

JULY 2024

Powering Next-generation Networks for GenAI Workloads With Dell Technologies PowerSwitch Z9864F-ON

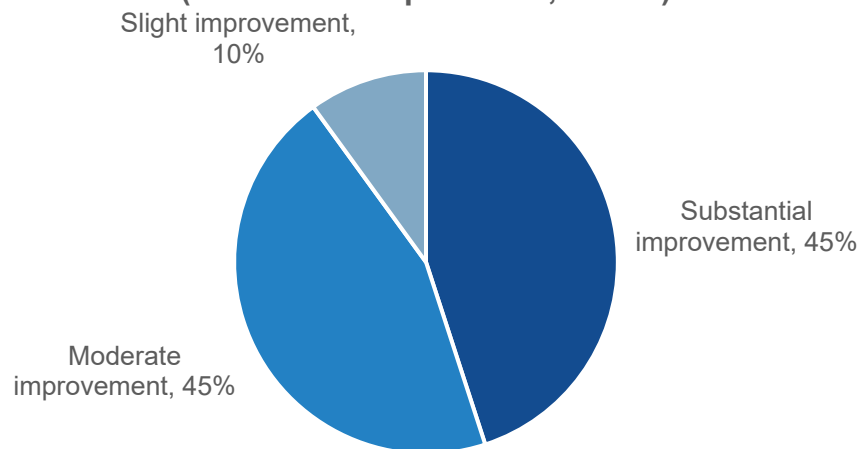
Alex Arcilla and Tony Palmer, Principal Analysts – Validation Services

Network Infrastructure Challenges for Supporting GenAI Workloads

Of all the artificial intelligence (AI) and machine learning (ML) technologies available, research from TechTarget's Enterprise Strategy Group found that organizations plan to make some of the most significant investments in generative AI (GenAI).¹ And organizations have high expectations for this and other AI technologies, as 90% of respondents expect moderate to substantial improvements in the speed and efficiency of their organizational processes and workflows (see Figure 1).² Yet the proper IT infrastructure must be implemented to support the high performance demands imposed by AI workloads. For GenAI, the demands originate from the need to process massive amounts of data for building, training, and fine-tuning AI models and inference, especially when needing to deliver highly accurate responses in a conversational flow.

Figure 1. Expected Impact of AI on Organizational Processes and Workflows

Which of the following best describes AI's impact on the speed and efficiency of your organization's processes or workflows?
(Percent of respondents, N=339)



Source: Enterprise Strategy Group, a division of TechTarget, Inc.

One approach to modernizing the current infrastructure—both front-end (application traffic, storage access, general network) and back-end (GPU-to-GPU connectivity)—has been to leverage proprietary technologies, such as InfiniBand, for GenAI deployments. Yet, the approach is a significant drawback for organizations that have

¹ Source: Enterprise Strategy Group Research Report, [2024 Technology Spending Intentions Survey](#), February 2024.

² Source: Enterprise Strategy Group Complete Survey Results, [Navigating the Evolving AI Infrastructure Landscape](#), December 2023.

This Enterprise Strategy Group Technical First Look was commissioned by Dell Technologies and is distributed under license from TechTarget, Inc.

standardized their networks on Ethernet. The skills and expertise to implement proprietary or other technologies aside from Ethernet, such as other open standard technologies, are scarce, leading to delayed deployments and longer time to value. Organizations need to ensure that the underlying network of any deployed GenAI infrastructure can ingest and process large data sets from distributed sources, orchestrate parallel computations across multiple GPU-enabled nodes, and deliver insights to end users in real time.

Solution – Dell Technologies Open Ethernet-powered Switch Fabric for GenAI

Dell Technologies has designed the PowerSwitch Z9864F-ON to deliver the throughput and low latency organizations need for AI-driven workflows. The latest in the Dell PowerSwitch Z-series has been designed as a high-performance, high-density, open networking, and multi-rate AI fabric switch that delivers the flexibility to support AI and ML fabric solutions, as well as other workload types that are characterized by intensive compute and storage traffic, such as IoT and streaming video.

The Dell PowerSwitch Z9864F-ON provides multi-rate speeds to help organizations enable denser switch footprints while simplifying migration to 800 Gbps as business demands dictate. With 64 ports of 800GbE contained in an OSFP112 form factor and 2U design, organizations that seek flexibility within their network infrastructure can use Dell PowerSwitch Z9864F-ON as a 100/200/400G switch via breakout cables for a maximum of 320 ports. For the increased throughput AI-driven workflows exhibit, this switch leverages the Broadcom Tomahawk 5. This silicon supports the 51.2 Tb/second throughput (half-duplex) ideal for supporting the high performance and low latency requirements of protocols ideal for AI workflows, such as RDMA over Converged Ethernet version 2 (ROCEv2).

With Dell's ongoing commitment to providing open networking solutions for network operating systems, orchestration, and monitoring, the Dell PowerSwitch Z9864F-ON supports the open-source Open Network Install Environment (ONIE) for zero-touch installation of Enterprise SONiC Distribution by Dell Technologies. To further maximize GenAI workload performance, Dell has developed advanced features into its SONiC distribution such as dynamic load balancing, adaptive routing, and RoCE congestion control. With SONiC, organizations can continue to expect the global support, scale, and features that organizations require, especially when supporting GenAI workloads.

As with all their other products and solutions, Dell Technologies provides the ongoing 24/7 global support and expertise in network design, implementation, and ongoing maintenance. Organizations can benefit when deploying GenAI, as Dell Technologies offers its extensive expertise and experience in network infrastructure and design via validated designs covering both inferencing and model customization. Dell also offers managed services to operate these AI environments should business needs dictate.

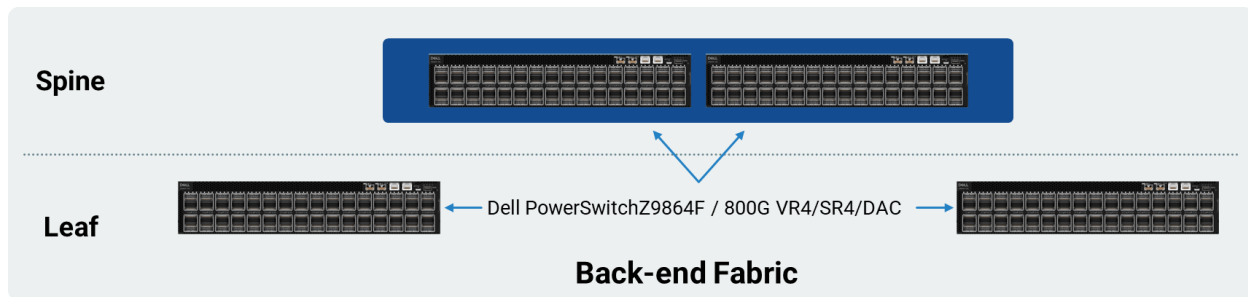
First Look

Enterprise Strategy Group took a closer look at the Dell PowerSwitch Z9864F-ON to examine how organizations can build out the network connectivity required for GenAI workloads as well as achieve denser data center footprints. Based on our review, we noted that the Dell PowerSwitch Z9864F-ON:

- Offers the following port breakdown to enable larger number of interfaces per switch system, including:
 - 64 ports of 800GbE (OSFP112 connectors with PAM-4 support).
 - 128 ports of 400GbE (100GbE PAM-4 support in breakout mode).
 - 256 ports of 200GbE (in breakout mode).
 - 320 ports of 100GbE.
 - 32 ports of 800G and 64 ports of 400G.

- Can be implemented as a high capacity Ethernet fabric for AI/ML training and inferencing clusters containing up to 8,000 GPU nodes in a single two-tier fabric.
- Can be used in leaf and spine topologies along with PowerSwitch Z- and S-series switches, enabling cost-effective aggregation of 100/200/400/800G uplinks (see Figure 2).

Figure 2. Dell PowerSwitch Z9864F-ON in Leaf and Spine Topologies



Source: Enterprise Strategy Group, a division of TechTarget, Inc.

- Demonstrates Dell Technologies' commitment to sustainability, as the switch's design facilitates power-efficient operation up to 45°C. Organizations can use the switch to reduce cooling costs in temperature-constrained data center environments.
- Incorporates a high radix of 200GbE ports—with 256 ports supported on a single chip—that can enable a flat AI/ML cluster for supporting low latency AI workloads.
- Supports features that can optimize network flows associated with AI workloads, such as versatile hashing, adaptive routing, and enhanced priority-based flow control queue configuration support (enabled by the Broadcom Tomahawk 5 chipset).

Conclusion

Using AI technologies, including GenAI, to improve process and workflow efficiency, a highly performant network is required for supporting the large amounts of compute and storage traffic indicative of AI workloads. Enterprise Strategy Group has evaluated multiple releases of the Dell PowerSwitch series and found that Dell Technologies has delivered the switches organizations need to build out end-to-end and high-performance networks for supporting both traditional and modern workloads. The Dell PowerSwitch Z9864F-ON is the latest evidence of this as AI workload adoption takes off. Based on our first look at this switch, Enterprise Strategy Group suggests seriously considering the Dell PowerSwitch Z9864E-ON as you leverage more AI into your operations.

©TechTarget, Inc. or its subsidiaries. All rights reserved. TechTarget, and the TechTarget logo, are trademarks or registered trademarks of TechTarget, Inc. and are registered in jurisdictions worldwide. Other product and service names and logos, including for BrightTALK, Xtelligent, and the Enterprise Strategy Group might be trademarks of TechTarget or its subsidiaries. All other trademarks, logos and brand names are the property of their respective owners.

Information contained in this publication has been obtained by sources TechTarget considers to be reliable but is not warranted by TechTarget. This publication may contain opinions of TechTarget, which are subject to change. This publication may include forecasts, projections, and other predictive statements that represent TechTarget's assumptions and expectations in light of currently available information. These forecasts are based on industry trends and involve variables and uncertainties. Consequently, TechTarget makes no warranty as to the accuracy of specific forecasts, projections or predictive statements contained herein.

Any reproduction or redistribution of this publication, in whole or in part, whether in hard-copy format, electronically, or otherwise to persons not authorized to receive it, without the express consent of TechTarget, is in violation of U.S. copyright law and will be subject to an action for civil damages and, if applicable, criminal prosecution. Should you have any questions, please contact Client Relations at cr@esg-global.com.