

# PT study shows that Dell PowerEdge XE9680 servers with AMD Instinct MI300X GPUs make a strong platform for GenAI

*In hands-on testing, a Dell PowerEdge XE9680 server with AMD Instinct MI300X Accelerators supported up to 136 simultaneous chatbot users*

ROUND ROCK, TX, UNITED STATES, May 13, 2025 /EINPresswire.com/ -- As more and more organizations are looking to use their own data to support GenAI chatbots to answer customer and employee questions, the need for real-world data about appropriate solutions increases. Principled Technologies (PT) tested a Dell PowerEdge XE9680 servers with AMD Instinct MI300X Accelerators to show the kind of chatbot performance an organization might expect.

According to the report, "Organizations have many ways to configure servers to meet their AI needs, but specific data to help them plan may be scarce. To assist these businesses, we used the PTChatterly service to showcase the AI chatbot performance of a Dell™ PowerEdge™ XE9680 server powered by AMD Instinct™ MI300X Accelerators with an industry-leading 192 GB of high bandwidth memory (HBM3) in two different use cases: with eight accelerators and then using just four of the accelerators, for those who wish to use their accelerators for multiple workloads. In fact, the huge memory size of the AMD Instinct MI300X



## Principled Technologies®

A Principled Technologies report: Hands-on testing. Real-world results.

### Dell PowerEdge XE9680 servers with AMD Instinct MI300X Accelerators: the power to host GenAI with Llama 3.1 405B LLMs

As generative AI (GenAI) adoption continues, organizations are increasingly using their specific in-house data in conjunction with large language models (LLMs) to provide AI chatbots to meet customer needs. In a recent Gartner poll, 38% of executives cited customer experience as the primary purpose for using GenAI—which means you need a hardware solution that's powerful enough to support the number of chatbot users you expect to serve and give them quick, accurate responses.

Organizations have many ways to configure servers to meet their AI needs, but specific data to help them plan may be scarce. To assist these businesses, we used the PTChatterly service to showcase the AI chatbot performance of a Dell™ PowerEdge™ XE9680 server powered by AMD Instinct™ MI300X Accelerators with an industry-leading 192 GB of high bandwidth memory (HBM3) in two different use cases: with eight accelerators and then using just four of the accelerators, for those who wish to use their accelerators for multiple workloads. In fact, the huge memory size of the AMD Instinct MI300X Accelerators is what makes it possible to run a very large LLM on only four accelerators; it has the most memory of any available GPU as of this writing. For testing, we used a very large Llama 3.1 405B LLM (which has 405 billion parameters) and FP8 precision, which requires very fast accelerators that have HBM. With this LLM, which we augmented with in-house data, we were able to quantify the number of supported users engaging in longer conversations to mirror what real-world users would experience.

Both Dell PowerEdge XE9680 server use cases supported substantial numbers of chatbot users: using four accelerators, the server supported 72 simultaneous users while leaving room for other workloads, while the standard eight accelerator-use case supported 136 simultaneous users. For organizations looking to support in-house chatbots using their own data and to use a high-precision very large LLM, these results show what's possible as you allocate resources for your new AI infrastructure. To help organizations understand how much a GenAI project might cost, we also calculated expected five-year TCO costs.

When your business requires the precision of a very large LLM, Dell PowerEdge XE9680 servers with AMD Instinct MI300X Accelerators offer the resources you need to make your in-house GenAI project a success.

**Support up to 72 simultaneous chatbot users**  
using only 4 accelerators; the other 4 are free for other applications

**Support up to 136 simultaneous chatbot users**  
with 8 accelerators

**Support 816 simultaneous users in a full rack for just \$7.6M over 5 years**

"By 2027, more than 50% of the GenAI models that enterprises use will be specific to either an industry or business function — up from approximately 1% in 2023."  
—Gartner<sup>1</sup>

Dell PowerEdge XE9680 servers with AMD Instinct MI300X Accelerators: the power to host GenAI with Llama 3.1 405B LLMs May 2025

Dell PowerEdge XE9680 servers with AMD Instinct MI300X Accelerators: the power to host GenAI with Llama 3.1 405B LLMs

Accelerators is what makes it possible to run a very large LLM on only four accelerators; it has the most memory of any available GPU as of this writing.”

PT used the very large Llama 3.1 405B LLM and FP8 precision, which requires very fast accelerators that have HBM. They augmented the LLM with in-house data and used the PTChatterly service to determine how many users could engage in long conversations. This approach mirrors what real-world LLM users would likely experience.

The report continues, “Both Dell PowerEdge XE9680 server use cases supported substantial numbers of chatbot users: using four accelerators, the server supported 72 simultaneous users while leaving room for other workloads, while the standard eight accelerator-use case supported 136 simultaneous users. For organizations looking to support in-house chatbots using their own data and to use a high-precision very large LLM, these results show what’s possible as you allocate resources for your new AI infrastructure. To help organizations understand how much a GenAI project might cost, we also calculated expected five-year TCO costs.

When your business requires the precision of a very large LLM, Dell PowerEdge XE9680 servers with AMD Instinct MI300X Accelerators offer the resources you need to make your in-house GenAI project a success.”

To learn more, read the full report at <https://facts.pt/ovSGH4f> or see the infographic at <https://facts.pt/6uHzsvR>.

About Principled Technologies, Inc.

Principled Technologies, Inc. is the leading provider of technology marketing and learning & development services.

Principled Technologies, Inc. is located in Durham, North Carolina, USA. For more information, please visit [www.principledtechnologies.com](http://www.principledtechnologies.com).

Sharon Horton

Principled Technologies, Inc.

[press@principledtechnologies.com](mailto:press@principledtechnologies.com)

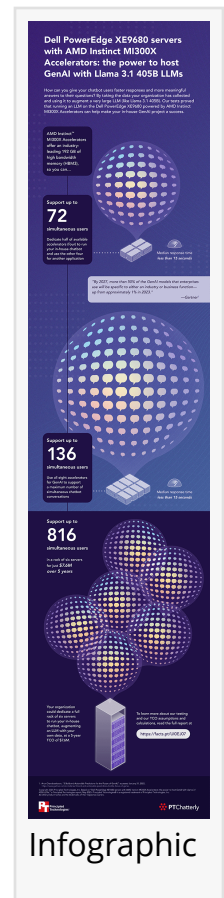
Visit us on social media:

[LinkedIn](#)

[Facebook](#)

[YouTube](#)

[X](#)



EIN Presswire's priority is source transparency. We do not allow opaque clients, and our editors try to be careful about weeding out false and misleading content. As a user, if you see something we have missed, please do bring it to our attention. Your help is welcome. EIN Presswire, Everyone's Internet News Presswire™, tries to define some of the boundaries that are reasonable in today's world. Please see our Editorial Guidelines for more information.

© 1995-2025 Newsmatics Inc. All Right Reserved.