# RDMA in Data Centers: Looking Back and Looking Forward

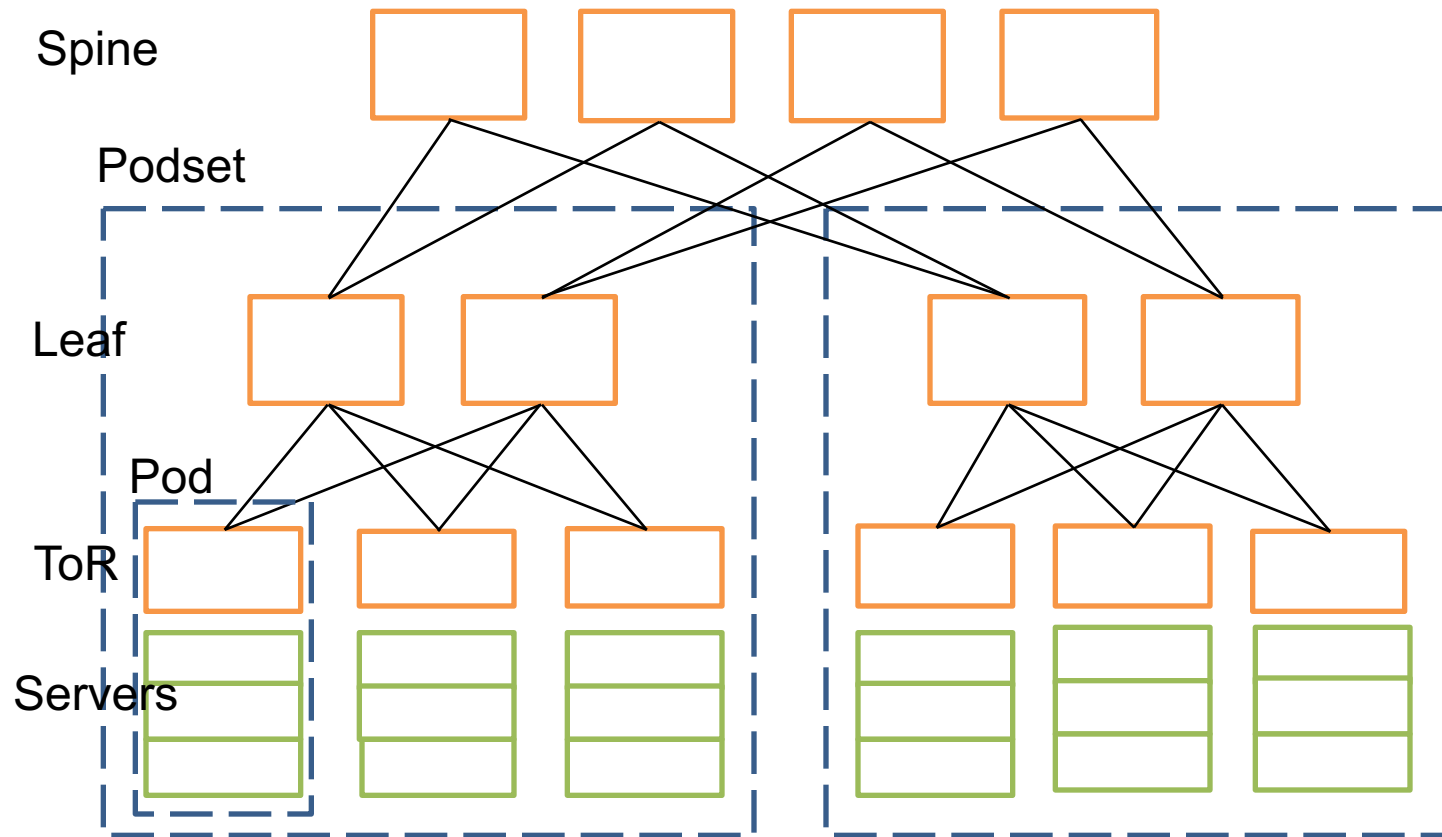Chuanxiong Guo

Microsoft Research

Data Centers

Data Centers

# Data center networks (DCN)

- Cloud scale services: IaaS, PaaS, Search, BigData, Storage, Machine Learning, Deep Learning

- Services are latency sensitive or bandwidth hungry or both

- Cloud scale services need cloud scale computing and communication infrastructure
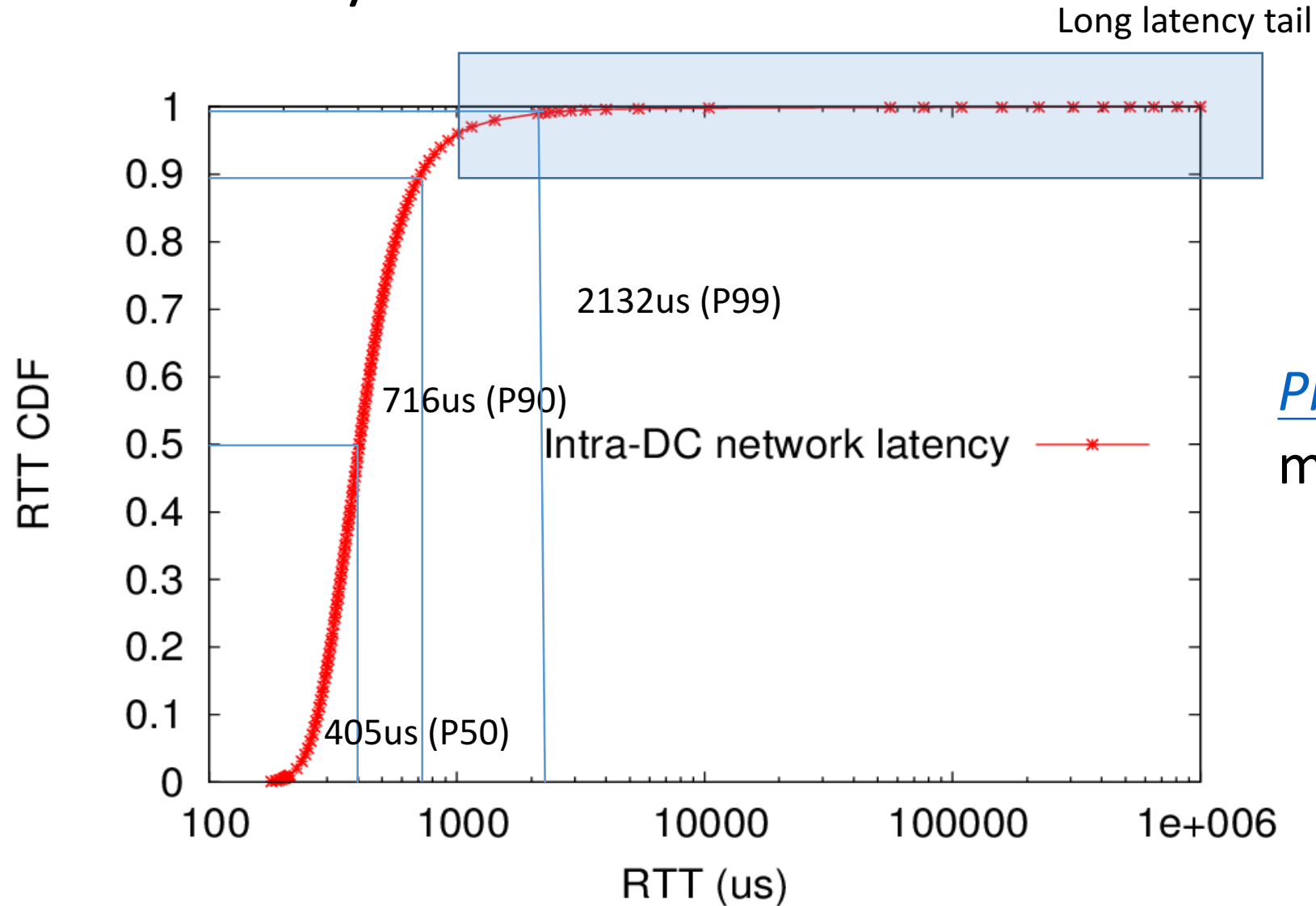
# Data center networks (DCN)

Spine

Podset

Leaf

Pod

ToR

Servers

- Single ownership
- Large scale
- High bisection bandwidth
- Commodity Ethernet switches
- TCP/IP protocol suite
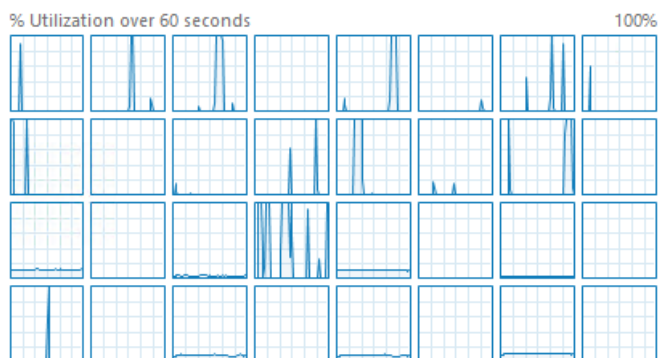
6

*But TCP/IP is not doing well*

# TCP latency



Long latency tail

2132us (P99)

716us (P90)

Intra-DC network latency

405us (P50)

RTT CDF

RTT (us)

*Pingmesh*
measurement results

# TCP processing overhead (40G)

Sender

Receiver

8 tcp connections

40G NIC

Send 70.0 Mbps
Receive 39.4 Gbps

9

*An RDMA renaissance story*

Virtual Interface
Architecture Spec
1.0        1997

*Infiniband Architecture
Spec*
1.0        2000
1.1        2002
1.2        2004
1.3        2015
RoCE      2010
RoCEv2 2014

# RDMA

- Remote Direct Memory Access (RDMA): Method of accessing memory on a remote system **without** interrupting the processing of the CPU(s) on that system

- RDMA offloads packet processing protocols to the NIC

- RDMA in Ethernet based data centers

# RoCEv2: RDMA over Commodity Ethernet



- RoCEv2 for Ethernet based data centers
- RoCEv2 encapsulates packets in UDP
- OS kernel is not in data path
- NIC for network protocol processing and message DMA

# RDMA benefit: latency reduction

| Msg size | TCP P50 (us) | TCP P99 (us) | RDMA P50 (us) | RDMA P99 (us) |
|---|---|---|---|---|
| 1KB | 236 | 467 | 24 | 40 |
| 16KB | 580 | 788 | 51 | 117 |
| 128KB | 1483 | 2491 | 247 | 551 |
| 1MB | 5290 | 6195 | 1783 | 2214 |

- For small msgs (<32KB), OS processing latency matters
- For large msgs (100KB+), speed matters

# RDMA benefit: CPU overhead reduction



Sender

Receiver

One ND connection

40G NIC

37Gb/s goodput

# RDMA benefit: CPU overhead reduction

Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz, two sockets 28 cores



RDMA: Single QP, 88 Gb/s, 1.7% CPU

TCP: Eight connections, 30-50Gb/s,
Client: 2.6%, Server: 4.3% CPU

# *RoCEv2 needs a lossless Ethernet network*

- PFC for hop-by-hop flow control
- DCQCN for connection-level congestion control

# Priority-based flow control (PFC)

- Hop-by-hop flow control, with eight priorities for HOL blocking mitigation

- The priority in data packets is carried in the VLAN tag or DSCP

- PFC pause frame to inform the upstream to stop

- PFC causes HOL and colleterial damage

Egress port

Data packet

Ingress port

p0

p1

p7

p0

p1

PFC pause frame
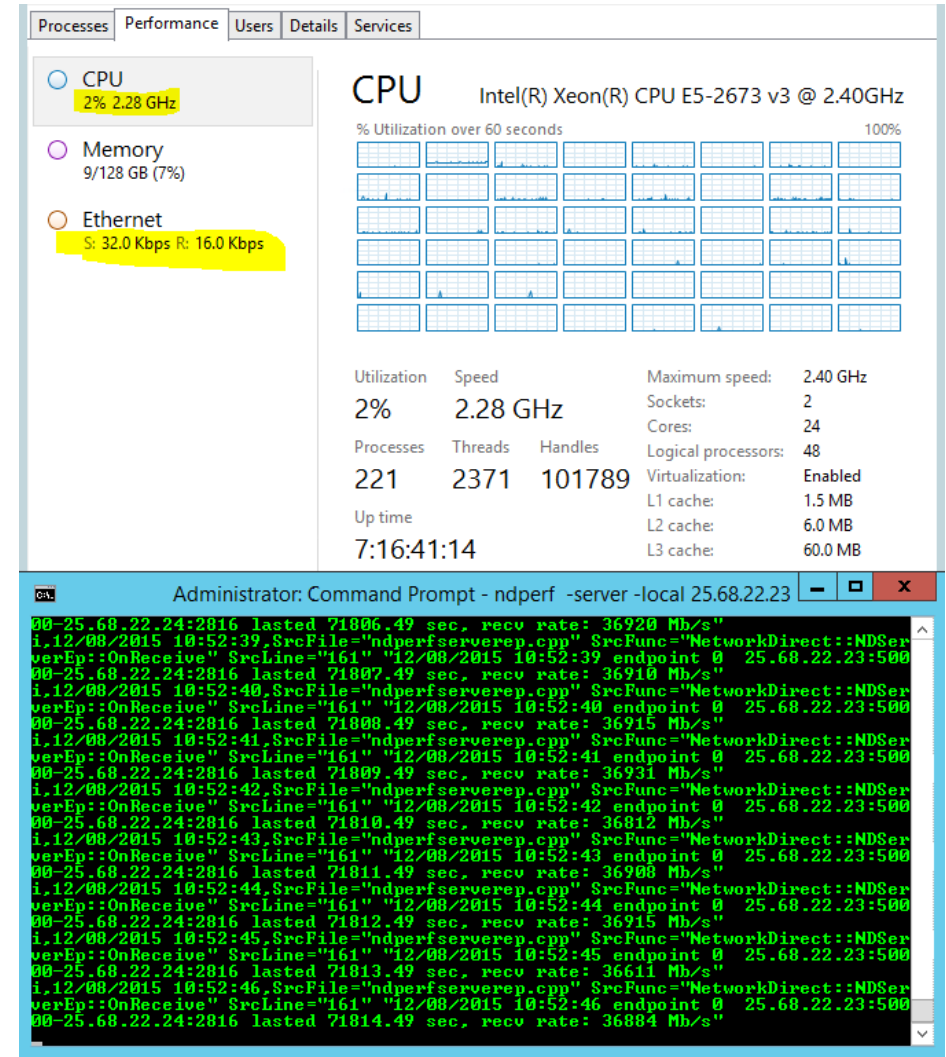
p0

p1

p7

XOFF threshold

| DMAC | SMAC | VLAN TAG | EtherType | Payload | FCS |
|------|------|----------|-----------|---------|-----|

| 16bits | 3bits | 1bits | 12bits |
|--------|-------|--------|--------|
| TPID 0x8100 | PCP | DEI | VID |
| | TCI | | |

Data packet

PFC pause frame

| DMAC 01-80-C2-00-00-01 | SMAC | EtherType 88-08 | OpType 01-01 | Priority-enable-vector (2B) | Time-vector (16B) | Padding | FCS |
|------------------------|------|-----------------|--------------|------------------------------|--------------------|---------|-----|

# DCQCN



Sender NIC
Reaction Point
(RP)

Switch
Congestion Point
(CP)

Receiver NIC
Notification Point
(NP)

> DCQCN = Keep PFC + Use ECN +
> hardware rate-based congestion control

- CP: Switches use ECN for packet marking
- NP: periodically check if ECN-marked packets arrived, if so, notify the sender
- RP: adjust sending rate based on NP feedbacks

# *The lossless requirement causes safety and performance challenges*

- RDMA transport livelock
- PFC deadlock

- PFC pause frame storm
- Slow-receiver symptom

# RDMA transport livelock

Switch

Pkt drop rate 1/256

Sender

Receiver

Sender — Receiver

RDMA Send 0

RDMA Send 1

RDMA Send N+1

RDMA Send N+2

NAK N

RDMA Send *0*

RDMA Send *1*

RDMA Send *2*

Go-back-0

Sender — Receiver

RDMA Send 0

RDMA Send 1

RDMA Send N+1

RDMA Send N+2

NAK N

RDMA Send *N*

RDMA Send *N+1*

RDMA Send *N+2*

Go-back-N   21

# PFC deadlock

- Our data centers use Clos network
- Packets first travel up then go down
- No cyclic buffer dependency for up-down routing -> no deadlock
- But we did experience deadlock!



Spine

Podset

Leaf

Pod

ToR

Servers

# PFC deadlock

- Preliminaries
  - ARP table: IP address to MAC address mapping
  - MAC table: MAC address to port mapping
  - If MAC entry is missing, packets are flooded to all ports

Input

ARP table

| IP | MAC | TTL |
|----|-----|-----|
| IP0 | MAC0 | 2h |
| IP1 | MAC1 | 1h |

MAC table

| MAC | Port | TTL |
|-----|------|-----|
| MAC0 | Port0 | 10min |
| MAC1 | - | - |

Dst: IP1

Output

23

# PFC deadlock

La

Lb

p0  p1 (La)

p0  p1 (Lb)

Path: {S1, T0, La, T1, S3}

Path: {S1, T0, La, T1, S5}

Path: {S4, T1, Lb, T0, S2}

PFC pause frames

② ③ ① ④

Congested port

p2 p3 (T0)

p3 p4 (T1)

T0

T1

p0 p1 (T0)

p0 p1 p2 (T1)

Packet drop

Ingress port

Egress port

Dead server

PFC pause frames

Server S1 S2 ✖

S3 ✖ S4 S5

24

# PFC deadlock

- The PFC deadlock root cause: the interaction between the PFC flow control and the Ethernet packet flooding

- Solution: drop the lossless packets if the ARP entry is incomplete

- Recommendation: do not flood or multicast for lossless traffic

# *Tagger*: practical PFC deadlock prevention



- Concept: Expected Lossless Path (ELP) to decouple Tagger from routing

- Strategy: move packets to different lossless queue before CBD forming

- Tagger Algorithm works for general network topology
- Deployable in existing switching ASICs

# NIC PFC pause frame storm



Spine layer

Podset 0                    Podset 1

Leaf layer

ToRs

0  1  2  3    4  5  6  7    servers

Malfunctioning NIC

- A malfunctioning NIC may block the whole network

- PFC pause frame storms caused several incidents

- Solution: watchdogs at both NIC and switch sides to stop the storm

27

# The slow-receiver symptom

- ToR to NIC is 40Gb/s, NIC to server is 64Gb/s

- But NICs may generate large number of PFC pause frames

- Root cause: NIC is resource constrained

- Mitigation
  - Large page size for the MTT (memory translation table) entry
  - Dynamic buffer sharing at the ToR

Server

CPU   DRAM

PCIe
Gen3 8x8 64Gb/s

WQEs   MTT

QPC

QSFP
40Gb/s

ToR

NIC

Pause frames

*Deployment experiences and lessons learned*

# Latency reduction

- RoCEv2 deployed in Bing world-wide for two and half years

- Significant latency reduction

- Incast problem solved as no packet drops

# RDMA throughput



- Using two podsets each with 500+ servers
- 5Tb/s capacity between the two podsets

- Achieved 3Tb/s inter-podset throughput
- Bottlenecked by ECMP routing
- Close to 0 CPU overhead

31

# Latency and throughput tradeoff



- RDMA latencies increase as data shuffling started

- Low latency vs high throughput

# Lessons learned

- Providing lossless is hard!
- Deadlock, livelock, PFC pause frames propagation and storm did happen
- Be prepared for the unexpected
  - Configuration management, latency/availability, PFC pause frame, RDMA traffic monitoring
- NICs are the key to make RoCEv2 work

*What's next?*

*Applications*

- RDMA for X (Search, Storage, HFT, DNN, etc.)

*Architectures*

- Software vs hardware
- Lossy vs lossless network
- RDMA for heterogenous computing systems

*Technologies*

- RDMA programming
- RDMA virtualization
- RDMA security
- Inter-DC RDMA

*Protocols*

- Practical, large-scale deadlock free network
- Reducing colleterial damage
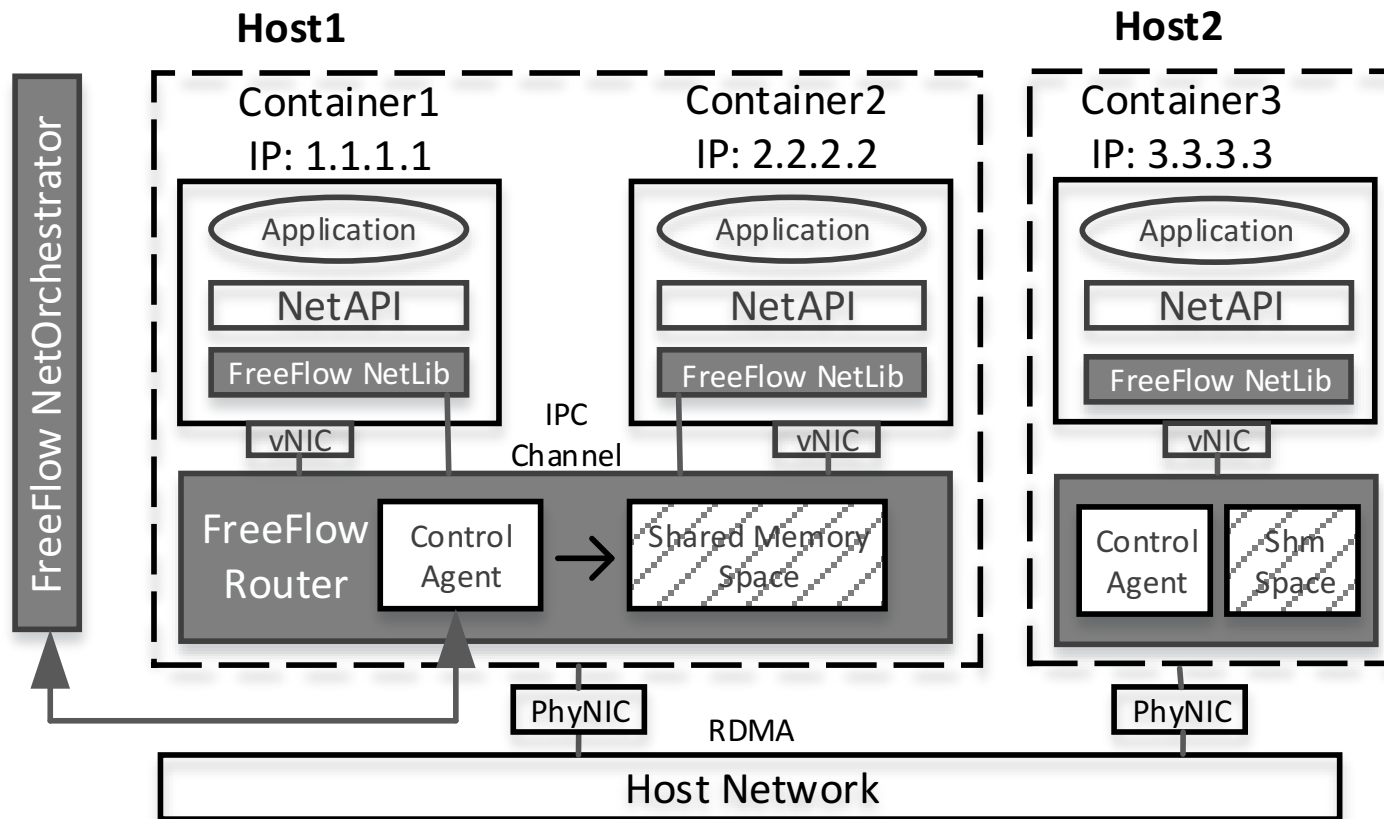
# *Will software win (again)?*

- Historically, software based packet processing won (multiple times)
  - TCP processing overhead analysis by David Clark, et al.
  - Non of the stateful TCP offloading took off (e.g., TCP Chimney)
- The story is different this time
  - Moore's law is ending
  - Accelerators are coming
  - Network speed keep increasing
  - Demands for ultra low latency are real

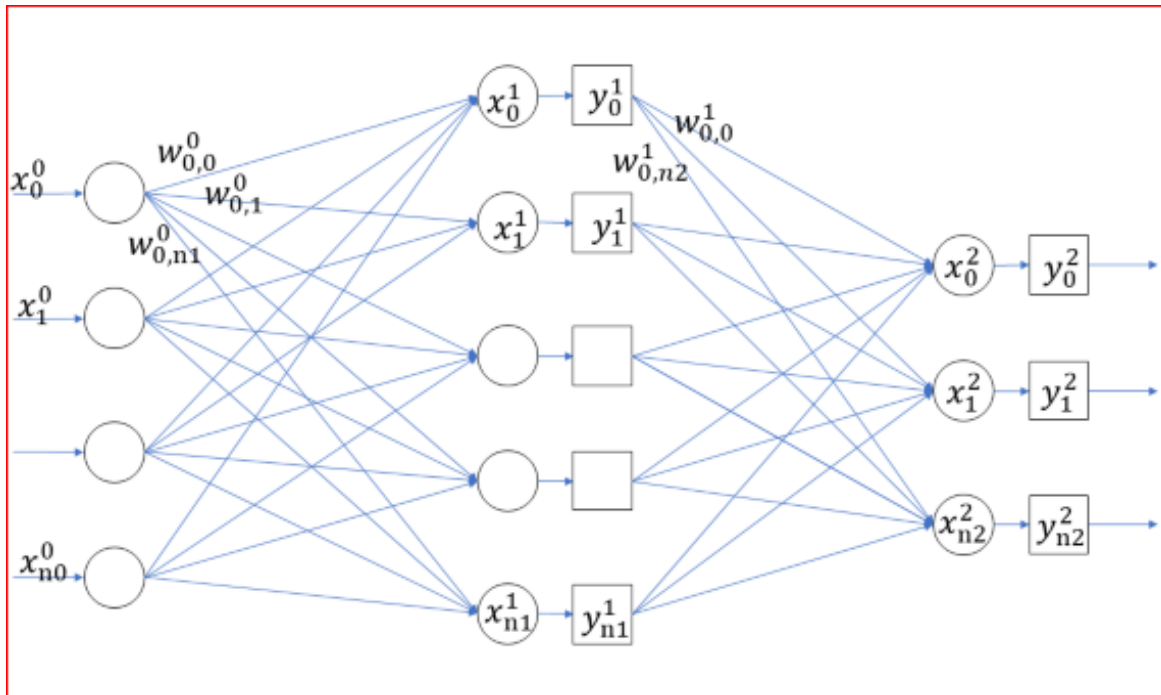# *Is lossless mandatory for RDMA?*

- There is no binding between RDMA and lossless network
- But implementing more sophisticated transport protocol in hardware is a challenge

# RDMA virtualization for the container networking



- A router acts as a proxy for the containers
- Shared memory for improved performance
- Zero copy possible

# *RDMA for DNN*



- TCP does not work for distributed DNN training
- For 16-GPU, 2-host speech training with CNTK, TCP communications dominant the training time (72%), RDMA is much faster (44%)

# RDMA Programming

- How many LOC for a "hello world" communication using RDMA?

- For TCP, it is 60 LOC for client or server code

- For RDMA, it is complicated ...
  - IBVerbs: 600 LOC
  - RCMA CM: 300 LOC
  - Rsocket: 60 LOC

# RDMA Programming

- Make RDMA programming more accessible
  - Easy-to-setup RDMA server and switch configurations
  - Can I run and debug my RDMA code on my desktop/laptop?
  - High quality code samples
- Loosely coupled vs tightly coupled (Send/Recv vs Write/Read)

# *Summary: RDMA for data centers!*

- RDMA is experiencing a renaissance in data centers
  - RoCEv2 has been running safely in Microsoft data centers for two and half years
- Many opportunities and interesting problems for high-speed, low-latency RDMA networking
- Many opportunities in making RDMA accessible to more developers

# *Acknowledgement*

*Questions?*