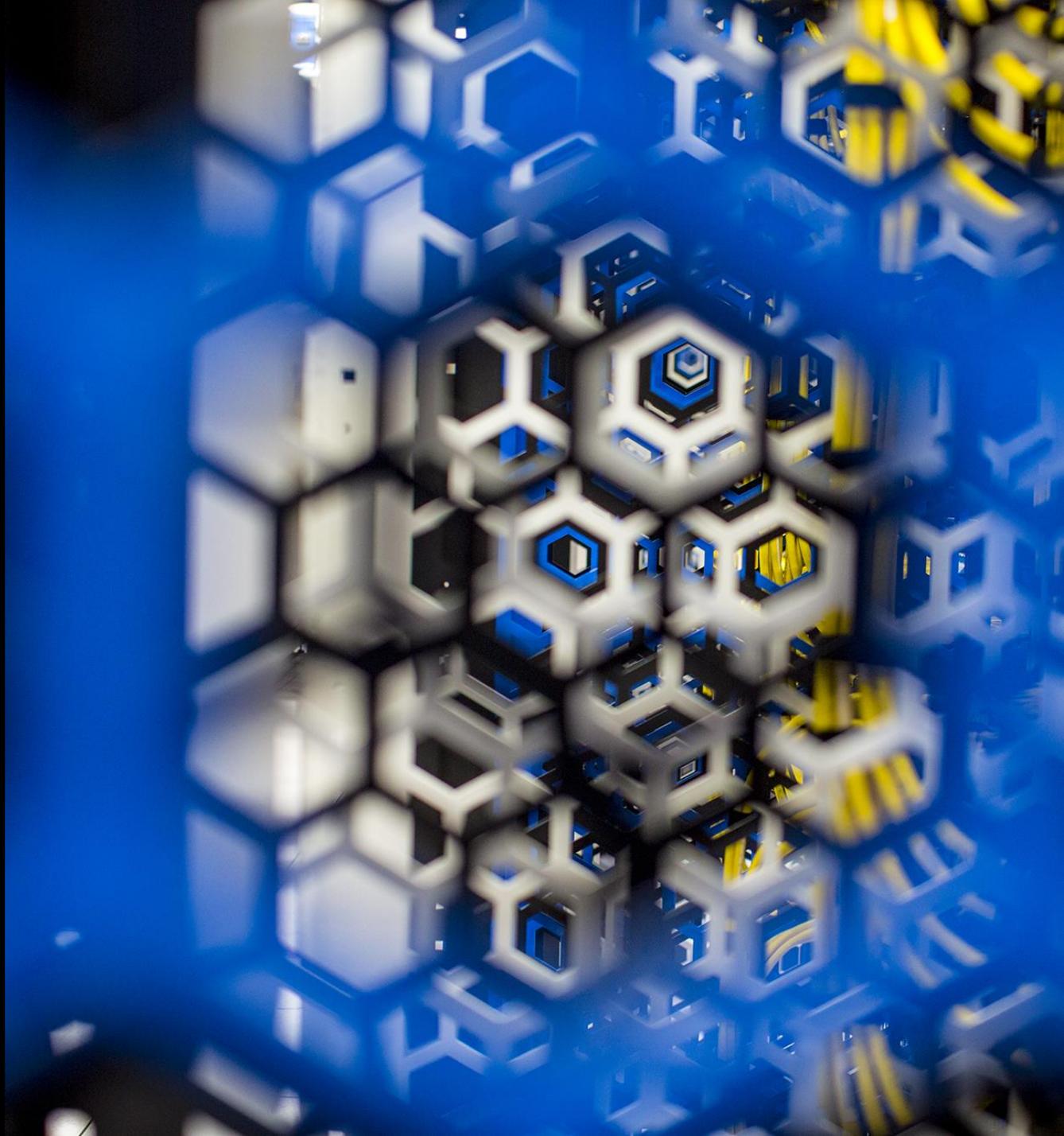




SuperComputing Japan! 2024

生成 AI & クラウド
スーパーコンピューティングを
支える Azure HPC のご紹介

日本マイクロソフト株式会社
クラウドソリューションアーキテクト
五十木 秀一 (Shuichi Gojuki)
2024/03/12



Azureにおける AI & 機械学習サービス

Azure AI services



Bot Service



Cognitive Search



Form Recognizer



Video Indexer



Metrics Advisor



Immersive Reader



Vision



Speech



Language



Decision



Azure OpenAI Service

Machine Learning



Azure Machine Learning platform



Deep learning frameworks

Azure AI Infrastructure

Azure AI & MLを支えるインフラストラクチャ



Transformative AI Services

開発者とデータサイエンティスト向けに設計されたAIサービスのポートフォリオ
[Azure AI Services](#)



Machine learning platform

マネージド エンドツーエンドの機械学習プラットフォームと OSS フレームワーク
[Azure Machine Learning](#) [PyTorch](#) [ONNX Runtime](#) OSS frameworks



Workload orchestration

使い慣れたツールとプロセスを使用したエンドツーエンドのワークフロー
[Azure Machine Learning](#) [VM Scale Sets](#) [Azure Batch](#) [Azure CycleCloud](#) OSS frameworks



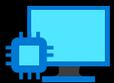
Fast, secure networking

高速なノード間相互接続とエッジからクラウドへの接続
[InfiniBand](#) [ExpressRoute](#)



High-performing Storage

用途に応じて選択可能な、幅広いストレージ機能
[Azure Blob](#) [Azure Managed Lustre](#)



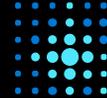
Optimized compute

80K+コアまで拡張可能なフルレンジのGPUおよびCPU
[N series virtual machines](#)

あらゆる規模のAIをサポートする確かな実績

Microsoft は、最も複雑なディープラーニングモデルの一部に対して、AIに最適化されたインフラストラクチャとスーパーコンピューティングスケールを提供します。

2023



マイクロソフトは、Azure OpenAI Service と新しい Bing を発表。これらはAzure AI インフラストラクチャにてトレーニング

2022



Microsoft が NVIDIA と最新の GPU に関して提携し、スーパーコンピューティングのリーダーシップを拡大



Microsoft は、530B パラメーターの NeMo Megatron ベンチマークを175台のVMで実行

2021



Microsoft は、クラウドにおいてTOP500スパコンで #1 にランクイン



Microsoft がクラウドで初めて記録的な MLPerf の結果を提供

2020



Microsoft はTOP500スパコンのトップ5に相当する環境をOpenAI 向けに構築



MicrosoftがAI スーパーコンピューティング VMを発表

最も包括的で信頼できるAIインフラストラクチャ



Autonomous Driving



Smart Health



Process Automation



Natural Language



Climate Science



Sustainable Energy

Azure AI Infrastructure

GPU に最適化された VM の Azure の包括的なポートフォリオは、あらゆる規模の AI 需要を満たし、Azure AI サービス、オープンソース ソリューション、ツールチェーンとシームレスに統合します。

Model Inference

リアルタイム推論、バッチ推論...



Mid-Range Training

スマート在庫管理、創薬...



High-End Training

深層学習モデル、生成AIモデル、自然言語処理...



パフォーマンスと拡張性の要件の増大



Azure Datacenters



Globally distributed,
trusted and secure



Redundancy,
Backup & Recovery



Compliance
Certifications

Microsoft is powered by Azure AI infrastructure

Security
Copilot

Edge
Bing Chat
Teams

Microsoft
365 Copilot

Windows
Copilot

Dynamics
365 Copilot

Azure
OpenAI API

Microsoft runs on Azure AI

Azure AI runs on Azure HPC infrastructure

Real time inference & low-cost
compute



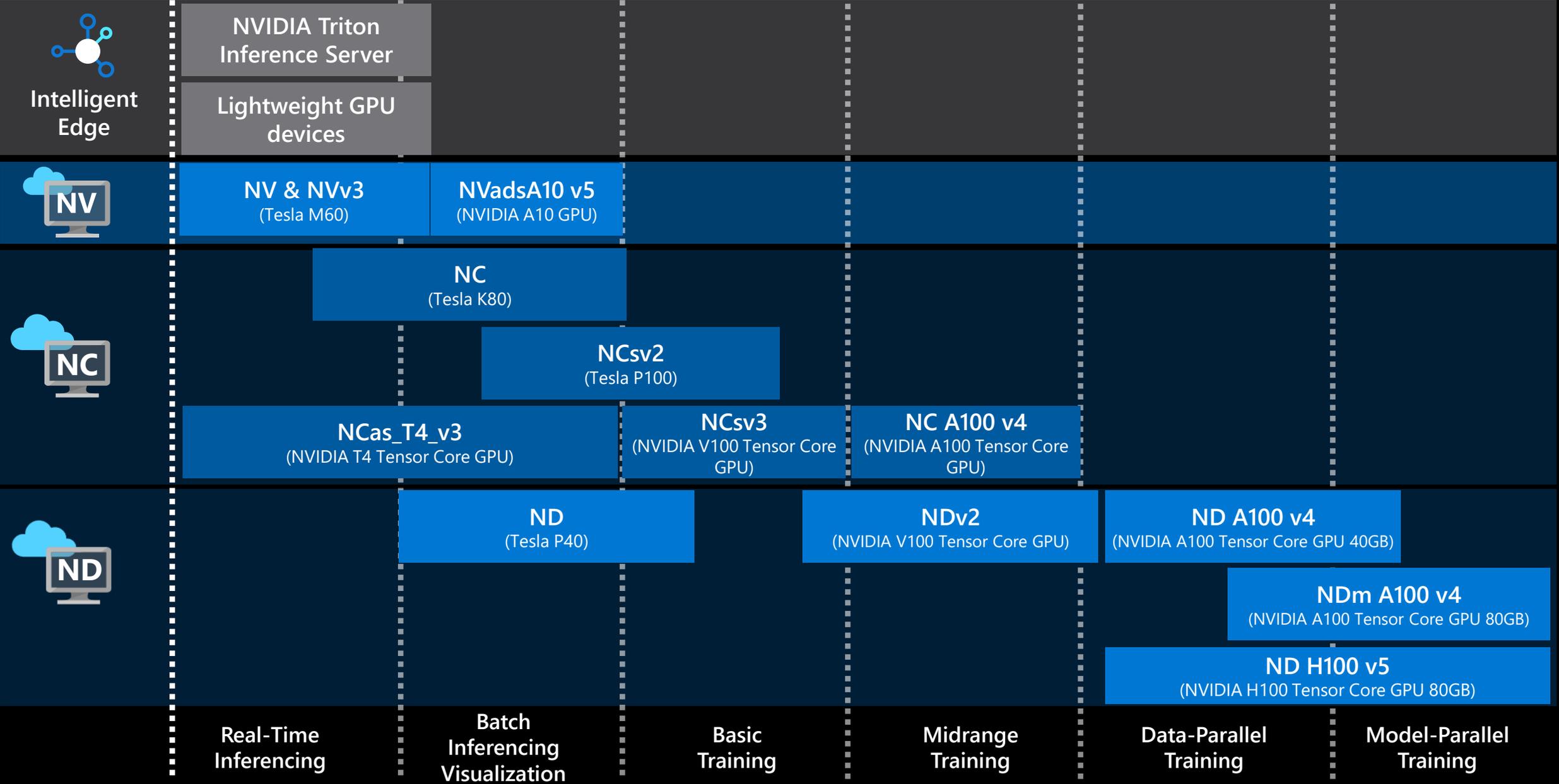
Mid-range training & dense
inference



Distributed training & generative
inference



Azure では、目的に応じた GPU 利用のための最適な選択肢が提供されます



AI イノベーターは Azure AI インフラ で実行



"NVIDIA と Microsoft Azure は、複数の世代の製品を通じて協力し、世界中の企業に最先端の AI イノベーションをもたらしてきました。NDv5 H100 仮想マシンは、ジェネレーティブ AI アプリケーションとサービスの新時代を後押しするだろう」と語った。

"

Ian Buck, Vice President of hyperscale and high-performance computing at NVIDIA



"スーパーコンピューターを Azure と共同で設計することは、要求の厳しい AI トレーニングのニーズを拡大し、ChatGPT のようなシステムでの研究と調整作業を可能にするために不可欠でした。"

Greg Brockman, President and Co-Founder of OpenAI



"会話型 AI に焦点を当てるには、最も複雑な大規模言語モデルのいくつかを開発してトレーニングする必要があります。Azure の AI インフラストラクチャは、これらのモデルを大規模で確実に効率的にトレーニングするために必要なパフォーマンスを提供します。Azure 上の新しい VM と、それらが AI 開発の取り組みにもたらすパフォーマンスの向上に興奮しています。"

Mustafa Suleyman, CEO, Inflection



175B

GPT-3 parameters



530B

Parameters in Megatron-Turing



5400

GPUs to advance Meta's AI research

New Azure ND H100 v5 VM series



NVIDIA H100 GPUを1000台規模まで
スケール可能



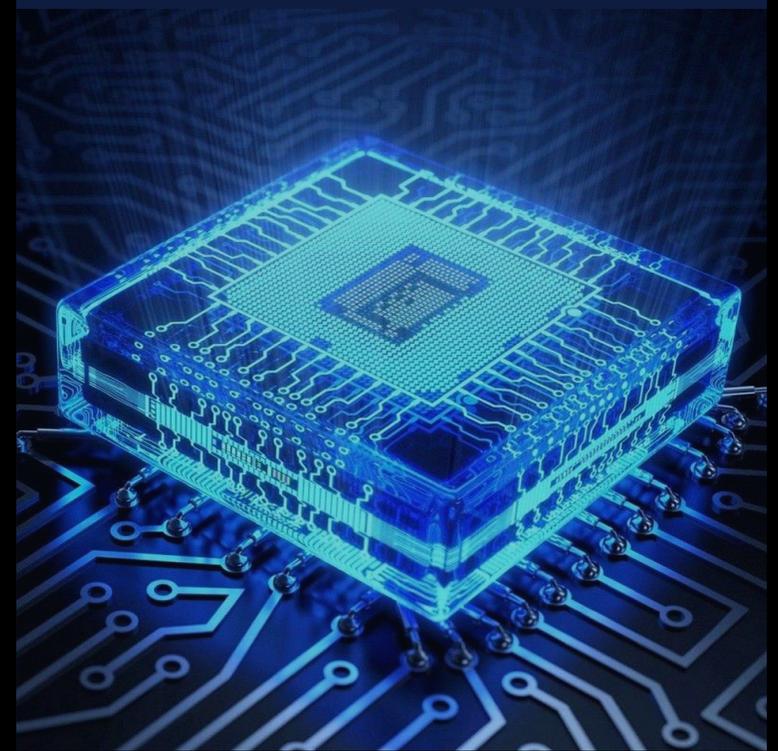
NVIDIA Quantum-2 InfiniBand で
ノード間接続 (NDR)



生成AIアプリケーション向けに設計

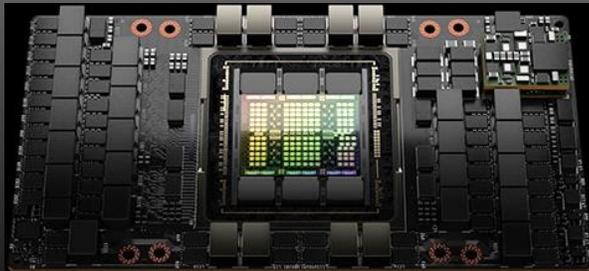


新しいクラスの大規模AIモデル対応



ND H100 v5 > 大規模スケーラブルAIスーパーコンピュータ

シングル H100 GPU



NCCL
+
NVLink

NVIDIA H100 Tensor Core GPU

- 80 GB of HBM3 Memory
- 2x – 30x A100 performance
- PCIe Gen 5, Intel Remote Host
- 8 per NDv5 VM

マルチ GPU

シングル ND_H100_v5 VM内の
NVLINKで接続された8GPU



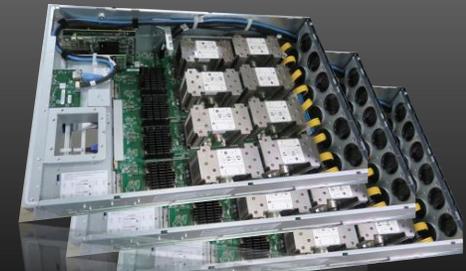
NVSwitch + NVLink 4.0

- Between 8 local GPUs within each VM
- 3.6 TB/s Bisection BW
- 450 GB/s AllReduce

マルチ GPU VM

Quantum-2 InfiniBandで接続された
複数台のND_H100_v5 VM

数百台のNDv5, 数千台のH100 GPUまでスケールアップ



NCCL
+
NDR

Quantum-2 InfiniBand

- 400 Gigabit dedicated link per GPU (3.2 Terabits/VM)
- Any to any, all to all, not over subscribed up to thousands of GPUs
- Dynamically provisioned via VMSS
- GPUDirect RDMA

ND H100 v5

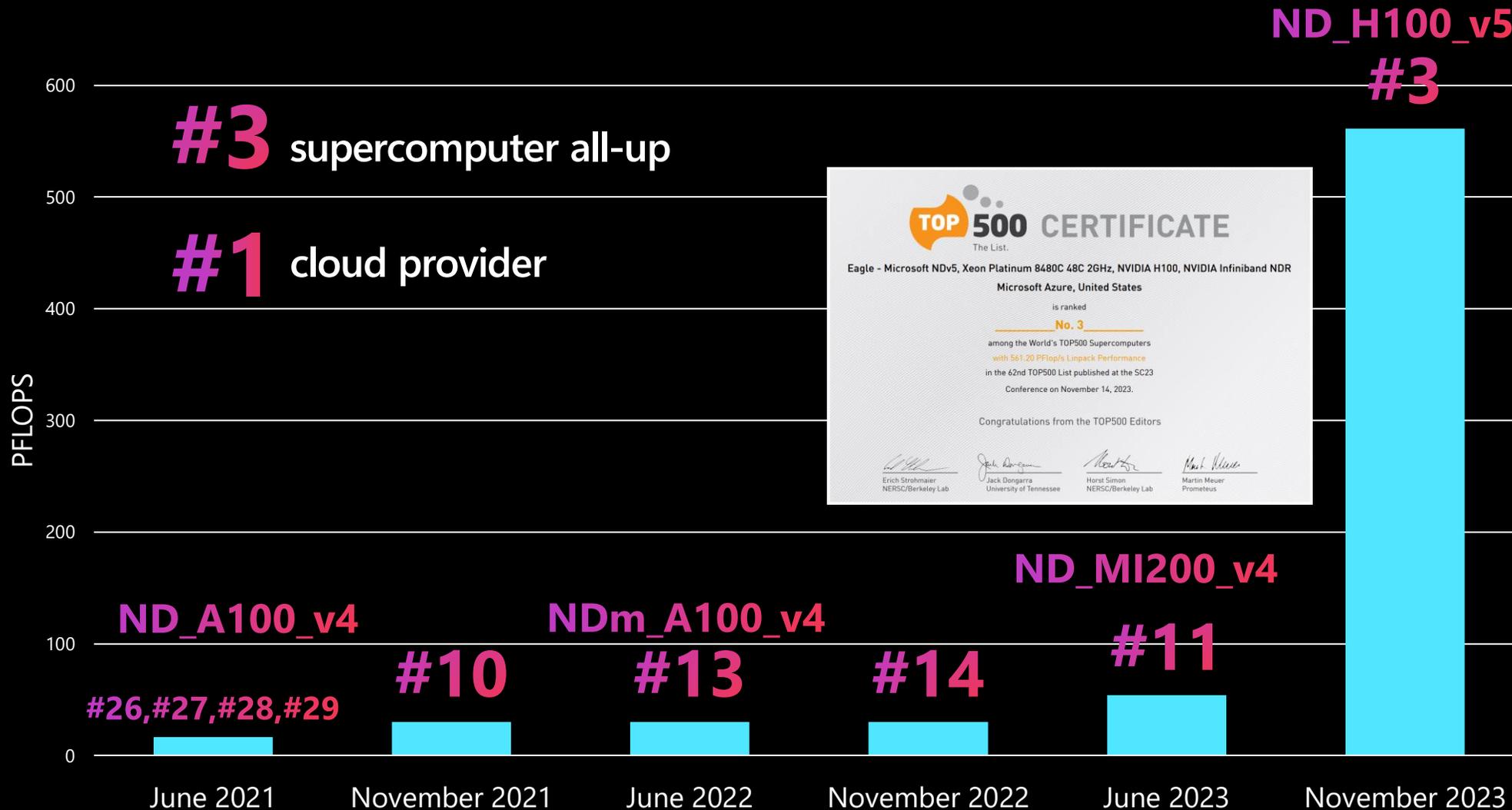
ND H100 v5 シリーズの仮想マシンは、ハイエンドのディープラーニングの学習と密結合のスケールアップおよびスケールアウトの生成AIおよびHPCワークロード向けに設計されています。ノードあたりNVIDIA H100 × 8で構成されており、3.2Tb/sのMellanox HDR InfiniBand × 8を用いて数千GPUまでスケールアップしてデプロイすることが可能です。多くのAI、MLのフレームワークを活用して優れたパフォーマンスを発揮し、さらにシームレスなGPUクラスタリングのためNVIDIAのNCCL通信ライブラリに対応したAIおよびHPCツールによって、InfiniBandインターコネクで優れたスケラビリティを実現します。

- ✓ NVIDIA H100 × 8 (NVLINK) 搭載
- ✓ Intel 第4世代 Xeon Scalable Processor (Sapphire Rapids) (96コア/ノード) 搭載
- ✓ 400Gbps NDR InfiniBand × 8 で ノードあたり3.2Tb/sの相互接続帯域幅を提供



	
CPU	4 th Gen Intel Xeon Scalable Processor (Sapphire Rapids)
コア数	96
GPU	8 x NVIDIA H100 (next gen NVSwitch and NVLink 4.0)
メモリ容量	1900 GiB (DDR5 DIMMS)
ローカルディスク	1000 GiB SSD
InfiniBand	400 Gbps NDR InfiniBand (NVIDIA Quantum-2 CX7) x 8 (3.2Tb/s per VM in a non-blocking fat-tree network)

Azure ND_H100_v5がTop500で3位にランクイン



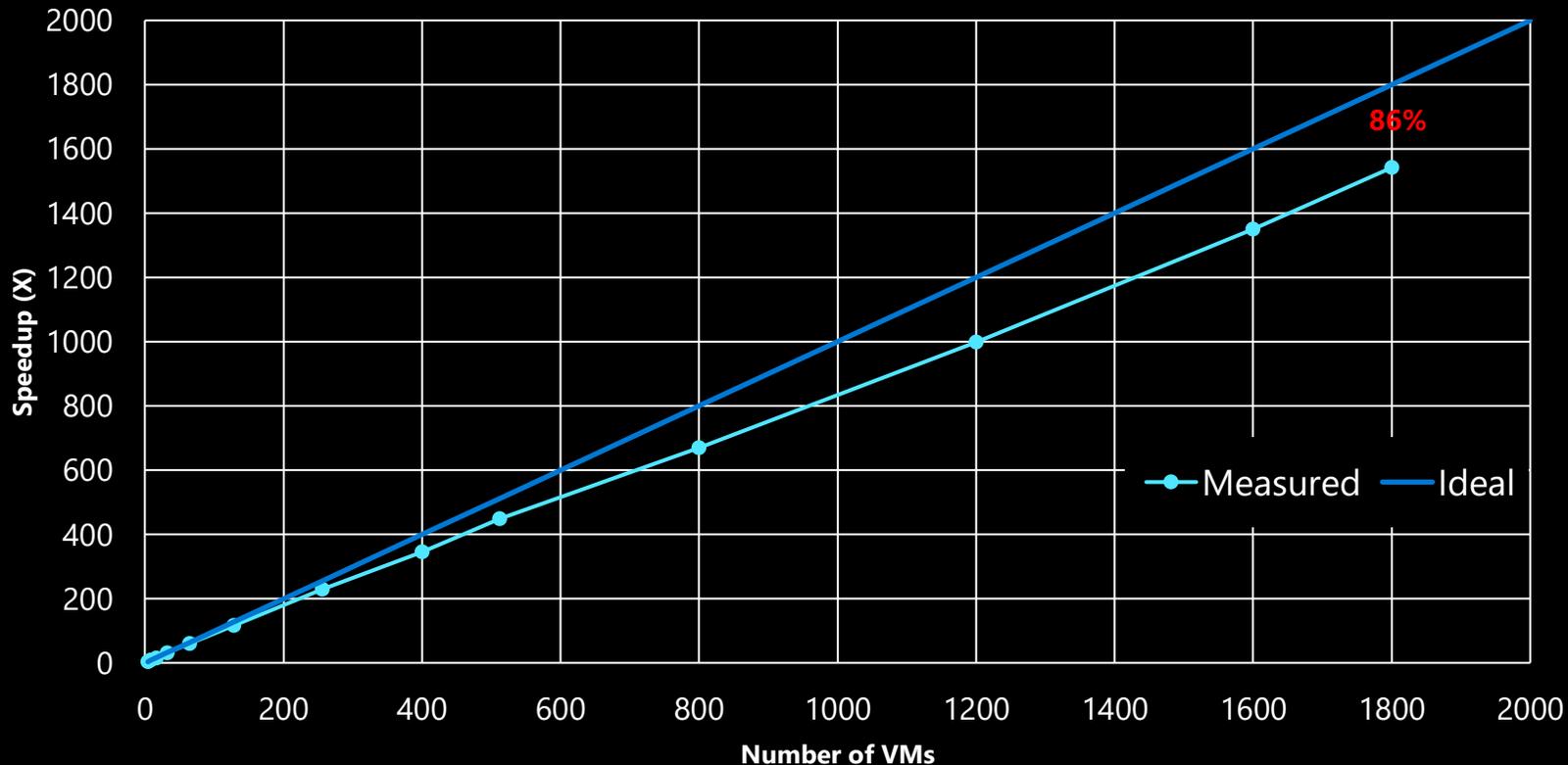
2023

Microsoft Azure

Ranking #3 in the #Top500 Supercomputing list

14,400 GPUs におけるスケーリング効果 (HPL)

HPL Scaling efficiency with 14,400 H100 GPUs



32 から 14,400 の NVIDIA H100 Tensor コア GPU で 86% のスケーリング効率

2023

TOP 500
The List.

Microsoft Azure

Ranking #3 in the #Top500
Supercomputing list

#3

MLPerf Training v3.1 benchmark

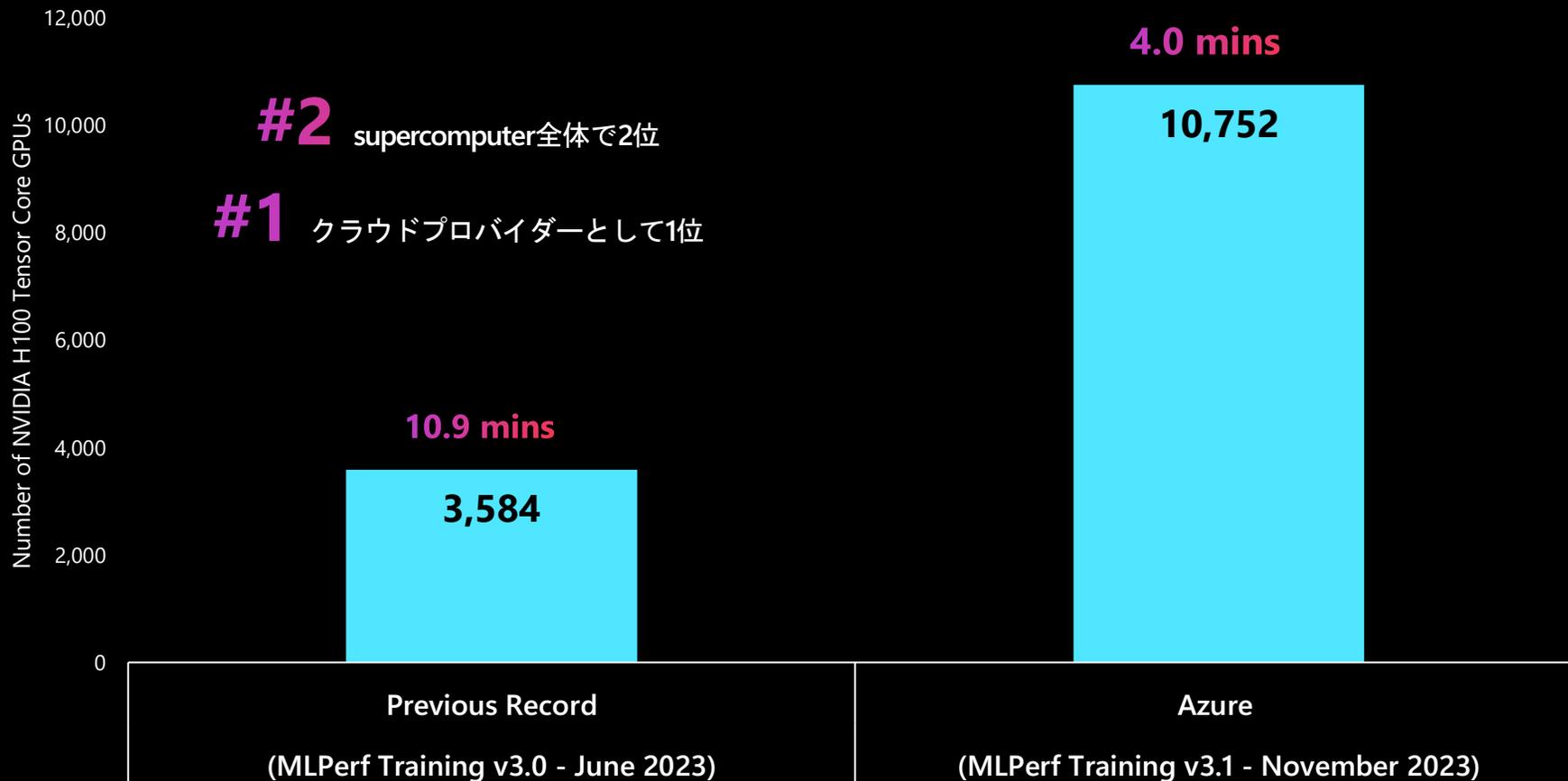
**MLPerf-LLM 175B training scale record by
Azure ND H100 v5-series**
as of November 2023

2023

MLPerf
on
Microsoft Azure

#1 highest performing cloud provider in the world

#2 highest performing supercomputer in the world



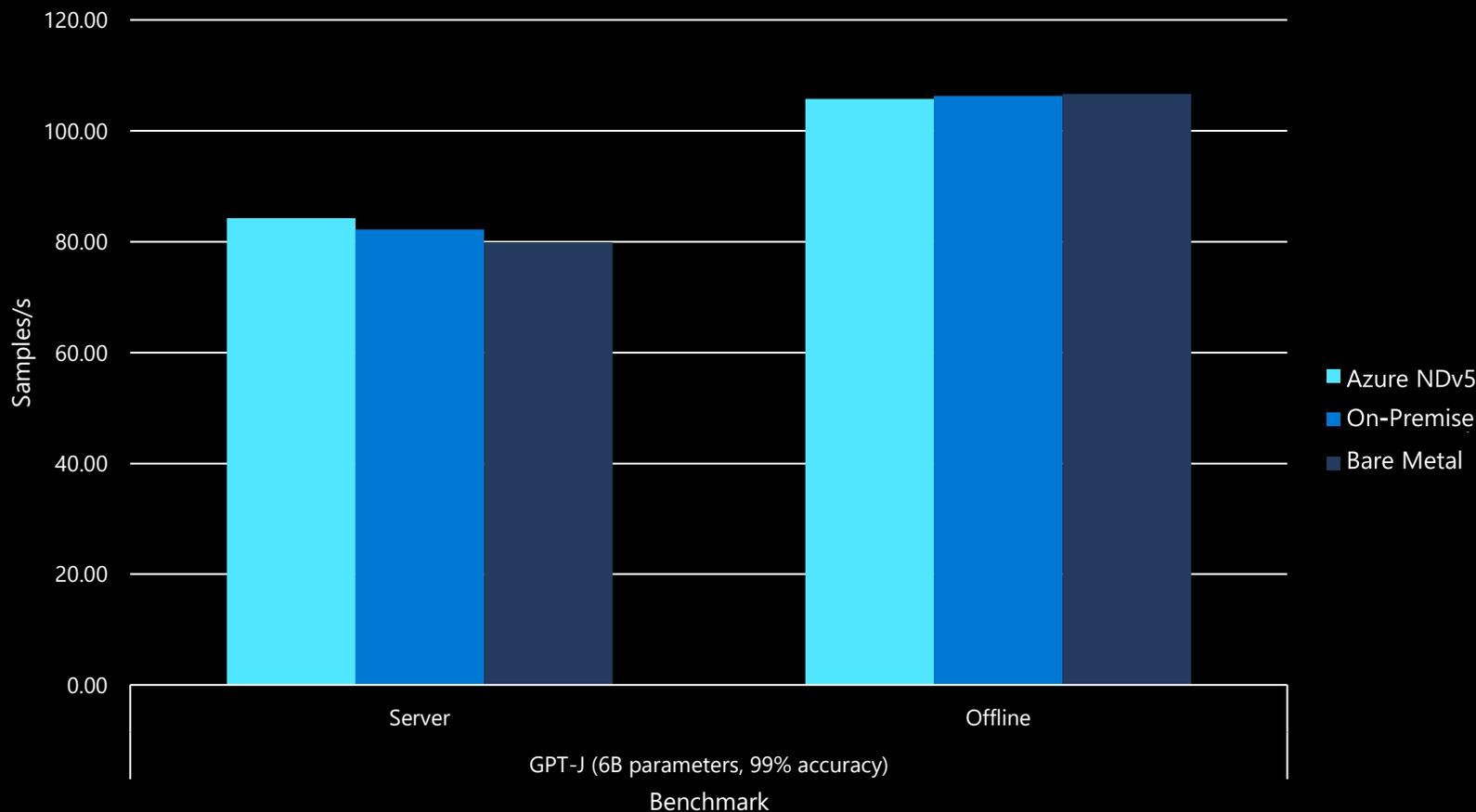
Azure's submission, the largest in the history of MLPerf Training, demonstrates the extraordinary progress we have made in optimizing the scale of training." said David Kanter, Executive Director of MLCommons

aka.ms/AzureBlog/MLPerf3.1

MLPerf Inference v3.1 benchmark

Performance on LLM GPT-J - MLPerf Inference v3.1

as of September 2023



オンプレミス、ベアメタルと同等性能

0.99x-1.05x relative performance

NC_H100_v5 Preview

NCads H100 v5 シリーズは、NVIDIA H100 NVL GPU と第 4 世代 AMD EPYC™ Genoa プロセッサを搭載しています。この VM には、最大 2 個の NVIDIA H100 NVL GPU (それぞれに 94 GB のメモリを装備)、最大 80 個の非マルチスレッド AMD EPYC Milan プロセッサコア、640 GiB のシステムメモリが搭載されています。これらの VM は、次のような実際の Applied AI ワークロードに最適です。

- ✓ GPU で高速化された分析とデータベース
- ✓ 大量の前処理と後処理があるバッチ推論
- ✓ 自律性モデルのトレーニング
- ✓ 石油とガスの貯留層シミュレーション
- ✓ 機械学習 (ML) 開発
- ✓ ビデオの処理
- ✓ AI/ML Web サービス

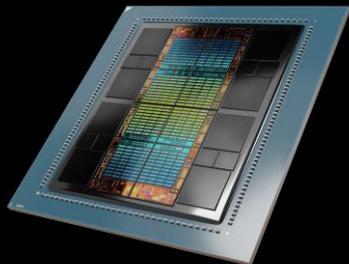


	
CPU	4 th Gen AMD EPYC Genona
コア数	40/ 80
GPU	1 or 2 x NVIDIA H100 NVL (NVLink 4.0)
GPUメモリ	94/ 188 GiB
メモリ容量	320 / 640 GiB (DDR5 DIMMS)
ローカルディスク	3576 / 7152 GiB SSD

ND_MI300X_v5 Preview

ND_MI300X_v5シリーズは、厳しいAIとHPCのワークロードに最適化された仮想マシンです。AMD Instinct MI300X GPUを8基と第4世代 Intel Xeon Scalableプロセッサプロセッサを搭載しています。また、GPUあたり1本の400GbpsのNVIDIA Quantum-2 CX7 (4x HDR InfiniBand)、ノードあたり3.2Tbpsのスループット性能のインターコネクタで大規模にスケールアップしてデプロイすることが可能です。

- ✓ AMD Instinct MI300X × 8 (Infinity Fabric 3.0) 搭載
- ✓ Intel 第4世代 Xeon Scalable Processor 搭載
- ✓ 400Gbps NDR InfiniBand × 8 で ノードあたり3.2Tb/sの相互接続帯域幅を提供

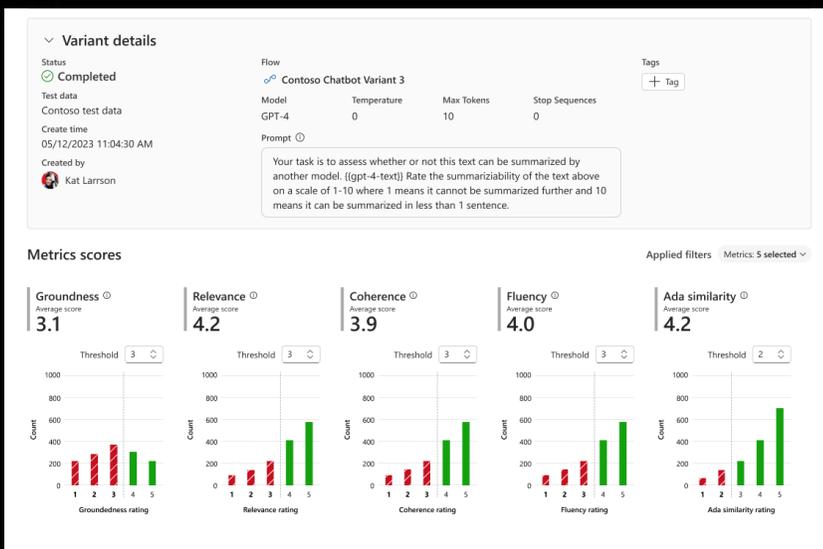


	
CPU	4 th Gen Intel Xeon Scalable Processor
コア数	??
GPU	8 x AMD Instinct MI300X (Infinity Fabric 3.0)
メモリ容量	??? GiB (16Ch. DDR5 DIMMS)
ローカルディスク	???? GiB SSD
InfiniBand	400 Gbps 4x NDR InfiniBand (NVIDIA Quantum-2 CX7) x 8 (3.2Tb/s per VM)

生成 AI を支援するさまざまなサービス (Azure AI Studio, Azure Machine Learning, Azure CycleCloud)

- 大量ノードの構成や起動/停止
- 推論エンドポイントの管理
- スクリプトによる自動化・トリガ実行
- GPU リソースなどの権限管理
- 生成 AI 専用のモデル評価と監視
- 生成 AI アプリケーションのフロー構築

Microsoft Azure Machine Learning Studio interface showing a Bing Grounded QA flow. The flow includes nodes for 'extract_query_from_quest...', 'search_on_bing', 'process_search_result', and 'augmented_qna'. The 'process_search_result' node contains Python code for processing search results. The 'augmented_qna' node shows deployment settings like 'max_tokens: 256' and 'temperature: 0.5'.



Index	Input	Expected response	Output	Groundedness	Relevance	Reasoning
1	Can you tell me if I can return this ABC I bought 2 days ago?	Your warranty for ABC product is 90 days.	The warranty for ABC line of products is 60 days.	1	4	Your warranty for ABC product is 90 days but the output is 60 days.
2	Can you tell me if I can return this ABC I bought 2 days ago?	Your warranty for ABC product is 90 days.	The warranty for ABC line of products is 60 days.	4	5	Your warranty for ABC product is 90 days but the output is 60 days.
3	Can you tell me if I can return this ABC I bought 2 days ago?	Your warranty for ABC product is 90 days.	The warranty for ABC line of products is 60 days.	1	5	Your warranty for ABC product is 90 days but the output is 60 days.
4	Can you tell me if I can return this ABC I bought 2 days ago?	Your warranty for ABC product is 90 days.	The warranty for ABC line of products is 60 days.	4	5	Your warranty for ABC product is 90 days but the output is 60 days.
5	Can you tell me if I can return this ABC I bought 2 days ago?	Your warranty for ABC product is 90 days.	The warranty for ABC line of products is 60 days.	1	4	Your warranty for ABC product is 90 days but the output is 60 days.

まとめ



生成 AI & クラウド HPC を支える Azure HPC インフラ

✓ 仮想マシン

- ✓ ND_H100_v5, NDm_A100_v4, ND_A100_v4
- ✓ ND_MI300X_v5 (Preview)
- ✓ NC_H100_v5 (Preview)
- ✓ GPU / InfiniBand

✓ ストレージ

- ✓ Azure Managed Lustre Filesystem

✓ 生成 AI を支援するサービス

- ✓ Azure AI Studio
- ✓ Azure Machine Learning
- ✓ Azure CycleCloud

Microsoft Azure で

生成 AI もクラウド HPC も対応可能！

■ 第3回 Azure HPC ユーザー会 (計画中)

■ 営業お問い合わせ先

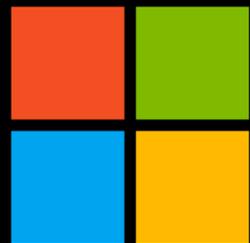
Tel: 0120-167-400 時間: 9:00 - 17:30 (月～金)

URL: <https://azure.microsoft.com/contact>



お問い合わせはこちら





Microsoft

- 本書に記載した情報は、本書各項目に関する発行日現在の Microsoft の見解を表明するものです。Microsoftは絶えず変化する市場に対応しなければならないため、ここに記載した情報に対していかなる責務を負うものではなく、提示された情報の信憑性については保証できません。
- 本書は情報提供のみを目的としています。Microsoft は、明示的または暗示的を問わず、本書にいかなる保証も与えるものではありません。
- すべての当該著作権法を遵守することはお客様の責務です。Microsoftの書面による明確な許可なく、本書の如何なる部分についても、転載や検索システムへの格納または挿入を行うことは、どのような形式または手段（電子的、機械的、複写、レコーディング、その他）、および目的であっても禁じられています。これらは著作権保護された権利を制限するものではありません。
- Microsoftは、本書の内容を保護する特許、特許出願書、商標、著作権、またはその他の知的財産権を保有する場合があります。Microsoftから書面によるライセンス契約が明確に供給される場合を除いて、本書の提供はこれらの特許、商標、著作権、またはその他の知的財産へのライセンスを与えるものではありません。

© 2024 Microsoft Corporation. All rights reserved.

Microsoft, Windows, その他本文中に登場した各製品名は、Microsoft Corporation の米国およびその他の国における登録商標または商標です。

その他、記載されている会社名および製品名は、一般に各社の商標です。