# Enhancing Distributed Computing with Optical Interconnects in Data Centers and High-Performance Computing (HPC) Systems

Abidur Rahman, Yuyang Wang, Brian Wu, Alex Meng

Lightwave Research Laboratory

Department of Electrical Engineering, Fu Foundation School of Engineering and Applied Science, Columbia University

**COLUMBIA | ENGINEERING**
The Fu Foundation School of Engineering and Applied Science
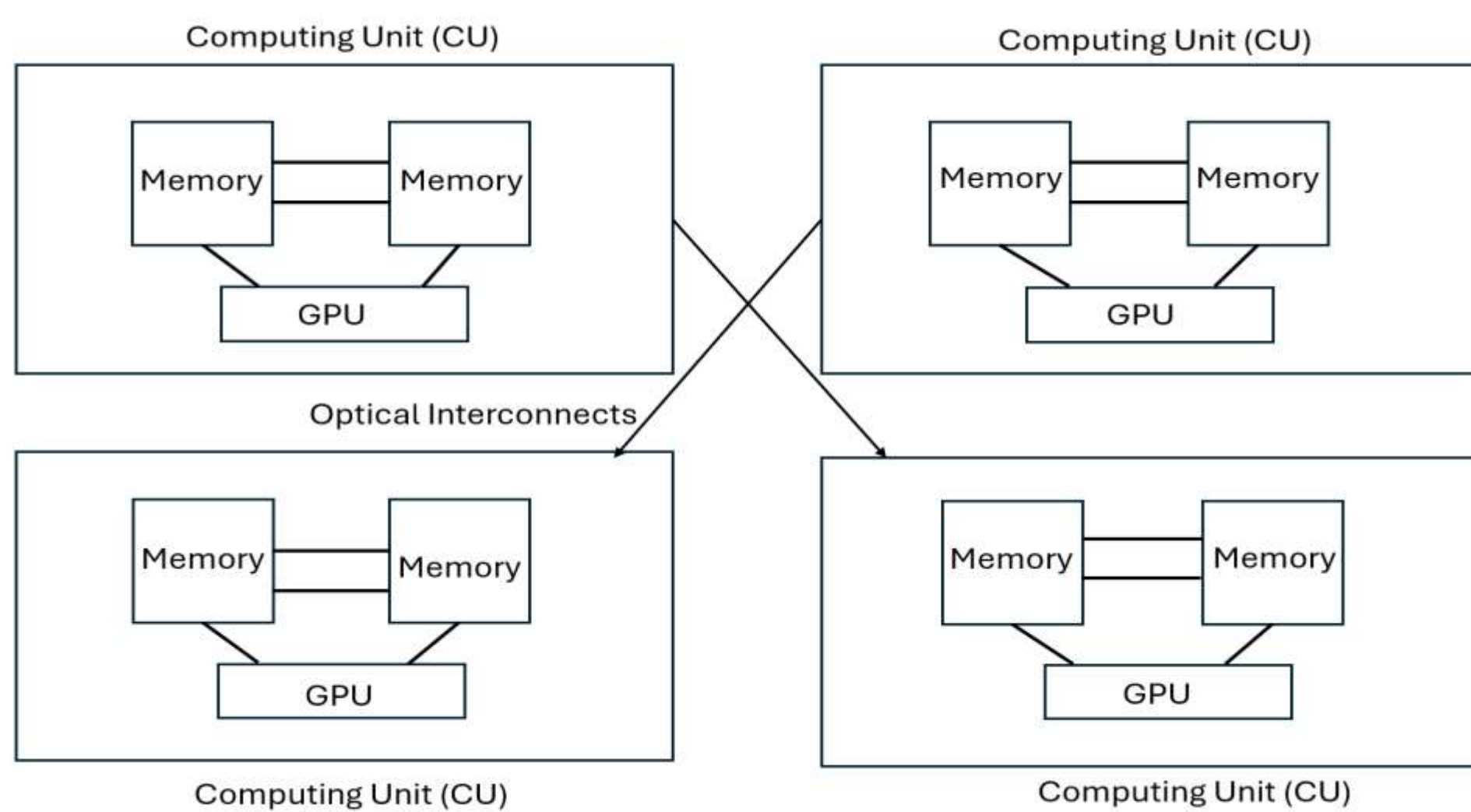
**amazon**

## Introduction

**Problem:** Off-chip electrical bandwidth links are not keeping up with compute power

**High Bandwidth:** Optical interconnects offer significantly higher bandwidth compared to traditional electrical interconnects, allowing for the transmission of larger volumes of data simultaneously.

**Low Latency:** Optical signals travel faster than electrical signals, resulting in reduced communication delay and improved synchronization in distributed systems.

**Scalability:** Optical interconnects can efficiently scale to support a larger number of computing units without the performance degradation seen in electrical interconnects.

**Energy Efficiency:** Optical interconnects consume less power, especially over longer distances, making them more energy-efficient and reducing the overall power consumption of data centers.



Computing Unit (CU) Group
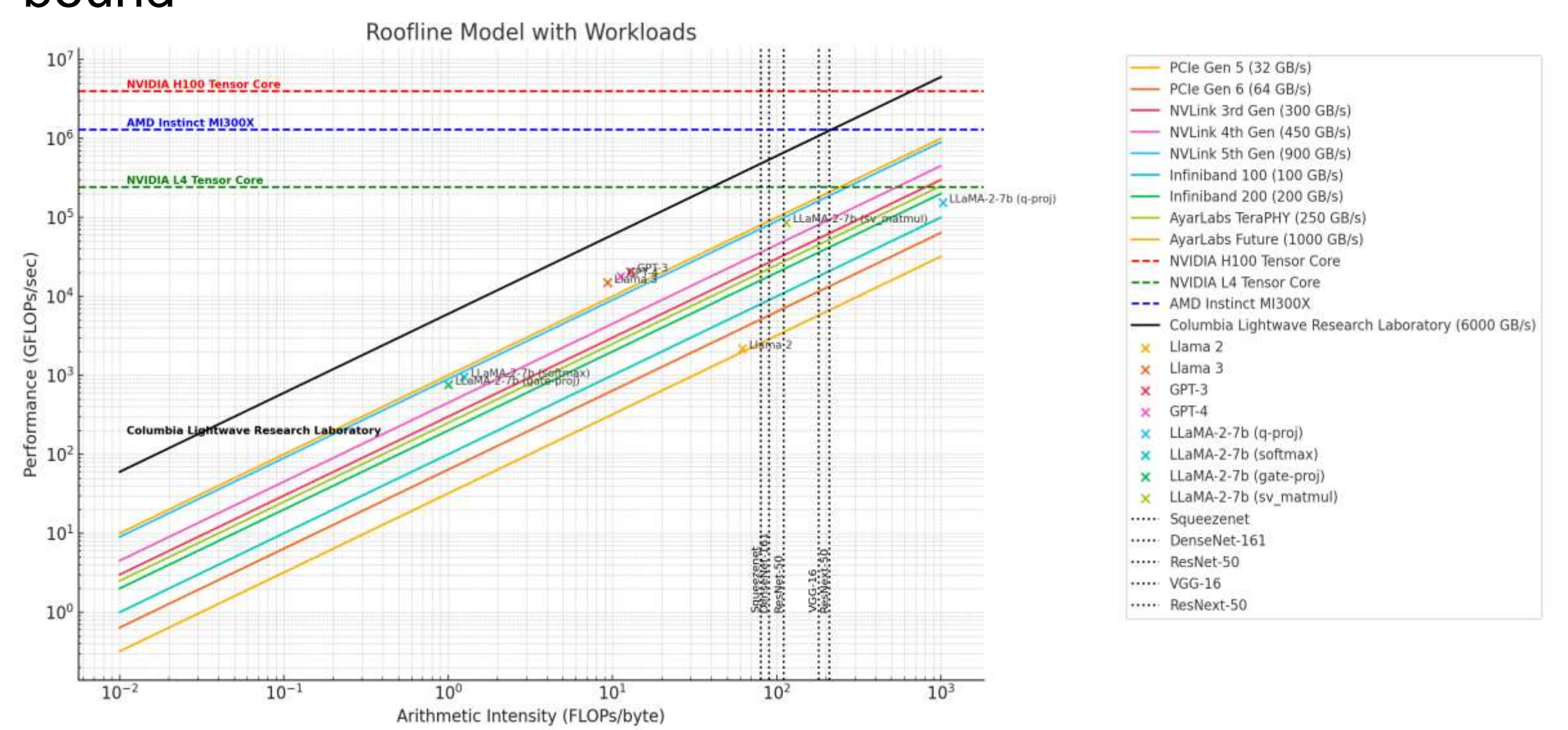


GPU bandwidth trend

## Methods

**Roofline Model**
The Roofline model is a visual representation used to analyze and optimize the performance of software on modern processors

**Purpose:** The Roofline model helps to identify the maximum achievable performance of a computing system for a given application, considering both the computational capability of the processor and the memory bandwidth.
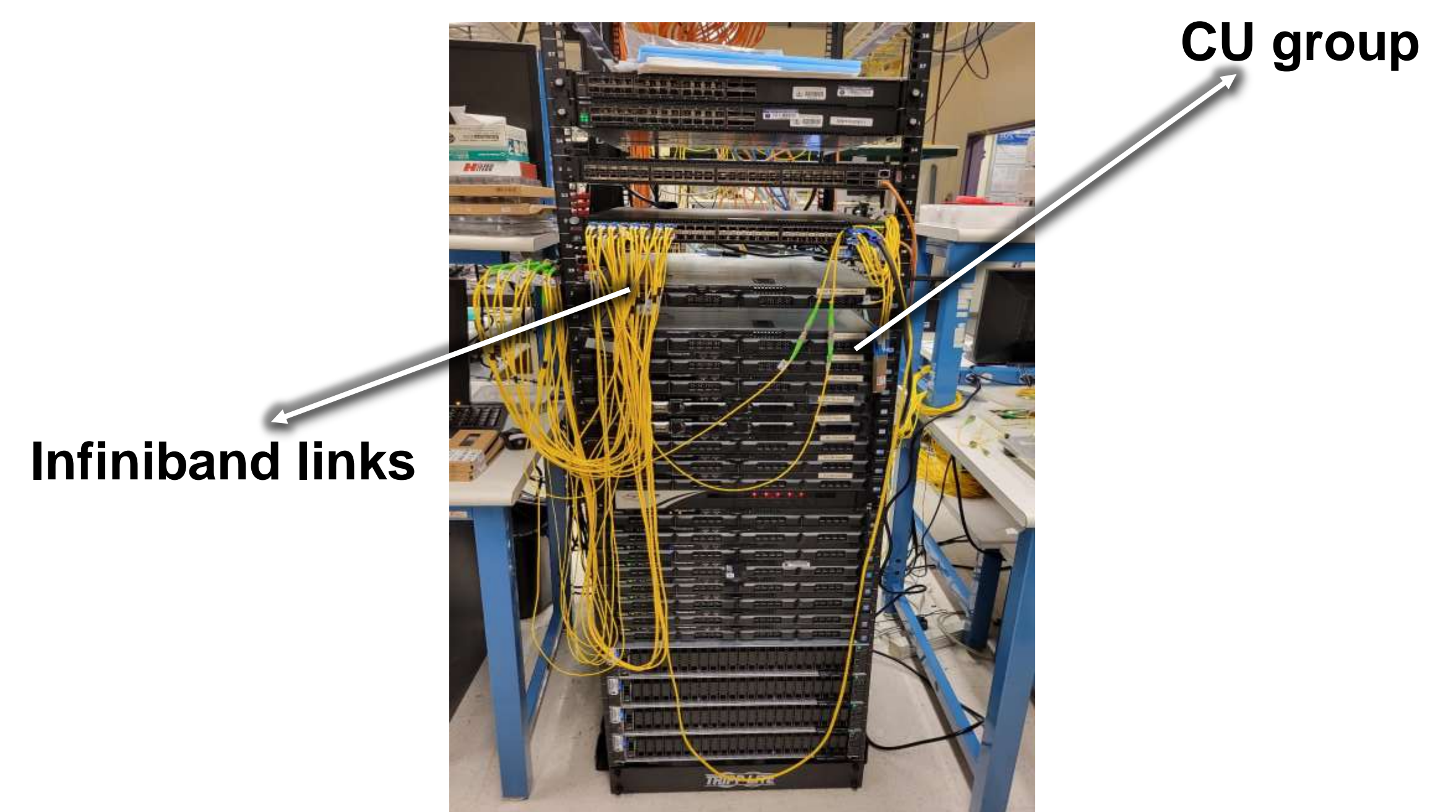
## Results

Left side: Workload Performance is Bandwidth-bound
Right side: Workload Performance is Compute-bound



Roofline Model

## Lightwave Research Laboratory Compute Cluster



CU group

Infiniband links

## Conclusions

The transition to optical interconnects in data centers and HPC systems addresses critical challenges related to bandwidth, latency, scalability, and energy efficiency, making them an optimal choice for supporting the next generation of high-performance computing and AI workloads.

## Acknowledgments