# Revisit RoCEv2 issues in large scale deployment and the future that UEC promise

AMD and Edgecore

# Agenda

**01**

**Problem Statement**

**02**

**Product**

**03**

**Solution**

**04**

**Performance**

**05**

**Q&A**

# AI Scale-out Networking Challenges

| Network Utilization | Reliability | Scalability | Operations | TCO |
|---|---|---|---|---|
| Inefficient GPU-to-GPU communication | Link, NIC and Switch failure | PFC & Queue Pair stalls Elephant flows sharing | Poor telemetry and lack of network state at CCL | Require deep buffer switches, lack of multi-plane/rail networks |

# RoCEv2 Requires Improvements for modern GenAI & HPC deployments

**PFC**
- PFC requires at least BW*RTT+MTU buffering for fully lossless transmission
- Blocked victim flows
- PFC storms

**Security**
- Flexibility for End-to-End confidentiality and service protection. Large session state (keys)

**Congestion Control**
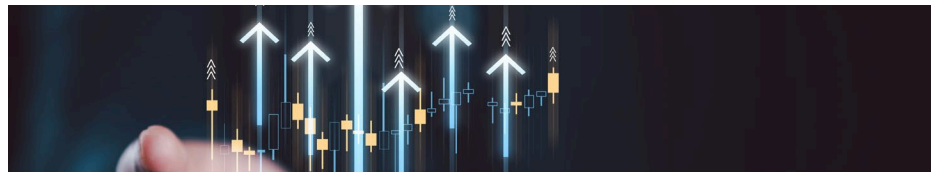- Different DCQCN implementations

**Link Level Reliability or Network Reliability**
- Delays become more significant as scale increases – Requires error handling at link layer

**Different traffics co-exists**
- RoCEv2 core design natively does not support different transport protocols for different services.

# Edgecore AIS800 Tomahawk 5 AI Switch

- 51.2Tbps while <1W per 100Gbps

- Best-in-Class SerDes that enable LPO

- (OSFP, QSFP) (AFO, AFI) complete portfolio

- Adaptive Routing & Cognitive Routing for all traffic types Improved Network Utilization ⇒ Lowest Tail Latency

- Programmable out-of-band telemetry (6 ARM cores) and Programmable inband telemetry ⇒ Minimized Packet Drops and Latency Jitter
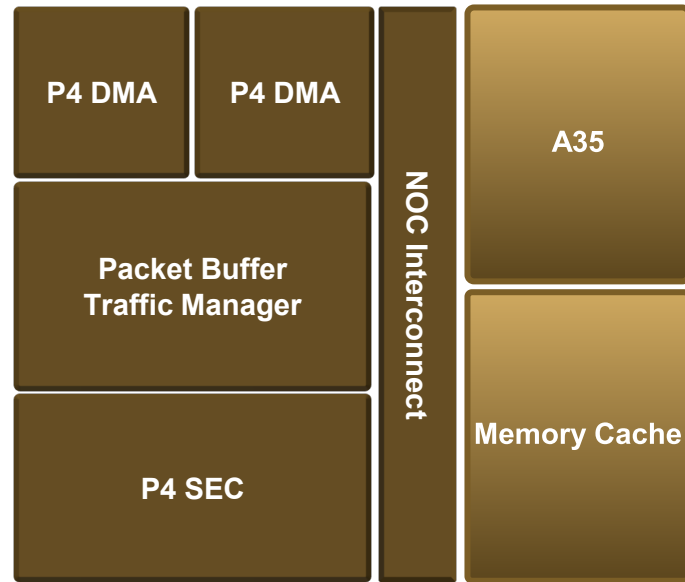
# AMD Pensando™ Pollara 400 AI NIC

- ❯ Fully Programable Customizable Transports

- ❯ Offload and Acceleration

- ❯ PCIe® Gen5, 400G

- ❯ Scale-Out Choice No Fabric Dependency

# AMD Pensando™ Pollara 400 AI NIC

- ▶ P4-based architecture - 72 MPU

- ▶ ATS and RDMA translation services to P4DMA

- ▶ High PPS / message rate and low latency RDMA services

- ▶ RDMA transport datapath with P4DMA Programmability

Pollara AI NIC RDMA Architecture

**High performance and Scale with the Flexibility of a FULLY P4 Programmable System**

# AMD AI Networking Solution

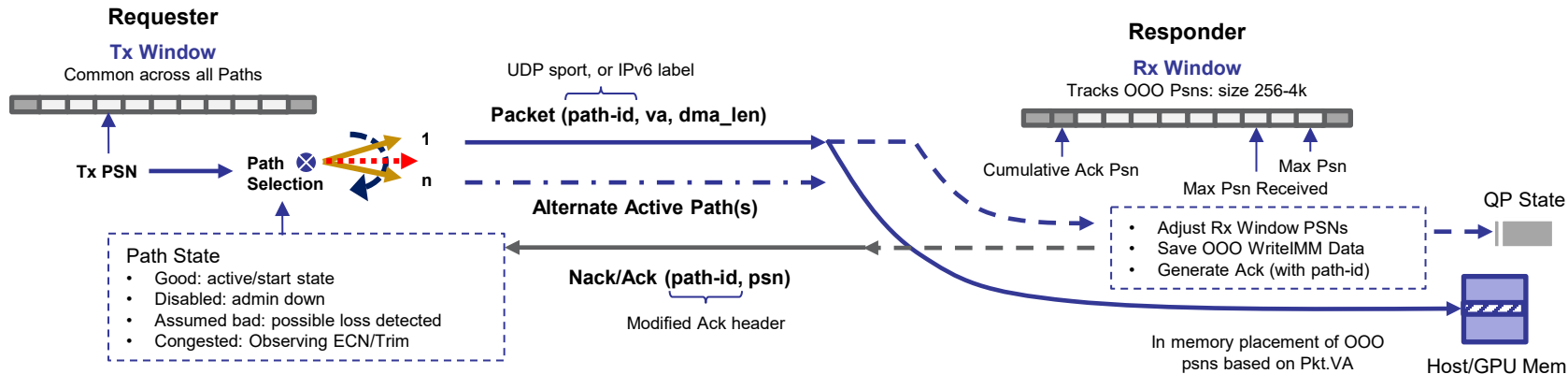| Network Utilization | Reliability | Scalability | Operations | TCO |
|---|---|---|---|---|
| Reliable Multi-path packet Spray, Out-of-order data placement, Flexible source routing | Fast data loss recovery with SACK and probes | Programmable Transport, Multi-path aware Congestion Management | Visibility into network paths and granular transport layer functions with extensible API. Easier to debug | No Fabric dependency, multi-plane network with fault isolation, redundancy and scale |

# Network Load-Balancing

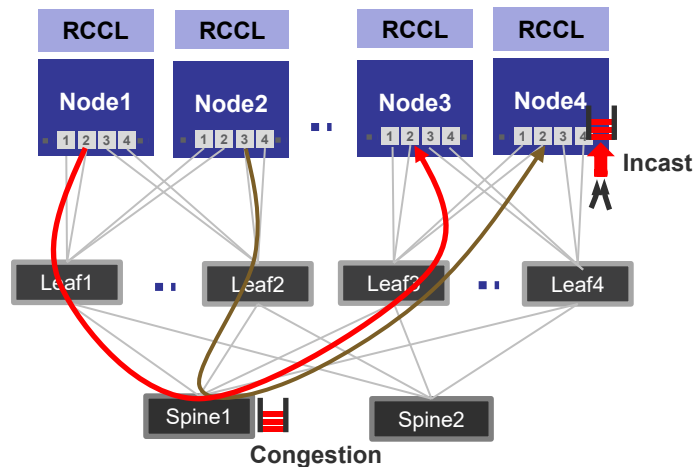| | |
|---|---|
| **Pollara supports multiple solutions for efficient network load-balancing for AI workload traffic** | **NIC** |

**NIC**
- **Multi-path Packet-Spray** for efficient load-balancing of individual RDMA Qpairs for ECMP-based Networks
- **Source-Routing** for traffic-engineering of RDMA traffic over multiple available network paths
  - Segment-routing / PBR based solutions with control-plane driven Entropy-value sets per RDMA Queue-pair

**SW**
- Network Switch driven dynamic-load-balancing (DLB) solutions with NIC Out-of-order data delivery
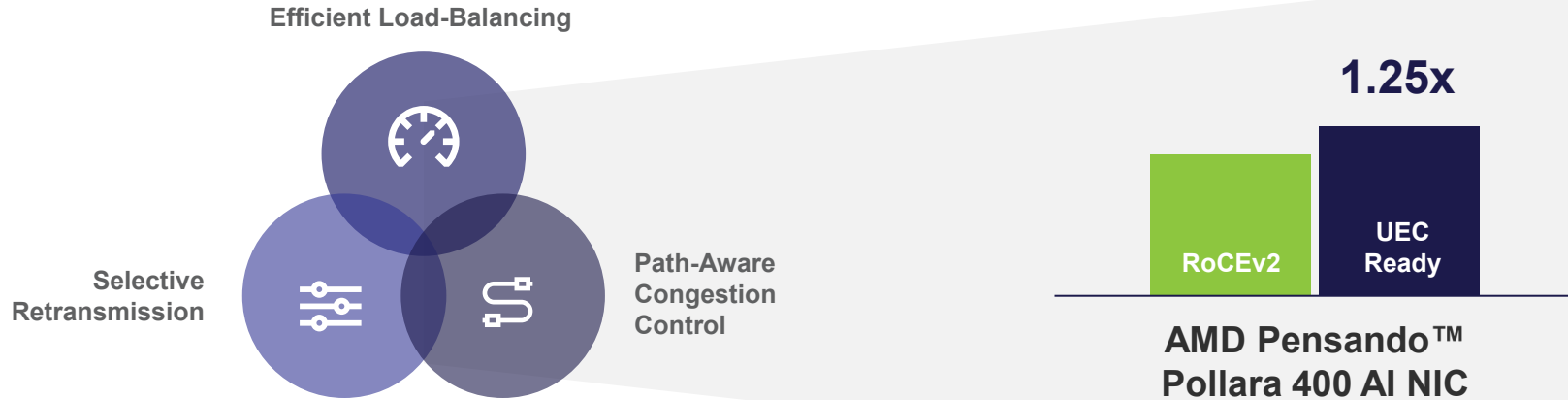
# Handling Out-of-order Data Delivery

AMD Pensando™ Pollara supports inline **Out-of-order data reception** (No NIC buffering) and delivery to handle packet re-ordering in network due to Packet Spray / DLB

- ❯ In-order completions of RDMA messages, ensuring packet re-ordering is transparent to AI workload applications

- ❯ Configurable Out-of-order Rx Window buffer for enhanced tolerance to re-ordering in network

- ❯ Lossy Network support UEC-NSCC based congestion control
  - No PFC requirement in the network
  - Window-based congestion-control algorithm with multiple congestion signals
  - Switch drop congestion notification based Fast-Retransmissions
  - Selective Acknowledgement and Selective Retransmissions for Fast data-loss recovery (SACK, SLEEK Algorithm)
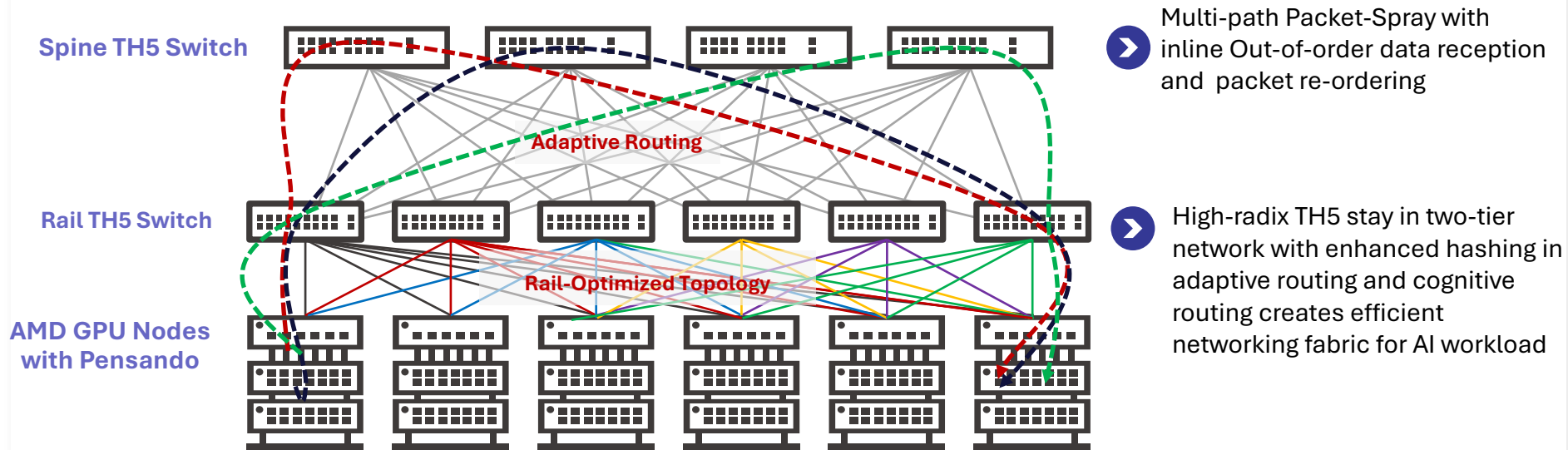
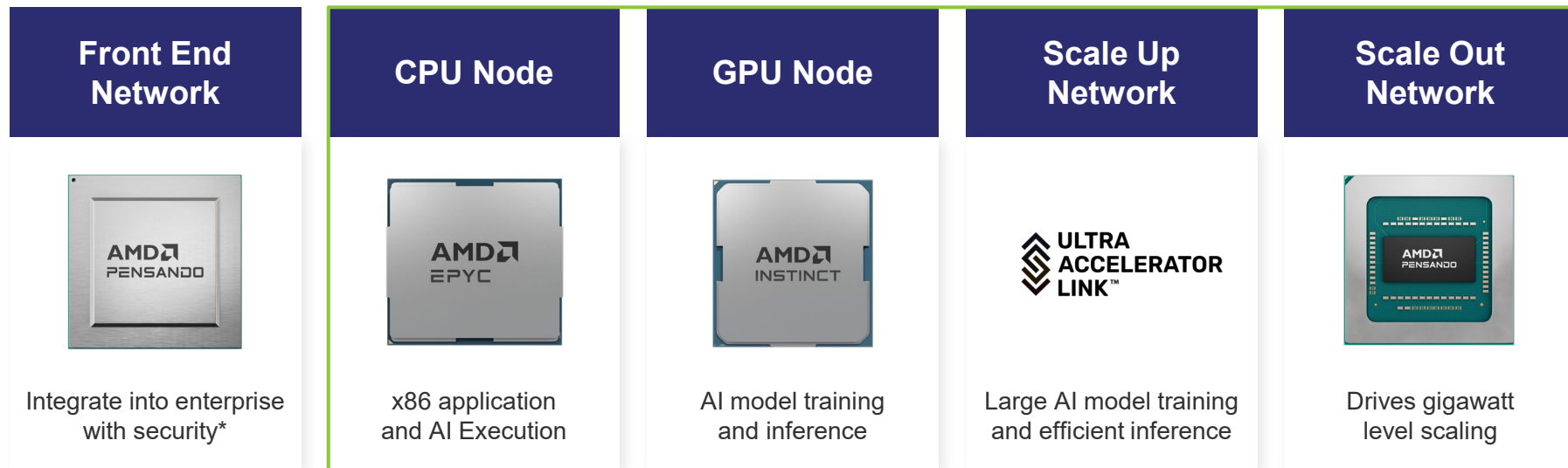# ~25% Higher Performance due to UEC Software differentiation

**Efficient Load-Balancing**

**Selective Retransmission**

**Path-Aware Congestion Control**

**1.25x**

RoCEv2

UEC Ready

**AMD Pensando™ Pollara 400 AI NIC**

# AMD Pensando™ Pollara NIC and Edgecore Tomahawk 5 AI Switch provides UEC-Ready Infrastructure that enable modern GenAI Deployment



Spine TH5 Switch

Rail TH5 Switch

AMD GPU Nodes with Pensando

Adaptive Routing

Rail-Optimized Topology

**Multi-path Packet-Spray with inline Out-of-order data reception and packet re-ordering**

**High-radix TH5 stay in two-tier network with enhanced hashing in adaptive routing and cognitive routing creates efficient networking fabric for AI workload**

# Complete AI System

| Front End Network | CPU Node | GPU Node | Scale Up Network | Scale Out Network |
|---|---|---|---|---|
| AMD PENSANDO | AMD EPYC | AMD INSTINCT | ULTRA ACCELERATOR LINK™ | AMD PENSANDO |
| Integrate into enterprise with security* | x86 application and AI Execution | AI model training and inference | Large AI model training and efficient inference | Drives gigawatt level scaling |

## ROCm™ Infrastructure Software

*No technology or product can be completely secure

# Call to Action

**Ultra Ethernet Consortium (UEC)**

**Validated Reference Guide**

**Where to find additional information (URL links)**
- https://www.amd.com/pensando

# Thank You!

**PoWen Tsai**

✉ powen_tsai@edge-core.com

🏠 http://www.edge-core.com

**Azeem Suleman**

✉ azeem.suleman@amd.com

🏠 http://www.amd.com

OCP APAC SUMMIT | 2025