

FAST TRACK GENERATIVE AI WITH DELL™ POWEREDGE™ XE9680

MADE POSSIBLE WITH NVIDIA® H100 TENSOR CORE GPUS
& BROADCOM 100 GIG-ETHERNET



Introduction

Unlock the potential of AI and transform your business with Dell™ PowerEdge™ XE9680 server - A machine that will elevate your AI performance to new heights.

The field of generative AI is experiencing an explosive growth, with cutting-edge developments in image, video, and audio media creation revolutionizing the creative industry. This remarkable technology is driving innovation in diverse sectors and opening up new frontiers for creative expression. Moreover, the immense promise of fine-tuned enterprise LLMs is bringing unparalleled insights to businesses, enabling them to protect their proprietary information, comply with data sovereignty issues, and improve their internal and external effectiveness. With the potential to process vast amounts of data in real-time, these finely-tuned models offer a decisive advantage in the fast-paced world of modern business. So, if you want to unlock the full potential of your data assets and stay ahead of the competition, enterprise LLMs are the way to go!

Dell™ PowerEdge™ XE9680 server with 8 NVIDIA® H100 GPUs and Broadcom BCM57508 100G Ethernet offers unmatched performance for high-performance AI training and inference. Broadcom BCM957508-P2100G is a dual-port 100 Gb/s PCI Express 4.0 x16 Network Interface Card that supports QSFP56/QSFP28 optical modules and copper direct-attach cables which makes the card a perfect choice for network-intensive AI applications. Our whitepaper evaluates its effectiveness on common AI workloads like language and image recognition.

PERFORMANCE

Dell™ PowerEdge™ XE9680 Server
(NVIDIA® H100 vs A100)

1.8x

While Training GPT-2

3.15x

with the NVIDIA® Transformer engine enabled (Float8)

➤ **Optimize your AI workloads and stay ahead of the competition with Dell™ PowerEdge™ XE9680**



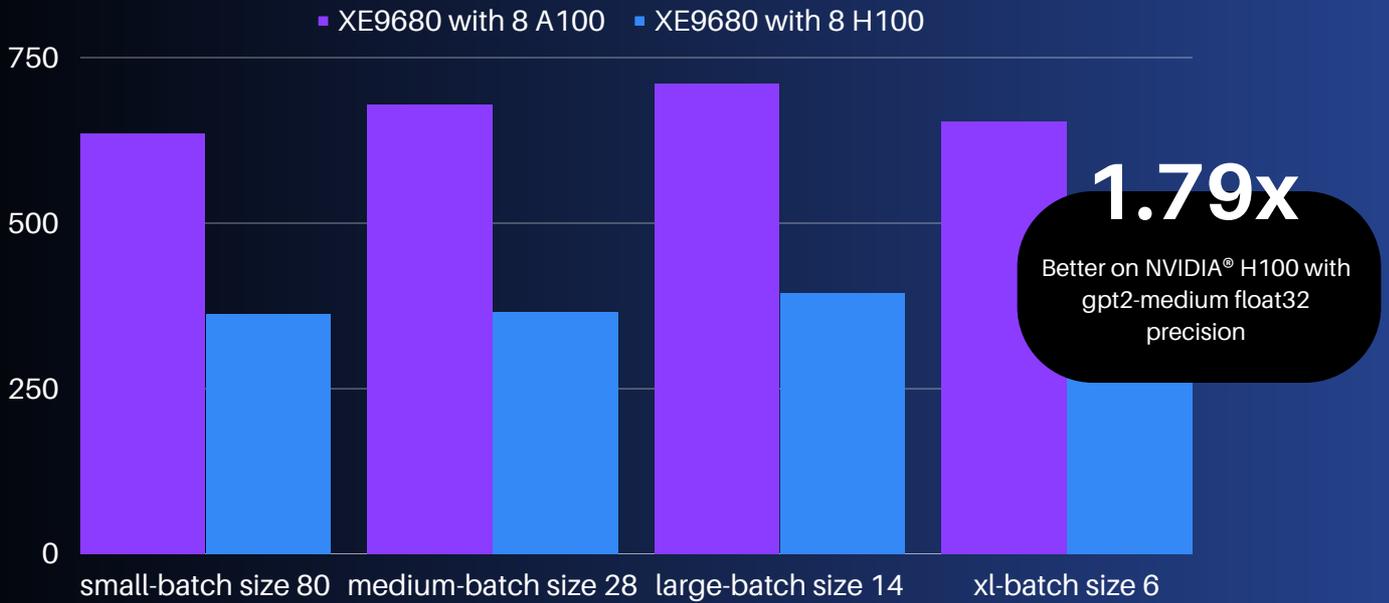
Dell™ PowerEdge™ XE9680 is purpose-built for generative AI workloads, including language and generative AI models. We found it to demonstrate compelling performance in both training and fine-tuning of LLMs.

- Chetan Gadgil, CTO at Scalers AI™

The Results | Performance Summary

(Training of GPT-2)

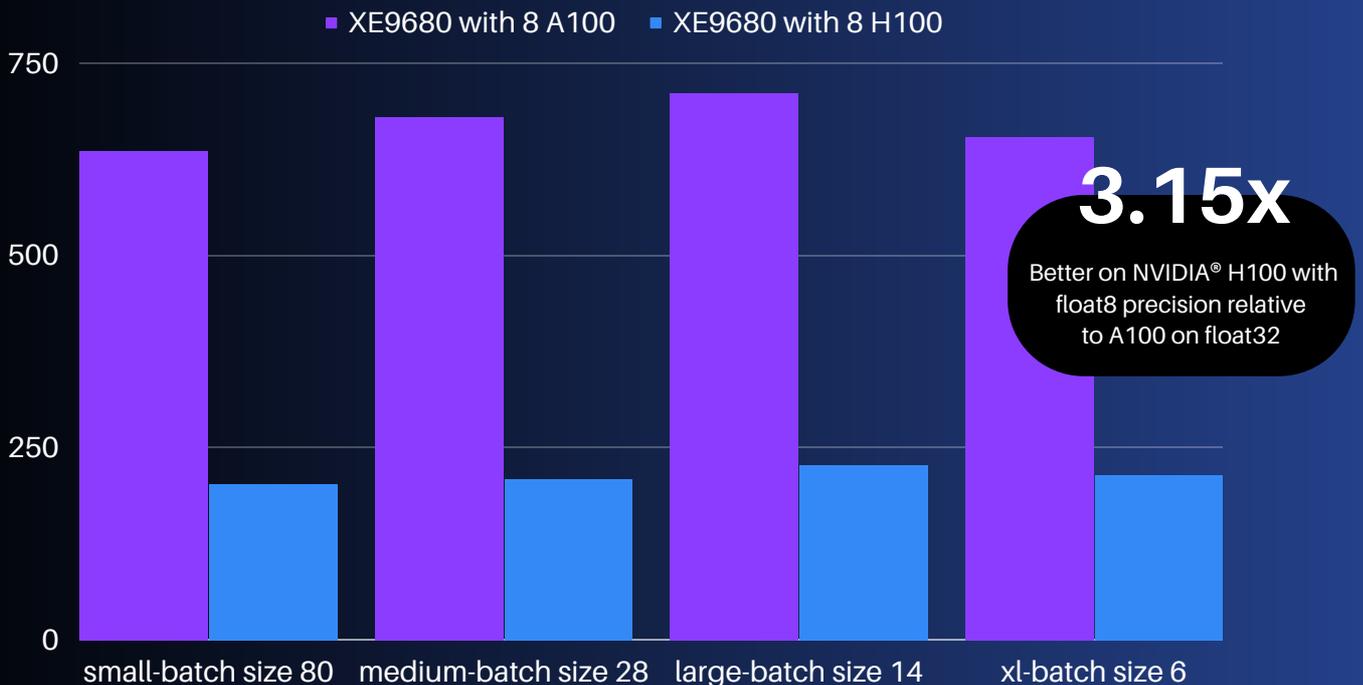
- Dell™ PowerEdge™ XE9680 server with 8 NVIDIA® A100 GPUs and Broadcom BCM57508
- Dell™ PowerEdge™ XE9680 server with 8 NVIDIA® H100 GPUs and Broadcom BCM57508



GPT-2 float32 - Training Step Time (ms)

Lower is Better

> With FP8 on Dell™ PowerEdge™ XE9680 server with 8 NVIDIA® H100 GPUs and Broadcom BCM57508



GPT-2 float8 - Training Step Time (ms)

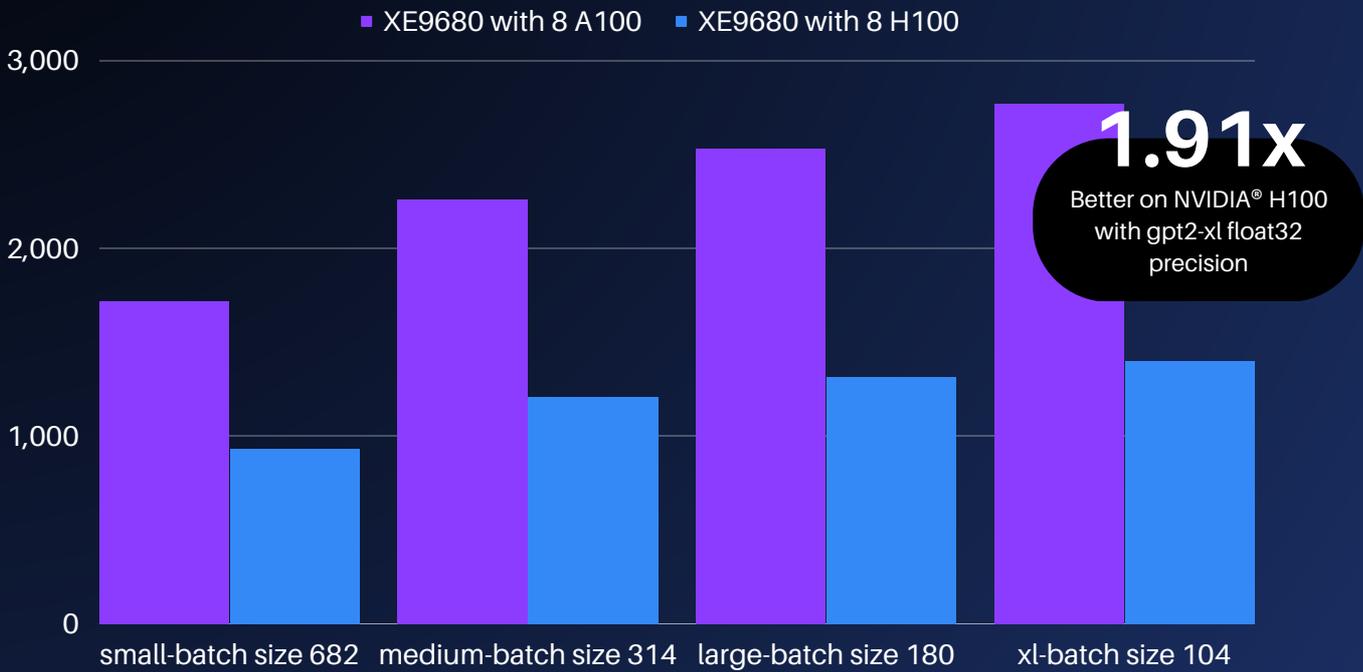
Lower is Better

*Accuracy of FP8 models was not tracked. FP8 using the NVIDIA® Transformer Engine is available only on H100

The Results | Performance Summary

(Inference of GPT-2)

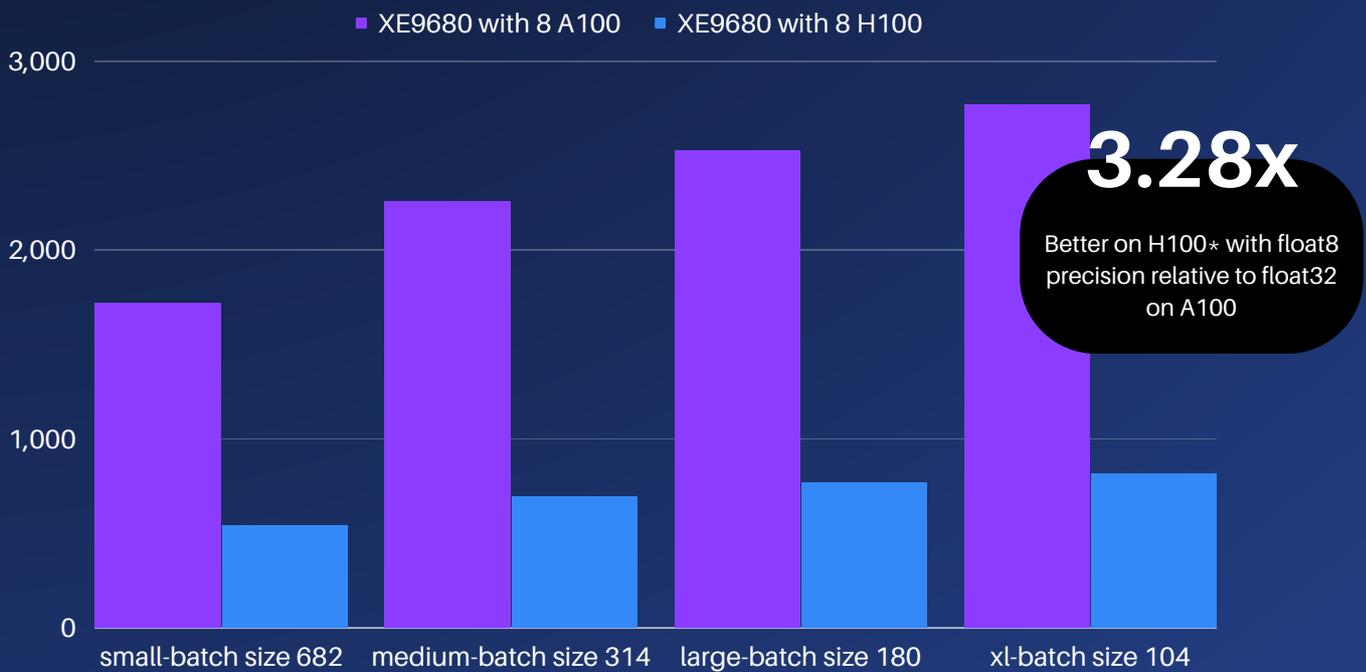
➤ Dell™ PowerEdge™ XE9680 server with 8 NVIDIA® H100 GPUs and Broadcom BCM57508



GPT-2 float32 - Inference Step Time (ms)

Lower is Better

➤ With FP8 on Dell™ PowerEdge™ XE9680 server with 8 NVIDIA® H100 GPUs and Broadcom BCM57508



GPT-2 float8 - Inference Step Time (ms)

Lower is Better

*Accuracy of FP8 models was not tracked. FP8 is available only on NVIDIA® H100

*Disclaimer - Performance varies by use case, model, application, hardware & software configurations, the quality of the resolution of the input data, and other factors. This performance testing is intended for informational purposes and not intended to be a guarantee of actual performance of an AI application.

| Generative AI

Generative AI is a fascinating technology that has revolutionized the field of artificial intelligence. Through its remarkable ability to generate realistic images and coherent text, it can create entirely new worlds that seem to come alive before your very eyes. From generating lifelike portraits to producing complex narratives, generative AI has ushered in a new era of innovation and discovery.

| Example Real World Applications

Custom Chatbot

Dell™ PowerEdge™ XE9680 offering for Generative AI workloads, specifically Large Language Models are very well suited. On Dell™ PowerEdge™ XE9680 server, we successfully built, deployed, and operated a chatbot designed to answer queries about Dell Technology World 2023. Our approach began by generating embeddings using a pre-trained Large Language Model (LLM) that processed text from Dell Technology World 2023 website. These embeddings were then loaded into a vector database, forming a solid foundation for the chatbot's knowledge.

Next, we developed a micro-service specifically designed to handle question-and-answer interactions. By integrating this service with the vector database, we ensured the chatbot could efficiently retrieve relevant responses to user queries. To facilitate seamless user interaction, we created a user-friendly interface for the chatbot, enabling users to ask questions and receive answers with ease.

Finally, we deployed and executed the chatbot on Dell™ PowerEdge™ XE9680 with NVIDIA® H100 server, taking full advantage of its exceptional performance and capabilities. As a result, users could effectively engage with the chatbot to gather valuable information about Dell Technology World 2023, showcasing the server's versatility and power in real-world applications.

| Stable Diffusion

Harnessing the impressive capabilities of Dell™ PowerEdge™ XE9680 with NVIDIA® H100 server, we successfully deployed and executed a pre-trained Stable Diffusion v2 model to generate intriguing images based on text prompts. NVIDIA® H100 server's remarkable performance and advanced features facilitated the efficient handling of the computationally demanding tasks inherent in the Stable Diffusion v2 model.

We deployed the pre-trained model onto Dell™ PowerEdge™ XE9680 with NVIDIA® H100 server, ensuring that it could fully utilize the server's processing power for real-time image generation. With the model in place, users could provide text prompts, prompting the Stable Diffusion v2 model, running on NVIDIA® H100 server, to generate visually appealing and contextually relevant images in response.

This process effectively demonstrated the server's ability to manage complex tasks and large-scale computations, highlighting its value for both creative and technical applications in the field of AI-driven image generation.

EXAMPLE OUTPUT



PROMPT

"a photo of an astronaut with red shirt riding a horse on moon"

OUTPUT

(Stable Diffusion 2)



PROMPT

"future city with artificial general intelligence"

A screenshot of the Dell Technologies World website. The header includes the logo, the event name "MANDALAY BAY, LAS VEGAS | MAY 22-25, 2023", and a "Register Now" button. Below the header, there are navigation links: "What to Expect", "Speakers", "Session Catalog", "Global Partner Summit", "Sponsors", and "Registration Details". The main content area features a large image of a man in a suit, James Cameron, with the text "Interact, inspire and ideate" and "The future of technology belongs to thought leaders, trailblazers and trendsetters like you. Join the Dell Technologies World community of forward thinkers and innovate how we live, work and play." A "Register Now" button is also present. On the right side, there is a chatbot interface titled "DTW Assistant" with a conversation log showing a user asking for session details and the chatbot providing information about James Cameron's session.

Custom LLM chatbot running on Dell™ PowerEdge™ XE 9680
(powered by NVIDIA® H100)

| Chatbot Architecture

Fined tuned Language Models (LLMs) are an incredibly potent tool for enterprises to safeguard their proprietary information, adhere to data sovereignty issues, and enhance their internal and external (customer) effectiveness. The targeted fine-tuning with internal data empowers organizations to unlock the full potential of their data assets, gaining unprecedented insights into their business operations. With the ability to rapidly process vast amounts of data, fined tuned LLMs offer a decisive competitive edge in the dynamic business environment of today. So, if you are looking to remain ahead of the curve and leverage the full potential of your data assets, fined tuned LLMs are the way to go.

Dell™ PowerEdge™ XE9680 is a great platform to fine tune LLMs with your private and proprietary enterprise information rather than rely on bespoke, generic public APIs.

| Use Cases and Benefits

Dell™ PowerEdge™ XE9680 server has several potential use cases in various industries, including healthcare, finance, and retail. The improved performance and energy efficiency of the server can enable businesses to train AI models more quickly and accurately, leading to better predictions and insights. Dell™ PowerEdge™ XE9680 can also help businesses reduce their operational costs by enabling them to use less energy to perform AI workloads.

| Enterprise Knowledge Base

Fined tuned Language Models (LLMs) are an incredibly potent tool for enterprises to safeguard their proprietary information, adhere to data sovereignty issues, and enhance their internal and external (customer) effectiveness. The targeted fine-tuning with internal data empowers organizations to unlock the full potential of their data assets, gaining unprecedented insights into their business operations. With the ability to rapidly process vast amounts of data, fined tuned LLMs offer a decisive competitive edge in the dynamic business environment of today. So, if you are looking to remain ahead of the curve and leverage the full potential of your data assets, fined tuned LLMs are the way to go.

Dell™ PowerEdge™ XE9680 is a great platform to fine tune LLMs with your private and proprietary enterprise information rather than rely on bespoke, generic public APIs.

Chatbot Application

CHAT WEB APP
(Browser)

CHAT SERVER (PORT 9000)
(HTML/CSS/JS)

DIALOG MANAGEMENT
(PORT 5005)
(Rasa Server)

ACTION HANDLER
(PORT 5005)
(Rasa Action Endpoint)

API SERVER (PORT 8000)

CONTEXT PROCESSING + TEXT GENERATION
(LangChain)

EMBEDDINGS
(HuggingFace)

VECTORDB
(ChromaDB)

LLM
(StabilityLM)

| Conclusion

The Dell™ PowerEdge™ XE9680 server is poised to become a vital tool for training large deep learning models due to its exceptional performance and capabilities. The significance of minimizing latency in the training process cannot be overstated, especially considering the substantial amount of data required for training such models. This is where the network interface speed plays a crucial role. The inclusion of the super-fast Broadcom PEX89000 family of PCIe Gen 5.0 switches which provide the flexibility to create a wide range of systems, including simple PCIe connections or complex, high-performance, low-latency, scalable, and cost-effective PCIe fabrics for composable hyper-scale compute systems supporting a variety of ML/AI and Server/Storage applications. 100G Ethernet interfaces in the server enables enterprises to efficiently move massive volumes of data across the network, ensuring swift preparation for training. The combination of the server's 8 NVIDIA® H100 GPUs and the high-speed Ethernet interfaces results in a compounding effect, maximizing the value and efficiency of the investment. Enterprises can fully leverage the Dell™ PowerEdge™ XE9680 server's power and capabilities to propel their Generative AI initiatives forward, unlocking new possibilities to fast track business transformation.

| About Scalers AI™

Scalers AI™ specializes in creating end-to-end artificial intelligence (AI) solutions to fast track industry transformation across a wide range of industries, including retail, smart cities, manufacturing, insurance and healthcare. Scalers AI™ industry offering include predictive analytics, generative AI chatbots, stable diffusion, image and speech recognition, and natural language processing. As a full stack AI solutions company with solutions ranging from the cloud to the edge, our customers often need versatile common off the shelf (COTS) hardware that works well across a range of workloads.

| Fast track development with access to the solution code

Save development time with the sample code.

As part of this effort Scalers AI™ is making the solution code available.



Reach out to your Dell™ representative or contact Scalers AI™ at contact@scalers.ai for access.



This project was commissioned by Dell Technologies™ and conducted by Scalers AI, Inc.
Scalers AI™ and Scalers AI™ logos are trademarks of Scalers AI, Inc.
Copyright © 2023 Scalers AI, Inc.
All rights reserved.
Other trademarks are the property of their respective owners.