

# Validating and Tuning High-Throughput Ethernet Fabrics for AI and HPC

Alex Bortok, Lead Product Manager, Hyperscale Infrastructure  
Ankur Sheth, Senior Director, Strategic Projects

# A Brief History of Keysight



**1939–1998:**  
**Hewlett-Packard years**



**1999–2013:**  
**Agilent Technologies years**



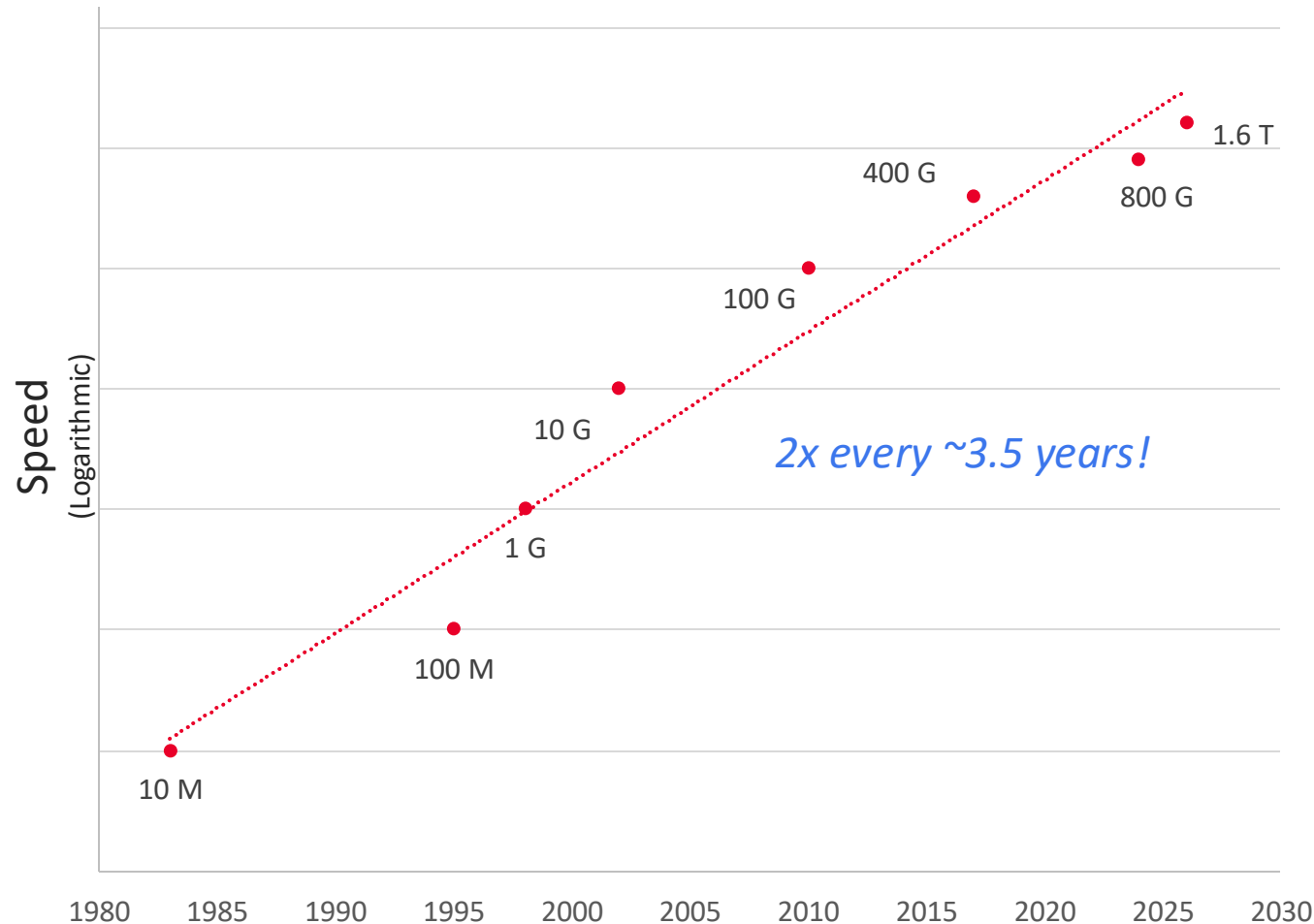
**2014+:**  
**Keysight years**



**2017:**  
**Keysight acquires Ixia**

# Ethernet

1983 - Present



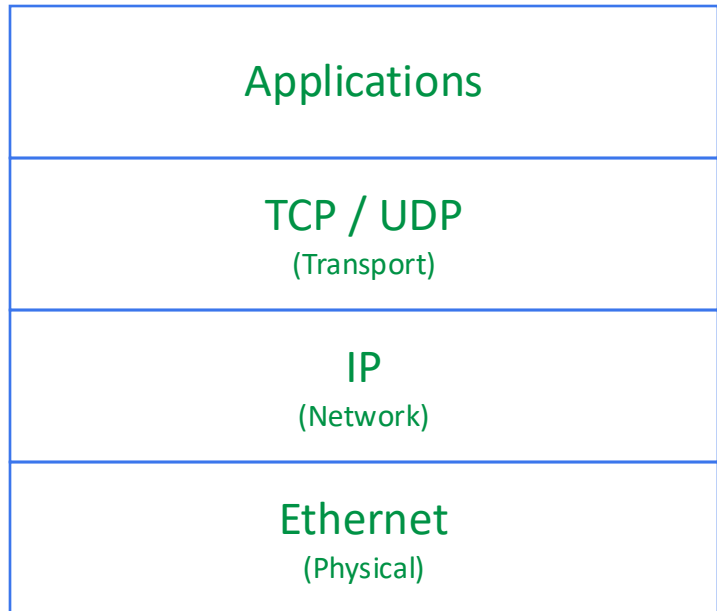
- Where?

- LANs & Internet (~1990s and onwards)
- Data Centers (~2010s and onwards)
- AI Fabrics (~2020s and onwards)

- Why?

- Standards Driven
- Versatile
- Large Eco-System
  - Vendors
  - Practitioners

# TCP/IP & Ethernet: Made for Each Other!



- Both Ethernet & IP

- ✗ Reliability

- ✗ Ordering

- ✗ Performance

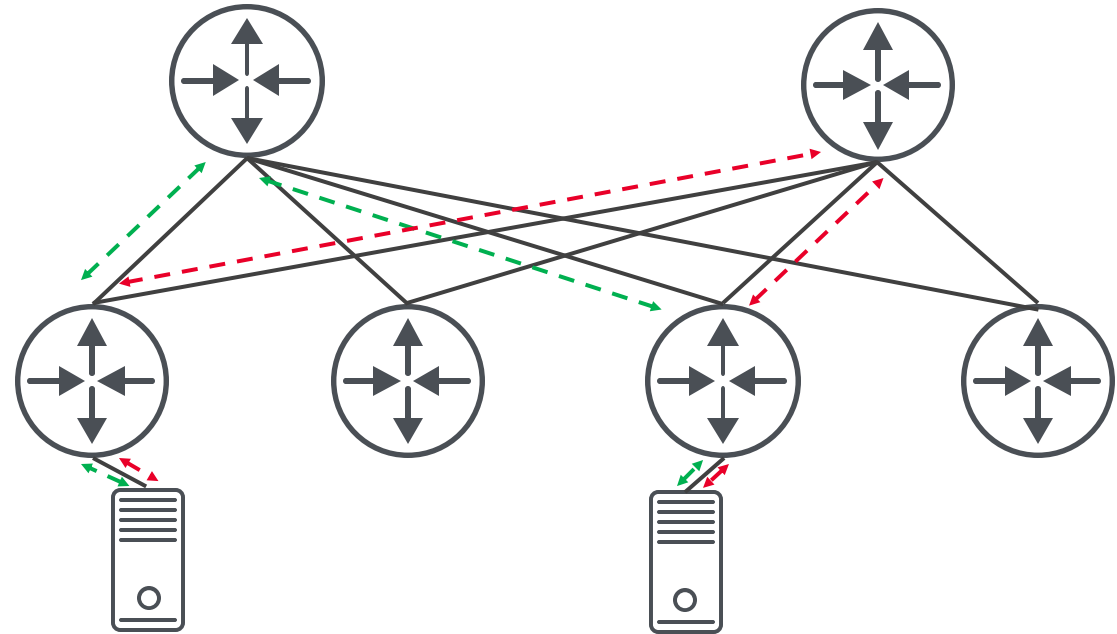
} Left to Upper Layers

- Perfect For Applications like:

- Web Browsing
  - Business Applications
  - Streaming

# Ethernet & The Data Center

- East  $\leftrightarrow$  West Traffic
- Small Bursty Flows
- 10,000s of Servers
- Resilience To Failures

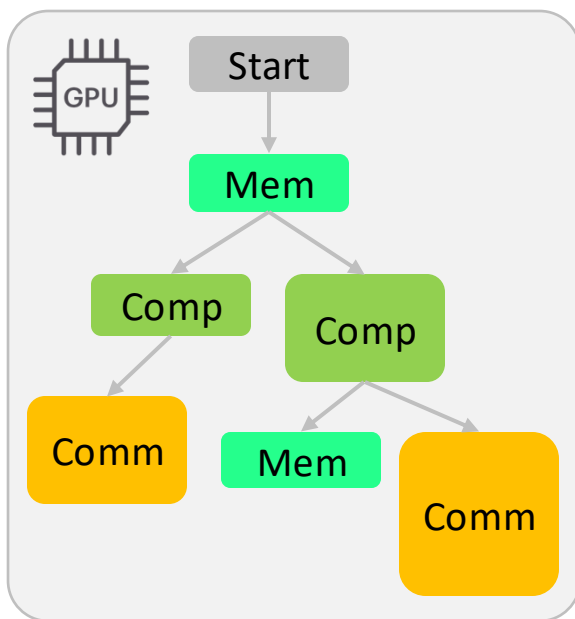


## Clos Topology

- Full Bi-Sectional Bandwidth
- Multiple Paths

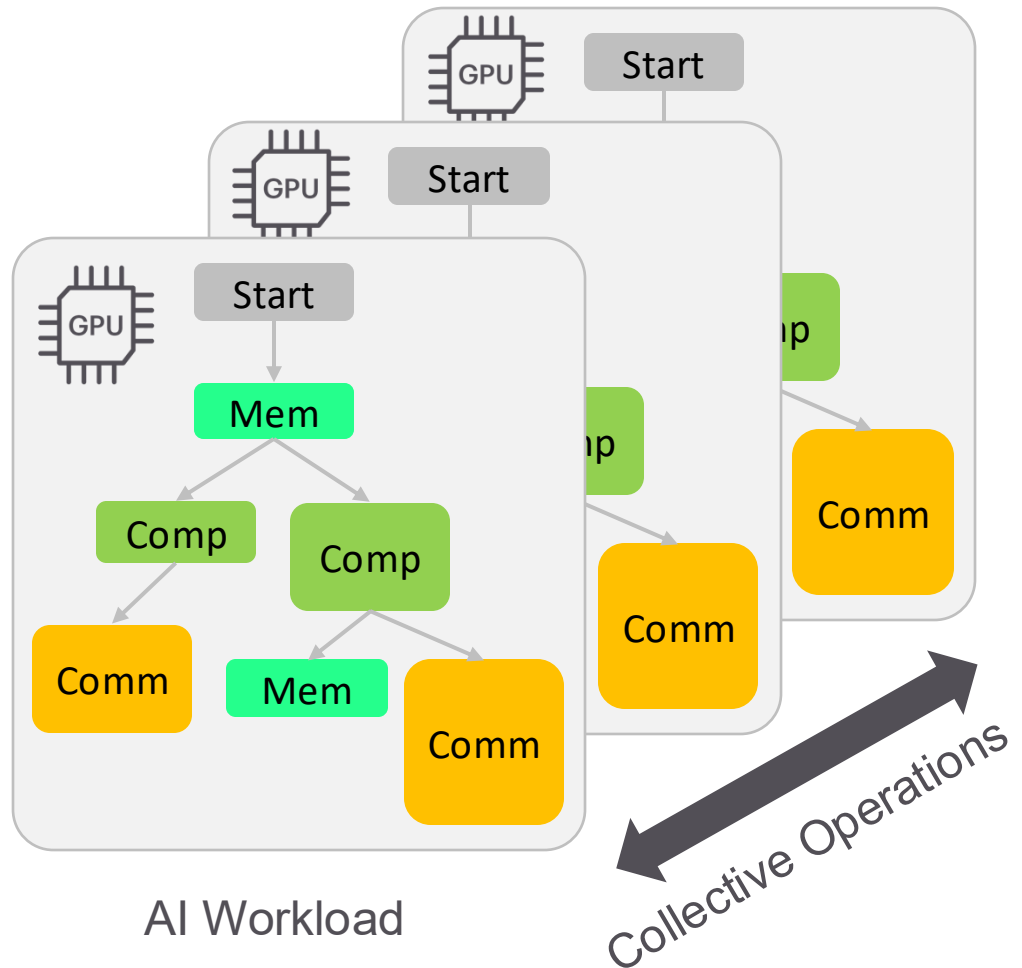
--> Equal Cost Multi Path

# Model Training & RDMA



AI Workload

# Model Training & RDMA



## Characteristics

- Reliable
- High Bandwidth
- Low Latency
- “Elephant” Flows
- Low Entropy
- RDMA : Best Performance over Reliable Interconnects
  - e.g. InfiniBand



# RDMA on Ethernet/IP



- iWARP : RDMA over TCP!
  - Host Networking Stacks → High Latencies
  - CPU Overhead
- RoCE: RDMA Over Converged Ethernet v1/v2
  - Simplified Protocol
  - RDMA Stack in NIC → No CPU Overhead



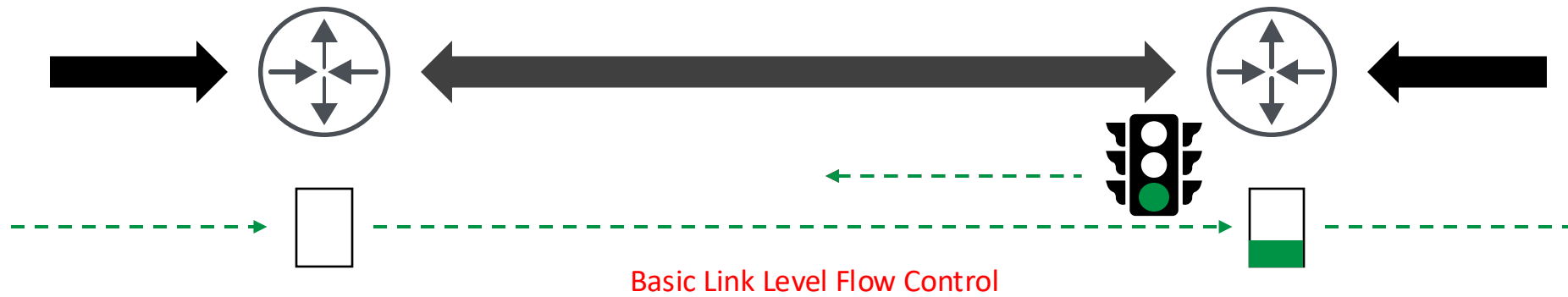
# Lossless Ethernet



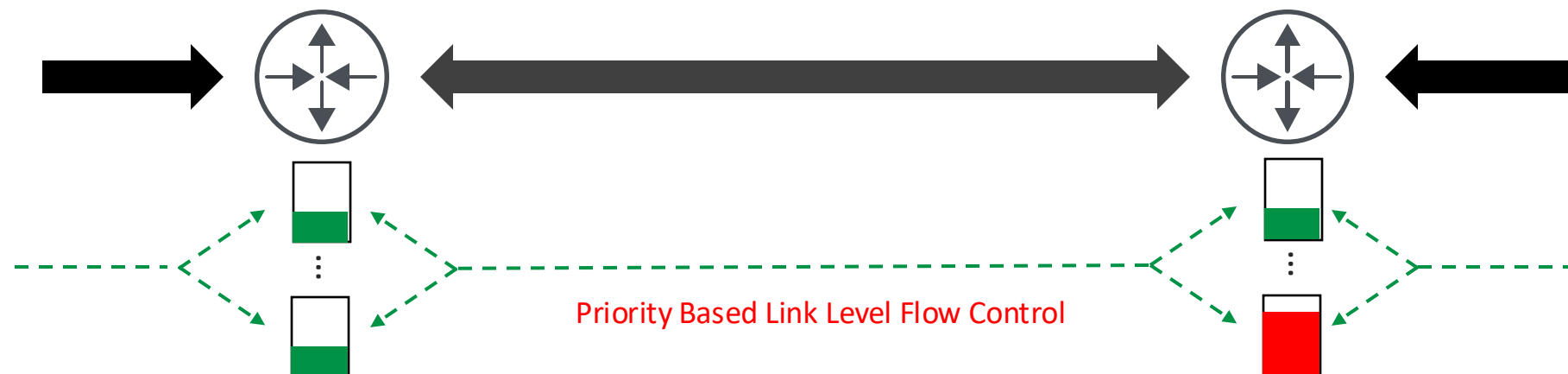
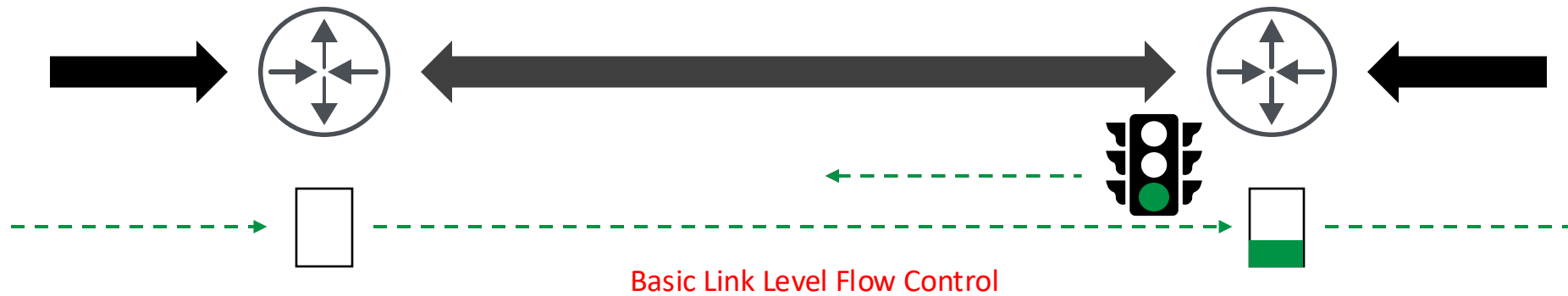
# Lossless Ethernet



# Lossless Ethernet



# Lossless Ethernet

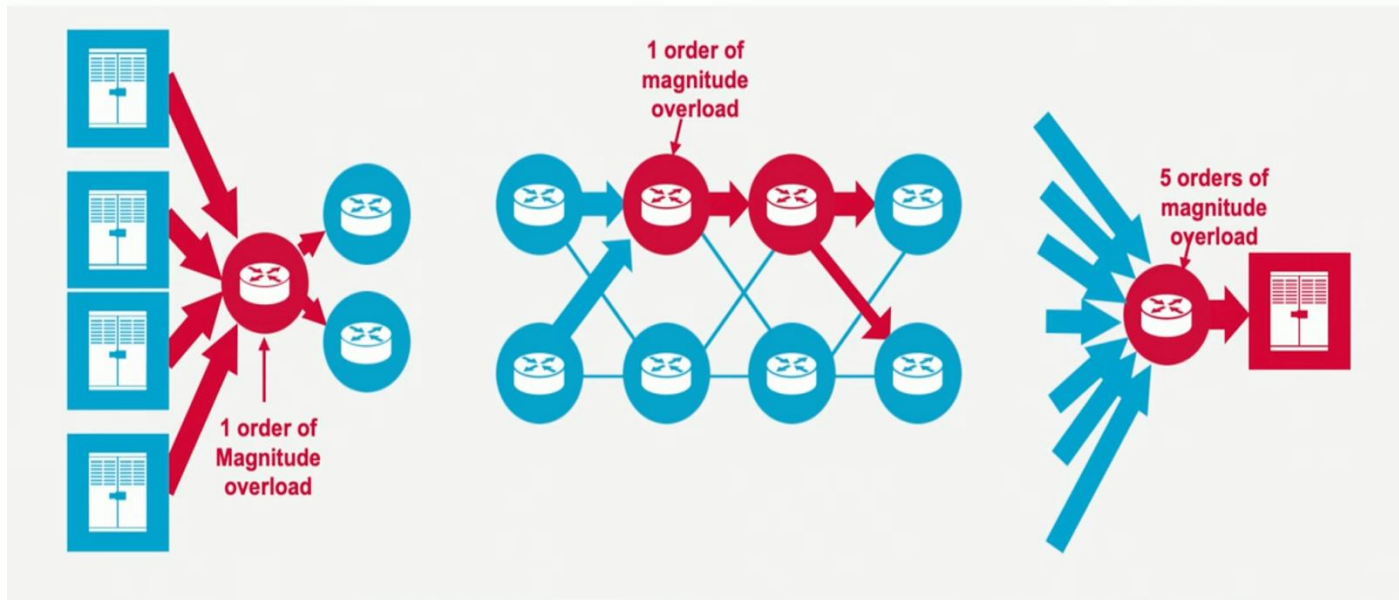


Issues:

- Head-of-Line Blocking
- Congestion Spreading
- Occasional Deadlocks

## Other Fabric Challenges

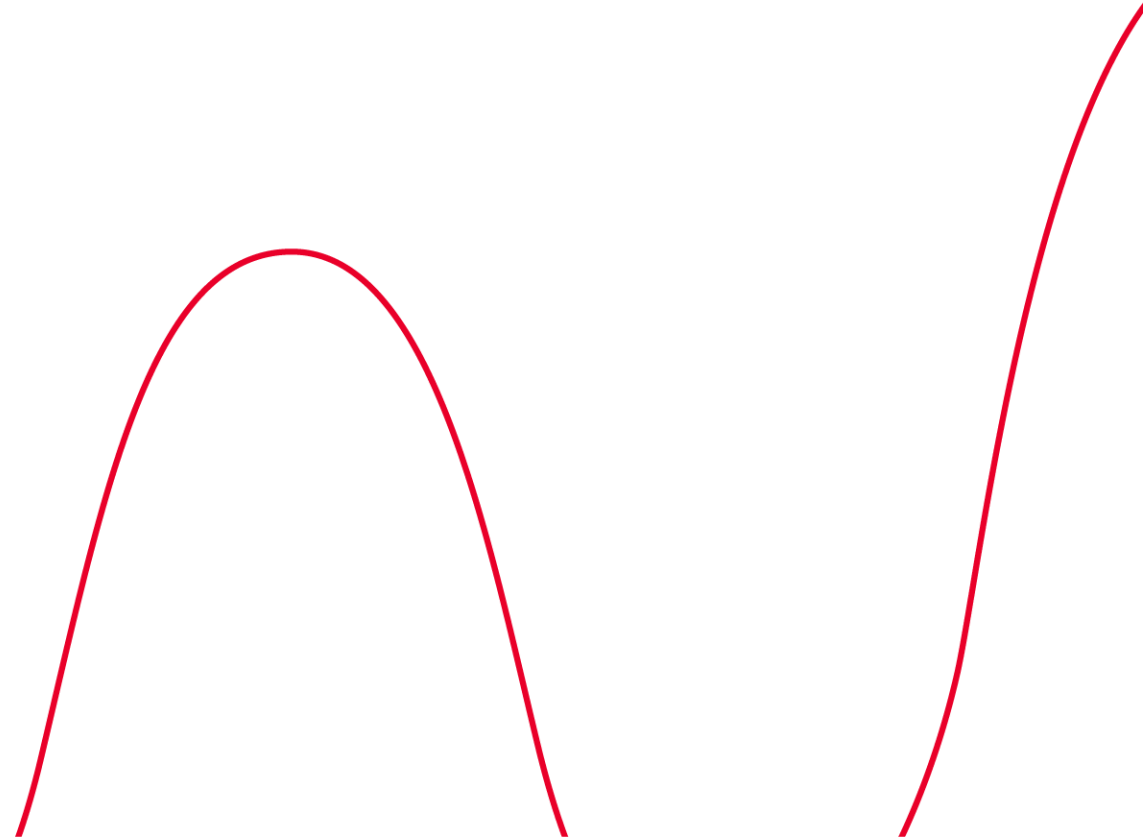
**To Minimize Latency, We Need Better Traffic Control to Keep Packets Out of Network Buffers**



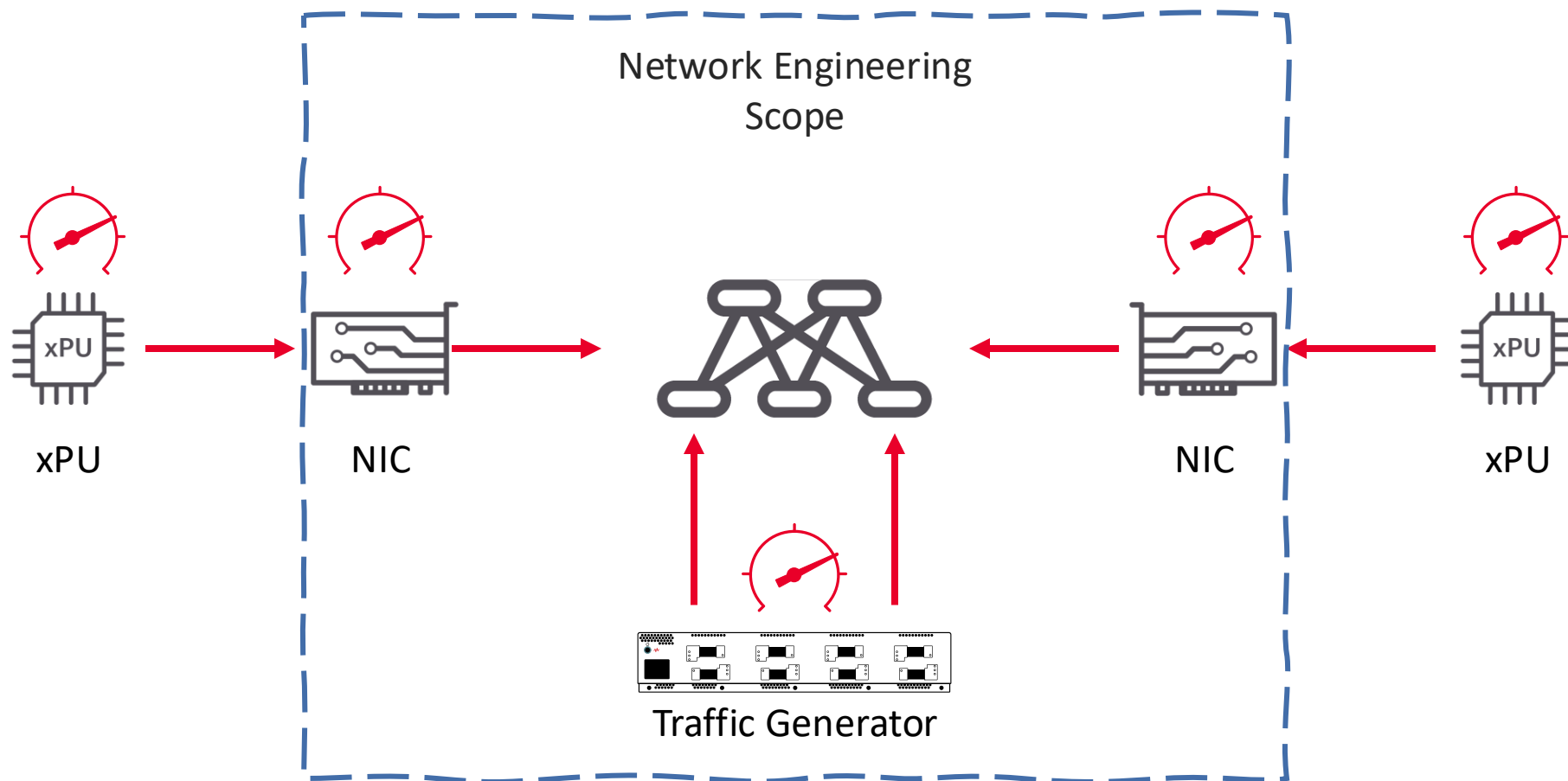
Source: Broadcom talk "Ethernet Fabric for High Performance Computing and AI/ML Workloads" @ OCP Summit 2022

- An Overwhelmed Receiver (In-cast)
- Uneven Link Utilization (Unequal Load Balancing)
- Re-ordering

# Network Engineer Perspective on Ethernet Validation for AI & HPC



# Layers of System Benchmarking





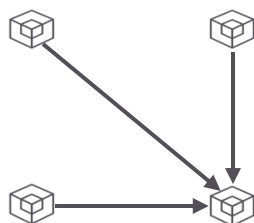
# Types and Tools of System Benchmarking

## Traditional Network Test Scenarios

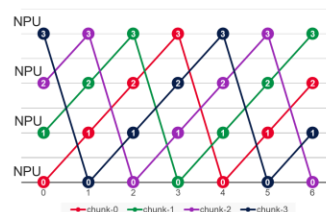
### P2P



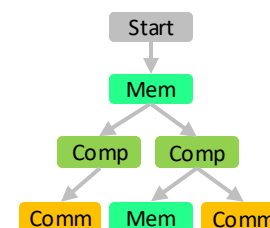
### Incast & Fanout



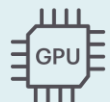
### Collective



### Model Parallelism



### Workload



gdr\_copy -----> ?

nccl-tests

PARAM

MLPerf



perftest -----> ?

OSU Micro-Benchmarks

?

?

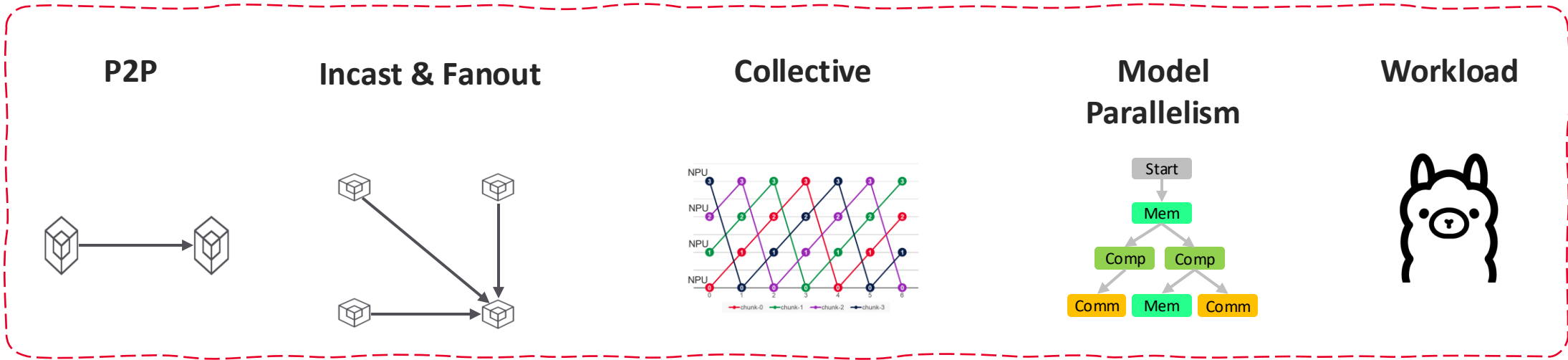


← Traffic Generator -----> ?

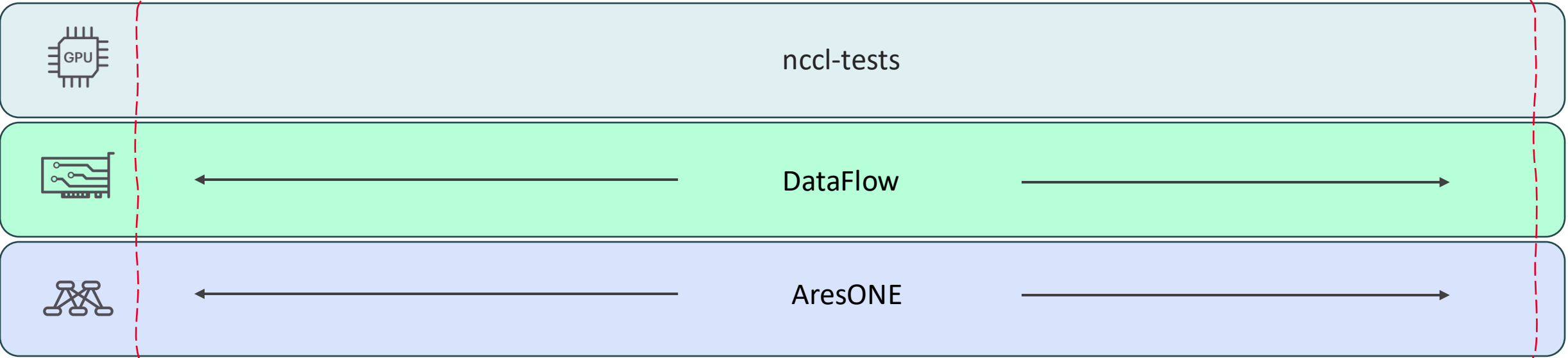
?

?

# Keysight's Approach



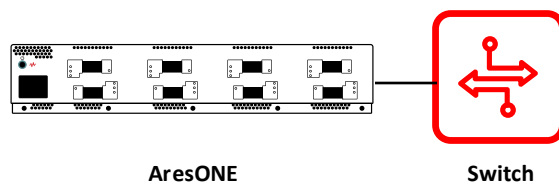
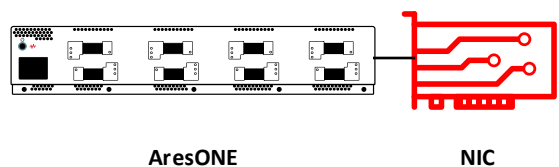
Apps



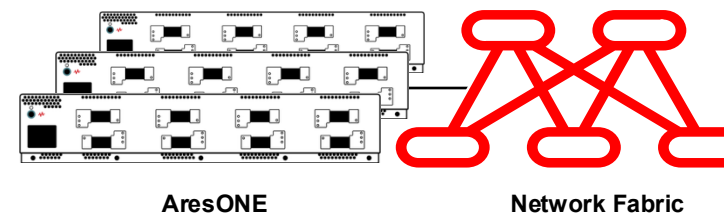
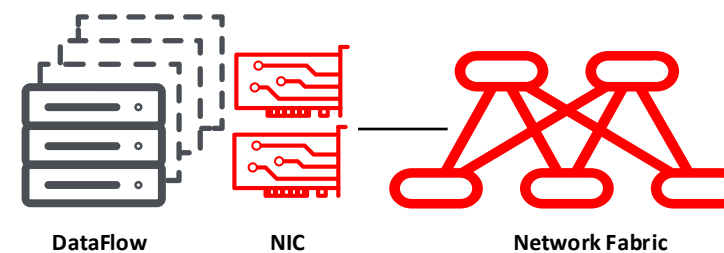
Test Platforms

# Testbed Topologies

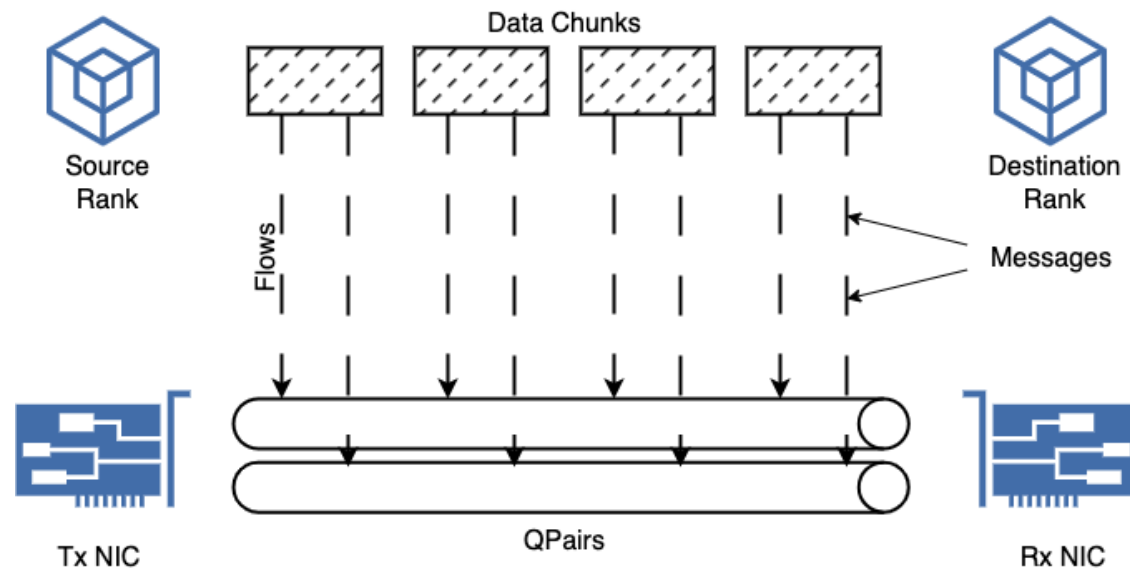
## Device level



## System level



# What can we measure in RoCEv2?



## Aggregate metrics per job

- Where: Ranks, Ports, **QPairs**
- What: Frames, Bytes, **Messages**, **Loss**

## Individual metrics per transfer

- Where: Chunks, Flows
- What: Timestamps

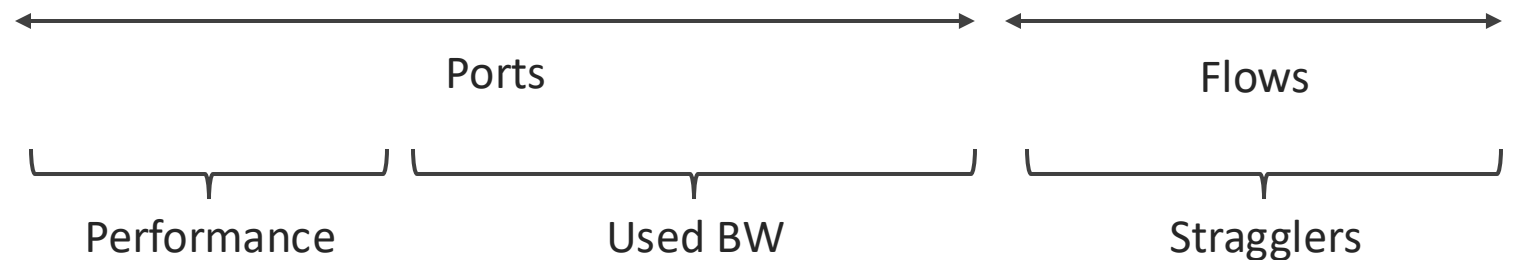
## Statistical distribution metrics

- Where: **Messages**, **Packets**
- What: **Completion**, **Latency**, **Reordering**

*(red) Traffic Generator metrics*

# Characterizing Efficiency

Benchmark		Busbw	Ideal %	Fabric Utilization (%)	P95 FCT Skew (%)
Type of Workload	Algorithm Specific		$\frac{Time_{min}}{Time}$	$\frac{\sum Bytes_{Ingress}}{Ports * Speed * Time}$	$\frac{FCT_{p95} - FCT_{p50}}{FCT_{p50}}$
AlltoAll-32 1GB		42.76	86.83	86.91	3.91
AlltoAll-32 1GB		46.67	94.76	94.81	0.66
AlltoAll-8 2GB		47.10	95.64	96.15	2.40

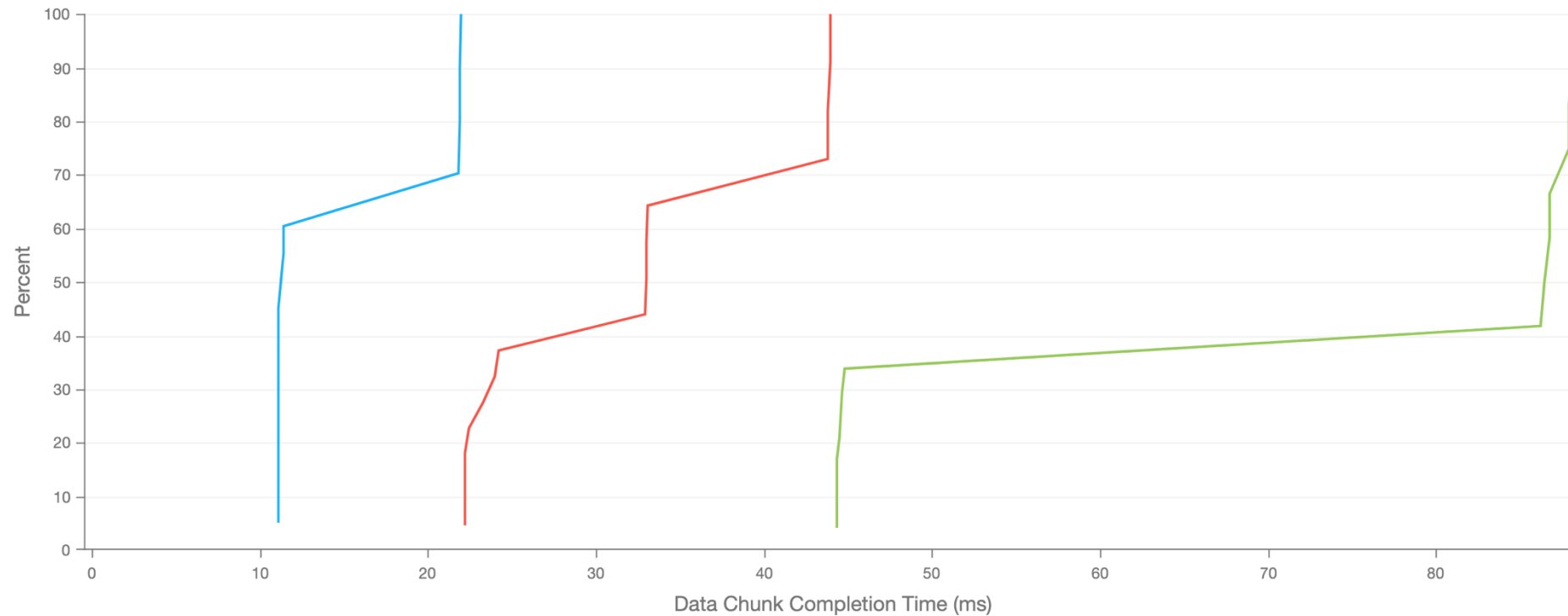


## Signals of Troubles

Benchmark	PFC Rx	ECN CE	ReTx Frames	Reordering
Type of Workload	$\sum Port(PFC Rx)$	$\sum QP(ECN CE)$	$\sum QP(ReTx)$	WIP
AlltoAll-32 1GB	546	311,645	0	—
AlltoAll-32 1GB	96	16,593	0	—
AlltoAll-8 2GB	0	0	10,257	—



# Visualizing Chunk & Flow Data

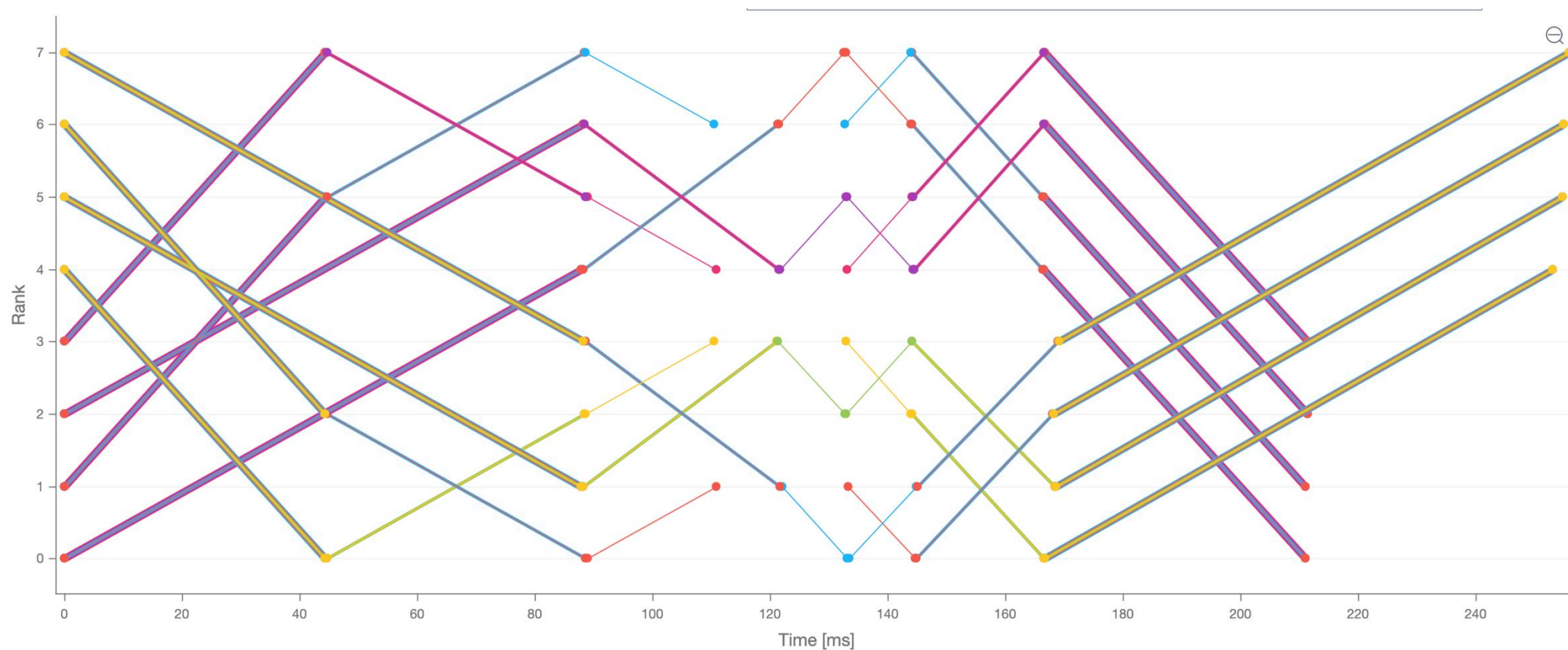


## Cumulative Distribution Function (CDF)

Severity and uniformity of issues



# Visualizing Chunk & Flow Data

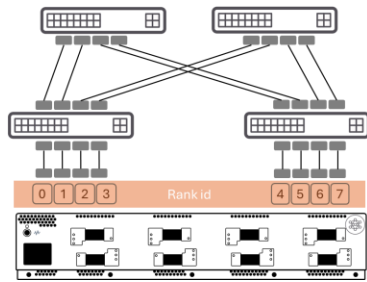


## Data Movement Timeline

Localize place and time of issues

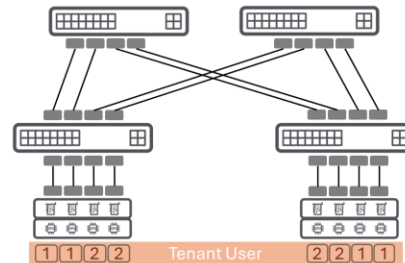
# KAI Data Center Fabric Test Methodology

Available for download



## Job Completion Time

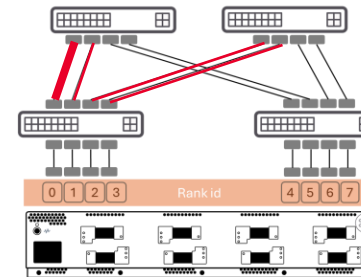
- Topologies
- Algorithms
- Data sizes
- RDMA message sizes



## Performance Isolation

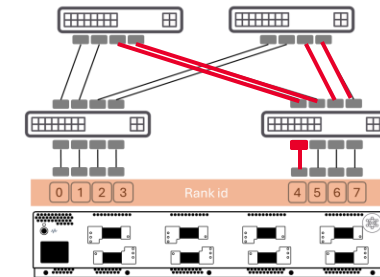
- Noisy neighbors
- Parallel collectives

Upcoming 2<sup>nd</sup> edition with  
new benchmark type:  
**Mix of Collectives**



## Load Balancing

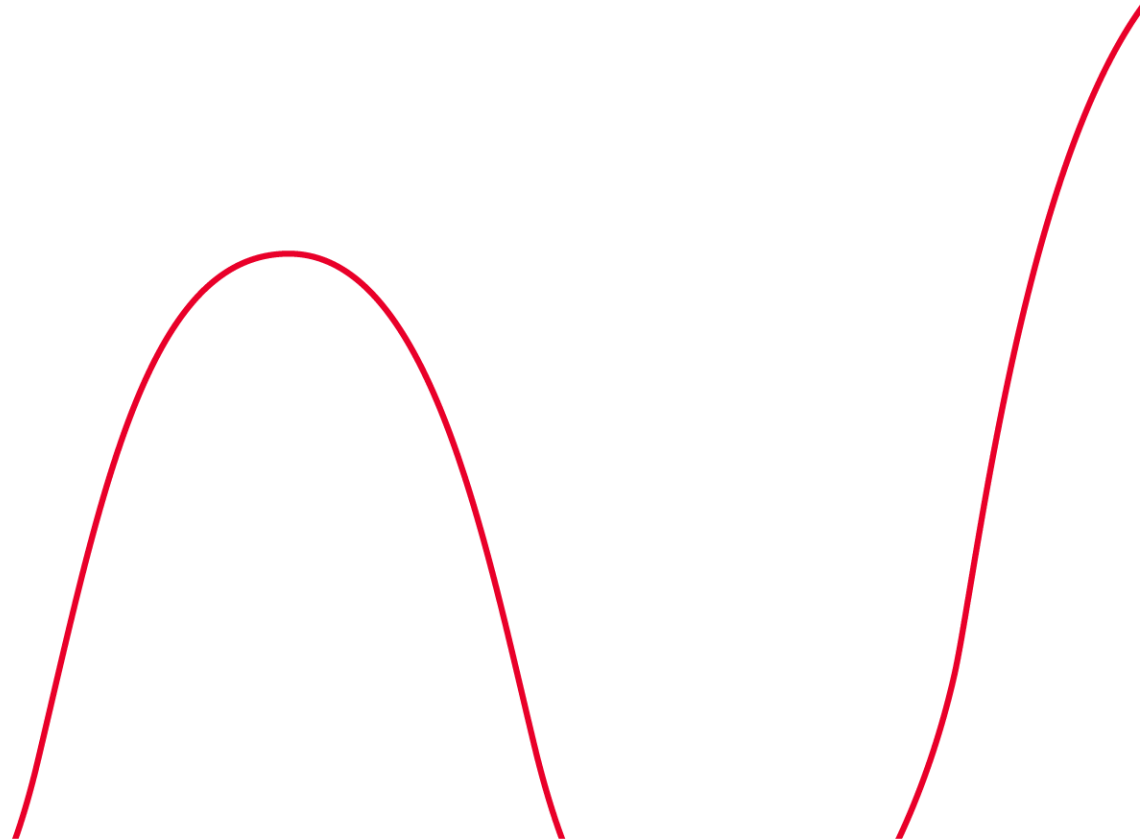
- ECMP hashing
- Traffic Engineering
- Q-Pair awareness
- Parallel Q-Pairs
- Dynamic Load Balancing



## Congestion Control

- PFC
- ECN/DCQCN
- Fine-tuning

# Closing Remarks



# Ethernet Future

- Congestion Control Algorithms
- Multi-Path
- Retransmission at Link Level
- Flow Control at Link Level
- Source Routing
- ...

*Holy Grail: One Interconnect To Rule Them All!!!*



# Learn More About the Keysight AI Data Center Builder



PDL-KAI-DC-Builder@keysight.com



Product Page



Methodology



Solution Brief



Documentation

## Meet Keysight Team

Dan Mihailescu, Winston Liu

### **Chakra – Standardized AI Workload Traces for System Co-Design**

**MUG '25**

Tue, August 19, 4:00 PM

Visit us in booth C1 for in-person demos

### **2025 OCP Global Summit**

---

**San Jose, California**

October 13-16, 2025

