# Azure Al Infrastructure

### Why Al Infrastructure **Matters**

Al is changing the way the world works, and advancements in AI infrastructure are at the center of that revolution. Al infrastructure powers your AI innovations—from optimizing hardware, processors, and accelerators, to networking and virtualization.

Requirements for AI infrastructure are advancing in two fundamental ways. First, the ever-increasing complexity of AI workloads requires having an infrastructure that's optimized to provide massive scale with parallel computation. Second, the increasing availability of AI infrastructure in the cloud is making AI innovation more accessible to organizations of all sizes.

## Accelerate your **AI** innovations with Al infrastructure

## ∠X throughput¹ Per GPU compared to competitors

performance<sup>2</sup> among Top 500 supercomputing rankings

for cloud

Performance at scale

Azure AI infrastructure provides access to large-scale AI models, GPU-based virtual machines (VMs), AI supercomputing, and Al services, allowing customers to build, train, test, and deploy their AI applications with speed and efficiency.

### Microsoft is on a continuous journey to

Microsoft: optimized for Al

revolutionize computing advancements by leading the development of AI-optimized infrastructure.

Customers that want to lead in the AI era

need cloud computing infrastructure optimized specifically for AI workloads. With highly scalable interconnected GPUs powering virtual machines (VMs), Azure Al infrastructure helps remove bottlenecks and deliver rapid data processing throughput and low latency—a must for Al.

services, like Azure Machine Learning, to operationalize models with MLOps seamlessly, and Azure AI services for pre-trained foundation models. It also natively integrates with open-source frameworks and optimization technologies such as PyTorch and ONNX Runtime. This integration allows developers and data scientists to build and deploy high-performing AI models quickly and easily.

The infrastructure is engineered for Azure Al



**Azure Al** 











**Open source** 



### Azure Al infrastructure offers proven supercomputing performance for the most advanced Al initiatives. It supports cross-platform interoperability and collaboration, letting you use

Proven workload performance

their preferred tools and languages across different environments and devices.

Microsoft runs on Azure Al infrastructure



## Microsoft Edge



WAYVE

Azure supercomputing

is enabling **Al-based** 

in Microsoft Teams



in Microsoft 365

## Partnerships that power ingenuity



research, products,

and API services.

**S**OpenAI

Azure is the exclusive

cloud provider for

**Purpose-built for innovation** 



Your considering AI as part of your innovation strategy, and that means you'll need secure,

GPU-based VMs that can handle different workloads and scenarios, including training large

high-performance infrastructure to support those workloads. Azure AI offers a range of

the flexibility required to strike the optimal balance between cost and performance.

Al models and running high-performance computing (HPC) simulations.



**Elekta** 

Azure HPC provides the

agility to scale storage

### From current and future-generation transformer-based models for natural language processing (NLP), to scalable implementations for midrange training projects and inference, Azure provides

A track record you can trust

### 2020 2021

and services. Here's a look at some of the most important milestones from the last few years:

Azure has been scaling supercomputing capabilities since 2020, deploying new infrastructure

**Built for OpenAl** 

Microsoft

ranks **Top 5** in supercomputing<sup>3</sup>

Microsoft announces Al supercomputing VMs available to customers

Microsoft

ranks **#1** for

First in cloud

Azure delivered #1 performance by

a cloud provider

in MLPerf results<sup>5</sup>

Microsoft partners

with NVIDIA on

2022

**Reaching new** 

Microsoft runs a 530B-parameter

supercomputing capabilities supercomputing milestones performance

> supercomputing<sup>4</sup> latest GPUs

> > NeMo Megatron benchmark<sup>6</sup> on **175** VMs

Azure OpenAl Service and the new Bing trained

Microsoft

announces

2023

**Extending Al** 

with Azure Al infrastructure

Connect with Azure Sales to learn more >

Take the next steps

<sup>1</sup>Microsoft performance benchmark for 1-Trillion parameter model with DeepSpeed, compared to another published cloud vendor running 1-Trillion parameter model with DeepSpeed optimized for throughput. As of July 2022 <sup>2</sup> Microsoft Azure ranks #10 on supercomputing performance and #1 for cloud-based performance—Top500 List. Published November 2022 <sup>3</sup> Microsoft supercomputer for OpenAI, completed in May 2020, was a single system with more than 285,000 CPU cores, 10,000 GPUs and 400 gigabits per second of network connectivity for each GPU server. Compared with other machines listed on the TOP500 supercomputers in May 2020, ranks in the top five.

<sup>5</sup> MLPerf 1.1 results show a debut performance by Azure delivering the #2 performance overall and the #1 performance by a cloud provider. <sup>6</sup> Microsoft Azure published <u>NVIDIA NeMo Megatron benchmarking</u> demonstrating Scale to 530B Parameter GPT-3 Model ©2023 Microsoft Corporation. All rights reserved. This document is for informational purposes only. Microsoft makes no warranties,

express or implied, with respect to the information presented. You may copy and use this document for your internal, reference purposes.

<sup>4</sup>Microsoft Azure ranks #10 on supercomputing performance and #1 for cloud-based performance—<u>Top500 List</u> November 2021