

# Dell PowerSwitch Z9664F-ON 400GbE Ethernet AI Fabric Performance

## EXECUTIVE SUMMARY

AI training involves massive amounts of data being processed in parallel by GPUs across a collection of high-end servers. Network designers need to have confidence that their Ethernet networking fabric can handle the demands of AI processing and perform at or near wire speed. The Dell PowerSwitch line offers high-performance and high-port density that can serve as a network fabric for AI.

Dell Inc. commissioned Tolly to benchmark a 400GbE network fabric consisting of 10 Dell PowerSwitch Z9664F-ON switches with Enterprise SONiC Distribution by Dell Technologies version 4.3.0 in a RAIL optimized topology. Eight Dell PowerEdge servers, each outfitted with eight GPUs and eight 400GbE network interfaces, were connected to the switch fabric and were tasked with running various AI workloads across the network fabric.

The Dell PowerSwitch Ethernet Fabric provided zero-loss handling of AI traffic from the 64 connected NICs/GPUs delivering ~391Gbps of Perfetest RDMA throughput and ~390GBps of NCCL benchmark inter-node throughput. See Figure 1.

### THE BOTTOM LINE

The Dell Z9664F-ON PowerSwitch Ethernet Fabric demonstrated support for carrying:

- 1 ~391Gbps RDMA perftest throughput
- 2 ~390GBps inter-node bandwidth for xCCL test (8 x 400GbE NICs)
- 3 AI training via LLM fine-tuning with zero packet loss

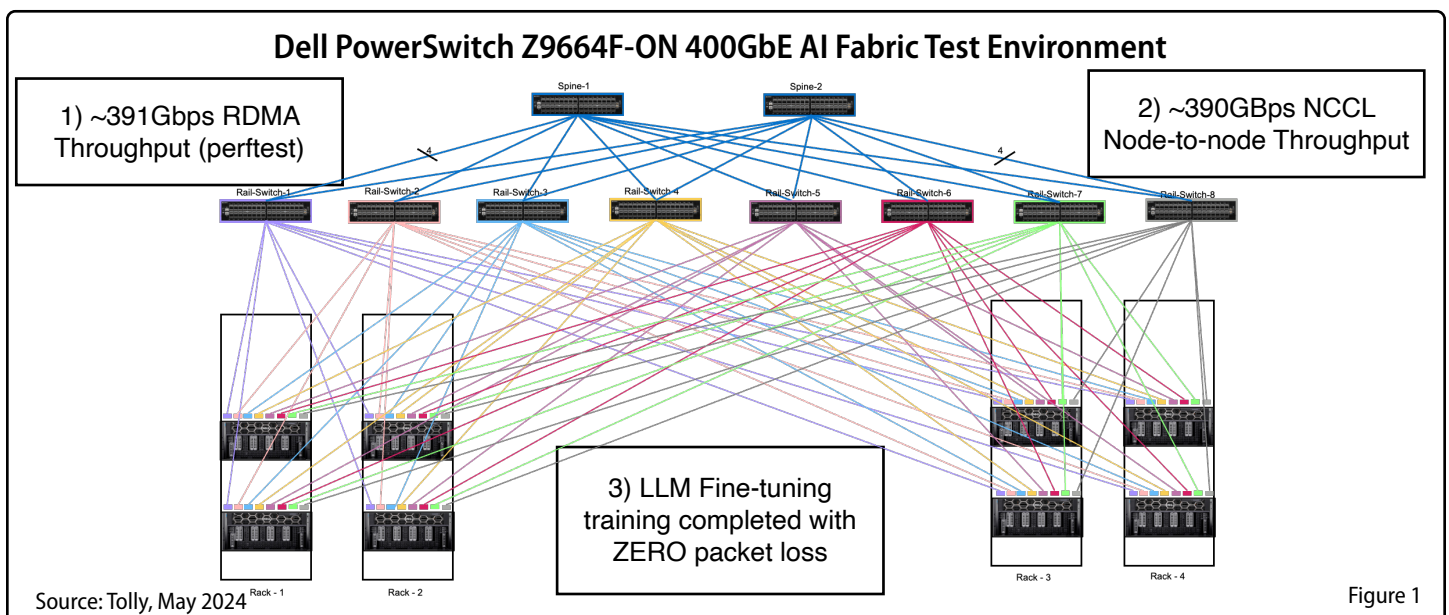


Figure 1

## Goal

While the test configuration was complex, the test goal was simple - to demonstrate that the Dell PowerSwitch 400G Ethernet network fabric could transport RDMA and NCCL traffic at or near line rate across the 400GbE network fabric and without packet loss. As noted briefly above, the test was successful.

## Configuration

While many test environments are small and represent a microcosm of a real-world environment, this test bed was a real-world environment, running actual AI micro benchmarks and applications, one of the most extensive ever used in a Tolly test. The scale was chosen to provide a realistic load and traffic flow and testers used recognized AI micro benchmarks to generate the workload on the systems and, thus, the

traffic across the network fabric. Configuration details are found in Tables 1 through 3 later in the report.

## Switch Fabric

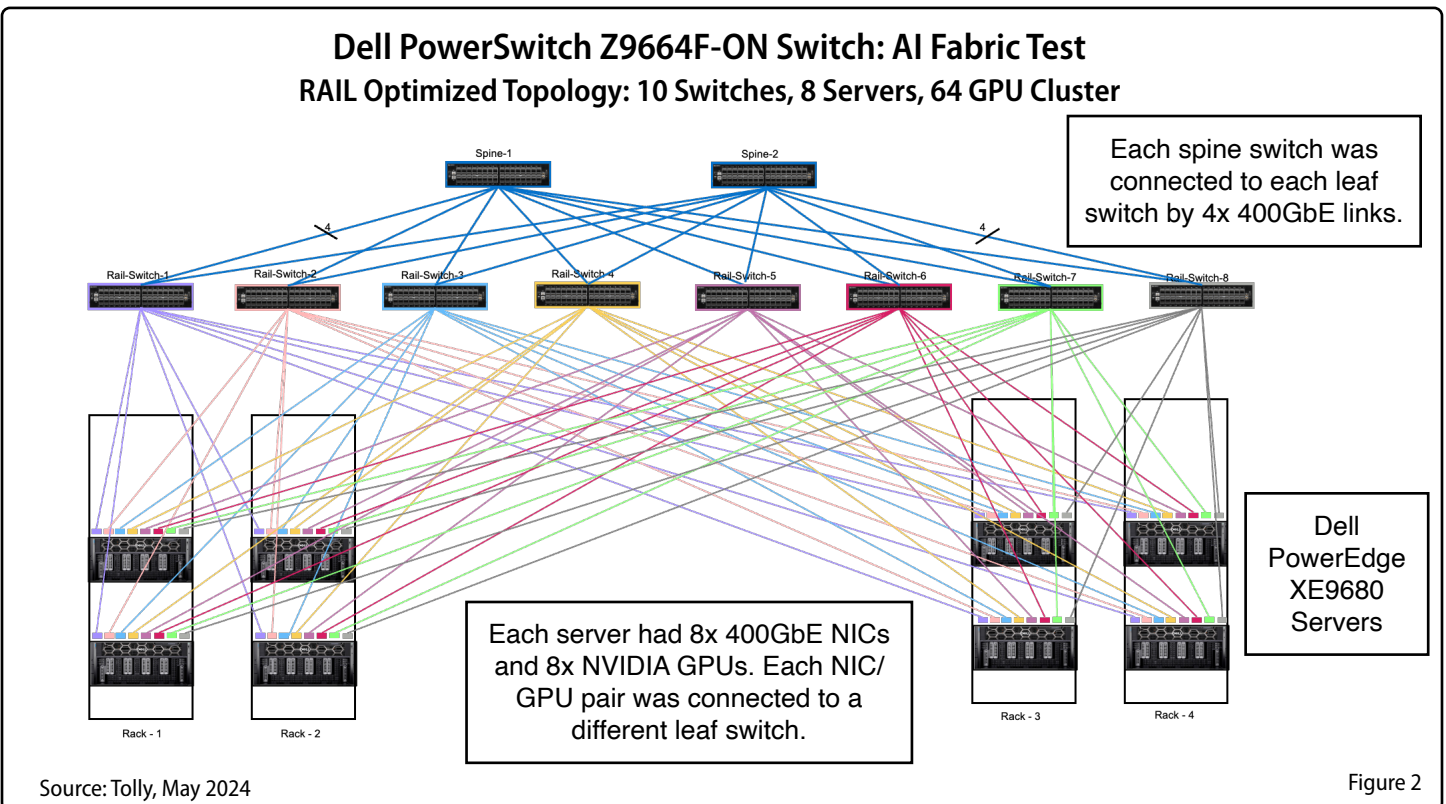
All LAN switches were Dell PowerSwitch Z9664F-ON devices with all ports set to 400Gbps. The fabric topology consisted of two spine switches and eight leaf switches with BGP using RAIL optimized architecture. The leaf switches used adaptive routing and switching to implement the fabric.

There were four, 400GbE links between each of the two spine switches and each of the eight leaf switches totaling 32x 400GbE uplinks for each spine switch. (See Figures 1 and 2, same diagram with different notations.) Dell has developed a “roce enable” command that configures its switches to provide optimized settings for

RDMA over converged Ethernet (RoCE) traffic.

## AI Server Nodes

The test used eight Dell PowerEdge XE9680 servers shown, two to a rack, in Figures 1 and 2. Each server was outfitted with 8x 400GbE NICs and 8 GPUs. To optimize traffic flow for GPU-to-GPU parallel processing, each GPU used a dedicated NIC. NIC 1 in each server node was connected to switch 1, NIC 2 to switch 2, and so forth. This is known as a RAIL optimized configuration. This configuration provides for GPUs in the same server node to communicate either internally or via the spine switches.





# Tests & Results

## RDMA Throughput

RDMA throughput was demonstrated using the perftest package. Perftest is an open source micro benchmark, architected as a client-server model, that drives RDMA reads/writes.

For this test, one node was configured as the target server and another was configured as clients. RDMA "reads" were issued from the clients to the server using the `ib_read_bandwidth` command. Test was run on all eight NICs simultaneously on both client and server nodes.

The test was run for two minutes and the network throughput for each of the eight NICs averaged 390.89Gbps, close to the 400Gbps maximum. Tolly engineers confirmed that the switch links being traversed by the pretest traffic were showing utilization of 99% or 100%.

## NVIDIA Collective Communication Library (NCCL) Test

These tests are provided by the GPU vendor, in this case NVIDIA and consist of various test routine. For this test, the "All reduce" routine was used. The algorithm, chosen by the benchmark, was "ring" for the eight nodes involved in the test. RoCE v2 was used.

The NCCL test was run using an open source "Pytorch" container and the test data size was 16GBytes. Test run time was unlimited.

The test showed node-to-node (i.e., eight NICs/ports to eight NICs/ports) bus bandwidth to be 389.7Gbytes. This means that each of the eight, 400GbE NICs and associated switch ports were running at wire speed.

## LLM Fine-Tuning

This test also used all eight nodes. Key elements:

- Model: Llama-2 70B
- Precision: bf16

- Fine-tuning Technique: PEFT - LoRA
- Data Set: Databricks-Dolly-15K
- Steps: 1010, Global Batch Size: 128, Micro Batch Size =1, TP=8, PP=1
- LLM Training Framework: Nemo

The purpose of this test was to illustrate that AI training could take place across the Ethernet fabric without packet loss caused by congestion.

The fine-tuning test ran for 2 hours and 6 minutes and completed successfully. The test consisted of multiple functions as the model was loaded and the test was run. Thus, the load on the network fabric varied during the test. The Dell Ethernet fabric was more than sufficient for fine-tuning workloads

After the test completed, Tolly engineers reviewed the counters on the Dell PowerSwitch fabric and confirmed that there was zero packet loss.

With any AI model, the training process or training loss is another key parameter that requires attention when benchmarking any AI model relative to the infrastructure. Training loss is a measure of how well a model is performing on the training data. The end goal is to decrease over time and be at 1 or under.

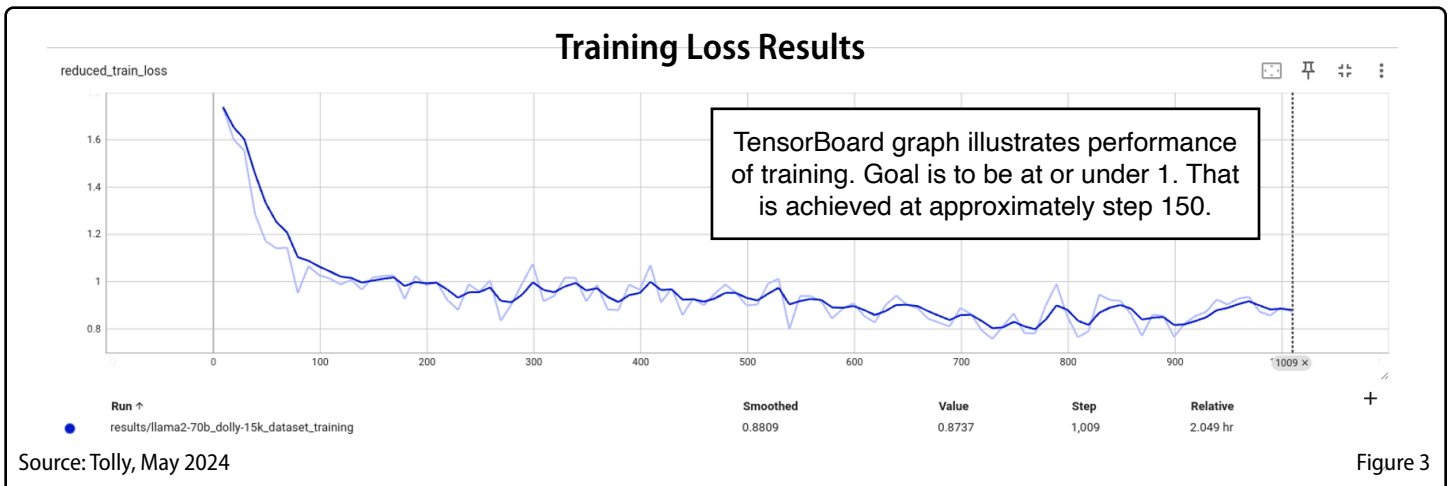


Figure 3



Figure 3 shows how the Dell Z9664F-ON switch fabric performed optimally. At the beginning of the model training, one can see how training loss is high, but eventually, as the training continues, the training loss decreases as expected and is equal to or less than 1. If the fabric infrastructure does not perform well, if for example, there is packet loss, congestion, or low throughput, training loss and overall completion time would be adversely affected.

## Test Setup & Methodology


The test methodology and setup for each test was explained, as appropriate, in the "Test Results" section of the report. Switch hardware and software as well as server hardware and software components and versions levels can be found in Tables 1 through 3.

Dell customers interested in additional details of server BIOS settings and specific switch configuration parameters should contact their Dell technical team.

**Dell Inc.**

**Dell PowerSwitch Z9664F-ON 100/400GbE Aggregation Switch**

**AI Ethernet Fabric Performance Evaluation**



*Tested  
May  
2024*

### Dell Ethernet AI Fabric Test Bed Details

#### Switch Fabric (Same for All 10 Switches)

Vendor	Switch	OS	Version	Notes
Dell	Z9664F-ON	Enterprise SONIC Distribution by Dell Technologies	4.3.0 GA	RAIL Optimized switch-server configuration: Two spine, eight leaf switches. All ports configured for 400Gbps. <span style="float: right;">Table 1</span>

#### Server Configuration (Same for All 8 Servers)

<b>Vendor/System</b>	Dell PowerEdge XE9680 (48c, 350W)
<b>CPU</b>	2x Intel Xeon 8468 (48c, 350W)
<b>Memory</b>	16x 64GB 4800 GHz RDIMMS
<b>GPU</b>	8x NVIDIA H100 module
<b>BOSS M.2</b>	2x 480GB
<b>Drives</b>	4x 1.92TB Samsung U.2 NVMe SSD
<b>PCIe cards</b>	8x NVIDIA ConnectX-7 VPI 400G 2x CX-6Dx 100GbE
<b>Onboard LOM</b>	Broadcom BCM5720 2x 1GbE

Table 2

#### Server Software

<b>Node Operating System</b>	Ubuntu Linux 22.04.3 LTS
<b>NVIDIA CUDA</b>	12.4
<b>NVIDIA GPU Driver</b>	550.54.15
<b>NVIDIA ConnectX-7 Firmware</b>	28.39.1002
<b>OFED Driver</b>	MLNX_OFED_LINUX-24.01-0.3.3.1
<b>NCCL</b>	2.21.5
<b>NVIDIA Fabric Manger</b>	550.54.15
<b>NGC Pytorch Container</b>	24.03

Table 3

Source: Tolly, May 2024

Note: Each server used eight 400GbE NICs and eight GPUs for running the tests.



## About Tolly

The Tolly Group companies have been delivering world-class IT services for more than 35 years. Tolly is a leading global provider of third-party validation services for vendors of IT products, components and services.

You can reach the company by E-mail at [sales@tolly.com](mailto:sales@tolly.com), or by telephone at +1 561.391.5610.

Visit Tolly on the Internet at:

<http://www.tolly.com>

## Terms of Usage

This document is provided, free-of-charge, to help you understand whether a given product, technology or service merits additional investigation for your particular needs. Any decision to purchase a product must be based on your own assessment of suitability based on your needs. The document should never be used as a substitute for advice from a qualified IT or business professional. This evaluation was focused on illustrating specific features and/or performance of the product(s) and was conducted under controlled, laboratory conditions. Certain tests may have been tailored to reflect performance under ideal conditions; performance may vary under real-world conditions. Users should run tests based on their own real-world scenarios to validate performance for their own networks.

Reasonable efforts were made to ensure the accuracy of the data contained herein but errors and/or oversights can occur. The test/audit documented herein may also rely on various test tools the accuracy of which is beyond our control. Furthermore, the document relies on certain representations by the sponsor that are beyond our control to verify. Among these is that the software/hardware tested is production or production track and is, or will be, available in equivalent or better form to commercial customers. Accordingly, this document is provided "as is", and Tolly Enterprises, LLC (Tolly) gives no warranty, representation or undertaking, whether express or implied, and accepts no legal responsibility, whether direct or indirect, for the accuracy, completeness, usefulness or suitability of any information contained herein. By reviewing this document, you agree that your use of any information contained herein is at your own risk, and you accept all risks and responsibility for losses, damages, costs and other consequences resulting directly or indirectly from any information or material available on it. Tolly is not responsible for, and you agree to hold Tolly and its related affiliates harmless from any loss, harm, injury or damage resulting from or arising out of your use of or reliance on any of the information provided herein.

Tolly makes no claim as to whether any product or company described herein is suitable for investment. You should obtain your own independent professional advice, whether legal, accounting or otherwise, before proceeding with any investment or project related to any information, products or companies described herein. When foreign translations exist, the English document is considered authoritative. To assure accuracy, only use documents downloaded directly from Tolly.com. No part of any document may be reproduced, in whole or in part, without the specific written permission of Tolly. All trademarks used in the document are owned by their respective owners. You agree not to use any trademark in or as the whole or part of your own trademarks in connection with any activities, products or services which are not ours, or in a manner which may be confusing, misleading or deceptive or in a manner that disparages us or our information, projects or developments.