

AI Infrastructure Spending

Researched by Hey Lefty

Automated research briefings on topics you choose — heylefty.com

TL;DR

The AI physical infrastructure race is shifting from grid-negotiation delays to self-sufficient microgrid deployments Oracle Bloom Energy Fuel Cell Partnership and radical hardware redesigns AMD Helios Platform. To bypass severe power transmission bottlenecks, major operators are deploying gigawatt-scale, on-site fuel cell networks to power their newest supercomputing clusters Oracle Bloom Energy Fuel Cell Partnership. Simultaneously, chipmakers are rewriting the physical constraints of the data center, introducing double-wide server racks AMD Helios Platform and co-packaged optical switches to handle the extreme thermal and electrical demands of next-generation silicon NVIDIA Vera Rubin Platform.

On-Site Microgrid Independence

Hyperscalers are abandoning traditional utility grid negotiations entirely in favor of massive, self-sufficient on-site microgrids to bypass electrical transmission delays.

"This distributed, on-site generation strategy allows Oracle to bypass local utility transmission constraints in high-demand areas." — Bloom Energy Press Release

"By utilizing Bloom's combustion-free fuel cells, Project Jupiter will reduce NOx emissions... and use a negligible amount of water..." — Oracle Press Release

This pattern matters because waiting years for traditional grid interconnections is no longer a viable strategy for tech firms racing to deploy frontier AI systems. By committing to up to 2.8 gigawatts of fuel cells for its rapid AI buildout, Oracle is accepting the massive capital burden of building its own clean power generation directly on-site to bring its largest clusters online years ahead of schedule Oracle Bloom Energy Fuel Cell Partnership.

What to watch: Watch whether other Stargate-affiliated data center designs adopt similar off-grid fuel cell architectures to bypass local utility approval timelines.

Physical Redesign of the Rack-Scale Layout

The physical constraints of next-generation silicon are forcing a fundamental redesign of data center layouts, shifting the industry standard toward double-wide rack architectures.

"...the Helios platform introduces a major physical design shift: the transition to a double-wide rack architecture." — ServeTheHome

This pattern matters because competing with dominant hardware providers now requires optimizing the physical space and power delivery of the entire server cabinet, not just the individual processor. By leveraging its acquisition of ZT Systems, AMD is engineering massive, integrated rack systems designed to pack 72

GPUs into a single double-wide footprint, a strategy that has already garnered high-profile endorsement from OpenAI AMD Helios Platform.

What to watch: Watch whether competing rack-scale integrators adopt double-wide standards to accommodate the extreme thermal and power demands of next-generation silicon.

The Optical Transition in Scale-Out Networks

Network bottlenecks in massive-scale clusters are driving a rapid transition toward co-packaged optics to prevent power delivery from choking compute scaling.

"By integrating optical components directly with the switching silicon, Spectrum-X Ethernet Photonics delivers... 5x better power efficiency compared to traditional networks using conventional transceivers." — NVIDIA Press Release

This pattern matters because as AI workloads transition to complex tasks requiring multi-step reasoning, the network interconnects themselves threaten to consume an unsustainable share of a data center's power budget. By integrating optical elements directly onto the silicon, hardware providers can dramatically improve cluster uptime and free up vital electricity for the actual GPUs NVIDIA Vera Rubin Platform.

What to watch: Watch how quickly early adopters like CoreWeave and Oracle transition their production clusters to co-packaged optics once volume shipments begin in the fall.

What surprised us

- **Oracle is spending far beyond its cash generation to secure power.** Oracle's capital expenditures have ballooned to \$18.64 billion, dramatically outstripping its \$7.15 billion in operating cash flow, highlighting the extreme financial strain hyperscalers are accepting to build out AI microgrids Oracle Bloom Energy Fuel Cell Partnership.
- **AMD is redesigning the physical layout of the data center.** Rather than fitting new chips into standard server racks, AMD's new Helios platform forces a shift to double-wide racks to manage the power density of its upcoming Instinct accelerators AMD Helios Platform.
- **Co-packaged optics are arriving in volume production sooner than expected.** Long viewed as a future-looking lab technology, co-packaged optics are debuting in Nvidia's Spectrum-X Ethernet Photonics platform as a production-ready solution to deliver massive power efficiency gains NVIDIA Vera Rubin Platform.

Open threads worth a vote

- Track Oracle Q4 FY2026 Earnings and Capex Guidance

Appendix: Findings

Oracle and Bloom Energy Partner for Up to 2.8 GW of On-Site Fuel Cells to Power AI Microgrid Campuses

Oracle and Bloom Energy Partner for Up to 2.8 GW of On-Site Fuel Cells to Power AI Microgrid Campuses

Oracle and Bloom Energy (BE) have announced a massive expansion of their strategic partnership, under which Oracle intends to procure up to 2.8 gigawatts (GW) of Bloom's solid oxide fuel cell systems to power its rapid AI and cloud computing infrastructure buildout. An initial 1.2 GW of capacity has already been contracted, with deployments currently underway and continuing into 2027.

This distributed, on-site generation strategy allows Oracle to bypass local utility transmission constraints in high-demand areas. In just 55 days, Bloom previously delivered a fully operational fuel cell system to Oracle, demonstrating the speed-to-power advantage of modular fuel cells over traditional grid interconnections.

Overhauling Project Jupiter's Power Design

As part of this partnership, Oracle and developer BorderPlex Digital Assets announced on April 27, 2026, that "Project Jupiter"—a massive \$165 billion AI data center campus in Santa Teresa, Doña Ana County, New Mexico, which is part of the broader Stargate AI initiative with OpenAI—will be fully powered by up to 2.45 GW of Bloom's fuel cells. This replaces previously planned natural gas turbines and diesel generators, consolidating the entire facility into a single, self-sufficient microgrid.

By utilizing Bloom's combustion-free fuel cells, Project Jupiter will reduce NOx emissions by roughly 92% compared to gas turbines and use a negligible amount of water, responding to local community concerns about air quality and water scarcity in the southern New Mexico desert. Oracle will bear all energy costs, shielding local residents from electricity rate increases or grid instability.

Sources

- Bloom Energy and Oracle Expand Strategic Partnership to Deploy up to 2.8 GW to Accelerate AI Infrastructure Build-Out
- Oracle, BorderPlex, and Bloom Energy to Power Project Jupiter with Cleaner, Water-Efficient Fuel Cell Technology

AMD Helios Platform and MI400 Series Accelerators Move AI Data Centers to Double-Wide Racks

AMD Helios Platform and MI400 Series Accelerators Move AI Data Centers to Double-Wide Racks

At Computex 2026, AMD unveiled its next-generation AI compute roadmap, highlighted by the **AMD Helios** rack-scale AI platform and its upcoming **Instinct MI400 series** accelerators (including the MI430X, MI440X, and MI455X). Slated for deployment starting in late 2026, the Helios platform introduces a major physical design shift: the transition to a **double-wide rack** architecture.

This double-wide rack configuration is designed to handle the massive power density and scale of next-generation hardware. Forrest Norrod, General Manager of AMD's Data Center Solutions Business Group, confirmed that the Helios platform is a direct output of **ZT Systems**, the rack-scale engineering firm AMD acquired in late 2024.

Technical Specifications and Open Standards

The AMD Helios platform is built to deliver up to 10x more AI compute performance compared to the previous-generation MI355X. Key technical highlights of the Helios and MI400 architecture include:

- **HBM4 Memory Integration:** The platform will scale up to 31TB of HBM4 memory and 1.4PB/s of memory bandwidth per rack. This points to approximately 440GB of memory per GPU across 72 GPUs in a rack, utilizing 12 stacks of 36GB HBM4 memory (which Micron began shipping in early 2026).
- **Next-Gen Processors:** The rack integrates AMD's upcoming 2nm **"Venice" EPYC CPUs**, which feature up to 256 cores, PCIe Gen6 support, and 1.6TB/s of memory bandwidth.
- **Networking and Interconnects:** Helios will adopt open industry standards, utilizing the **AMD Vulcano 800G NIC** and incorporating UALink (Ultra Accelerator Link) and Ultra Ethernet Consortium (UEC) standards.

During the announcement, OpenAI CEO Sam Altman appeared on stage to highlight OpenAI's close collaboration with AMD on the Helios platform, demonstrating growing hyperscaler support for AMD as an alternative to NVIDIA.

Sources

- Not Just for Oreos and Trailers AMD Helios Next-Gen AI Racks Go Double-Wide
- AMD launches Instinct MI350 GPUs, unveils double-wide Helios AI rack-scale system

NVIDIA Launches Vera Rubin Platform and Spectrum-X Ethernet Photonics for Agentic AI Factories

NVIDIA Launches Vera Rubin Platform and Spectrum-X Ethernet Photonics for Agentic AI

Factories

At GTC Taipei / Computex 2026 on May 31, 2026, NVIDIA CEO Jensen Huang announced that the next-generation **NVIDIA Vera Rubin** computing platform is ramping into full production, with volume shipments scheduled to begin in the fall of 2026.

Positioned as the successor to the Grace Blackwell platform, Vera Rubin is designed specifically for "agentic AI" workloads, which involve complex multi-step reasoning, tool use, and code execution. According to NVIDIA, the platform delivers up to **10x more agent throughput** at scale compared to Grace Blackwell.

Vera Rubin Architecture and Component Ecosystem

The Vera Rubin platform represents NVIDIA's most extensive POD-scale system, integrating five purpose-built racks that operate as a single supercomputer. Highly integrated hardware components include:

- **NVIDIA Vera CPU:** Built specifically to handle heavy CPU-bound agentic tasks and reinforcement learning. The Vera CPU will utilize high-bandwidth memory chips from **SK Hynix**.
- **NVIDIA Rubin GPUs** and **Vera Rubin NVL72** rack systems.
- **NVIDIA Spectrum-6 SPX Ethernet** and **Vera BlueField-4 STX storage** racks.
- **NVIDIA Confidential Computing:** Full-stack hardware-level encryption across high-speed NVLink interconnects, providing a trusted execution environment at rack scale.

Spectrum-X Ethernet Photonics: A Co-Packaged Optics Breakthrough

To support scale-out networks for million-GPU AI factories, the platform introduces **NVIDIA Spectrum-X Ethernet Photonics**, the world's first co-packaged-optics (CPO)-based switches with 200Gb/s SerDes.

By integrating optical components directly with the switching silicon, Spectrum-X Ethernet Photonics delivers:

- **5x better power efficiency** compared to traditional networks using conventional transceivers.
- **5x longer AI cluster uptime** and **1.3x faster deployment times** by simplifying physical design and freeing up power for compute.

NVIDIA announced that cloud infrastructure providers **CoreWeave**, **Lambda**, and **Oracle Cloud Infrastructure** are among the first ecosystem partners and early adopters of Spectrum-X Ethernet Photonics.

Sources

- NVIDIA Vera Rubin Ramps Into Full Production to Power Agentic AI Factories Worldwide
- Next Gen Data Center CPU | NVIDIA Vera CPU