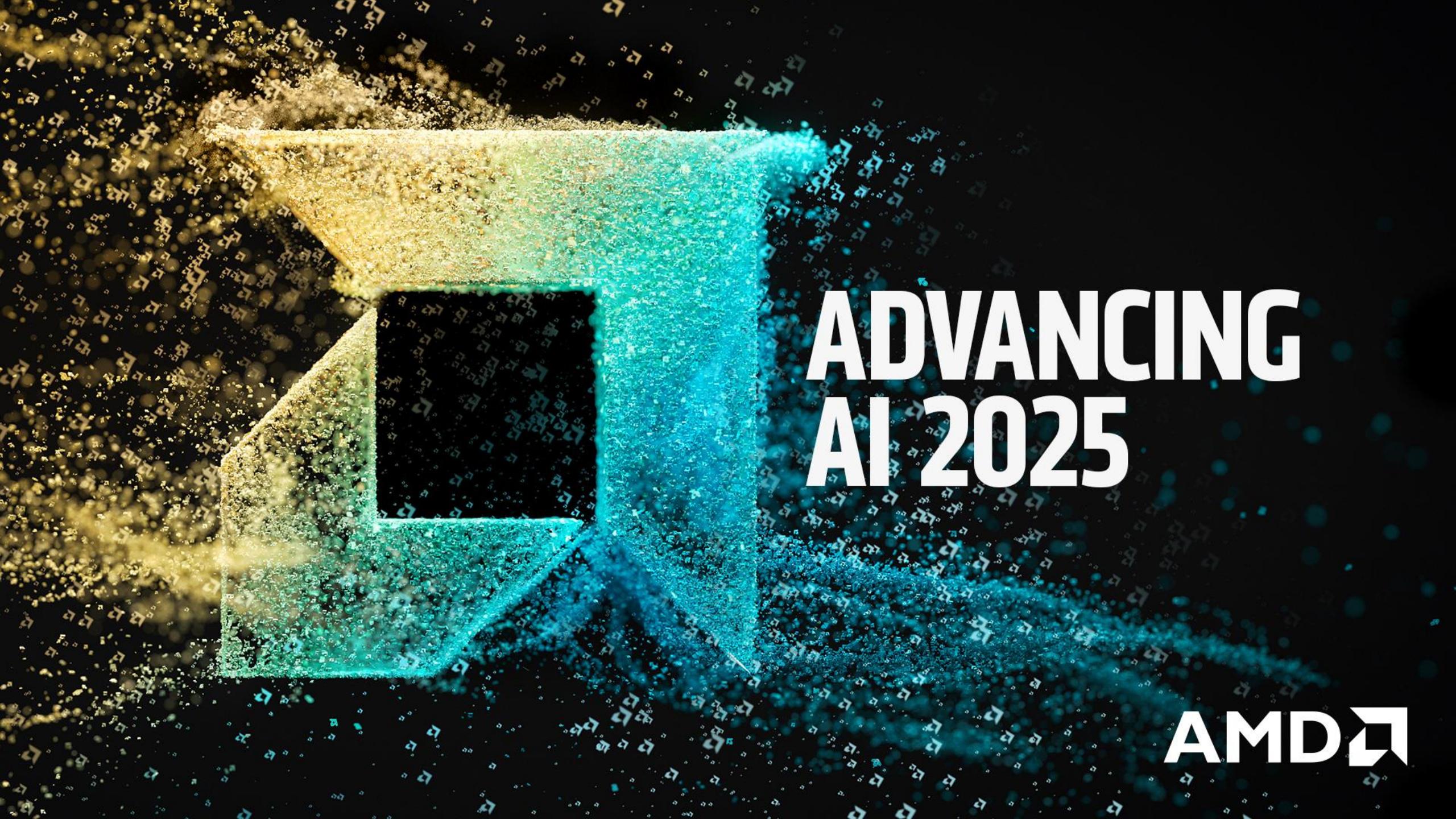


CAUTIONARY STATEMENT

This presentation contains forward-looking statements concerning Advanced Micro Devices, Inc. (AMD) such as the features, functionality, performance, availability, timing and expected benefits of AMD products and product roadmaps including the AMD Instinct™ MI350 Series, AMD ROCm 7 and the associated products with AMD "Helios" and Next-Gen AI Racks; projected data center AI accelerator TAM in 2028; AMD's AI strategy; and AMD's ability to accelerate momentum across its AI capabilities, which are made pursuant to the Safe Harbor provisions of the Private Securities Litigation Reform Act of 1995. Forward-looking statements are commonly identified by words such as "would," "may," "expects," "believes," "plans," "intends," "projects" and other terms with similar meaning. Investors are cautioned that the forward-looking statements in this presentation are based on current beliefs, assumptions and expectations, speak only as of the date of this presentation and involve risks and uncertainties that could cause actual results to differ materially from current expectations. Such statements are subject to certain known and unknown risks and uncertainties, many of which are difficult to predict and generally beyond AMD's control, that could cause actual results and other future events to differ materially from those expressed in, or implied or projected by, the forward-looking information and statements. Investors are urged to review in detail the risks and uncertainties in AMD's Securities and Exchange Commission filings, including but not limited to AMD's most recent reports on Forms 10-K and 10-Q.

AMD does not assume, and hereby disclaims, any obligation to update forward-looking statements made in this presentation, except as may be required by law.



ADVANCING AI 2025

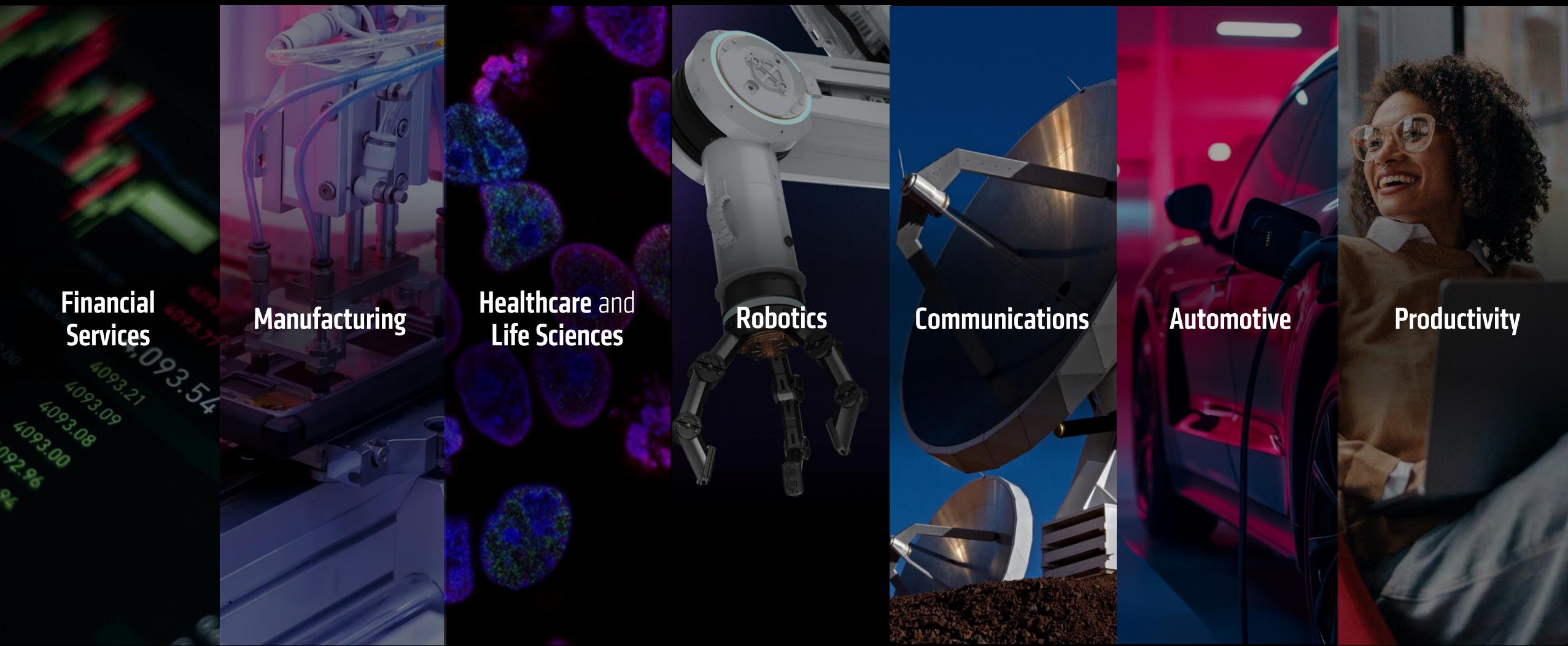
AMD



High Performance & Adaptive Computing
Solving the World's Most Important Challenges



AI Everywhere, for Everyone



Financial
Services

Manufacturing

Healthcare and
Life Sciences

Robotics

Communications

Automotive

Productivity

AI Innovation is Accelerating



Training is Evolving



Inference Scaling
Accelerates



Explosion of Models



Reasoning &
Agents Surge

Reasoning & Agents Fuel Compute Surge



Leadership GPU
Lowers TCO



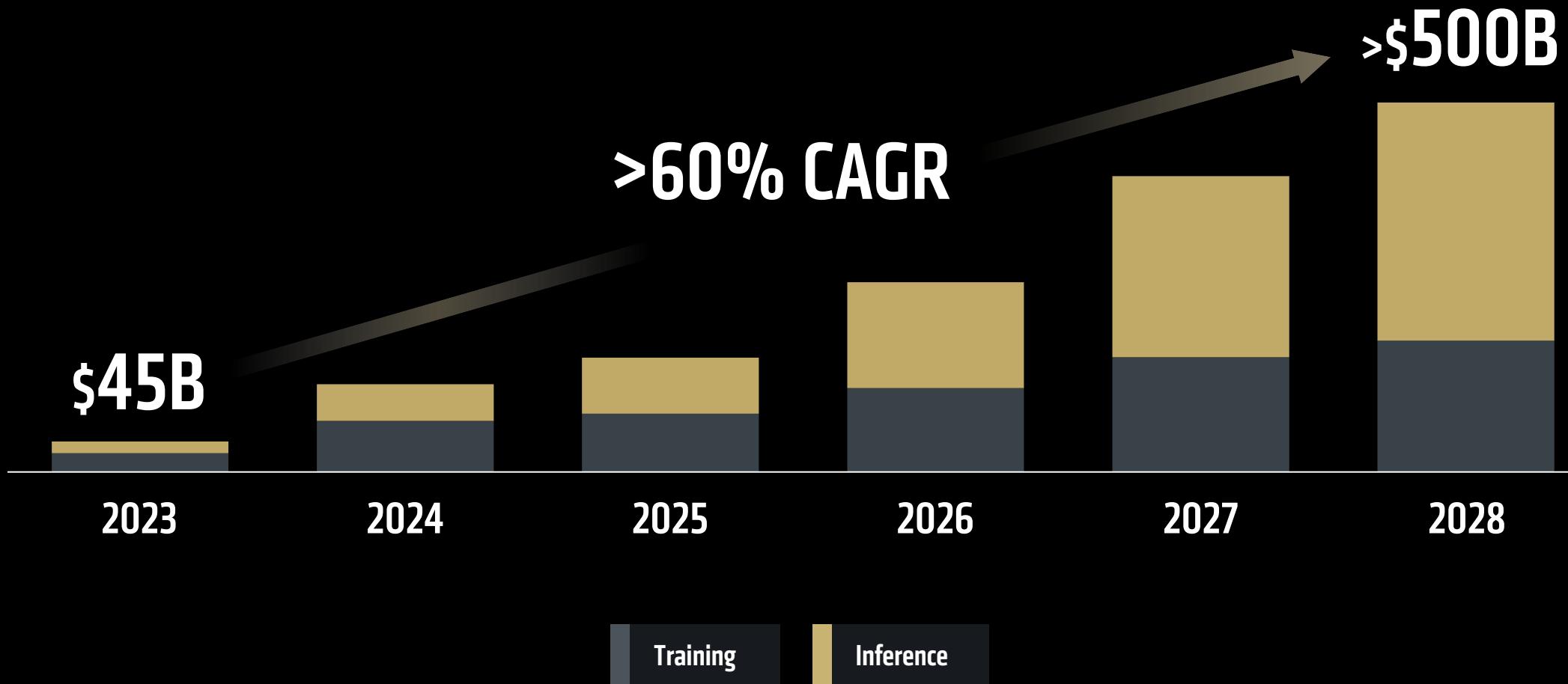
Leadership CPU
Powers apps



Openness
Accelerates Innovation

Data Center AI Accelerator TAM

Inference Growing at >80% CAGR



AI Beyond the Data Center

Inference Scaling Across Cloud to Edge to Client



Cloud



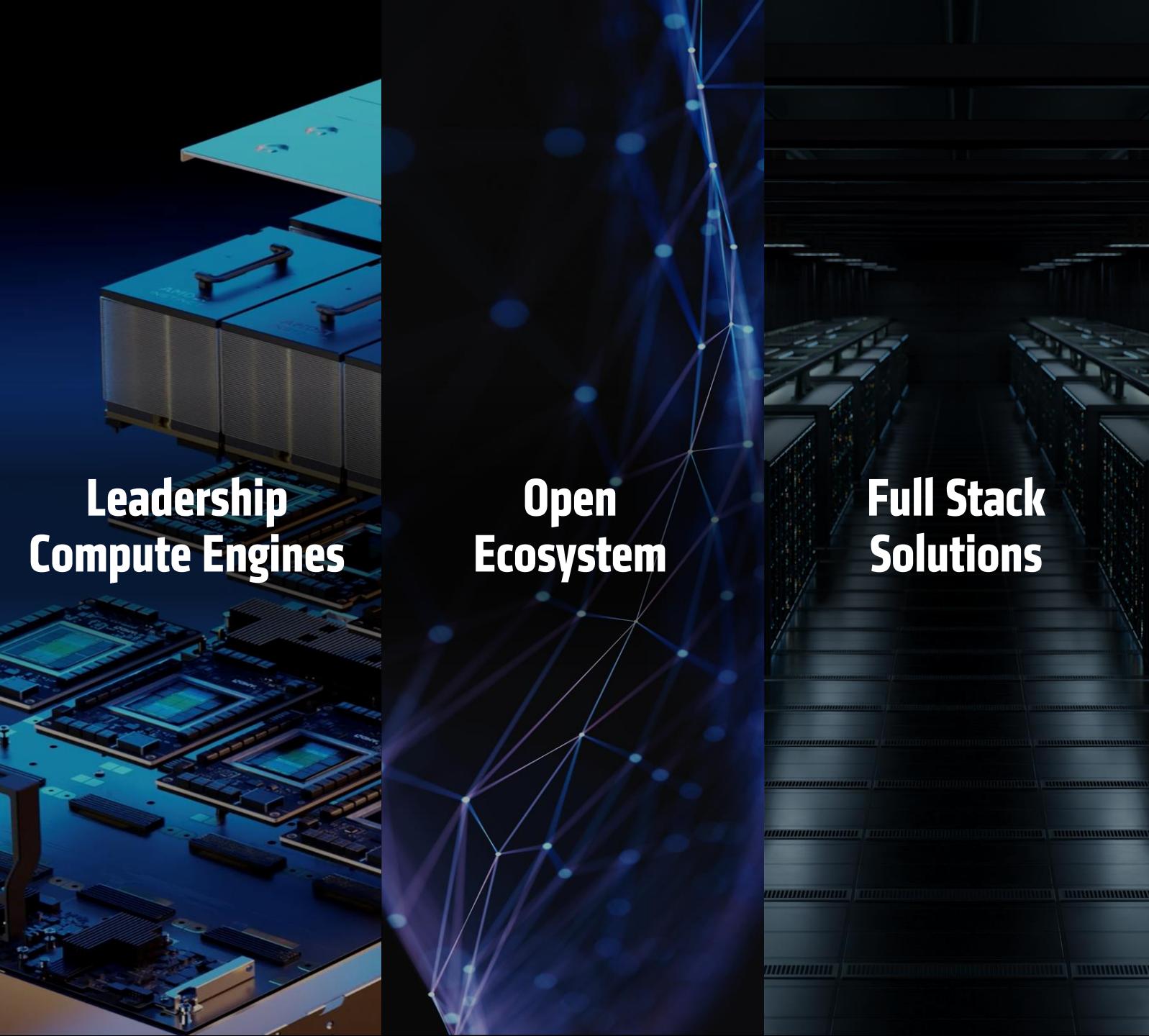
Edge



Client

Driven by domain-specific compute engines & open software stack

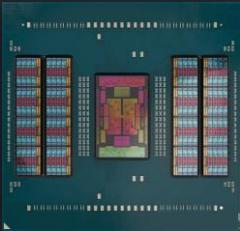
AMD AI Strategy





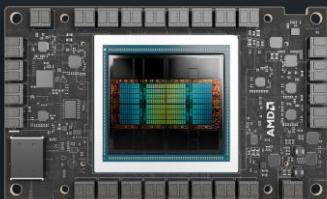
Best End-to-End AI Compute Portfolio in the Industry

AMD EPYC™ Processors



Leading server CPU

AMD Instinct™ Accelerators



World's best GPU accelerator

AMD Pensando™ Networking



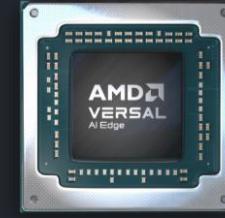
Premier programmable
DPUs & AI NICs

AMD Ryzen™ AI AMD Radeon™ AI Processors



Most powerful client
AI processors

AMD Versal™ Adaptive SOCs



Leadership AI
Processing at the edge

Open Development Drives Value & Innovation

Open Hardware



Open Software



Open Ecosystem



OPEN
Compute Project™



Ultra Ethernet
Consortium



Hugging Face



PyTorch



Choice

Flexibility

Rapid Co-Innovation

Portability

Proven

Investing in Full-Stack Solutions

Acquisitions Span Entire AI Value Chain



ENOSEMI

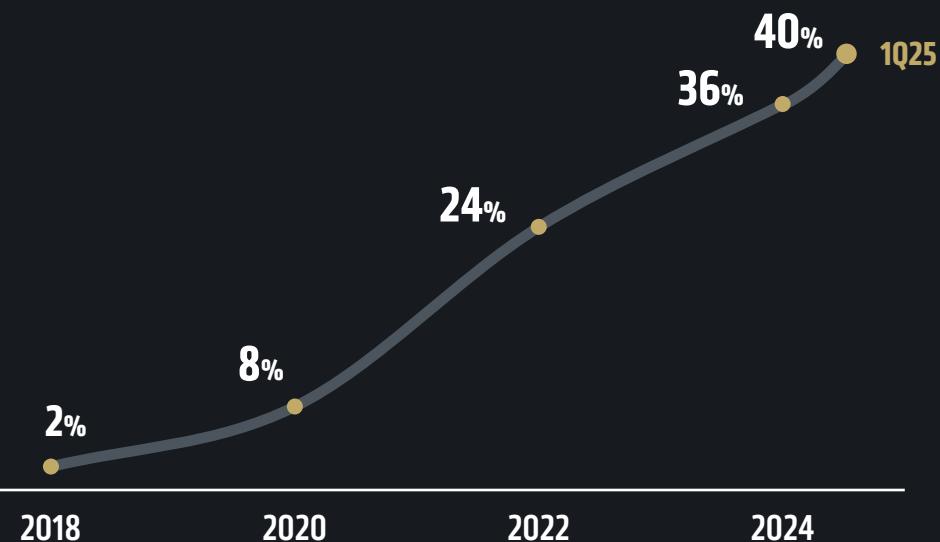
PENSANDO

XILINX.

Over 25 AI Acquisitions & Investments in the Last Year Alone

EPYC Momentum Accelerates...

>18x Server CPU Market Share Growth



Industry Leaders Run on EPYC™

Cloud



Digital



Enterprise



OEM





Delivering on Leadership GPU Commitment

Powering Top Supercomputers

Trusted by Leading AI Practitioners



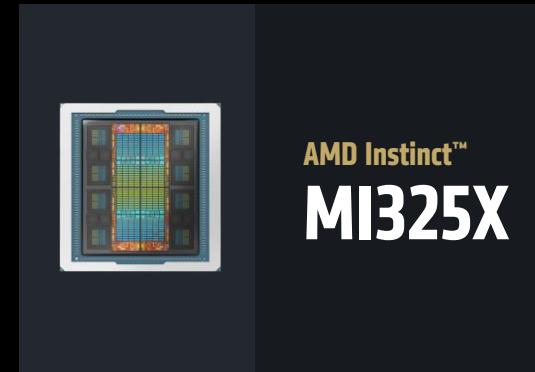
AMD Instinct™
MI250X



AMD Instinct™
MI300A



AMD Instinct™
MI300X



AMD Instinct™
MI325X

2021

2025



Growing Industry Adoption

7 of 10 Largest AI Companies Use AMD Instinct





Rapidly Advancing Open Software Capabilities

**Day-0 Support
for Leading Models**

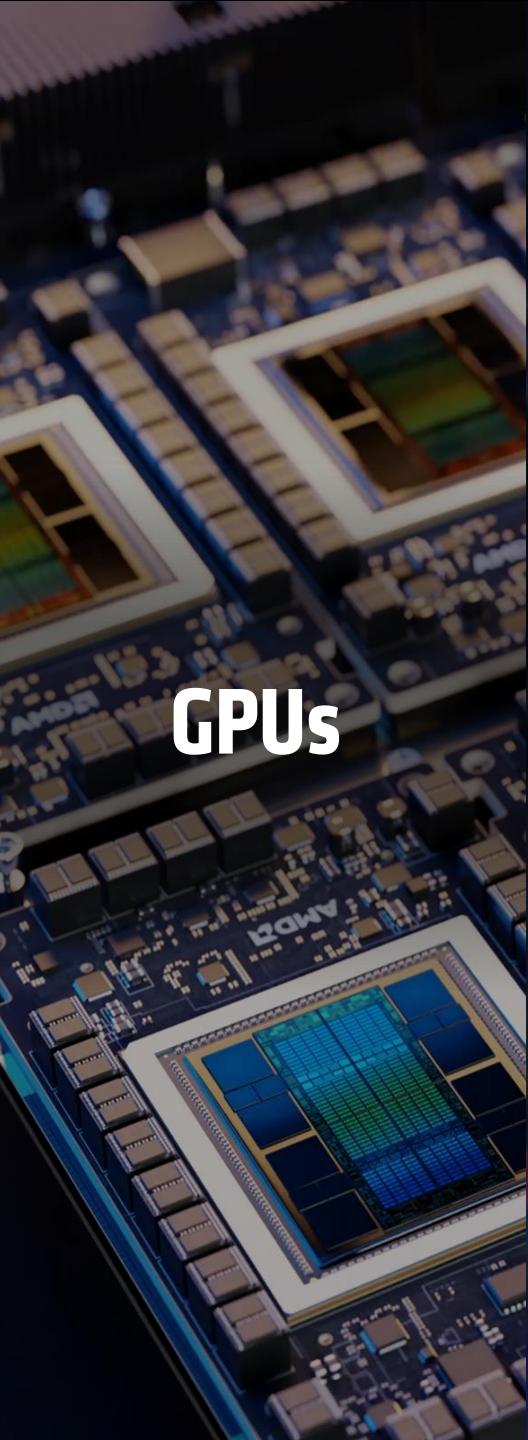
**Accelerated
Pace of Innovation**

**Broadening
Ecosystem Partnerships**

**Developer First
Approach to Enablement**



Today at **Advancing AI**



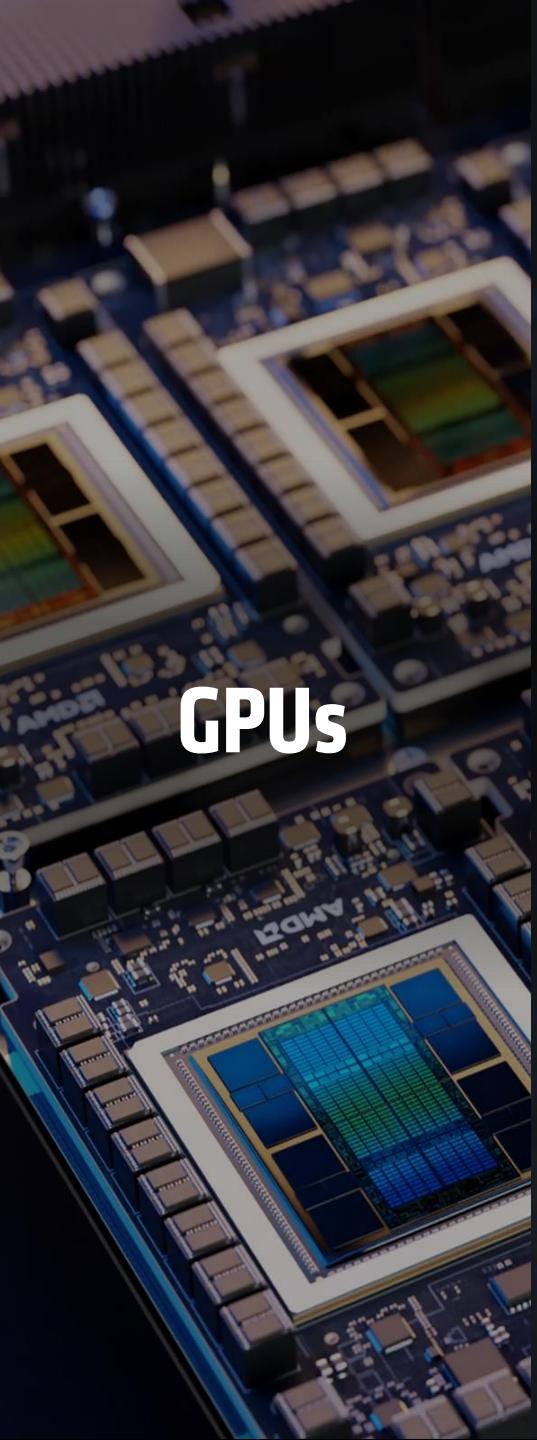
Software

```
You, 7 months ago | 1 author
import VueRouter from 'vue-router'
import routes from './routes'
import store from './store'
import vuexI18n from 'vuex-i18n'
import enLangFile
```

Solutions



Today at **Advancing AI**



GPUs

Software

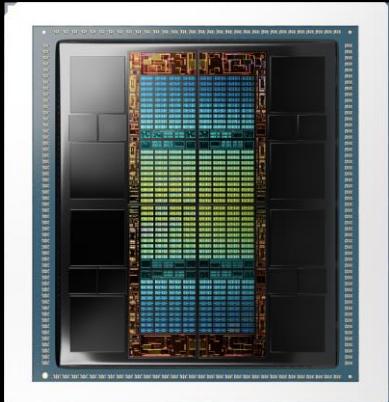
Solutions

You, 7 months ago | 1 author
import VueRouter from 'vue-router'
import routes from './routes'
import store from './store'
import vuexI18n from 'vuex-i18n'
import enLangFile from 'enLangFile'



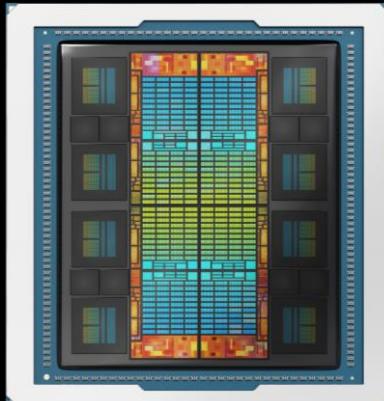
Delivering on Annual Roadmap Commitment

AMD Instinct™
MI300A/X



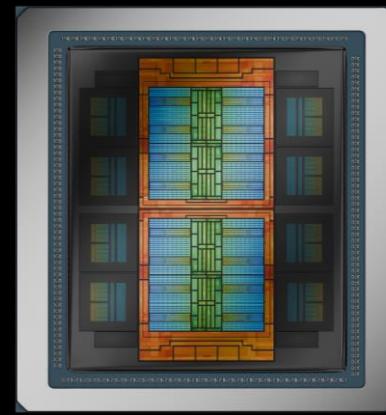
2023

AMD Instinct™
MI325X



2024

AMD Instinct™
MI350 SERIES



2025

AMD Instinct™
MI400 SERIES



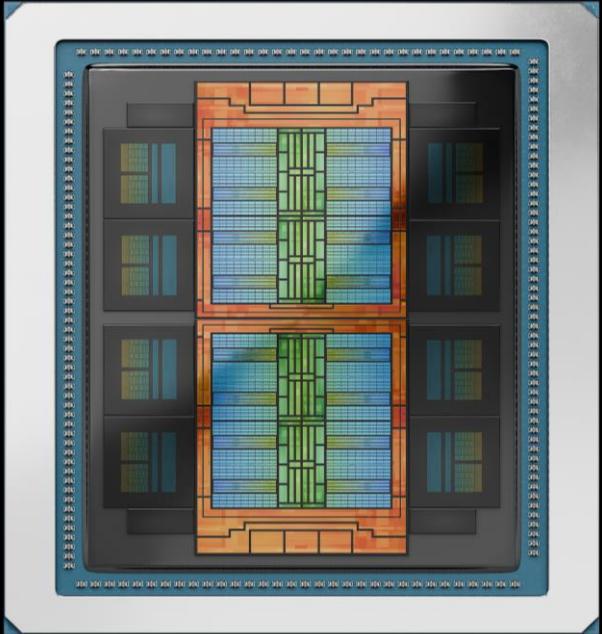
2026

Roadmap subject to change

Launching Today at **Advancing AI 2025**

AMD Instinct™ MI350 Series

Continued Generative AI Leadership



AMD Instinct™ **MI350 Series**

4th Gen Instinct™ architecture

3nm process node

185 billion transistors

FP4 & FP6 new gen AI datatypes

HBM3E leadership capacity

MI350 Series Accelerates Your Gen AI Outcomes

Faster AI Inference & Training

20PF

FP4 & FP6

4x Gen-on-Gen AI
Compute Increase

Larger AI Model Support

288GB

HBM3E

Supports up to 520B
Parameter AI Model

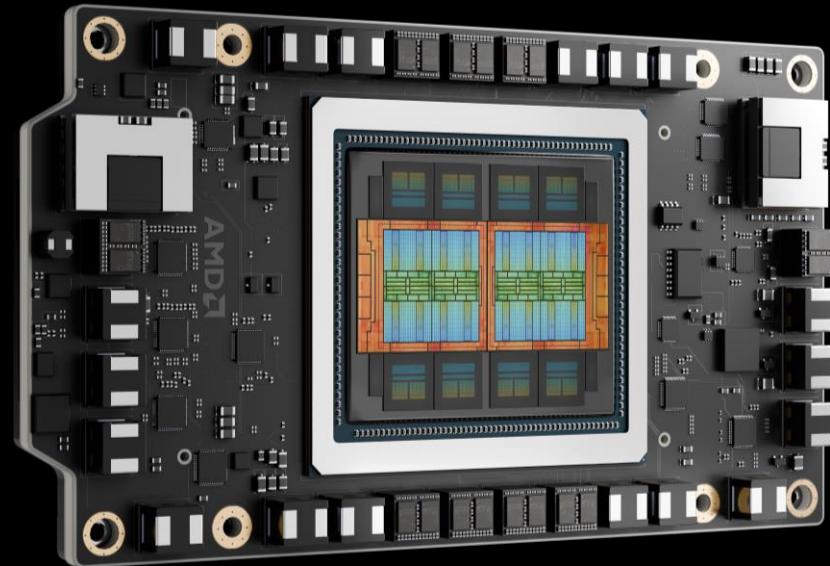
Rapid AI Infrastructure Deployment

UBB8

Industry Standard GPU Node

Available in Air Cooled
& Direct Liquid Cooled

AMD Instinct™ MI350 Series



Instinct™ MI355X

MEMORY **288 GB HBM3E**

MEMORY BANDWIDTH **8 TB/s**

FP64 **79 TF**

FP16 **5 PF**

FP8 **10 PF**

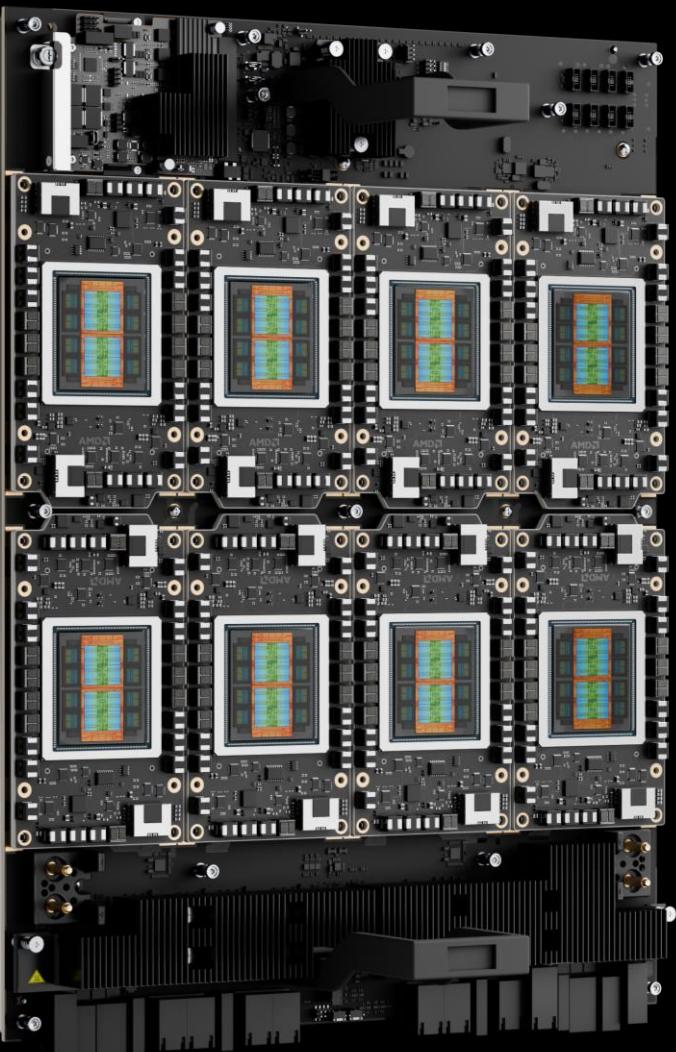
FP6 **20 PF**

FP4 **20 PF**

TBP **1400W**

Instinct™ MI350 Series Advantage

	vs. GB200	vs. B200
MEMORY	1.6x	1.6x
MEMORY BANDWIDTH	1.0x	1.0x
FP64	2.0x	2.1x
FP16	1.0x	1.1x
FP8	1.0x	1.1x
FP6	2.0x	2.2x
FP4	1.0x	1.1x



AMD Instinct™
MI350 Series Platform

Instinct™ MI355X • 8x

MEMORY **2.3 TB HBM3E**

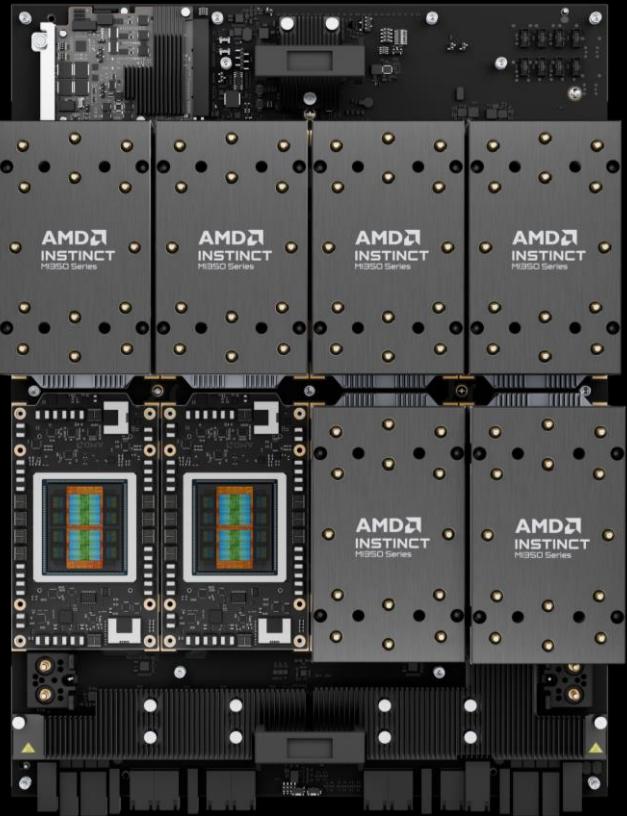
MEMORY BANDWIDTH **64 TB/s**

FP64 **0.63 PF**

FP8 **81 PF**

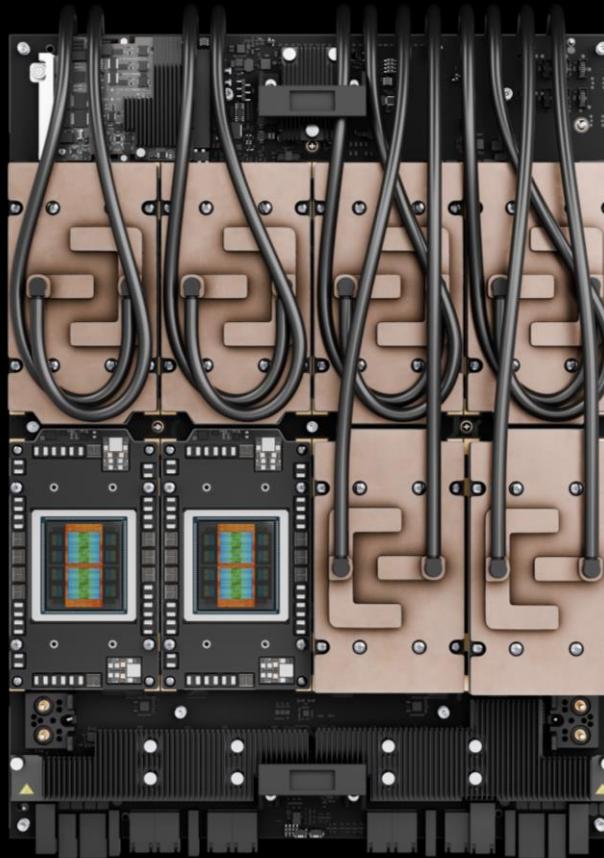
FP6 **161 PF**

FP4 **161 PF**



Air Cooled

AMD Instinct™ MI350 Series



Liquid Cooled

AMD Instinct™ MI350 Series

MI355X Delivers Generational Performance Leap For Ultra Low Latency Inference

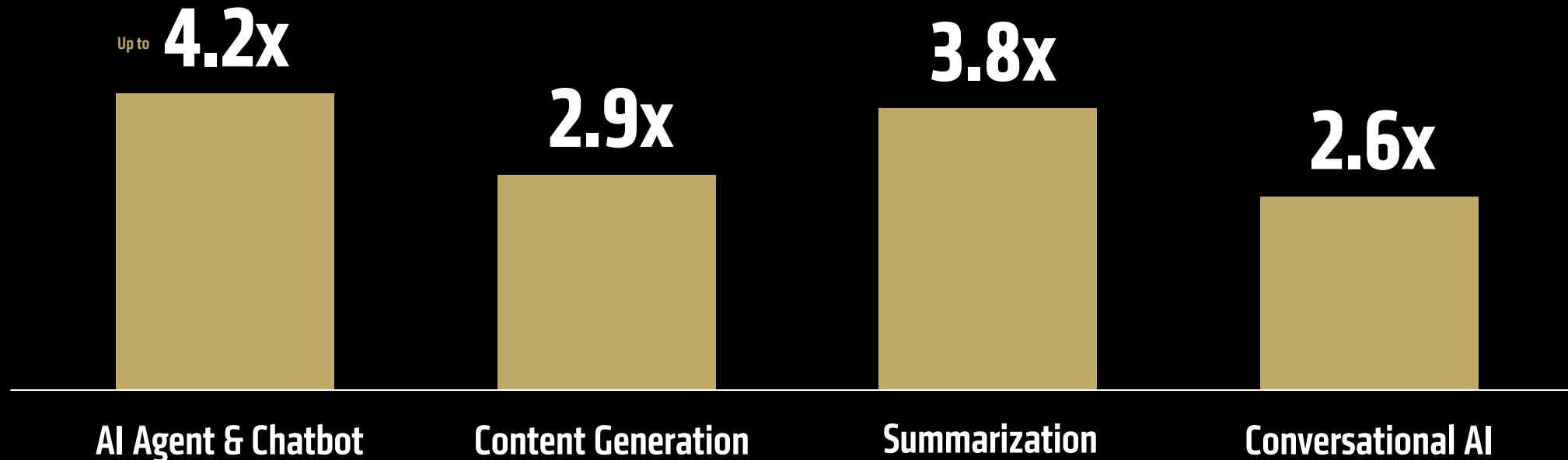


Llama 3.1 405B

Inference performance, throughput • MI355X (FP4) and MI300X (FP8)

Over 3x Generational Inference Improvement

For Broad AI Use Cases



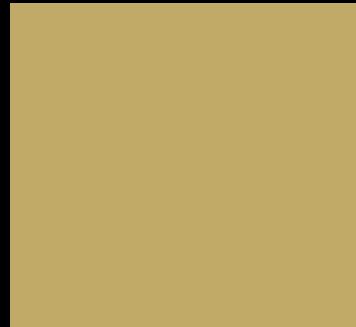
Llama 3.1 405B • MI355X vs. MI300X

Inference performance, throughput • MI355X (FP4) and MI300X (FP8)

Meeting Next-Gen Inference Demands

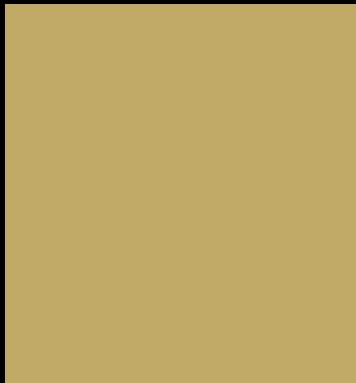
For the Most Popular AI Models

Up to
3.0x



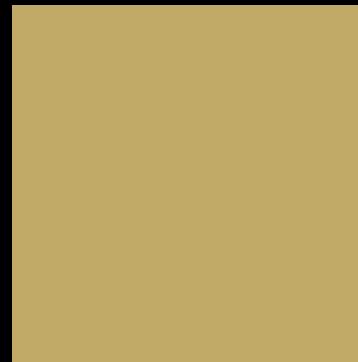
DeepSeek R1

3.3x



Llama 4 Maverick

3.2x

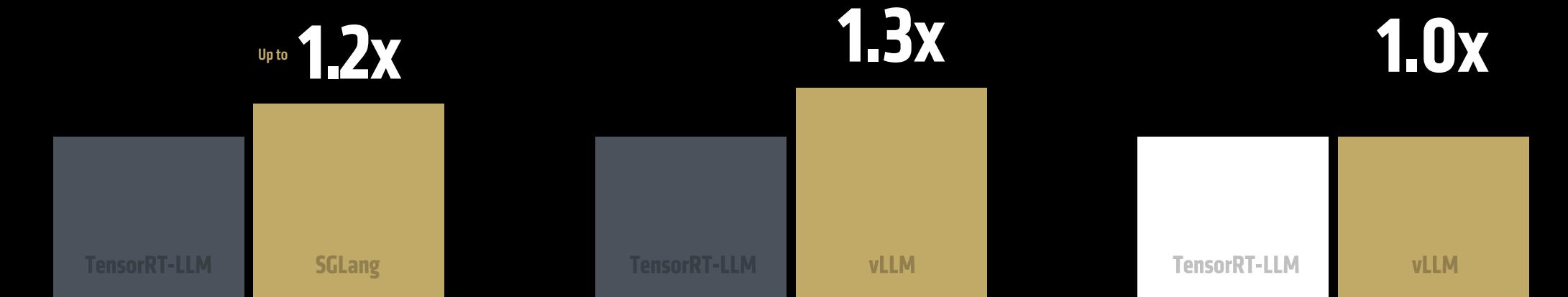


Llama 3.3 70B

MI355X vs. MI300X

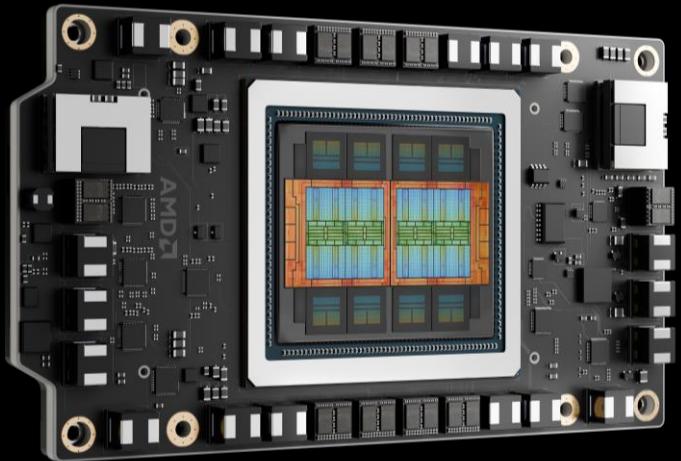
Inference performance, throughput • MI355X (FP4) and MI300X (FP8)

MI355X Delivers the Highest Inference Throughput For Large Models



Inference performance, throughput

See endnote: MI350-038, 039, 040



Up to 40% More Tokens / \$

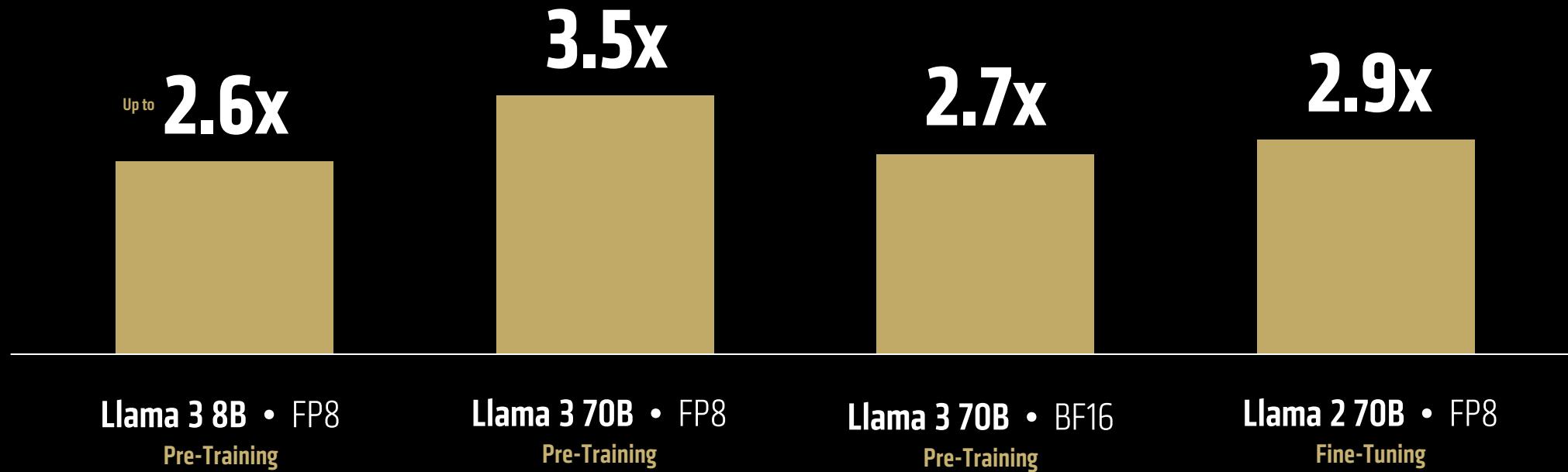
Using AMD Instinct™ MI355X vs. B200

*In addition to Llama Inference, Meta's
AI Recommendation Inference & Training
models run on AMD MI300X GPUs*



Accelerating Model Training by 3x

For Pre-Training & Fine-Tuning



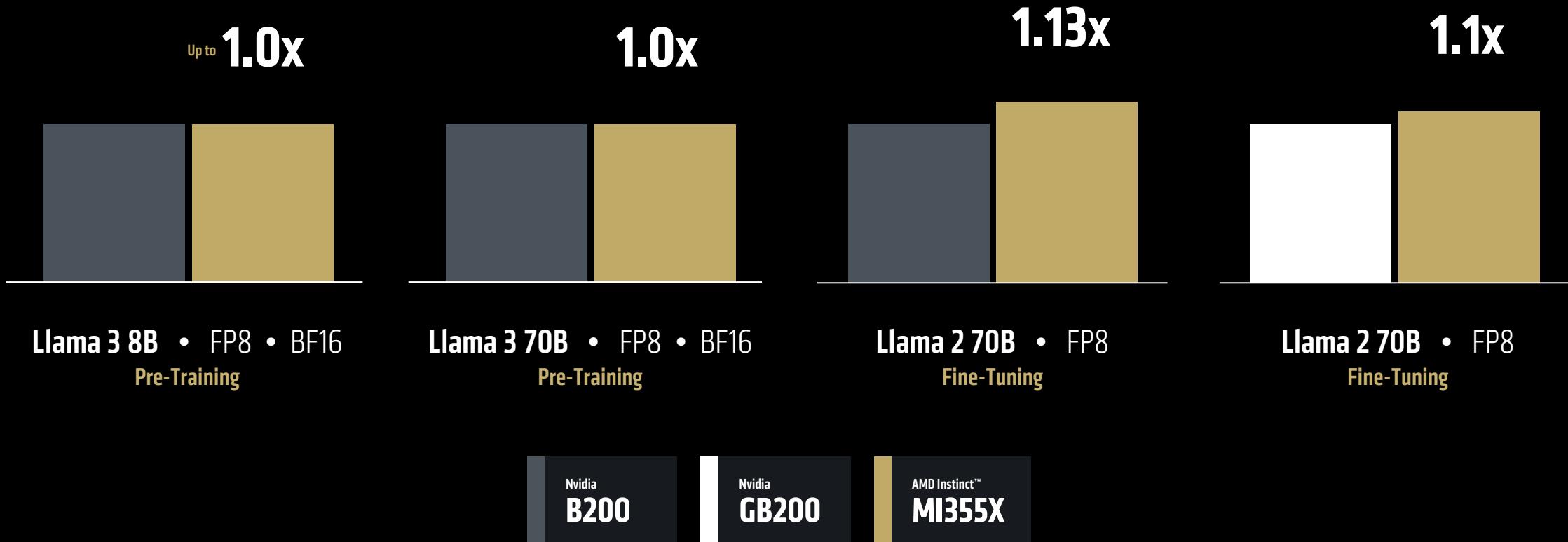
MI355X vs. MI300X

Training: Throughput • Fine Tuning: Time To Train

See endnote: MI350-034, 035

World Class Training & Fine-Tuning Performance

Across Models & Data Types



Pre-Training: Throughput • Fine Tuning: Time To Train • Unofficial MLPerf 5.0 MI355X vs. MLPerf 5.0 B200 and GB200

See endnote: MI350-031, 032, 030, 033

AMD Instinct™ MI350 Series

Proven Liquid Cooled Infrastructure

128/96

Up to **36 TB**

Up to **2.6 EF**

Up to **1.3 EF**

GPUs

HBM3E

FP4

FP8



“Turin”

MI350 Series

Pollara 400

Specs available on 2ou/3ou server



AMD Instinct™ MI350 Series

Proven Air Cooled Infrastructure

64

GPUs

18 TB

HBM3E

1.3 EF

FP4

0.6 EF

FP8

AMD
EPYC

AMD
INSTINCT

AMD
PENSANDO

"Turin"

MI350 Series

Pollara 400

Specs available on 6u server



AMD × ORACLE

 cohere

absci.

 seekr

Modular

 Fireworks AI

Driving AI Innovation

Uber

 vodafone

 PayPal

Trusted by Enterprises

 8x8

Global Cloud Communications

 Palantir

 zoom

 TANIUM

 Xactly

 phenix

 ORACLE
E-BUSINESS SUITE

 ORACLE
Cloud ERP

Powering Enterprise Software

AMD Instinct MI350 Series Solution Partners

ORACLE

DELL
Technologies



Hewlett Packard
Enterprise

CISCO

EVIDEN
an atos business

GIGABYTE™

ASUS®

ASRock
Rack

Cirrascale®

core42

Crusoe

HUMAIN
THE END OF LIMITS.

TENSORWAVE

VULTR

HOT AISLE

EVERGRID

Scaleway

Celestica

COMPAL

ingrasys®

Inventec

MITAC



PEGATRON

wistron®

QCT

wiwynn®

Available starting Q3

Harnessing AI Across Nations

End-to-End Infrastructure

Diverse Ecosystems

Open Architectures

Advancing National Economies

Lawrence Livermore
National Laboratory

OAK RIDGE
National Laboratory

USA

NATIONAL
ENERGY
TECHNOLOGY
LABORATORY

SciNet

OPEN
EURO
LLM

LUMI

USA

Canada

EU

cscs

cea

CiNES
GENCI

MAX PLANCK
GESELLSCHAFT

HLR8
High Performance Computing Center | Stuttgart

Switzerland

France

France

Germany

Germany

eni

MINISTERSTVO VNITRA
CESKE REPUBLIKY

KTH
VETENSKAP
OCH KONST

BBG
BUNDES
BESCHAFFUNG

SG
NSCC
National
Supercomputing
Centre
SINGAPORE

Italy

Czech Republic

Sweden

Austria

Singapore

G42

METI
Agency for Natural Resources and Energy

CENIA
CENTRO NACIONAL DE INTELIGENCIA ARTIFICIAL
NLHPC
National Laboratory
for High Performance
Computing
Chile

Pawsey

HUMAIN
THE END OF LIMITS.

UAE

Japan

Chile

Australia

Saudi Arabia

AMD Silo AI: Europe's AI Solution Factory

Working with European AI Stakeholders Across Public & Private Sectors



OpenEuroLLM



Deploy AI



EuroHPC
Joint Undertaking



LumiOPEN



PORO



Cyber
Valley



appliedAI
institute
for europe



MINISTRAL
AI_



ALEPH ALPHA



WASP

Partnerships

AMD Instinct™ MI350 Series

Generative AI Leadership

Today at **Advancing AI**



GPUs

Software

Solutions



Enabling Open Innovation at Scale

Relentless Progress, Focused on Developers

Accelerated inference
capabilities

Expanded support
for training

Richer out-of-the-box
experience

Developer-first approach
to enablement

Accelerated release
cadence

Day-0 support for
leading models

Deepening ecosystem
partnerships

Industry benchmarks



Deepening Ecosystem Collaboration



Pytorch

Day 0 support
daily performance CI



Triton
v3.3

Performance focus



Hugging Face
1.8 million models

Nightly CI/CD,
finetuning support



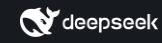
SGL



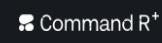
Serving leadership
Distributed
inference



Gemma 3



QwQ-32B



Command R⁺



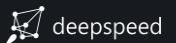
Grok



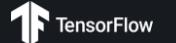
MISTRAL
AI



ONNX



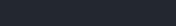
deepspeed



TensorFlow

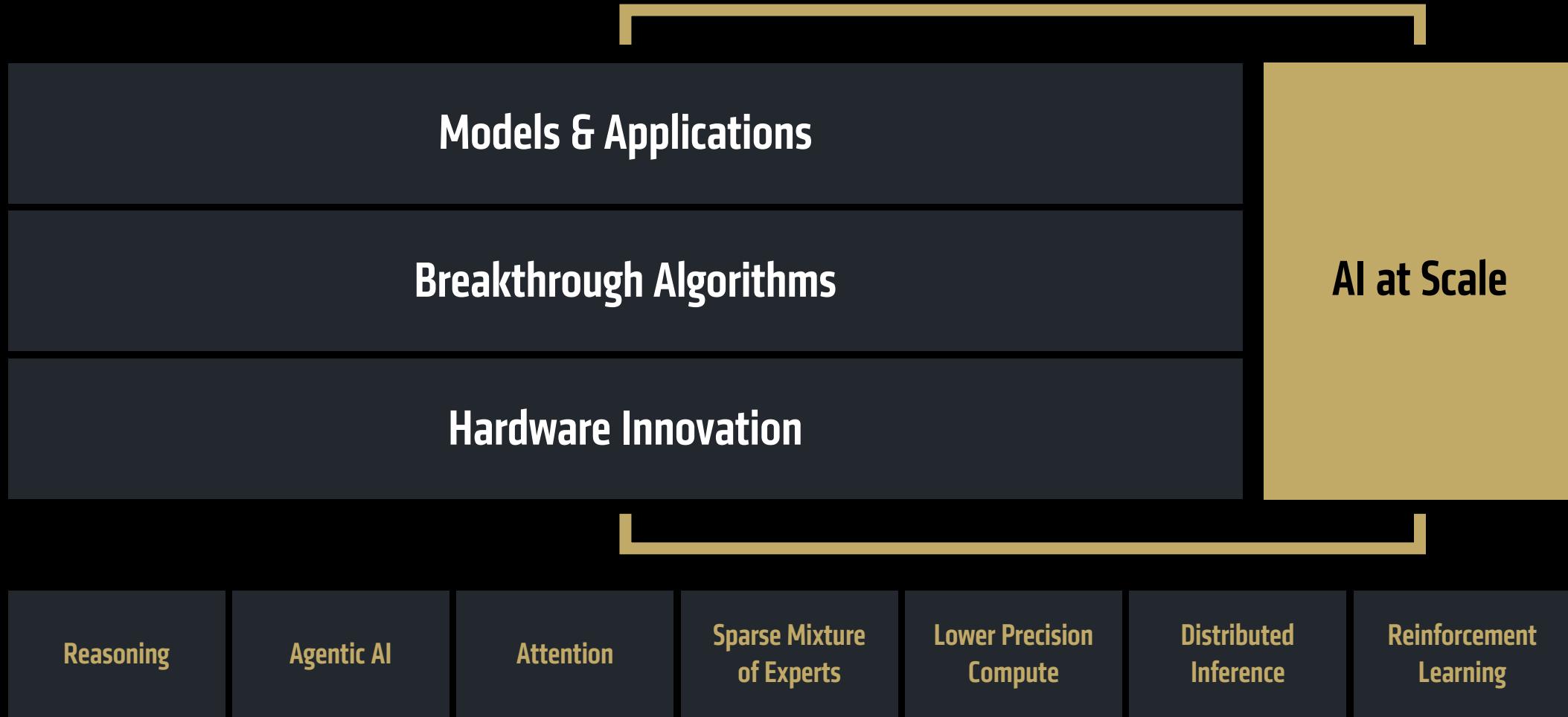


OpenXLA



Expanding open-
source footprint

AI Software Innovation Continues at Rapid Pace



Introducing AMD ROCm™ 7

Accelerating AI Innovation & Developer Productivity

Latest Algorithms
& Models

Advanced Features
for Scaling AI

MI350 Series
Support

Cluster
Management

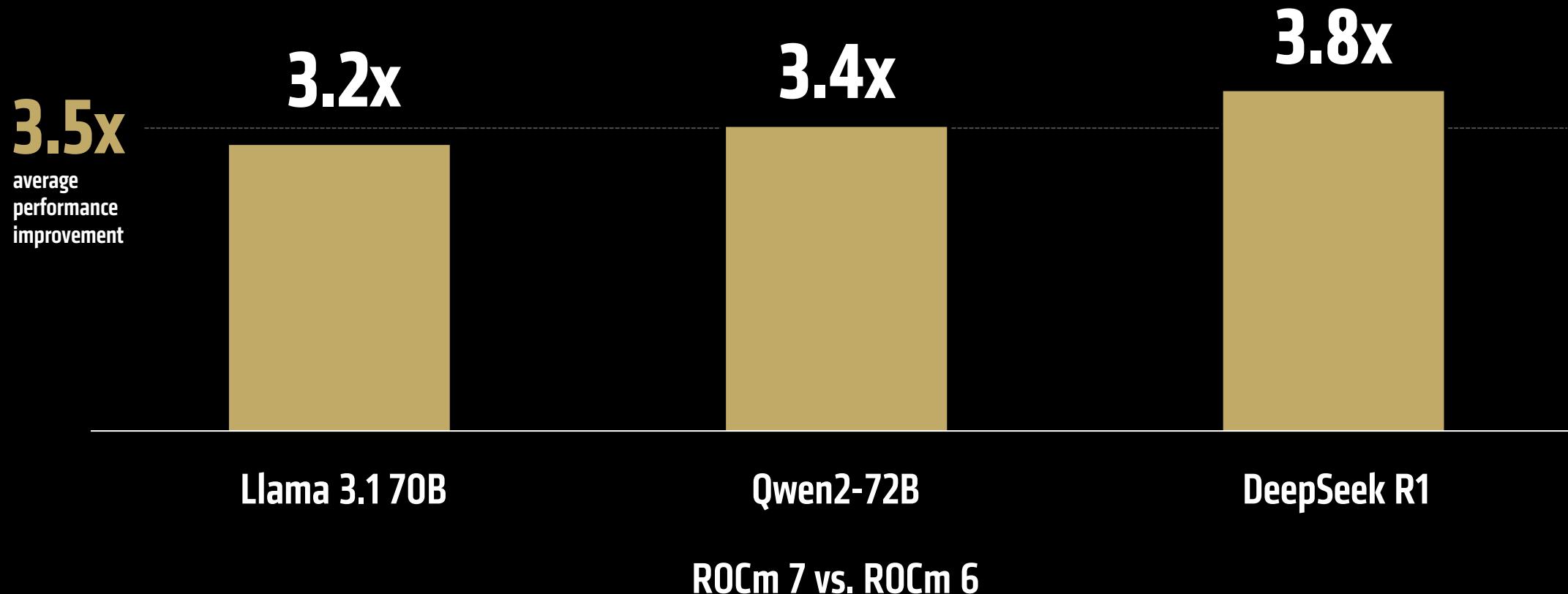
Enterprise
Capabilities

Growing Inference Capabilities

New AMD ROCm Features

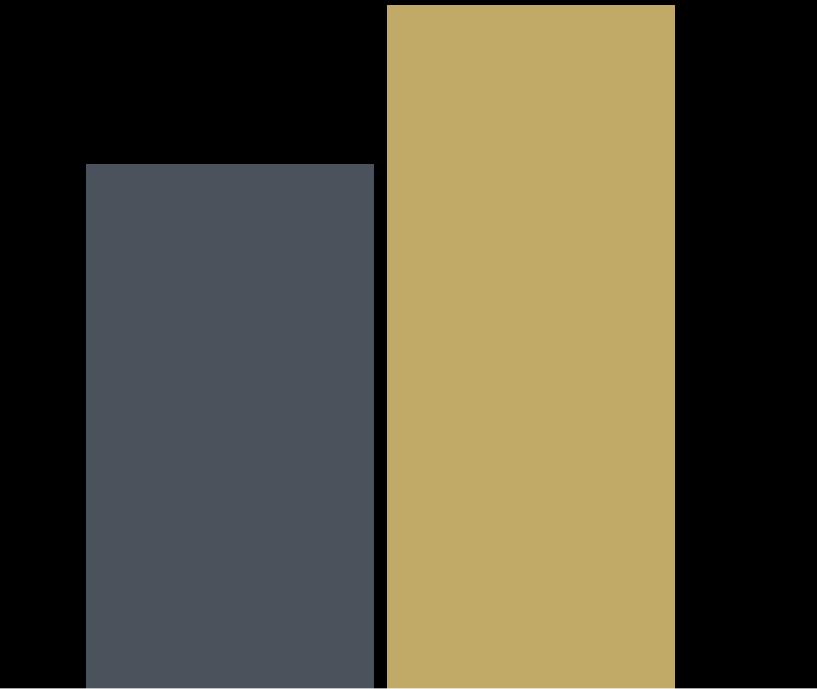
Enhanced Frameworks	vLLM v1	llm-d	SG Lang	
Serving Optimization	Distributed Inference	Prefill	Disaggregation	
Kernels & Algorithms	GEMM Autotuning	MoE	Attention	Python-Based Kernel Authoring
Communication	rocSHMEM		GPU Direct Access	RCCL
Advanced Data Types	FP8	FP6	FP4	Mixed

Accelerating Inference Performance



Open Source: Feature Velocity & Leadership Performance

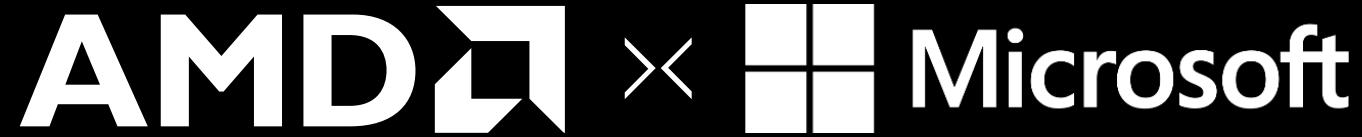
Up to **1.3x**



DeepSeek R1 FP8 Throughput

FP8 Model Support			
vLLM	✓	SGL	✓

Nvidia B200	AMD Instinct MI355X
-----------------------	-------------------------------



M365 Copilot App

AI Assistants



Azure OpenAI Service



Agentic Solutions



Azure AI Foundry

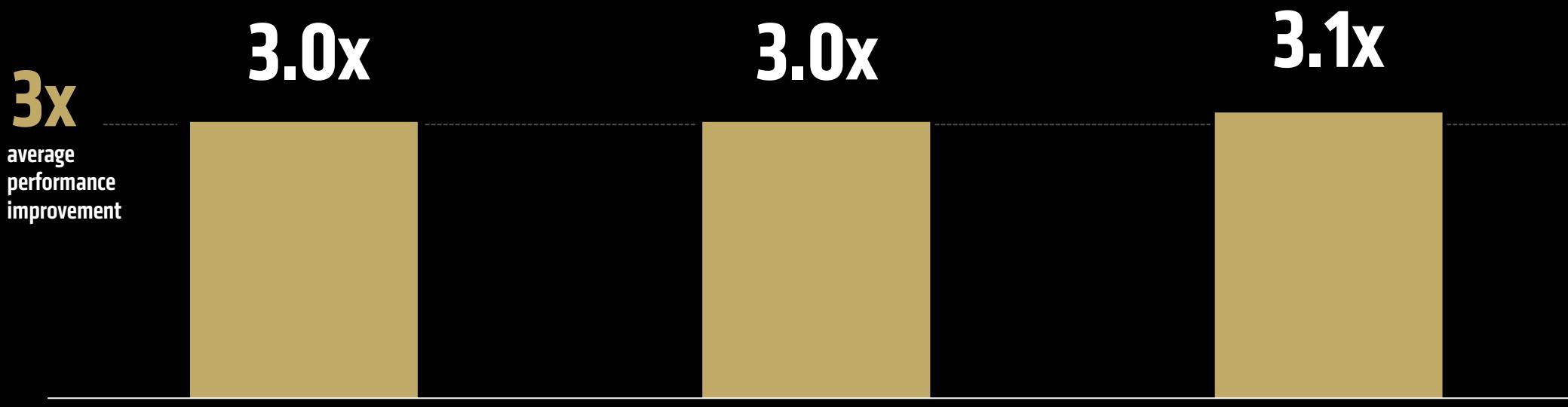
Multi-Modal Training

Growing Training Capabilities

New AMD ROCm Features

AMD Open-Source Models	Text-to-Text	Text-to-Image	European Models	Multimodal	Game Agent
Enhanced Frameworks	PyTorch	JAX Maxtext	Torchtune	Torch-titan	
Parallelization	DP	PP	TP	FSDP	CP
Kernels and Algorithms	GEMM		Attention		
Advanced Data Types	BF16		FP8		

Accelerating Training Performance



Llama 2 70B

Llama 3.1 8B

Qwen 1.5 7B

ROCM 7 vs. ROCM 6

Distributed Inference at Scale with Open Ecosystem

AI Serving Throughput, Multiplied

Orchestration Framework



Key Functions

PD KVcache Transfer

Cross-node Communication

Cross-PD Group Schedule

Key Technologies

Mooncake

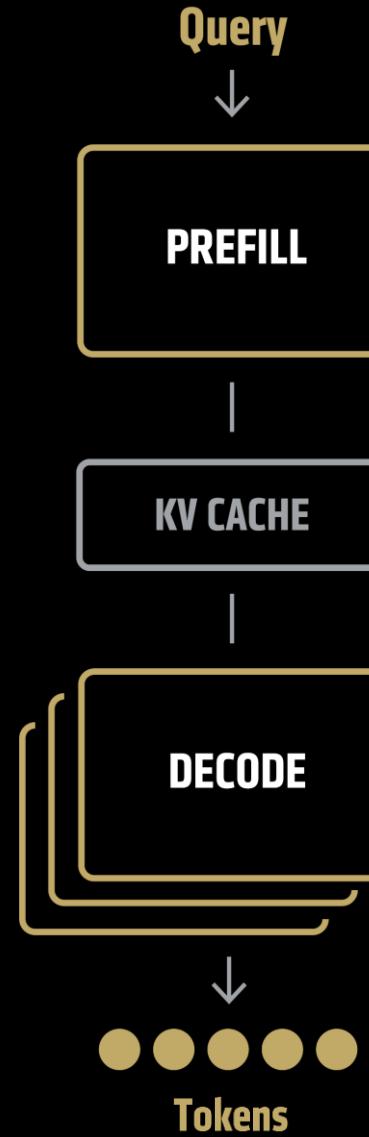
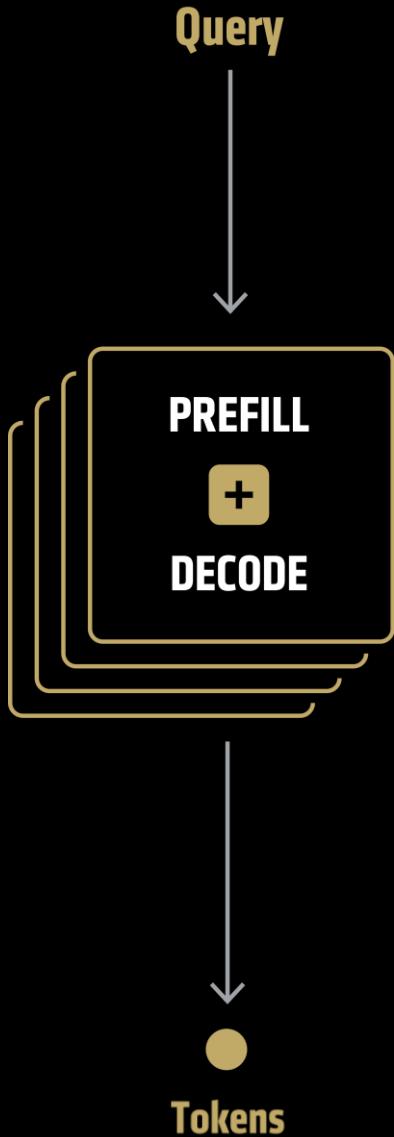
GPU Direct Access

DeepEP

Distributed Triton

SHMEM

Lowering the Cost of Token Generation



Extending AMD ROCm for Enterprise AI

Enterprise Ready

End-to-End Solutions

Secure Data Integration

Ease of Deployment



CLEAR | ML



MLOps

AI Workload & Quota Management

Kubernetes & Slurm Integration



AMD ROCm Enterprise AI

Operations Platform | Cluster Management

Cluster Provisioning & Telemetry

Compiler

Libraries

Profiler

Runtime

AMD ROCm 7

GPUs

CPUs

DPUs

Data Center Infrastructure



Empowering Developers

Ease of Use Collateral for Rapid Adoption

Tutorials

Blogs

Videos

Engaging Developers

Developer-Focused Events to Strengthen the Community

Hackathons

Hands-on Workshops

Developer Contests

Meet-Ups

Community CI





Announcing AMD Developer Cloud & Developer Credits

Available for Developers & Open-Source Contributors

Expanding AMD ROCm on Client

AI-Assisted Coding

Customization

Automation

Advanced Reasoning

Model Fine-Tuning

In-Box Linux Support

2H 2025

Red Hat EPEL



2H 2025

Ubuntu



NEW

OpenSUSE



Fedora



Full Windows Support

NEW

PyTorch

Preview Q3 2025



NEW

ONNX-EP

Preview July 2025



HIP SDK

Linux in Windows WSL

The Ultimate Client AI Solutions for Every Need



AMD Ryzen™ AI 300

Up to **24B** parameters



AMD Ryzen™ AI Max

Up to **70B** parameters



**AMD Threadripper™
+ Radeon™ AI**

Up to **128B** parameters

Developer Track: Meet the AI Experts



Andrew NG

Co-Founder, Coursera,
Landing AI, DeepLearning AI



Chris Lattner

CEO, Modular Inc.



Ashish Vaswani

CEO, Essential AI



Robert Shaw

Director of Engineering, Red Hat AI



Mattias Reso

AI Partner Engineer, Meta



Daniel Han

CEO, Unslot AI



Kunlun Zhu

OpenManus Project, Open Manus



Mark Saroufim

Co-Founder of GPU Mode, Meta



Joe Chau

VP of Engineering, Microsoft Azure



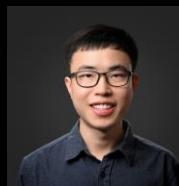
Simon Mo

Co-Lead of vLLM, vLLM



Peng Cheng

Sr. Principal Researcher, Microsoft



Lianmen Zheng

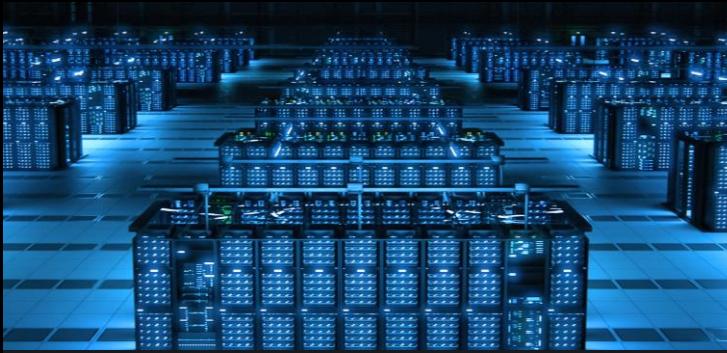
Member of Technical Staff, xAI



Accelerating Innovation



Open by Design

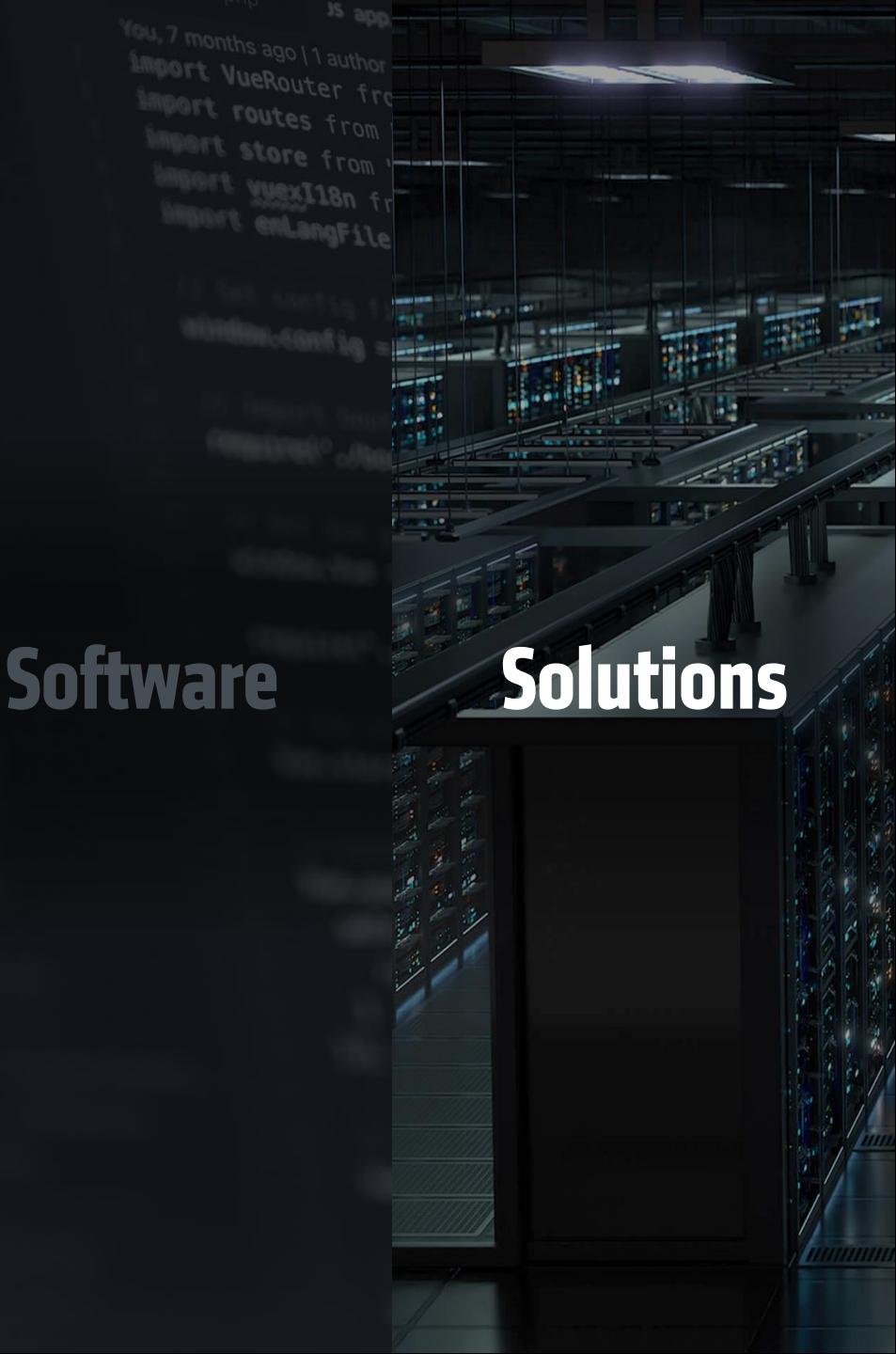
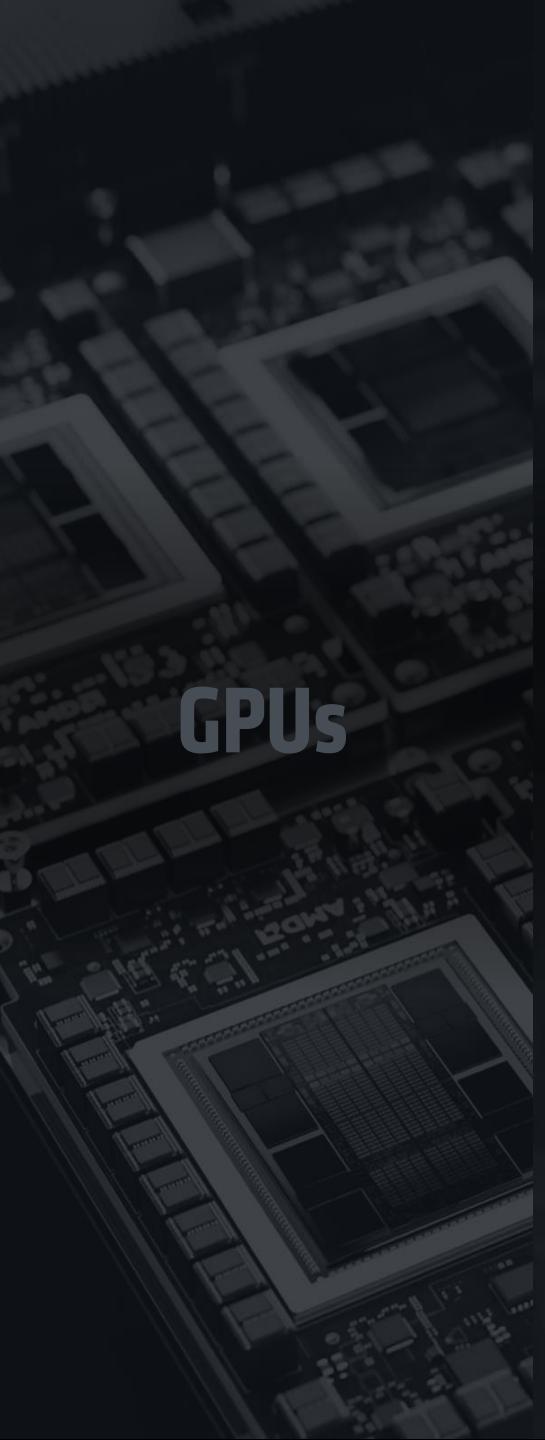


Trusted at Scale



Empowering Developers

Today at **Advancing AI**



Solutions

Transforming Business & Society

Optimize Business Processes

Legal Contracts



Inventory Management

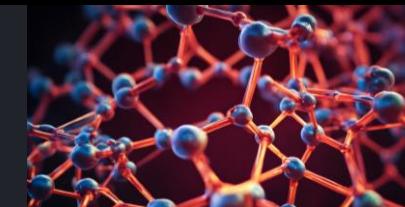


Product Pricing



Accelerate Innovation

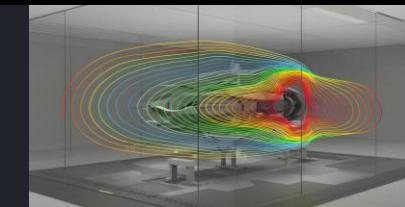
Material Science

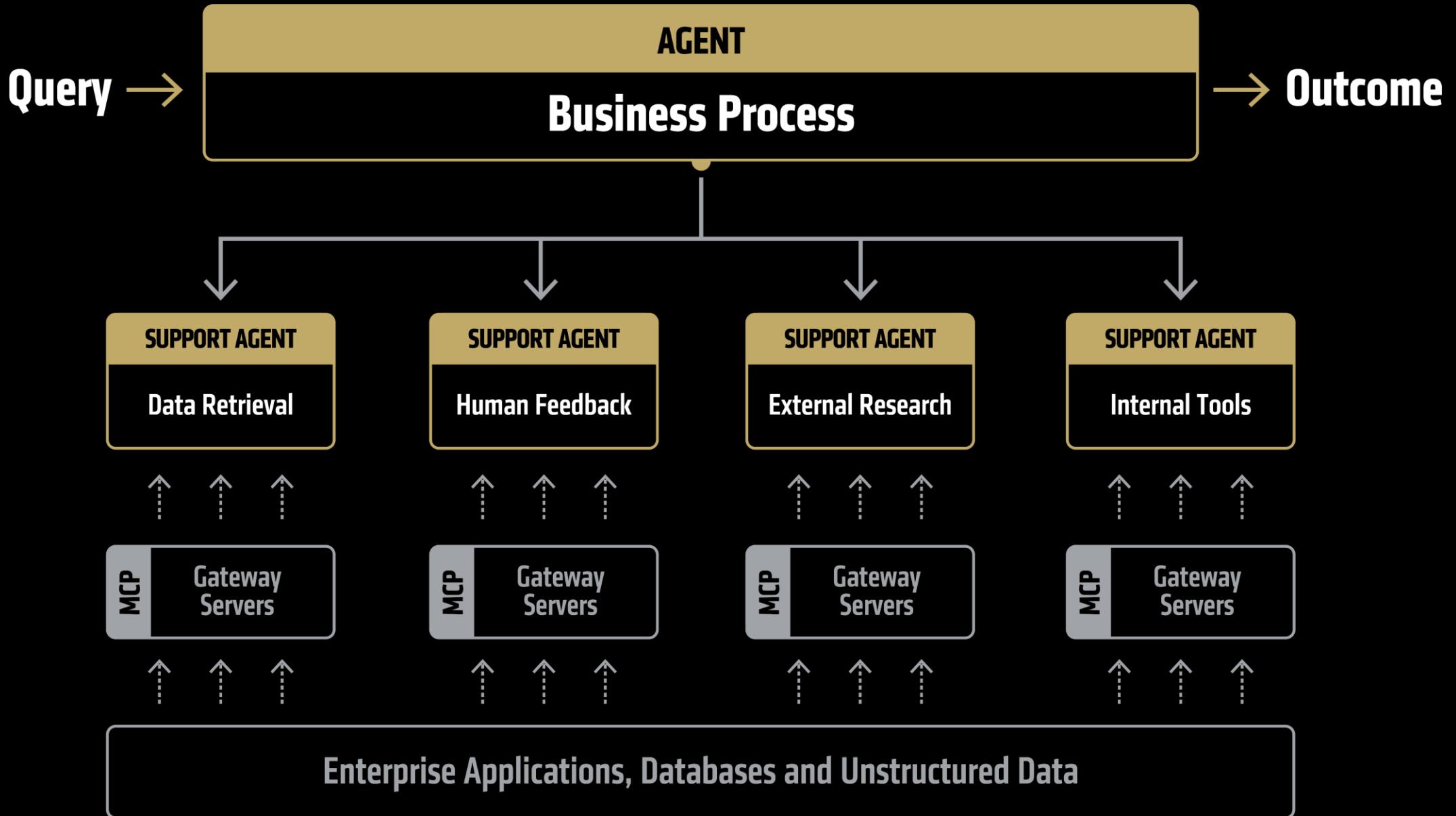


Drug Discovery

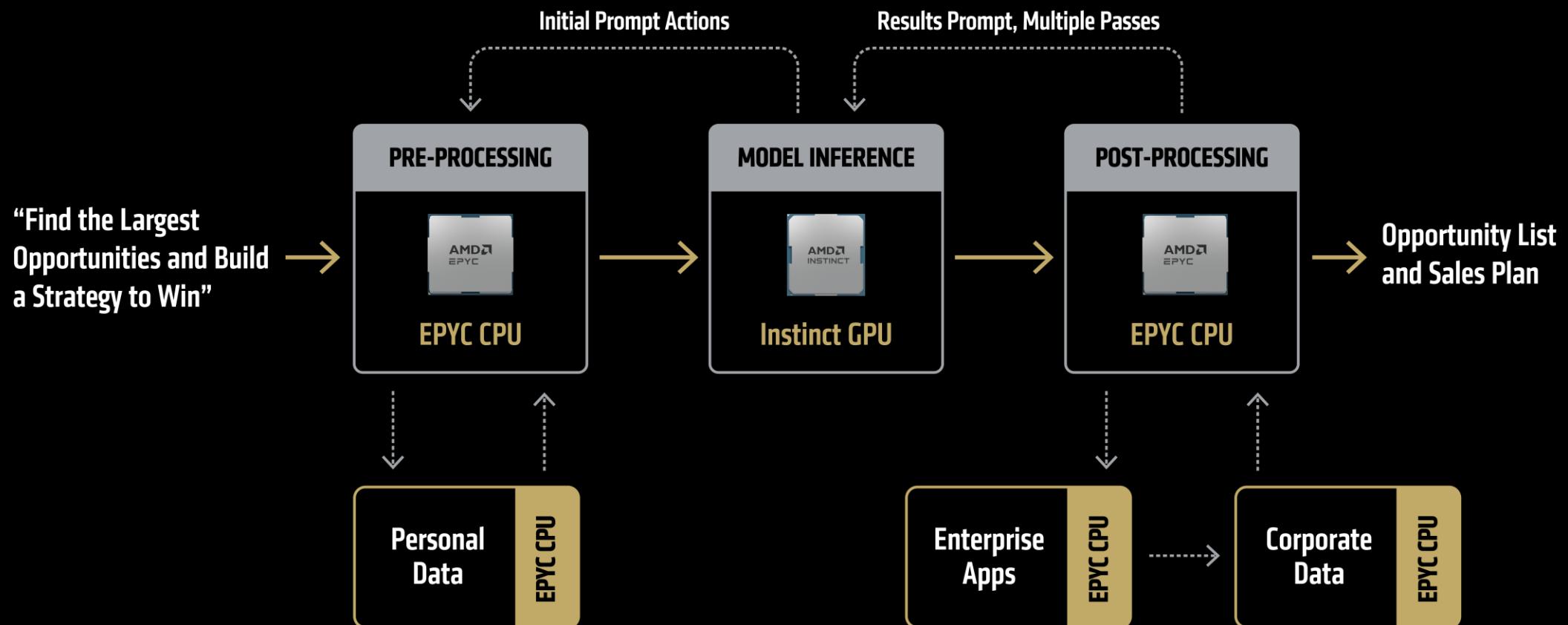


Mechanical Design

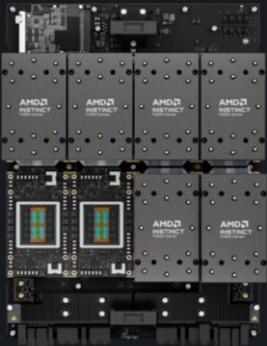




Agentic Execution



End-to-End Integrated AI Platform

Front End Network	CPU Node	GPU Node	Scale Up Network	Scale Out Network
 The AMD Pensando logo, featuring a stylized blue square with the text "AMD PENSANDO" in white.	 The AMD EPYC logo, featuring a silver rectangular shape with the text "AMD EPYC" in black.	 A photograph of an AMD Instinct GPU Node board, showing multiple green and grey circuit boards stacked vertically.	 The Ultra Accelerator Link logo, consisting of three vertical arrows pointing upwards, followed by the text "ULTRA ACCELERATOR LINK™".	 The AMD Pensando logo, featuring a stylized blue square with the text "AMD PENSANDO" in white. Below it, the text "Ultra Ethernet Consortium" is written in a smaller font.

Securely Integrate
into Enterprise

X86 Applications
& AI Execution

AI Model Training
& Inference

Connectivity for Training
& Distributed Inference

Enable Gigawatt
Level Scaling

Front-End Networks Feed AI

**Securely bridge AI servers
to enterprise data & apps**

**Boost AI server performance with
network & security offloads**

**DPU accelerates network functions
up to 40x versus CPU only**

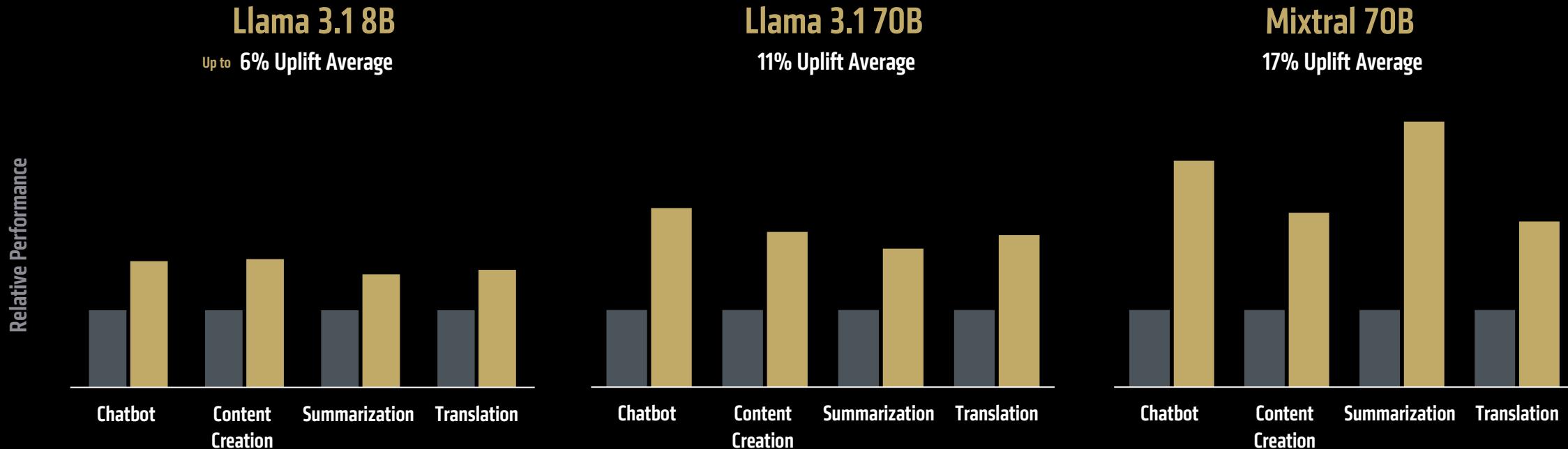
High Performance CPUs Drive Agentic AI

Leverages x86 ecosystem to power agentic use-cases

5Ghz operation boosts GPU operating efficiency

Robust Enterprise class reliability, availability & serviceability

AMD EPYC™ Driving End-to-End System Performance



Optimized for Performance

High efficiency protocol for
low power & low latency <1uS

Up to 800 Gbps per port with
multiple ports per accelerator

Leverages 200 G ethernet
physical layer & infrastructure



Protocol Interface

Transaction Layer

Data Link Layer

Industry Standard Ethernet PHY

Cabling

Connectors

Re-Timers

Management
Software

Engineered to Scale

Scales from SMB to
Gigawatt Data Center

Supports virtual POD
partitioning

Robust & reliable
with fault resiliency

Accelerating AI models with
in-network collectives



ULTRA ACCELERATOR LINK™



Ultra Accelerator Link, the Truly Open Standard

	UALINK	NVLINK FUSION
LOW ROUND TRIP LATENCY	Yes	Yes
HIGH SPEED I/O	224 Gbps	224 Gbps
MAXIMUM SCALABILITY	1024 GPUs	576 GPUs
CPU USE	Any	Only to Nvidia GPU
GPU USE	Any	Only Nvidia GPU
MANAGEMENT SOFTWARE	Open	Nvidia Proprietary
SPECIFICATION	Fully Open	Closed



ULTRA
ACCELERATOR
LINK™

Any CPU

Any Accelerator

Any Switch



Open Standard Drives Innovation & Choice

100+ Consortium Members



Google

Meta

AMD

Microsoft

synopsys[®]

intel.

AsteraLabs



CISCO

Alibaba.com

Hewlett Packard
Enterprise

Promoters

MARVELL[®]

auradine

ALPHAWAVE SEMI

DELL Technologies

XCONNTECH

enfabrica

cādence[®]

BROADCOM[®]

Qualcomm

Lenovo

XCONNTECH

MEDIATEK

arm

CREDO
we connect.

Select Contributors



Switch Fabric

Signal Conditioning

Controllers

More...

Chiplets

ICs

Modules

Boards

Software

Copper

Optical

Purpose-Built Connectivity for AI Infrastructure





MARVELL™

Optics

Connectivity

Switching

CPO

Scale-Up Fabric

CPC

Custom Silicon

Rack-Scale AI



Open Architecture for Performant AI Scale-Out

Modern RDMA transport with
automatic configuration & tuning

Optimized protocol for high fabric
utilization & congestion control

Unparallel scalability to
over a million GPUs



ARISTA

BROADCOM®



intel.



Hewlett Packard
Enterprise

EVIDEN

Meta

ORACLE

Steering Committee



Google

Lenovo



Qualcomm

SYNOPSYS®

cadence®

NOKIA



SAMSUNG

Select General Members

100+ Consortium Members

Scaling Out to Gigawatt Data Centers

Programmability enables
collective acceleration

Load balancing reduces
networking cost

Line-rate engines boost
flexibility & performance

Improved reliability
& resiliency

World Class Open Platform

Highest memory bandwidth & capacity drives performance

Industry standard x86 for easy integration to enterprise applications

Ready to deploy rack infrastructure

AMD Instinct™ GPU

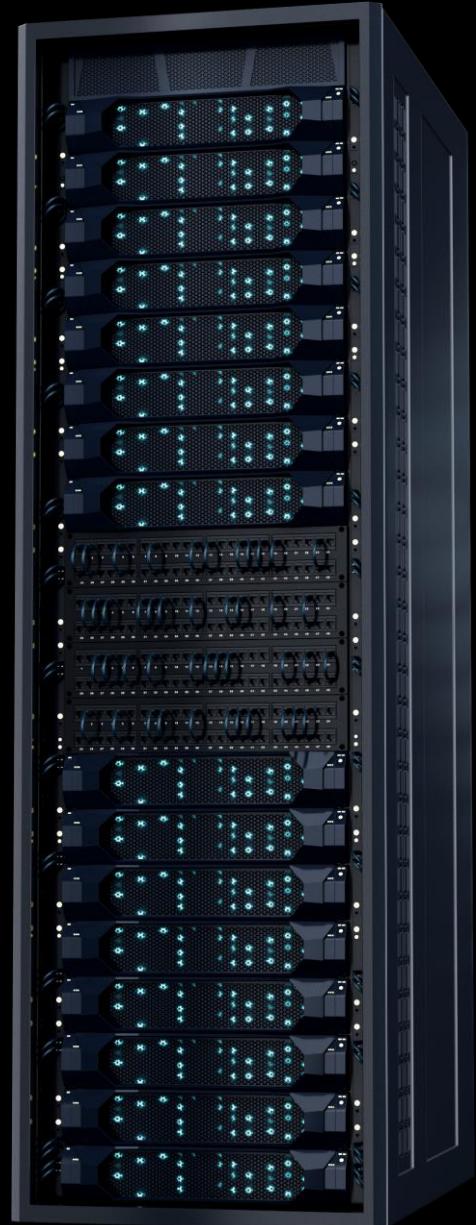
MI350 SERIES

AMD Server CPU

EPYC™ “TURIN”

AMD Pensando™ NIC

POLLARA 400



96 / 128 GPUs

Up to 1.8x more than GB200 NVL72

36 TB HBM3E

Up to 2.8x more than GB200 NVL72

2.6 EF FP4 1.3 EF FP8

Up to 1.6x more than GB200 NVL72

2.6 EF FP6

Up to 3.3x more than GB200 NVL72

x86 CPU

5th Gen EPYC "Turin"

Ultra Ethernet Consortium Scale-Out NIC

AMD Pollara NIC

OCP Compliant
Industry standard design

Ultra Ethernet
Consortium



Advancing AI Infrastructure Solutions

2024

AMD EPYC™
“GENOA”

AMD Instinct™
MI300 SERIES



2025

AMD EPYC™
“TURIN”

AMD Instinct™
MI350 SERIES

AMD Pensando™
POLLARA 400



2026

AMD EPYC™
“VENICE”

AMD Instinct™
MI400 SERIES

AMD Pensando™
“VULCANO”



Previewing Today at **Advancing AI 2025**

The World's Best AI Rack Solution

For At-Scale Training & Distributed Inference

Previewing Today at **Advancing AI 2025**

AMD “Helios” Optimized AI Rack Solution

AMD
EPYC

AMD
INSTINCT

AMD
PENSANDO

AMD
ROCm

Available in 2026

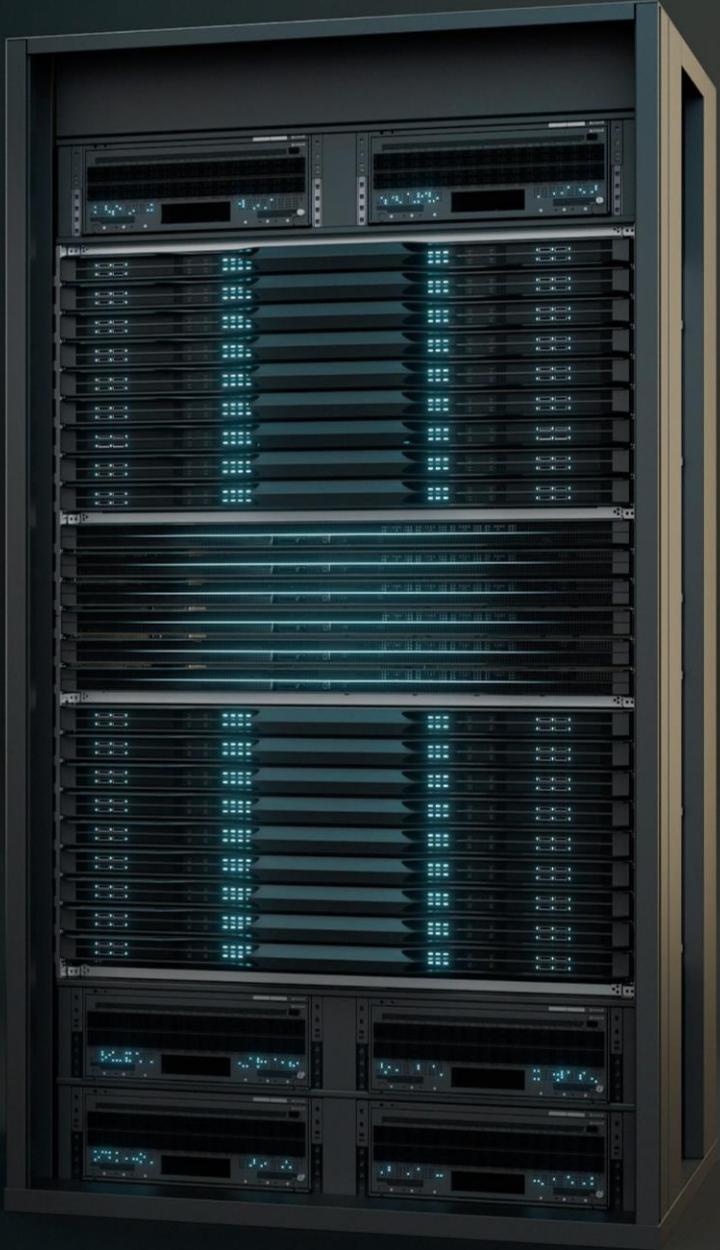


OPEN

Compute Project®



ULTRA
ACCELERATOR
LINK™



AMD “Helios” AI Rack

Rack Scale AI Performance Leadership

AMD Instinct™ MI400 Series vs. Vera Rubin

“Helios”

Oberon

GPU DOMAIN	72	1.0x
SCALE UP BANDWIDTH	260 TB/s	1.0x
FP4 • FP8 FLOPS	2.9 EF • 1.4 EF	1.0x
HBM4 MEMORY CAPACITY	31 TB	1.5x
MEMORY BANDWIDTH	1.4 PB/s	1.5x
SCALE OUT BANDWIDTH	43 TB/s	1.5x

AMD EPYC™ “Venice”

Highest Performance Server CPU

Up to **256 cores**

2nm • Zen 6

2.0x

CPU to GPU Bandwidth

1.7x

Gen vs. Gen Performance

1.6 TB/s

Memory Bandwidth

Coming in 2026

AMD Instinct™ MI400

Leadership Gen AI Accelerator

40 PF | 20 PF

FP4 • FP8 Flops

432 GB

HBM4 Memory Capacity

19.6 TB/s

Memory Bandwidth

300 GB/s

Scale Out Bandwidth / GPU

Coming in 2026

AMD Pensando™ “Vulcano”

Next Gen NIC for AI Clusters

3nm

Process Node

800G

Network Throughput

Up to 8x

Scale Out Bandwidth per GPU

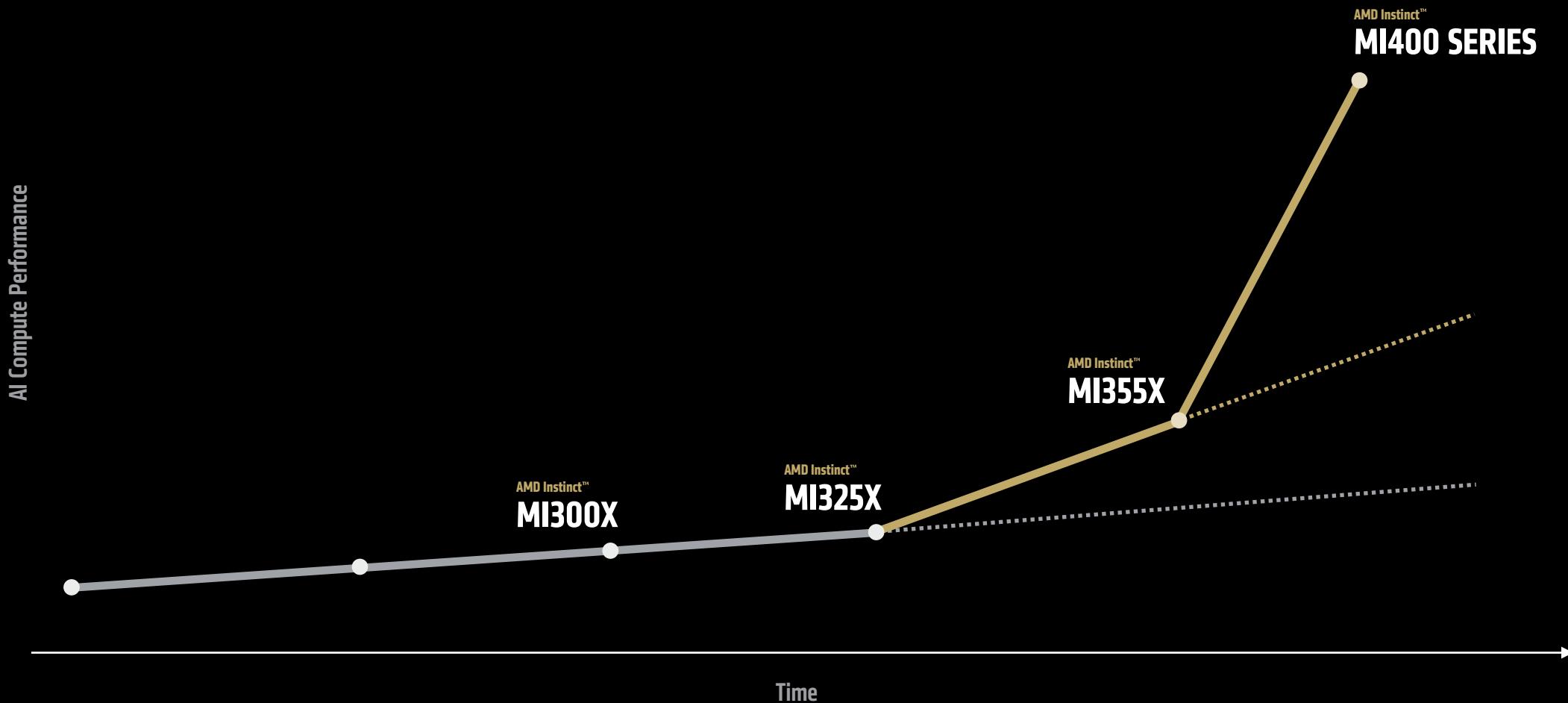
UAL | PCIe®

Host Interface

Ultra Ethernet
Consortium

Coming in 2026

AI Compute Performance



**Up to 10x More Performance
with AMD Instinct™ MI400 Series**

Compared to Instinct™ MI355X



Advancing AI Infrastructure on an Annual Cadence

2025

AMD EPYC
"TURIN"

AMD Instinct
MI350 SERIES

AMD Pensando
POLLARA 400



2026

AMD EPYC
"VENICE"

AMD Instinct
MI400 SERIES

AMD Pensando
"VULCANO"



2027

AMD EPYC
"VERANO"

AMD Instinct
MI500 SERIES

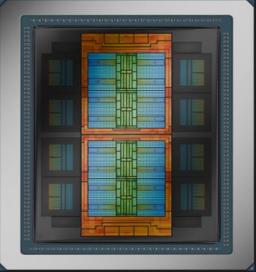
AMD Pensando
"VULCANO"



Next Gen AI Rack

Announced Today at **Advancing AI 2025**

Accelerating Momentum Across Our Broad AI Capabilities



AMD
ROCm

AMD
Developer Cloud



AMD Instinct™
MI350 Series

AMD ROCm 7 Optimized
for Reasoning & Agents

Empowering
Developers

AMD "Helios"
AI Rack

AMD Threadripper™ and
Radeon™ AI solutions

In Production Now
Systems Q3

Available August

Available Now

Coming in 2026

Available July



AI Innovation is a Global, Collective Effort

Endnotes

SHO-06: Testing as of Dec 2024 using the following benchmark scores compared to Intel Core Ultra 9 288V and Qualcomm Snapdragon X Elite X1E-84-100. Cinebench 2024 nT, 3Dmark Wildlife Extreme, and Blender.. Next gen AI PC defined as a Windows PC with a processor that includes a NPU with at least 40 TOPS. Configuration for AMD Ryzen™ AI Max+ 395 processor: AMD reference board, Radeon™ 8060S graphics, 32GB RAM, 1TB SSD, VBS=ON, Windows 11. Configuration for Qualcomm Snapdragon X Elite X1E-84-100 processor: Samsung Galaxybook, Adreno Graphics, 16GB RAM, Microsoft Windows 11. Configuration for Intel Core Ultra 9 288V: ASUS Zenbook X 14, Intel Arc Graphics, 32GB RAM, 1TBSSD, Microsoft Windows 11 Home. Laptop manufacturers may vary configurations yielding different results.

MI300-080: Testing by AMD Performance Labs as of May 15, 2025, measuring the inference performance in tokens per second (TPS) of AMD ROCm 6.x software, vLLM 0.3.3 vs. AMD ROCm 7.0 preview version SW, vLLM 0.8.5 on a system with (8) AMD Instinct MI300X GPUs running Llama 3.1-70B (TP2), Qwen 72B (TP2), and Deepseek-R1 (FP16) models with batch sizes of 1-256 and sequence lengths of 128-204. Stated performance uplift is expressed as the average TPS over the (3) LLMs tested. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of the latest drivers and optimizations.

MI300-081: AMD Instinct MI300X platform (8x GPUs) and AMD ROCm 7.0 preview version software running Llama2-70B, Qwen1.5-14B, Llama3.1-8B, Megatron-LM using the FP16 and FP8 datatypes, shows a combined average of 3.04x or average of 304% better training performance (TFLOPS) vs. AMD Instinct MI300X platform (8x GPUs) with ROCm 6.0 SW.

MI350-004: Based on calculations by AMD Performance Labs in May 2025, to determine the peak theoretical precision performance of eight (8) AMD Instinct™ MI355X and MI350X GPUs (Platform) and eight (8) AMD Instinct MI325X, MI300X, MI250X and MI100 GPUs (Platform) using the FP16, FP8, FP6 and FP4 datatypes with Matrix. Server manufacturers may vary configurations, yielding different results. Results may vary based on use of the latest drivers and optimizations.

MI350-008: Based on measurements taken by AMD Performance Labs in May 2025, of the peak theoretical precision performance of an AMD Instinct™ MI355X GPU with FP64 datatype with Matrix vs. Nvidia Grace Blackwell GB200 accelerator with FP64 datatype with Tensor; MI355X: FP32 with Matrix vs. GB200: FP32 datatype with Vector; and MI355X: FP6 datatype with Sparsity vs. GB200: FP6 datatype with Sparsity. Results may vary based on configuration, datatype. **MI350-008**

MI350-009: Based on calculations by AMD Performance Labs in May 2025, to determine the peak theoretical precision performance for the AMD Instinct™ MI350X / MI355X GPUs, when comparing FP64, FP32, TF32, FP16, FP8, FP6 and FP4, INT8, and bfloat16 datatypes with Vector, Matrix, Sparsity or Tensor with Sparsity as applicable, vs. NVIDIA Blackwell B200 accelerator. Server manufacturers may vary configurations, yielding different results.

MI350-025: Testing by AMD Performance Labs as of May 25, 2025, measuring the inference performance in tokens per second (TPS) of the AMD Instinct MI355X platform with ROCm 7.0 pre-release build 16047, running DeepSeek R1 LLM on SGLang versus NVIDIA Blackwell B200 platform with CUDA version 12.8. Server manufacturers may vary configurations, yielding different results. Performance may vary based on hardware configuration, software version, and the use of the latest drivers and optimizations.

MI350-030: Based on calculations by AMD internal testing as of 6/4/2025. Using 8 GPU AMD Instinct™ MI355X Platform for overall GPU-normalized Training Throughput (processed tokens per second) for text generation using the Llama3-70B chat model running TorchTITAN (FP8) when using a maximum sequence length of 8192 tokens compared to published 64 GPU Nvidia B200 Platform performance running NeMo (FP8) when using a maximum sequence length of 8192 tokens. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

Endnotes

MI350-031: Based on calculations by AMD internal testing as of 6/4/2025. Using 8 GPU AMD Instinct™ MI355X Platform for overall GPU-normalized Training Throughput (processed tokens per second) for text generation using both LLaMA3-70B and LLaMA3-8B chat models running TorchTITAN (BF16) or Megatron-LM (BF16) where applicable when using a maximum sequence length of 8192 tokens compared to 8 GPU Nvidia B200 Platform performance running NeMo (BF16) when using a maximum sequence length of 8192 tokens. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-032: Based on calculations by AMD internal testing as of 6/4/2025. Using 8 GPU AMD Instinct™ MI355X Platform for overall GPU-normalized Training Throughput (processed tokens per second) for text generation using both LLaMA3-70B and LLaMA3-8B chat models running TorchTITAN (BF16) or Megatron-LM (BF16) where applicable when using a maximum sequence length of 8192 tokens compared to 8 GPU Nvidia B200 Platform performance running NeMo (BF16) when using a maximum sequence length of 8192 tokens. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-033: Based on calculations by AMD internal testing as of 6/5/2025. Using 8 GPU AMD Instinct™ MI355X Platform for overall GPU-normalized Training Throughput (time to complete) for fine-tuning using the Llama2-70B LoRA chat model (FP8) compared to published 8 GPU Nvidia B200 and 8 GPU Nvidia GB200 Platform performance (FP8). Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-034: Based on AMD internal testing as of 6/4/2025, using an (8) GPU AMD Instinct™ MI355X Platform for overall GPU-normalized Training Throughput (processed tokens per second) for text generation using the LLaMA3-70B and LLaMA3-8B chat models running TorchTITAN or Megatron-LM (FP8 and BF16) as applicable, using a maximum sequence length of 8192 tokens, compared to an (8) GPU AMD Instinct™ MI300X Platform using Megatron-LM (FP8 and BF16). Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations

MI350-035: Based on AMD internal testing as of 6/5/2025. Using 8 GPU AMD Instinct™ MI355X Platform for overall GPU-normalized Training Throughput (time to complete) for fine-tuning using the Llama2-70B LoRA chat model (FP8) compared 8 GPU AMD Instinct™ MI300X Platform performance with (FP8). Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations

MI350-038: Based on testing by AMD internal labs as of 6/6/2025 measuring text generated throughput for LLaMA 3.1-405B model using FP4 datatype. Test was performed using input length of 128 tokens and an output length of 2048 tokens for AMD Instinct™ MI355X 8xGPU platform compared to NVIDIA B200 HGX 8xGPU platform published results. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-039: Based on Lucid automation framework testing by AMD internal labs as of 6/6/2025 measuring text generated throughput for LLaMA 3.1-405B model using FP4 datatype. Test was performed using 4 different combinations (128/2048) of input/output lengths to achieve a mean score of tokens per second for AMD Instinct™ MI355X 4xGPU platform compared to NVIDIA DGX GB200 4xGPU platform. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations

MI350-040: Based on testing (tokens per second) by AMD internal labs as of 6/6/2025 measuring text generated online serving throughput for DeepSeek-R1 chat model using FP4 datatype. Test was performed using input length of 3200 tokens and an output length of 800 tokens with concurrency up to 64 looks, serviceable with 30ms ITL threshold for AMD Instinct™ MI355X 8xGPU platform median total tokens compared to NVIDIA B200 HGX 8xGPU platform results. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-041: Based on AMD internal testing as of 6/5/2025. Using 8 GPU AMD Instinct™ MI355X Platform measuring text generated offline inference throughput for Llama4 Maverick chat model (FP4) compared 8 GPU AMD Instinct™ MI300X Platform performance with (FP8). MI355X ran 8xTP1 (8 copies of model on 1 GPU) compared to MI300X running 2xTP4 (2 copies of model on 4 GPUs). Tests were conducted using a synthetic dataset with different combinations of 128 and 2048 input tokens, and 128 and 2048 output tokens. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

Endnotes

MI350-042: Based on AMD internal testing as of 6/5/2025. Using 8 GPU AMD Instinct™ MI355X Platform measuring text generated offline inference throughput for Llama 3.1-405B chat model (FP4) compared 8 GPU AMD Instinct™ MI300X Platform performance with (FP8). MI355X ran 8xTP1 (8 copies of model on 1 GPU) compared to MI300X running 2xTP4 (2 copies of model on 4 GPUs). Tests were conducted using a synthetic dataset with different combinations of 128 and 2048 input tokens, and 128 and 2048 output tokens. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-043: Based on AMD internal testing as of 6/5/2025. Using 8 GPU AMD Instinct™ MI355X Platform measuring text generated online serving inference throughput for DeepSeek-R1 chat model (FP4) compared 8 GPU AMD Instinct™ MI300X Platform performance with (FP8). Test was performed using input length of 3200 tokens and an output length of 800 tokens with concurrency set to maximize the throughput on each platform, 128 for MI300X and 2048 for MI355X platforms. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations

MI350-044: Based on AMD internal testing as of 6/9/2025. Using 8 GPU AMD Instinct™ MI355X Platform measuring text generated online serving inference throughput for Llama 3.1-405B chat model (FP4) compared 8 GPU AMD Instinct™ MI300X Platform performance with (FP8). Test was performed using input length of 32768 tokens and an output length of 1024 tokens with concurrency set to best available throughput to achieve 60ms on each platform, 1 for MI300X (35.3ms) and 64ms for MI355X platforms (50.6ms). Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations. Based on AMD internal testing as of 6/9/2025. Using 8 GPU AMD Instinct™ MI355X Platform measuring text generated online serving inference throughput for Llama 3.1-405B chat model (FP4) compared 8 GPU AMD Instinct™ MI300X Platform performance with (FP8). Test was performed using input length of 32768 tokens and an output length of 1024 tokens with concurrency set to best available throughput to achieve 60ms on each platform, 1 for MI300X (35.3ms) and 64ms for MI355X platforms (50.6ms). Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations

MI350-047: Based on engineering projections by AMD Performance Labs in June 2025, to estimate the peak theoretical precision performance of seventy-two (72) AMD Instinct™ MI400X GPUs (Rack) vs. an 8xGPU AMD Instinct MI355X platform using the FP6 Matrix datatype. Results subject to change when products are released in market.

MI350-048: Based on AMD internal testing as of 6/9/2025. Using 8 GPU AMD Instinct™ MI355X Platform measuring text generated offline inference throughput for Llama 3.3-70B chat model (FP4) compared 8 GPU AMD Instinct™ MI300X Platform performance with (FP8). MI355X ran 8xTP1 (8 copies of model, one per GPU) compared to MI300X running 8xTP1 (8 copies of model, one per GPU). Tests were conducted using a synthetic dataset with different combinations of 128 and 2048 input tokens, and 128 output tokens. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-049: Based on performance testing by AMD Labs as of 6/6/2025, measuring the text generated inference throughput on the LLaMA 3.1-405B model using the FP4 datatype with input length of 128 tokens and an output length of 2048 tokens on the AMD Instinct™ MI355X 8x GPU, and published results for the NVIDIA B200 HGX 8xGPU. Performance per dollar calculated with current pricing for NVIDIA B200 and Instinct MI355X based cloud instances. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations. Current customer pricing as of June 10, 2025, and subject to change

MI400-001: Performance projection as of 06/05/2025 using engineering estimates based on the design of a future AMD Instinct MI400 Series GPU compared to the Instinct MI355x, with 2K and 16K prefill with TP8, EP8 and projected inference performance, and using a GenAI training model evaluated with GEMM and Attention algorithms for the Instinct MI400 Series .Results may vary when products are released in market.

Endnotes

MI350-54: Based on calculations by AMD internal testing as of 6/10/2025. Using 8 GPU AMD Instinct™ MI355X Platform for overall GPU-normalized Training Throughput (time to complete) for fine-tuning using the Llama2-70B LoRA chat model (FP8) compared to published 8 GPU Nvidia H200 Platform performance (FP8). Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI350-055: Based on engineering projections by AMD Performance Labs in June 2025, to estimate the peak theoretical precision performance of seventy-two (72) AMD Instinct™ MI400X GPUs (Rack) using the FP4 Matrix datatype vs. an 8xGPU AMD Instinct MI00 platform using the FP16 Matrix datatype. Results subject to change when products are released in market.

MI350-056: Based on calculations by AMD Performance Labs in June 2025, to determine the peak theoretical precision performance of 8x GPU AMD Instinct MI355X platform with FP6 Matix datatype vs. an 8x GPU AMD Instinct MI325X/MI300X platforms with FP8 Matrix datatype. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of the latest drivers and optimizations.

MI350-057: *Based on calculations by AMD Performance Labs in June 2025, to determine the peak theoretical precision performance of 8x GPU AMD Instinct MI325X/MI300X platform with FP8 Matix datatype vs. an 8x GPU AMD Instinct MI250X platform with FP16 Matrix datatype. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of the latest drivers and optimizations.

MI350-058: *Based on calculations by AMD Performance Labs in June 2025, to determine the peak theoretical precision performance of 8x GPU AMD Instinct MI325X/MI300X platform with FP8 Matix datatype vs. an 8x GPU AMD Instinct MI250X platform with FP16 Matrix datatype. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of the latest drivers and optimizations.

VEN-003: PCIe Gen comparison based on PCI-SIG published statements, <https://pcisig.com/pci-express-6.0-specification>. 2P 6th Gen EPYC CPU with 128 lanes of PCIe Gen 6 and 5th Gen EPYC with 128 lanes of PCIe Gen 5 as of 6/3/2025. PCIe is a registered trademark of PCI-SIG Corporation