

# Why China remains hungry for AI chips despite US restrictions

Waters, Richard; Liu, Qianer

[ProQuest document link](#)

---

## ABSTRACT (ENGLISH)

The impact of soaring global demand for Nvidia's products is likely to underpin the chipmaker's second-quarter financial results due to be announced on Wednesday. Besides reflecting demand for improved chips to train the internet companies' latest large language models, the rush has also been prompted by worries that the US might tighten its export controls further, making even these limited products unavailable in future. Washington set a cap on the maximum processing speed of chips that could be sold in China, as well as the rate at which the chips can transfer data — a critical factor when it comes to training large AI models, a data-intensive job that requires connecting large numbers of chips together. Cost-benefit Many Chinese tech companies are still at the stage of pre-training large language models, which burns a lot of performance from individual GPU chips and demands a high degree of data transfer capability.

## FULL TEXT

The US acted aggressively last year to limit China's ability to develop artificial intelligence for military purposes, blocking the sale there of the most advanced US chips used to train AI systems.

Big advances in the chips used to develop generative AI have meant that the latest US technology on sale in China is more powerful than anything available before. That is despite the fact that the chips have been deliberately hobbled for the Chinese market to limit their capabilities, making them less effective than products available elsewhere in the world.

The result has been soaring Chinese orders for the latest advanced US processors. China's leading internet companies have placed orders for \$5bn worth of chips from Nvidia, whose graphical processing units have become the workhorse for training large AI models.

The impact of soaring global demand for Nvidia's products is likely to underpin the chipmaker's second-quarter financial results due to be announced on Wednesday.

Besides reflecting demand for improved chips to train the internet companies' latest large language models, the rush has also been prompted by worries that the US might tighten its export controls further, making even these limited products unavailable in future.

However, Bill Dally, Nvidia's chief scientist, suggested that the US export controls would have greater impact in future.

"As training requirements [for the most advanced AI systems] continue to double every six to 12 months," the gap between chips sold in China and those available in the rest of the world "will grow quickly," he said.

### Capping processing speeds

Last year's US export controls on chips were part of a package that included preventing Chinese customers from buying the equipment needed to make advanced chips.

Washington set a cap on the maximum processing speed of chips that could be sold in China, as well as the rate at which the chips can transfer data — a critical factor when it comes to training large AI models, a data-intensive job that requires connecting large numbers of chips together.

Nvidia responded by cutting the data transfer rate on its A100 processors, at the time its top-of-the-line GPUs,

creating a new product for China called the A800 that satisfied the export controls.

This year, it has followed with data transfer limits on its H100, a new and far more powerful processor that was specially designed to train large language models, creating a version called the H800 for the Chinese market. The chipmaker has not disclosed the technical capabilities of the made-for-China processors, but computer makers have been open about the details. Lenovo, for instance, advertises servers containing H800 chips that it says are identical in every way to H100s sold elsewhere in the world, except that they have a transfer rate of only 400 gigabytes per second.

That is below the 600 GB/s limit the US has set for chip exports to China. By comparison, Nvidia has said its H100, which it began shipping to customers earlier this year, has a transfer rate of 900 GB/s.

The lower transfer rate in China means that users of the chips there face longer training times for their AI systems than Nvidia's customers elsewhere in the world—an important limitation as the models have grown in size.

The longer training times raise costs since chips will need to consume more power, one of the biggest expenses with large models.

However, even with these limits, the H800 chips on sale in China are more powerful than anything available anywhere else before this year, leading to the huge demand.

The H800 chips are five times faster than the A100 chips that had been Nvidia's most powerful GPUs, according to Patrick Moorhead, a US chip analyst at Moor Insights & Strategy.

That means that Chinese internet companies that trained their AI models using top-of-the-line chips bought before the US export controls can still expect big improvements by buying the latest semiconductors, he said.

"It appears the US government wants to not shut down China's AI effort, but make it harder," said Moorhead.

#### **Cost-benefit**

Many Chinese tech companies are still at the stage of pre-training large language models, which burns a lot of performance from individual GPU chips and demands a high degree of data transfer capability.

Only Nvidia's chips can provide the efficiency needed for pre-training, say Chinese AI engineers. The individual chip performance of the 800 series, despite the weakened transfer speeds, is still ahead of others on the market.

"Nvidia's GPUs may seem expensive but are, in fact, the most cost-effective option," said one AI engineer at a leading Chinese internet company.

Other GPU vendors quoted lower prices with more timely service, the engineer said, but the company judged that the training and development costs would rack up and that it would have the extra burden of uncertainty.

Nvidia's offering includes the software ecosystem, with its computing platform Compute Unified Device Architecture, or Cuda, that it set up in 2006 and that has become part of the AI infrastructure.

Industry analysts believe that Chinese companies may soon face limitations in the speed of interconnections between the 800-series chips. This could hinder their ability to deal with the increasing amount of data required for AI training and they will be hampered as they delve deeper into researching and developing large language models. Charlie Chai, a Shanghai-based analyst at 86Research, compared the situation with building many factories with congested motorways between them. Even companies that can accommodate the weakened chips might face problems within the next two or three years, he added.

Richard Waters in San Francisco and Qianer Liu in Hong Kong

## **DETAILS**

<b>Subject:</b>	Internet; Artificial intelligence; Semiconductors; Training; Processing speed; Engineers; Export controls
-----------------	---

<b>Business indexing term:</b>	Subject: Artificial intelligence Training Export controls; Corporation: NVidia Corp; Industry: 54133 : Engineering Services
--------------------------------	---

<b>Location:</b>	China; United States--US
<b>Company / organization:</b>	Name: NVidia Corp; NAICS: 334413, 513210
<b>Classification:</b>	54133: Engineering Services
<b>Identifier / keyword:</b>	Lenovo Group Ltd; US Government; Asia-Pacific companies; US-China relations; US-China trade dispute; Washington D.C.; GPU; Nvidia; Moor Insights & Strategy, Inc.; Artificial intelligence; Chinese business & finance
<b>Publication title:</b>	FT.com; London
<b>Publication year:</b>	2023
<b>Publication date:</b>	Aug 21, 2023
<b>Publisher:</b>	The Financial Times Limited
<b>Place of publication:</b>	London
<b>Country of publication:</b>	United Kingdom, London
<b>Publication subject:</b>	Business And Economics
<b>Source type:</b>	Trade Journal
<b>Language of publication:</b>	English
<b>Document type:</b>	News
<b>ProQuest document ID:</b>	2853776405
<b>Document URL:</b>	<a href="https://unh.idm.oclc.org/login?url=https://www.proquest.com/trade-journals/why-china-remains-hungry-ai-chips-despite-us/docview/2853776405/se-2?accountid=14612">https://unh.idm.oclc.org/login?url=https://www.proquest.com/trade-journals/why-china-remains-hungry-ai-chips-despite-us/docview/2853776405/se-2?accountid=14612</a>
<b>Copyright:</b>	Copyright The Financial Times Limited 2023
<b>Full text availability:</b>	This publication may be subject to restrictions within certain markets, including corporations, non-profits, government institutions, and public libraries. In those cases records will be visible to users, but not full text.
<b>Last updated:</b>	2023-08-21
<b>Database:</b>	Global Newsstream

## LINKS

[Linking Service](#), [Find Full Text at UNH](#)

Database copyright © 2023 ProQuest LLC. All rights reserved.

[Terms and Conditions](#) [Contact ProQuest](#)