

Fujitsu's Approach to Using Cutting-Edge AI Supercomputers

● Tsuguchika Tabaru ● Akihiko Kasagi ● Takashi Arakawa ● Yuji Yoshioka
● Hidenori Hayashi ● Kinji Ito

With the Japanese government's current policy to accelerate the research and development of AI technology and its social implementation of research results, supercomputers designed for AI processing (hereafter, AI supercomputers) are rapidly being developed. Unlike conventional supercomputers, AI supercomputers have the following new requirements: system sizing based on standard performance evaluations that focus on learning accuracy and learning speed in deep-learning processes, and a software arrangement where frequent updates to deep-learning frameworks can be dealt with in a timely manner. Fujitsu and Fujitsu Laboratories have devised new performance evaluation methods and operational technologies for AI supercomputers. We combined these with the technologies for conventional supercomputers and deployed AI supercomputers to the National Institute of Advanced Industrial Science and Technology (AIST) and RIKEN Center for Advanced Intelligence Project (AIP). These systems will be a reference point in terms of future AI supercomputers. This paper first outlines AI supercomputers and their trends, then describes the approach to operating issues and the performance evaluation methods required by AI supercomputers. Finally, it describes the application of the operating technologies to AI supercomputers and their effects.

1. Introduction

Artificial intelligence (AI) has been attracting considerable attention in recent years. On receiving instructions from the Prime Minister, the government of Japan established the Strategic Council for AI Technology in 2016.¹⁾ This council has a supervisory role in accelerating the research and development of AI technology and the social implementation of research results by coordinating the AI activities of the Ministry of Internal Affairs and Communications (MIC), Ministry of Education, Culture, Sports, Science and Technology (MEXT), and Ministry of Economy, Trade and Industry (METI). In addition, the deployment of supercomputers oriented to AI processing (hereafter, AI supercomputers) is progressing rapidly to provide a platform for accelerating AI research.

There are key differences between the requirements for AI supercomputers and the characteristics required for conventional supercomputers. In deep-learning processing, which features a particularly high computational load even within AI technology, there is

no need for high-precision numerical calculations such as double-precision floating-point arithmetic (hereafter, double-precision arithmetic) as required for conventional supercomputers. Of more importance here is the ability to perform high-volume processing using single-precision floating-point arithmetic (hereafter, single-precision arithmetic) or half-precision floating-point arithmetic (hereafter, half-precision arithmetic). Consequently, for AI supercomputers, there is a need to supplement the performance evaluation methods used for conventional supercomputers with standard ones that focus on learning accuracy and learning speed in deep-learning processing.

At the same time, deep-learning frameworks (hereafter, DL frameworks) required for AI research are many and varied and updated frequently. Thus, in relation to the operation of AI supercomputers, there is also a need for a new software environment that can provide such application execution frameworks to many researchers.

Fujitsu and Fujitsu Laboratories made an early

start in studying performance evaluation methods for AI supercomputers and setting up operating environments. As a result, they have already received orders and made deliveries of AI supercomputers to the National Institute of Advanced Industrial Science and Technology (AIST) and RIKEN Center for Advanced Intelligence Project (AIP).

This paper begins by outlining AI supercomputers and trends in this field. It then describes Fujitsu's approach to operating issues and performance evaluation methods required by AI supercomputers. Finally, it describes the application of Fujitsu's operating technologies to AI supercomputers and their effects.

2. Features of AI supercomputers

2.1 Overview of AI supercomputers

The aim of an AI supercomputer is to perform high-speed execution of matrix operations, the dominant processing in machine learning. An AI supercomputer constitutes a platform dedicated to AI processing by joining several hundred computing nodes with large-capacity storage by a high-speed interconnect technology such as InfiniBand.²⁾ These computing nodes are equipped with accelerators for AI processing (hereafter, AI accelerators) to achieve high-speed processing of large-scale, matrix sum-of-product operations.

2.2 Differences with conventional supercomputers

Conventional supercomputers assume the use of double-precision floating-point numbers to obtain more accurate results. This simplifies the description of application algorithms since only a single data type of double precision is being used.

On the other hand, AI supercomputers used in deep learning do not necessarily require double-precision arithmetic. Examining the characteristics of deep-learning processing reveals that a major portion is occupied by large-scale matrix multiplication operations, so there is a need here for high-performance sum-of-product operations. However, since the data processed in deep learning has a relatively small number of digits, even computations at less than single precision do not make for any significant differences in learning results. For this reason, single- and half-precision floating-point arithmetics and even fixed-point arithmetic at 16 bits or 8 bits has come to be used in AI

supercomputers, and even new operation formats are being considered.

Conventional supercomputers and AI supercomputers differ also from the viewpoint of parallel processing. Conventional supercomputers often apply large-scale parallel processing to execute applications and find solutions. However, while it is generally true that computational throughput in deep-learning processing targeted by AI supercomputers can be improved through large-scale parallel processing, it is not necessarily true that all-important learning accuracy can be improved. With this in mind, the parallel processing of applications in deep-learning processing is generally of a small-to-medium scale (several tens to several-hundred nodes) considering the difficulty of a large-scale implementation (several-thousand to several-ten-thousand nodes).

There are also differences in terms of application execution environment. At present, many types of deep-learning frameworks required for AI research are being developed, so there is a need for an environment that enables AI researchers to use whatever DL frameworks are optimal for their research whenever needed. As a result, the establishing of a detailed application execution framework that can immediately provide many DL frameworks to each user is an essential operating requirement of AI supercomputers.

Differences between conventional supercomputers and AI supercomputers are summarized in **Table 1**.

3. AI-supercomputer requirements and issues

As described above, the characteristics of AI supercomputers differ from those of conventional supercomputers, so there is a need to revise existing performance evaluation methods for AI supercomputers. It is also important that AI supercomputers be capable of flexible operation so that the latest DL frameworks can be used in a timely manner.

3.1 Issues in performance evaluation of AI supercomputers

Similar to conventional supercomputers, it is important when drafting a plan for installing an AI supercomputer to quantitatively clarify performance requirements that take into account both learning accuracy and learning speed as an AI supercomputer. For conventional supercomputers, there are a set of

Table 1
Differences between conventional supercomputers and AI supercomputers.

	Conventional supercomputers	AI supercomputers
Accelerator	Mounted on some nodes (used for some scientific and technical calculations)	Mounted on all nodes (used in various types of DL frameworks)
Arithmetic precision	Assumes double precision	Assumes single/half-precision or fixed-point (a tradeoff exists between arithmetic precision and learning accuracy)
Parallel processing	Assumes large-scale parallel processing	Small/medium-scale parallel processing dominates
Application execution environment	Standard computation packages	Many and diverse DL frameworks

standard benchmark programs for evaluating computer simulation performance and other characteristics and a uniform performance evaluation method. For AI supercomputers, however, no standard benchmark programs or uniform performance evaluation methods have been established. Therefore, the present situation is such that even system performance comparisons or system sizing are hardly performed.

The speed of deep learning having high computational complexity is considered to be a suitable performance index for AI supercomputers. Here, "speed" comes in two types: the amount of data that can be processed per unit time (throughput performance) and the degree to which recognition accuracy can be improved per unit time (learning speed). The former is basically constant and can be evaluated from the operation speed of existing applications, which is the same as that in conventional supercomputers. The latter, on the other hand, changes as learning progresses. In actual use, the time taken until reaching the target recognition accuracy is important, so an index expressing average learning speed across all processing of a benchmark program representing deep learning is also necessary.

Additionally, in terms of practical use, it is also important that an AI supercomputer be able to execute major DL frameworks and typical neural-network learning used in research and development. Devising an index expressing practical performance (as a standard performance evaluation method) by a comprehensive benchmark that includes these DL frameworks is an issue to be solved.

3.2 Flexible operating technologies and associated issues

An AI supercomputer requires a usage format in which multiple researchers make use of computing resources such as AI accelerators and storage resources through a time-sharing process. However, research themes differ among researchers and needs with respect to an AI supercomputer can be quite diverse. It must therefore be possible to deal flexibly with the needs of each researcher in the operation of an AI supercomputer. The following two issues, in particular, must be addressed in AI supercomputer operation.

1) Effective use of AI-supercomputer computing resources

Resource management using a job scheduler can be considered as one way of making effective use of the limited computing resources of an AI supercomputer. Users, however, are accustomed to the interactive use of computing resources in contrast to batch processing. Interactive use monopolizes computing resources, but at the same time, it cannot help causing free moments during operation to occur resulting in a wasteful use of computing resources. Satisfying the need for interactive computing while minimizing inefficient use of computing resources is an issue in system operation that must be solved if limited computing resources are to be efficiently used among multiple users.

2) Provision of new DL frameworks

A DL framework has many external libraries required for computing. This means that installation must be performed while resolving complicated dependency among various versions and levels, all of which makes implementation troublesome.

DL frameworks have a short technical-innovation cycle as new algorithms are constantly being thought

up. This means highly frequent updates that can result in multiple releases per month. For this reason, ensuring that the latest version of a DL framework is always provided places a heavy burden on the system manager. Resolving this problem is an issue in system operation.

Sections 4 and 5 introduce Fujitsu's approach to solving the above issues.

4. Performance evaluation methods for AI supercomputers

We consider the following four elements to be essential to the performance evaluation of AI supercomputers.

1) Evaluation by basic performance

Focusing on basic performance as a standard requirement in deep-learning processing, we evaluate computational performance using matrix multiplication libraries and storage performance using existing standard benchmark programs. This is the same as the approach used in the evaluation of conventional supercomputers.

2) Evaluation by learning speed

Considering that AI supercomputers are to be used for computationally intensive deep learning, the evaluation of learning speed becomes necessary. We adopted well-known open source software (OSS) that users frequently use to perform this evaluation (for example Caffe,³⁾ etc.).

As previously pointed out, there is a difference between learning speed and throughput performance in deep learning. Since learning progresses by the stochastic gradient descent (SGD) method, the learning process may get stuck at a local optimum, resulting in the learning accuracy failing to improve even though throughput performance improves by the use of many AI accelerators. We therefore consider that measuring the time taken to reach the specified learning accuracy makes it possible to both measure learning speed and ensure learning accuracy important in practice.

These methods take into account evaluation that allows for a practical degree of accuracy while keeping the time taken for the evaluation within a realistic level.

3) Evaluation by throughput performance

The throughput performance of an AI application is also an important evaluation index. Major AI applications have been prepared with single-node operation

in mind. To extract AI-supercomputer performance and execute these applications at high speed, it is necessary to use AI accelerators installed in many nodes.

Furthermore, to maximize the throughput performance of an AI application, it is necessary to conceal communication time so that the AI accelerators do not enter an idle state. This is accomplished by a tuning scheme that enables three types of communication—that between the CPU and AI accelerators, that between AI accelerators, and that between servers where communication bottlenecks easily occur—to be performed in parallel with AI-accelerator processing. Taking the above measures into account, we evaluated the extent to which the throughput performance of an application can be sped up.

4) Evaluation by multiple frameworks

Finally, we evaluated whether a number of major DL frameworks could be effectively used.

Since multiple DL frameworks and neural networks are currently in use, it is necessary to perform evaluations using combinations of the two. We therefore devised an evaluation method that targets a small number of such combinations while including all evaluation criteria at once. We did this by combining "evaluation by learning speed" and "evaluation by throughput performance" described above and setting an evaluation pattern that includes the frameworks essential to the evaluation and a pattern that enables the use of any open-source frameworks.

5. Operating technologies for a flexible AI-supercomputer platform

5.1 Resource allocation by job scheduler

We apply job-scheduler technology as used in conventional supercomputers to AI supercomputers to manage the execution units of AI applications as jobs and enable effective use of computing resources. A job scheduler secures computing resources including AI accelerators for each job and controls the order of job execution.

Jobs to be executed come in two types: batch and interactive. The job scheduler, however, is capable of allocating resources in a uniform manner without separating job types, thereby supporting the interactive-computing needs of users. However, when executing interactive jobs, the time taken up by user operation must also be added on, and as a result,

heavy use of interactive jobs will make for longer wait times for initiating job execution. For this reason, we reduced job-execution wait times by narrowing down the resources allocated to interactive jobs and limiting the execution time of interactive jobs to one hour.

5.2 Provision of DL environment by container technology

We have constructed an environment that provides DL frameworks in the form of Docker containers to enable system managers to provide users with the latest DL frameworks in a timely manner.

Docker⁴⁾ is a virtualization technology that uses the Linux Kernel of the host OS, instead of a hypervisor or guest OS as in conventional virtualization software (VMware, Hyper-V, etc.). It has consequently been attracting attention as a virtualization technology with little overhead. Applying Docker to AI supercomputers enables system managers to install the latest DL framework without degradation even if other DL frameworks are currently running.

6. Examples of application to AI supercomputers

In this section, we introduce leading AI supercomputers in Japan as examples of applying the technologies described in Sections 4 and 5.

6.1 Large-scale, low-power cloud platform for AI processing (AIST)

The National Institute of Advanced Industrial Science and Technology (AIST) has deployed a supercomputer system for AI applications centered about FUJITSU Server PRIMERGY CX2570 M4 servers to conduct cutting-edge research and development related to AI technologies. This system achieves a theoretical peak performance in half-precision arithmetic of 550 petaflops and features computing nodes each equipped with four NVIDIA Tesla V100 graphics processing units (GPUs).

At the time of the public offering for facility construction, individual ICT vendors were given a hearing on performance evaluation criteria. Fujitsu Laboratories offered its views on performance evaluation methods as described above, which were adopted in some benchmark tests. Furthermore, when it came time to make a proposal, Fujitsu was able to propose

a system that satisfied the requirements set forth by AIST by leveraging these views on performance evaluation. This proposal was accepted, and Fujitsu went on to receive an order for this supercomputer system and complete its delivery.

From the viewpoint of operation, DL frameworks are executed under privileged authority with Docker, which opens a system to security risks. To mitigate these risks we are planning to use Singularity, a new container technology oriented to high performance computing (HPC), in place of Docker so that the latest parallel-distributed DL frameworks can be executed under user authority. Going forward, we plan to incorporate more advanced technologies in this way.

6.2 Computer system for AI research (RIKEN Center for AIP)

Fujitsu deployed an AI supercomputer system for AI research (RAIDEN) at the RIKEN Center for AIP in March 2017. It then upgraded RAIDEN in March 2018 boosting its performance to a theoretical peak performance in half-precision arithmetic of 54 petaflops. This system applies the operating technologies described above enabling the use of Docker containers and the timely provision of frequently updated DL frameworks. RAIDEN also provides users with both a highly convenient interactive environment and a batch job environment, thereby achieving an AI computing environment that reduces user wait time and enables the effective use of computing resources 24 hours a day.

7. Conclusion

This paper introduced performance evaluation methods and operating technologies devised to accelerate AI research by meeting the needs for evaluation methods and operating environments specifically for AI supercomputers. It described the application of these methods and technologies to AI supercomputers and examined their effects. The AI supercomputer systems at AIST and RIKEN Center for AIP that apply the technologies introduced in this paper will serve as reference systems for future AI supercomputers.

Looking forward, we can expect requirements and issues surrounding AI supercomputers to change as new AI-specific processors and innovative AI techniques appear while incorporating the latest technologies. We intend to propose optimal solutions to these issues

by leveraging the advanced AI technologies of Fujitsu Laboratories and the delivery and operating experience of large-scale AI supercomputers of Fujitsu.

References

- 1) Ministry of Education, Culture, Sports, Science and Technology (MEXT), Information Science and Technology Committee, 94th Materials Distribution, "Materials 2-1 Strategic Council for AI Technology, Progress Report on AIP: AI, Big Data, IoT, Cyber Security Integration Project," pp. 1–3, 2016/6/2 (in Japanese).
http://www.mext.go.jp/b_menu/shingi/gijyutu/gijyutu2/006/shiryo/_icsFiles/afieldfile/2016/08/09/1374745_002_1.pdf
- 2) InfiniBand Trade Association: About InfiniBand.
<https://www.infinibandta.org/about-infiniband/>
- 3) Berkeley Artificial Intelligence Research Lab website.
<http://caffe.berkeleyvision.org/>
- 4) Docker website.
<https://www.docker.com/what-docker/>



Tsuguchika Tabaru

Fujitsu Ltd.

Mr. Tabaru is currently engaged in research and development of deep-learning parallel-distributed processing and HPC compilers.



Akihiko Kasagi

Fujitsu Ltd.

Mr. Kasagi is currently engaged in research and development of many-core optimization and deep-learning algorithms.



Takashi Arakawa

Fujitsu Ltd.

Mr. Arakawa is currently engaged in research and development of high-speed algorithms in the fields of HPC and deep learning.



Yuji Yoshioka

Fujitsu Ltd.

Mr. Yoshioka is currently engaged in HPC and AI-business negotiations in the field of science.



Hidenori Hayashi

Fujitsu Ltd.

Mr. Hayashi is currently engaged in HPC and AI-business negotiations in the field of science.



Kinji Ito

Fujitsu Ltd.

Mr. Ito is currently engaged in HPC and AI-business negotiations in the field of science.