

Microsoft used tens of thousands of chips for OpenAI supercomputer

When Microsoft invested \$1 billion in OpenAI in 2019, it agreed to build a massive, cutting-edge supercomputer for the artificial intelligence research startup. The only problem: Microsoft didn't have anything like what OpenAI needed and wasn't totally sure it could build something that big in its Azure cloud service without it breaking.

OpenAI was trying to train an increasingly large set of artificial intelligence programs called models, which were ingesting greater volumes of data and learning more and more parameters, the variables the AI system has sussed out through training and retraining.

That meant OpenAI needed access to powerful cloud computing services for long periods of time.

To meet that challenge, Microsoft had to find ways to string together tens of thousands of Nvidia's A100 graphics chips — the workhorse for training AI models — and change how it positions servers on racks to prevent power outages.

Scott Guthrie, the Microsoft executive vice president who oversees cloud and AI, wouldn't give a specific cost for the project, but said "it's probably larger" than several hundred million dollars.

"We built a system architecture that could operate and be reliable at a very large scale. That's what resulted in ChatGPT being possible," said Nidhi Chappell, Microsoft general manager of Azure AI infrastructure.

"That's one model that came out of it. There's going to be many, many

others."

The technology allowed OpenAI to release ChatGPT, the viral chatbot that attracted more than 1 million users within days of going public in November and is now getting pulled into other companies' business models, from those run by billionaire hedge fund founder Ken Griffin to food-delivery service Instacart. As generative AI tools such as ChatGPT gain interest from businesses and consumers, more pressure will be put on cloud services providers like Microsoft, Amazon and Google to ensure their data centers can provide the enormous computing power needed.

Now Microsoft uses that same set of resources it built for OpenAI to train and run its own large artificial intelligence models, including the new Bing search bot introduced last month. It also sells the system to other customers. The software giant is already at work on the next generation of the AI supercomputer, part of an expanded deal with OpenAI in which Microsoft added \$10 billion to its investment.

"We didn't build them a custom thing — it started off as a custom thing, but we always built it in a way to generalize it so that anyone that wants to train a large language model can leverage the same improvements," said Guthrie in an interview.

"That's really helped us become a better cloud for AI broadly."

Training a massive AI model requires a large pool of connected graphics processing units in one place like the AI supercomputer Microsoft assembled.

Once a model is in use, answering all the queries users pose — called inference — requires a slightly different set up. Microsoft also deploys graphics chips for inference but those processors — hundreds of thousands of them — are geographically dispersed throughout the company's more

than 60 regions of data centers. Now the company is adding the latest Nvidia graphics chip for AI workloads — the H100 — and the newest version of Nvidia's Infiniband networking technology to share data even faster, Microsoft said in a blog post.

The new Bing is still in preview with Microsoft gradually adding more users from a waitlist. Guthrie's team holds a daily meeting with about two dozen employees they've dubbed the "pit crew" after the group of mechanics that tune race cars in the middle of the race. The group's job is to figure out how to bring greater amounts of computing capacity online quickly, as well as fix problems that crop up.

"It's very much a kind of a huddle, where it's like, 'Hey, anyone has a good idea, let's put it on the table today, and let's discuss it and let's figure out OK, can we shave a few minutes here? Can we shave a few hours? A few days?' " Guthrie said.

A cloud service depends on thousands of different parts and items — the individual pieces of servers, pipes, concrete for the buildings, different metals and minerals — and a delay or short supply of any one component, no matter how tiny, can throw everything off.

Recently, the pit crew had to deal with a shortage of cable trays — the basket-like contraptions that hold the cables coming off the machines. So they designed a new cable tray that Microsoft could manufacture itself or find somewhere to buy. They've also worked on ways to squish as many servers as possible in existing data centers around the world so they don't have to wait for new buildings, Guthrie said.

When OpenAI or Microsoft is training a large AI model, the work happens at one time. It's divided across all the GPUs and at certain points, the units need to talk to each other to share the work they've done. For the AI supercomputer, Microsoft had to make sure the networking gear that

handles the communication among all the chips could handle that load, and it had to develop software that gets the best use out of the GPUs and the networking equipment. The company has now come up with software that lets it train models with tens of trillions of parameters.

Because all the machines fire up at once, Microsoft had to think about where they were placed and where the power supplies were located.

Otherwise you end up with the data center version of what happens when you turn on a microwave, toaster and vacuum cleaner at the same time in the kitchen, Guthrie said.

The company also had to make sure it could cool off all of those machines and chips, and uses evaporation, outside air in cooler climates and high-tech swamp coolers in hot ones, said Alistair Speirs, director of Azure global infrastructure.

Microsoft is going to keep working on customized server and chip designs and ways to optimize its supply chain in order to wring any speed gains, efficiency and cost-savings it can, Guthrie said.

"The model that is wowing the world right now is built on the supercomputer we started building couple of years ago. The new models will be built on the new supercomputer we're training now, which is much bigger and will enable even more sophistication," he said.

Microsoft Philanthropies underwrites some Seattle Times journalism projects.