



Leading AI Scalability Benchmarks with Microsoft Azure

AUTHOR

Russ Fellows

VP, Labs | Signal65

IN PARTNERSHIP WITH



NOVEMBER 2024

Overview

Businesses understand that AI has rapidly moved from being an interesting phenomenon to a tool that can be leveraged to provide firms with competitive advantages. Companies are now exploring options for developing and deploying AI applications both on premises and in public clouds. This trend is occurring across nearly all vertical segments and enterprise classes, with virtually all businesses affected by the drive to adopt AI implementations.

The ability to experiment, and rapidly grow or shrink infrastructure has always been one of the most compelling benefits of cloud computing. As companies experiment with AI training, fine-tuning and other resource intensive workloads, the benefits of cloud computing are clear. The ability to scale infrastructure to precisely match rapidly changing resource needs is vital.

Futurum Research found that the use of public clouds for AI workloads will grow at a 30% compound annual rate over the next five years.¹ Additionally, Futurum Research found that Azure was the most commonly mentioned cloud vendor in its AI cloud research survey.²

Our clients see value in the ability to quickly scale their AI investments without requiring significant up front capital. As opportunities for growth present themselves, cloud AI infrastructure enable users; then afterwards reducing resources is easily achieved.

Microsoft Azure asked Signal65 to analyze their AI portfolio, along with the latest set of MLPerf Training benchmark results. The MLPerf benchmarks provide a way to effectively compare alternative offerings. While this test measures large-scale AI training, the results may be analyzed to understand how various infrastructure options could be utilized for workloads many enterprises are evaluating, including AI fine-tuning and inferencing workloads.

In analyzing Microsoft Azure AI's capabilities, together with the MLPerf 4.1 training benchmark results³, we found the following:

- 1. Microsoft Azure's MLPerf outperforms a leading cloud competitor by 28%:** Using same number of GPUs to run Llama70B fine-tuning
- 2. Microsoft Azure Delivers AI at Scale:** Azure is one of three cloud vendors to demonstrate Llama 70B performance at scale, beyond 100 GPUs
- 3. Azure AI has Industry Leading Options:** Azure AI optimized software stacks, together with leading hardware options provide the highest performing, cloud computing AI workload options (based upon MLPerf DC Training Results).
- 4. Azure AI has Leading Price / Performance:** As a top public cloud provider, Azure provides highly competitive price / performance options for AI workloads.

Microsoft Azure AI

Microsoft Azure has established itself as a leading cloud platform for AI training and inferencing with innovative hardware investments including the Azure Maia 100 chip and GPUs from industry leaders NVIDIA and AMD. This helps Azure to provide comprehensive AI infrastructure, with options including NVIDIA A100, H100 and the new H200 GPUs to cost effectively support a range of AI training and inferencing workloads.

Key Advantages of Azure AI Infrastructure

- **Top-Tier Performance Accelerators:** Access the best performing and most efficient AI software stacks along with state-of-the-art hardware.
- **Optimized Software Infrastructure:** Continuous optimizations help keep AI applications running smoothly.
- **Research and Innovation:** Microsoft's continuing investments in AI research and open-source software benefit all AI researchers.
- **Industry-Leading Service:** Implement AI solutions quickly without sacrificing performance, backed by Azure's robust infrastructure.
- **Scalability:** From single nodes to thousands, Azure scales with your needs, ensuring seamless growth.

The combination of industry leading AI hardware accelerators, together with optimized AI software solution stacks and Microsoft Azure's ability to deliver scalable infrastructure make Azure a leading destination for building, training and running AI workloads. This scalability extends beyond GPUs to encompass data, networking, and more, ensuring that companies can grow infrastructure to meet their needs.

Signal65 Comments: Microsoft Azure's global resources enable companies to leverage this massive scale to rapidly expand their AI resources required for training, and then down again for production inferencing. Azure's scalability and flexibility is extremely difficult to achieve with on-premises solutions and has few competitors among cloud providers.

Azure: Optimized Software and Hardware

Microsoft Azure AI provides many advantages with its optimized AI software stacks compared to general-purpose cloud providers and home-grown on-premises deployments. These advantages stem from Azure's purpose-built infrastructure, advanced software optimizations, and integration with a wide range of AI tools and services.

Few cloud providers can offer GPU clusters at the scale demonstrated by Azure, nor provide the variety or software stacks of Azure. In contrast, Azure's demonstrated ability to scale to over 500 GPU nodes, along with their optimized AI software stacks provide a cost-effective option for a variety of AI workloads.

Signal65 Comments: In speaking with enterprises regarding AI, there are several challenges with on-premises deployments, including the significant capital resources required, and aligning equipment availability with their particular AI needs. For these reasons, many companies are looking to leverage cloud providers significant hardware and software resources for some portion of their AI development and deployments.

Azure Optimized AI Hardware

The latest generation NVIDIA H200 GPUs provide several advantages over the previous NVIDIA H100 GPUs for training, fine-tuning and inferencing, due to several factors. The H200 has more GPU memory, and the H200 uses a new, faster type of memory that has higher bandwidth⁴.

The larger memory capacity of the H200 GPUs enables running larger models and datasets with fewer GPUs, particularly when running production AI inferencing workloads. According to Signal65 analysis, H200 GPUs can support running the Llama 3.1 405B parameter model with 8 or fewer GPUs, compared to 14 – 16 GPUs using H100 GPUs when utilizing the default BF16 datatype .

Additionally, more memory reduces paging when inferencing large batches or long context results. Along with improved price / performance, there are additional power and cooling benefits from using fewer GPUs to perform the same work. Beyond memory enhancements, the H200 includes optimization when running lower precision models, significantly enhancing performance of FP8, FP16 and INT8 data types.

An overview of the H200 improvements compared to the H100 include:

- **Bandwidth improvement of 37%:** The H200's new HBM3e memory bandwidth is 4.8 TB/s vs. 3.35 TB/s for the H100.
- **Expanded Memory of 76%:** Along with improved bandwidth, the memory capacity of the H200 is 141 GB, compared to typical 80 GB configs for the H100.
- **Enhanced Tensor Core Performance:** The H200's tensor cores handle mixed-precision operations more efficiently with FP8, FP16, and INT8 datatypes.
- **Improved Energy Efficiency:** The H200 is more energy-efficient, while using a similar power envelope the H200 delivers up to 2x better inferencing performance.
- **Enhanced GPU scaling:** The H200 supports both PCIe Gen5 and enhanced NVLink, for faster communications with other GPUs and systems to enhance distributed AI workloads.
- **Superior Performance:**
 - **Training:** The H200's performance across a range of workloads is between 20 and 50% faster than the H100.
 - **Inferencing:** The H200 achieved a 90% increase in inference speed on the Llama 2 (70 billion parameters) compared to the H100.



Nvidia H200, source: [Nvidia.com](https://www.nvidia.com)

Signal65 Comments: The NVIDIA H200 GPU offers substantial advantages over the previous best in class H100 for AI workloads, due to increased memory capacity and bandwidth along with enhanced processing performance. These benefits translate into the ability to utilize fewer GPUs to achieve the same results or process workloads faster. As a result, users can improve energy efficiency and price performance when using the H200 for training, fine-tuning and inferencing.

Microsoft Azure: Research, Innovation and Scale for AI

Microsoft Azure AI has been a leader of artificial intelligence advancements, through their investments in hardware and software, and in particular Azure's commitment to AI research and innovation as evidenced by partnerships and in-house developments. Azure's AI software investments are demonstrated by its cutting-edge research initiatives such as OpenAI, DeepSpeed, and Phi. These projects highlight Azure's role as a thought leader in the AI space, pushing the boundaries of what's possible.

- **DeepSpeed:** Is an AI/ML deep learning optimization library developed by Microsoft, designed to make distributed training easy, efficient, and effective. DeepSpeed enables high performance during training and inferencing of models using billions of parameters, while efficiently scaling to thousands of GPUs.
- **ONNX Runtime:** Is an AI/ML inference optimization library developed by Microsoft. ONNX libraries are compatible with different hardware, drivers, and operating systems, and provides optimal performance by leveraging hardware accelerators when possible, alongside graph optimizations and transforms. ONNX provides an efficient solution for running SLMs\LLMs on edge devices.
- **Phi:** Are a family of small language, AI models (SLMs) developed by Microsoft, balancing performance and efficiency for commercial applications and research. These models support language and visual capabilities making them versatile for applications including code generation, OCR imaging, summarization and reading comprehension. Additionally, these models are designed to be compact while producing good results, making them appropriate for deployment on edge devices.

Azure Benefits Summary

Microsoft Azure AI offers several distinct advantages over other cloud vendors' AI offerings:

1. **Comprehensive AI Services:** Azure provides a wide range of AI services to more than 60,000 customers⁵, including Azure Machine Learning, Azure AI Services, and Azure OpenAI Service. This extensive suite allows developers to build, deploy, and manage AI solutions tailored to their specific needs, whether they require custom models or pre-built AI functionalities.
2. **Research and Innovation:** Microsoft's commitment to AI research is evident through initiatives like DeepSpeed and Phi. These projects push the boundaries of AI capabilities, offering advanced optimization techniques and efficient, powerful models.
3. **Scalable Performance:** Azure's infrastructure is designed to scale efficiently, supporting everything from small applications to large-scale AI deployments, including a #3 ranking on the top 500 supercomputer list.⁷ This scalability extends beyond GPU resources to include data and networking, ensuring enterprises can grow with high performance. Additionally, Azure offers a variety of hardware accelerators, providing cost and price / performance options.
4. **Integration and Ecosystem:** Azure AI includes a model catalog with over 1600 models⁸ (Foundational, open source, multimodal and SLMs) that enables practitioners the ability to easily select and compare various LLM models and run those models on scalable, optimized hardware instances within Azure.
5. **Cost Efficiency:** Azure offers competitive pricing and cost management tools, making it a cost-effective choice for businesses of all sizes. The economic benefits of migrating to Azure for AI readiness include lower costs, increased innovation, and better resource allocation.
6. **Leading Scale:** With a global network of data centers, Azure ensures low latency and high availability for AI applications, with over 300 datacenters in 60 regions.⁹ Additionally, Azure complies with a wide range of international standards and regulations, providing a secure and compliant environment for AI development.

The combination of industry leading GPUs from NVIDIA, AMD along with the Azure Maia chip position Microsoft Azure as a leading cloud platform for AI training and inferencing. Azure's scalable infrastructure, advanced hardware, and comprehensive AI services provide excellent performance, scalability, and efficiency, making it a top choice for AI practitioners and developers.

MLPerf – Training Results

Microsoft Azure, along with other industry leaders has recently participated in the MLCommons AI benchmarks, releasing results for the MLPerf Training workloads. In analyzing the results of the latest training results¹⁰, Signal65 found that AI customers have an advantage running the Llama 70B fine-tuning workload on the Azure AI Platform compared to alternatives. These advantages include overall performance, performance on a per GPU basis and Azure's ability to deploy and run at massive scale.



Performance Analysis

- **Azure's Performance:** Azure's ND_H200_v5 x64 system, with 512 GPUs, achieved a performance of **1.6085** minutes for the Llama 70B fine-tuning workload.
- **Comparison with Alternatives:** Azure's results were one of only 4 entries using 100 or more GPUs for the Llama 70B fine-tuning workload, and one of only two cloud vendors to submit results at this scale.
- **Leading Scale:** Azure's previous 10,752 GPU, NVIDIA H100 results on the MLPerf Training 3.1 results remains as one of the largest GPU results submitted.¹¹

Advantages of Azure

- **Scalability:** Azure's infrastructure can efficiently scale to hundreds of GPUs, providing consistent performance improvements as the number of GPUs increases.
- **Efficiency:** Azure's performance per GPU demonstrates efficient utilization of hardware resources as solutions scale.
- **Infrastructure:** Azure's supporting, server, storage and networking infrastructure is designed to support high-performance workloads, ensuring scalable AI deployments.

Azure's ND_H200_v5 x64 system shows strong GPU performance for the Llama 70B workload, making it a competitive choice for practitioners looking to deploy AI workloads efficiently in the public cloud. Azure's ability to scale and its efficient use of hardware resources provide significant advantage over other cloud vendors solutions, as evidenced by the latest MLPerf v 4.1 Training Results.

Within the continuum of AI workloads, fine-tuning is a common approach to creating large language models that are designed to solve specific problems. Fine-tuning typically utilizes fully trained base models as the starting point. The specific MLPerf workload that was measured uses a fine-tuning optimization method known as Low Rank Adaptation (LoRA) to help reduce the computational requirements needed to fine-tune large language models such as Llama 70B.

“ Azure’s MLPerf 4.1 results outperform a leading cloud competitor by 28%, using 512 GPUs to fine-tune Llama-70B ”

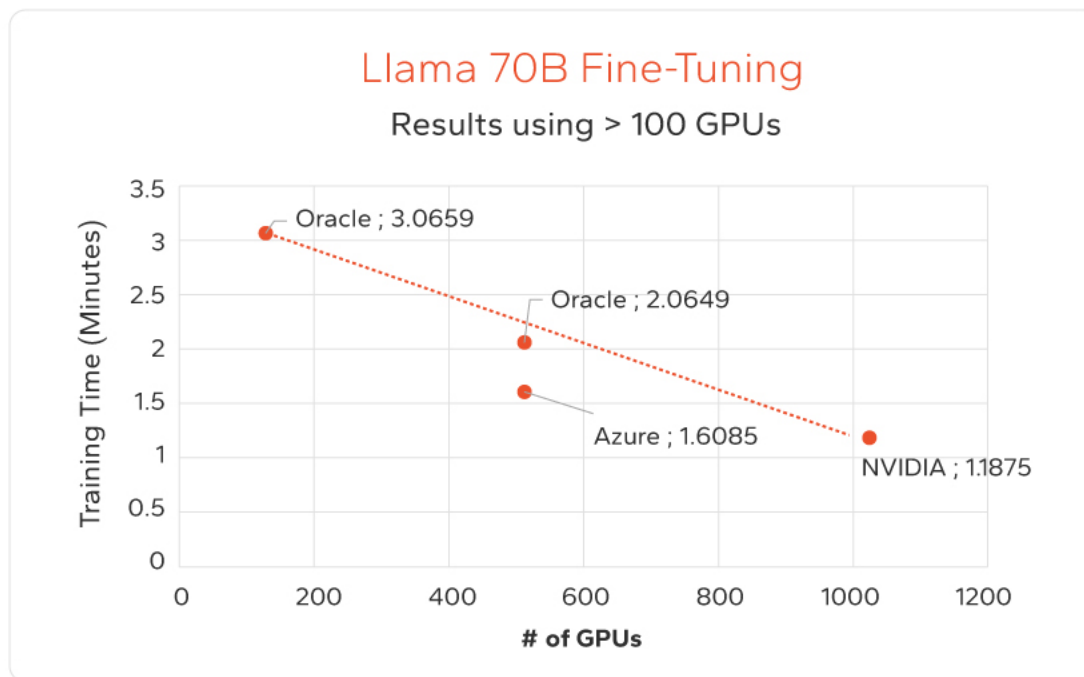


Figure 1: MLPerf DC Training – Llama70B fine-tuning results > 100 GPUs

As seen above, Azure’s result using 512 GPUs is significantly better, being below the trend line. By requiring less time compared to Oracle’s 512 GPU results, implementations will have reduced training time to achieve the same results. This provides Azure clients with a choice: either use fewer GPUs to achieve results in the same amount of time, or use the same number of GPUs and achieve results faster.

Relative Advantage of Azure

Azure’s ND_H200_v5 x64 system shows a total time of **1.6085**, which is significantly lower (better) compared to the alternatives using H100 GPUs. Here are the comparisons:

- **Compared to Oracle 512xBM.GPU.H100.8:**

- **Azure’s Advantage:** Approximately 1.28x the speed, or 28% faster

◊ Math: $(2.0649 / 1.6085) = 1.28x$; $(1.28 - 1.0) = 28\%$ faster

- **Compared to Oracle 128xBM.GPU.H100.8:**

- **Azure’s Advantage:** Approximately 87.1% better

Moreover, Azure achieved 70% of the improvements of the larger NVIDIA cluster, without requiring any additional GPUs.

Overall, the MLPerf results demonstrate that Azure’s H200 GPUs provide a clear performance advantage on a per GPU basis for the Llama 70B fine-tuning workload, achieving significantly lower times compared to alternatives using H100 GPUs, making it a competitive choice for AI practitioners looking for scalable and efficient AI solutions.

Final Thoughts

As companies evaluate their options for developing and deploying AI workloads, there are several factors that are important considerations, including the ability to match their infrastructure requirements to their application needs. As AI workloads continue to evolve rapidly, many companies find that they have neither the skills, capital or expertise required to deploy the hardware necessary for running AI workloads.

Microsoft Azure AI offers several distinct advantages over other cloud vendors' AI offerings:

- **Comprehensive AI Services:** Azure provides a wide range of AI services, including Azure Machine Learning, Cognitive Services, and the Azure OpenAI Service, allowing developers to build, deploy, and manage AI solutions tailored to their needs.
- **Scalable Performance:** Azure's infrastructure is designed to scale efficiently, supporting everything from small to large-scale AI deployments. The ability to deliver scalable infrastructure, requires extensive expertise in deploying systems at scale, along with architecture, GPU's, systems, storage and networking.
- **Integrated AI Ecosystem:** Azure AI integrates seamlessly with other Microsoft products and services, enhancing productivity and enabling a cohesive workflow across different platforms.
- **Cost Efficiency:** Azure offers competitive pricing and cost management tools, making it a cost-effective choice for businesses of all sizes.

As AI workloads mature, and businesses understand their long-term utilization needs for production AI inferencing, many will find the need to deploy additional AI resources for experimentation, custom training or fine-tuning workloads. In many cases, the capital costs and timelines required to build out large GPU clusters on premises will not be cost effective compared to cloud options, such as Azure AI.

Based upon the MLPerf and other data points, Signal65 believes that Microsoft Azure AI is a leading option for companies who want access to optimized AI software and hardware stacks, without requiring the significant investments in time, capital equipment and human resources necessary to build this infrastructure themselves. Azure's proven ability to scale clusters to hundreds of GPUs, with industry leading performance clearly demonstrates Azure AI's leadership in this space.

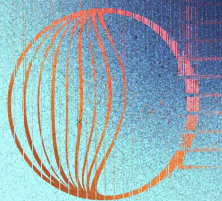
Appendix

Key data points comparing the MLPerf, data center Llama 70B fine-tuning time required, for solutions at 100 or more GPUs:

	Azure ND_H200_v5 x63	Oracle 512xBM.GPU.H100.8	Oracle 128xBM.GPU.H100.8	NVIDIA Eos_n128
Number of GPUs	512	512	128	1024
Performance (Llama 70B)	1.6085	2.0649	3.0659	1.1875

Footnote References:

- ¹ Futurum Intelligence AI Dashboard (AI growth projections): <https://www.futurumgroup.com>
- ² Futurum Intelligence AI Dashboard (AI cloud services survey): <https://www.futurumgroup.com>
- ³ ML Perf, Data Center Training v4.1 results: <https://mlcommons.org/2024/11/mlperf-train-v41-results/>
- ⁴ NVIDIA H200 Product Page: <https://www.nvidia.com/en-us/data-center/h200/>
- ⁵ MS Azure 2024 Updates: <https://azure.microsoft.com/en-us/blog/announcing-the-availability-of-azure-openai-data-zones-and-latest-updates-from-azure-ai>
- ⁶ Microsoft Phi SLM models: <https://azure.microsoft.com/en-us/blog/introducing-phi-3-redefining-whats-possible-with-slms/>
- ⁷ Top 500 Supercomputer; June 2024: <https://www.top500.org/lists/top500/2024/06/>
- ⁸ MS Tech Community: <https://techcommunity.microsoft.com/blog/aiplatformblog/ai21-jamba-instruct-launches-on-azure-ai-models-as-a-service/4170102>
- ⁹ MS Azure Infrastructure: <https://datacenters.microsoft.com/>
- ¹⁰ MLCommons, Data Center Training, v 4.1 Results: <https://mlcommons.org/benchmarks/training/>
- ¹¹ MLPerf Data Center Training, version 3.1 results, Public Identifier: 3.1-2002



Important Information About this Report

CONTRIBUTORS

Russ Fellows

VP, Labs | Signal65

Mitch Lewis

Performance Analyst | Signal65

PUBLISHER

Ryan Shrout

President and GM | Signal65

INQUIRIES

Contact us if you would like to discuss this report and Signal65 will respond promptly.

CITATIONS

This paper can be cited by accredited press and analysts, but must be cited in-context, displaying author's name, author's title, and "Signal65." Non-press and non-analysts must receive prior written permission by Signal65 for any citations.

LICENSING

This document, including any supporting materials, is owned by Signal65. This publication may not be reproduced, distributed, or shared in any form without the prior written permission of Signal65.

DISCLOSURES

Signal65 provides research, analysis, advising, and consulting to many high-tech companies, including those mentioned in this paper. No employees at the firm hold any equity positions with any companies cited in this document.

IN PARTNERSHIP WITH



Microsoft

ABOUT SIGNAL65

Signal65 is an independent research, analysis, and advisory firm, focused on digital innovation and market-disrupting technologies and trends. Every day our analysts, researchers, and advisors help business leaders from around the world anticipate tectonic shifts in their industries and leverage disruptive innovation to either gain or maintain a competitive advantage in their markets.



signal65



signal65

CONTACT INFORMATION

Signal65 | signal65.com