# AI Data Centers

Michal Styszynski  &  Mahesh Subramaniam
Juniper Networks

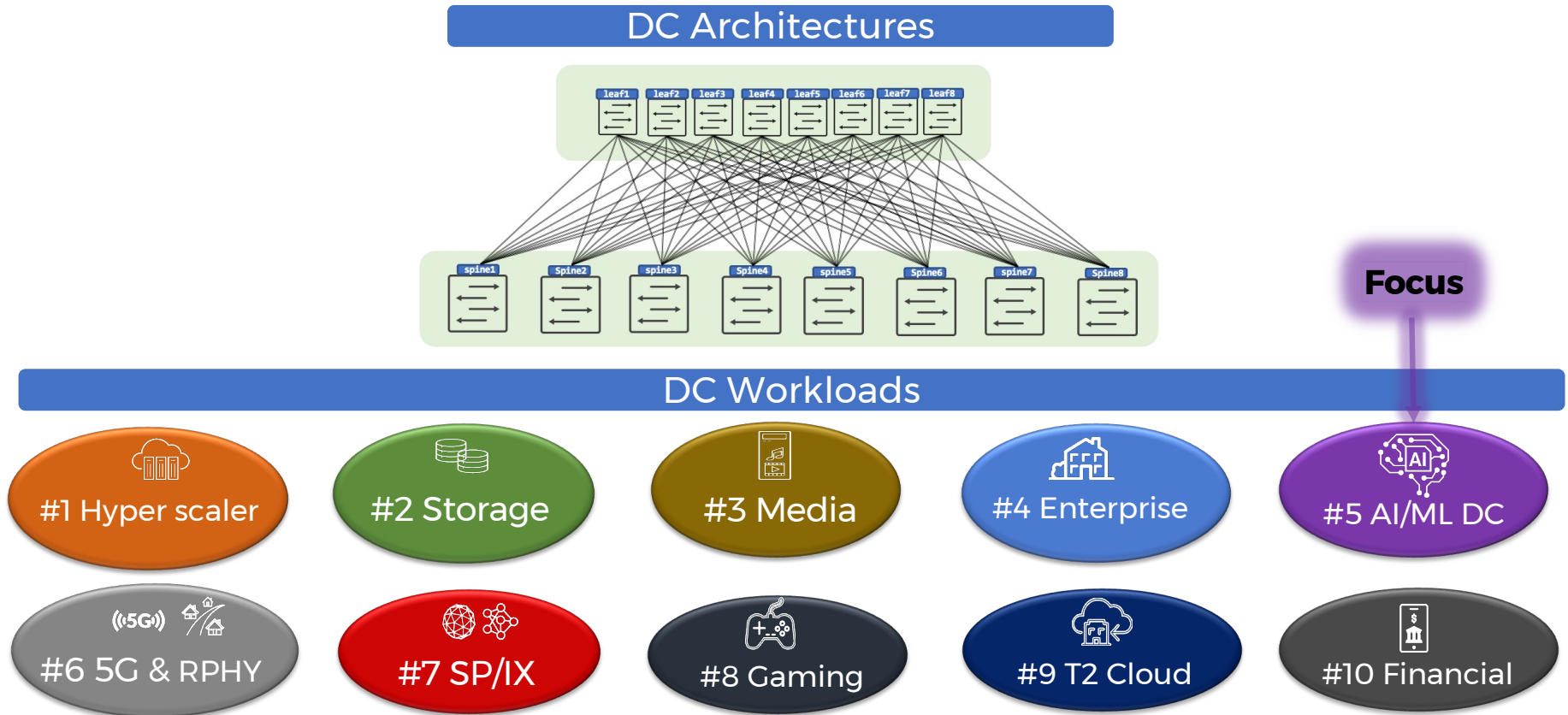ID 4950

February 13th, 2024

NANOG 90 – Charlotte, NC

# AGENDA

- DC architectures for Existing & Modern workloads
- Lifecycle of an AI DC Network
- AI DC technologies
- Key takeaways

NANOG™

# DC Arch. : Existing & Modern Workloads

# Why is AI DC now?

- Maturity of AI ML models development:
  - AIML models became more accurate, more fluent, and more creative
  - Availability of opensource AI models increased recently

- The increasing availability of data:
  - As the amount of data available to AI models grows, so does the ability of those models to learn and improve
  - The more data an AIML model must learn from, the better it will be at generating natural language responses
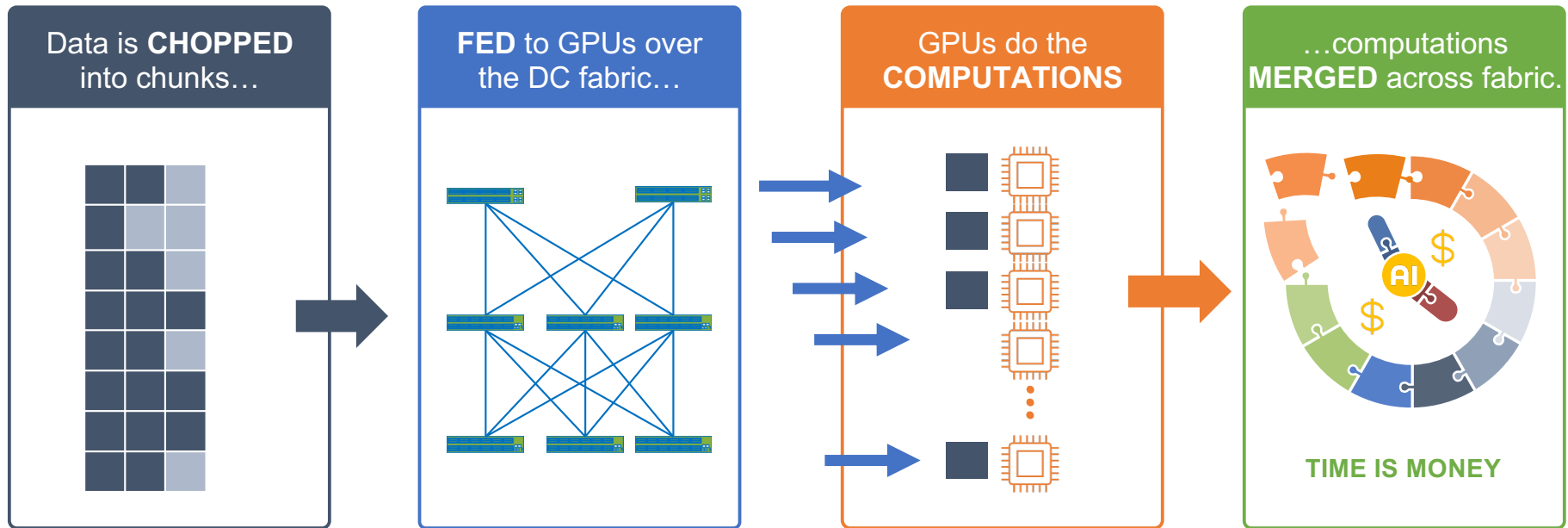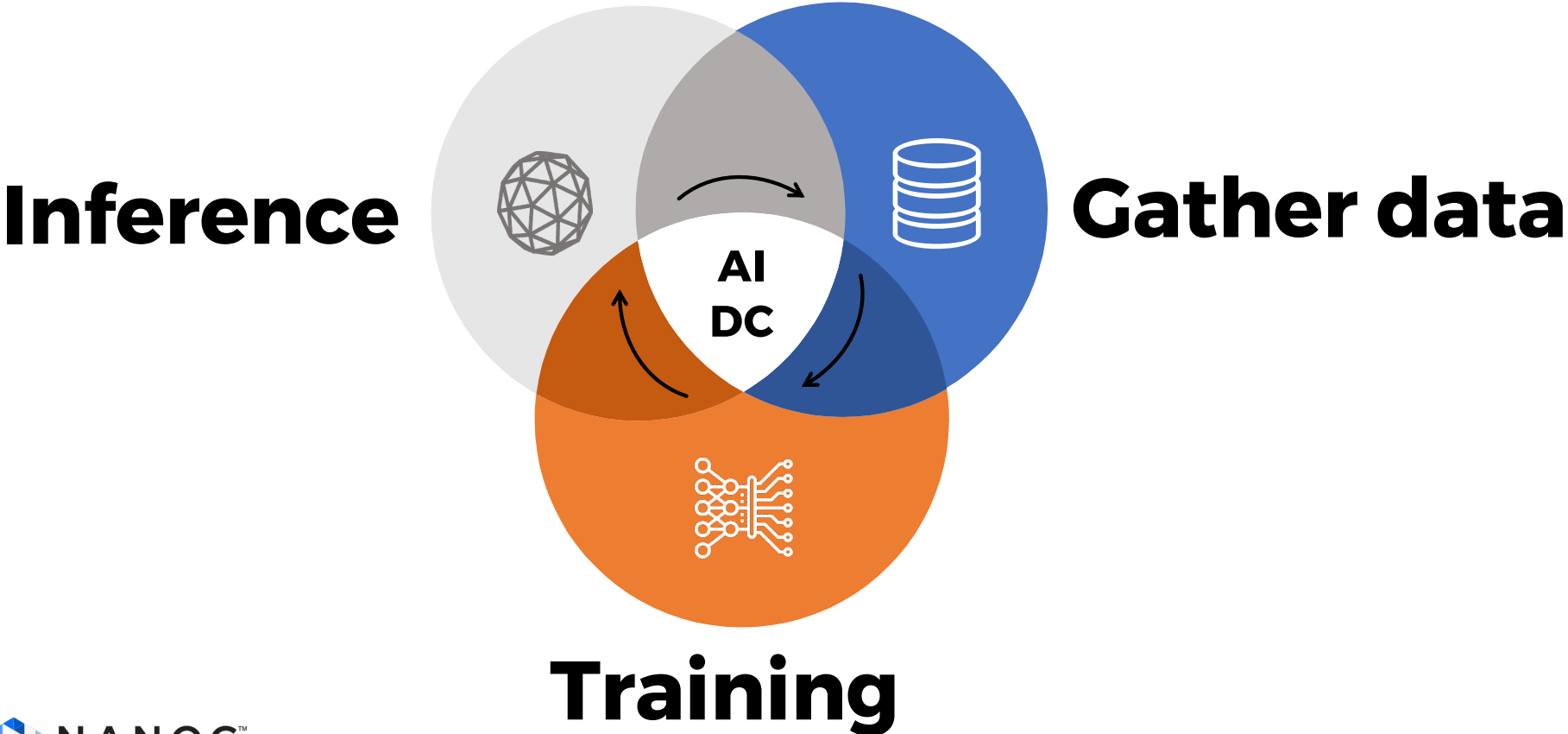
NANOG™

# Why is AI DC now?

- Technological advancements at the servers:
  - parallel processing of the data requirement
  - use GPU instead of serialized CPU processing

- Quick adoption of Generative AI applications by the users

**NANOG**™

# AI Model – Lifecycle

**Data is CHOPPED** into chunks…

**FED** to GPUs over the DC fabric…

GPUs do the **COMPUTATIONS**

…computations **MERGED** across fabric.

**TIME IS MONEY**

NANOG

# AI Model – Lifecycle

**Inference**

**Gather data**

AI
DC

**Training**
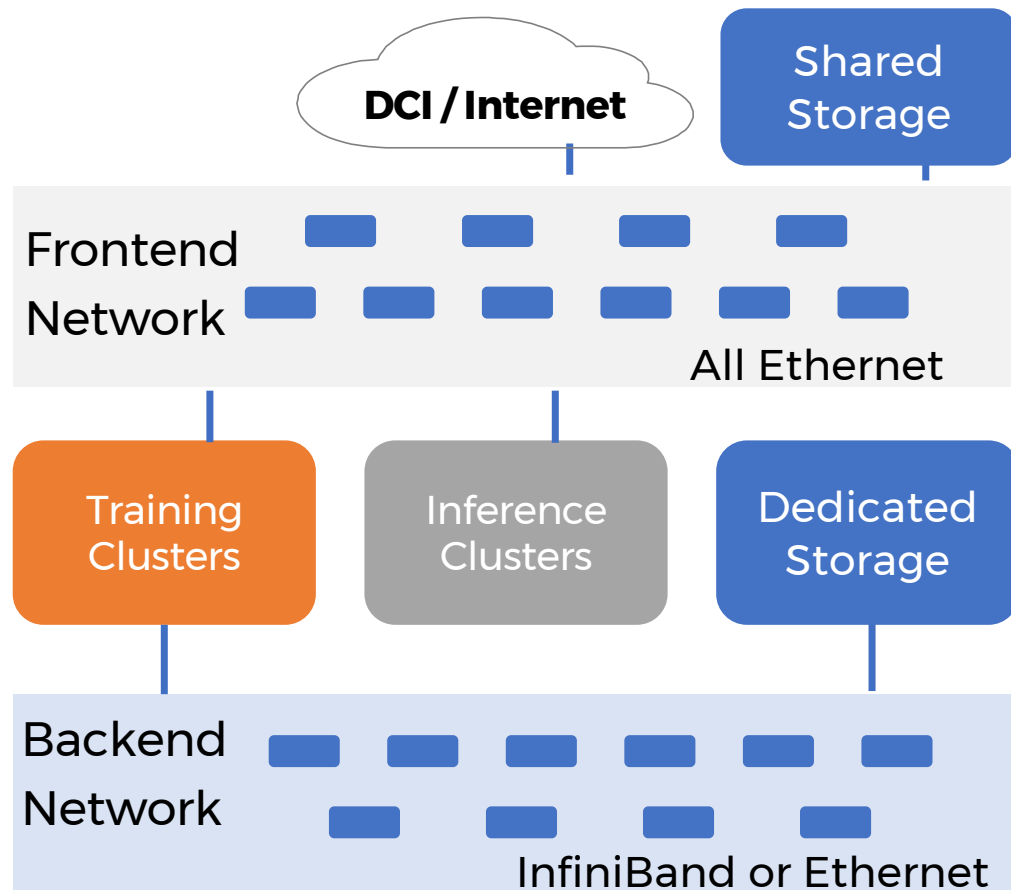
NANOG™

# Anatomy of an AI DC Network

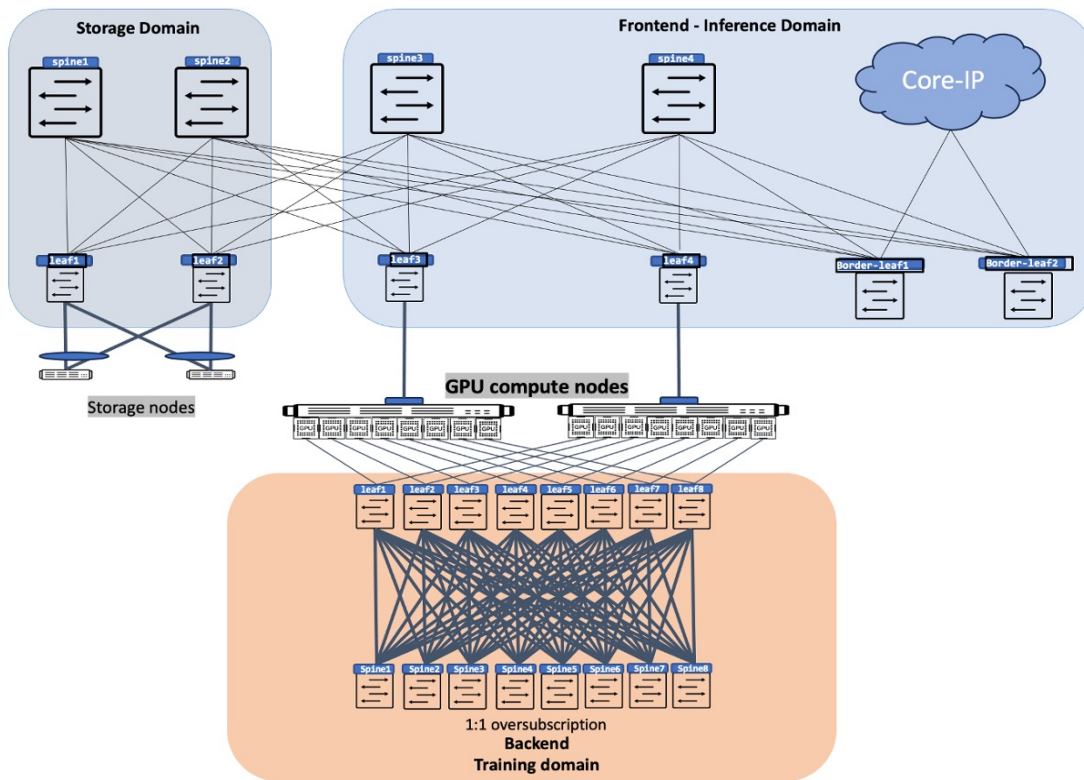## AI Cluster Networks

**"Frontend"**

- Inference clusters
- Shared storage pools
- Management

**"Backend"**

- GPU Compute Fabric
- Dedicated Storage Fabric

# AI DC - Architectures



**AI DC: Key capabilities:**

- Efficient Load Balancing
- ROCEv2 Transport
- Congestion Mgmt
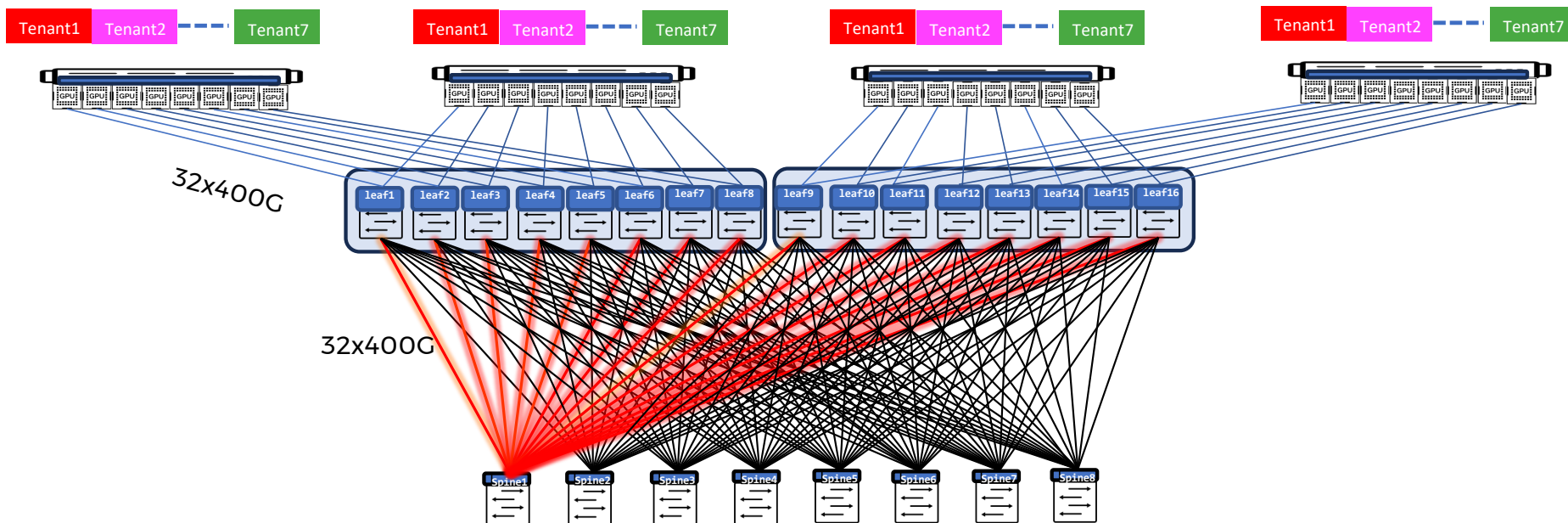- Adaptive IP Routing
- Monitoring

# AI DC – Stripe Optimized Design – SOD

GPU compute nodes rail optimized – Stripe 1
Each server with 8 x GPU

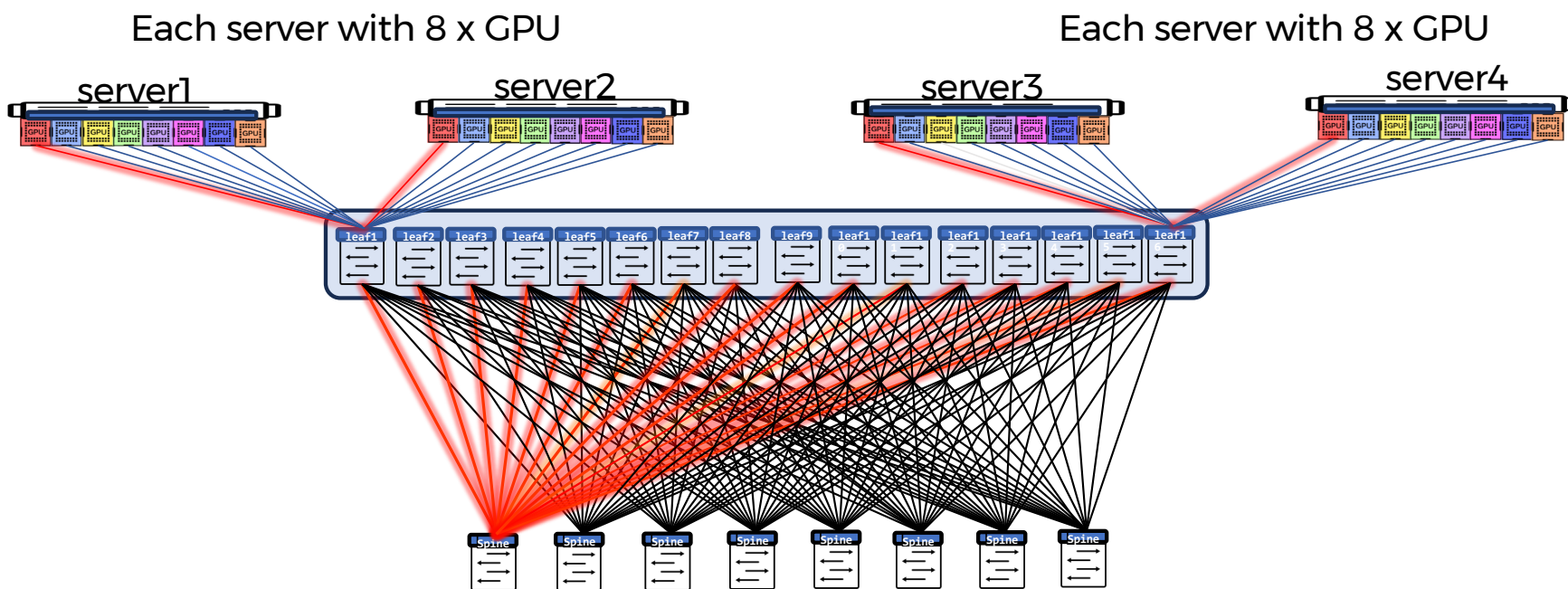GPU compute nodes rail optimized – Stripe 2
Each server with 8 x GPU

| Tenant1 | Tenant2 | – – – – | Tenant7 |

| Tenant1 | Tenant2 | – – – – | Tenant7 |

| Tenant1 | Tenant2 | – – – – | Tenant7 |

| Tenant1 | Tenant2 | – – – – | Tenant7 |

32x400G

leaf1 leaf2 leaf3 leaf4 leaf5 leaf6 leaf7 leaf8    leaf9 leaf10 leaf11 leaf12 leaf13 leaf14 leaf15 leaf16

32x400G

Spine1 Spine2 Spine3 Spine4 Spine5 Spine6 Spine7 Spine8

## Stripe Optimized Design - SOD
+ Multi-tenancy

NANOG

# AI DC – Stripe Unified Design - SUD

Each server with 8 x GPU

Each server with 8 x GPU

server1  server2  server3  server4

leaf1 leaf2 leaf3 leaf4 leaf5 leaf6 leaf7 leaf8 leaf9 leaf1 leaf1 leaf1 leaf1 leaf1 leaf1 leaf1

Spine Spine Spine Spine Spine Spine Spine Spine

**Stripe Unified Design - SUD**

NANOG

# AI DC: Requirements

**THIS SESSION FOCUS**

| | |
|---|---|
| ROCEv2 Transport | Congestion Control |
| Efficient load Balancing | IP routing in AI DC |
| Thermal Management | Liquid cooling |
| Optics | 400G – 800G |

NANOG

# RDMA Workload : AI DC

AI server-1

**RDMA Send**

**Initiate Transfer**

Memory

GPU-1..8

PCIE

NIC

Encap Data ROCEv2 Frame

ROCEv2 Transport

AI server-201

Memory

GPU-1..8

PCIE

NIC

**RDMA Receive**

Write directly to the GPU memory without going through the CPU of the server

400G/200G ToR switch

**Stripe-A Leaf-1**

**Stripe-B Leaf-2**

**Spines**

IP Fabric

NANOG™

# ROCEv2 - Transport for AI DC



Ethernet | IP | UDP | IB BTH | IB Payload | ICRC | FCS

Destination port 4791 indicates next header is RoCEv2

**Base Transport Header – BTH:**

Opcode: Transport type

Destination queue pair

Packet Sequence number
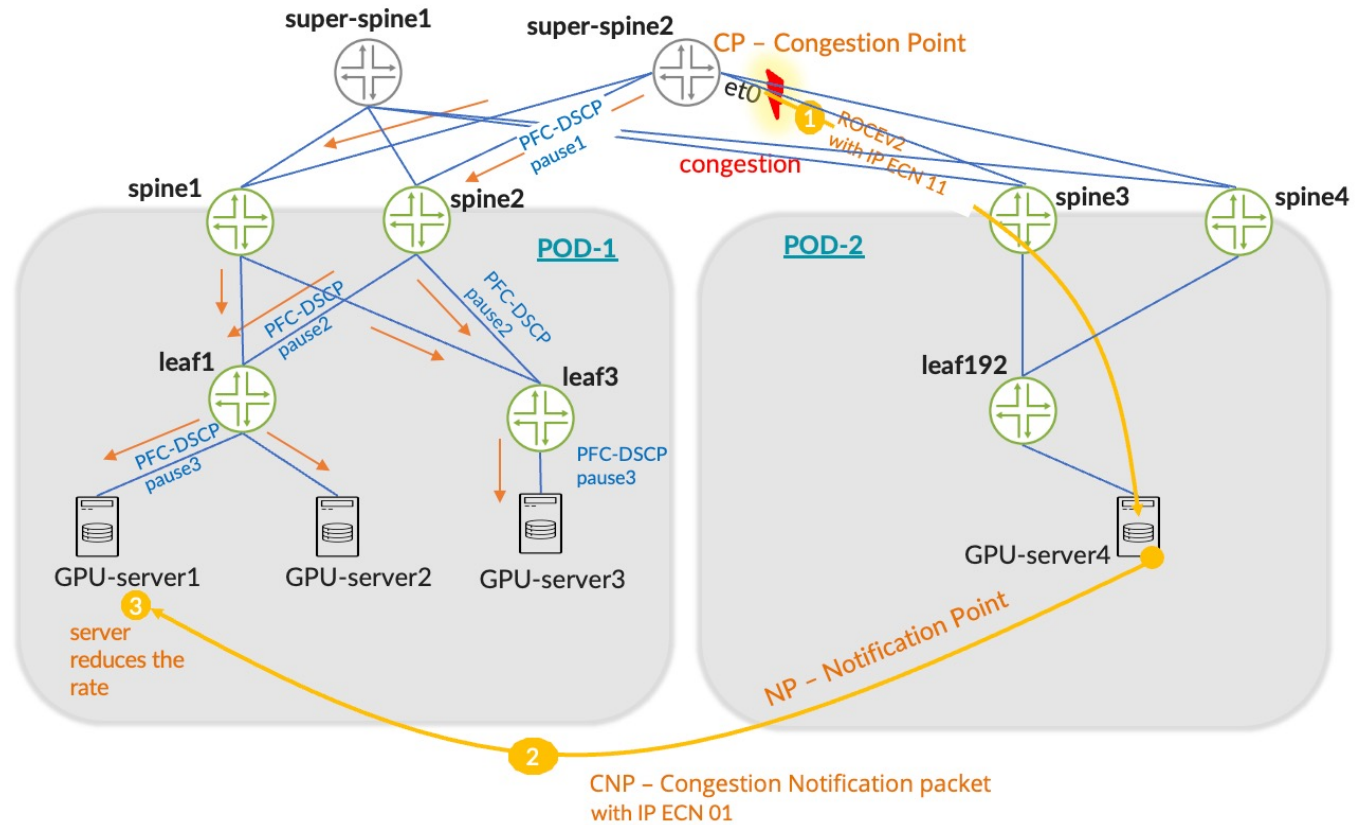
ROCE 32-bit end to end FCS= z_

NANOG™

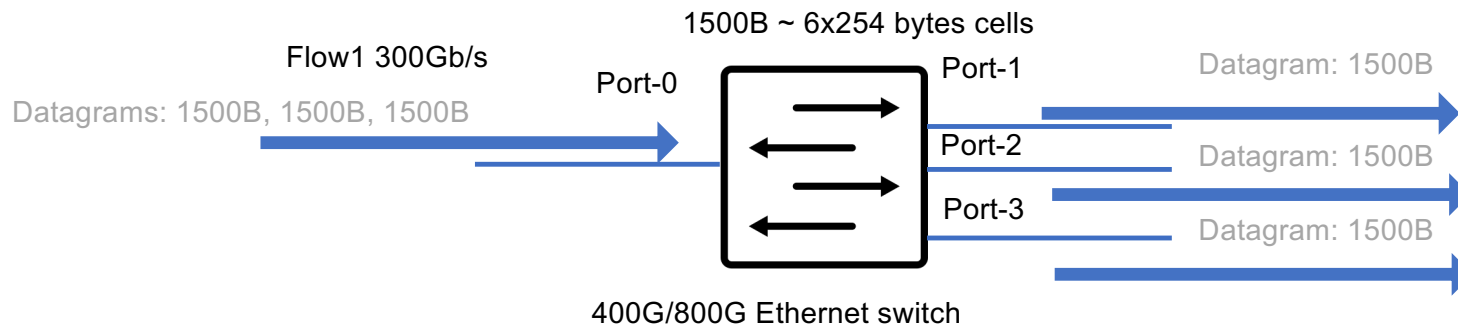# ROCEv2 – session establishment

# DCQCN – PFC-DSCP vs ECN

PFC-DSCP
Pause-level-1

↓

PFC-DSCP
Pause-level-2

↓

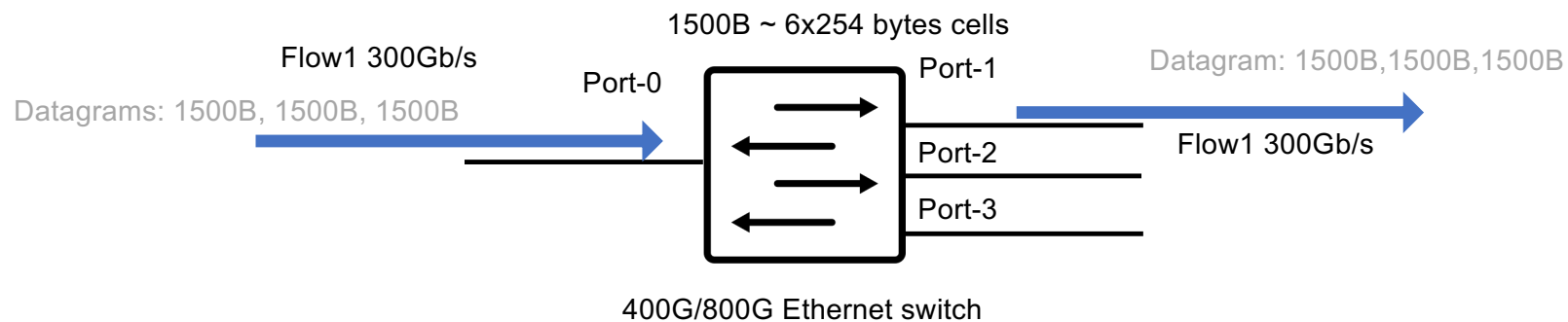PFC-DSCP
Pause-level-3



NANOG

# AI DC - Dynamic Load Balancing

## DLB (Dynamic Load Balancing) - per-packet optimal spraying



1500B ~ 6x254 bytes cells

Flow1 300Gb/s

Datagrams: 1500B, 1500B, 1500B

Port-0

Port-1    Datagram: 1500B

Port-2    Datagram: 1500B

Port-3    Datagram: 1500B

Datagram: 1500B

400G/800G Ethernet switch

Packet **re-ordering may happen** at the destination NIC Card connected to the GPU

NANOG

# AI DC - Dynamic Load Balancing

## DLB (Dynamic Load Balancing) – "flowlet" mode

1500B ~ 6x254 bytes cells

Flow1 300Gb/s

Datagrams: 1500B, 1500B, 1500B

Datagram: 1500B,1500B,1500B

Port-0

Port-1

Port-2

Port-3

Flow1 300Gb/s
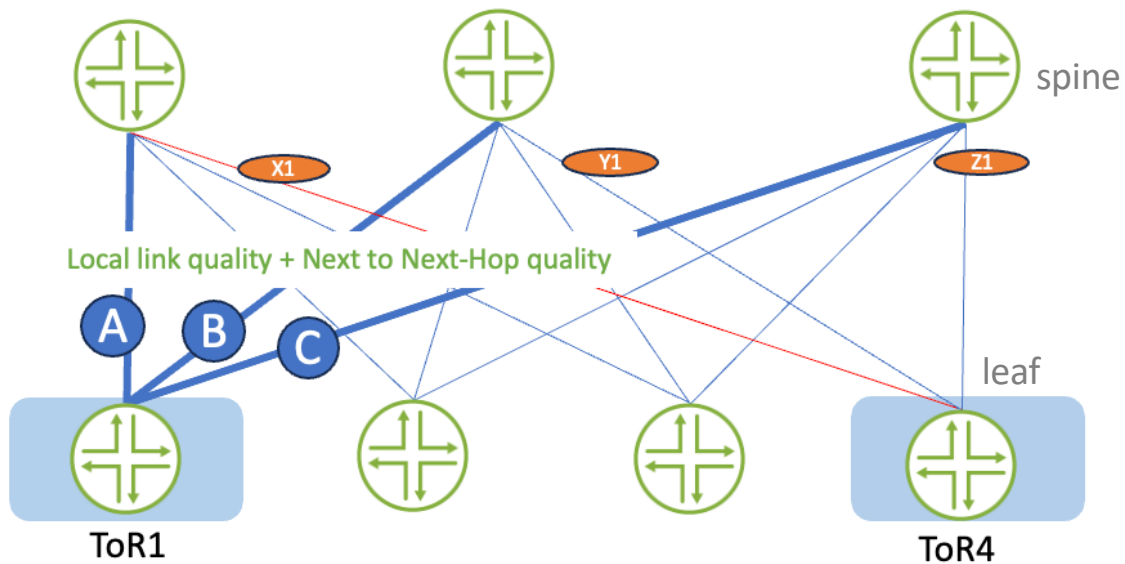
400G/800G Ethernet switch

Packet **re-ordering won't happen** at the destination NIC Card connected to the GPU

NANOG™

# AI DC – Selective Load Balancing



- Ability to selectively enable DLB via access lists for read/write operations

- It can handle out-of-order packets and enable DLB per packet mode for just that service
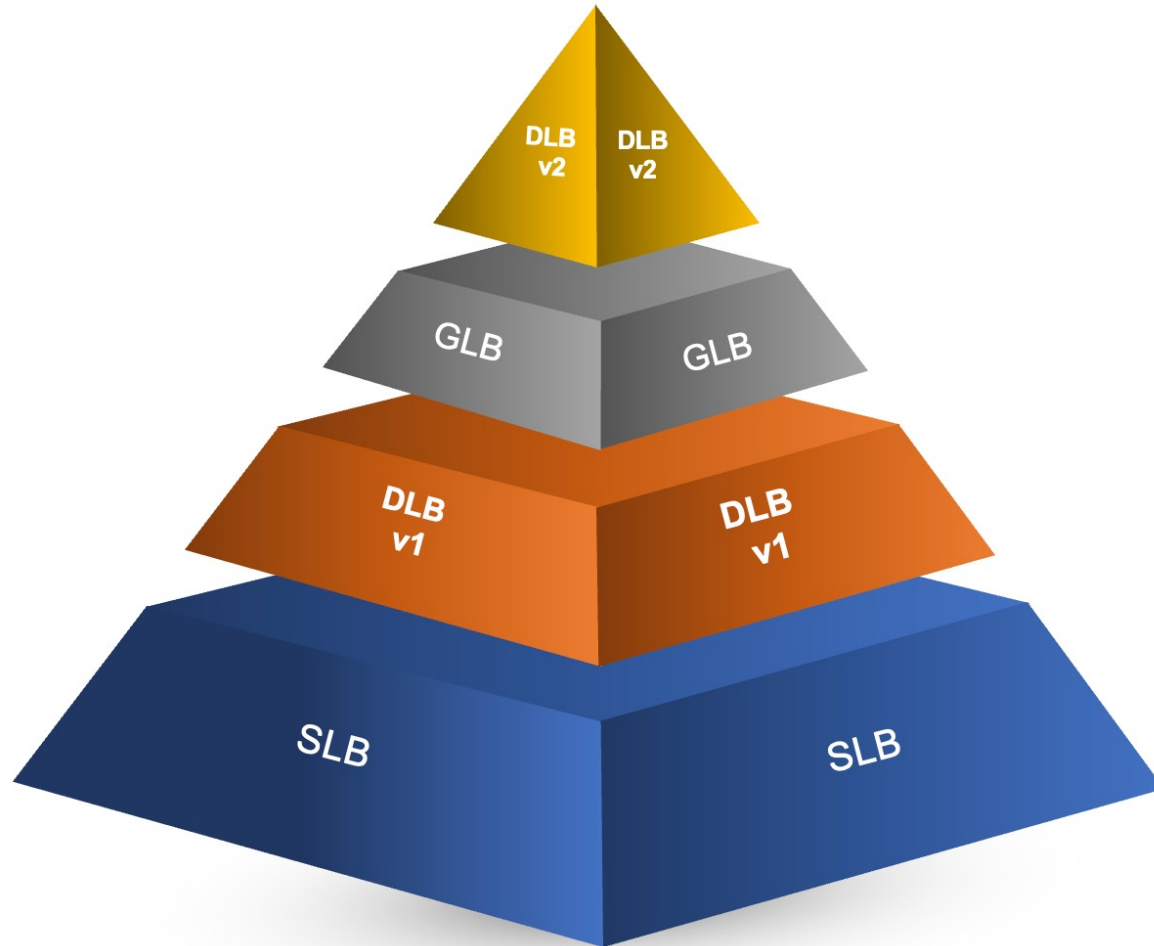
# AI DC – Global Load Balancing



Local link quality + Next to Next-Hop quality

- Global Load Balancing (GLB) uses **path quality**

- GLB selects a better end-to-end path.

| Route in HW at ToR1 | |
|---|---|
| NEXTHOP ID | Remote Quality Profile ID |
| ECMP_NH_ID1 | SwitchID.ToR4 |

| Local Port | Quality | | Quality | Remote Port |
|---|---|---|---|---|
| A | Q(A) | | Q(X1) | X1 |
| B | Q(B) | | Q(Y1) | Y1 |
| C | Q(C) | | Q(Z1) | Z1 |

# AI DC: Efficient load Balancing summary

# IP Routing for AI DC

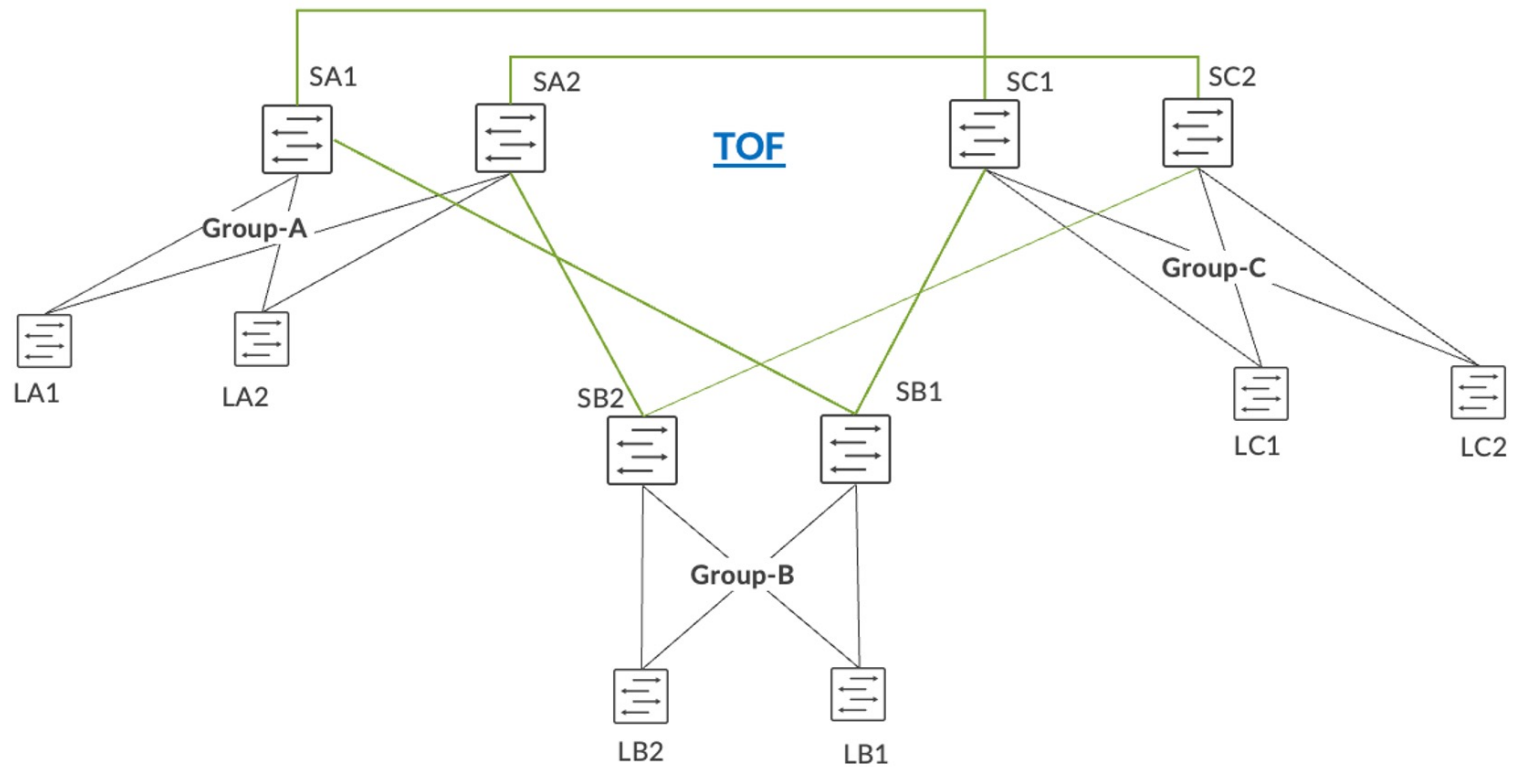| Frontend | eBGP underlay/overlay:<br>• underlay:  eBGP unnumbered / RFC5549<br>• overlay: EVPN-VXLAN |
| --- | --- |

| Backend | • BGP unnumbered/RFC5549<br>• IGP protocols: RIFT or ISIS |
| --- | --- |

NANOG™

# BGP unnumbered / RFC5549

# RIFT routing for backend network

# AI DC key takeaways

- The number of new AI applications is increasing over time.

- Dedicated AI DC infrastructures are built to accelerate parallel data processing.

- Ethernet 400G/800G adoption is increasing thanks to AI

- Congestion Management & Load Balancing efficiency are the key network components in AI DC

# Thank you

Feb 12-14, 2024

**NANOG**™