

Networking for AI Cloud – Why do we need a different network

Barak Gafni

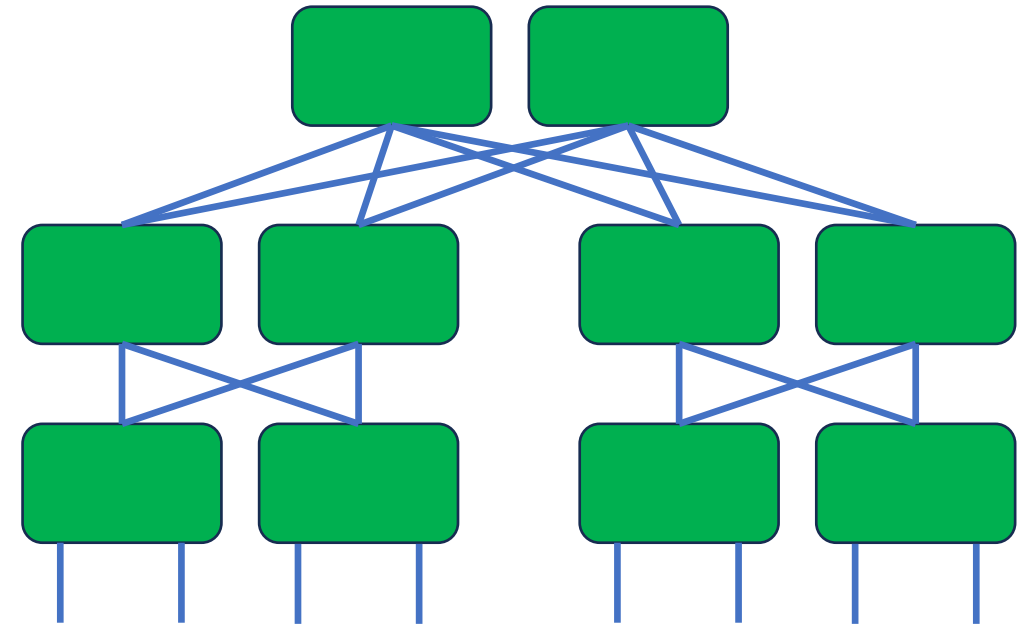
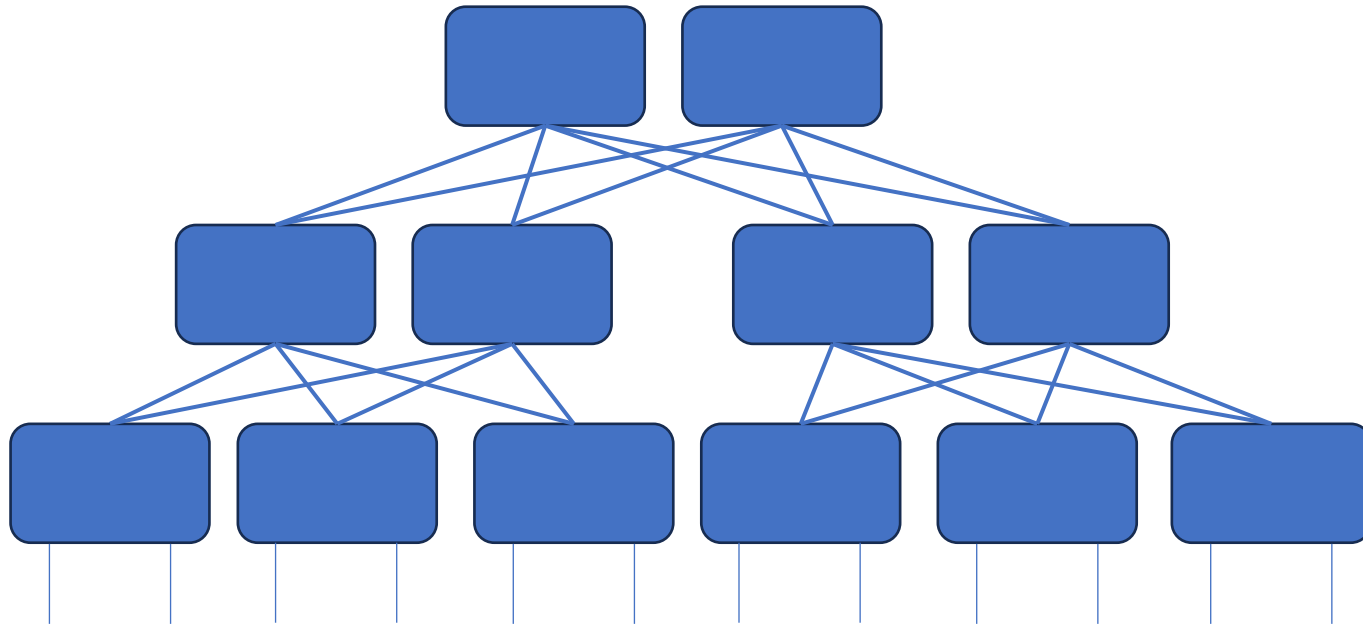
IETF-117

July 2023

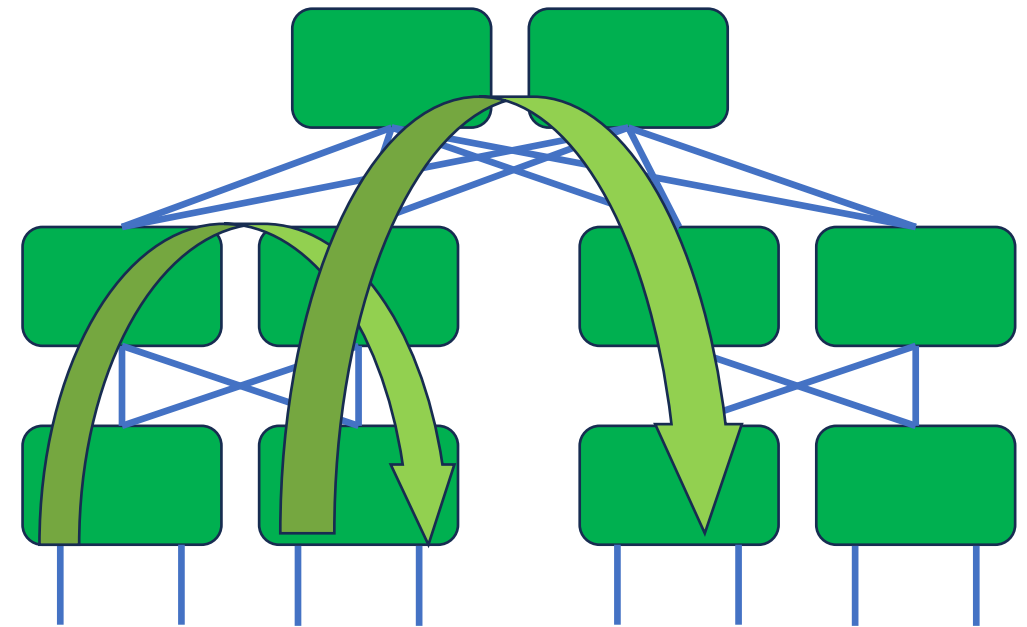
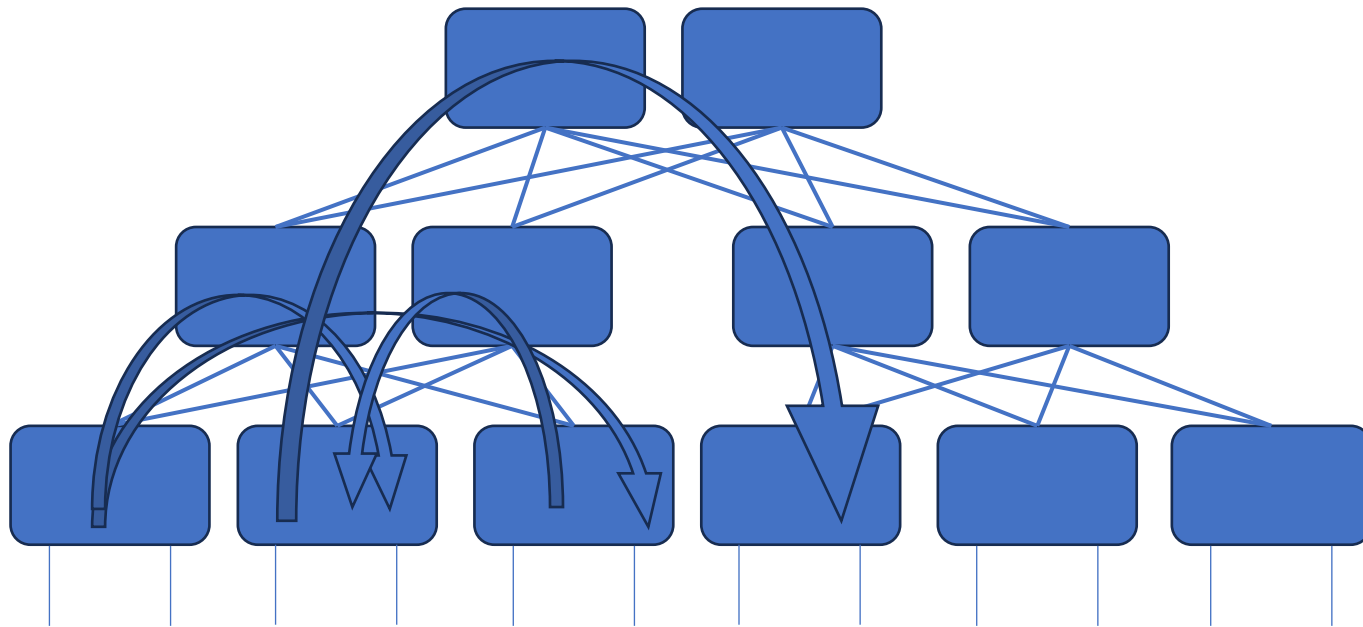
AI Cloud Networking Requirements

- High BW
- Low latency
- Lossless network
- Predictable performance
- Robust network

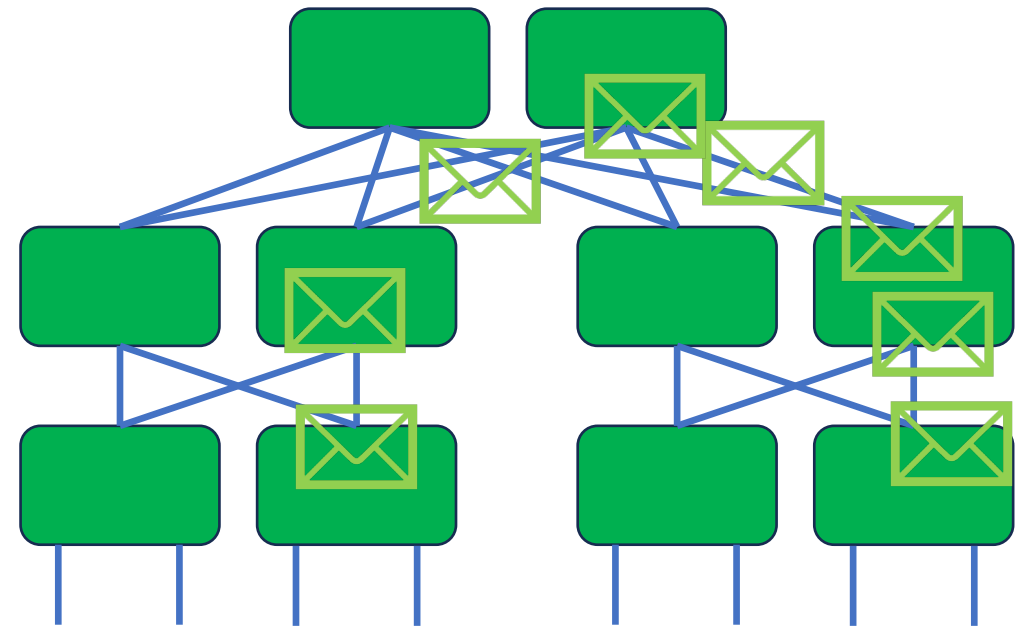
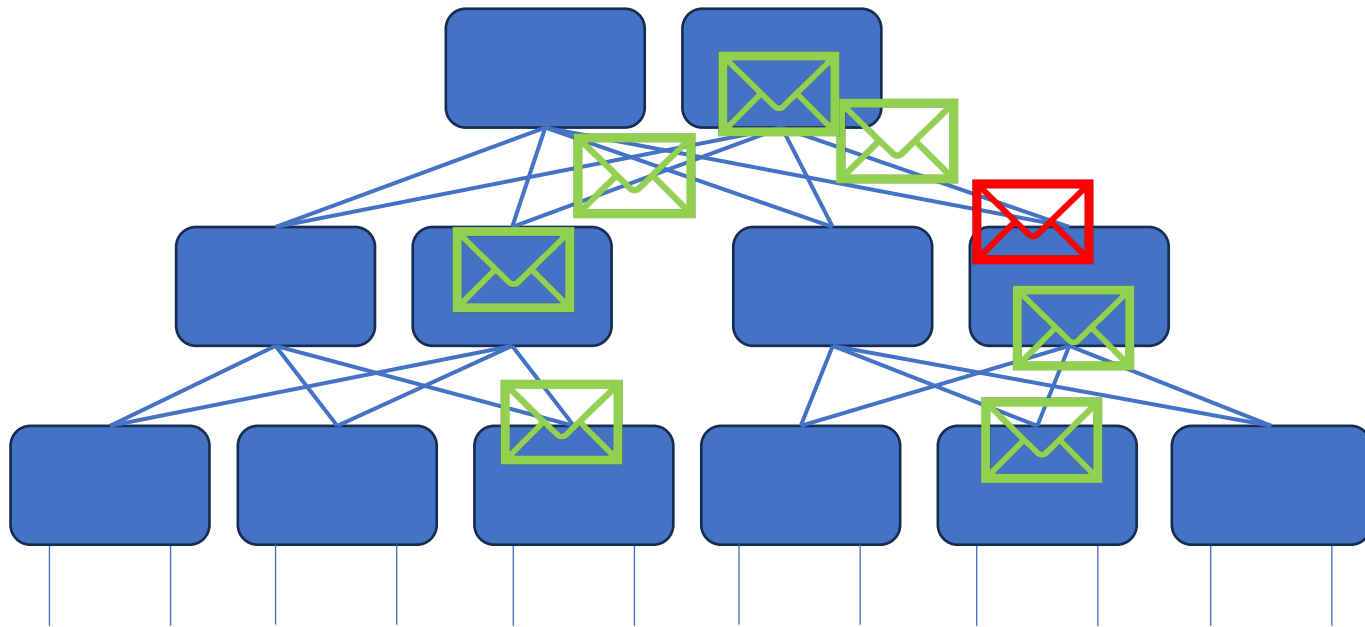
Oversubscribed vs non-blocking network



Low BW TCP/UDP vs High BW RDMA

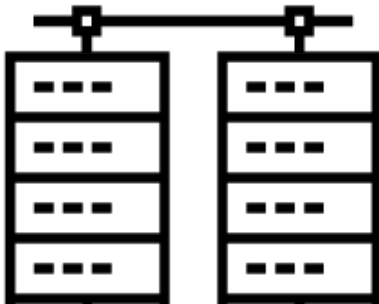


Lossy vs lossless



Additional observations on AI networking

General purpose cloud network	AI cloud network
Designed for Loosely-Coupled Applications	Designed for tightly-coupled distributed applications
TCP/UDP-Focused Network, Low Effective Bandwidth	Highly Effective RDMA Bandwidth
High Jitter Tolerance	Low Jitter Tolerance
Switch/Link failure impact is moderate	Switch/Link failure impact is dramatic
Designed for best effort	Designed for High Performance at scale, at load



Reminder, there are many networks to operate an AI cloud

- GPU Compute
- General compute
- Storage
- Management
- More...

Thank You