

Memory in AI Servers



Overview

This guide provides a practical, data-driven framework to determine RAM requirements for AI workloads, including AI server memory planning, GPU RAM requirements, and large-scale LLM infrastructure design. AI workloads differ fundamentally from traditional enterprise. As a trusted U. Micron Technology has announced the sampling of its new 256-GB DDR5 registered dual in-line memory module (RDIMM) to key server ecosystem partners, targeting next-generation AI and. Local AI inference means running an already trained model on your own server. The model is not trained from scratch; it is used to answer questions, analyze documents, generate text, recognize speech, classify tickets, search a knowledge base or process images. SK Hynix officially begins mass production of its 192GB SOCAM M2 memory, “establishing a new benchmark for memory performance for AI servers. We will explore their architectural differences, their respective strengths and weaknesses in handling various AI tasks, and how to optimally configure them.

Article Content

SQL Server 2025 In-Memory: New Cleanup Features & SQLBits 2026

SQL Server 2025 finally allows dropping In-Memory filegroups, a breakthrough analyzed here alongside expert performance insights from SQLBits 2026.

Riding the AI Supercycle: Navigating the 2026 Memory

The memory and storage market has entered a multi year, AI driven supercycle, as suppliers shift aggressively toward HBM and server class DRAM

3 AI Memory Stocks to Watch in 2026 (Besides Micron)

AI memory stocks are booming as demand outpaces supply. Here are three worth watching in 2026 - beyond the obvious Micron trade.

The Hidden Crisis in AI Right Now: Server Memory Is In

AI teams are running into a problem the market isn't built to solve: server memory prices are up more than 300 percent this year thanks to supply shortages and

Memory & Flash Crisis: March 2026 Update

The global memory market entered 2026 in a state of structural supply constraint. AI infrastructure demand has reallocated

Memory | Awesome MCP Servers

Knowledge graph-based persistent memory system Knowledge Graph Memory Server
A basic implementation of persistent memory using a local knowledge

High-Bandwidth Memory Solution for AI Servers

A memory module is set to power AI servers with higher speed, lower energy use, and smoother performance for large AI workloads.

Micron Samples 256-GB DDR5 RDIMM for AI Servers

Micron is now sampling its new 256-GB DDR5 registered dual in-line memory module (RDIMM) for AI and HPC platforms.

2026 Market Outlook: SK hynix's HBM to Fuel AI Memory Boom

In its 2026 semiconductor market outlook, SK hynix forecasts that demand for its HBM3E and HBM4 products will fuel the AI memory supercycle.

Nvidia Collaborates with Major Memory Makers on New SOCAMM

This new modular memory form factor is designed to unlock the full potential of AI platforms and has been developed exclusively for Nvidia's Grace Blackwell platform. SOCAMM, or

Local AI Inference Server 2026: How to Choose GPU, CPU and VRAM

Learn how to size VRAM, CPU, PCIe lanes, memory, power and cooling for a reliable local AI inference server. A practical guide for avoiding GPU overkill and planning around real workloads

AI Server Bottlenecks: Memory, Packaging & Power Limits

Explore AI server bottlenecks and how memory, advanced packaging, and power constraints impact data center growth and investment strategy.

Nvidia shift to smartphone-style memory could double

Nvidia's move to use smartphone-style memory chips in its artificial intelligence servers could cause server-memory prices to double by late 2026,

How Much RAM Do AI Workloads Really Need?

Learn how much RAM for AI workloads your organization really needs. A detailed guide for CTOs and AI teams covering AI server memory, GPU

AI Server Industry Analysis

In 2025, global AI chips focus on high-end HBM memory; NVIDIA's new Blackwell platform drives growth, amid geopolitical limits and steady AI

Samsung targets Q4 rollout of next-gen CXL memory for AI servers

Samsung Electronics is preparing to mass-produce next-generation Compute Express Link (CXL) memory modules as early as the fourth quarter, to capture rising demand for flexible AI data

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.

Winbond Electronics stock (TW0002344009): Memory chip maker rides AI ...

Winbond Electronics shares gain as Taiwan's memory sector surges 309% year-over-year in April 2026, driven by AI server demand and DRAM price strength.

2026 Memory Price Forecast | TrendForce

TrendForce forecasts a structural memory price hike by 2026. Strong AI server demand and profit-first supplier strategies lead to capacity shifting

Unihost: Choosing the Right Server Specs for AI Workloads - CPU vs

A comprehensive guide to selecting the right server specifications (CPU, GPU, RAM) for AI workloads, covering deep learning, inference, and data processing."

AI Memory Requirements: Why Memory — Not

CXL memory contributes to this infrastructure revolution by providing a flexible way to expand and deploy memory — one of the most critical resources

AI Server Demand to Drive Memory Contract Price Increases in 2Q26

TrendForce's latest memory pricing survey reveals that DRAM suppliers are reallocating capacity toward HBM and server applications in 2Q26, while implementing catch-up pricing to narrow

AI Memory: Enabling The Next Era Of High

Investment in AI infrastructure is accelerating, with hyperscale data centers expanding their AI server capacity to accommodate increasingly complex

Contact Us

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

Website: <https://www.buglerdental.co.za>

Email: sales@buglerdental.co.za

Phone: +27 71 549 2836

Address: 22 Impala Crescent, Waterfall Business Estate, Midrand, 1685, South Africa

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

