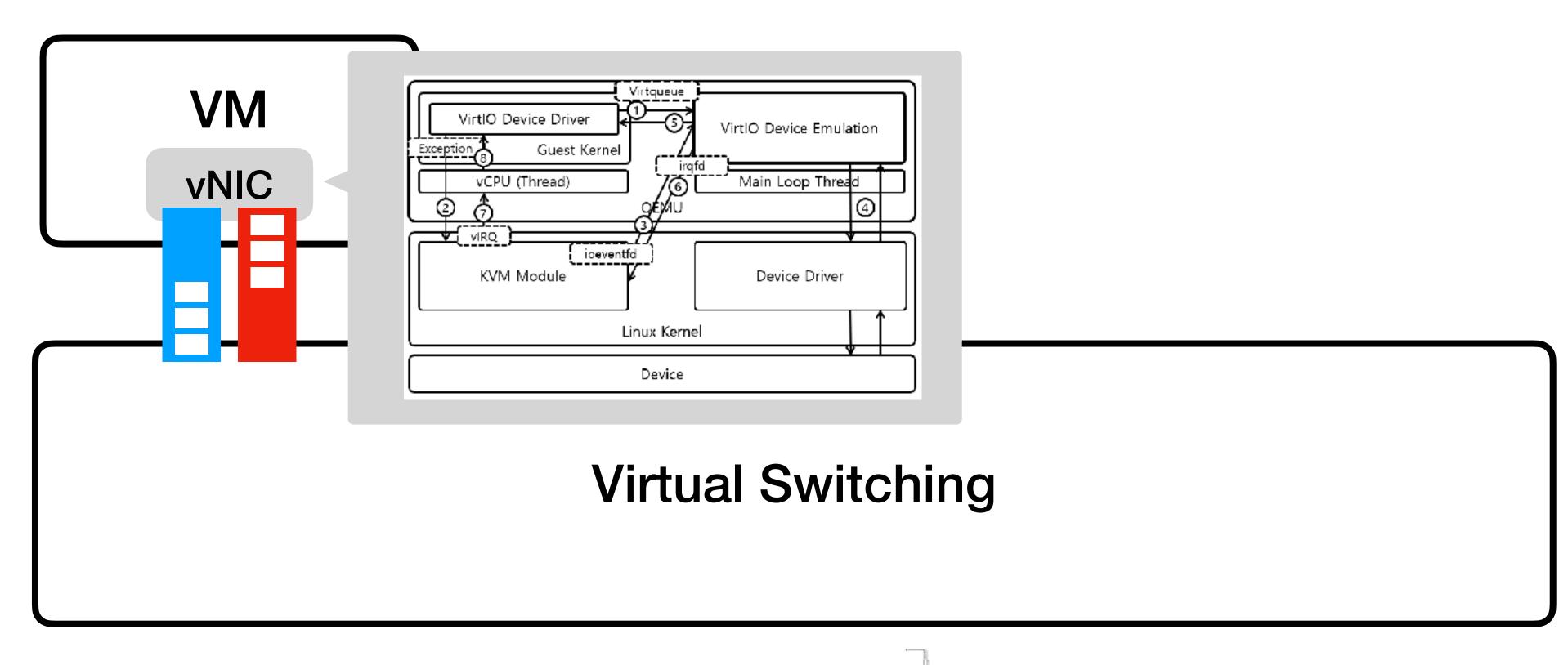
- Announcements
 - Lab2 due 11/05/2025 11:59 PM
 - Midterm report due 11/04/2025 11:59 PM



Virtual Switching

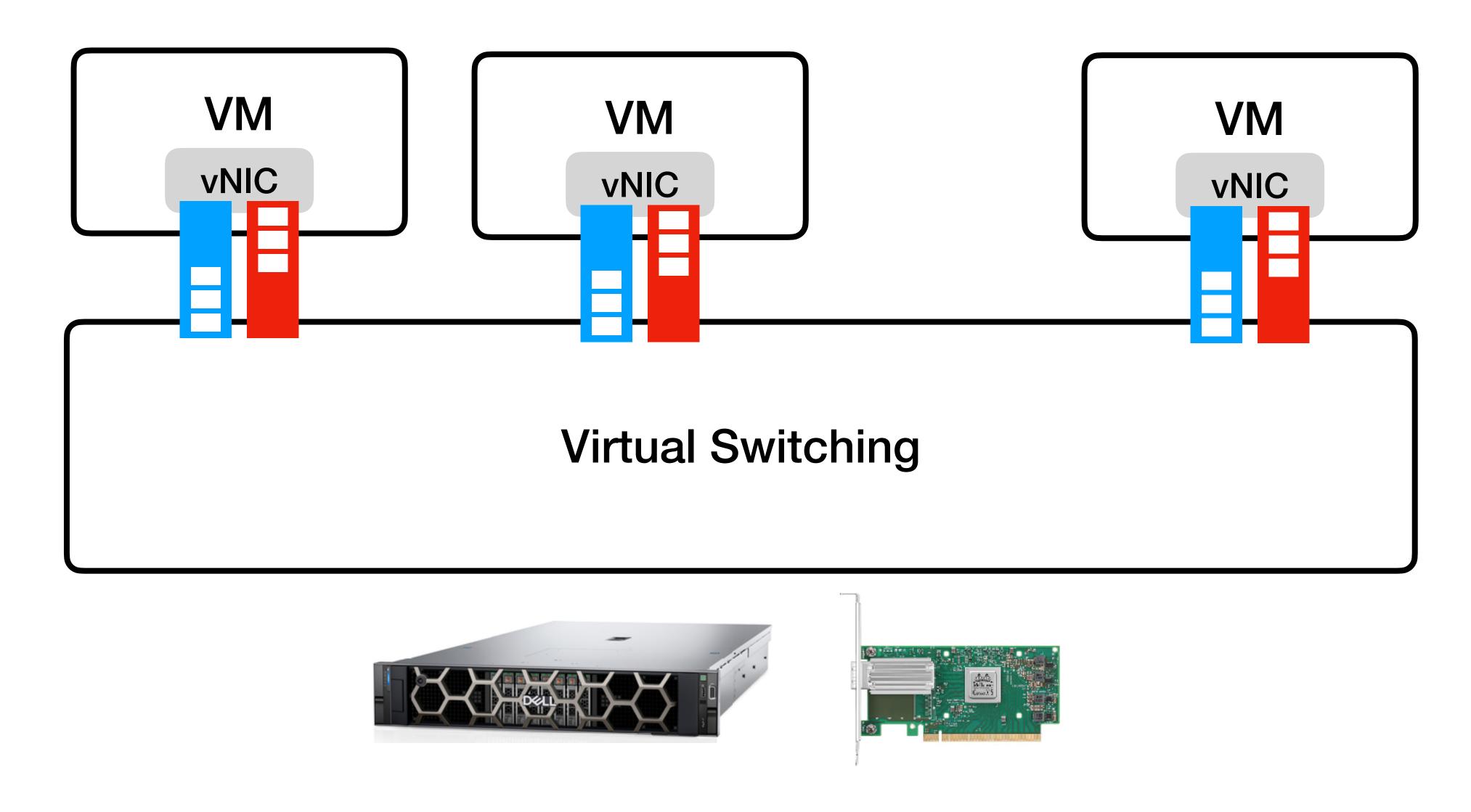


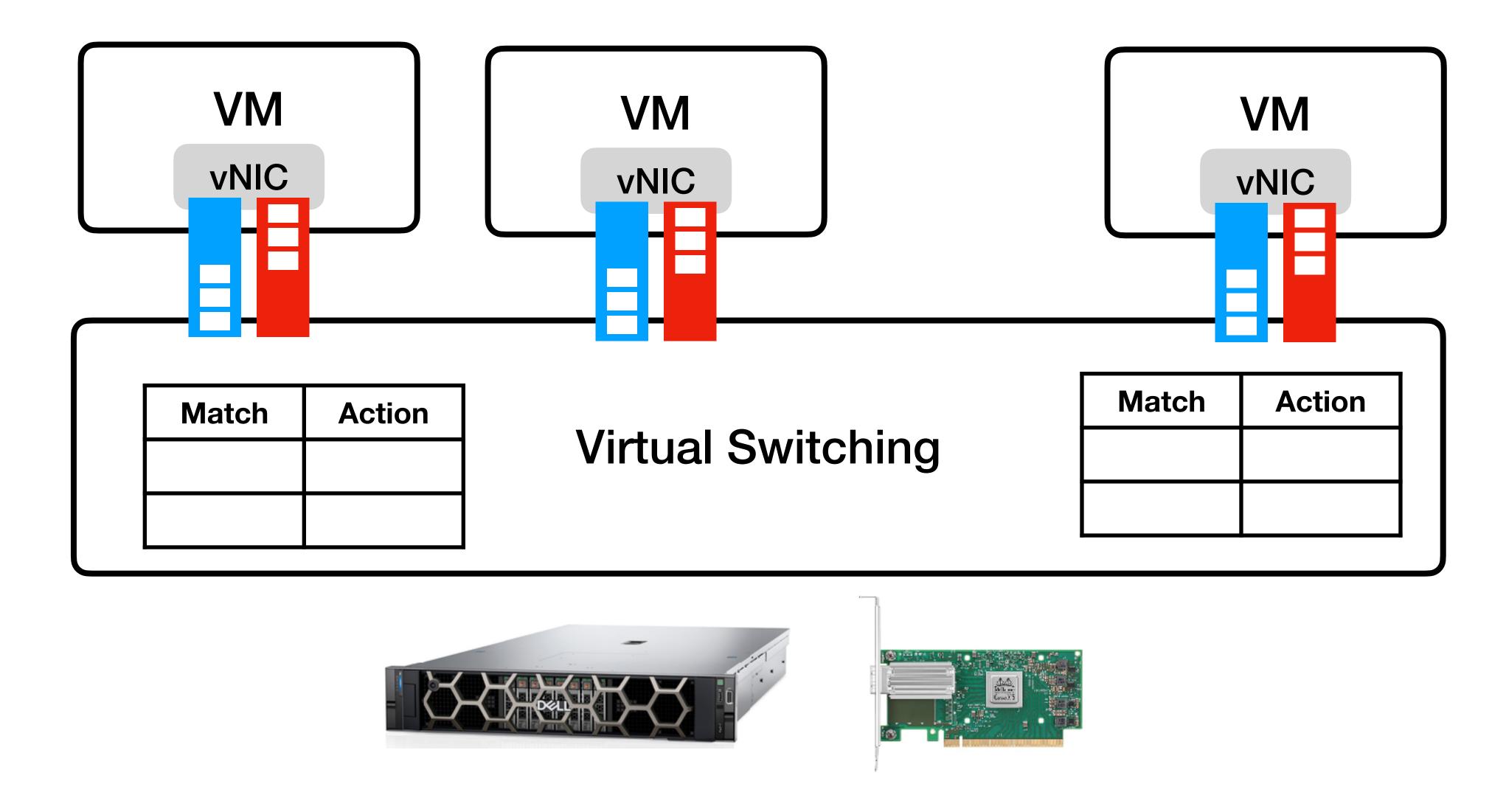


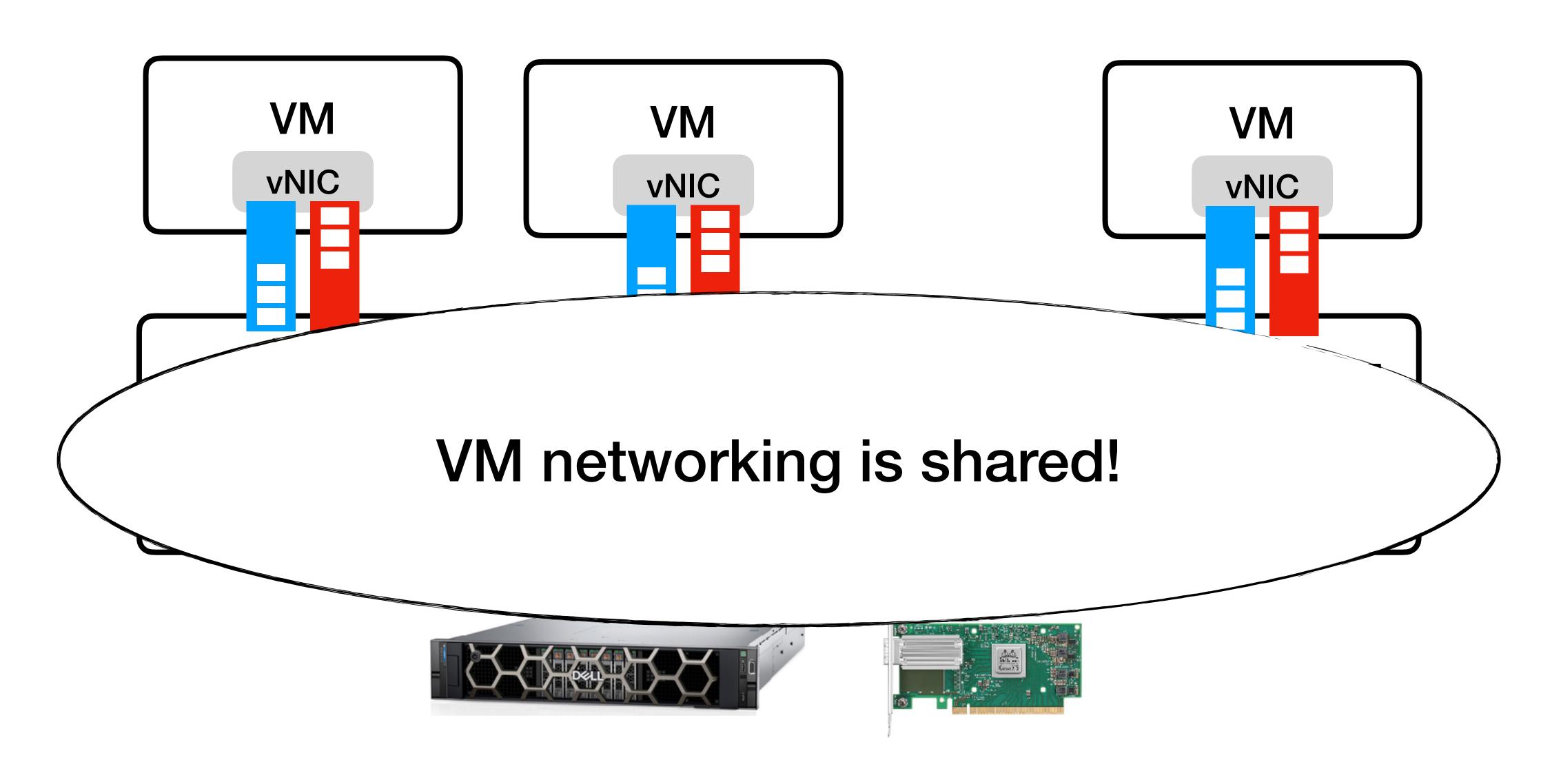








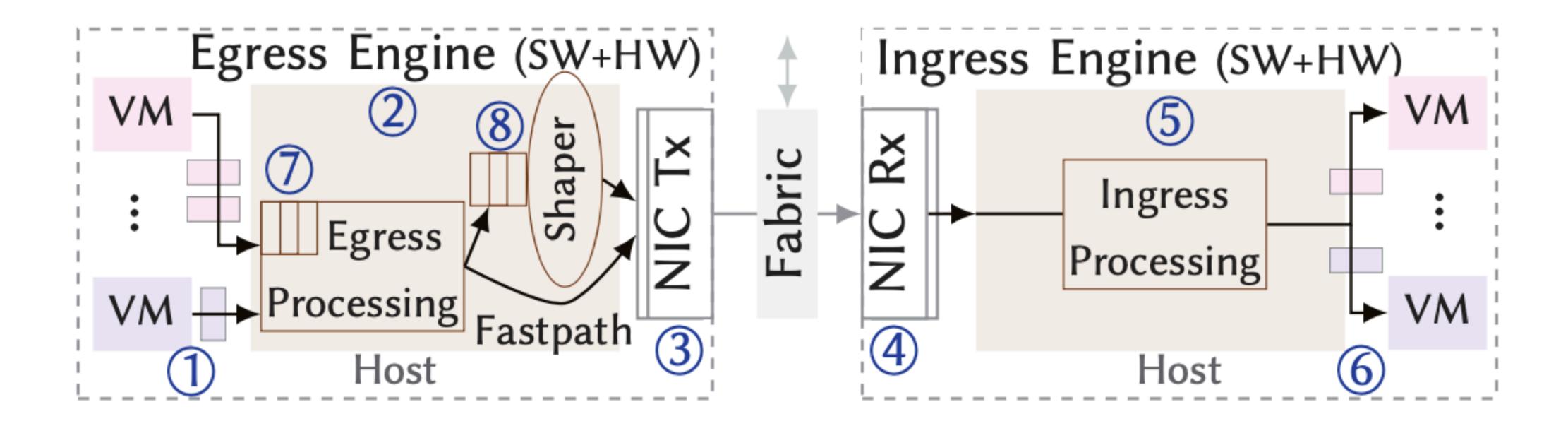




How can we achieve network performance isolation at the endhost?

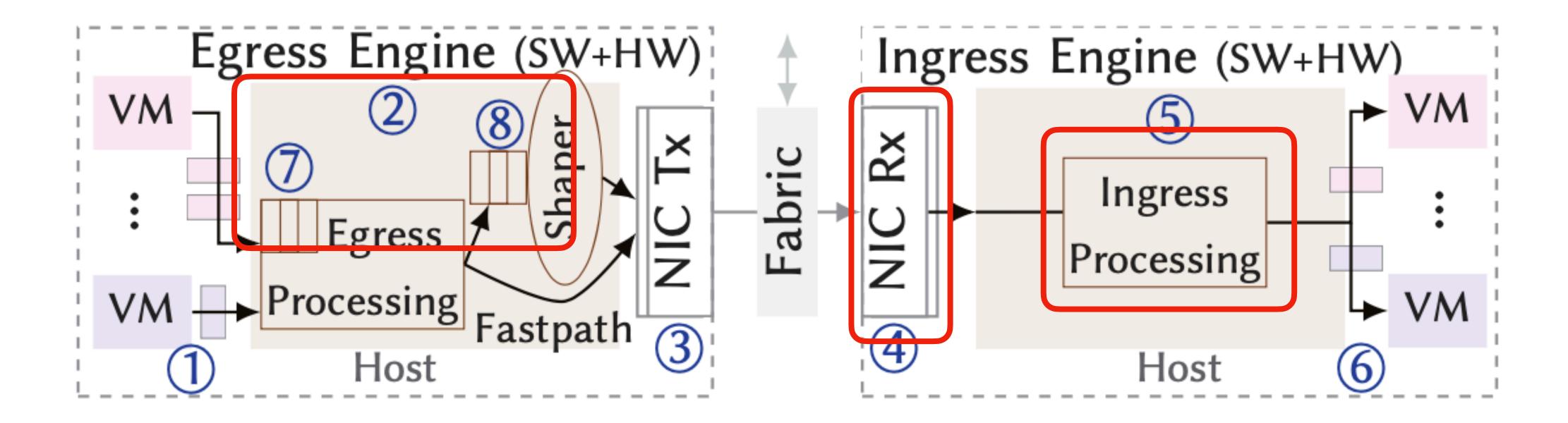
Packet Path Between Two VMs

• Per-packet processing costs are not const.



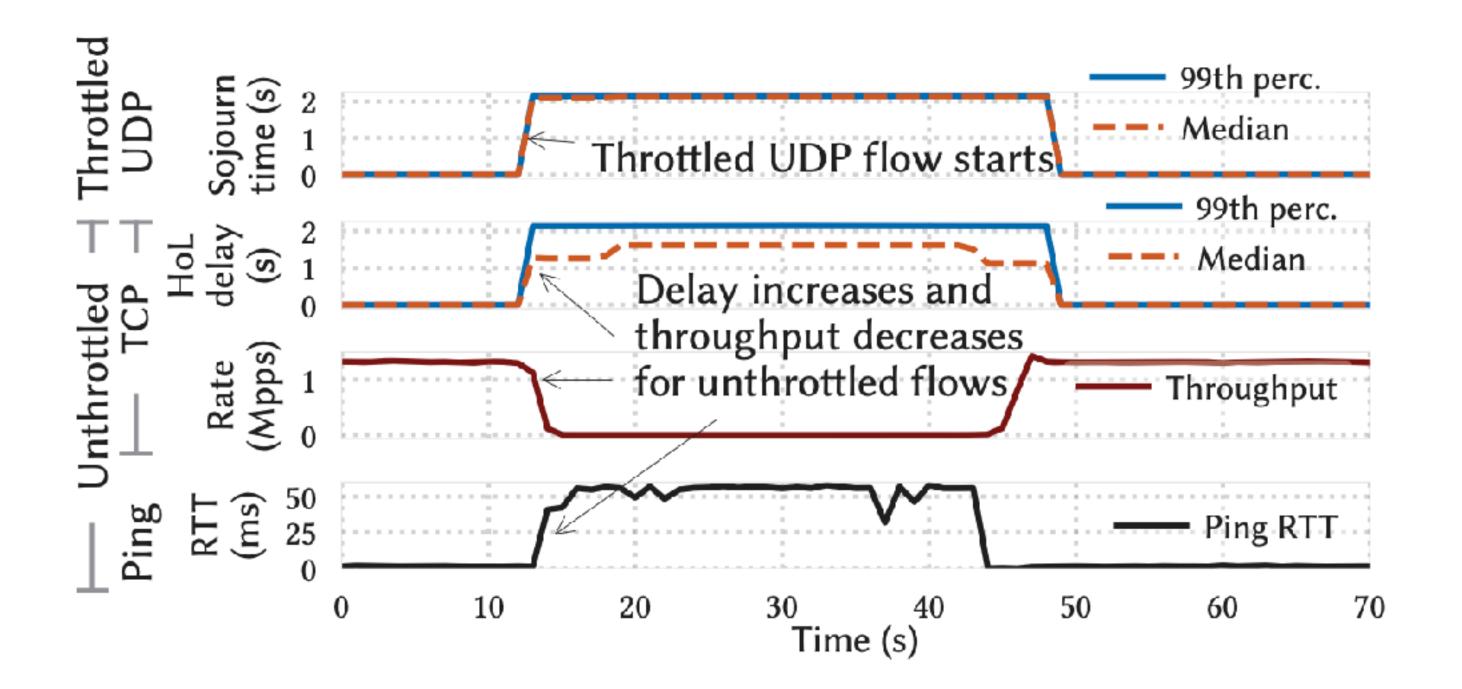
Packet Path Between Two VMs

• Per-packet processing costs are not const.



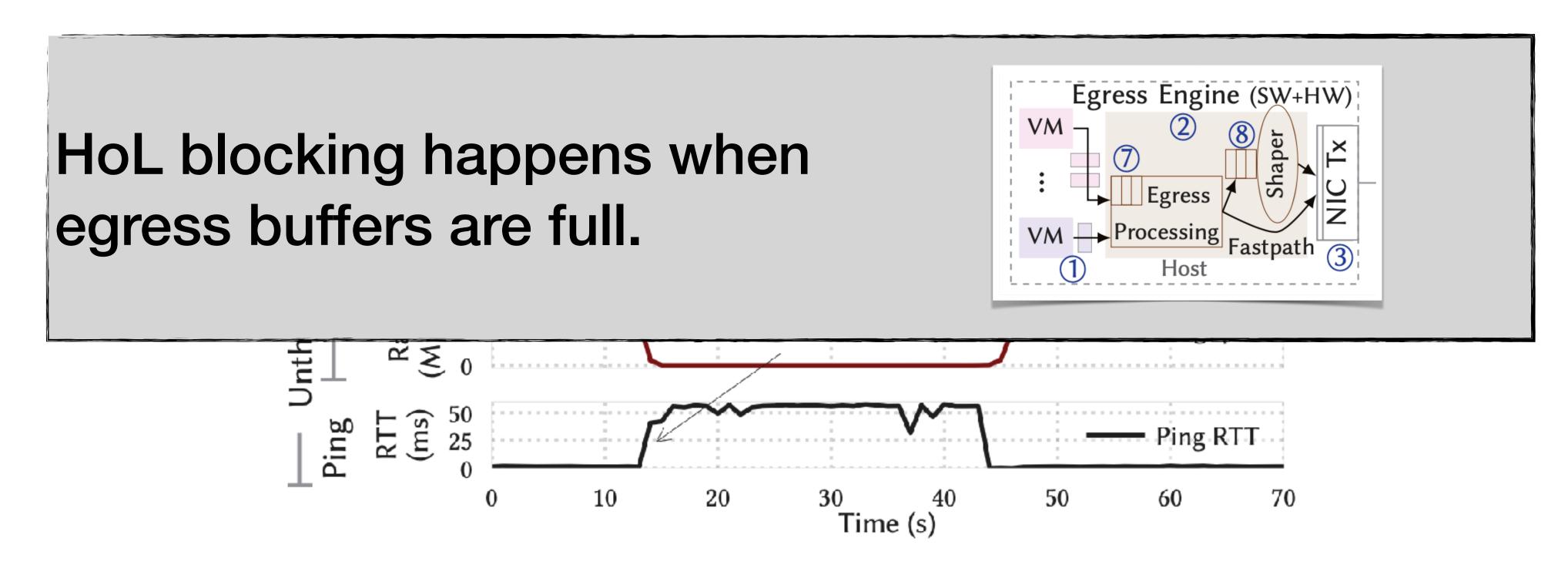
Issue #1: Egress Buffer Contention

- Experiment setup:
 - VM1 —> VM2: unthrottled TCP flow
 - VM1 —> VM3: throttled UDP flow @ 10Mbps



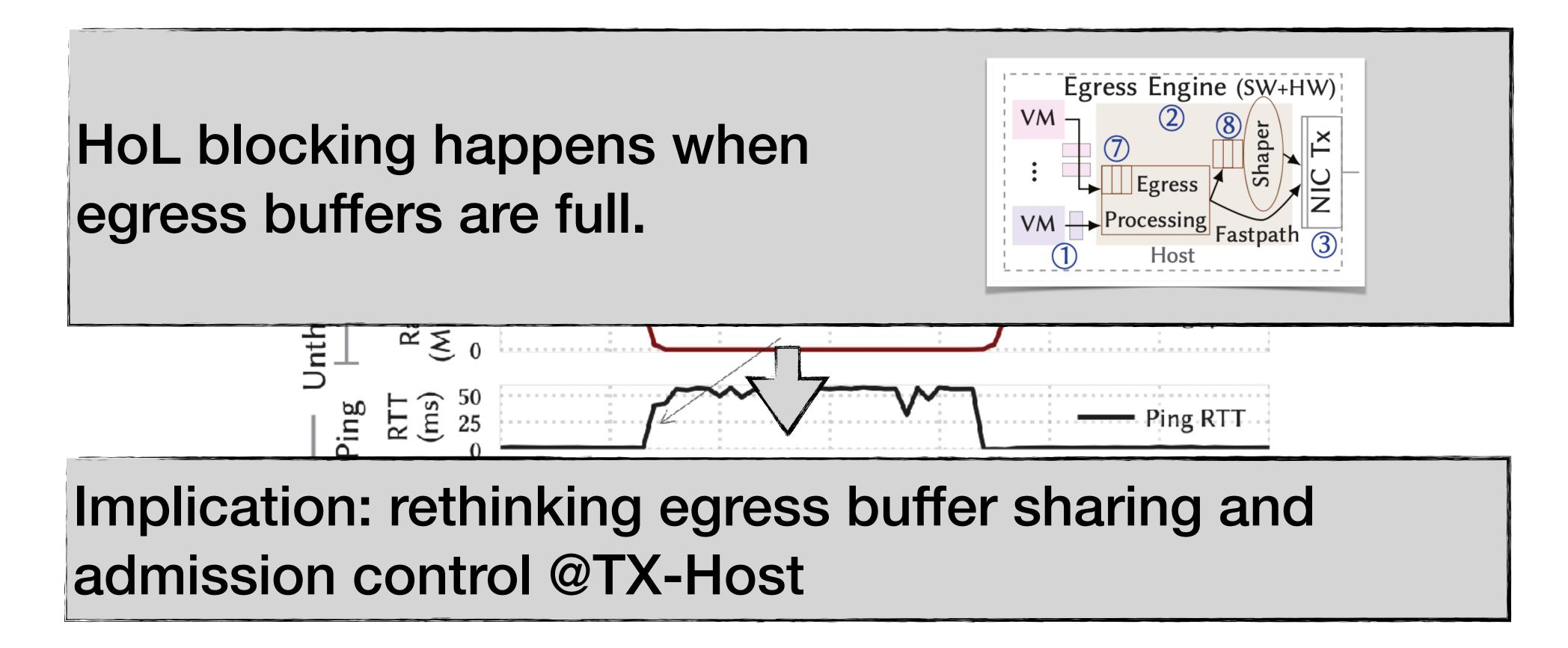
Issue #1: Egress Buffer Contention

- Experiment setup:
 - VM1 —> VM2: unthrottled TCP flow
 - VM1 —> VM3: throttled UDP flow @ 10Mbps



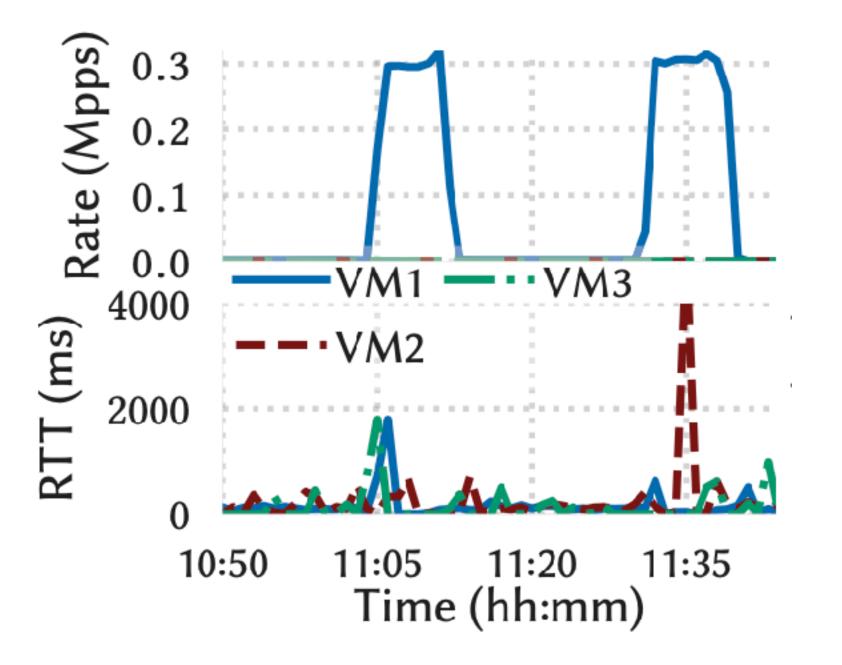
Issue #1: Egress Buffer Contention

- Experiment setup:
 - VM1 —> VM2: unthrottled TCP flow
 - VM1 —> VM3: throttled UDP flow @ 10Mbps



Issue #2: Ingress NIC Contention

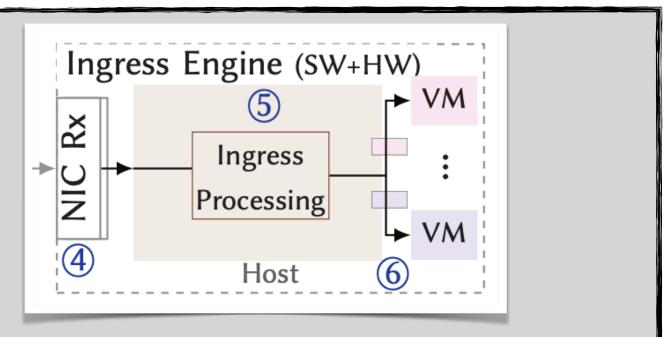
- Experiment setup:
 - Three co-located: VM1, VM2, VM3
 - VM1 receives a packet burst from small packets

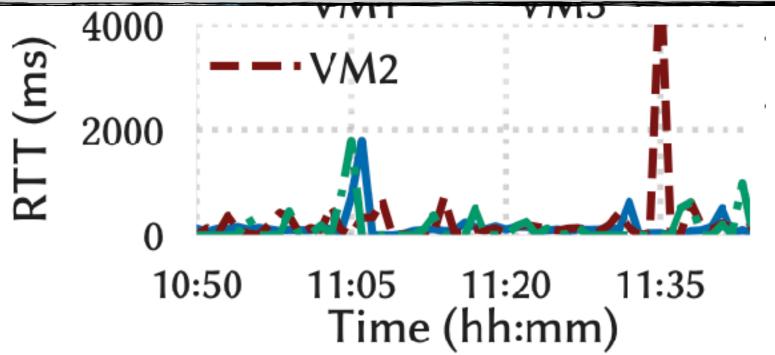


Issue #2: Ingress NIC Contention

- Experiment setup:
 - Three co-located: VM1, VM2, VM3
 - VM1 receives a packet burst from small packets

A storm of small packets can exceed the NIC packet processing capacity.





Issue #2: Ingress NIC Contention

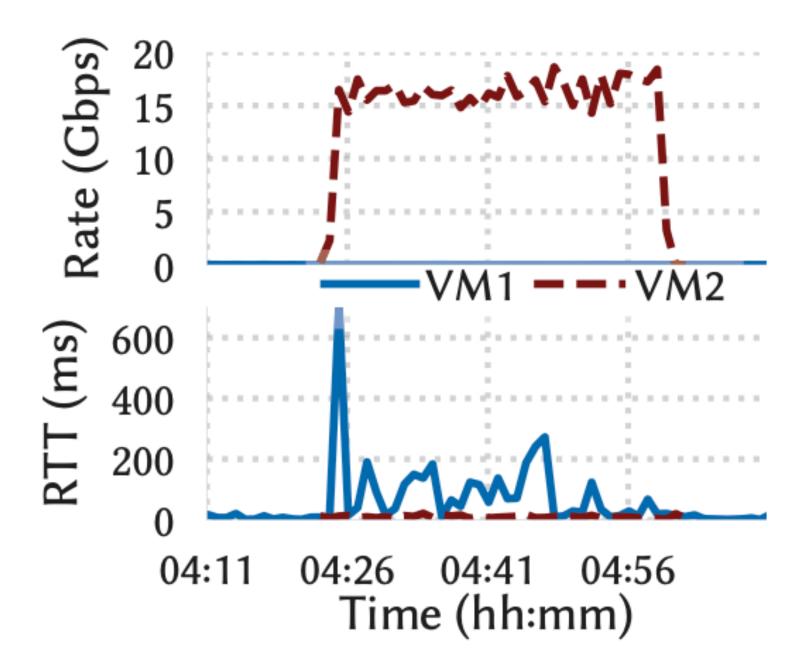
- Experiment setup:
 - Three co-located: VM1, VM2, VM3
 - VM1 receives a packet burst from small packets

Implication: rethinking packet dropping policy @ RX-NIC

Time (hh:mm)

Issue #3: Ingress Engine Contention

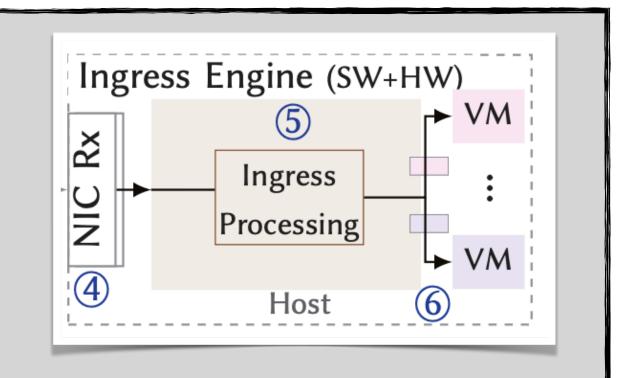
- Experiment setup:
 - Two co-located: VM1, VM2
 - VM2 receives a burst of ~16Gbps traffic

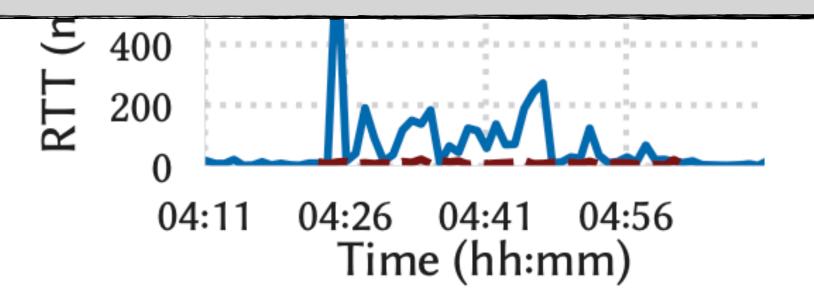


Issue #3: Ingress Engine Contention

- Experiment setup:
 - Two co-located: VM1, VM2
 - VM2 receives a burst of ~16Gbps traffic

A packet burst can exceed the ingress engine's packet processing capacity.

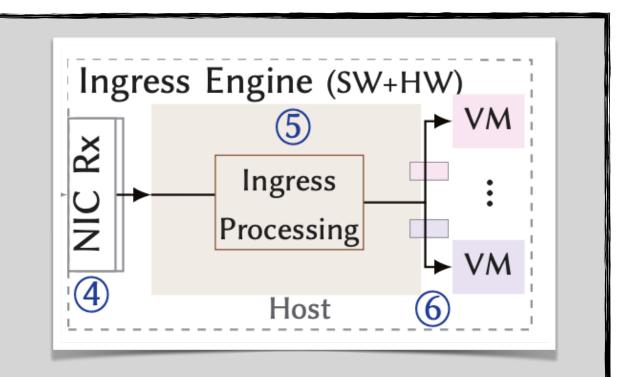


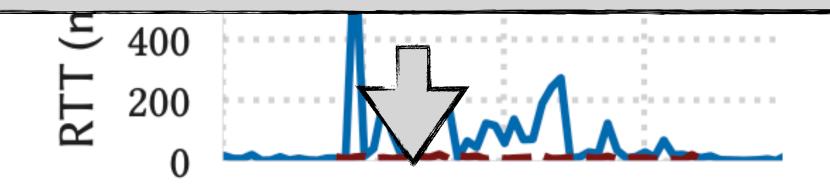


Issue #3: Ingress Engine Contention

- Experiment setup:
 - Two co-located: VM1, VM2
 - VM2 receives a burst of ~16Gbps traffic

A packet burst can exceed the ingress engine's packet processing capacity.





Implication: rethinking computation allocation@RX-Host

Design principles behind PicNIC:

- P1: SLO-based resource sharing => Sacrifice utilization
- P2: Backpressure and early drops => Traffic regulation

- #1: Bandwidth
 - [MIN_BPS, MAX_BPS]
 - Consider PPS (packets per second)

- #1: Bandwidth
 - [MIN_BPS, MAX_BPS]
 - Consider PPS (packets per second)
- #2: Delay
 - Ingress delay: delay in NIC + delay in engine
 - Egress delay: delay from the guest OS Tx queue to the wire

- #1: Bandwidth
 - [MIN_BPS, MAX_BPS]
 - Consider PPS (packets per second)
- #2: Delay
 - Ingress delay: delay in NIC + delay in engine
 - Egress delay: delay from the guest OS Tx queue to the wire
- #3: Loss rate
 - Well-behaved VMs: no drop
 - Uncooperative VMs: drop

- CPU-fair weighted fair queues
 - Per-VM packet queues

- CPU-fair weighted fair queues
 - Per-VM packet queues
- Enqueue
 - Packet classification

- CPU-fair weighted fair queues
 - Per-VM packet queues
- Enqueue
 - Packet classification
- Dequeue
 - Monitor the packet processing time for each VM and apply EMWA
 - Allocate dequeueing CPU works based on SLO

- CPU-fair weighted fair queues
 - Per-VM packet queues
- Enqueue
 - Packet classification
- Dequeue
 - Monitor the packet processing time for each VM and apply EMWA
 - Allocate dequeueing CPU works based on SLO

Target: issue #3

Benefits: delay and drops

Receiver-driven

PCCB: throughput mode

• PCCP: latency mode

- Receiver-driven
 - PCCB: throughput mode
 - PCCP: latency mode

PCCB

- Receivers use the Max-min fairness to determine the rate
- Senders apply an RCP-like approach to control the rate

- Receiver-driven
 - PCCB: throughput mode
 - PCCP: latency mode

PCCB

- Receivers use the Max-min fairness to determine the rate
- Senders apply an RCP-like approach to control the rate

PCCP

```
Input: delay_in
delay \leftarrow EWMA(delay, delay_in)
                                   ▶ multiplicative decrease (MD)
if delay > threshold then
  rate \leftarrow (1 - \beta \cdot (1 - \frac{threshold}{delay})) \cdot rate
                              enter fast recovery (FR)
  counter \leftarrow 0
  target\_rate \leftarrow rate
                                      ▶ default: fast recovery (FR)
else
                                          ▶ additive increase (AI)
  if N_{AI} < counter \le N_{HAI} then
    target\_rate \leftarrow target\_rate + \delta
  if counter > N_{HAI} then
                                    ▶ hyper-active increase (HAI)
    target\_rate \leftarrow target\_rate + (counter - N_{HAI}) \cdot \delta
  rate \leftarrow \frac{rate + target\_rate}{2}
  counter \leftarrow counter + 1
```

- Receiver-driven
 - PCCB: throughput mode
 - PCCP: latency mode

PCCB

- Receivers use the Max-min fairness to determine the rate
- Senders apply an RCP-like approach to

PCCP

```
Input: delay\_in
delay \leftarrow EWMA(delay, delay\_in)
if delay > threshold then \Rightarrow multiplicative decrease (MD)
\begin{vmatrix} rate \leftarrow (1 - \beta \cdot (1 - \frac{threshold}{delay})) \cdot rate
counter \leftarrow 0 \qquad \Rightarrow \text{ enter fast recovery (FR)}
target\_rate \leftarrow rate
```

Target: cause #2

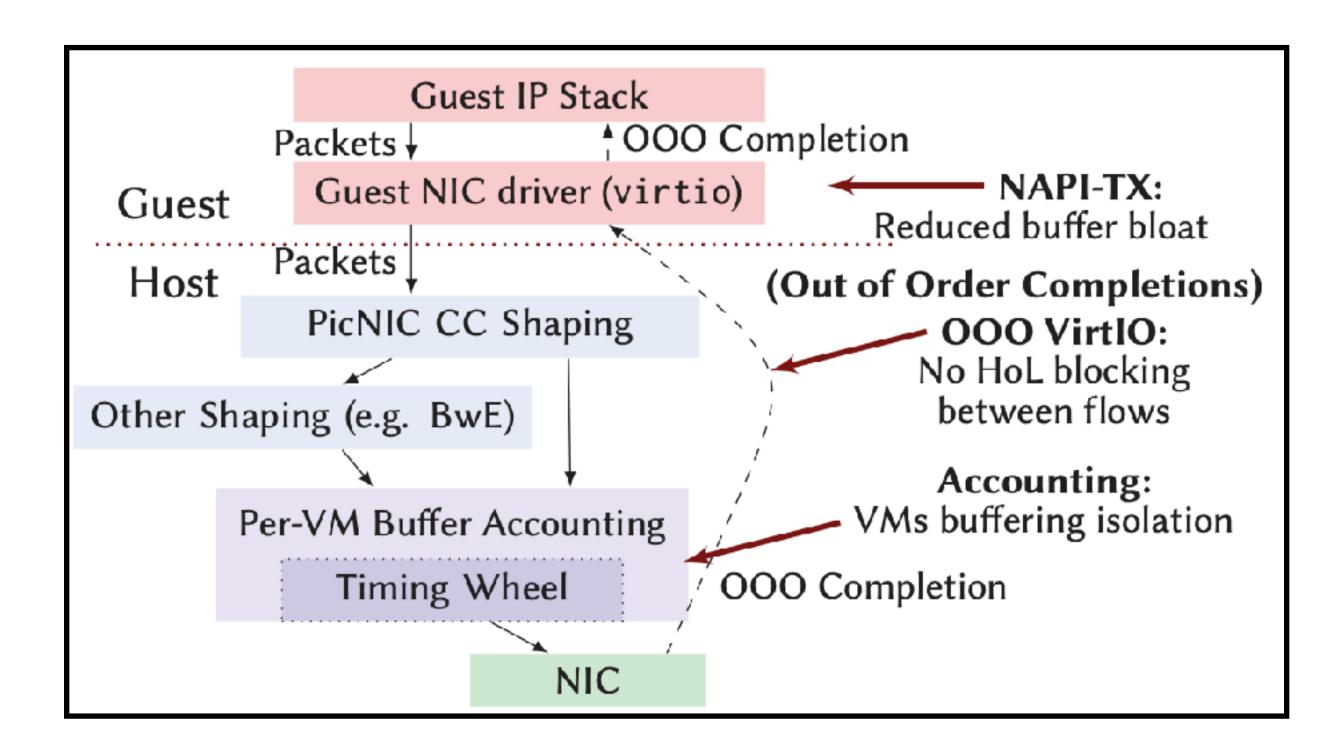
Benefits: bandwidth, delay, and drops

Technique #4: Sender-side Admission Control

- Smart traffic shaper
 - Per-VM packet counting
 - Backpressure to guest OS

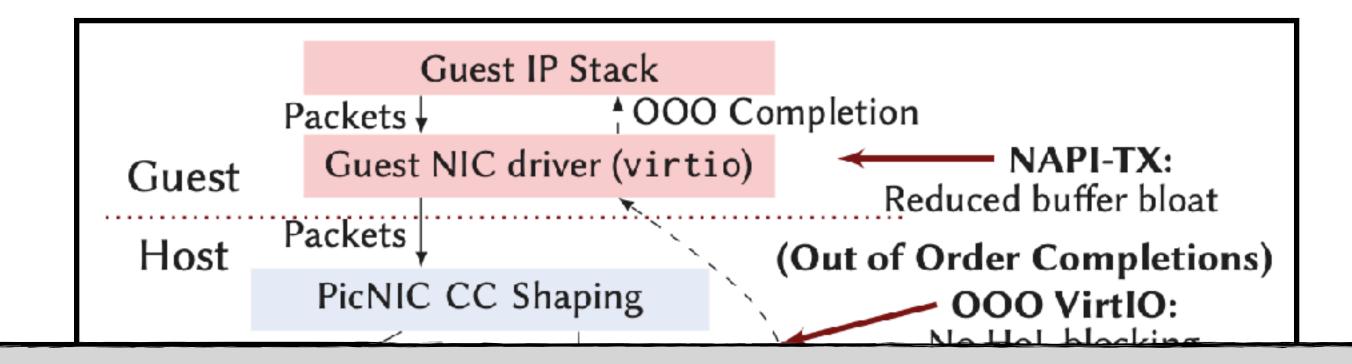
Technique #4: Sender-side Admission Control

- Smart traffic shaper
 - Per-VM packet counting
 - Backpressure to guest OS



Technique #4: Sender-side Admission Control

- Smart traffic shaper
 - Per-VM packet counting
 - Backpressure to guest OS



Target: issue #1

Benefits: bandwidth, delay, and drops

NIC

Network performance isolation requires considering (1) the entire communication path; (2) the interaction between network with CPU/memory.

Summary

- Today
 - Network virtualization in data center networks (II)

- Next topic: SDN and programmable networks
 - Ethane (Sigcomm'07) and OpenFlow (Sigcomm CCR'08)
 - RMT (Sigcomm'13) and AFQ (NSDI'18)
 - AccelNet (NSDI'18) and iPipe (Sigcomm'19)