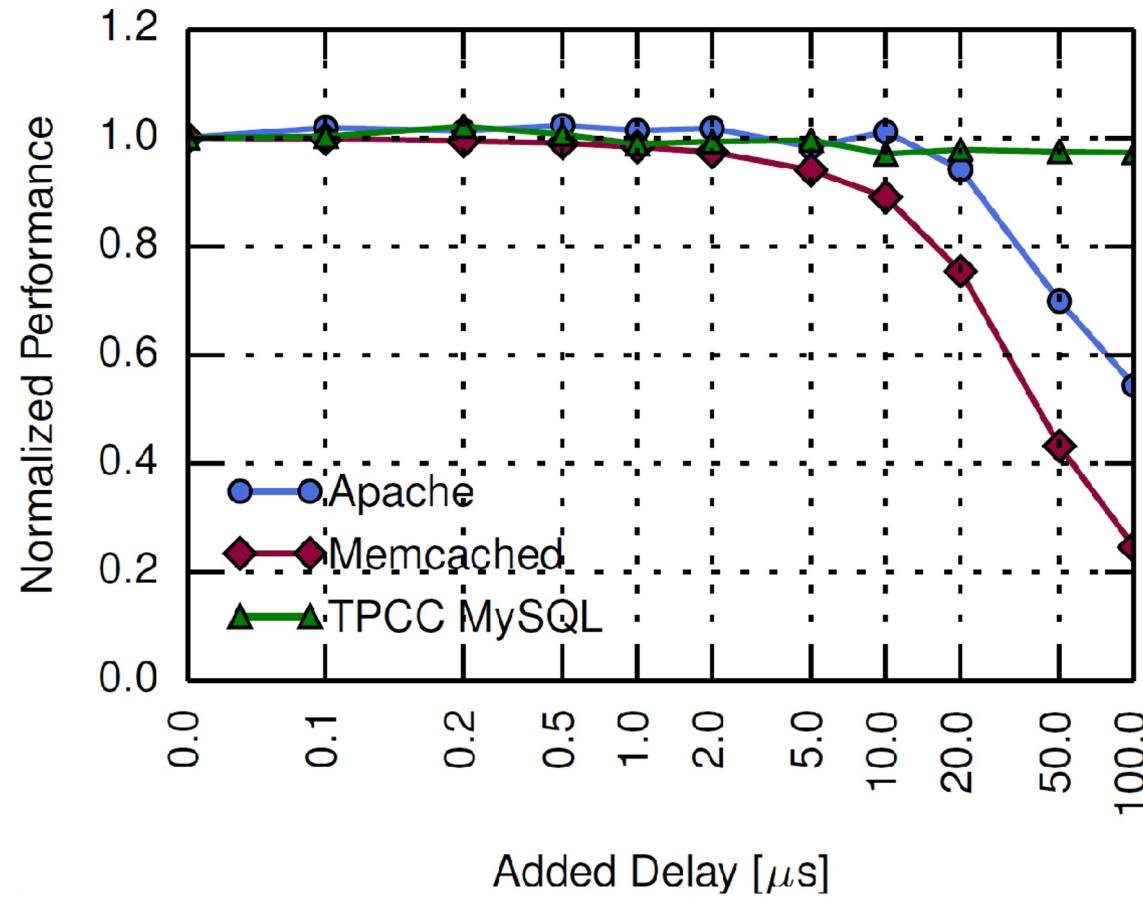


Where Has My Time Gone?

Noa Zilberman, Matthew Grosvenor, Diana Andreea Popescu, Neelakandan Manihatty-Bojan,
Gianni Antichi, Marcin Wójcik, Andrew W. Moore

Latency matters



It's Time For Low Latency + Low Variance!

2011

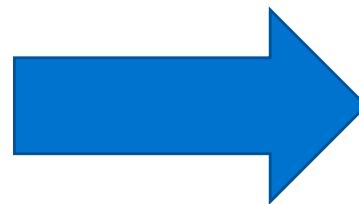
It's Time for Low Latency

Stephen M. Rumble, Diego Ongaro, Ryan Stutsman,
Mendel Rosenblum, and John K. Ousterhout
Stanford University

community has ignored network latency in the past, speed-of-light delays and unoptimized network hardware make round-trip times impossible. In years datacenters will be dominated by Ethernet. Without the burden of the datacenter campus and network devices, it will be up to us to take advantage of this benefit through application researchers must lead the charge to push the boundaries of low-latency communication.

Component	Delay	Round-Trip
Network Switch	10-30 μ s	100-300 μ s
Network Interface Card	2.5-32 μ s	10-128 μ s
OS Network Stack	15 μ s	60 μ s
Speed of Light (in Fiber)	5ns/m	0.6-1.2 μ s

Table 2: Factors that contribute to latency in TCP datacenter communication. “Delay” indicates the cost of a single traversal of the component, and “Round-Trip” indicates the total impact on round-trip time. Messages typically traverse 5 switches in each direction in a large datacenter network and must pass through the OS stack 4 times.



2016

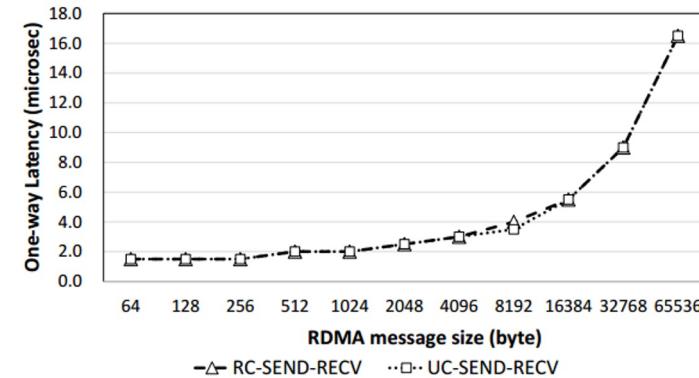


Fig. 4: Median one-way latency of RoCE RC and UC transport types. “Exploring Low-latency Interconnect for Scaling Out Software Routers”, Ma, Kim, and Moon

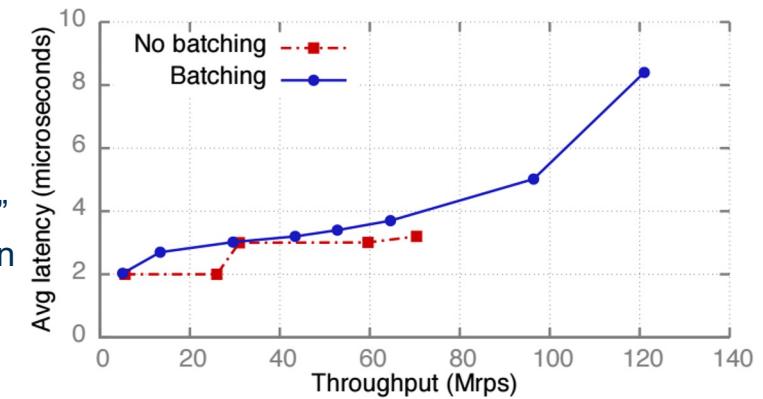


Figure 8: Impact of response batching on Spec-S0 latency

It's Time For Low Latency + Low Variance!

2011

It's Time for Low Latency

Stephen M. Rumble, Diego Ongaro, Ryan Stutsman,
Mendel Rosenblum, and John K. Ousterhout
Stanford University

Community has ignored network latency in the past, speed-of-light delays and unoptimized network hardware make round-trip times impossible. In years datacenters will be dominated by Ethernet. Without the burden of the datacenter campus and network devices, it will be up to us to take advantage of this benefit through application researchers must lead the charge to push the boundaries of low-latency communication.

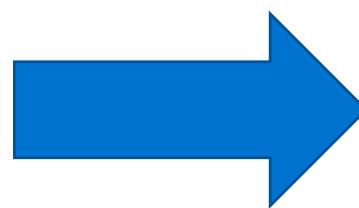
	1983	2011	Improved
CPU Speed	1x10Mhz	4x3GHz	> 1,000x
Memory Size	$\leq 2\text{MB}$	8GB	$\geq 4,000\text{x}$
Disk Capacity	$\leq 30\text{MB}$	2TB	> 60,000x
Net Bwidth	3Mbps	10Gbps	> 3,000x
RTT	2.54ms	$80\mu\text{s}$	32x

Table 1: Network latency has improved far more slowly over the last three decades than other performance metrics for commodity computers. The V Distributed System [5] achieved round-trip RPC times of 2.54ms. Today, a pair of modern Linux servers require $80\mu\text{s}$ for 16-byte RPCs over TCP with 10Gb Ethernet.

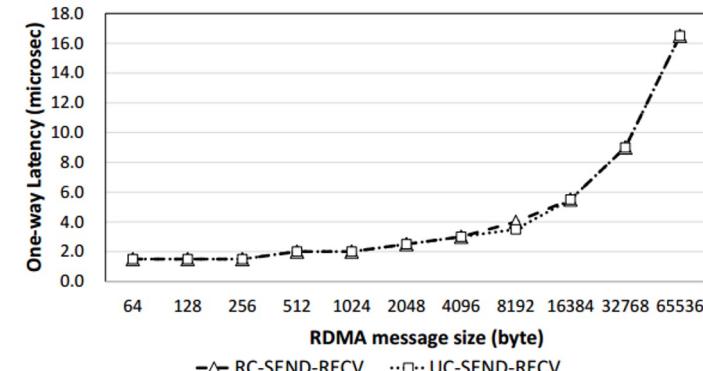
Component	Delay	Round-Trip
Network Switch	$10\text{-}30\mu\text{s}$	$100\text{-}300\mu\text{s}$
Network Interface Card	$2.5\text{-}32\mu\text{s}$	$10\text{-}128\mu\text{s}$
OS Network Stack	$15\mu\text{s}$	$60\mu\text{s}$
Speed of Light (in Fiber)	5ns/m	$0.6\text{-}1.2\mu\text{s}$

Table 2: Factors that contribute to latency in TCP datacenter

It Usually Works



2016



The Unavoidable Latency

I will talk about:

- The **essential** latency contributions
- **Commodity** hardware, standard coding
- Set a (reproducible) **baseline** for research

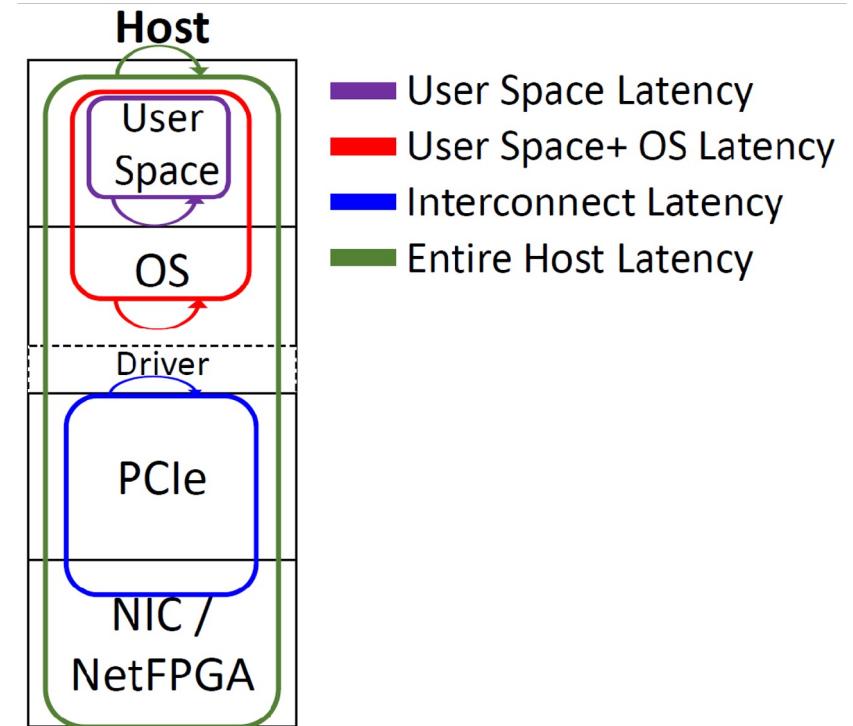
I will not talk about:

- TCP, DCTCP, MPTCP etc.
- Congestion, Buffer bloat, In-cast, etc.

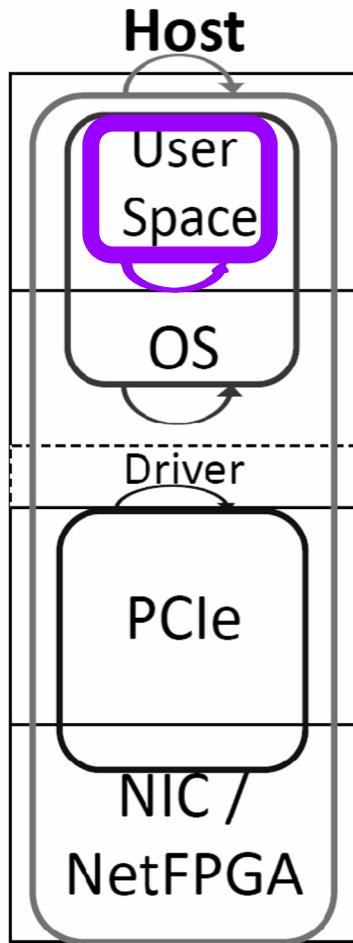


Methodology

- Breakdown overall system latency into **component contributions**
- Use decompositional analysis to report component distributions e.g. min, median, 99.9th, max

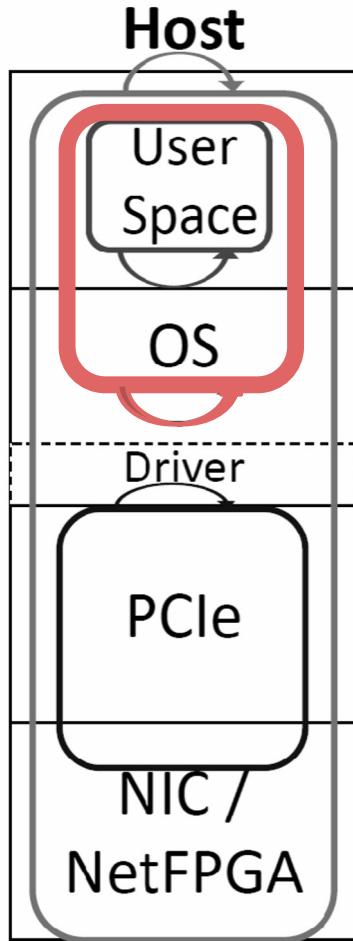


Time Stamp Counter (TSC) Measurement Accuracy



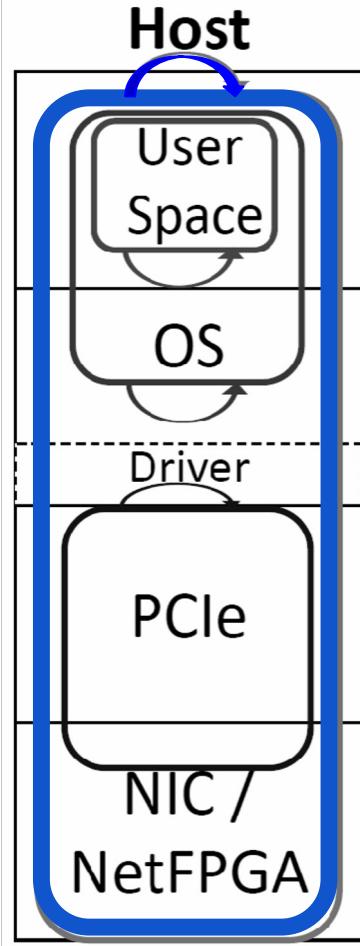
Experiment	Minimum	Median	99.9th	Tail
TSC - From User Space	9 ns	10 ns	11 ns	49 µs
TSC - Kernel	9 ns	9 ns	9 ns	6.9 µs
TSC - Kernel Early Boot	7 ns	7 ns	7 ns	11 ns
TSC - From VM User Space	12 ns	12 ns	13 ns	64 ms

A Breakdown of Basic Latency Components: User space + OS



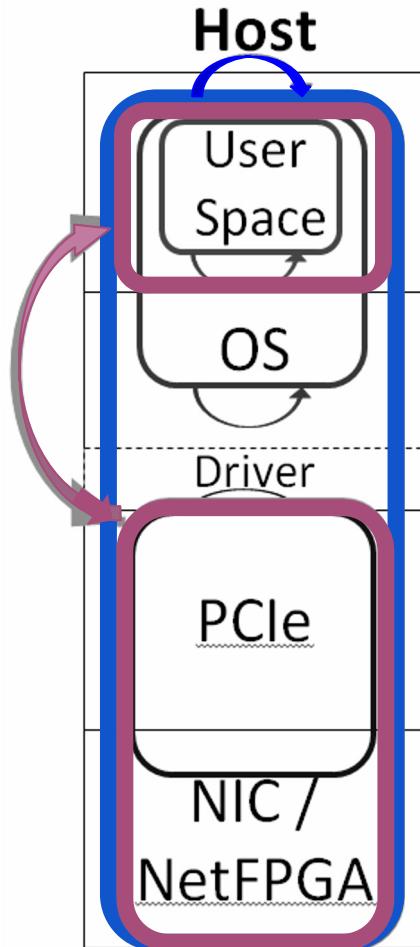
Experiment	Minimum	Median	99.9th	Tail
User space + OS (same core)	2 μ s	2 μ s	2 μ s	68 μ s
User space + OS (other core)	4 μ s	5 μ s	5 μ s	31 μ s

A Breakdown of Basic Latency Components: Host



Experiment	Minimum	Median	99.9th	Tail
Host	3.9 μ s	4.5 μ s	21 μ s	45 μ s

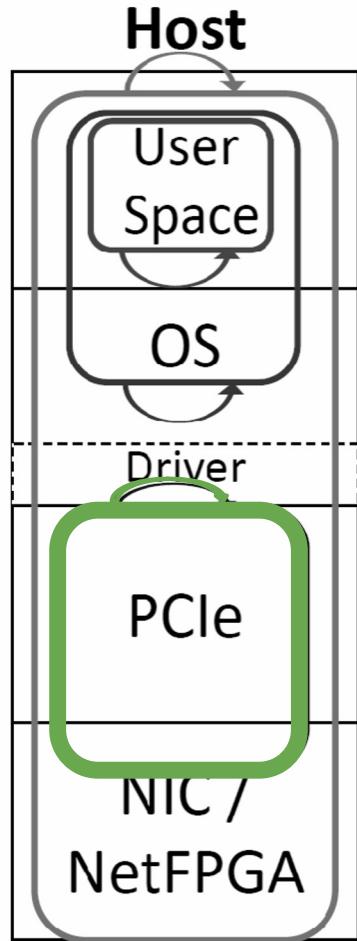
A Breakdown of Basic Latency Components: Host



Experiment	Minimum	Median	99.9th	Tail
Host	3.9 μ s	4.5 μ s	21 μ s	45 μ s
Kernel Bypass	0.89 μ s	0.94 μ s	1.1 μ s	5.4 μ s



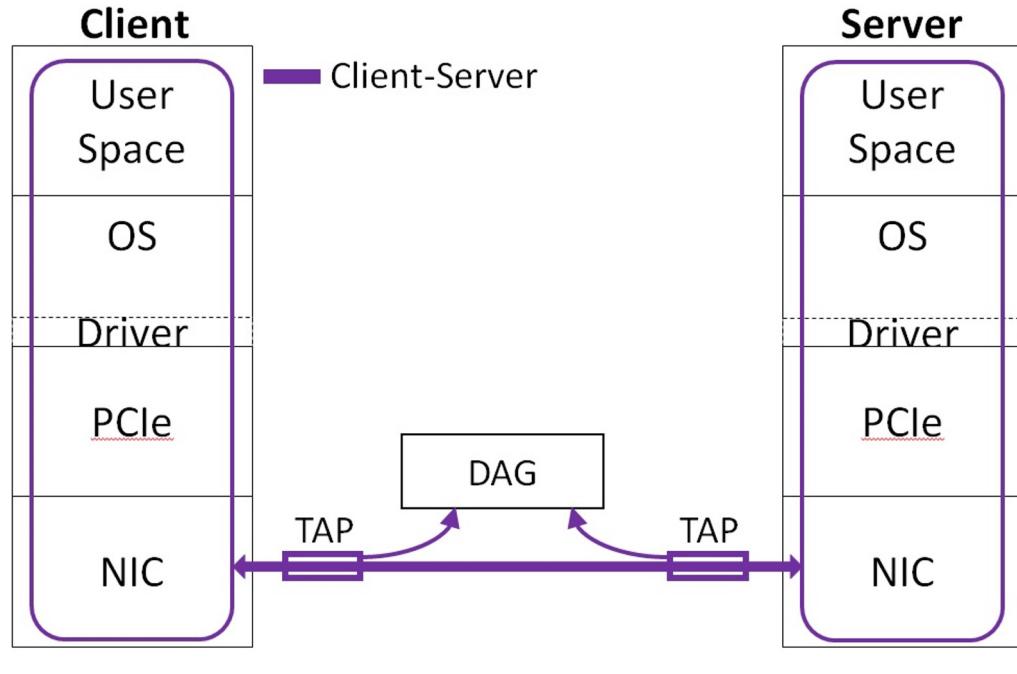
A Breakdown of Basic Latency Components: PCIe



Experiment	Minimum	Median	99.9th	Tail
Host	3.9 μ s	4.5 μ s	21 μ s	45 μ s
Kernel Bypass	0.89 μ s	0.94 μ s	1.1 μ s	5.4 μ s

Experiment	Minimum	Median	99.9 th	Tail
PCIe Interconnect (64B)	0.55 μs	0.57 μ s	0.59 μ s	0.6 μs
PCIe Interconnect (1536B)	0.97 μs	0.98 μ s	1.02 μ s	1.03 μs

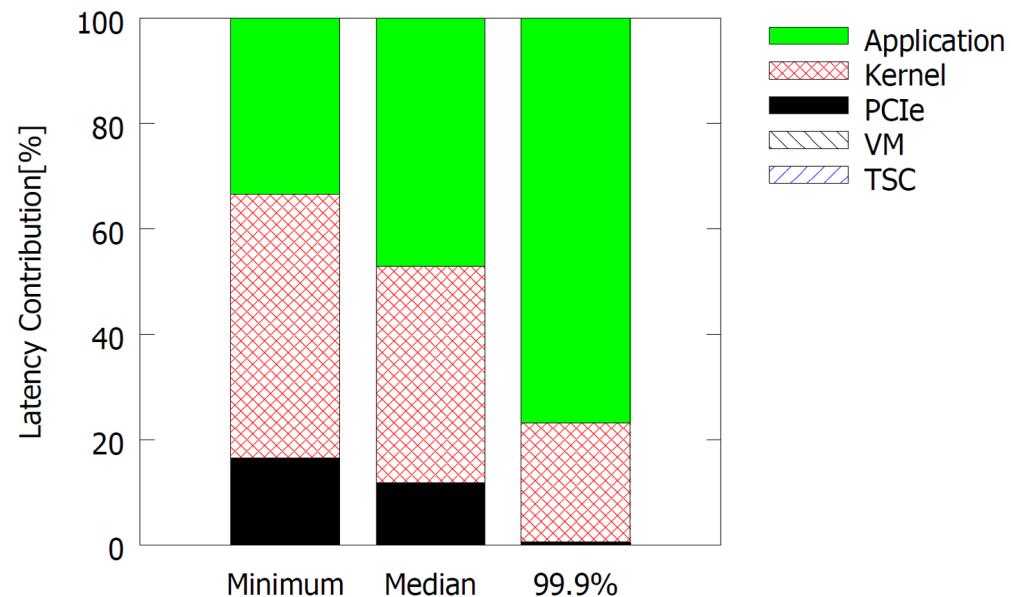
A Breakdown of Basic Latency Components: Client-Server



Experiment	Minimum	Median	99.9th	Tail
Client-Server (UDP)	7 μ s	9 μ s	107 μ s	203 μs
Client-Server (Memcached)	10 μ s	13 μ s	240 μ s	20.3 ms

Summary end host results

- All experiments bounded by application variance (TSC)
- PCIe latency variance is low

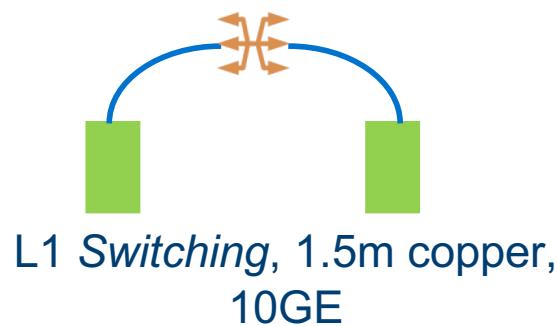


Basic Latency Components of the Network

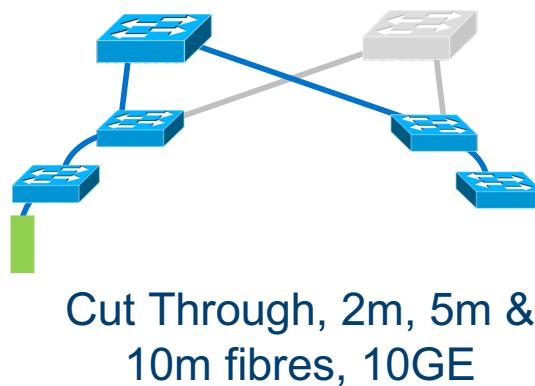
Single Rack



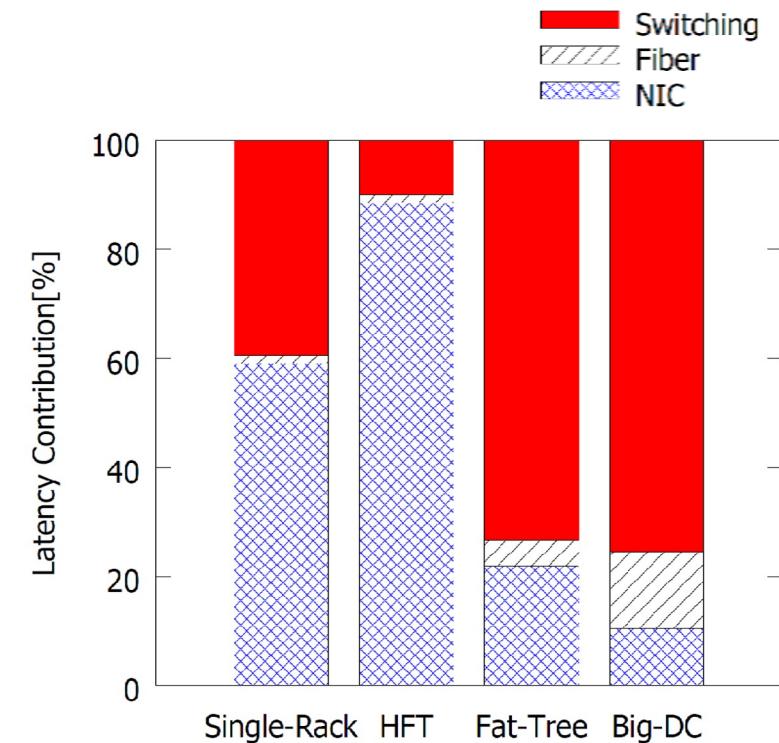
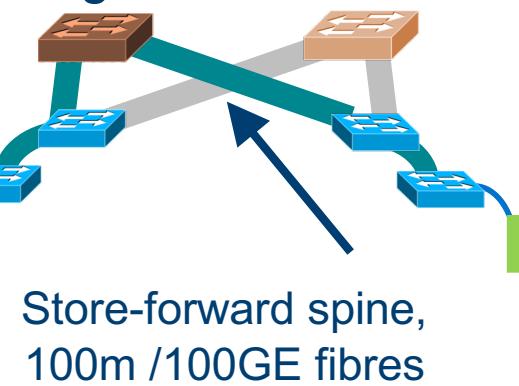
HFT



Fat tree



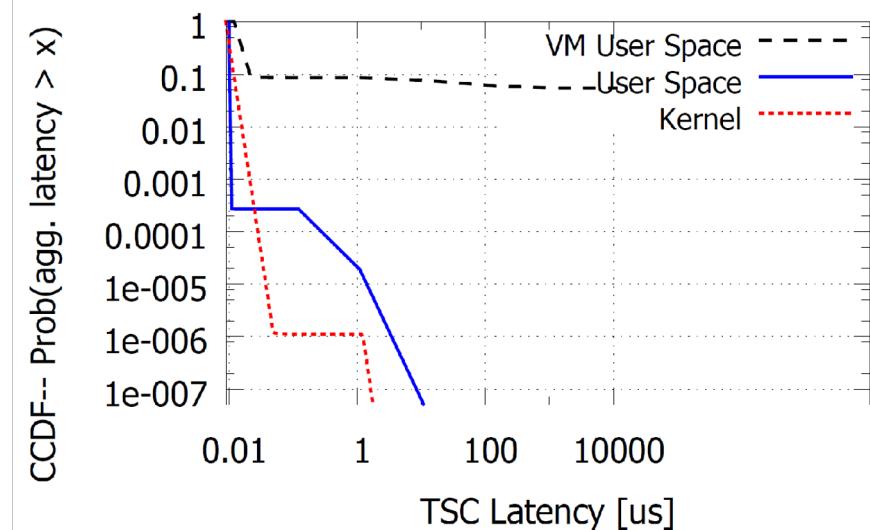
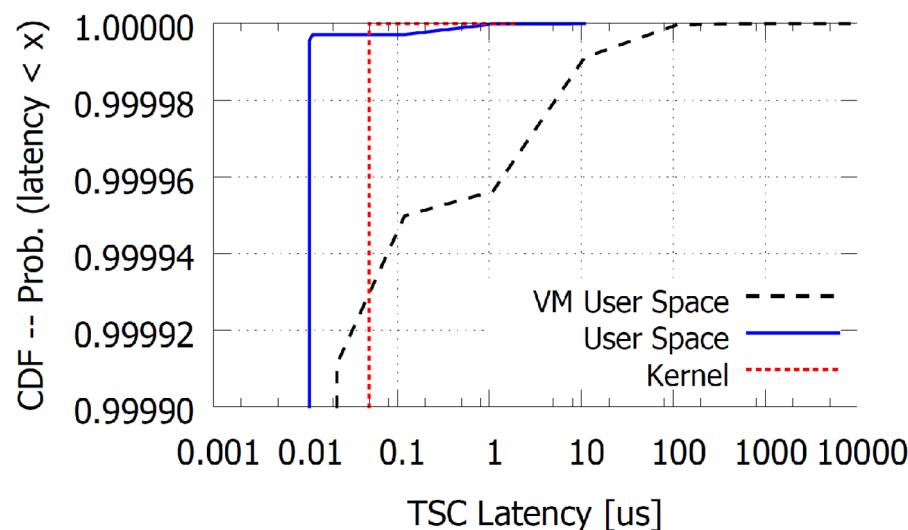
Big DC Fat tree



Based on median results

Tail Latency

- The long latency tail can take much more than you think
- For VM, events of 1ms or longer take almost 5% of the time.



The Good, The Bad and The Ugly

Good ($\leq 1\mu\text{s}$)

- Simple kernel and user space operations
- PCIe
- Single through-switch latency (no queueing)

Bad ($1\mu\text{s}-100 \mu\text{s}$)

- Sending packets over user space+OS
- Host latency and client-server latency
- Multi-stage network topology
 - RTT over 100m fibers

Ugly ($>100\mu\text{s}$)

- The far end of the tail, i.e. the “variance”
- Mostly in user space and within a VM.

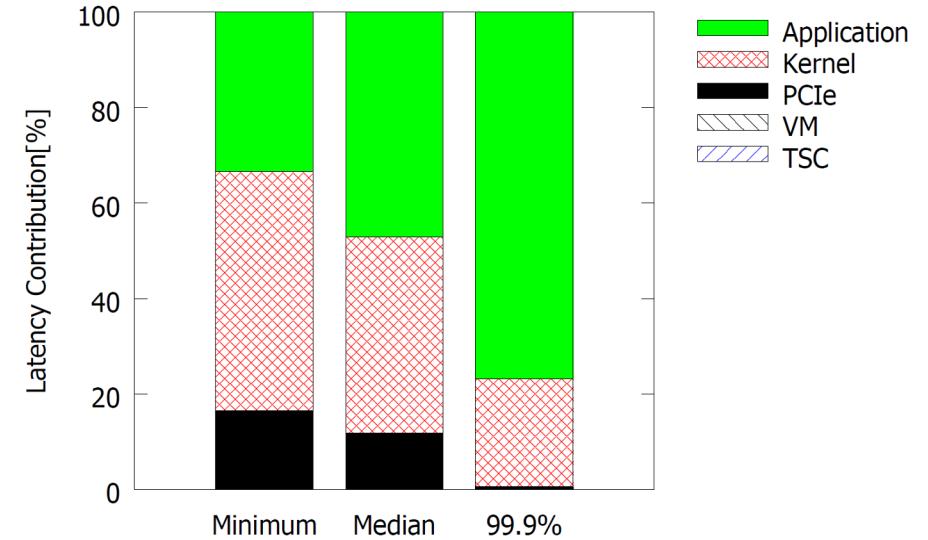
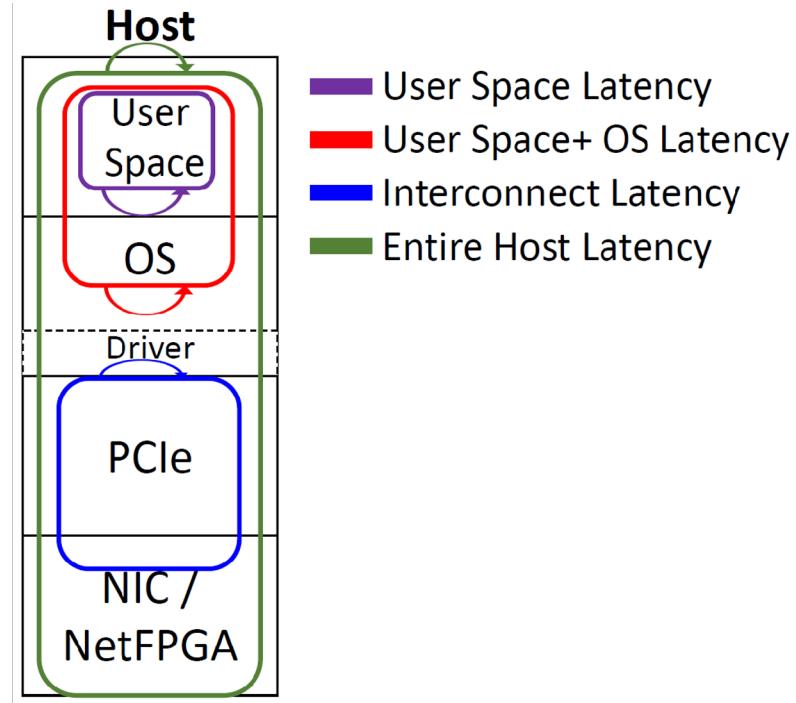
Summary

- End host latency has significantly improved over half a decade, but
 - Apps and VMs are the source of significant variance
 - Fibre length matters
 - “Variance” is no longer negligible
 - Fewer opportunities to improve end-host networking & interconnect
- The toolkit and reproduction environment are available!

www.lowlatencylab.org/data/pam2017



Questions?



Acknowledgements



UNIVERSITY OF
CAMBRIDGE



The Leverhulme Trust

EPSRC

Pioneering research
and skills



UNIVERSITY OF
CAMBRIDGE

Improving Tail Latency

```
1  while (!done)
2  {
3      //Read TSC once
4      do_rdtscp(tsc, cpu);
5      //If the gap between the current and the previous
6      //TSC value is above a certain threshold, save
7      //it
8      if ((tsc - last > threshold) && (cpu == lastcpu))
9          buffer[samples++] = tsc-last;
10     last = tsc;
11     lastcpu = cpu;
12 }
```

65% Min latency
50% Max latency

- Running in real time
- Pinning to a core
- Inhibiting interrupts
- Coding practices

```
1  while (!done)
2  {
3      //Read TSC twice, one immediately after the other
4      do_rdtscp(tsc, cpu);
5      do_rdtscp(tsc2, cpu2);
6      //If the gap between the two reads is above a
7      //certain threshold, save it
8      if ((tsc2 - tsc > threshold) && (cpu == cpu2))
9          buffer[samples++] = tsc2-tsc;
10 }
```

Application Only

- What we do: Read TSC
- Min: 9ns
- Median: 10ns
- 99.9%: 11ns
- Max: 10's to 100's of μ s
- 50-100 events/second > 1 μ s

