# CacheCloud
## Towards Speed of Light Datacenter Communication
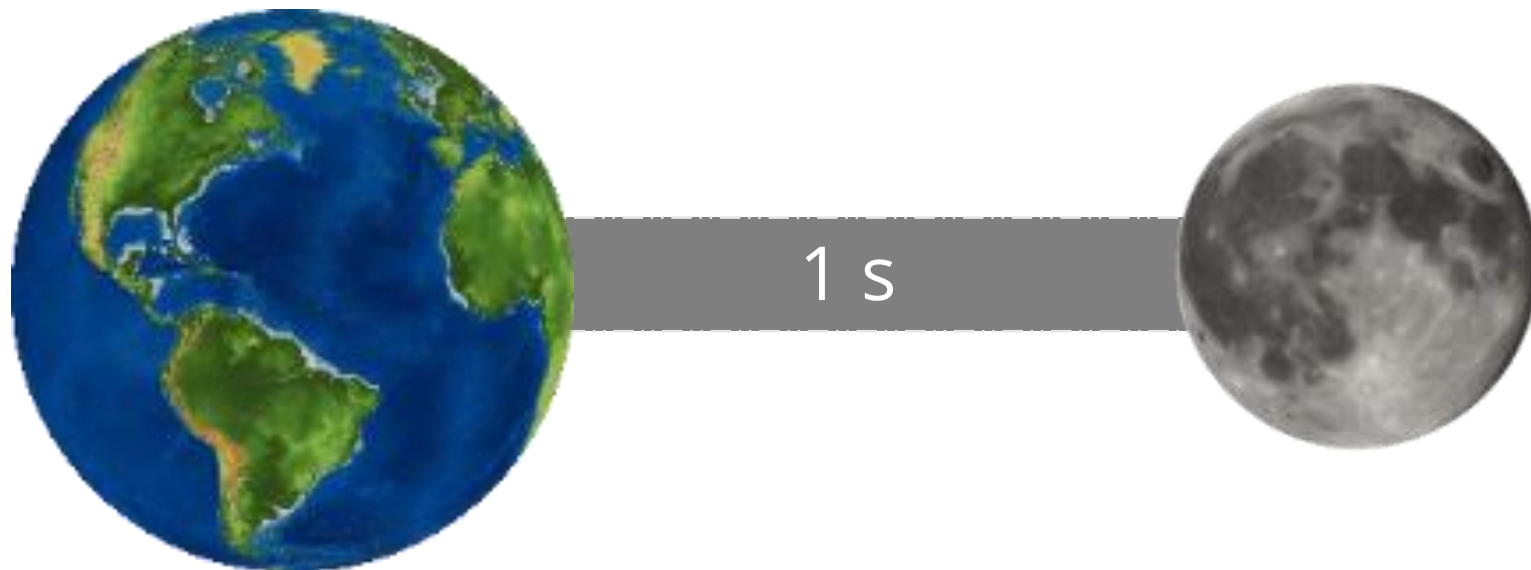
**Shelby Thomas**

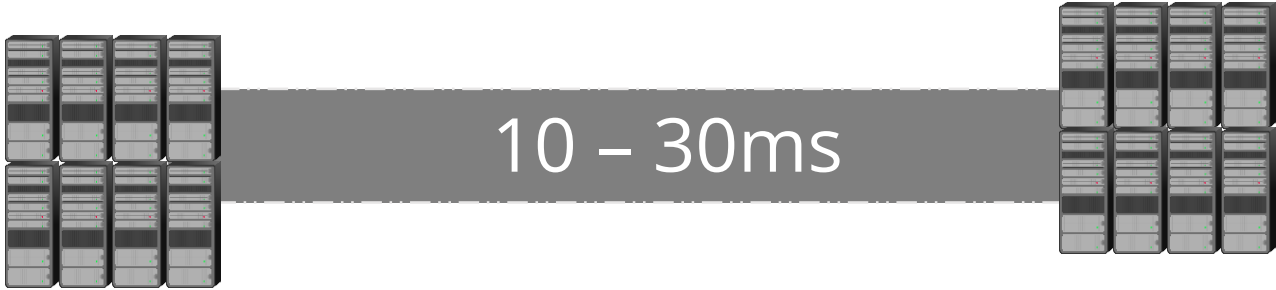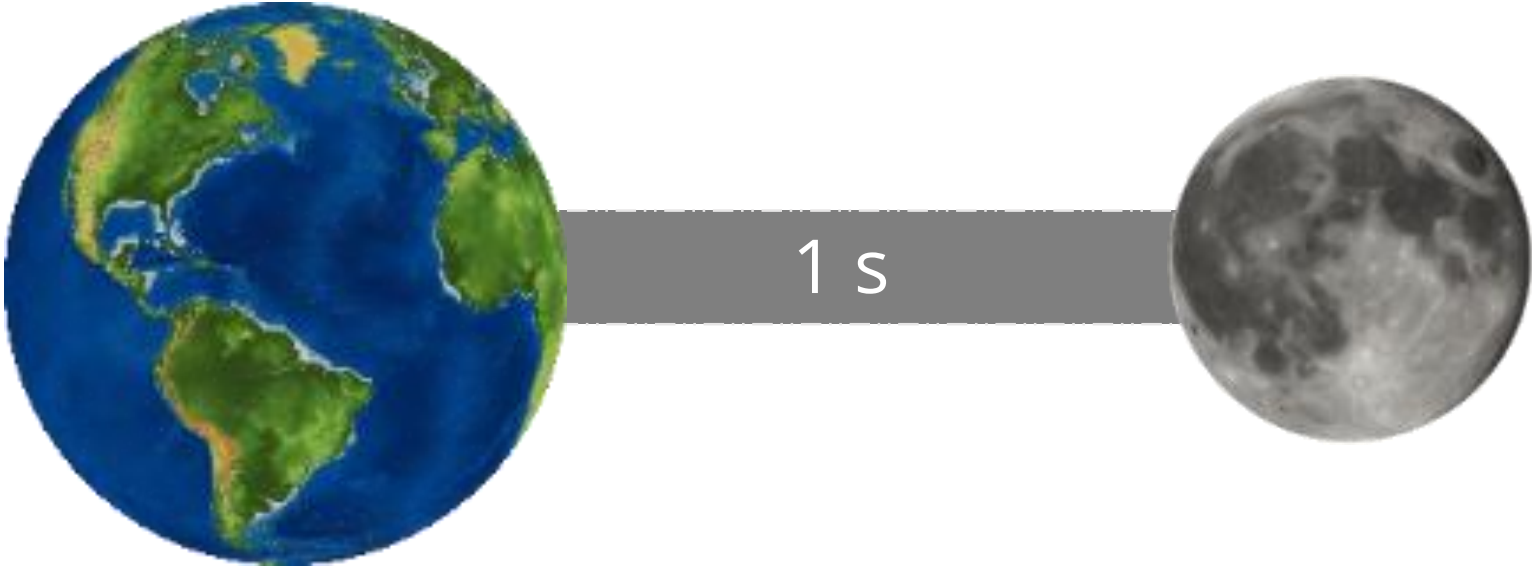Geoff Voelker

George Porter

University of California, San Diego

# The Speed of Light Baseline

# The Speed of Light Baseline



1 s

10 – 30ms

[Singla Et. Al.]

# Speed of Light in the **Wide Area**

Bufferbloat and congestion

DNS resolution

Distance

10x – 1000x slower than speed of light propagation in WAN

[Singla Et. Al.]

# Speed of Light in the Wide Area

Bufferbloat and congestion

## What about the Local Area?

Distance

10x – 1000x slower than speed of light propagation in WAN

[Singla Et. Al.]

# Speed of Light in the **Local** Area

~~Bufferbloat and congestion~~

~~DNS resolution~~

~~Distance~~

**Expectation:** Datacenter applications are closer to the speed of light than wide area networks.

# Speed of Light in the **Local** Area

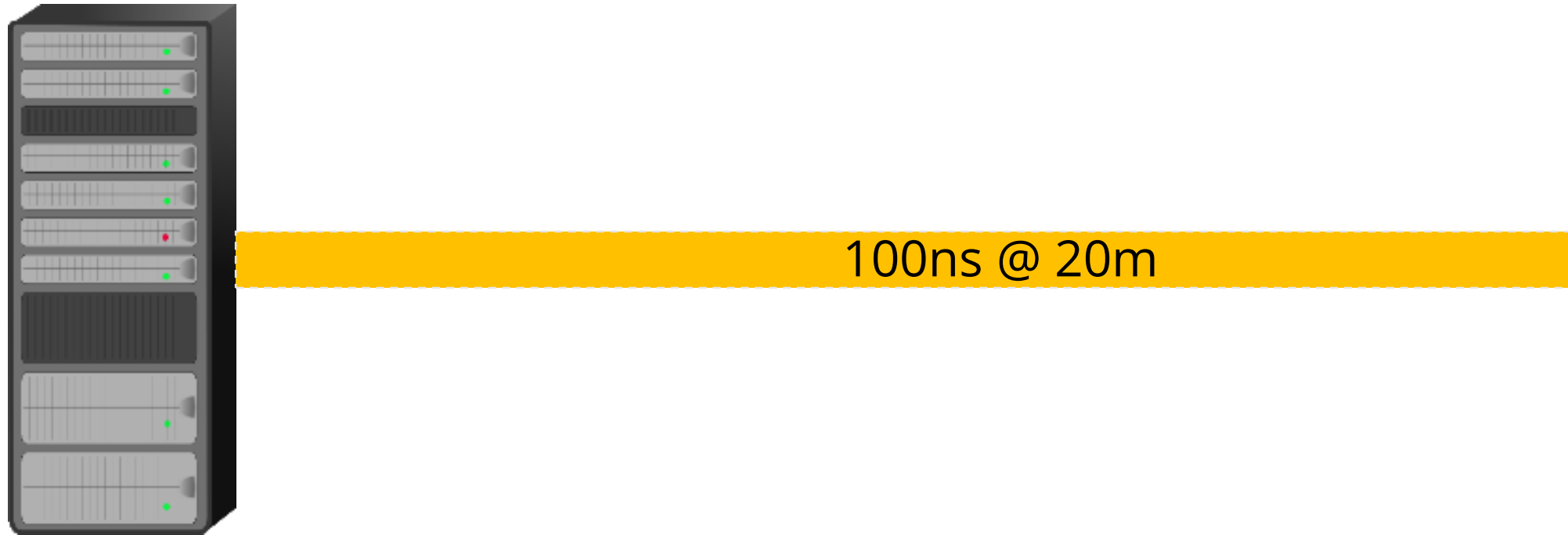~~Bufferbloat and congestion~~

~~DNS resolution~~

~~Distance~~

**Reality:** Datacenter applications are **10x – 1000x slower** than speed of light propagation in the LAN
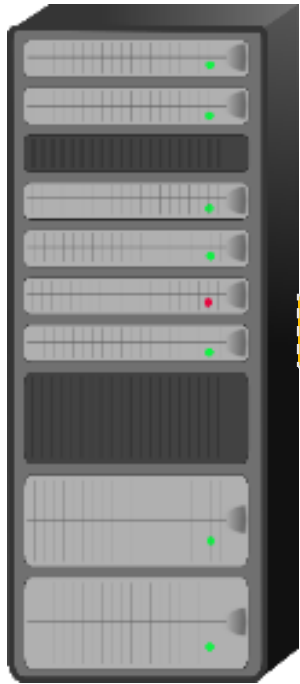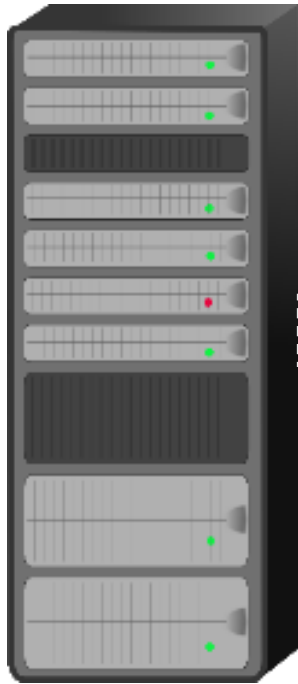
# Different Distances Same Performance Gap

# Datacenter Propagation Delay

100ns @ 20m

# Datacenter Propagation Delay

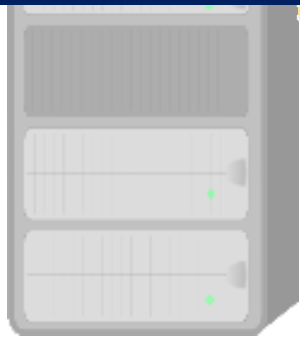Datacenter applications do not run at nanosecond scale

100ns @ 20m

# Datacenter Propagation Delay

Datacenter applications do not run at nanosecond scale

100ns @ 20m

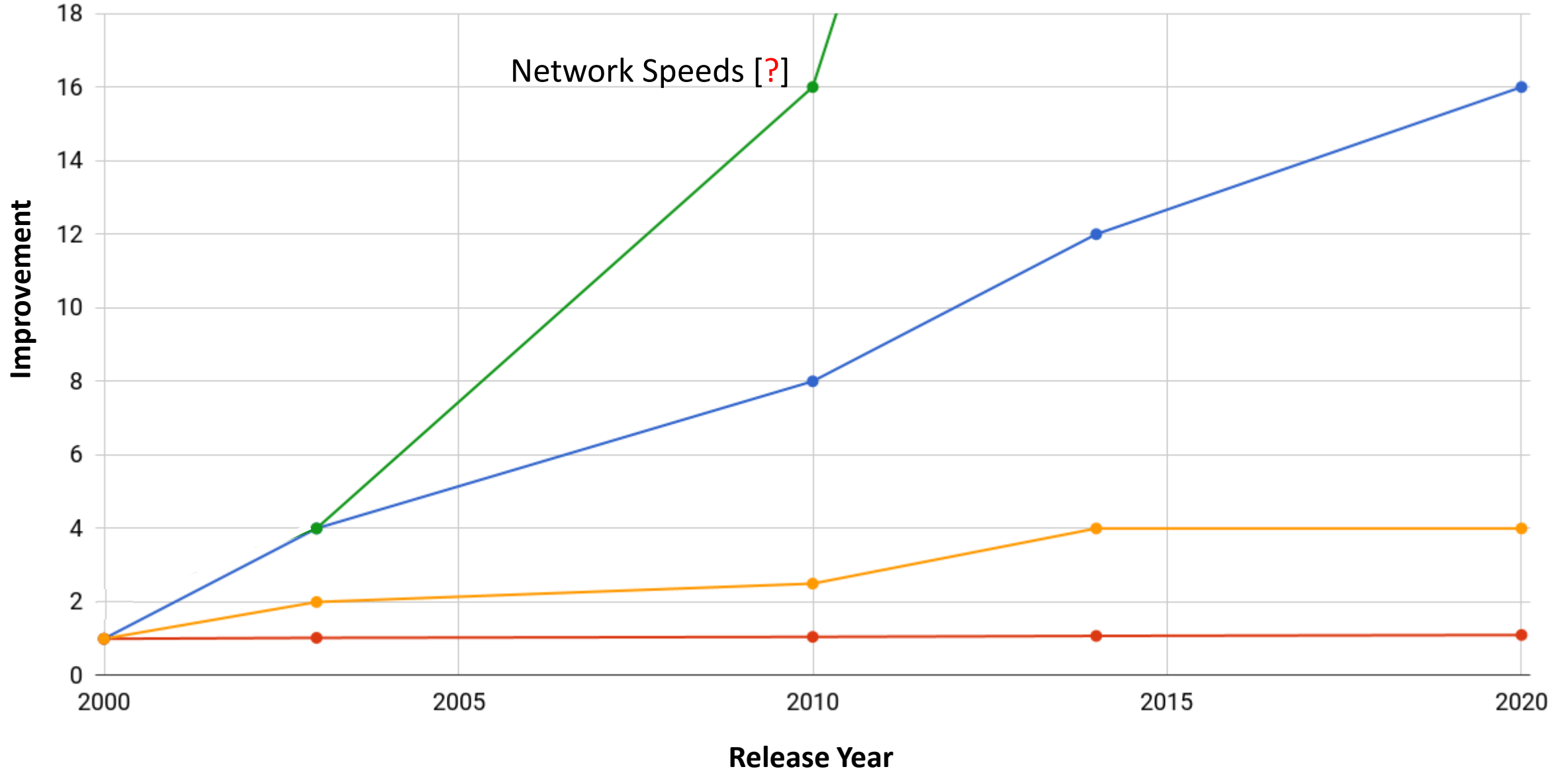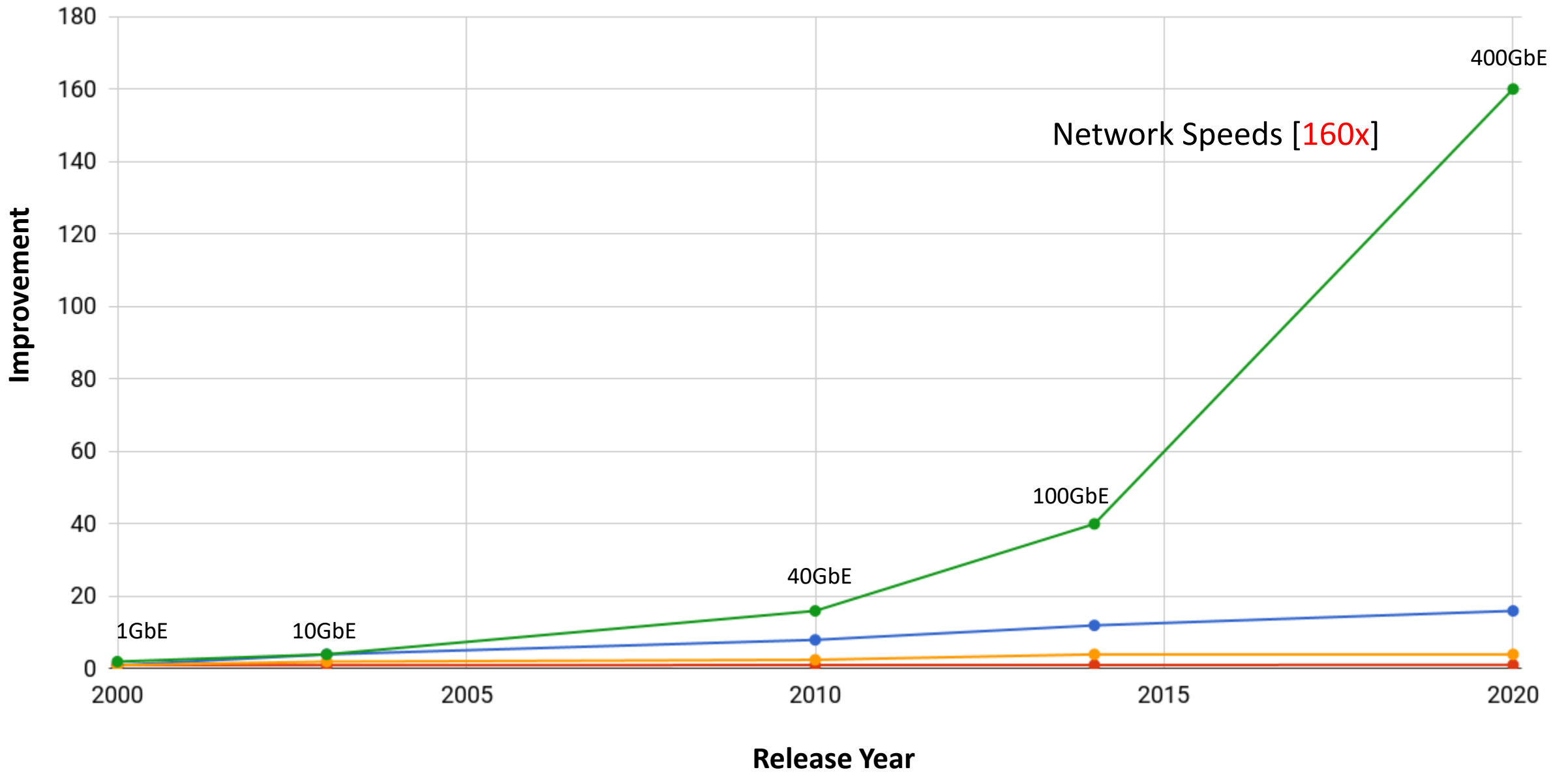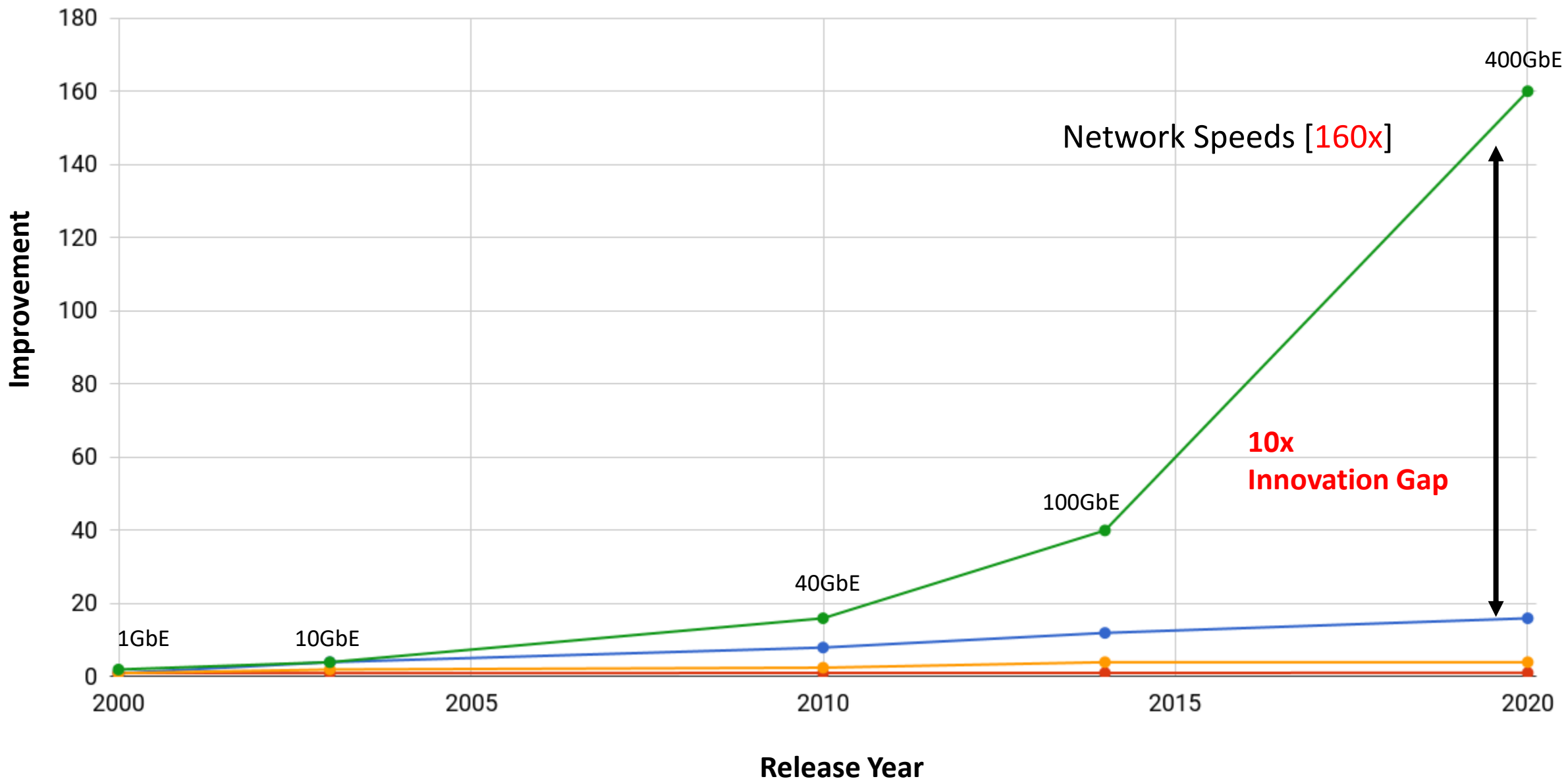| Application | Latency | Gap |
|---|---|---|
| Key Value Stores | 4.3us – 100us | **43x** |
| Industry End-to-End RPC | 75us | **750x** |
| Consensus | 100us – 300us | **1000x** |

# Datacenter Propagation Delay

Datacenter applications do not run at nanosecond scale

**How has our infrastructure evolved?**

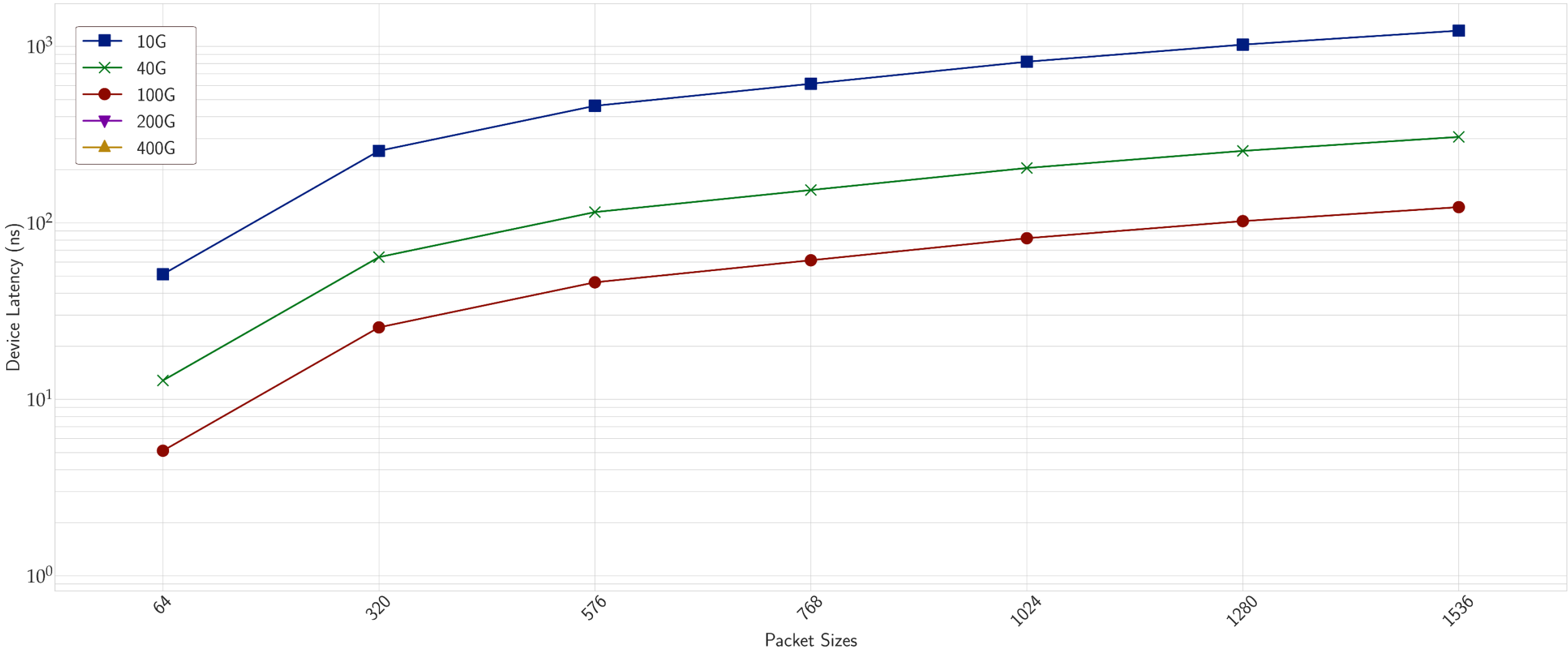| Application | Latency | Gap |
|---|---|---|
| Key Value Stores | 4.3us – 100us | **43x** |
| Industry End-to-End RPC | 75us | **750x** |
| Consensus | 100us – 300us | **1000x** |

DRAM Latency [1.1x]
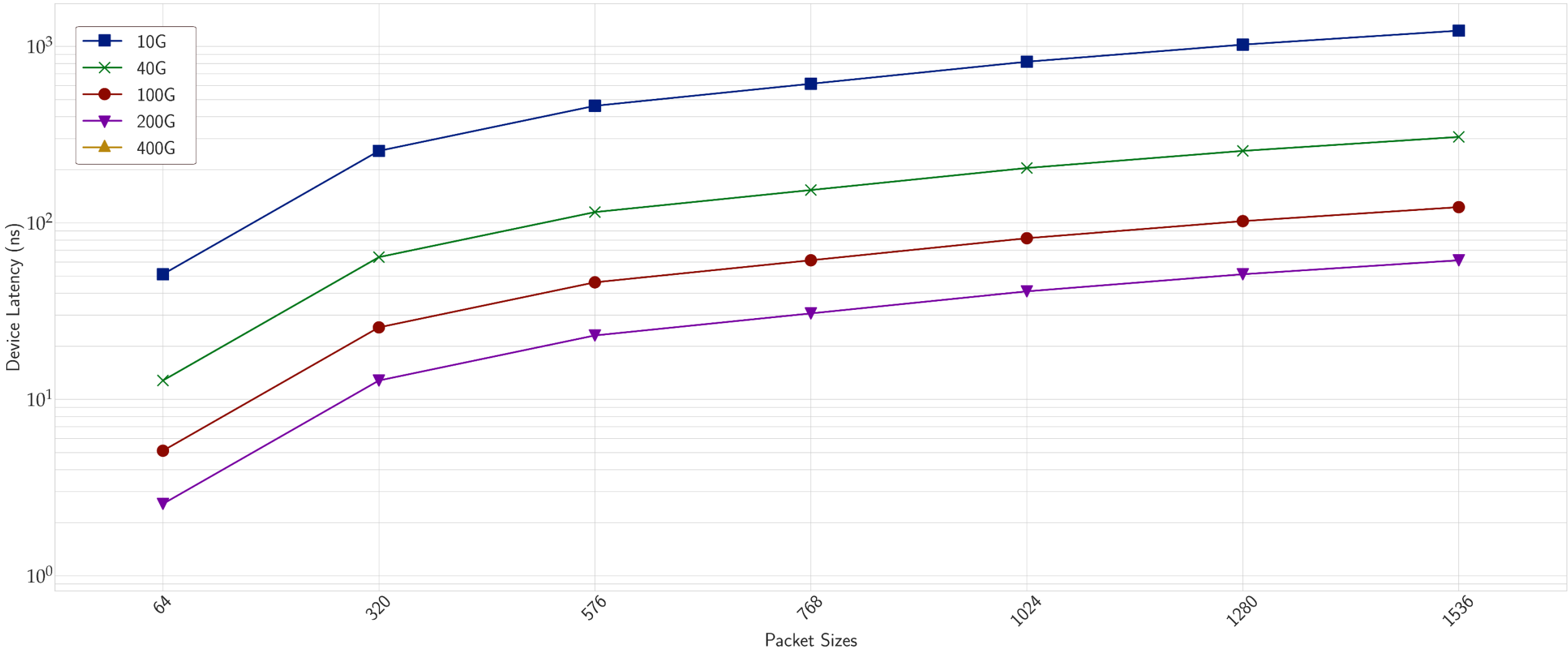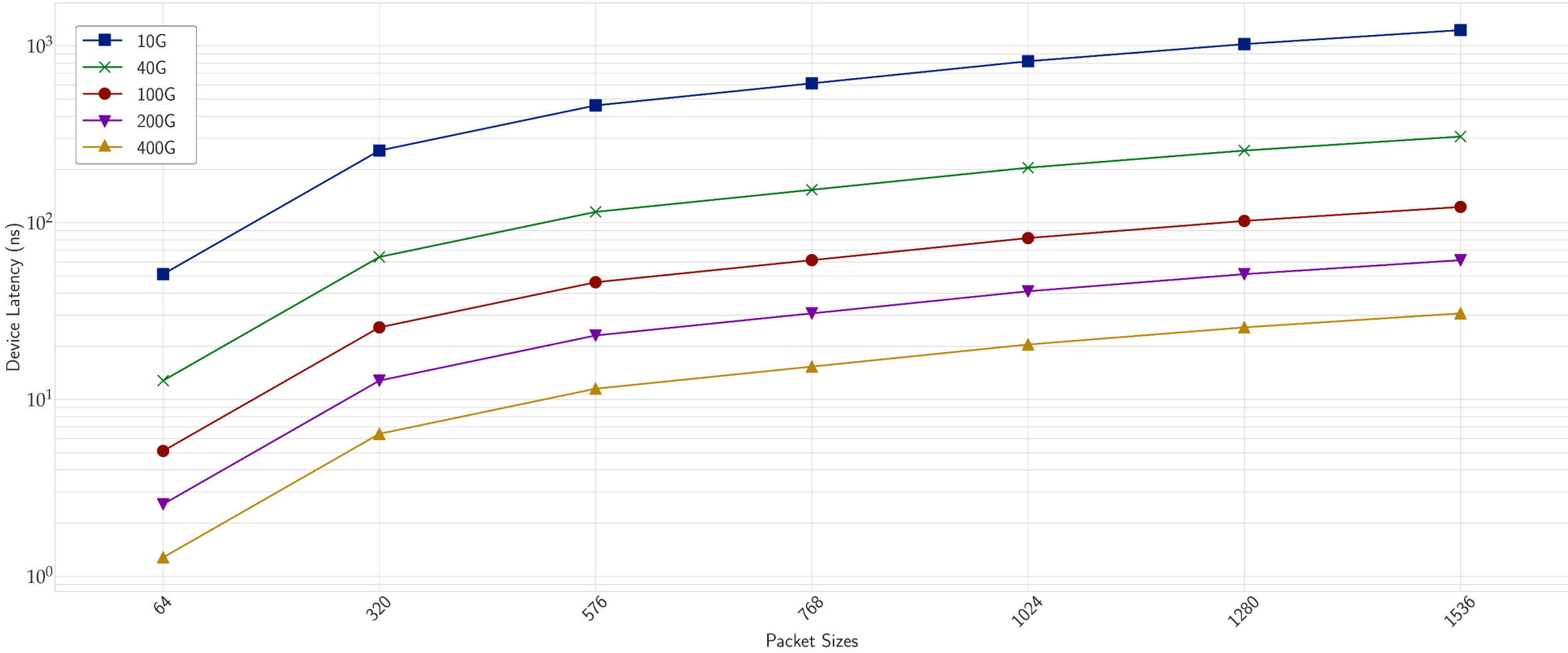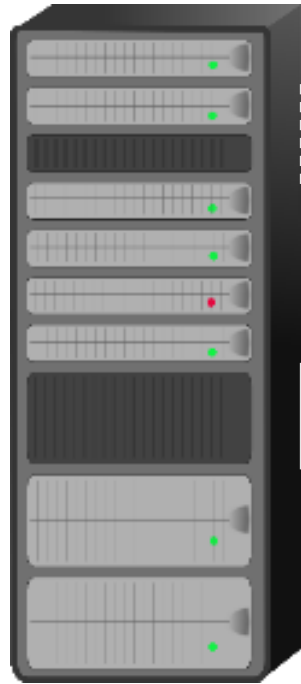
[Mutlu Et. Al.]

Network Speeds [?]

# Interpacket Gap Decreased Drastically

Speed of Light Propagation (100ns)

1GbE [10000ns] (MTU)

# Interpacket Gap Decreased Drastically

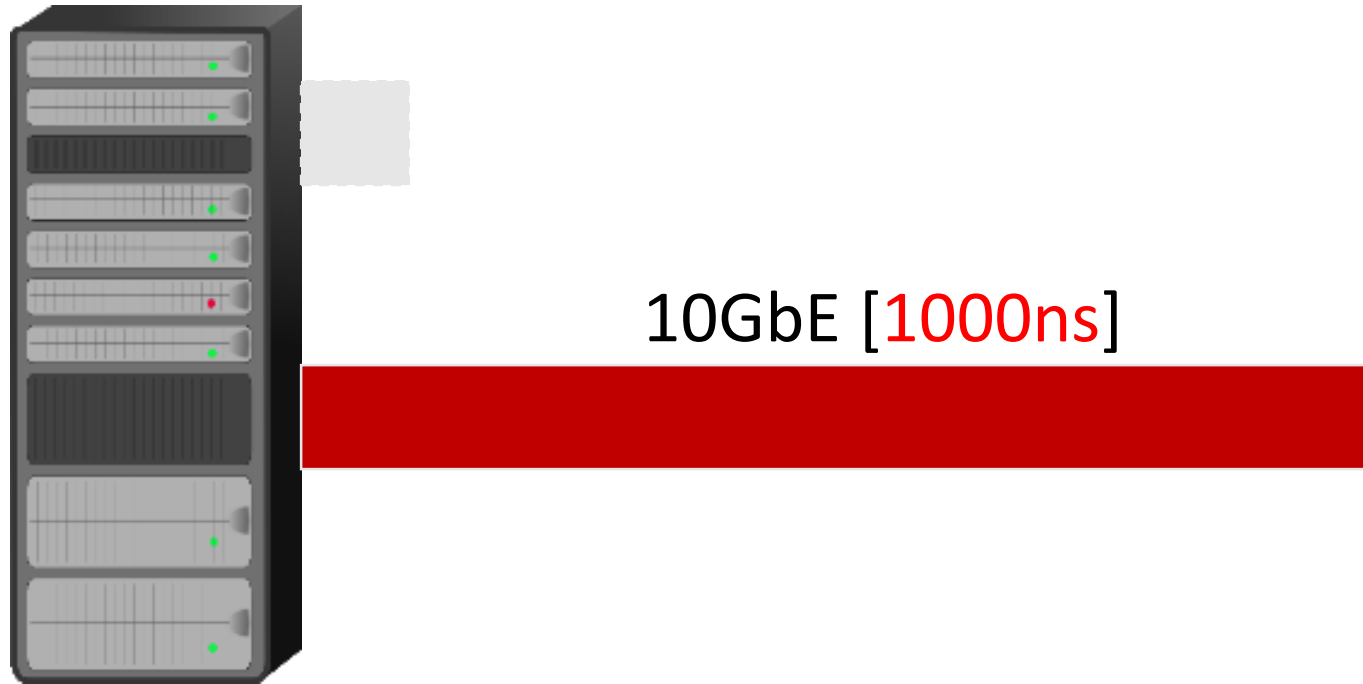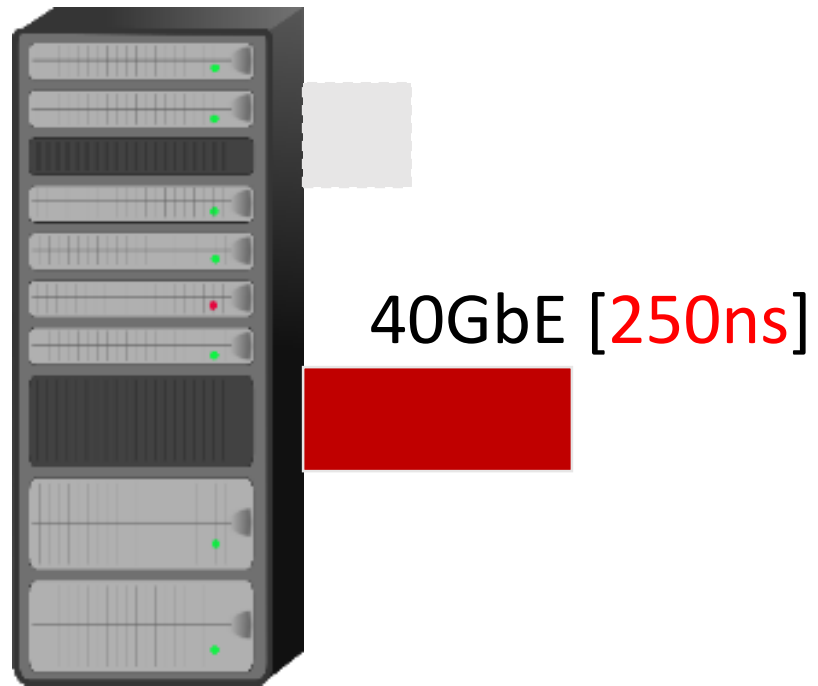# Interpacket Gap Decreased Drastically

10GbE [1000ns]

# Interpacket Gap Decreased Drastically

# Interpacket Gap Decreased Drastically
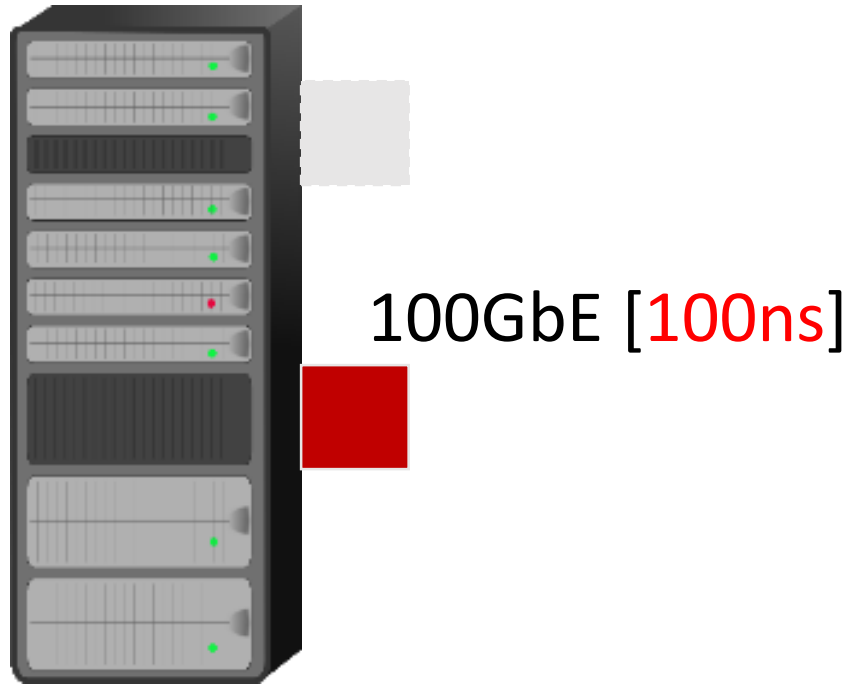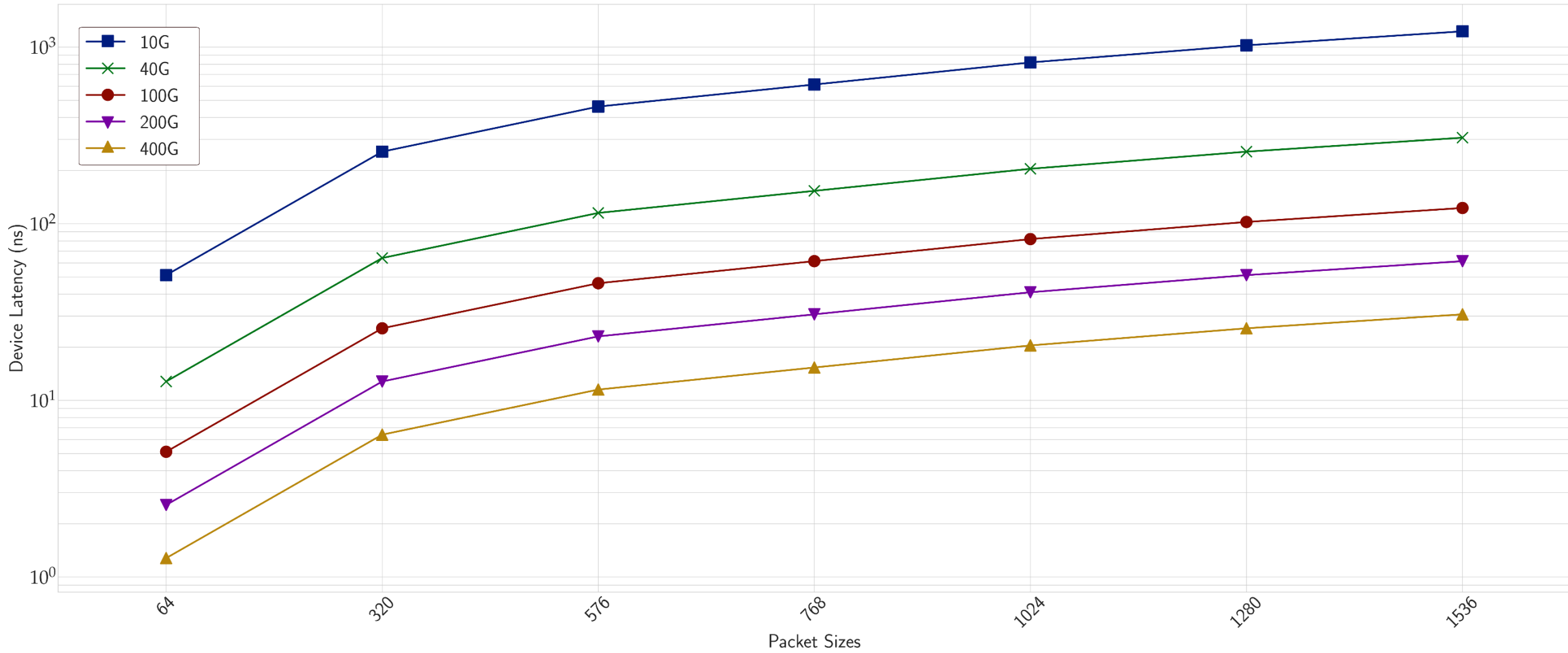
40GbE [250ns]

# Interpacket Gap Decreased Drastically
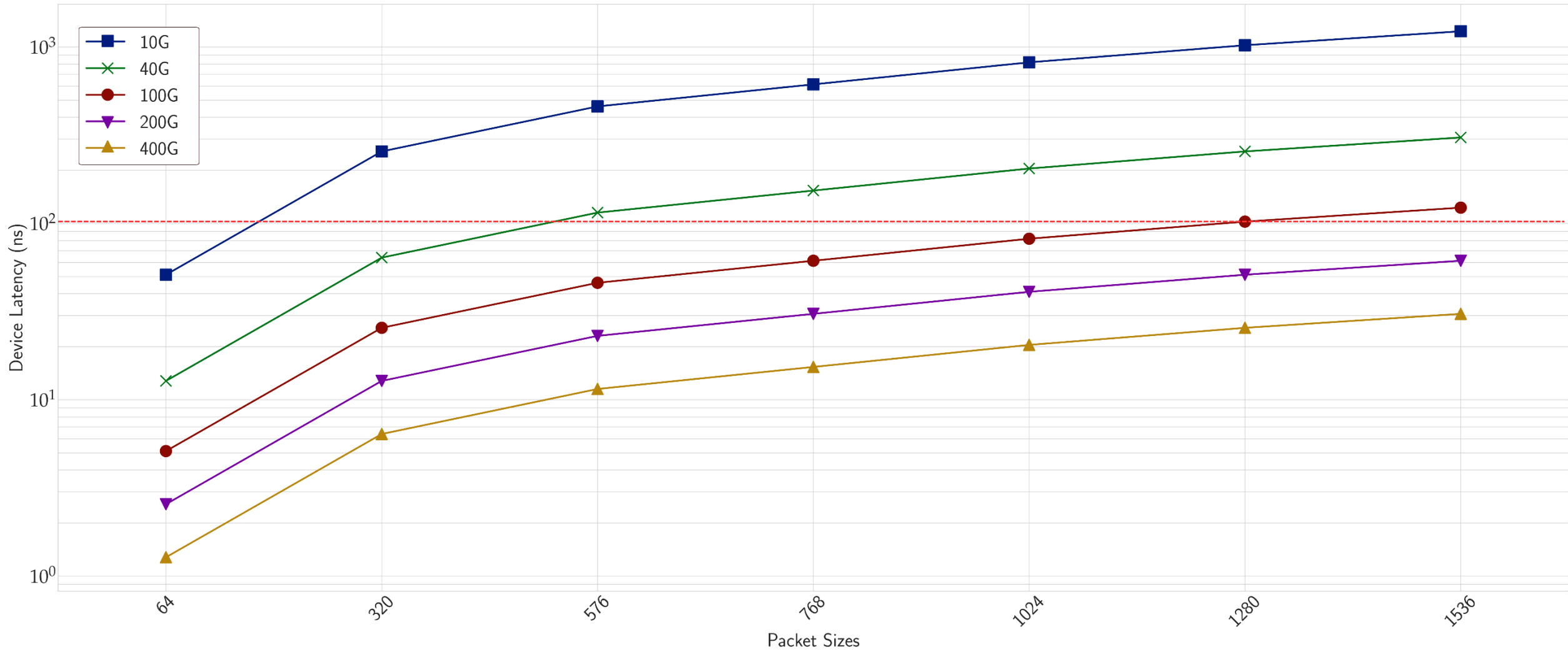
# Interpacket Gap Decreased Drastically



100GbE [100ns]
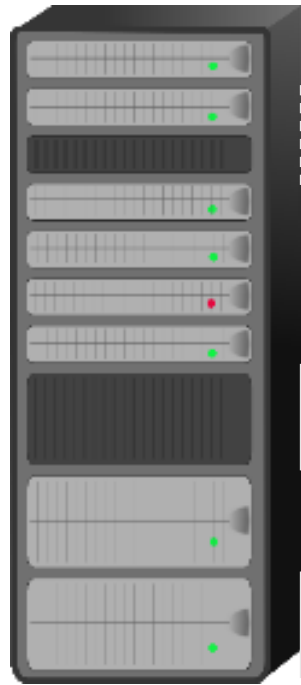
DRAM Latency ~ 100ns
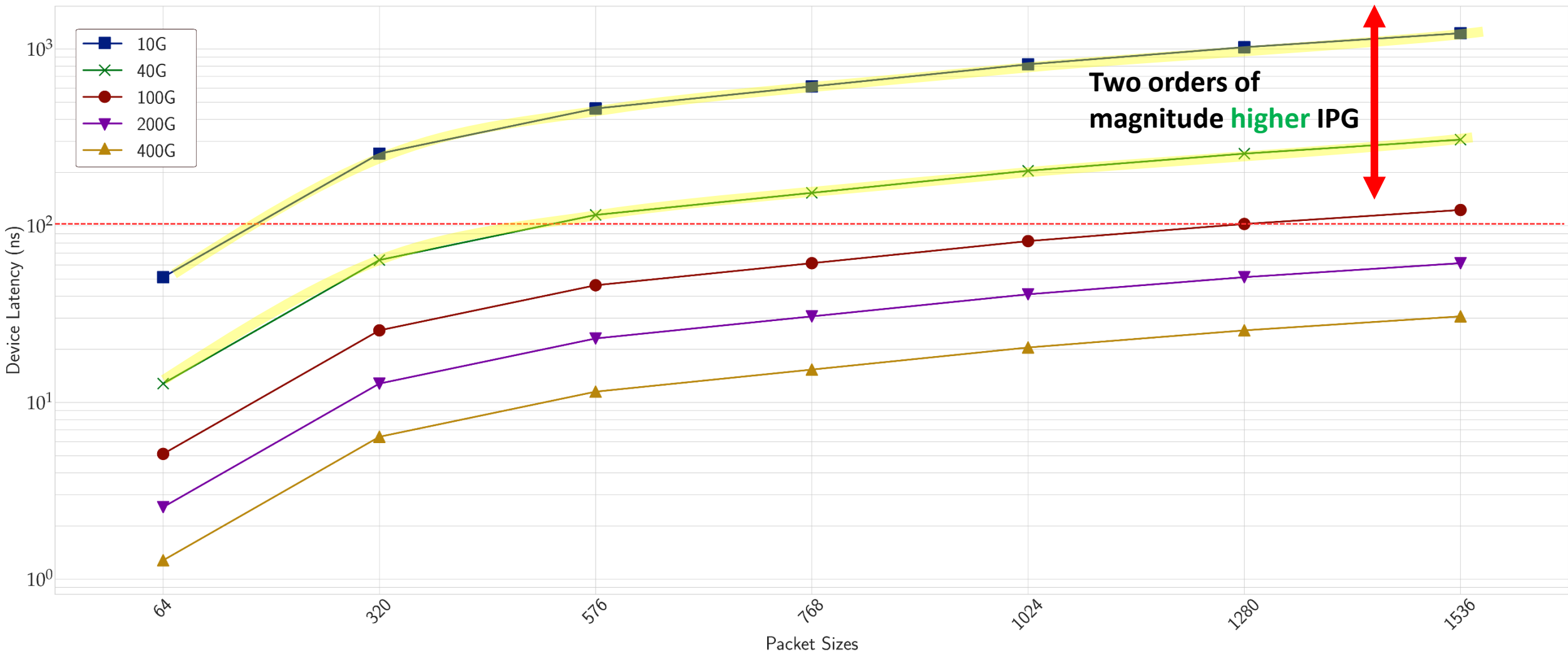
Interpacket Gap < DRAM Latency

# Latency Convergence



Speed of Light Propagation [100ns]

40GbE/100GbE Interpacket gap [100ns]

DRAM Latency [100ns+]

DRAM accesses need to be minimized or eliminated

DRAM is the new Disk

How do we overcome DRAM limitation?

How do we overcome DRAM limitation?  → **Pivot to SRAM**

How do we overcome DRAM limitation? → **Pivot to SRAM**

# Faster than Light SRAM

Speed of Light Propagation [100ns]
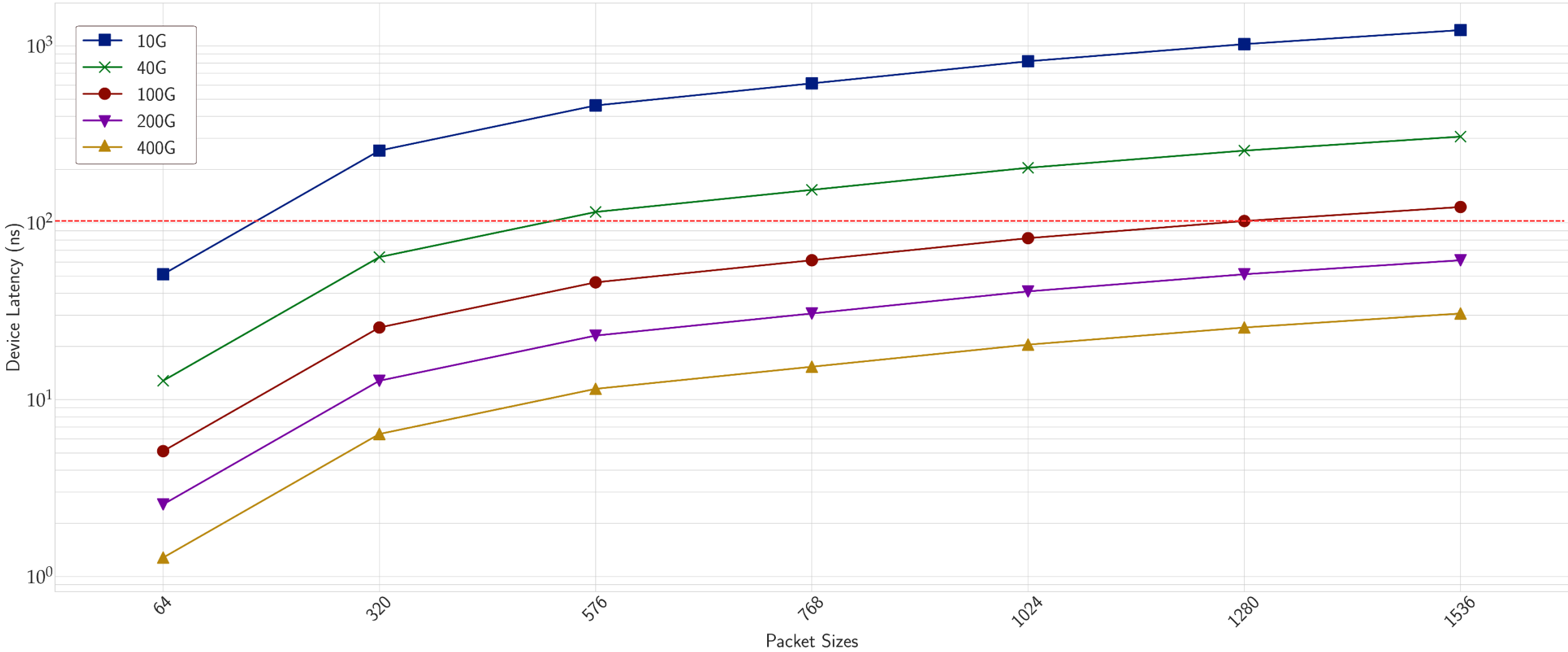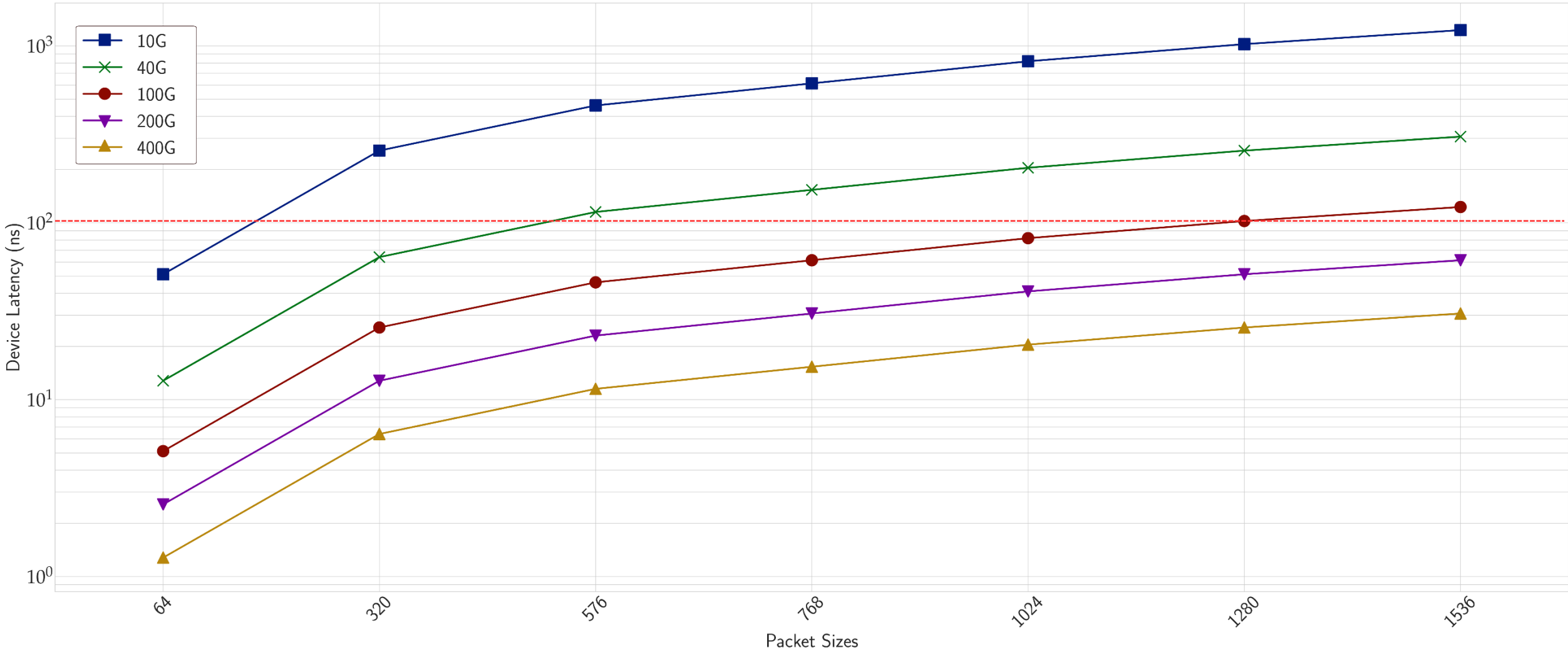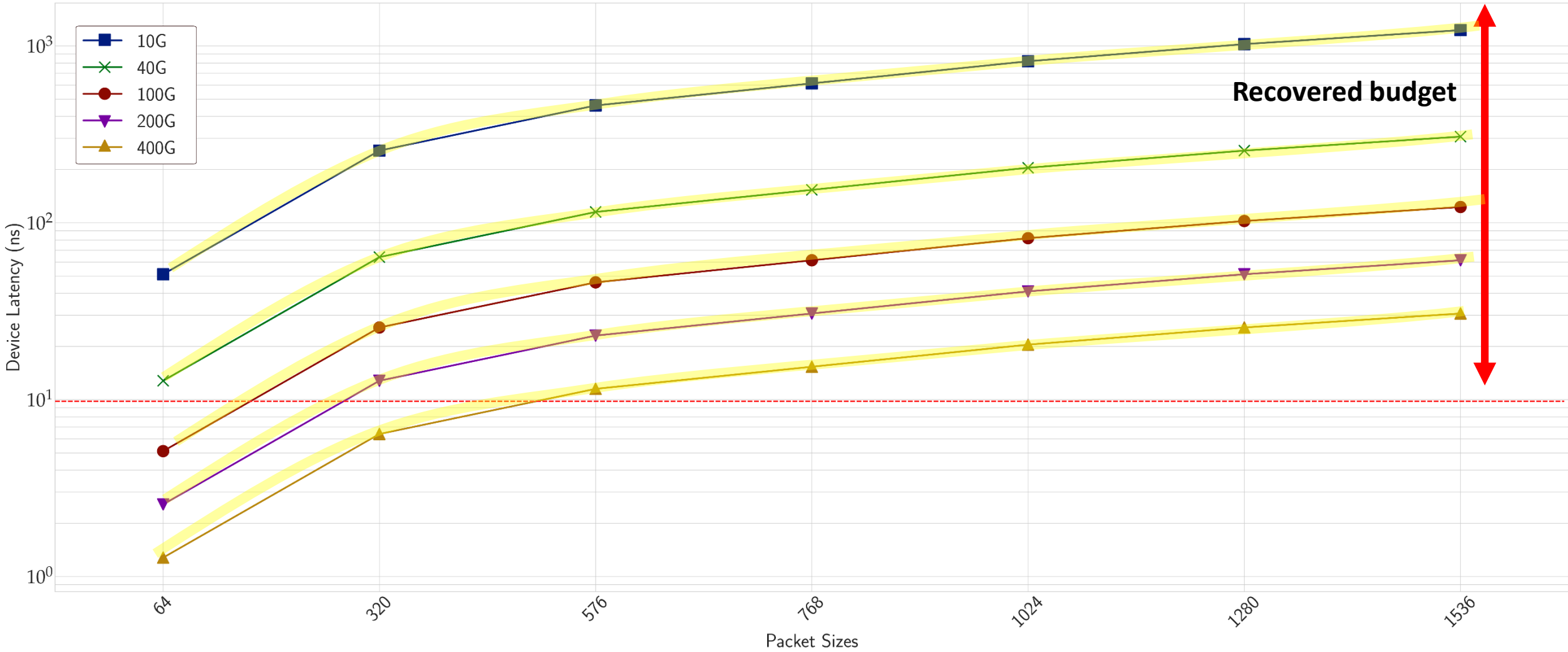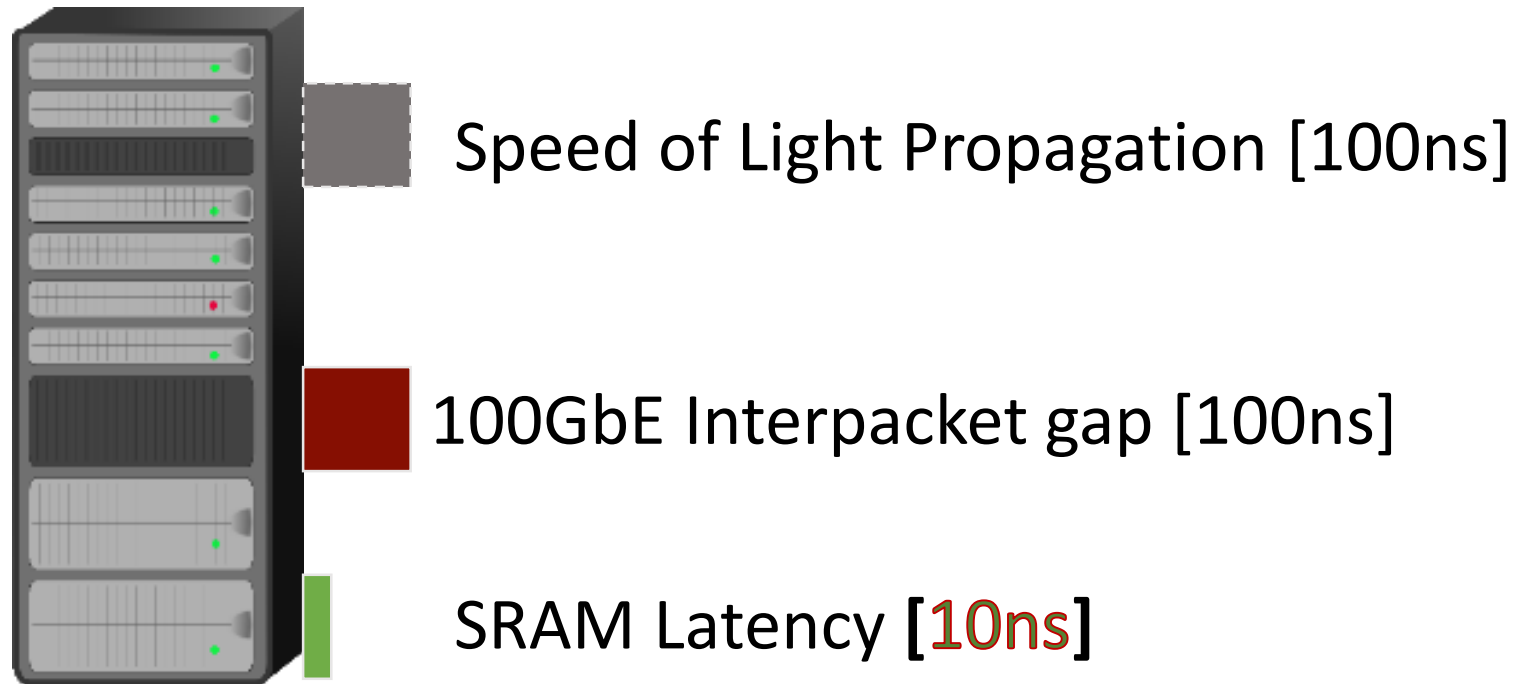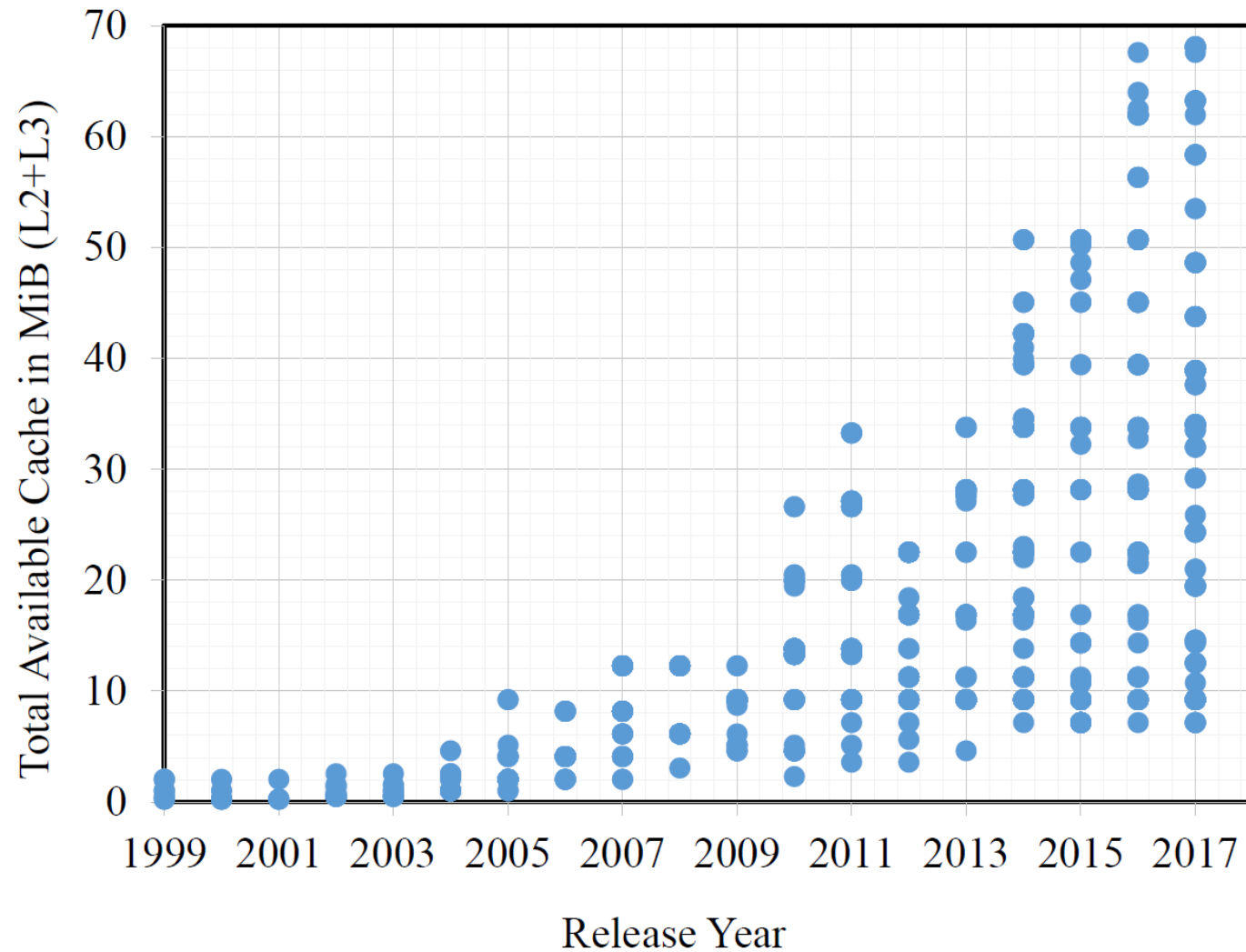
100GbE Interpacket gap [100ns]

SRAM Latency [10ns]
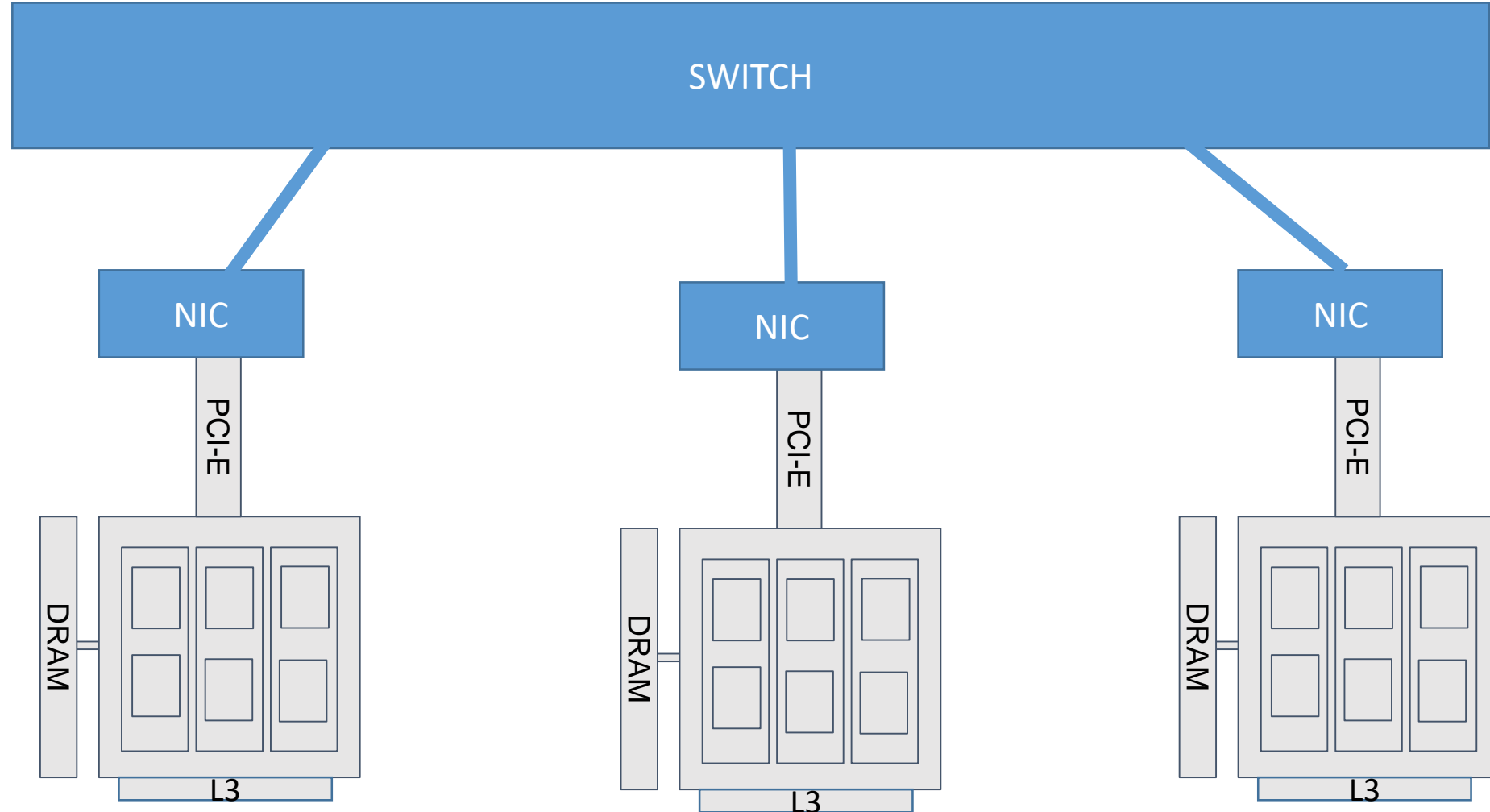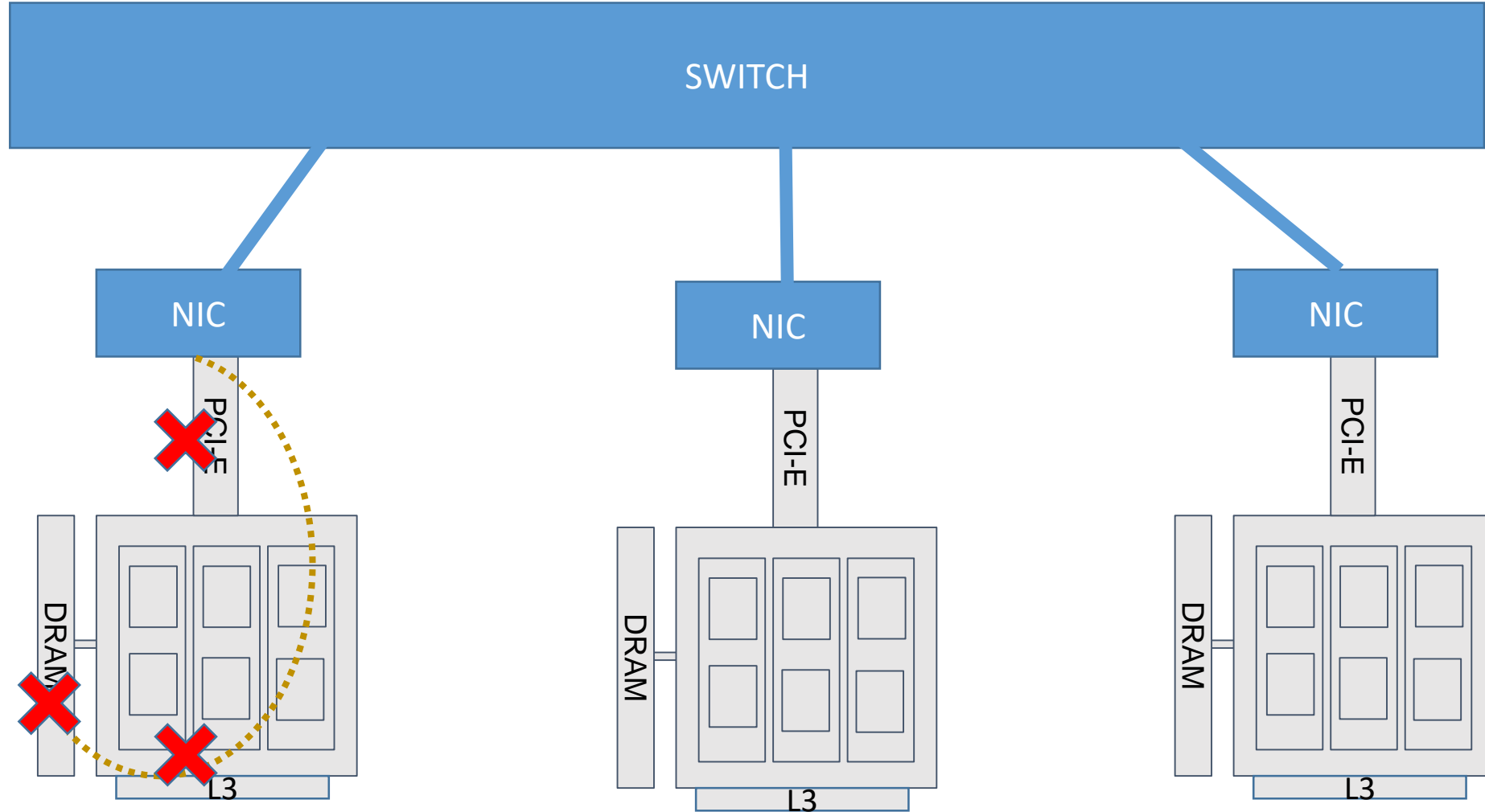
# SRAM Capacity Continues to Scale

# Conventional System Model
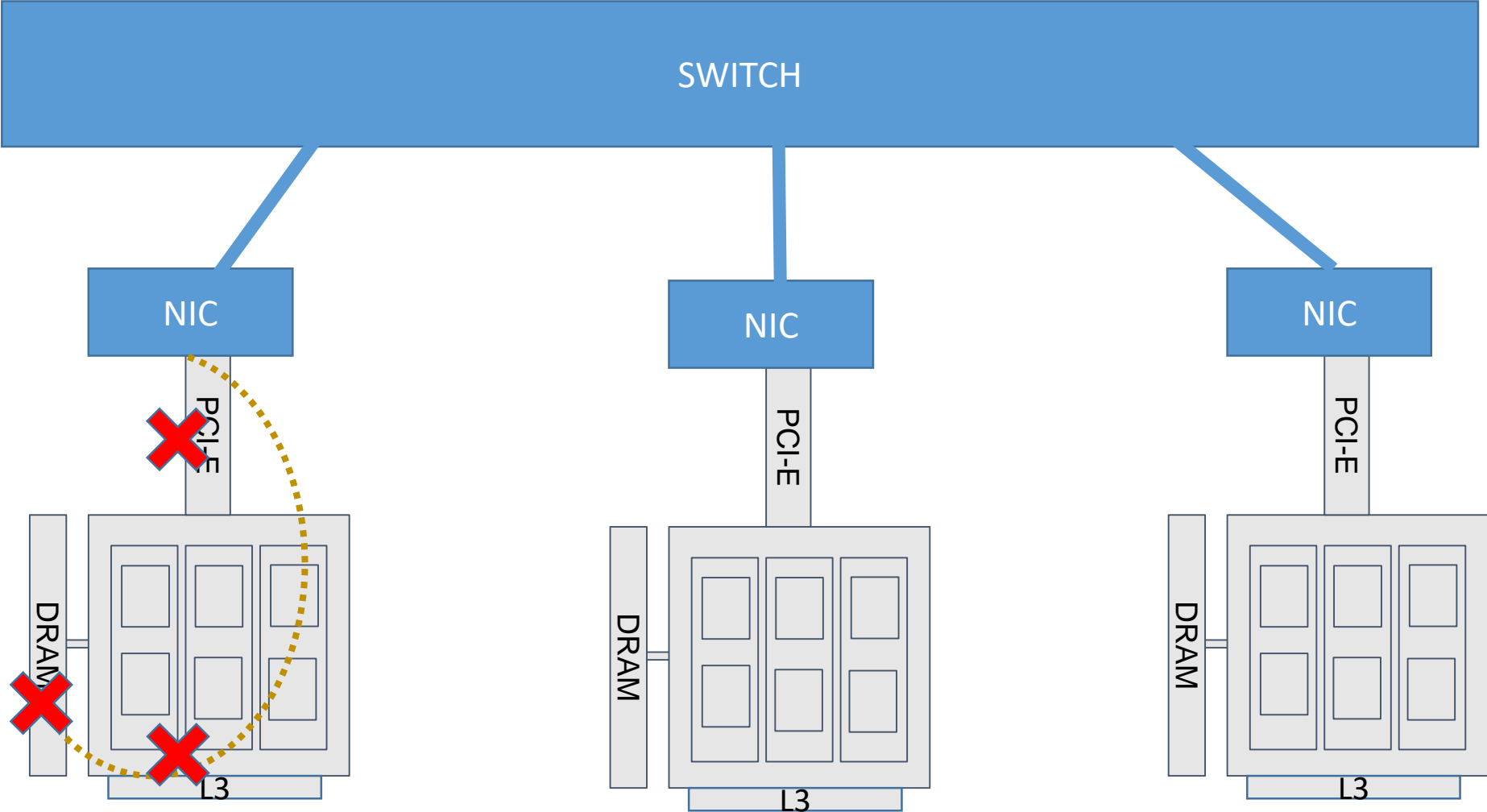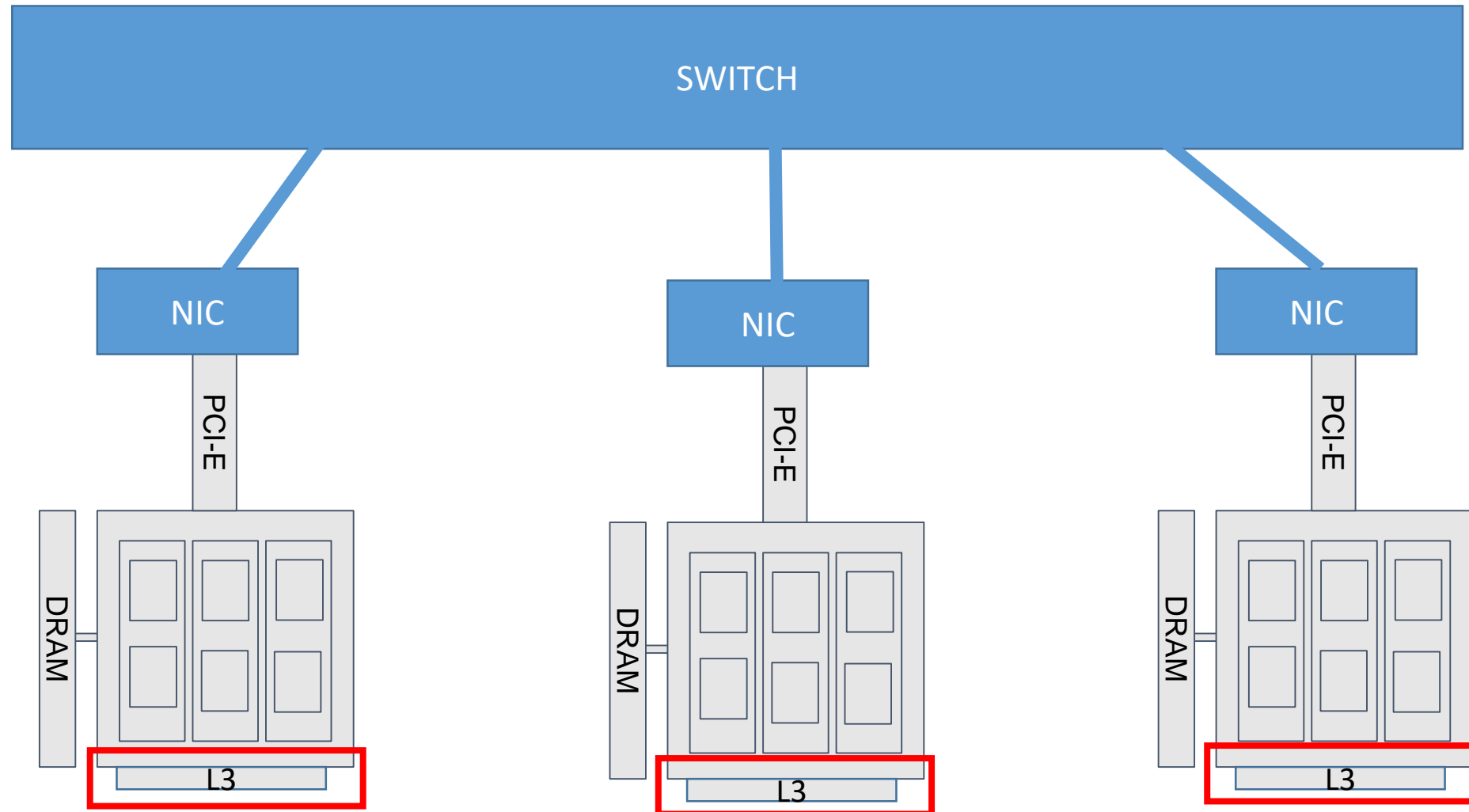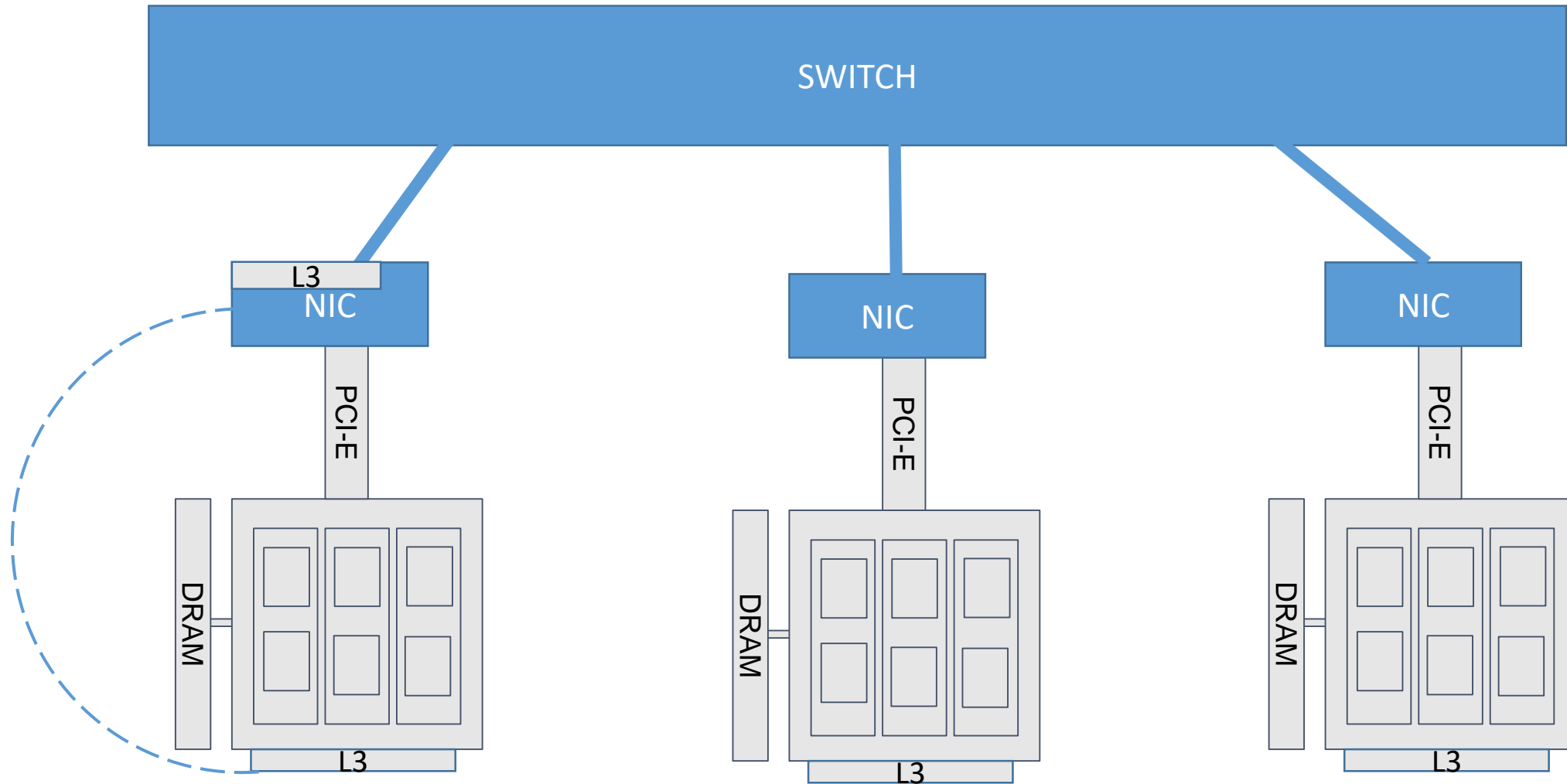
# Missing is Expensive

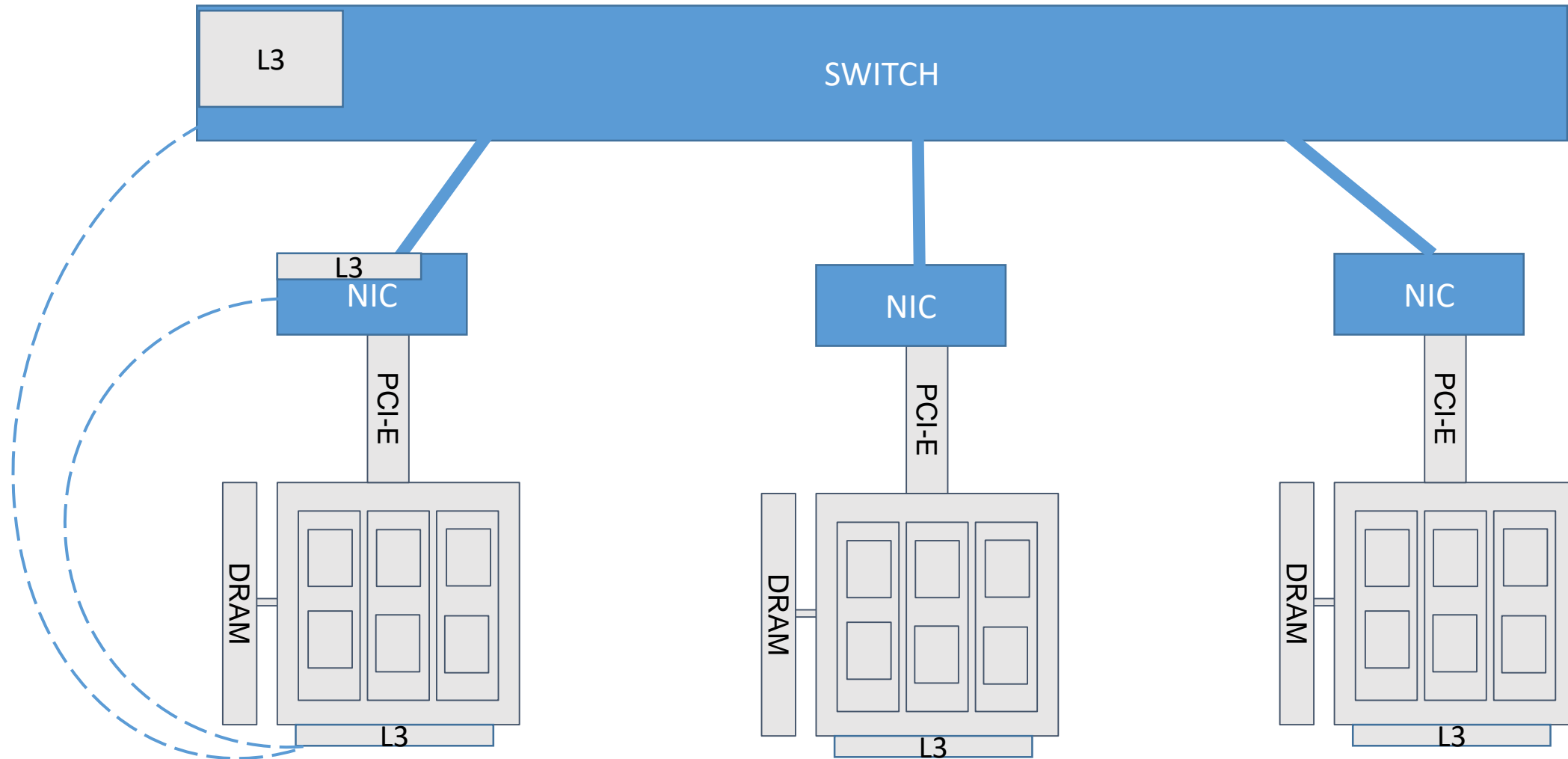# Can we Predict L3 Misses *Before They Happen?*

# Idea: Global View of L3 Caches Across Cluster

# Closer Integration with Networks and Caches

# Reroute Packets Based on In-Network Data

# Improving Interconnects

# Interconnects are Significant Bottlenecks



Speed of Light Propagation

PCIe[2000ns]

CacheCloud
Co-Packaging Network Interfaces and Cache

# CacheCloud: Takeaways

We are more than an **order of magnitude away** from speed of light propagation

Most hardware components **have stopped scaling** while networks **scale up**

**Cache is the new DRAM, DRAM is the new Disk**

Architecture, Network and OS **integration** is crucial

# Issues and Feedback for CacheCloud

Q: Won't this be solved with new hardware?

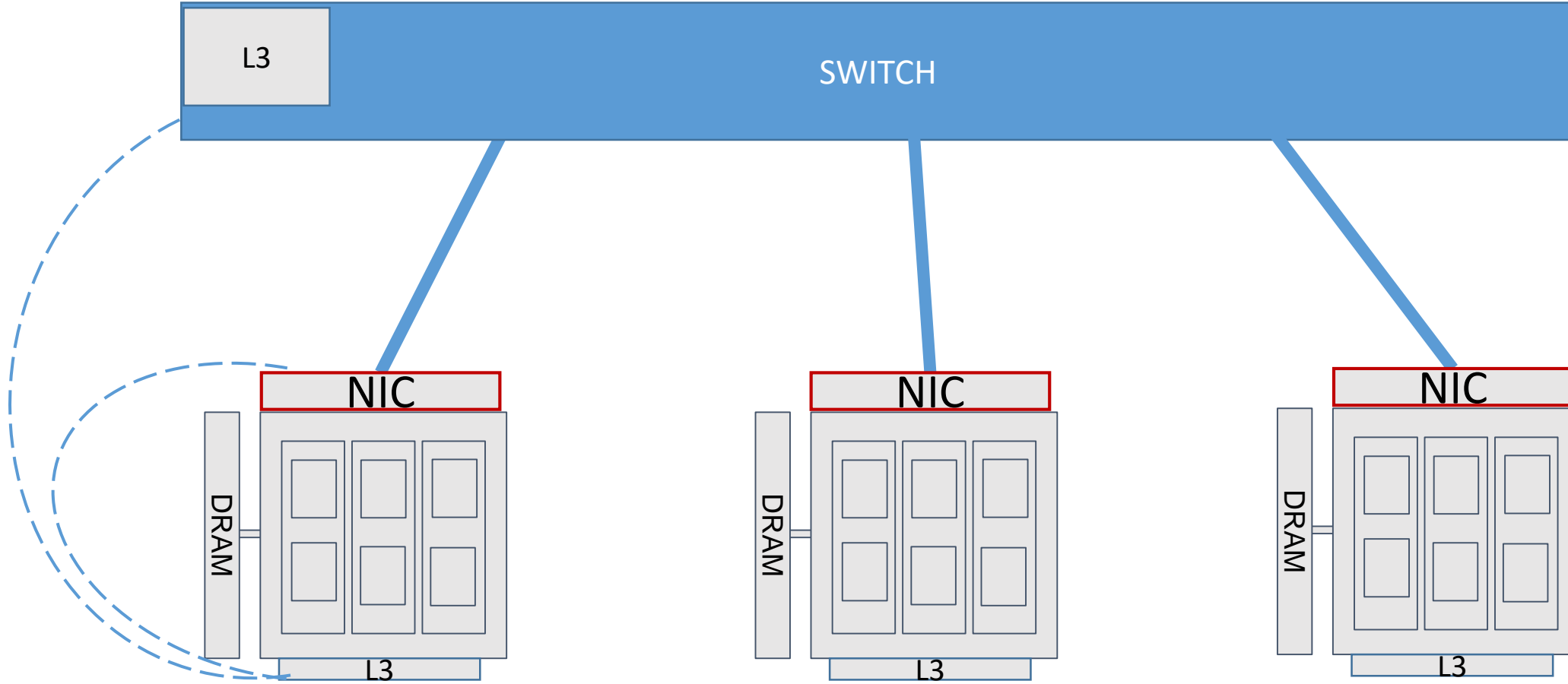- SmartNICs, Accelerators, FPGAs, etc.
- Coherency, programmability, size, and cost still a problem

Q: Datasets are too large….but do all applications have large dataset?

- Network function virtualization (NFV)
- Coordination services

Feedback:

- How do hardware accelerators (TPU) change network and end-host design?
- Starting from scratch what architecture should we build?
- Starting from scratch what network protocols should we build?

# Thank You