

## Unlock high-value insights faster with GenAI

Rapidly deploy a full-stack solution for Generative Artificial Intelligence (GenAI) large language model inferencing

### Increase productivity and insights

This joint architecture delivers a modular and flexible design supporting a multitude of use cases and computational requirements. Components can be mixed and matched, and independently scaled depending on your application needs.

Leverage the potential of GenAI for key use cases:

**Content Creation:**

Across marketing, sales, back-office operations, and more

**Digital Assistants:**

A tailored self-service experience in almost any language

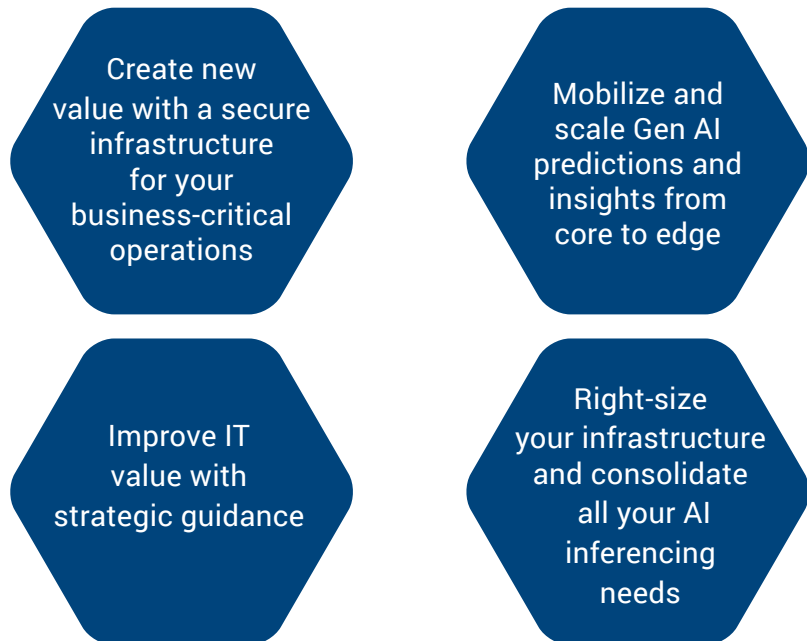
**Code Development:**

Assistance in generating initial code drafts for example

Generate higher quality, faster time-to-value predictions and outputs, while accelerating decision-making with a powerful GenAI solution from Dell Technologies and NVIDIA. This jointly engineered solution addresses Inferencing challenges such as latency, responsiveness, and computational demands helping transform enterprise data into high value, smarter outcomes.

With innovative technologies, comprehensive professional services, and a broad partner ecosystem, your organization can accelerate GenAI at an enterprise-wide scale. Now IT organizations, data scientists, and AI devOps can easily deliver a modular and scalable platform for GenAI and LLM inferencing.

**Over 340k engineering hours spent on design, development and validation on GenAI solutions<sup>1</sup>**



### Reduce time-to-results with a proven solution

Quickly build on-premises infrastructure for your application needs using pre-tested solutions made to simplify adoption. With a proven blueprint, you can reduce risk while avoiding design, planning, and adoption pitfalls.

## Learn more

- [See Design Guide](#)
- [The Power of the PowerEdge XE9680 Server on the GPT-J Model from MLPerf™ Inference](#)

## What is Inferencing?

Inferencing in AI refers to the process of using a trained model to generate predictions, make decisions, or produce outputs based on input data. It involves applying the learned knowledge and patterns acquired during the model's training phase to new, unseen data.

During inferencing, the trained model takes input data and processes it through its computational algorithms or neural network architecture to produce an output or prediction. The model applies its learned parameters, weights, or rules to transform the input data into meaningful information or actions.

Inferencing is a crucial stage in the lifecycle of an AI system. After training a model on labeled or unlabeled data to learn patterns and correlations, inferencing allows the model to generalize its knowledge and make predictions or generate responses on real-world or unseen data.

## Deliver outcomes faster with our help

Dell Services experts help you realize the value of GenAI for your data more quickly with a portfolio of services to assist you at every stage of your GenAI journey:

- **Strategize** - build your roadmap to achieve the innovation objectives of your IT and business stakeholders
- **Implement** - establish your platform, leveraging Dell Validated Designs to implement GenAI inferencing hardware and software
- **Adopt** - accelerate the value of your GenAI use cases by implementing a pre-trained inferencing model
- **Scale** - Manage your GenAI innovation portfolio with resident technical experts and training offers to develop the skills of your team

## Technical Specifications

The Validated Design configurations are based on the newest, AI-acceleration-optimized Dell [PowerEdge XE](#) and [rack servers](#), leveraging the latest NVIDIA GPUs and NVIDIA AI Enterprise, with Triton Inference Server and the NeMo framework. Fast, ample data lake storage for Generative AI and large language models is provided by [Dell PowerScale](#) all-flash or hybrid storage arrays.

Compute	Networking	Software
<ul style="list-style-type: none"><li>• PowerEdge XE9680 server equipped with eight NVIDIA H100 SXM GPUs with NVSwitch</li><li>• PowerEdge XE8640 server equipped with four NVIDIA H100 SXM GPUs with NVLink</li><li>• PowerEdge R760xa servers supporting up to four NVIDIA H100 PCIe GPUs with NVLink Bridge</li><li>• PowerEdge R760xa servers supporting up to four NVIDIA L40S PCIe GPUs</li><li>• Management: PowerEdge R660 servers</li></ul>	<p>NVIDIA Networking, Dell PowerSwitch S5232F-ON or S5248F-ON</p> <p><b>Storage</b></p> <p>Supported by Dell PowerScale storage</p>	<p>Dell OpenManage Enterprise, Power Manager, CloudIQ. NVIDIA AI Enterprise with NeMo Framework for LLMs and Triton Inference Server; NVIDIA Base Command Manager Essentials</p>

## Dell Technologies and NVIDIA

Dell Technologies and NVIDIA work together to enable and accelerate Generative AI adoption, deliver engineering-validated hardware and software to accelerate AI, ML and DL workloads to meet customer needs across all businesses and verticals. With this Validated Design for LLM inferencing, you can accelerate your digital transformation with real-time data that improves key decision-making at-scale, with solutions optimized for fastest time to value from your AI initiatives.



[Learn more](#) about Dell solutions



[Contact](#) a Dell Technologies Expert



[View more](#) resources



[Join the conversation](#) with #HashTag

<sup>1</sup> Based on internal analysis, October 2023

© 2024 Dell Inc. or its subsidiaries. All Rights Reserved. Dell and other trademarks are trademarks of Dell Inc. or its subsidiaries.