**HPE**
**GreenLake**

# Why enterprises need supercomputers to train their most complicated AI models



The complexity of generative AI demands more computing muscle than most enterprises can muster. Enter the supercomputer.

For nearly 50 years, supercomputers have been tackling the biggest and most difficult problems on the planet — like predicting long-term weather patterns, simulating the effects of nuclear fusion, enabling the discovery of lifesaving medicines, and tracing the very origins of the universe.

Increasingly, however, these massively powerful machines are finding a home in the enterprise. And with the emergence of artificial intelligence (AI) as an essential tool for driving business decisions, supercomputers are poised to have a much more prominent role moving forward.

As their name implies, supercomputers are often room-sized beasts consisting of tens of thousands of CPUs and GPUs, millions of gigabytes of memory and storage, and thousands of nodes, each one its own immensely powerful computer. The most advanced versions of these, so-called exascale supercomputers, can handle 1 billion x 1 billion ($10^{18}$) operations per second.

Supercomputers are built to run single, massive applications that require fast communications among thousands of nodes operating in parallel, notes Paolo Faraboschi, vice president of the AI research lab at Hewlett Packard Enterprise. That makes them ideal for training the large language models (LLMs) that serve as the foundation for today's generative AI applications.

"Training a generative AI model is a lot like simulating weather patterns or creating a digital twin of an engine," Faraboschi says. "It's essentially one enormous application with millions of parameters. AI training workloads are exactly where supercomputing hardware shines from the point of view of performance."

## How supercomputers surpass the cloud in raw computing power

While training LLMs is possibly the most visible application for high-performance computing (HPC), it's far from the only one. Supercomputers are already hard at work managing high-frequency trading systems for Wall Street brokers, maintaining digital twins of supply chain networks for manufacturers, and helping pharmaceutical firms quickly evaluate millions of drug candidates in the race to uncover newer and more effective vaccines.

The factors that make these scenarios ideal for supercomputing applications — the need to run millions of tightly coupled processes in parallel — also make them a poor fit for hyperscale cloud solutions, explains Faraboschi. Cloud computing excels at running discrete applications for millions of users, but it's not designed for the high-speed internal communications needed for a single, highly parallel application.

"These applications need to work in lockstep," he says. "Think about a weather simulation. Every node takes a little chunk of the world and analyzes it. But it also needs to communicate with its neighboring nodes because the weather in one area affects the weather in others. Fast node-to-node data exchange is the key to making that work."

Server farms in the cloud typically don't have the same ability to exchange data rapidly and predictably amongst themselves, Faraboschi explains. If one of those servers is slower than the others, the rest must wait for it to catch up.

"The bus can't leave until the last person is on it," he says. "That's where cloud architecture tends to break down."

Faraboschi says it's possible to order a cluster of tightly interconnected and GPU-accelerated servers in the cloud, but the cost can be prohibitive. These servers are in limited supply, and unlike most multitenant cloud architectures, they're typically dedicated to a single customer. That means cloud providers can charge a huge premium for them.

"A lot of organizations start by using the cloud to run small AI training experiments," he adds. "But once they realize how much it will cost to scale the infrastructure to a full training job, they usually bring it back in-house. If I need 50,000 nodes for a week or a month in order to train my LLM, the cloud business model falls apart."

## Liquid cooling makes supercomputers more sustainable

Another surprising advantage supercomputers have over the cloud and most on-premises data centers is energy efficiency. An average corporate data center has a power usage effectiveness (PUE) rating of 1.58 or higher. This means that for every kilowatt consumed by computing resources, another 580 watts (58%) is expended delivering power and cooling to the data center.[1]

Hyperscale cloud providers typically have a PUE of around 1.1 or 1.2, says Faraboschi, making them considerably more efficient than on-premises centers. But supercomputers, especially liquid-cooled models like Frontier, the exascale supercomputer at Oak Ridge National Laboratory, have a PUE of 1.03, meaning they waste only 3% of energy. That makes them three to six times more efficient than typical cloud data centers.[2]

Liquid cooling is what makes these systems more sustainable, Faraboschi notes. Because water dissipates heat more efficiently than air, computing nodes can run at a higher temperature and be packed more closely together. This greater density results in significantly higher performance over a much smaller footprint. And unlike with air-cooled data centers, the waste heat can be used for heating other facilities or even growing greenhouse fruits and vegetables in frigid climates.

"There are many ways to do liquid cooling," says Faraboschi. "But when you do it for real, the cooling fluid gets into every element — memory, GPUs and CPUs, even capacitors, and power converters. That means designing the whole cold-plate infrastructure and liquid delivery systems from scratch."

## Enterprises desperately need AI expertise

Enterprises need to consider several factors in deciding whether they require a supercomputer's power. One is the nature and scale of the workload they're hoping to run. As noted earlier, supercomputers excel at tackling big jobs that involve many disparate computing elements. If you're a global retailer who wants to understand and predict customer behavior, an oil and gas company looking to identify untapped reservoirs of petroleum without the expense of drilling test holes, or an engineering firm that designs commercial aircraft and other large, complex machines, you'll probably use a supercomputer to do it.

Organizations looking to create a competitive advantage by building and training their own AI models in house will also want to employ a supercomputer to do the job more efficiently, says Faraboschi.

"You need to understand where in the pyramid of AI players you want to be," he says. "Can you afford to be at the top of the pyramid, pretraining a model from scratch? Then, that's a supercomputing problem."

Another factor is how long you might need access to massive computing resources. If you're maintaining a digital twin of an engine or a supply chain, or you need to do regular, periodic forecasting, your need for supercomputing is constant and you probably will want this resource on-premises. But if you're looking to train an LLM once (using your proprietary data before deploying it to the field), on-demand HPC as a service would likely be a more efficient way to go than purchasing a supercomputer outright.

"The difference is really economics and whether you have an IT group with expertise in HPC and AI," says Faraboschi. "If you don't, and you only need to train your LLM once or cyclically, then you probably want to go with supercomputing as a service."

In all use cases, however, expertise is key. It's hard to keep any application running for a month without crashing, especially an enormously complex application running across thousands of nodes. Today, only 10% of enterprises have the expertise and resources in-house to train their own AI models.[3] The rest are using publicly available LLMs, or they're bringing in outside expertise from systems integrators.

"The complexity of training a large AI model is very, very high," warns Faraboschi. In many cases, we've seen customers move into AI and then realize that it's not their core business, and their money is better spent doing what they're supposed to be doing, whether retail, finance, or engineering. That's when they usually decide to leave the expertise to us or their other system integration partners.

[1] Why has PUE Remained Flat for So Long After Years of Progress?

[2] olcf.ornl.gov/wp-content/uploads/2022-OLCF-User-Meetig-Overview-of-Frontier-Whitt.pdf

[3] "2023: The State of Generative AI in the Enterprise," Menlo Ventures, Nov 13, 2023

## Why the AI future is going to be "super"

Just over half of enterprises use some form of AI, and together they plan to spend some $70 billion on AI technology and services in 2024.[4] Today, generative AI represents a relatively small portion of that investment, but it's likely to grow steadily over time.

Organizations looking to capitalize on the potential of generative AI first need to get their data houses in order, says Faraboschi. Without good, clean, reliable data, any predictions an AI model makes won't be worth much. Then, they need to decide how big a role AI will play inside their companies and how deeply they intend to commit to procuring the necessary hardware and expertise.

"Organizations are only right now starting to appreciate the complexity of the AI environment," he adds. "You need to have efficient resilience mechanisms because some element of the infrastructure is going to fail every few hours, and you cannot afford to restart a large AI training job from scratch. You need to provision the networking and the storage in a way that maximizes the value of this expensive hardware. The whole system design space is very complex, which is where system integrators can help."

Faraboschi says the size and scale of today's AI models demand high-performance computing solutions. Ten or fifteen years ago, you could use a supercomputer to train a million of the largest AI models in a week, he notes. Today, incredibly powerful supercomputers like Frontier would need a week to train just one of the largest LLMs. In other words, as Faraboschi says, "People need to realize that training large AI models is now a supercomputing problem."

[4] Ibid

## Learn more at
HPE.com/AI

Visit **HPE GreenLake**

**Chat now**

**Hewlett Packard Enterprise**