



From Datacenter to AI Center

Building the networks that build AI

Tyler Conrad
Principal Engineer

Datacenter Design Principles

Speeds

10/25/100/400G

Redundancy

MLAG / ESI

Services

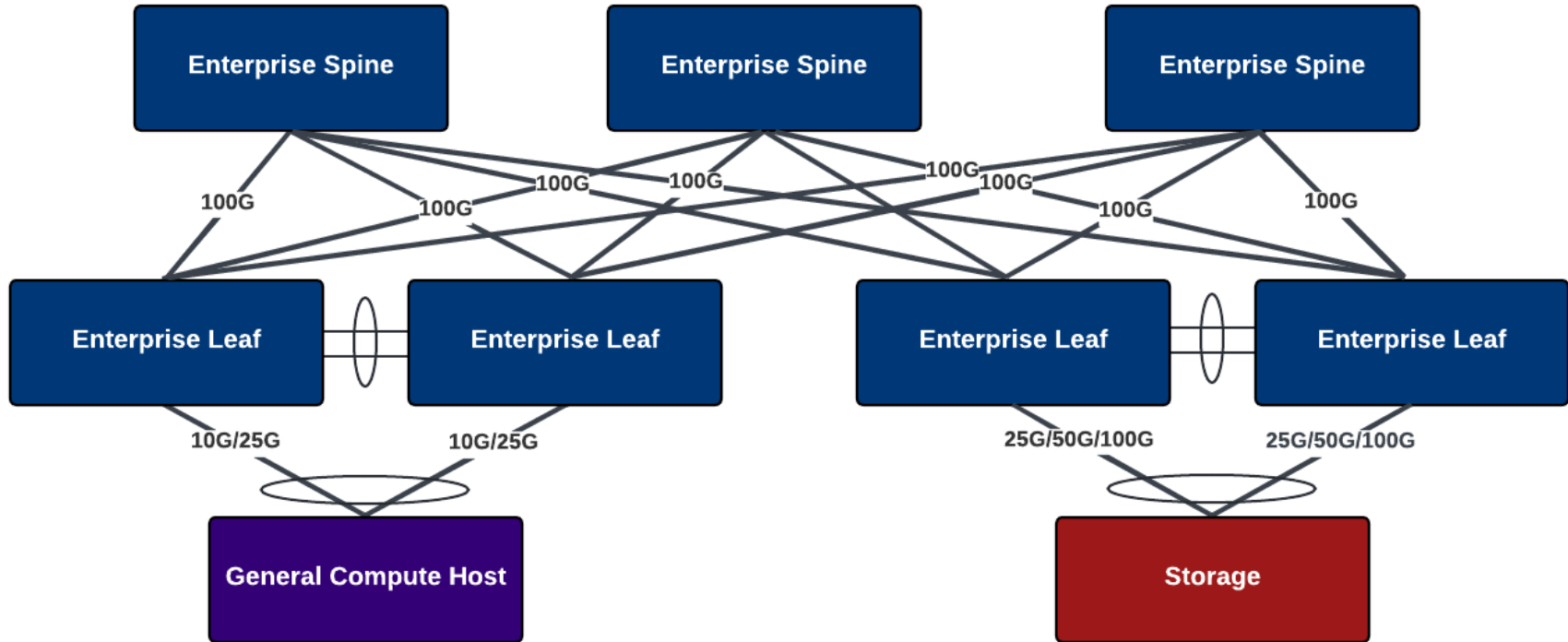
VXLAN-EVPN

- VLAN & VRF Extension

Capacity

~3:1 Oversubscription

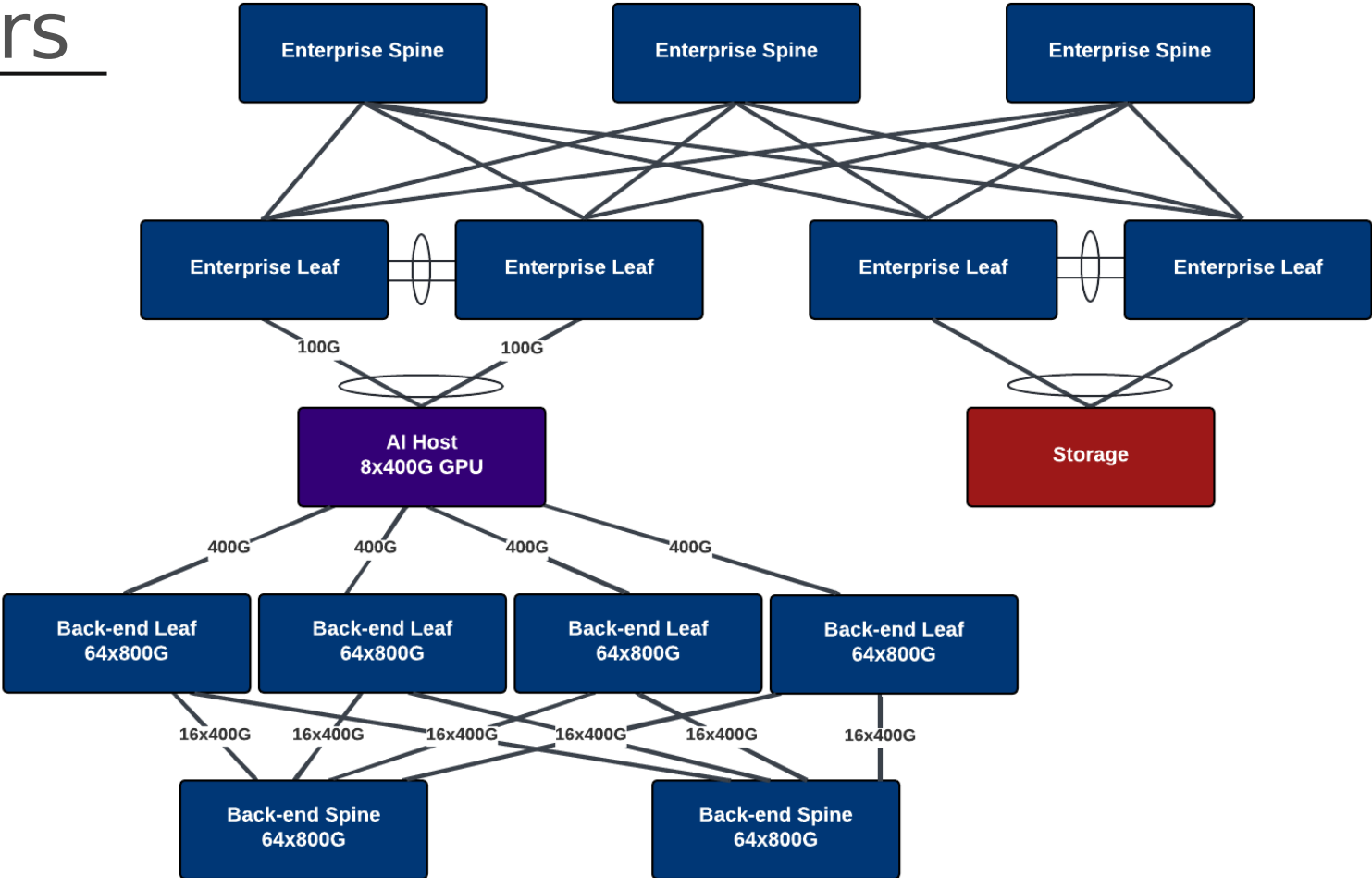
Datacenter Networking



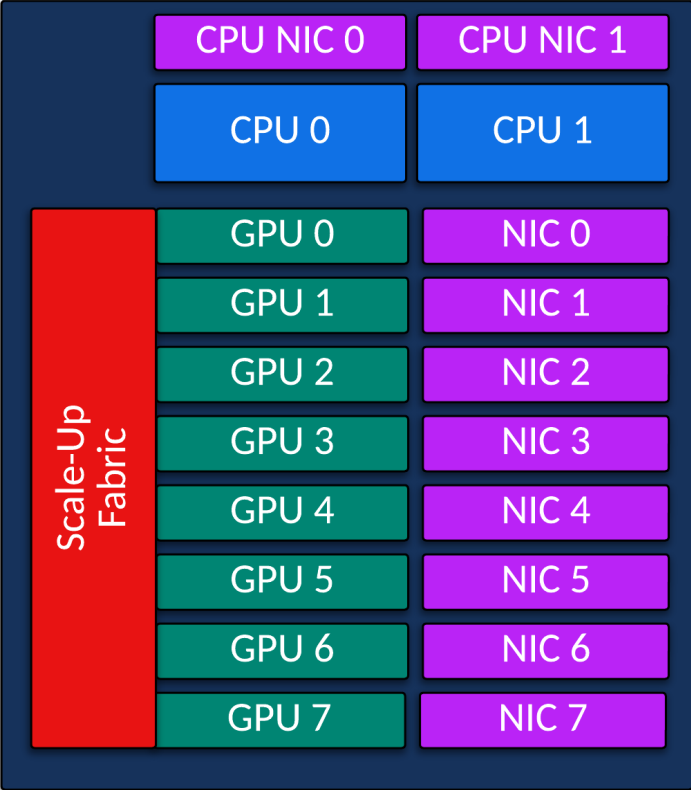
AI Centers

Front End
CPU

Back End
GPU



What Do These Servers Look Like?

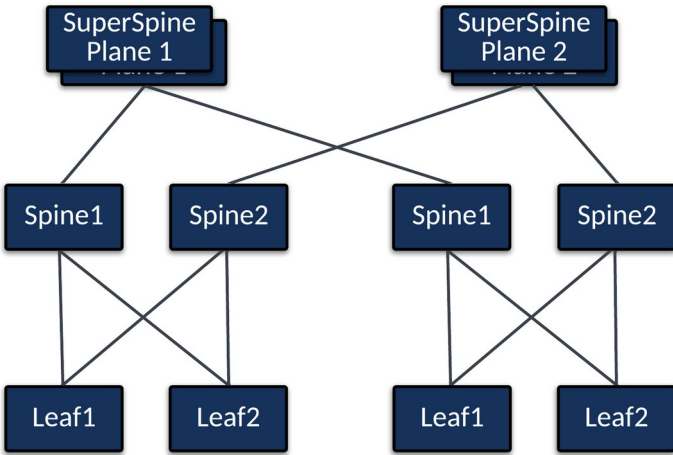
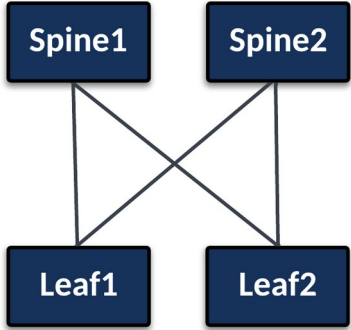


Fixed vs Rack-scale Architecture

Topologies

Single-Box vs 2-Tier vs 3-Tier

Modular Leaf



<1K GPU

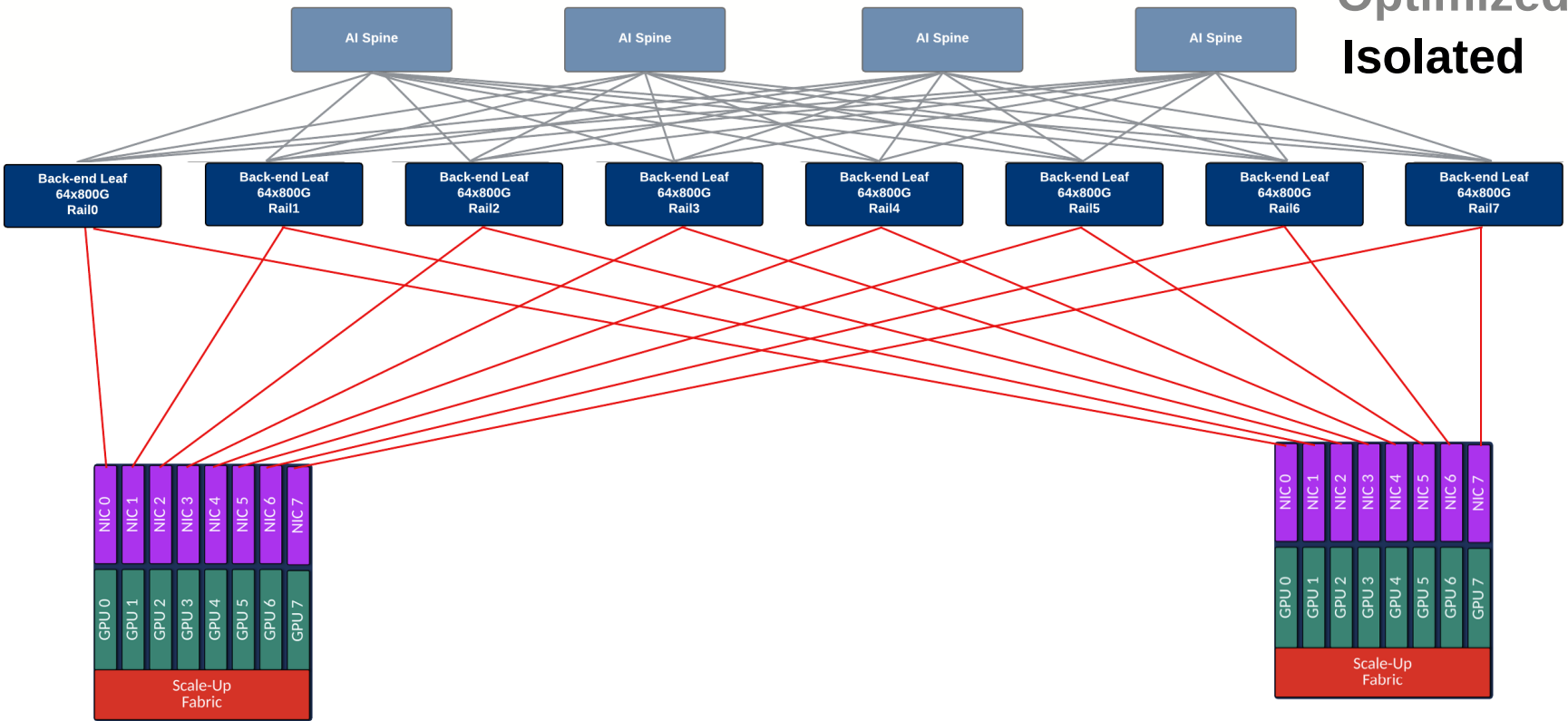
<8K Fixed
<73K Hybrid
<663K Modular

<524K Fixed
>100M Modular

Normalized for a 400G GPU

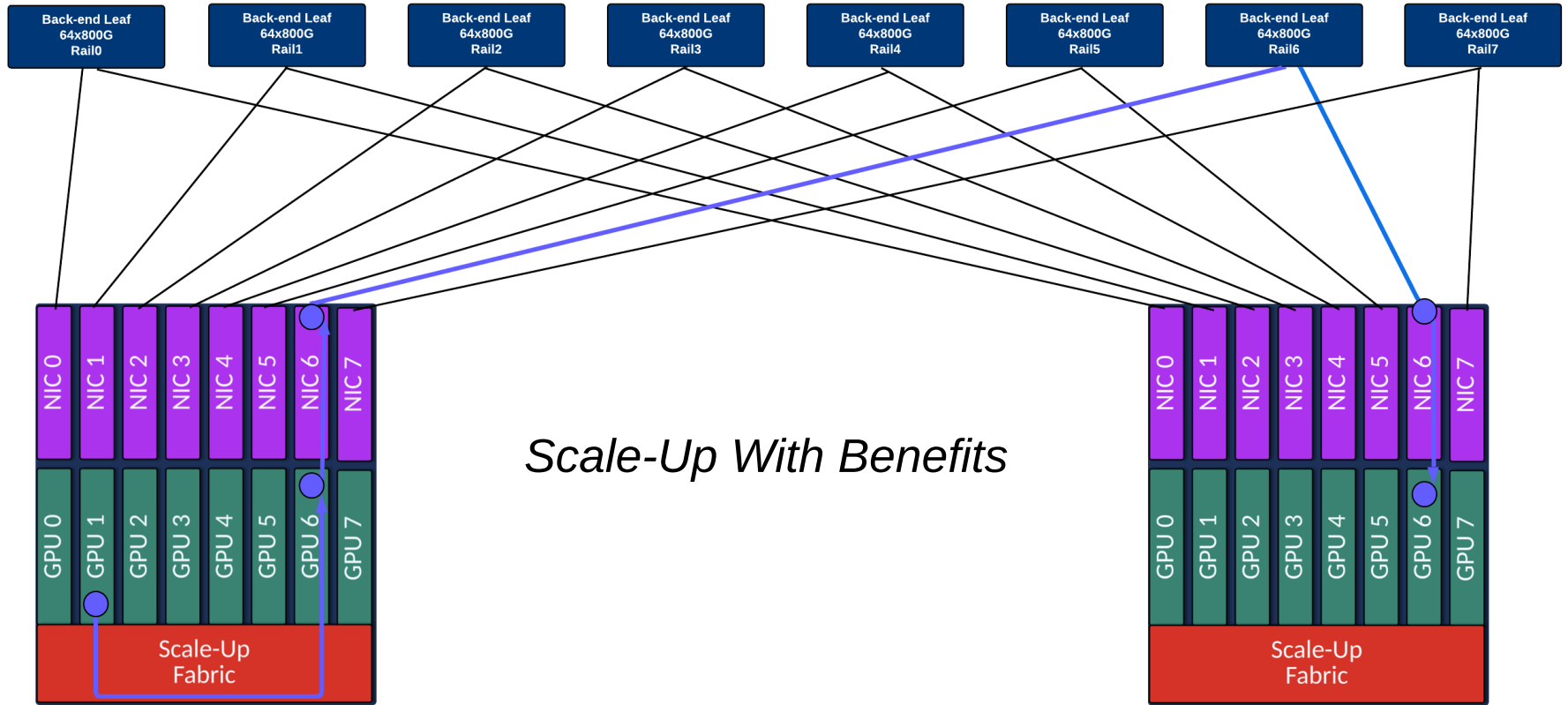
Rail Example

Optimized
Isolated



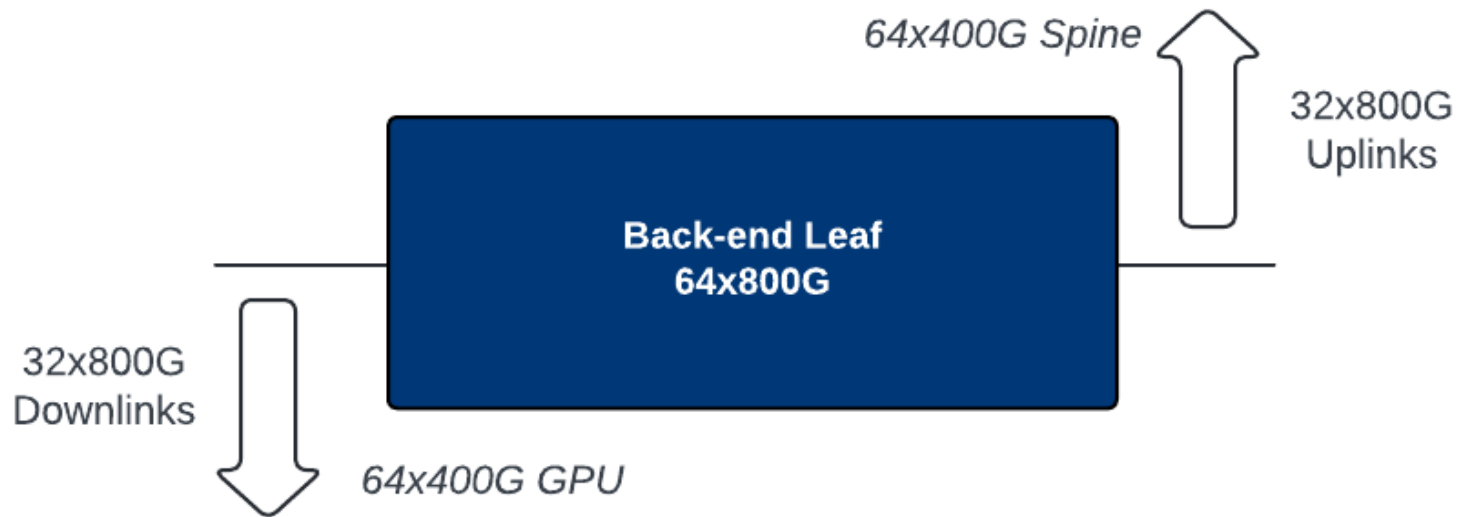
Each GPU in the chassis connects to one leaf.

PXN - PCI-E Extended Networking



Use another GPU NIC to reach other Rails

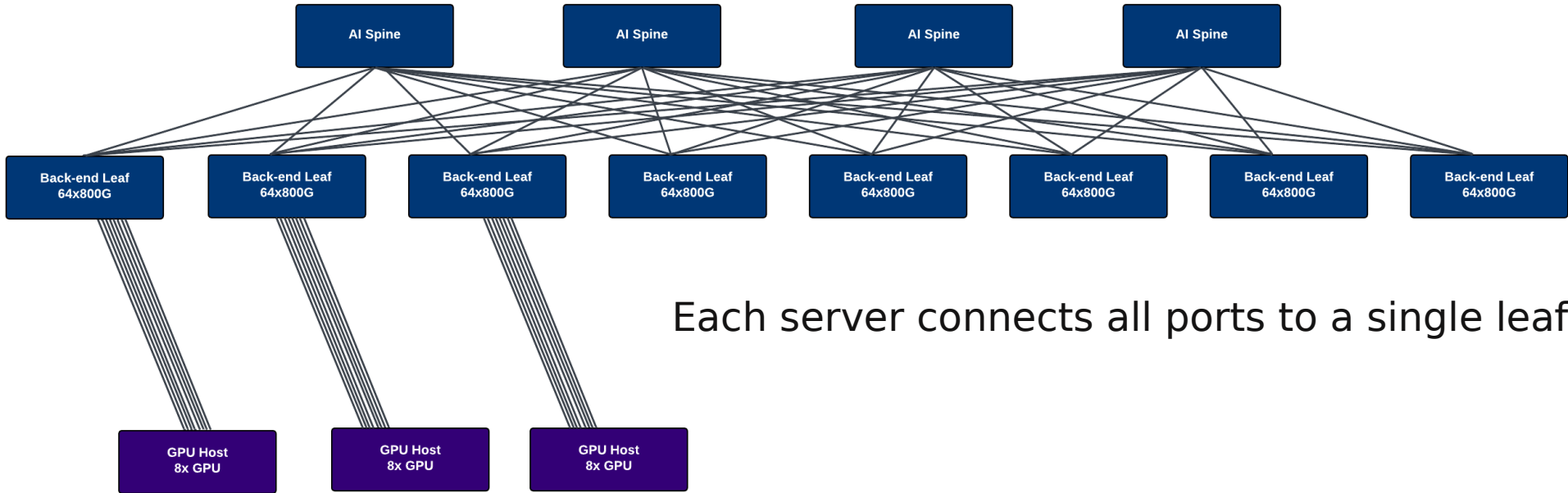
Subscription Ratios



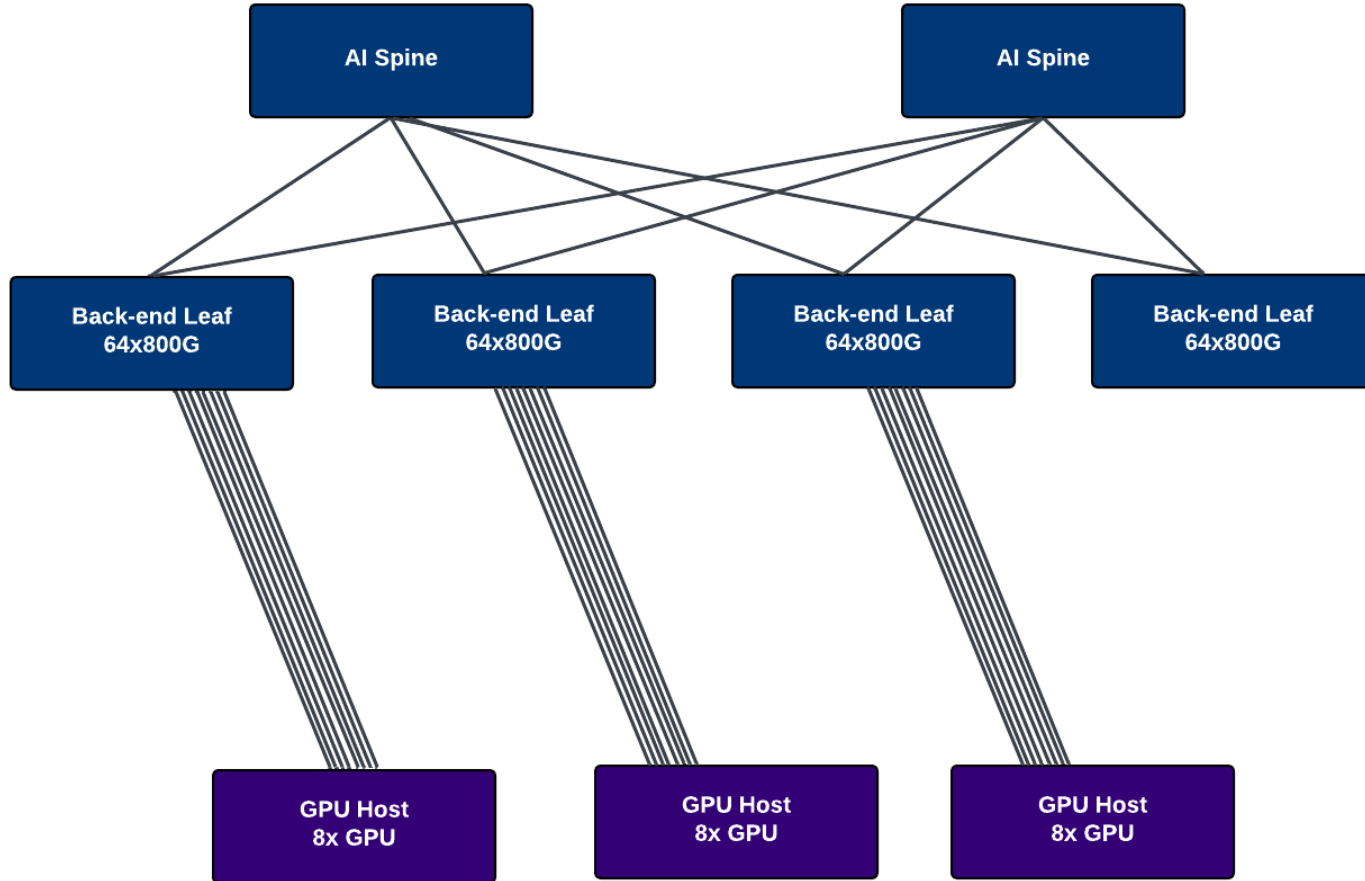
Target is 1:1 or better for uplink to downlink ratios*.

*Unless you really know your workloads

Non-Rail / “Fat Tree” Example



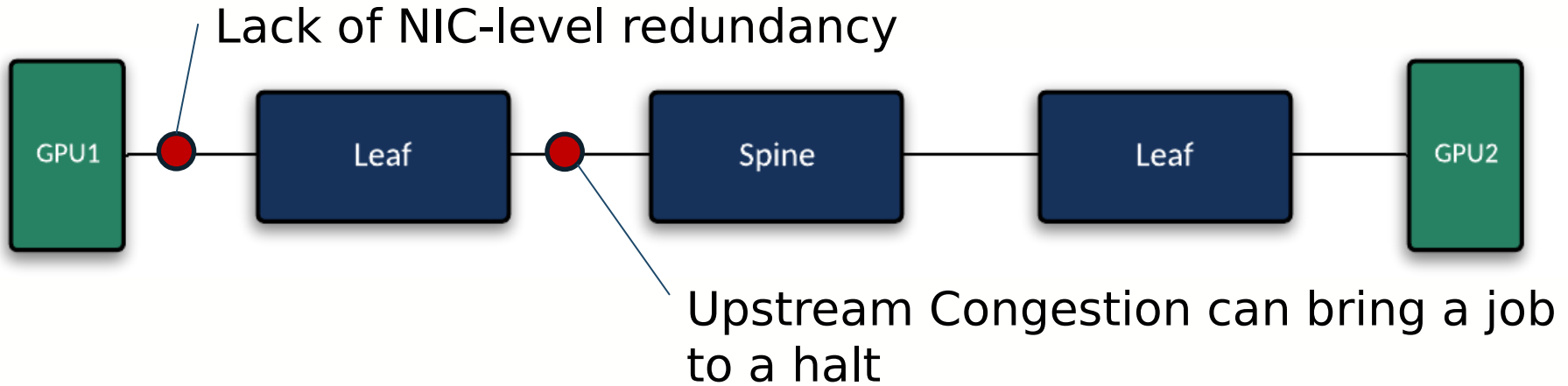
Scaling Topologies Down



Introduction to Multi-Plane

2-Tier vs 3-Tier:

- Lower Latency
- Less Power
- Fewer(ish) Optics/Switches

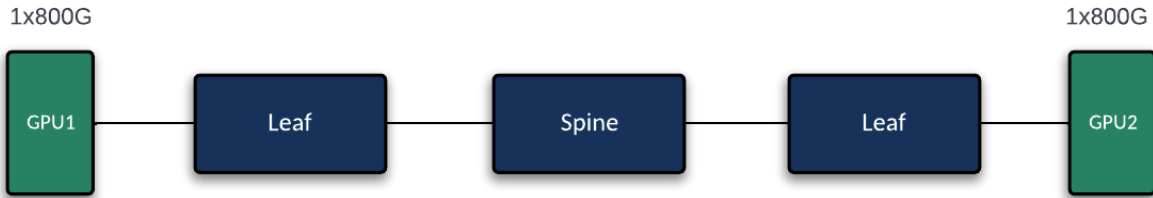


Higher speeds = lower radix.

64x800G Switch w/400G Uplinks in 2-tier = 8192 GPU

64x800G Switch w/800G Uplinks in 2-tier = 4096 GPU

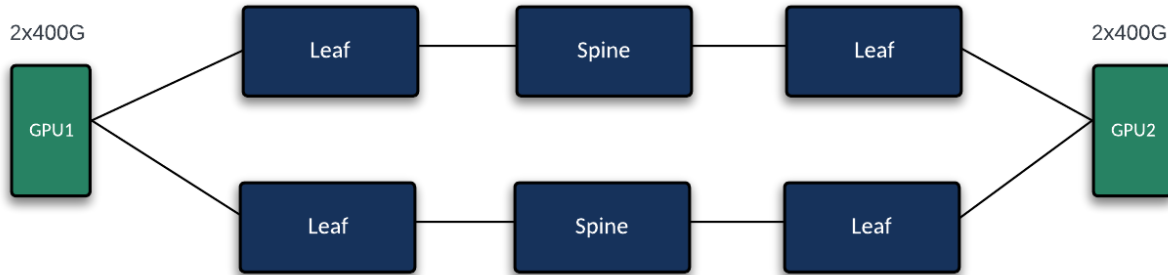
Make like a banana and...



64x800G

- 32 Spine per plane
- 64 Leaf per plane (32 downlinks)

1x800G Max Scale (64x800G 2-tier) = 2048 GPU

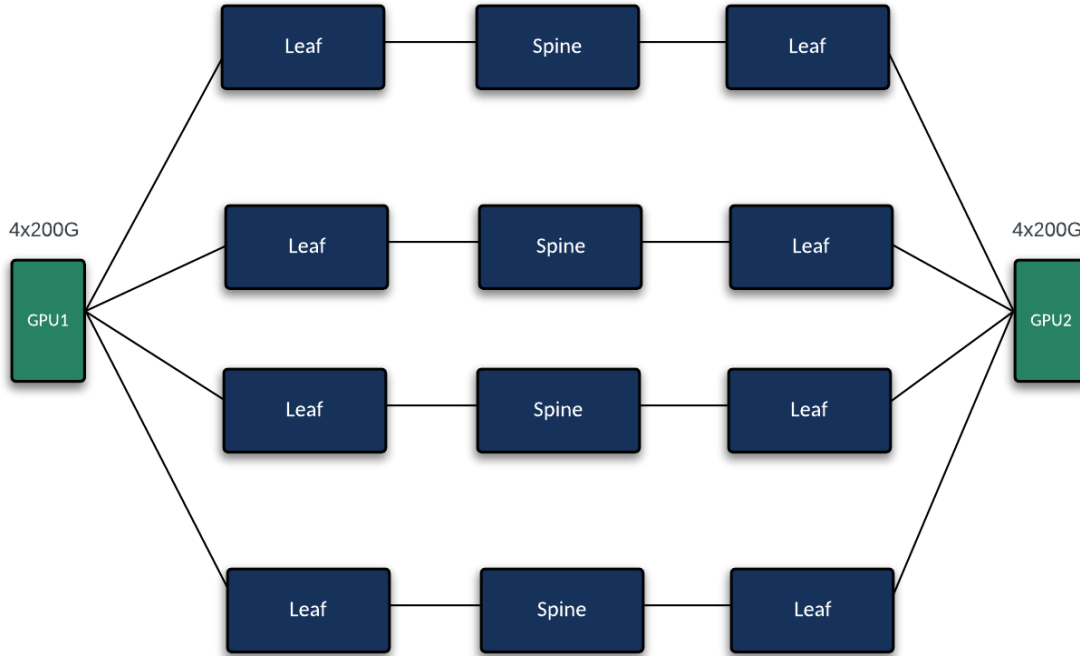


128x400G

- 64 Spine per plane
- 128 Leaf per plane (64 downlinks)

2x400G Max Scale (64x800G 2-tier) = 8192 GPU

Taking it a step further...



256x200G

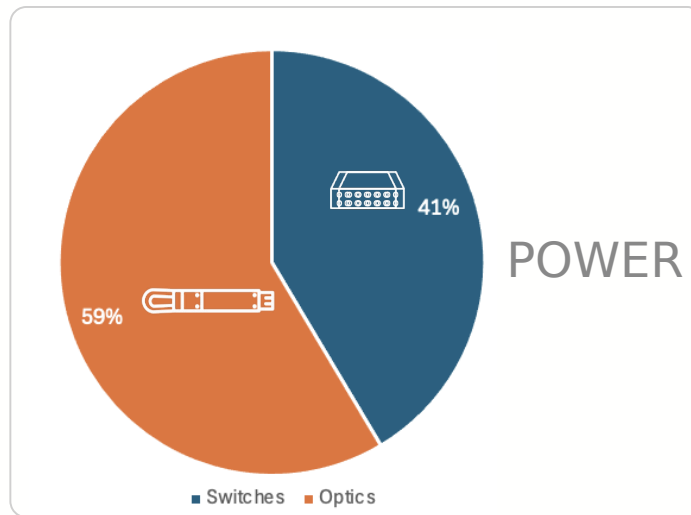
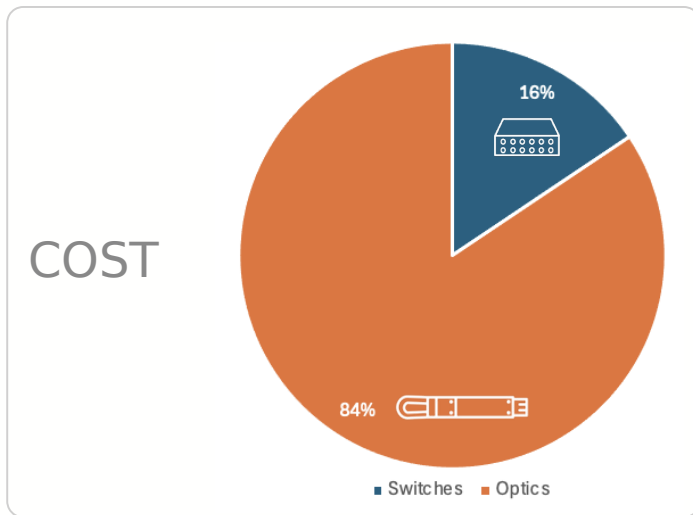
- 128 Spine per plane
- 256 Leaf per plane (128 downlinks)

4x200G Max Scale (64x800G 2-tier) = 32768
GPU

Optics

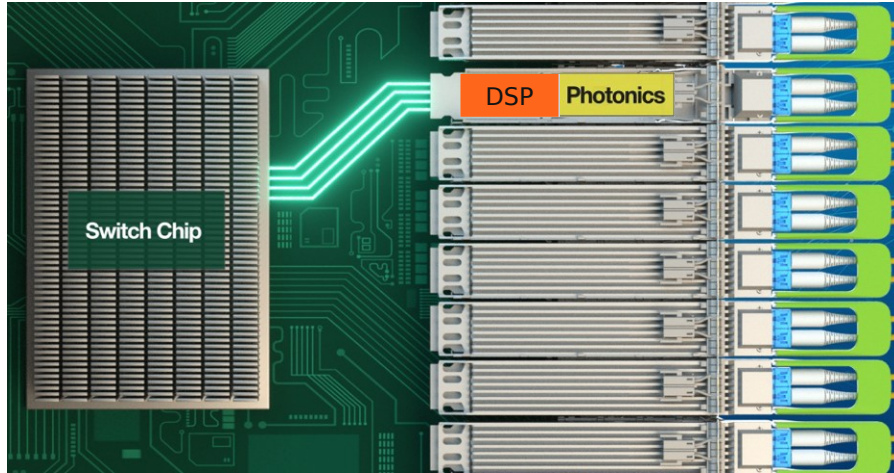
Optimizing optics

- Ex: 8k XPU ports / 2-tier AI-center infrastructure
 - 64x 64-port 800G AI Spine
 - 128x 64-port 800G AI Leaf
 - 16384x 800G optics
 - 8192x 800G optics (Fabric ports)
 - 8192x 800G optics (Host-facing ports)

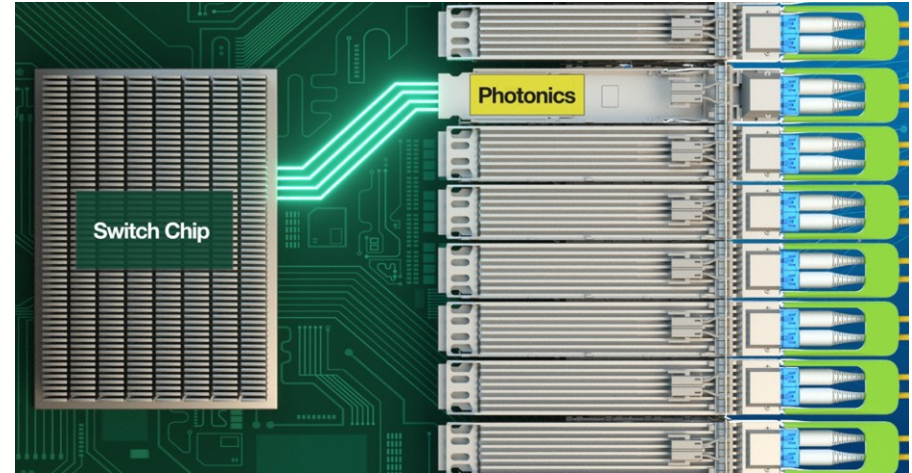


Linear-drive Pluggable Optics (LPO)

Traditional pluggable optical modules



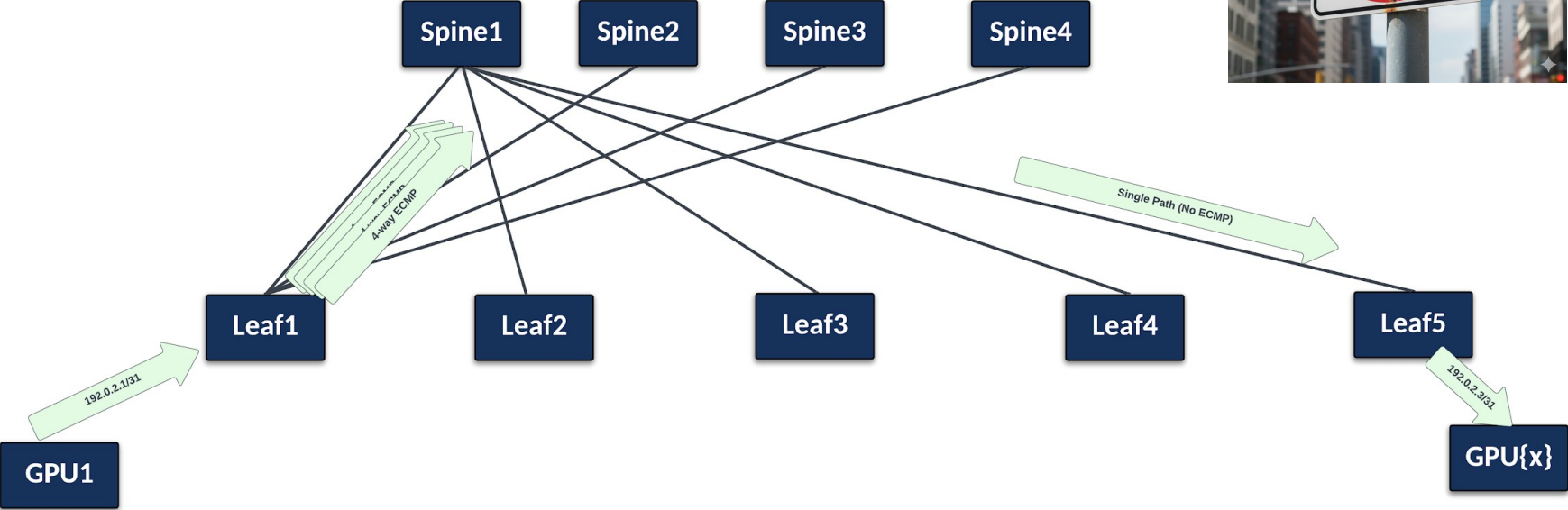
LPO modules



- LPO means no DSP in the optic module
- How is this possible ?
 - Certain switch silicon has advanced DSP technology on-chip
 - Requires careful system design and SerDes tuning
- Lower power ($\sim 0.5x$), cost ($\sim 0.7x$) and latency ($\sim 0.01x$) with higher reliability

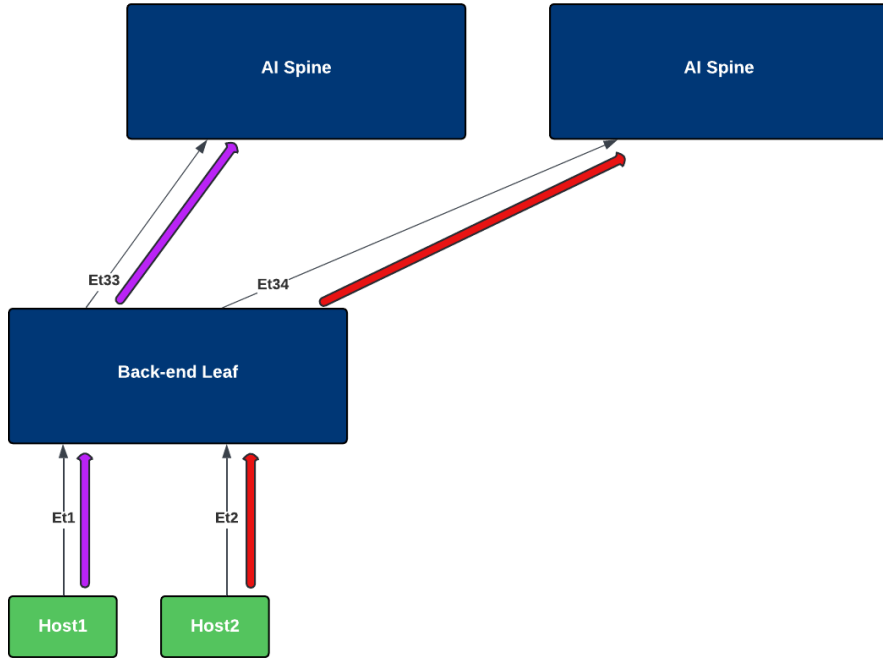
Load Balancing / Congestion Control

Equal Cost MultiPath (ECMP) / Fan-Out

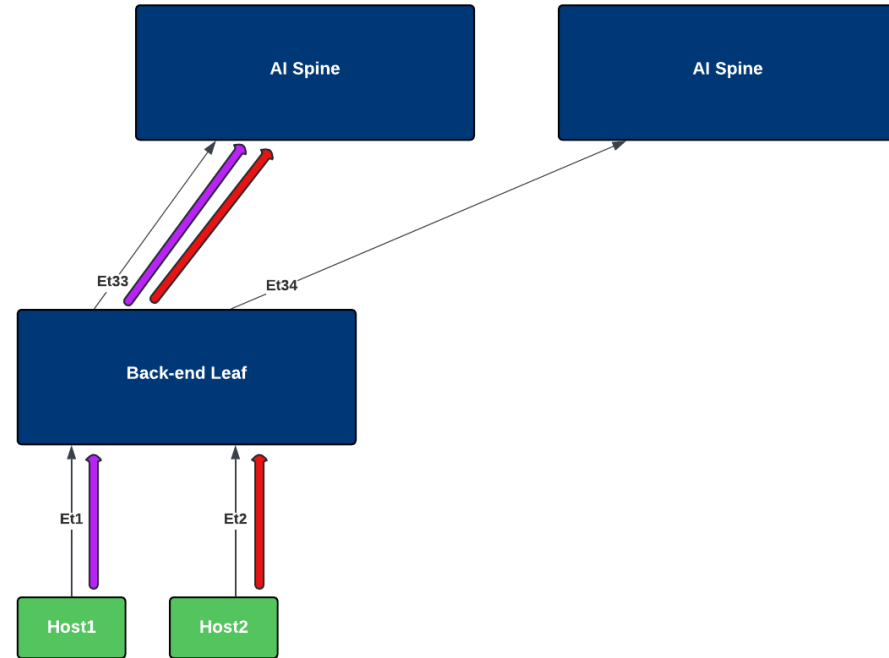


ECMP Hashing

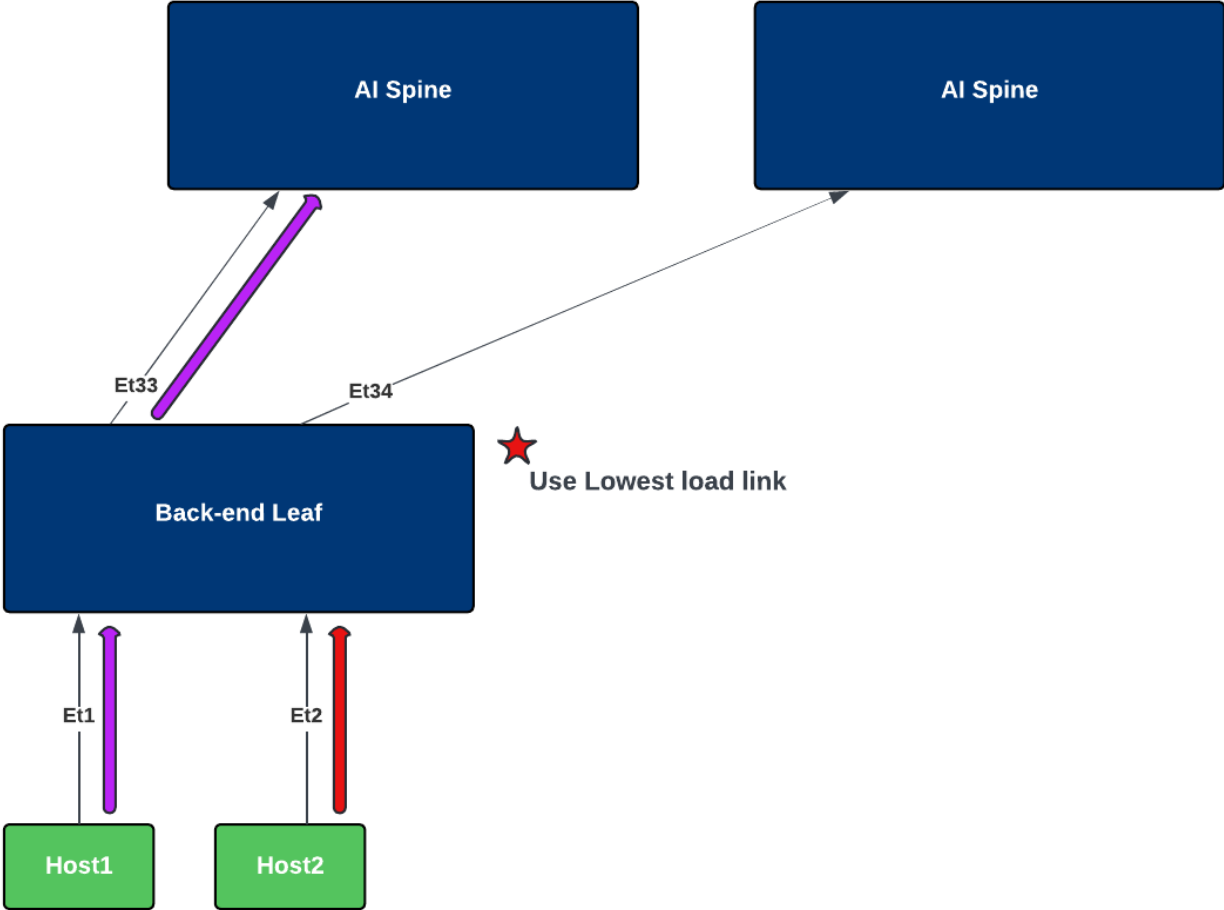
Good



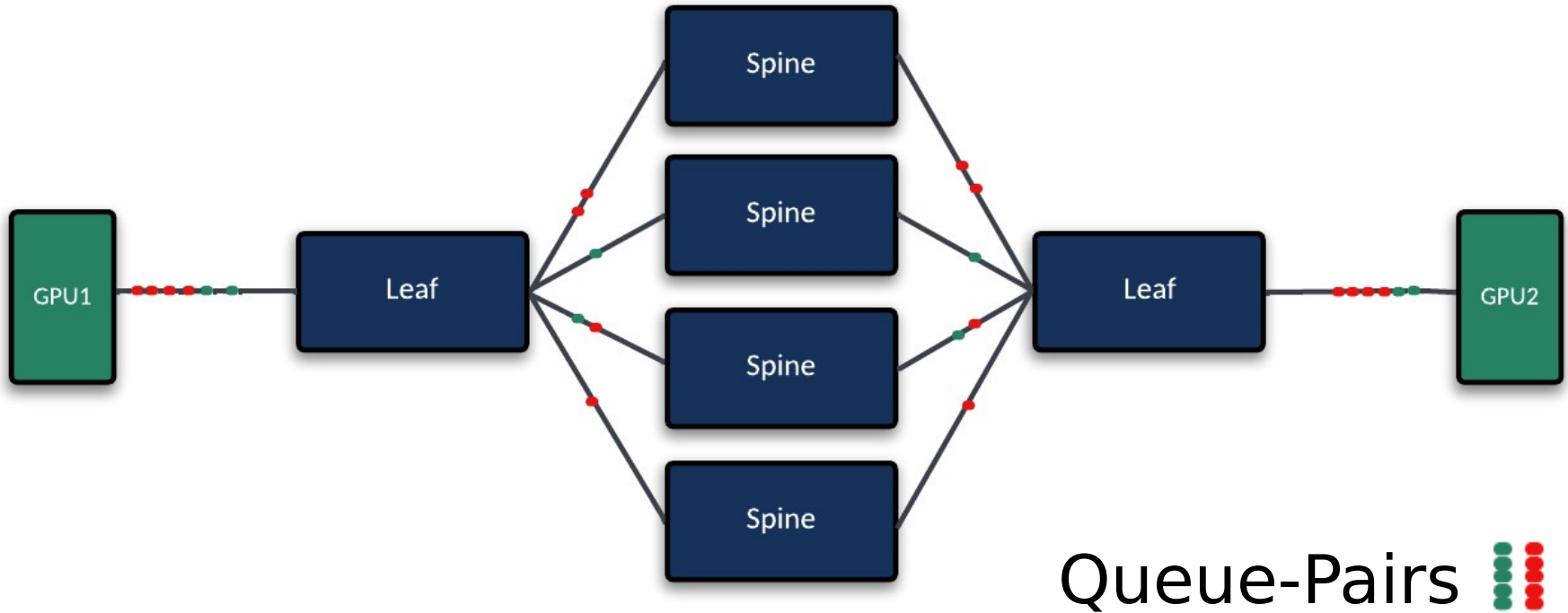
Bad



Dynamic Load-Balancing (DLB)

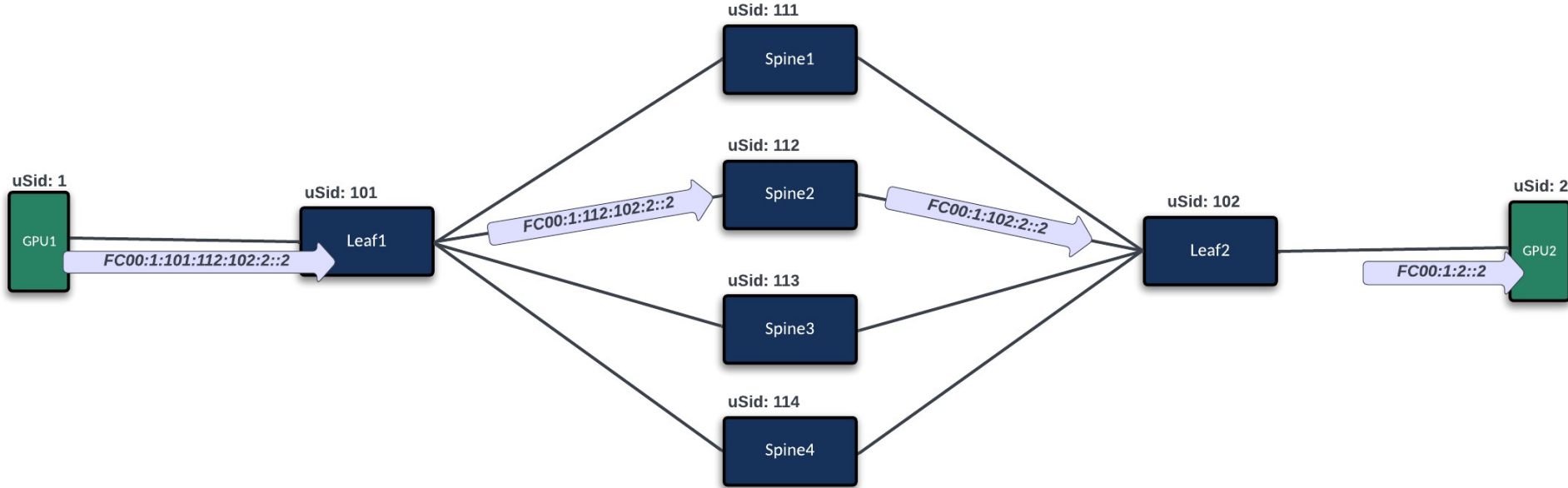


Packet Spraying



SRv6 Micro-SID

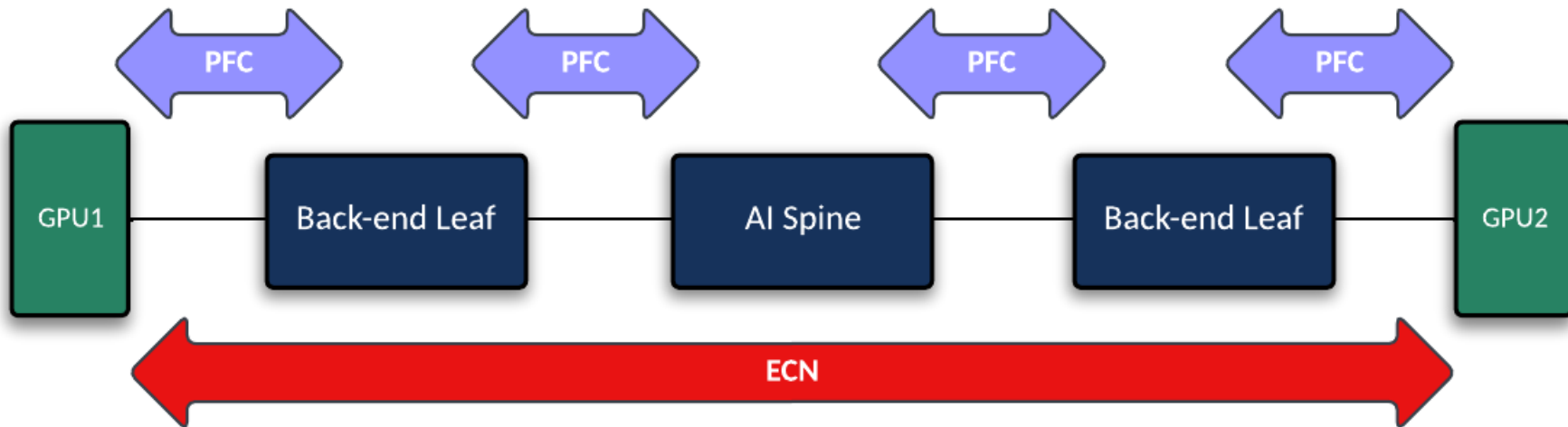
Load-Balancing & Segmentation



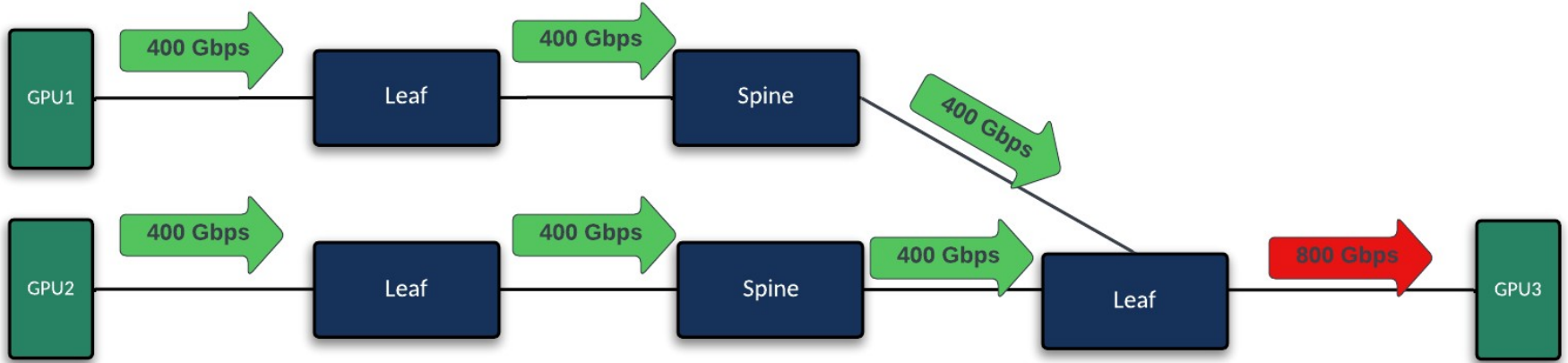
Source node determines network path. Usually requires a controller

Containerlab SRv6 Topology: https://github.com/brokenpackets/clab_Topos

RDMA / RoCEv2



Priority Flow Control



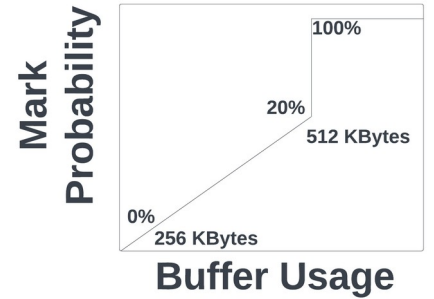
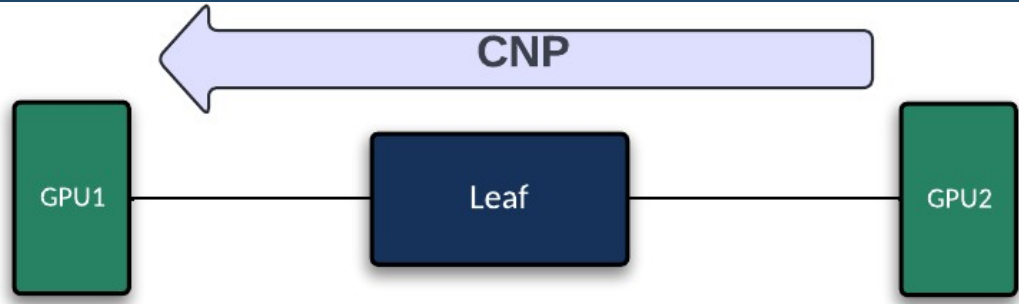
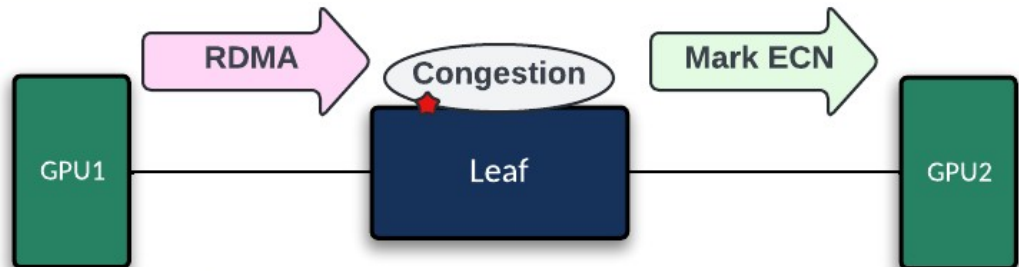
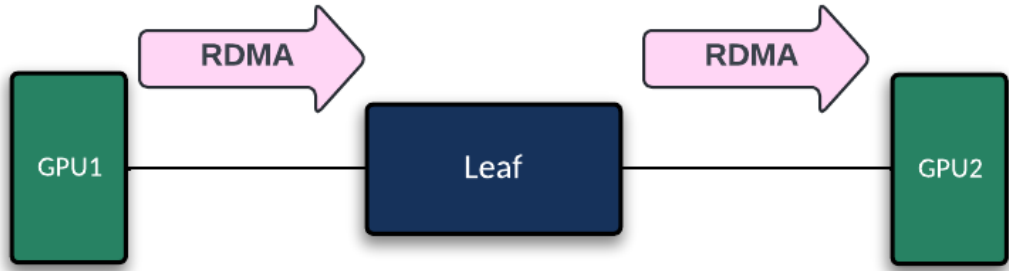
Ingress Queue



Egress Queue

Queue Backpressure in Egress causes Ingress to start building.

ECN/ECT



7	6	5	4	3	2	1	0
IP Precedence				Unused			
DiffServ Code Point (DSCP)				IP ECN			

Quality of Service

QoS Basics for AI



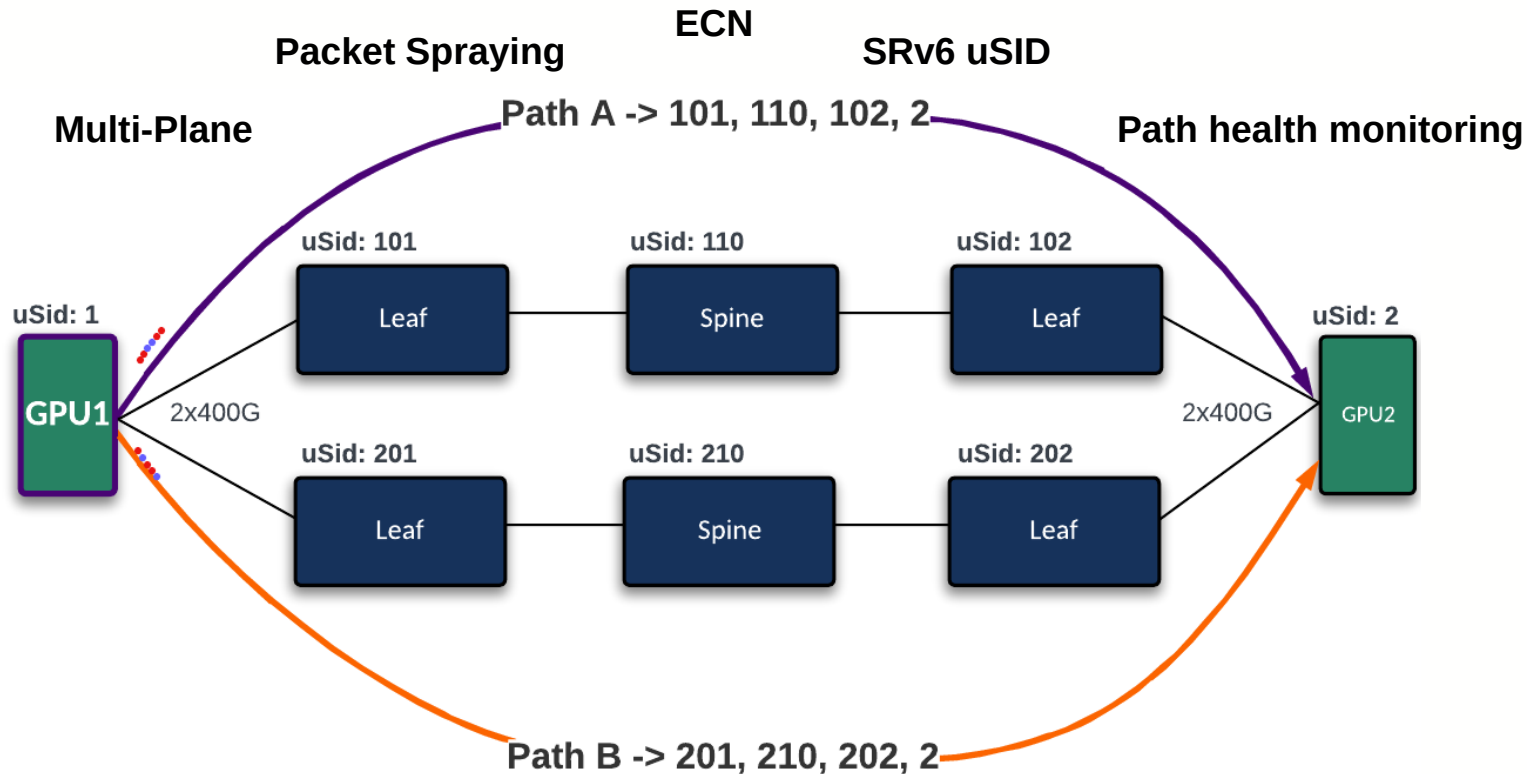
SP Queue 6/7 - CNP

SP Queue 4 - RDMA Control

WRR Lossless Queue 3 - RDMA Data

WRR Queue 1 - Everything Else

OCF MRC



No SRv6 Controller required. Just need initial topology information.

The Shift to AI Centers

Traditional Enterprise Datacenters focus on Services and Features.

In contrast, AI Centers prioritize Radix and Bandwidth above all else. Redundancy is usually handled at the application layer rather than with multihoming

- **Front-End:** Connecting CPUs and Storage to the Data Center and Users
- **Back-End:** High-radix fabric for GPU-to-GPU clusters.
- **Metrics:** Priority shifts from physical redundancy to subscription ratios.

If you are tasked with building an AI Center Network, plan for order-of-magnitude bandwidth increases, as well as a significant shift from the features that you likely leverage today.

ARISTA

THANK YOU