**AMD** 

# AMD Instinct™ MI300 Series Cluster Reference Architecture Guide

# Contents

**AMD**

# List of Figures

# List of Tables

# Chapter 1: Abstract

Artificial Intelligence (AI) and Machine Learning (ML) models continue to advance in capability and scale at an increasingly rapid rate. This advancement increases performance and efficiency requirements on every element of AI/ML infrastructure, including networking infrastructure. As AI/ML model scale has increased, the workload must be distributed across Graphics Processing Units (GPUs) operating in parallel at massive scale. Performance is increasingly dependent on the network that enables data movement between these GPUs.

AI/ML model training and inference operations require the movement and processing of massive volumes of data. The GPU-GPU communication network must support a wide range of requirements such as latency-sensitive inference operations and iterative, high-throughput, parallel mathematical training operations. The highest throughput, lowest latency GPU-GPU data movement occurs within a node's *accelerator network* that connects a group of GPUs.

Additional GPU-GPU data movement occurs across a *backend network* that connects multiples of these nodes into a large-scale network cluster. AI/ML model deployments require a wide range of network scaling. Therefore, the design of the accelerator network and the backend network is crucial as it scales in size from small clusters of inferencing nodes to much larger-scale training backend networks supporting thousands of nodes and beyond. The efficiency and performance of GPU-GPU communication fabrics is therefore of critical importance.

The *frontend network* in existing data centers supports a wide range of functions including AI/ML data ingestion, storage as well as management functions. Traditionally, the front-end network is connected directly to CPUs in the nodes.

This reference architecture document outlines the components required to build a backend network cluster of *MI300 Series GPUs*, which uses primarily Ethernet-based network interface cards (NICs) and switches that scale to meet the increasing scaling requirements of AI models. AMD MI300 Series products support a wide range of networking technologies and topologies beyond Ethernet via standard PCIe-based NICs. AMD is committed to the development and enhancement of open standards-based networks like the Ultra Ethernet Consortium (UEC), and the Ultra Accelerator Link Consortium (UALink). AMD works with partners to support an open ecosystem of multiple networking solutions including AMD networking products. This reference architecture document describes a wide range of networking topologies including fat tree and rail-based topologies.

## Key Terms

The following table defines the key terms used in this document.

**Table 1.1:** Key Terms

| Terminology | Description |
|---|---|
| AMDMI300 Series | • AMD Instinct™ MI300X Platform<br>• AMD Instinct™ MI325X Platform |
| Backend Network | Network forming the cluster with GPU NICs, also referred to as the scale-out network, backend scale-out network, and backside scale-out network. The NICs in this network are referred to as backend NICs. |
| Accelerator Network | Network connecting GPUs within a node in a mesh with Infinity Fabric™ links, also referred to as the scale-up network, backend scale-up network, and backside scale-up network. NICs are not used in the MI300 Series. |
| Frontend Network | Network with a different set of NICs (from the backend network), also referred to as the frontside network. Depending on the server design, this network can also support storage and in-band management operations. The NICs in this network are referred to as frontend NICs. |

# Chapter 2: Components of an MI300 Series Cluster

A cluster consists of several components:

- AMD MI300 Series Platform compute nodes,
- Network fabrics that are composed of at least three networks with NICs, switches and cables, and
- Software libraries, system and management components.

## AMD MI300 Series Platform Compute Node

The AMD MI300 Series platform comprises eight OCP Accelerator Module (OAM) form-factor MI300 Series GPUs in a Universal Baseboard (UBB) 2.0 design. The following figure shows the air-cooled platform.

**Figure 2.1:** AMD MI300 Series Platform

A compute node consists of the AMD MI300 Series Platform together with CPUs, memory, and NIC devices. Specifications of the AMD MI300 Series reference compute node are given in the following table. Compute nodes with AMD MI300 Series platforms are available from select vendors (see Vendor List for Cluster Networking).

**Table 2.1:** Compute Node Reference Design Specifications

| Component | Specification |
|---|---|
| CPU | 2 x 4th-gen AMD EPYC Processors |
| GPU | 8 x AMD Instinct™ MI300 Series Accelerators with AMD Universal Base Board (UBB 2.0) |
| Memory | Configurable; typical designs use 6 TB (24 x 256 GB DRAM) DDR5 |
| Drives | NVMe SSDs; typical designs use 8-16 2.5-inch drives, 1-2 OS drives, high performance scratch drives |
| Networking | 8 x PCIe 5.0 high-performance networking cards, 400 Gb Ethernet |
| Accelerator Interconnect | Incorporating the AMD Infinity Architecture platform with 128 GB/s bidirectional Infinity Fabric™ bandwidth between each GPU for a peak aggregate bandwidth of 896 GB/s |
| Cooling | Air cooling or liquid cooling |

# Network Fabrics

Network fabrics are composed of at least three networks with NICs, switches and cables, as detailed in the following table. Several such components are listed in Vendor List for Cluster Networking.

**Table 2.2:** Network Hardware Components

| Hardware Component | Description |
|---|---|
| Backend scale-out Network | Fat-tree, or rail-optimized cluster topology with RDMA optimized Ethernet NICs and switches |
| Accelerator Network | Infinity Fabric™ mesh interconnecting 8 GPUs in the Compute Node |
| Storage Network (Frontend network) | Storage network topology connected through frontend NICs |

**Table 2.2:** Network Hardware Components (continued)

| Hardware Component | Description |
|---|---|
| In-band management network (Frontend network) | Management network connected through frontend NICs. Also provides services accessible by users. |
| Out-of-band management network | Separate network with its own NICs connecting BMC |

# Software Components

An MI300 Series Cluster requires the following software components.

**Table 2.3:** Software Components

| Software Component | Description |
|---|---|
| Data Center Management Software (RDC, SMI) | ROCm Data Center Tools and System Management Interface Libraries (see Cluster Management with RDC and SMI Tools) |
| System Management | Software and user interface for system management of nodes |

# Chapter 3: Design Requirements

AI/ML deployments have a wide range of cluster network scale requirements. The optimal system design should consider the node design, target NIC cards, switch capabilities, and target workloads to deliver required efficiency and performance.

This reference architecture provides a starting point with common usage models for AI/ML or High Performance Computing (HPC) workloads.

## System Design

The following hardware system design components are recommended:

- Scalable cluster architecture based on a scalable unit of 32 compute nodes
- Datacenter racks may have 1, 4, 8 or 16 compute nodes. The specific number of nodes is influenced by rack power and cooling requirements.
- Supported Networking adapters (NICs) and switches from AMD and partners supporting up to 400 Gb/s
- Storage Networking components to support storage servers and storage network

The following software system design components are recommended:

- Cluster management software from AMD and partners
- System management software from AMD and partners

## Compute Node

The compute node consists of eight MI300 Series GPUs interconnected by 4th gen AMD Infinity Fabric™ Links. A typical compute node also includes dual-socket CPUs, memory, and NICs connected via two PCIe 5.0 switches.

Performance-optimized designs have a specific mapping of MI300 Series GPU to frontend NICs and backend NICs as illustrated in Figure 3.1. Each CPU has one directly connected frontend NIC, so there are a total of two frontend NICs per compute node. Each MI300 Series GPU has one backend NIC that is connected through the PCIe switch, so there are a total of eight backend NICs per compute node.

## Frontend, Accelerator and Backend Networks

Cluster network fabric is composed of at least three networks, as illustrated in the following figure and discussed in the following subsections.

**Figure 3.1:** Frontend, Accelerator and Backend Networks



# Frontend Network

The frontend network is the traditional datacenter network comprised of switches and network adapters (or NICs) which support storage and management functions.

Storage network (part of the frontend network):

* As language models grow, it is important to consider fault tolerance for hardware failures and software errors. Creating and storing checkpoints are essential for recovery and reducing JCT (job completion time). The storage fabric connects compute and storage servers. .

* Independent of the backend network. RoCE is a prerequisite (RDMA over Converged Ethernet).

* If a separate network is designed, a different plane is ideal for storage needs.

In-band management network (part of frontend network):

* The in-band management fabric is used for node provisioning, data movements, SLURM, Kubernetes, and downloading from package repos such as pypi, docker repo, gcr.io, etc.

* Ethernet based, is used for provisioning of nodes and services that need to be accessed by users.

* If a separate network is designed by vendors, a 100 Gbps network is desirable.

# Accelerator Network

The accelerator network is a high bandwidth, low latency network that connects a group of GPUs and supports load/store transactions between the GPUs, as shown in the following figure. In MI300 Series based designs, this network is Infinity Fabric™ interconnecting 8 GPUs in a mesh topology within a compute node.

**Figure 3.2:** Accelerator Mesh Network with MI300 Series Accelerators

## Backend Network

The backend network connects a bigger set of GPUs (that are beyond the set available within the accelerator or scale-up network). Each compute node has eight NICs with a 1:1 GPU:NIC ratio, utilizing a PCIe switch between the GPUs and NICs. RDMA communication using RoCE protocol, congestion management and the support for UEC defined transport layer improvements are essential in this network. Communication between GPUs in this network is enhanced by NICs supporting acceleration of collective operations.

This network is designed to be highly scalable:

- The backend network topology can be either a fat tree or rail optimized.
- Ethernet switches form the 2- (leaf-spine) or 3-tier (leaf-spine-core) switching fabric.

NICs and switches:

- NICs connect compute nodes to the backend network through a set of switches in a well-defined topology.
- The number of required switches depends on the number of nodes in the cluster, the switch radix, and the cluster performance requirements.
- NICs and leaf (ToR) switches reside with the compute nodes in a rack.

## Out-of-band Management Network

The out-of-band management fabric is a separate, slow-speed (usually 1-10 Gbps) network that connects to the management ports of all nodes, storage servers, racks and switches in the whole cluster. Within a node, it connects to the Baseboard Management Controller (BMC), which is used to change BIOS settings, monitor and set the node health such as fan speed, voltage levels, temperatures, etc. Users can interact with the BMC through IPMI or Redfish API, or through the BMC web portal.

## System Management

For management and maintenance of a server, system vendors provide management software and interfaces that perform real-time health monitoring and management on each server, including firmware updates.

# Chapter 4: Cluster Architecture

A cluster consists of a group of racks, each of which consists of a group of servers. These servers are placed in a rack with backend switches for the backend network. In MI300 Series systems these racks are designed with qualified vendors (see Vendor List for Cluster Networking). The rack layouts are scalable and adjustable to meet the data center requirements. The following figure is a reference rack layout consisting of either 4, 8 or 16 GPU server nodes, each with 8 MI300 Series GPUs.

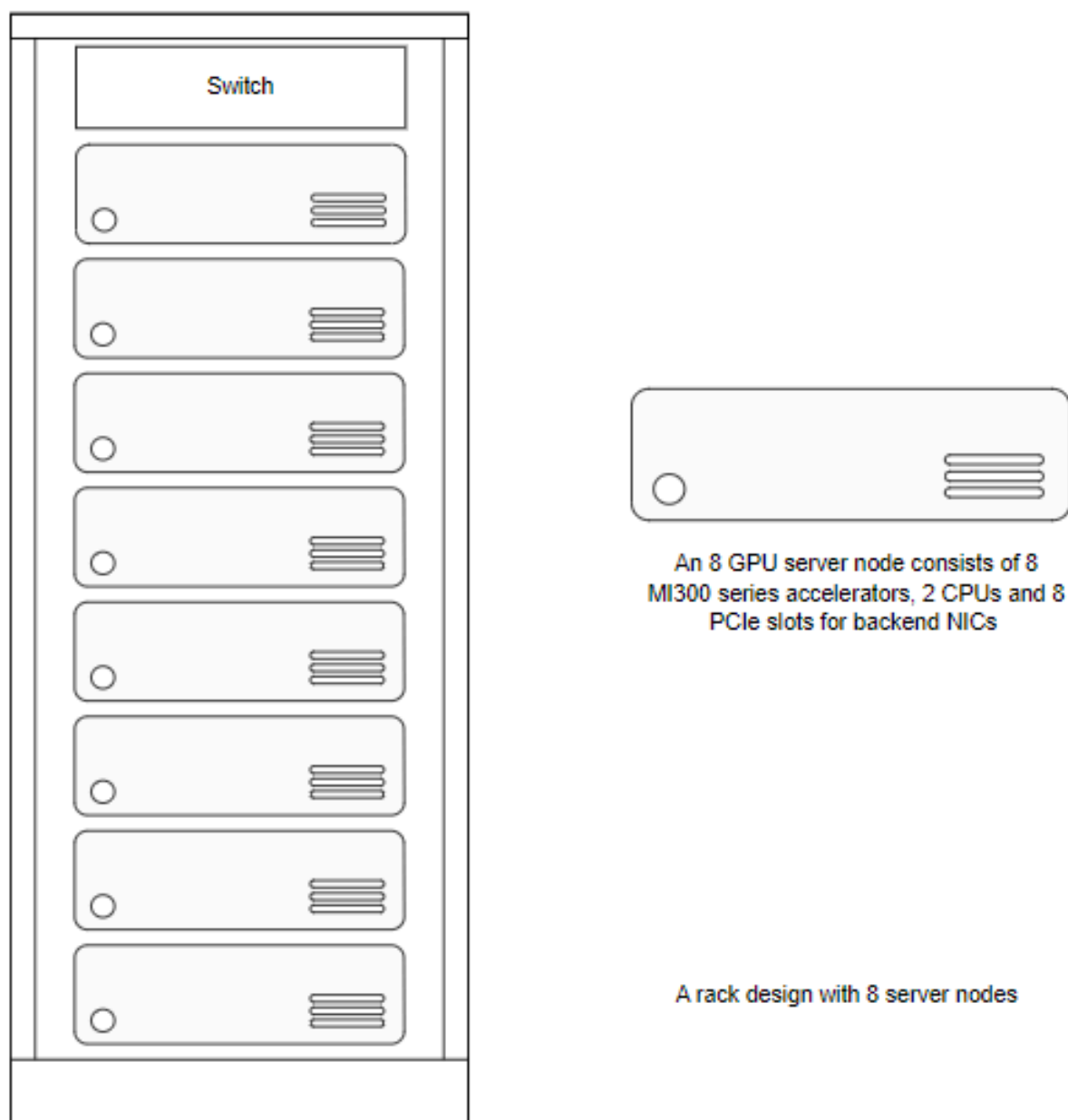**Figure 4.1:** Server and Rack Design with MI300 Series Accelerators



An 8 GPU server node consists of 8 MI300 series accelerators, 2 CPUs and 8 PCIe slots for backend NICs

A rack design with 8 server nodes

**Table 4.1:** Component Count with 64 x 400G Switch

| Node Count | GPU Count | Switch Count | | | | Cable Count | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Leaf | Spine | Core | Total switches | Nodes-Leafs | Leafs-Spines | Spines-Cores | Total |
| 128 | 1024 | 32 | 16 | – | 48 | 1024 | 1024 | – | 2048 |
| 256 | 2048 | 64 | 32 | – | 96 | 2048 | 2048 | – | 4096 |
| 512 | 4096 | 128 | 128 | 64 | 320 | 4096 | 4096 | 4096 | 12288 |
| 1024 | 8192 | 256 | 256 | 128 | 640 | 8192 | 8192 | 8192 | 24576 |
| 2048 | 16384 | 512 | 512 | 256 | 1280 | 16384 | 16384 | 16384 | 49152 |

**Table 4.2:** Component Count with 128 x 400G Switch

| Node Count | GPU Count | Switch Count | | | | Cable Count | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Leaf | Spine | Core | Total switches | Nodes-Leafs | Leafs-Spines | Spines-Cores | Total |
| 128 | 1024 | 16 | 8 | – | 24 | 1024 | 1024 | – | 2048 |
| 256 | 2048 | 32 | 16 | – | 48 | 2048 | 2048 | – | 4096 |
| 512 | 4096 | 64 | 32 | – | 96 | 4096 | 4096 | – | 8192 |
| 1024 | 8192 | 128 | 64 | – | 192 | 8192 | 8192 | – | 16384 |
| 2048 | 16384 | 256 | 256 | 128 | 640 | 16384 | 16384 | 16384 | 49152 |

# Chapter 5: Topology of Network Fabrics

A network fabric in a cluster design consists of the following fabrics:

- Compute fabric (backend network),
- Storage fabric (frontend network),
- In-band management fabric (frontend network), and
- Out-of-band management fabric.

## Backend Network Topology

There are two topologies that will be discussed for the backend network: fat tree, and rail.
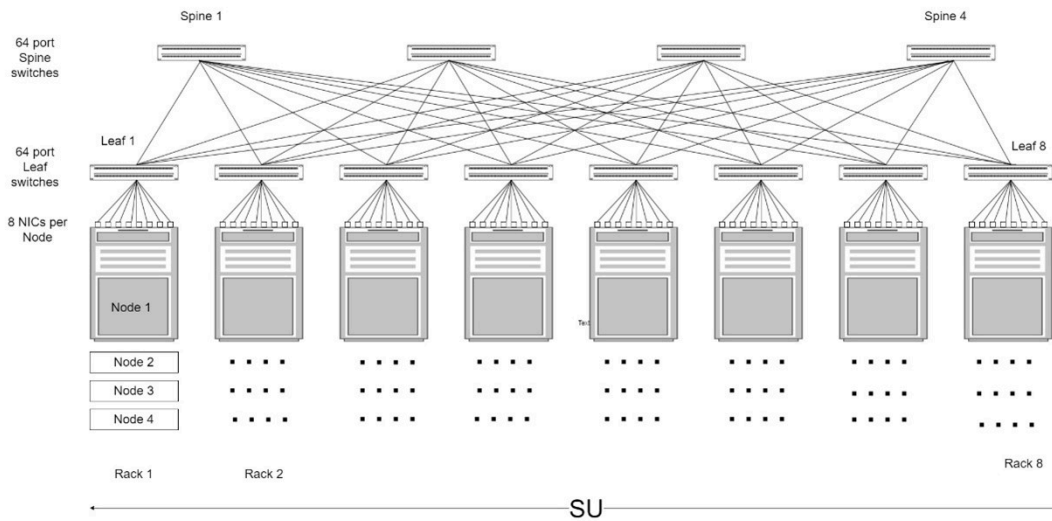
Consider the following when designing a network topology:

- Blocking factor: A switch has downlink and uplink ports; the blocking factor is defined as "downlink_port:uplink_port". For a 64-port switch using 32 uplink and 32 downlink ports, the blocking factor is 1:1. (A 1:1 blocking factor is defined as a non-blocking configuration).
- Undersubscription: A switch can have fewer downlink ports than uplink ports, for example a 64-port switch can have 24 downlink ports and 28 uplink ports, with 12 ports unused. This is referred to as 16% undersubscription. The safe approach is to have undersubscription (especially at higher switch tiers), but the cost effective approach is 1:1 which utilizes all switch ports.

### Fat Tree Non-blocking Topology

A 2-tier fat-tree consists of 2 layers of leaf-spine switches (T1, T2), with the T1 (leaf) switches connected to the NICs in the backend network. All NICs of a node are connected to the same T1 switch. A third tier adds a T3 layer of switches.

The fat tree topology is a familiar scalable design; some networks may require undersubscription to mitigate ECMP hash collisions (with a blocking design). The following diagram and table illustrate a 2-tier Fat Tree non-blocking topology.

**Figure 5.1:** A 32 Node 2-Tier Fat Tree Topology



**Table 5.1:** Fat Tree 2-Tier with Switch Radix = 64 (non-blocking 32 downlink, 32 uplink)

| GPUs, NICs (1:1) | Nodes | Leaf Switches | Spine Switches |
|---|---|---|---|
| 32 | 4 | 1 | 0 |
| 64 | 8 | 2 | 1 |
| 128 | 16 | 4 | 2 |
| 256 | 32 | 8 | 4 |
| 512 | 64 | 16 | 8 |
| 1024 | 128 | 32 | 16 |
| 2048 | 256 | 64 | 32 |

## Rail Topology

A 2-tier rail consists of 2 layers of leaf-spine switches (T1, T2), with the T1 (leaf) switches connected to the NICs in the backend network. Each NIC of a node is connected to one port of each T1 leaf switch. A 3rd tier adds a T3 layer of switches.

Rail topology benefits by containing traffic to rails, thereby minimizing probability of congestion. The communication libraries are dependent (aware) of the rail connections and the scale-up fabric. The following diagram and table illustrate a 2-tier rail topology.
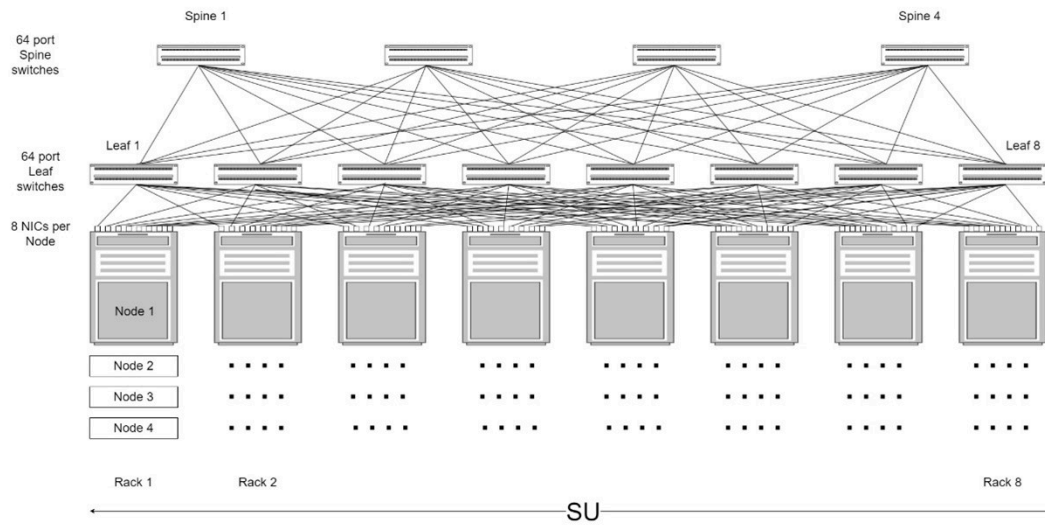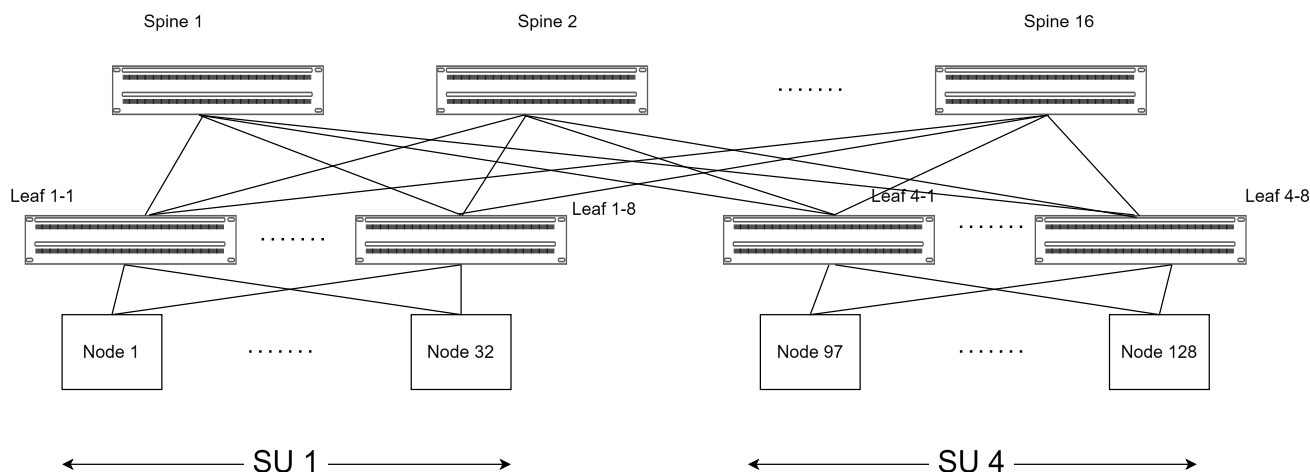
**Figure 5.2:** A 32 Node 2-Tier Rail Topology



**Table 5.2:** Rail Topology with 64 Port Switch (32 downink, 32 uplink)

| GPUs, NICs (1:1) | Nodes | Leaf Switches | Spine Switches |
|---|---|---|---|
| 256 | 32 | 8 | 4 |
| 512 | 64 | 16 | 8 |
| 768 | 96 | 24 | 12 |
| 1024 | 128 | 32 | 16 |
| 2048 | 256 | 64 | 32 |

The following figure is of a full 128 node cluster (where the nodes are in rail layout). Within a rail, a node is one hop from the other node. The layout can also be a fat tree where all the links from a rack terminate in a leaf switch, with similar number of leafs and switches for the same number of scalable units.

**Figure 5.3:** Layout of a Full 128 Node Cluster with 64 Port Switches



To build larger clusters, refer to the following table. The maximum numbers of nodes are dependent on switches and topology.

**Table 5.3:** Maximum Counts Based on Radix and Switch Tiers

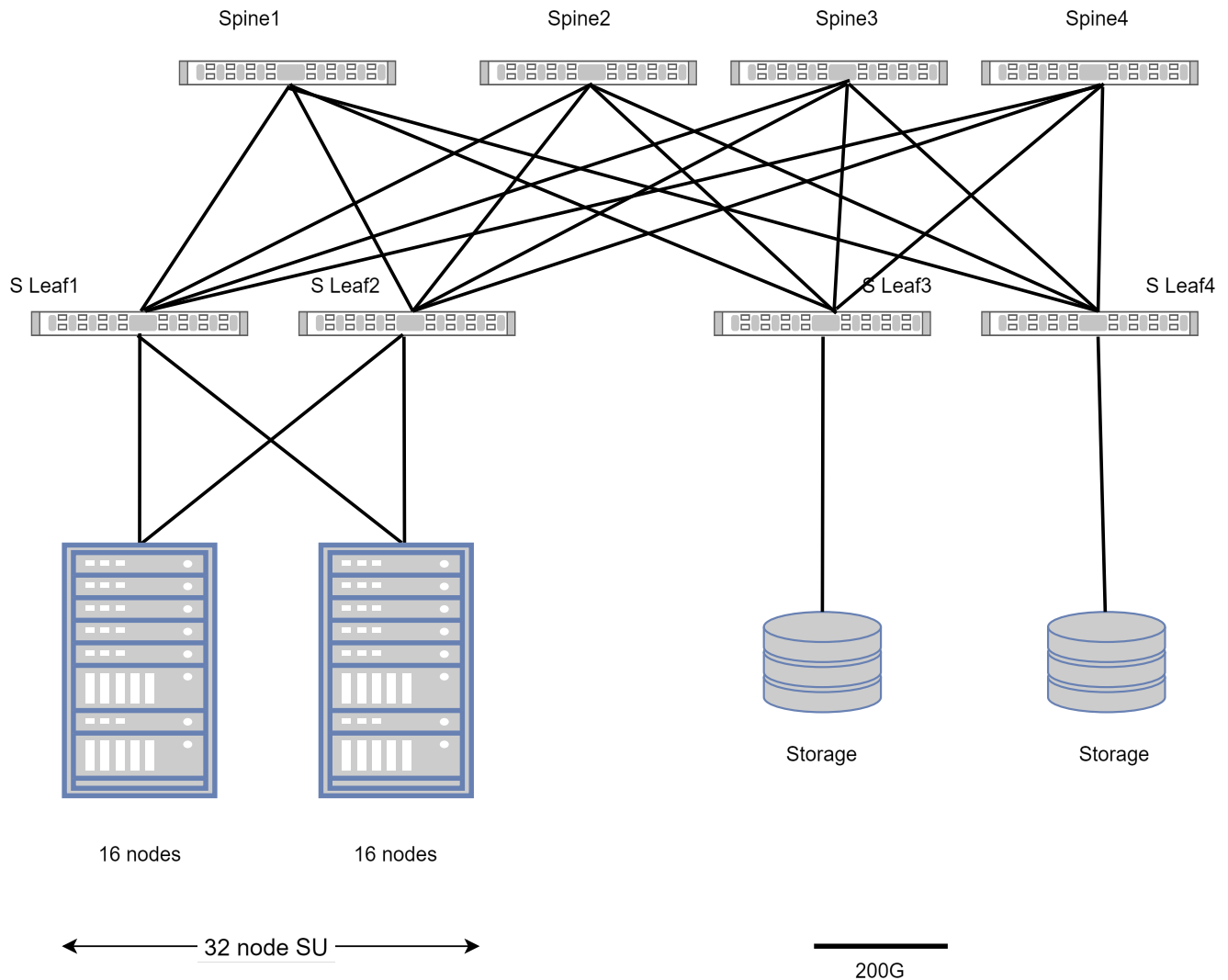| Parameter Count | 2 Tier Rail/Fat Tree (64 port, 400G) | 3 Tier Rail/ Fat Tree (64 port, 400G) | 2 Tier Rail/Fat Tree (128 port, 400G) | 3 Tier Rail/Fat Tree (128 port, 400G) |
|---|---|---|---|---|
| Switch Radix | 64 | 64 | 128 | 128 |
| NICs per Node | 8 | 8 | 8 | 8 |
| Max Leaf Switches | 64 | 2048 | 128 | 8192 |
| Max Spine Switches | 32 | 2048 | 64 | 8192 |
| Max Core Switches | -- | 1024 | -- | 4096 |
| Max NICs | 2048 | 65536 | 8192 | 524288 |
| Max GPUs | 2048 | 65536 | 8192 | 524288 |
| Max Nodes | 256 | 8192 | 1024 | 65536 |

# Frontend Network Topology

The frontend network composed of Ethernet NICs in 1:1 NIC:CPU organizations carries the storage and in-band communications, if a separate fabric is not provided. The in-band management network connects the cluster management services.

As datasets for AI workloads continue to expand in size, it is becoming increasingly critical that GPUs are not constrained by the I/O network and storage systems. The storage fabric provides the path between

GPU memory and the storage systems. Storage systems can be connected by the frontend network, but benefit by a separate plane of storage network.

**Figure 5.4:** Storage Network for a 32 Node Design



A separate optimized storage network as shown above provides benefits such as:

- Deep learning models accesses large datasets for training, a dedicated network provides frequent and iterative access to the data from the GPUs over the storage network.

- As datasets grow in size the capex and opex expenditures are kept separate from the compute needs.

# Chapter 6: AMD Software Tools for Cluster Operation and Management

The AMD open-source ROCm software platform, containers for AI/ML and Data Center Infrastructure empowers the accelerated computing community to innovate on top of a robust, flexible stack designed for scalability. These components work in concert to extract the full potential of heterogeneous architectures. The platform's open-source philosophy gives developers complete visibility while enabling customization and co-development. Users can optimize the ROCm software platform runtimes, programming models and utilities based on their workloads and scale requirements.

The AMD software components, shown below, consists of a collection of drivers, development tools, and APIs that enable GPU programming from low-level kernels to end-user applications. The ROCm Data Center tools (RDC), and AMD SMI (System Management Interface) are essential building blocks in cluster management and datacenter operation.

**Figure 6.1:** Software Stack with AMD ROCm, Container, and Infrastructure Blocks

| | |
|---|---|
| **AI Models and Algorithms**<br><br>Pytorch TensorFlow ONNX | **AI Ecosystem optimized for AMD** |
| **Workflow Orchestration and Job scheduling** | **Lamini, SLURM, Kubernetes** |
| **Cluster Management** | **Container Applications**<br>Redhat Openshift |
| **Data Center Management** | **AMD ROCm**<br>ROCm-SMI  ROCm Data Center Tool |
| **Hardware** | **AMD MI300 series GPU** |

## Cluster Management with RDC and SMI Tools

ROCm Data Center (RDC) enables GPU cluster administration with the capability of monitoring, validating and configuring policies. It enables full diagnostic and stress testing at cluster level. Administrators can use device monitoring, job statistics and error collection for a group of GPUs in a cluster and provides APIs for 3rd party integration. Full documentation and API reference are available at ROCm Data Center Tool documentation.

AMD System Managaement Interface (SMI) is a C library on linux providing user space interface to monitor and control AMD devices. The SMI libraries are available on AMD SMI Github Repository.

# Appendix A: Vendor List for Cluster Networking

AMD Instinct<sup>TM</sup> Accelerator powered servers are available from our partners. A [complete catalog](#) of qualified servers are available from [AMD Instinct Solutions](#).

**Table A.1:** NICs in Backend and Frontend Network

| Vendor | Link |
|---|---|
| AMD | [Pensando™ Giglio DPU 200G](#) |
| AMD | [Pensando™ Pollara 400](#) |
| AMD | [Pensando™ DSC3-400](#) |
| Broadcom | [Thor 2 400G](#) |
| Broadcom | [Thor 200G](#) |
| NVIDIA | [ConnectX®-7 400G](#) |

**Table A.2:** Switches in Backend Network

| Vendor | Link |
|---|---|
| Arista | [7060X6PE 51.2T Tomahawk 5](#) |
| Arista | [7060DX5-64S 25.6T Tomahawk 4](#) |
| Arista | [7800R4 Jericho 3 28.8T](#) |
| Arista | [Distributed Etherlink Switch](#)<br>• 7720R4-128PE (Ramon 102.T)<br>• 7700R4C-38PE (Jericho 3-AI 14.4T) |
| Cisco | [G200 51.2T](#) |
| Dell | [Z9864F-ON 51.2T Tomahawk 5](#) |
| Dell | [Z9664F-ON 25.6T Tomahawk 4](#) |
| Dell | [Z9432F-ON Trident4-X11](#) |
| Juniper | [QFX5240-64OD 51.2T Tomahawk 5](#) |
| Juniper | [QFX5230-64CD 25.6T Tomahawk 4](#) |
| Juniper | [PTX10008/10016 28.8T Express5](#) |
| Nokia | [7220 IXR-H4 Tomahawk 4](#) |
| Nokia | [7250 X1b/X3b Jericho 2C+](#) |

**Table A.2:** Switches in Backend Network (continued)

| Vendor | Link |
|---|---|
| Nokia | 7250 IXR-6e/10e/18e Jericho 2C+/Jericho 3 |

**Table A.3:** Switches in Frontend Network

| Vendor | Link |
|---|---|
| Arista | 7060X6-32PE 25.6T Tomahawk 5 |
| Arista | 7060DX5-64S 25.6T Tomahawk 4 |
| Arista | 7280R3A (up to 21.6T) Jericho 2C+ |
| Cisco | G200 51.2T |
| Dell | Z9864F-ON 51.2T Tomahawk 5 |
| Dell | Z9664F-ON 25.6T Tomahawk 4 |
| Dell | Z9432F-ON Trident4-X11 |
| Juniper | QFX5130-32CD 25.6T Trident4 |
| Juniper | QFX5220-32CD 12.8T Tomahawk 3 |
| Nokia | 7220 IXR-H4 Tomahawk 4 |
| Nokia | 7220 IXR-D5 Trident4 |

**Table A.4:** Switches in Storage Network

| Vendor | Link |
|---|---|
| Arista | 7050DX4-32S Trident4 |
| Arista | 7280R3A-72 (up to 21.6T) Jericho 2C+ |
| Dell | S5232F-ON Trident3-X7 |
| Juniper | QFX5230-64CD 25.6T Tomahawk 4 |
| Juniper | QFX5220-32CD Tomahawk 3 |
| Juniper | QFX5130-65CD Trident4 |
| Nokia | 7220 IXR-D5 Trident4 |
| Nokia | 7220 IXR-H4 Tomahawk 4 |
| Nokia | 7250 IXR-6e/10e/18e Jericho 2C+/Jericho 3 |

**Table A.5:** Switches in OOB Network

| Vendor | Link |
|---|---|
| Arista | 7010TX-48C Trident3 |
| Dell | S3248T-ON Trident3-X3 |
| Juniper | QFX5120 Trident3 |
| Juniper | EX4400 Trident3 |
| Nokia | 7220 IXR-D2L/D3L Trident3 |
| Nokia | 7215 IXS-A1 Mrvl AC5X |

Storage systems offer a high performance for data handling in AI workloads. Efficient storage systems allow GPUs to access data with low latency and prevent GPU stalls waiting for data completion. The following vendors offer storage systems in a cluster:

**Table A.6:** Storage Systems

| Vendor | Link |
|---|---|
| AMD-Supermicro | WEKAIO Reference Storage |
| Dell | Powerscale |
| HPE | HPE Greenlake |
| IBM | IBM Storage Scale System |

**Table A.7:** Validated Designs

| Vendor | Link |
|---|---|
| Dell | Dell Validated Design |
| Juniper | Juniper Validated Design |

📄 **Note:** This table will be updated as additional validated designs are made available.

# Appendix B: Acronyms

The acronyms used in this document are expanded in the following table.

**Table B.1:** Acronyms

| Acronym | Definition |
| --- | --- |
| AI | Artificial Intelligence |
| API | Application Programming Interface |
| BIOS | Basic Input/Output System |
| BMC | Baseboard Management Controller |
| CNP | Congestion Notification Packet |
| CPU | Central Processing Unit |
| DDR | Double Data Rate |
| DNS | Domain Name System |
| DRAM | Dynamic Random Access Memory |
| ECMP | Equal Cost MultiPath |
| GPU | Graphics Processing Unit |
| HPC | High Performance Computing |
| IP | Internet Protocol |
| IPMI | Intelligent Platform Management Interface |
| NIC | Network Interface Card |
| NVMe | Non Volatile Memory Express |
| OAM | OCP Accelerator Module |
| OCP | Open Compute Project |
| OOBM | Out Of Band Management Network |
| OS | Operating System |
| PCI | Peripheral Component Interconnect |
| PCIe | PCI Express |
| RDC | ROCm Data Center |
| RDMA | Remote Direct Memory Access |
| RoCE | RDMA over Converged Ethernet |

**Table B.1:** Acronyms (continued)

| Acronym | Definition |
|---------|------------|
| SMI | System Management Interface |
| SSD | Solid State Drive |
| ToR | Top of Rack |
| UALink | Ultra Accelerator Link Consortium |
| UBB | Universal Baseboard |
| UEC | Ultra Ethernet Consortium |

# Appendix C: Additional Resources and Legal Notices

## Revision History

The following table shows the revision history for this document.

| Revision Summary |
|---|
| November 2024 Version 1.00 |
| Initial release. |
| March 2025 Version 1.1 |
| • Cluster Architecture: Updated component counts in Table 4.1 and Table 4.2. <br> • Topology of Network Fabrics: Updated Table 5.2 with 2-tier rail. <br> • Vendor List for Cluster Networking: Updated vendor list and added Table A.7. <br> • Minor non-technical edits throughout to comply with AMD standards. |

## Notices

## Trademarks

AMD, the AMD Arrow logo, and combinations thereof are trademarks of Advanced Micro Devices, Inc.

Other product names used in this publication are for identification purposes only and may be trademarks of their respective companies.