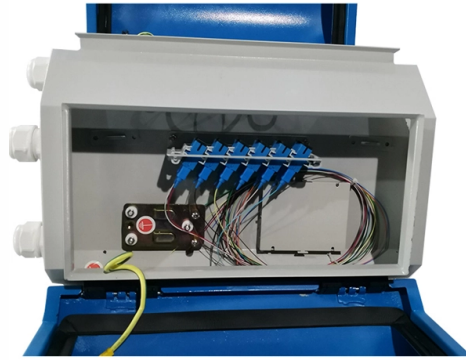


AI Server Performance Comparison



Overview

This study presents a systematic, empirical comparison of GPU- and NPU-based server platforms across key AI inference domains: text-to-text, text-to-image, multimodal understanding, and object detection. Compare specifications, pricing, support, and real-world performance to select the optimal infrastructure for your AI workloads. The enterprise AI server market reached \$245 billion in 2025 (ABI Research) and is projected to grow at 18% CAGR through 2030. The transition from NVIDIA Hopper. Dell's AI Factory platform (e. PowerEdge XE97xx/XE9712) provides high-density rack-scale clusters (72 GPUs per rack with NVLink, ~30× LLM inference speed-up and up to 25× energy efficiency advantage over prior-gen systems ()) with both liquid- and air-cooled options. HPE's Private Cloud AI. Live now: AA-AgentPerf is now open for submissions of configurations for benchmarking. The models supported at launch are gpt-oss-120b and DeepSeek V3. We'll be publishing results on a rolling basis. AA-AgentPerf has been shaped by our work with inference providers and engagement with AI. Artificial Intelligence (AI) server manufacturers have experienced surging demand as data center operators require significantly more computing power than before the advent of ChatGPT and other Generative Artificial Intelligence (Gen AI) tools. A year. This paper is a substantially extended version of a preliminary abstract presented at the 19th International Conference on Innovative Computing, Information and Control (ICICIC 2025), Kitakyushu, Japan, 29 August 2025.

Article Content

AI Server Comparison: Dell vs HPE vs Supermicro vs Lenovo | SLYD

Compare specifications, pricing, support, and real-world performance to select the optimal infrastructure for your AI workloads. The enterprise AI server market reached \$245 billion in 2025 (ABI Research)

Performance and Efficiency Gains of NPU-Based

This study presents a systematic, empirical comparison of GPU- and NPU-based server platforms across key AI inference domains: text-to-text, text-to-

Comparison of AI Models across Intelligence,

Comparison and analysis of AI models across key performance metrics including quality, price, output speed, latency, context window & others.

On-Prem AI Infrastructure: Comparing Dell, HPE, & More

Compare on-prem AI infrastructure from Dell, HPE, Lenovo, Supermicro & Cisco. Analyze NVIDIA GB200/GB300 NVL72 and Blackwell Ultra hardware specs,

Computers, Monitors & Technology Solutions | Dell USA

Dell provides technology solutions, services & support. Buy Laptops, Touch Screen PCs, Desktops, Servers, Storage, Monitors, Gaming & Accessories

AI Hardware Index

Track AI hardware prices across 24+ vendors. Daily updated pricing for GPU servers, workstations, and accelerators from \$109 to \$500k+.

Microsoft

Service Level Agreements (SLA) for Online Services. The Service Level Agreements (SLA) describe Microsoft's commitments for uptime and connectivity for Microsoft Online Services

AI M.2 Accelerator: Hailo-8 AI Module | Superior Edge Performance

Hailo-8 M.2 AI Module is an AI accelerator compatible with NGFF M.2 form factor. Based on a 26 TOPS processor with high power efficiency.

AI Server vs Traditional Server Comparison | Lenovo US

Compare AI servers and traditional servers across architecture, performance, scalability, and workloads. Understand which server type fits your business needs.

AISBench: an performance benchmark for AI server systems

Artificial intelligence (AI) server systems, including AI servers and AI server clusters, are widely utilized in AI applications. The performance of an AI server system determines the

Gartner | Delivering Actionable, Objective Insight to

Gartner provides actionable insights, guidance, and tools that enable faster, smarter decisions and stronger performance on an organization's mission-critical priorities.

GPU Server Benchmarks 2026: Training Speed & Cost Per Token

Check GPU Server Benchmarks with insights into AI training speed, cost per token, etc. Compare performance and scalability to choose the right infrastructure.

Mac Mini M4 AI Server: Local LLM + Agent Setup (2026)

Turn your Mac Mini M4 into a local AI server. Ollama for LLMs, OpenClaw for AI agents, Claude Code for dev workflows. Hardware tiers \$599-\$2,000 tested.

Top AI Server Companies & How to Compare Them (2025)

With numerous vendors vying for dominance, choosing the right AI server can be daunting. Understanding the key players and evaluation criteria is essential for making informed decisions.

AI Hardware Benchmarking & Performance Analysis

AI Hardware Benchmarking & Performance Analysis We measure real-world performance of AI accelerator systems during language model inference. For

Asus to compete for AI servers; not concerned about

Compared to the impressive performance of its gaming notebooks and boards, server business tends to not receive much attention from the market

AI Performance Comparison of Processors

Detailed comparison of AI performance across CPUs, GPUs, and dedicated AI accelerators. Covers key specs like FP64/FP32/FP16/FP8 FLOPS, INT16/INT8/INT4 TOPS, memory bandwidth, and capacity.

What is an AI server? Why artificial intelligence needs

AI servers are playing an increasingly pivotal role as enterprises across industries race to implement sophisticated gen AI tools and AI agents.

AISBench: an performance benchmark for AI server systems

In response to this need, this paper introduces AISBench, a performance benchmark for AI server systems. AISBench comprises standardized rules and a test toolkit that has been agreed

Claude vs ChatGPT: Why Users Are Switching and Which AI Is Better

Claude vs ChatGPT comparison in 2026—context limits, coding benchmarks, and creative writing quality to help choose the best AI chatbot or alternative.

Self-Hosting AI Models: Hardware Requirements, Model Selection,

A practical guide to self-hosting AI models on your own infrastructure. Covers hardware requirements, VRAM and quantisation, model selection for 2026, cost comparisons with cloud APIs,

Huawei Ascend 910C AI chip cluster dubbed

Huawei Ascend P910C AI chip cluster called CloudMatrix, will have superior performance to NVIDIA GB200 NVL72 AI servers, but with more power

Feature Comparison Of Windows Server 2025 Vs 2022

Feature Comparison of Windows Server 2025 Vs 2022 Vs 2019 Windows Server 2025 offers enhanced multi-layer security, Hyper-V, AI, and

Top Five AI Server Companies for Data Centers and

We evaluated server manufacturers based on performance, partner channels, workload optimization, environmental impact, future-readiness, and

AI Hardware Benchmarking & Performance Analysis

Comprehensive benchmarking of AI accelerator systems for language model inference. We test different chip configurations, inference software (vLLM vs.

Best AI Agent Memory Frameworks in 2026: Compared and Ranked

A comparison of the top AI agent memory frameworks in 2026 — Mem0, Zep, LangMem, Letta, and more — covering architecture, strengths, and enterprise fit.

Desktop-Class Performance For Gaming and AI Development: Razer

Crafted for peak performance, with up to 2.2× faster AI¹, uncompromising power for the most demanding gaming and AI workloads, and a brighter 18" dual-mode display.

Contact Us

For more information, pricing, or custom solutions, please contact us:

Website: <https://pvprojekt.com.pl>

Email: contact@pvprojekt.com.pl

Phone: +48 512 897 346

Address: ul. Tęczowa 17, 61-001 Poznań, Greater Poland Voivodeship, Poland

This document is for informational purposes only. Specifications subject to change without notice.

