# 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, the evolving AI landscape, AMD's ability to advance AI, and the growing AMD EPYC™ market share, 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.
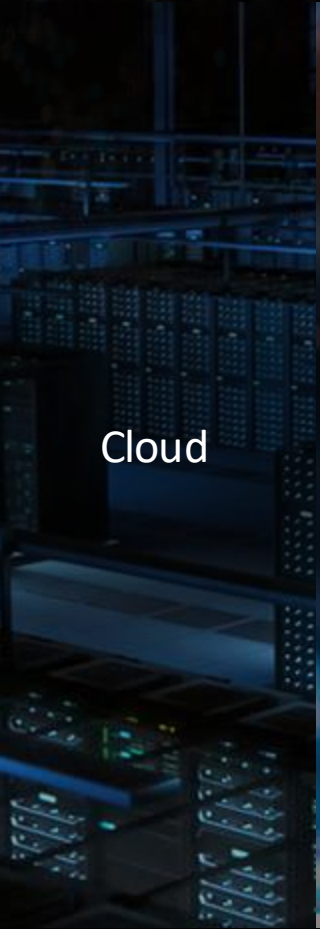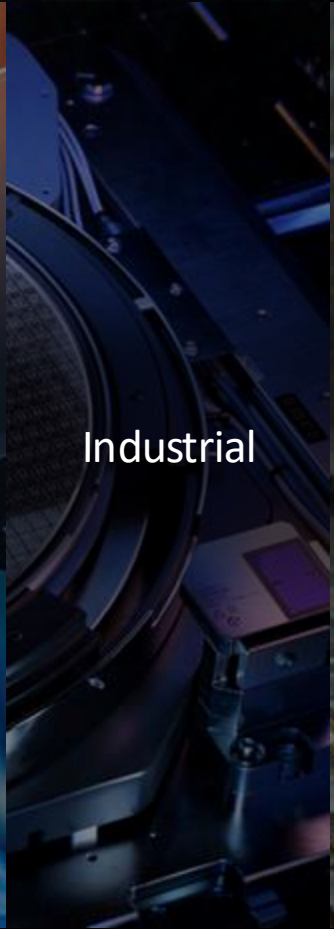
# AI

## Most transformational technology in 50 years

Agents

Smarter Cities

Robotics

Healthcare

Research

Supply Chain

**AMD**
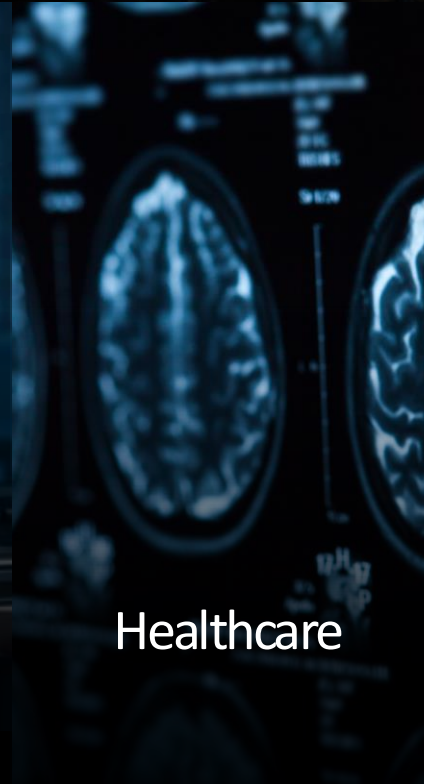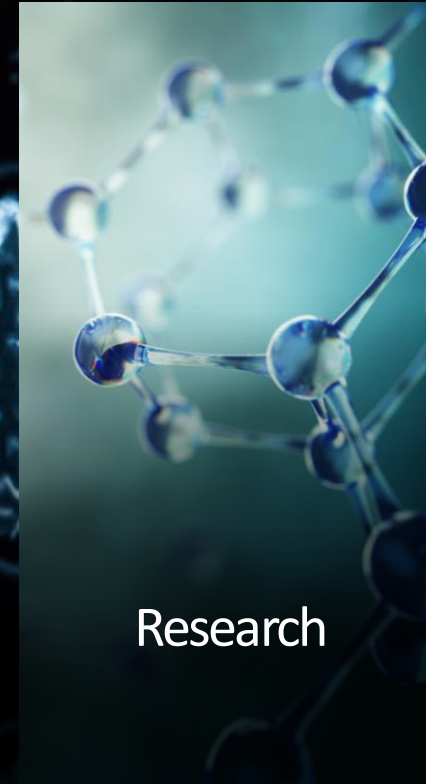AI Platforms

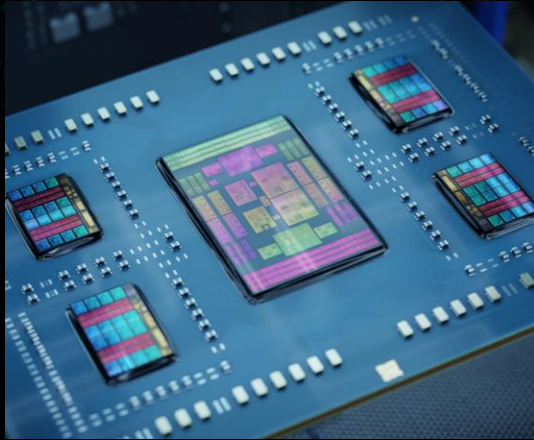Unmatched portfolio of training and inference compute engines

Open software solutions

AI ecosystem with deep co-innovation

Cluster level systems design

# Evolving AI Landscape

## Innovation moving from silicon to nodes to racks to clusters



**Silicon**



**Server**



**Rack**



**Data Center**

# Leadership Engines for Enterprise AI Workloads

**AMD EPYC**

x86 CPU

Data from x86-Based Systems

Data Input → Data Cleaning → Pre-Processing → Model Training → Deployment

**AMD INSTINCT**

GPU

LLM Training and Inference

◄ Agentic AI ►

From analytics to generative AI to agentic AI

# AMD EPYC™ record market share...and growing



| | 1st Gen Processor Family | 2nd Gen Processor Family | | 3rd Gen Processor Family | | 4th Gen Processor Family | |
|---|---|---|---|---|---|---|---|

Market share line chart:
- 2018: 2%
- 2019 — 2nd Gen Processor Family
- 2020: 8%
- 2021 — 3rd Gen Processor Family
- 2022: 27%
- 2023 — 4th Gen Processor Family
- 2024: 31%
- H1'24: 34%
- 2025

**>350 OEM Platforms** | **>950 Cloud Instances**

Source: Mercury Research Sell-in Revenue Shipment Estimates

# #1 CPU for hyperscalers

aws    Alibaba Cloud    Microsoft Azure    Google Cloud    IBM Cloud    ORACLE    Meta    Tencent

Hyperscale leaders power internal workloads
with AMD, serving billions worldwide
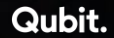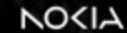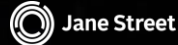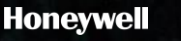
NETFLIX    Office 365    ORACLE EXADATA    salesforce    SAP    Uber    zoom

# Trusted by industry leaders on-prem

# AMD Instinct™ MI300 Series

## Powering the most popular Gen AI platforms

**Microsoft**   **OpenAI**   **Meta**

cohere   stability.ai   essential AI   LAMINI   Reka   LUMA AI   Lepton AI   Fireworks AI

databricks   310.AI   scale   MOREH   World Labs   anyscale   clarifai

UbiOps   FlexAI   OPEN INNOVATION   ZYPHRA   Rhymes   rapt.ai   CLEAR ML   NEURAL MAGIC

# AMD INSTINCT

## Solutions from leading OEMs and cloud

Launching Today: "Turin"

# 5th Gen AMD EPYC™

## World's best CPU for Cloud, Enterprise and AI

ZEN 5

3nm
4nm

150 billion
transistors

Up to 192 cores
384 threads

17% IPC uplift
Full AVX512

Up to 5 GHz

*~17% Across 36 cloud and enterprise workloads
As of 10/1/2024. See endnotes 9xx5-001, EPYC-029C

# "Turin" Continues AMD EPYC™ Leadership

**Consistent x86 ISA**
Consistent IPC

**SP5 Socket**
"Genoa" Compatible

**8 to 192 cores**
125W to 500W

Up to
**12Ch DDR5-6400**
128 PCIe® 5.0/CXL® 2.0

**Confidential Compute**
with Trusted I/O

# Scale-Up

Up to

16 "Zen 5" CCDs

128 Cores · 256 Threads

# Industry's Highest Performing Server CPU

**AMD EPYC™**
**5th Gen 9965**

**AMD EPYC™**
**4th Gen 9754**

**Intel™ Xeon®**
**5th Gen 8592+**

192 cores — 2.7

128 cores — 1.7

64 cores — 1.0

SPECrate®_2017_int_base

# 2.7x

vs. top-of-stack
"Emerald Rapids"

# 60% More Performance at the Same Licensing Cost

**AMD EPYC™**
5th Gen 9575F

64 cores                                    1.6

**AMD EPYC™**
4th Gen 9554

64 cores                          1.2

**Intel™ Xeon®**
5th Gen 8592+

64 cores                     1.0

Virtualized Infrastructure
VMmark® 4.0

up to **1.6x**

Performance per core in virtualized infrastructure

# Grow Your Database and Media Processing Capabilities



Open-Source Database
MySQL OLTP

- 1.9x
- 3.9x

Video Transcoding
FFMPEG

- 1.9x
- 4x

5th Gen Intel®
Xeon® 8592+ 64C

4th Gen AMD
EPYC™ 9654 96C

5th Gen AMD
EPYC™ 9965 192C

up to **4x**

Throughput performance

# Fastest CPU for the Most Challenging HPC Problems



up to

# 3.9x

Improved time to insight

**Dense Linear Solver**
HPL

1.1x  3.1x

**Molecular Dynamics**
GROMACS

1.9x  3.9x

5th Gen Intel®
Xeon® 8592+ 64C

4th Gen AMD
EPYC™ 9654 96C

5th Gen AMD
EPYC™ 9965 192C

As of 9/30/2024. See endnote 9xx5-038, 9xx5-039 and 9xx5-022.

# End-to-End AI and Inference Performance



up to **3.8x**

AI performance on CPU

**Machine Learning**
XGBoost (Higgs)

3.0x
1.2x

**End-to-End AI**
Workload Derived from TPCx-AI

3.8x
1.7x

5th Gen Intel®
Xeon® 8592+ 64C

4th Gen AMD
EPYC™ 9654 96C

5th Gen AMD
EPYC™ 9965 192C

# 1,000 legacy servers

## 2P Intel® Xeon® Platinum 8280 servers

To deliver 391,000 unit of integer performance

Servers required to achieve a total of 391,000 SPEC®rate_2017_ int_base performance score. AMD EPYC 9965 SPEC®rate_2017_int_base score is estimated

# 131 modern servers

## 2P AMD EPYC™ 9965

To deliver 391,000 unit of integer performance

# 7:1 consolidation

Use the savings, space and power to grow your business

~**87%** fewer servers | ~**67%** lower TCO | ~**68%** less power

2P EPYC™ 9965 vs 2P Intel Xeon® 8280 to deliver 391,000 unit of integer performance

Today at Advancing AI 2024

# AMD end-to-end AI infrastructure leadership

| Data Center CPUs | Data Center GPUs | Networking | AI PCs |
| --- | --- | --- | --- |
| 5th Generation AMD EPYC™ "Turin" | AMD Instinct™ MI325X and MI350 Series | AMD Pensando™ Pollara 400/Salina 400 | 3rd Generation AMD Ryzen™ AI PRO |

$45B

2023

>60% CAGR

Data Center AI Accelerators

Source: AMD

$500B

2028

AMD Instinct™ MI300 GPU

Fastest ramping product in AMD history

# Advancing ROCm™ performance and ecosystem

**~2x**

Improvement in inference and training performance

**~3x***

ROCm now supports 1M+ models out of box*

Deepening partnership with AI ecosystem

**AMD SILO AI**

Growing AI expertise and customer implementations

Demonstrated inference leadership at key customers

Nvidia
**H100**

Conversational AI — 1.3x

Content Generation — 1.3x

AI Agent & Chatbot — 1.2x

Summarization — 1.1x

Llama 3.1  •  405B Latency Improvement

up to **1.3x**

Higher performance across key workloads

Llama 3 ∞    Llama 2 ∞    Mistral    Mixtral

Qwen    CommandR    **Stable Diffusion**

See endnotes MI300-053, MI300-054, MI300-064.

Launching Today

# AMD Instinct™ MI325X GPU
## Extending generative AI leadership

| 256GB HBM3E | 6TB/s | 1.3 PF | 2.6 PF | AMD CDNA 3 |
| --- | --- | --- | --- | --- |
| 1.8x memory | 1.3x bandwidth | 1.3x FP16 | 1.3x FP8 | |

AMD Instinct™
# MI325X Platform

**2 TB** | HBM3E
1.8x memory vs. H200 HGX

**48 TB/s** | Memory Bandwidth
1.3x memory bandwidth vs. H200 HGX

**10.4 PF** | FP16
1.3x compute flops vs. H200 HGX

**20.8 PF** | FP8
1.3x compute flops vs. H200 HGX

See endnotes MI325-001A, MI325-002

*Dense flops

# AMD Instinct™ MI325X Platform
## Leadership Inference performance using 8x MI325X

~1.4x

~1.2x

up to 1.4x

Meta Llama-3.1
405B
Throughput

Meta Llama-3.1
70B
Latency improvement

Inference performance

Nvidia
H200 HGX

AMD Instinct™
MI325X Platform

See endnotes MI325-015, MI325-014

# World-Class Training Performance
## Single GPU and 8 GPU Training

1x GPU

8x GPU

~1.1x

~1x

| | |
|---|---|
| H200 | MI325X |

Meta Llama-2
7B

| | |
|---|---|
| H200 HGX | MI325X Instinct™ Platform |

Meta Llama-2
70B

| | |
|---|---|
| Nvidia H200 HGX | AMD Instinct™ MI325X Platform |

See endnotes MI325-013, MI325-012.

# AMD Instinct™ Annual Roadmap Cadence

**AMD** ◢
CDNA 3

AMD Instinct™
MI300X

2023

**AMD** ◢
CDNA 3

AMD Instinct™
MI325X

2024

**AMD** ◢
CDNA 4

AMD Instinct™
MI350
S E R I E S

2025

Roadmap subject to change

Previewing today

# AMD Instinct™ MI350 Series
## Continued Gen AI Leadership

| 3nm | Up to 288GB | FP4 / FP6 | AMD |
|-----|-------------|-----------|-----|
| Process Node | HBM3E | Datatype Support | CDNA 4 |

**Planned availability 2H 2025**

# AMD CDNA 4

| | |
|---|---|
| **New Datatypes** | **~3.5x** AI Flops vs. FP8 |
| **AI Compute** | **~1.8x** FP16 / FP8 |
| **HBM3E Memory** | **~1.5x** memory capacity<br>memory bandwidth |

# AMD Instinct™ MI355X Accelerator

## Leadership performance for Gen AI

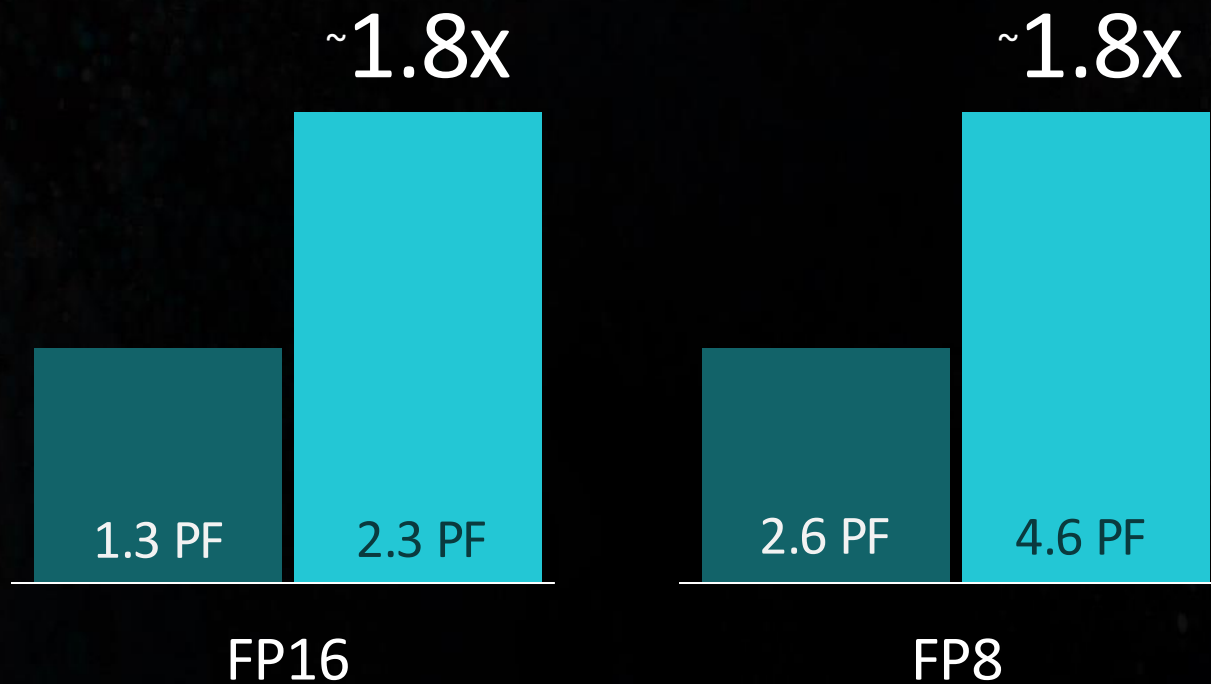~1.8x                    ~1.8x

1.3 PF    2.3 PF    2.6 PF    4.6 PF

FP16                      FP8

# 9.2 PF

## Introducing FP6 and FP4

AMD Instinct™ **MI325X**          AMD Instinct™ **MI355X**

# Open, modular software stack

**AMD ROCm**

**AI Models and Algorithms**

PyTorch  ONNX  JAX  TensorFlow

**Libraries**

**Compilers and Tools**

**Runtime**

**AMD GPUs** | AMD INSTINCT  AMD RADEON PRO W7600 AMD RADEON

**Support all major frameworks and models**

**Expanded Gen AI optimizations**

New algorithms

New libraries

Expanding platform support

**Extended developer support**

# Deep collaboration with developer community

**PyTorch**

Day 0 support for latest features

**Triton**

Vendor agnostic compiler support

**Hugging Face**

Nightly CI/CD ensuring all models work out-of-box

vLLM   SGLang

JAX   Tensor Flow

ONNX Runtime   OpenXLA

DeepSpeed   MLIR | IREE

Increasing open-source contributions and expanding footprint

# AMD Instinct™ MI300X Accelerator
## Expanding out-of-box support on popular generative AI models

GPT-4  GPT-4o  Llama 2  Llama 3  MISTRAL AI_  Mixtral  Grok 1  Stable Diffusion

DBRX  Qwen  Zamba  Yi  StarCoder  Phi  Flux  Hugging Face

CommandR  Gemma  BLOOM  MPM4  GPT-NeoX  OLMo  Falcon LLM  Aria  MPT

# Generational inference improvement
## ROCm™ 6.2 vs. ROCm 6.0

~2.4x
average
performance
improvement

~1.9x
~2.1x
~2.6x
~2.6x
~2.8x

Mixtral 8x22B
Mixtral 8x7B
Qwen2 72B
Llama3.1 70B
Llama3.1 8B

Runtime Optimization | Kernel Fusion | Collective Communication | Subgraph

# Generational training improvement
## ROCm™ 6.2 vs. ROCm 6.0

~1.7x

~1.9x

~1.9x

~1.8x

average
performance
improvement

| Llama2 70B | Qwen 1.5 14B | Llama2 7B |
|---|---|---|

Flash Attention 3 | Parallelization Strategy | Kernels | Scale out Efficiency

# Advancing Data Center Solutions

Data Center CPUs

Data Center GPUs

Networking

# CPU Performance Enhances GPU Performance

AMD EPYC™
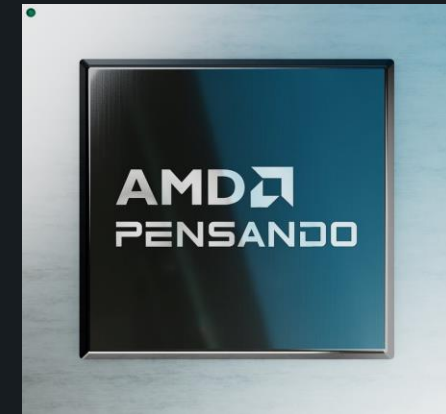5th Gen 9575F

5.0GHz

Intel™ Xeon®
5th Gen 8592+

3.9GHz

CPU Host Processor
Max Frequency

# 28%

Faster processing for GPU orchestration tasks

Speeds up data prep, memory copies, kernel launch and task orchestration

# AMD EPYC™ 9575F

## Purpose built for GPU host nodes

~**700,000**
more inference tokens/s

on 1K node AI cluster running Llama3.1-70B

Up to
**20%**
faster training

with Stable Diffusion XL V2

Llama 3.1: 8.8% more perf on 1000 Node Cluster of Turin + 8xMI300X vs Emerald Rapids + 8xMI300X on Llama3.1-70B with 128 Input tokens, 2048 output tokens, batch size 1000
Stable Diffusion XL V2: 20% better training time on Turin + 8xMI300X vs Emerald Rapids + 8xMI300X
As of 10/4/2024. See endnote 9xx5-087, 9xx5-059a.

# Programmable DPU
## Evolving front-end network

Enables:
### Faster Data Integration

Storage offload
and acceleration

Enables:
### SDN and Security

Evolving network
infrastructure services
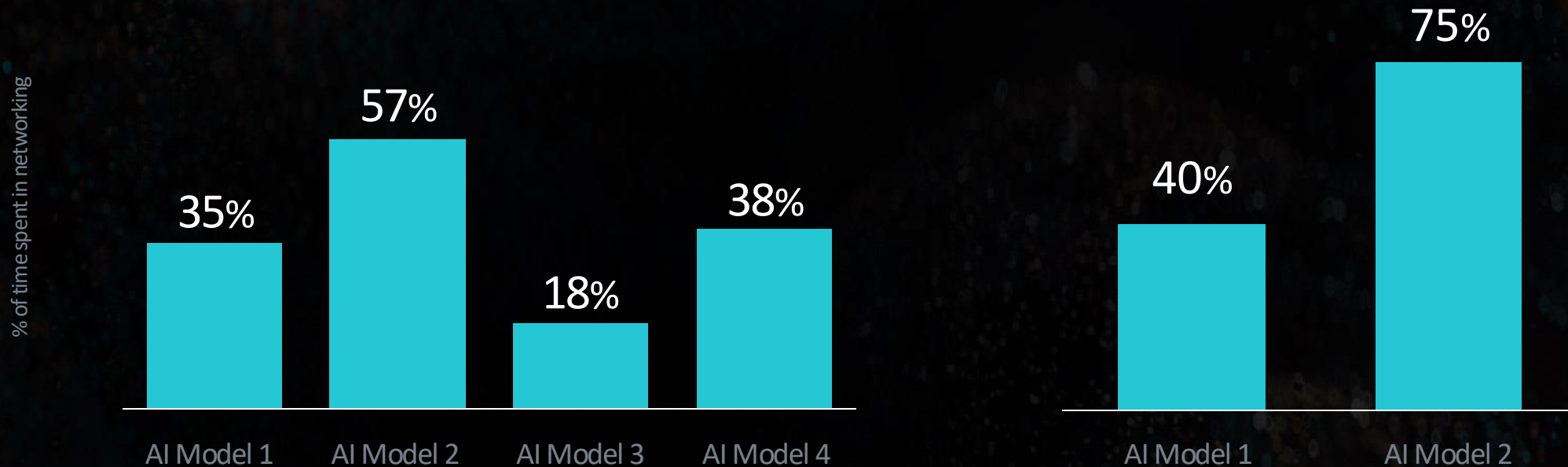
Secure multi-tenant access
Data privacy and integrity

Enables:
### Zero CPU Overhead

DPU accelerates infrastructure
services at line rate

Dedicated CPU for
AI workload processing

# Back-end networks drive AI system performance

% of time spent in networking

35%
57%
18%
38%

AI Model 1    AI Model 2    AI Model 3    AI Model 4

At an average 30% of training cycle time
is elapsed in waiting for networking

40%
75%

AI Model 1    AI Model 2

Communication accounts for 40%-75% of time
with Training and Distributed Inference Models[2]

Source: 1) 2022 OCP Keynote by Alexis Bjorlin, VP at Meta, 2) Computation vs. Communication Scaling for Future Transformers on Future Hardware  https://arxiv.org/pdf/2302.02825

# Advancing Data Center Solutions

| Data Center CPUs | Data Center GPUs | Networking |
|---|---|---|
| AMD EPYC | AMD INSTINCT | AMD PENSANDO |

# Ethernet is always the preferred choice

> **50%**
TCO Saving

InfiniBand

Ethernet
RoCEv2

**Total Cost of Ownership[1]**
Lower is better

1,000,000+ GPU

up to
**48,000** GPU

InfiniBand

Ethernet
RoCEv2

**Scalability**
Higher is better

~ **95**%

< **50**%

General
Purpose
Network

AI Back-End
Network

Network Utilization

# The challenge of high network utilization

## AI backend networks drive sustained data transfers

| Intelligent Load Balancing | Congestion Management | Fast Failover and Loss Recovery |

# Ultra Ethernet
## Consortium

Evolve ethernet as an open, interoperable, high performance, full-communications stack architecture to meet the growing network demands of AI and HPC at scale

UEC 1.0 Specification  - Q1CY25

Performant | Scalable | Cost Effective

# Ultra Ethernet
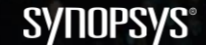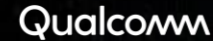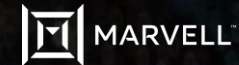## Consortium

## Steering Members

AMD · ARISTA · BROADCOM · CISCO · EVIDEN an atos business

Hewlett Packard Enterprise · intel · Meta · Microsoft · ORACLE

# Ultra Ethernet
## Consortium

## General Members

Alibaba Cloud | ARRCUS NETWORK DIFFERENT | Baidu 百度 | 世纪互联 VNET | ByteDance | cādence | CORNELIS NETWORKS | DELL Technologies | enfabrica

HUAWEI | IBM | JUNIPER NETWORKS | KEYSIGHT TECHNOLOGIES | Lawrence Livermore National Laboratory | Lenovo | MARVELL | H3C

NOKIA | NVIDIA | Preferred Networks | PURESTORAGE | Qualcomm | Spirent Promise. Assured. | SYNOPSYS | ZTE

• Total 97 Members •

# RDMA outperforms RoCEv2

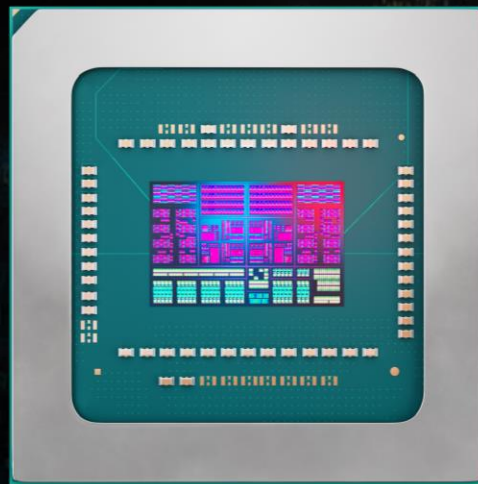## 6x faster
message completion time

## 5x faster
collective completion time

Intelligent packet spray and
in-order message delivery

Path aware
congestion avoidance

Selective retransmission
and fast loss recovery

* Reference: STrack: A Reliable Multipath Transport for AI/ML Clusters, July 2024

# 3rd Gen AMD P4 Engines

## Deliver network innovation at the speed of AI

| 120M Packets/s | 400Gb/s | 5M Connections/s |
| --- | --- | --- |
| Fully Programmable | Wire Rate | Concurrent Services<br>(SDN, Security, Storage Acceleration) |

Solutions you can rely on for your business

350+
server platforms

950+
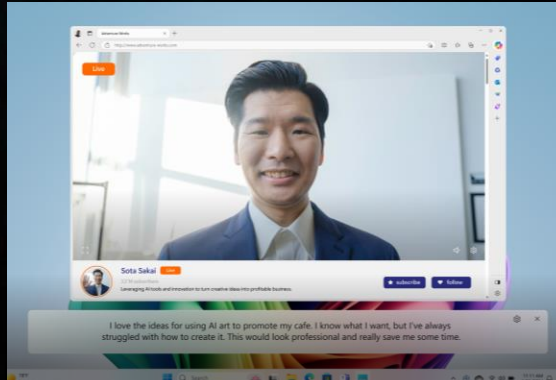cloud instances

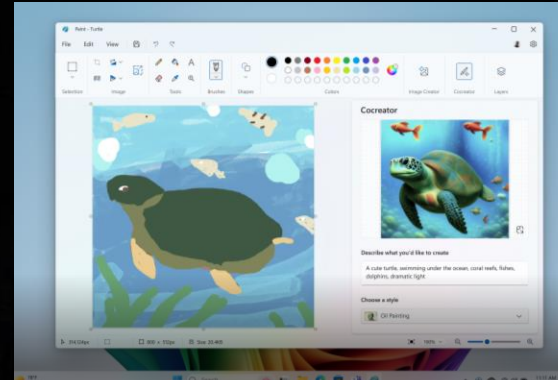# Unprecedented transformational experiences with next-gen AI PCs



## Enterprise Productivity

Business LLMs for data processing

Multilingual chatbots
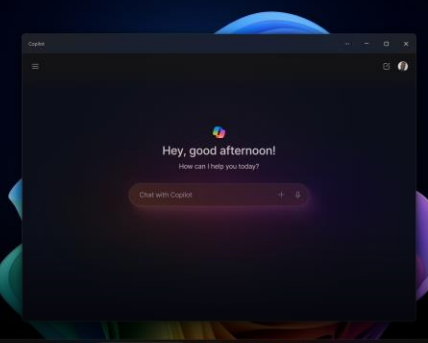
Software coding assistants

Real-time threat detection

## Immersive Collaboration

Live captions with translation

Speech recognition and transcription

Intelligent meeting assistance

Sentiment analysis

## Revolutionary Creation and Editing
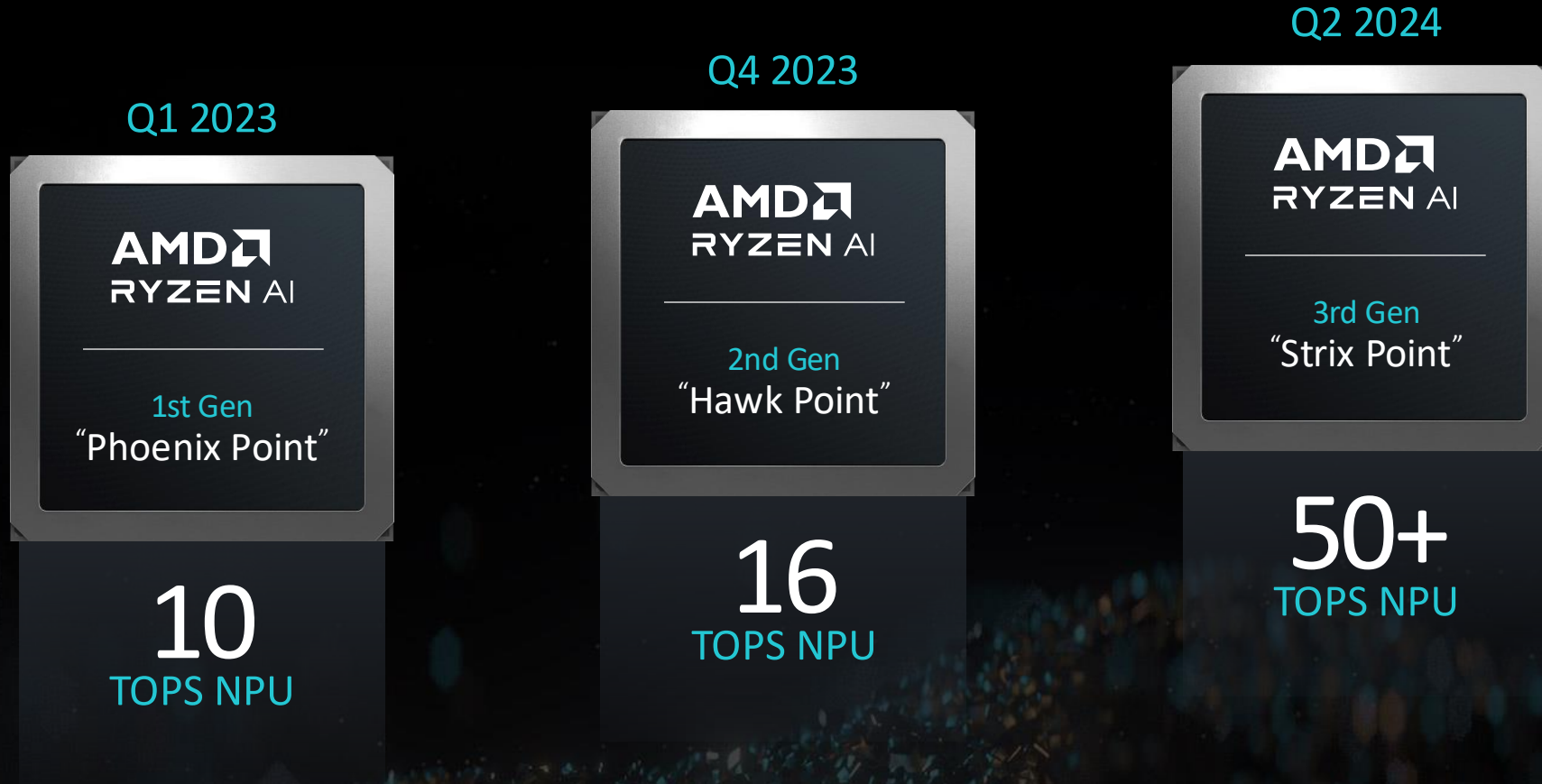
Automated content generation

Multimedia production

Art and design assist

Advanced video & audio effects

## Personal AI Assistance

Smart information retrieval

Document analysis

Calendar management

Travel planning

# Best Processors for Copilot+ Enterprise PCs

**AMD RYZEN AI PRO**

BUFFERZONE · zoom · splashtop · WHISPP · webex by CISCO · Microsoft · LM Studio · grammarly · MAXON

Blackmagicdesign · OBS Open Broadcaster Software · BORISFX · SOLIDWORKS · Avid · GoPro · Topaz Labs · blender · Adobe

**Enterprise AI Experiences**

Pretrain ❯ Quantize ❯ Deploy

Open Platform

**Simplified Development Framework**

Top-Tier Performance

Multi-Day Battery Life

Security
AMD PRO TECHNOLOGY

Reliability
AMD PRO TECHNOLOGY
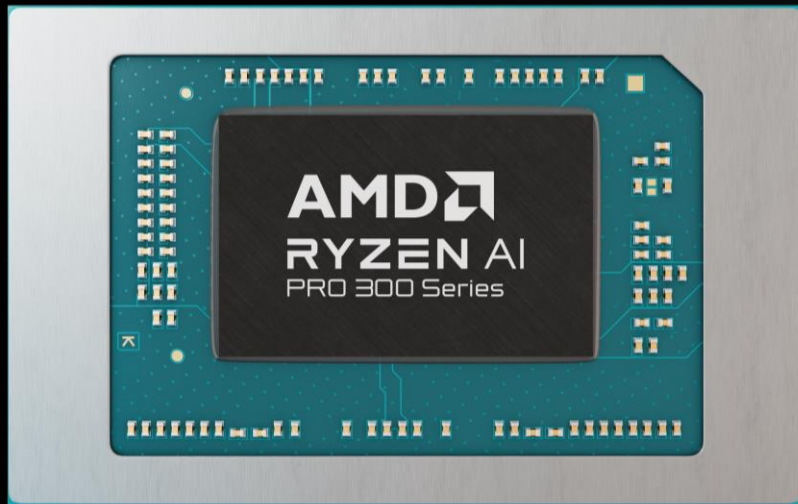
Manageability
AMD PRO TECHNOLOGY

**Ryzen™ AI PRO**
**Designed for Enterprise**

See endnote GD-173a, STXP-04.

Up to

# 1.4x

multithreaded performance

Up to

# 23hrs
## Multi day battery life
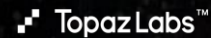
AMD Ryzen™ AI 9 HX PRO 375, VIdeo Playback. See endnote STXP-30

Up to

# 9hrs

battery life with Microsoft Teams

# Enterprise AI PC Application Ecosystem

Adobe   Microsoft   webex by CISCO   ZOOM

splashtop   WHISPP   DS SOLIDWORKS   BUFFERZONE   Camo   bitdefender secure your every bit   LM Studio   grammarly

Blackmagicdesign   nero   BORIS FX   Avid   Rhinoceros   GoPro Be a HERO.   Topaz Labs™   ARKRUNR   voicemy.ai
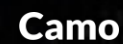
blender   RADiCAL AI-POWERED 2D ANIMATION   CyberLink   AFFINITY Photo 2   ACCA ACCA SOFTWARE   convai   OBS Open Broadcaster Software   MAXON

# Endnotes

# Endnotes

# Endnotes

# Endnotes

9xx5-048: AMD EPYC™ 9005 Series processors require OEM enablement and a BIOS update from your server or motherboard manufacturer if used with a motherboard designed for the SP5 socketed AMD EPYC™ 9004 Series processors. Contact your system manufacturer prior to purchase to determine compatibility.

9xx5-059A: Stable Diffusion XL v2 training results based on AMD internal testing as of 10/10/2024. SDXL configurations: DeepSpeed 0.14.0, TP8 Parallel, FP8, batch size 24, results in seconds 2P AMD EPYC 9575F (128 Total Cores) with 8x AMD Instinct MI300X-NPS1-SPX-192GB-750W, GPU Interconnectivity XGMI, ROCm™ 6.2.0-66, 2304GB 24x96GB DDR5-6000, BIOS 1.0 (power determinism = off), Ubuntu® 22.04.4 LTS, kernel 5.15.0-72-generic, 334.80 seconds. 2P Intel Xeon Platinum 8592+ (128 Total Cores) with 8x AMD Instinct MI300X-NPS1-SPX-192GB-750, GPU Interconnectivity XGMI, ROCm 6.2.0-66, 2048GB 32x64GB DDR5-4400, BIOS 2.0.4, (power determinism= off), Ubuntu 22.04.4 LTS, kernel 5.15.0-72-generic, 400.43 seconds. For 19.600% training performance increase. Results may vary due to factors including system configurations, software versions and BIOS settings.

9xx5-069A: SPECrate®2017_int_base comparison based on published scores from www.spec.org as of 10/10/2024. Generational scores are based on highest published scores from www.spec.org from respective launch years. 2P AMD EPYC 9965 (3000 SPECrate®2017_int_base, 384 Total Cores, https://www.spec.org/cpu2017/results/res2024q4/cpu2017-20240923-44837.pdf) 2P AMD EPYC 9654 (1790 SPECrate®2017_int_base, 192 Total Cores, https://www.spec.org/cpu2017/results/res2022q4/cpu2017-20221024-32607.html ) 2P AMD EPYC 7763 (861 SPECrate®2017_int_base, 128 Total Cores, https://www.spec.org/cpu2017/results/res2021q4/cpu2017-20211121-30148.html ) 2P AMD EPYC 7742 (701 SPECrate®2017_int_base, 128 Total Cores, https://www.spec.org/cpu2017/results/res2019q4/cpu2017-20191125-20001.html ) 2P AMD EPYC 7601 (275 SPECrate®2017_int_base, 64 Total Cores, https://www.spec.org/cpu2017/results/res2017q4/cpu2017-20171211-01594.html) SPEC®, SPEC CPU®, and SPECrate® are registered trademarks of the Standard Performance Evaluation Corporation. See www.spec.org for more information. Intel CPU TDP at https://ark.intel.com/. SPEC - Standard Performance Evaluation Corporation

9xx5-071: VMmark® 4.0.1 host/node FC SAN comparison based on "independently published" results as of 10/10/2024. Configurations: 2 node, 2P AMD EPYC 9575F (128 total cores) powered server running VMware ESXi8.0 U3, 3.31 @ 4 tiles, https://www.infobellit.com/BlueBookSeries/VMmark4-FDR-1003. 2 node, 2P AMD EPYC 9554 (128 total cores) powered server running VMware ESXi 8.0 U3, 2.64 @ 3 tiles, https://www.infobellit.com/BlueBookSeries/VMmark4-FDR-1002. 2 node, 2P Intel Xeon Platinum 8592+ (128 total cores) powered server running VMware ESXi 8.0 U3, 2.06 @ 2.4 Tiles, https://www.infobellit.com/BlueBookSeries/VMmark4-FDR-1001. VMmark is a registered trademark of VMware in the US or other. countries.

9xx5-083::5th Gen EPYC processors support DDR5-6400 MT/s for targeted customers and configurations. 5th Gen production SKUs support up to DDR5-6000 MT/s to enable a broad set of DIMMs across all OEM platforms and maintain SP5 platform compatibility.

9xx5-087: As of 10/10/2024; this scenario contains several assumptions and estimates and, while based on AMD internal research and best approximations, should be considered an example for information purposes only, and not used as a basis for decision making over actual testing. Referencing 9XX5-056A: "2P AMD EPYC 9575F powered server and 8x AMD Instinct MI300X GPUs running Llama3.1-70B select inference workloads at FP8 precision vs 2P Intel Xeon Platinum 8592+ powered server and 8x AMD Instinct MI300X GPUs has ~8% overall throughput increase across select inference use cases" and 8763.52 tokens/s (9575F) versus 8,048.48 tokens/s (8592+) at 128 input / 2048 output tokens, 500 prompts for 1.089x the tokens/s or 715.04 more tokens/s. 1 Node = 2 CPUs and 8 GPUs. Assuming a 1000 node cluster, 1000 * 715.04 = 715,040 tokens/s. For ~700,000 more tokens/s. Results may vary due to factors including system configurations, software versions and BIOS settings.

# Endnotes

99xx5TCO-002A: This scenario contains many assumptions and estimates and, while based on AMD internal research and best approximations, should be considered an example for information purposes only, and not used as a basis for decision making over actual testing. The AMD Server & Greenhouse Gas Emissions TCO (total cost of ownership) Estimator Tool - version 1.12, compares the selected AMD EPYC™ and Intel® Xeon® CPU based server solutions required to deliver a TOTAL_PERFORMANCE of 391000 units of SPECrate2017_int_base performance as of October 10, 2024. This estimation compares a legacy 2P Intel Xeon 28 core Platinum_8280 based server with a score of 391 versus 2P EPYC 9965 (192C) powered server with a score of 3000 (https://www.spec.org/cpu2017/results/res2024q4/cpu2017-20240923-44837.pdf) along with a comparison upgrade to a 2P Intel Xeon Platinum 8592+ (64C) based server with a score of 1130 (https://spec.org/cpu2017/results/res2024q3/cpu2017-20240701-43948.pdf). Actual SPECrate®2017_int_base score for 2P EPYC 9965 will vary based on OEM publications. Environmental impact estimates made leveraging this data, using the Country / Region specific electricity factors from the 2024 International Country Specific Electricity Factors 10 – July 2024, and the United States Environmental Protection Agency 'Greenhouse Gas Equivalencies Calculator'. For additional details, see https://www.amd.com/en/legal/claims/epyc.html#q=epyc4#SP9xxTCO-002A.

EPYC-029C: Comparison based on thread density, performance, features, process technology and built-in security features of currently shipping servers as of 10/10/2024. EPYC 9005 series CPUs offer the highest thread density [EPYC-025B], leads the industry with 500+ performance world records [EPYC-023F] with performance world record enterprise leadership Java® ops/sec performance [EPYCWR-20241010-260], top HPC leadership with floating-point throughput performance [EPYCWR-2024-1010-381], AI end-to-end performance with TPCx-AI performance [EPYCWR-2024-1010-525] and highest energy efficiency scores [EPYCWR-20241010-326]. The 5th Gen EPYC series also has 50% more DDR5 memory channels [EPYC-033C] with 70% more memory bandwidth [EPYC-032C] and supports 70% more PCIe® Gen5 lanes for I/O throughput [EPYC-035C], has up to 5x the L3 cache/core [EPYC-043C] for faster data access, uses advanced 3-4nm technology, and offers Secure Memory Encryption + Secure Encrypted Virtualization (SEV) + SEV Encrypted State + SEV-Secure Nested Paging security features. See the AMD EPYC Architecture White Paper (https://library.amd.com/l/3f4587d147382e2/) for more information.

MI300-53: Testing completed on 05/28/2024 by AMD performance lab attempting text generated throughput measured using Mistral-7B model comparison. Tests were performed using batch size 1 and 2048 input tokens and 2048 output tokens for Mistral-7B **Configurations: \**2P AMD EPYC 9534 64-Core Processor based production server with 8x AMD InstinctTM MI300X (192GB, 750W) GPU, Ubuntu® 22.04.1, and ROCm™ 6.1.1 Vs. 2P Intel Xeon Platinum 8468 48-Core Processor based production server with 8x NVIDIA Hopper H100 (80GB, 700W) GPU, Ubuntu 22.04.3, and CUDA® 12.2. Only 1 GPU on each system was used in this test. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI300-54: Testing completed on 05/28/2024 by AMD performance lab attempting text generated Llama3-70B using batch size 1 and 2048 input tokens and 128 output tokens for each system. **Configurations:** 2P AMD EPYC 9534 64-Core Processor based production server with 8x AMD InstinctTM MI300X (192GB, 750W) GPU, Ubuntu® 22.04.1, and ROCm™ 6.1.1 Vs. 2P Intel Xeon Platinum 8468 48-Core Processor based production server with 8x NVIDIA Hopper H100 (80GB, 700W) GPU, Ubuntu 22.04.3, and CUDA® 12.2 **8 GPUs on each system was used in this test.**
Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI300-62: Testing conducted by internal AMD Performance Labs as of September 29, 2024 inference performance comparison between ROCm 6.2 software and ROCm 6.0 software on the systems with 8 AMD Instinct™ MI300X GPUs coupled with Llama 3.1-8B, Llama 3.1-70B, Mixtral-8x7B, Mixtral-8x22B, and Qwen 72B models. ROCm 6.2 with vLLM 0.5.5 performance was measured against the performance with ROCm 6.0 with vLLM 0.3.3, and tests were performed across batch sizes of 1 to 256 and sequence lengths of 128 to 2048. Configurations: 1P AMD EPYC™ 9534 CPU server with 8x AMD Instinct™ MI300X (192GB, 750W) GPUs, Supermicro AS-8125GS-TNMR2, NPS1 (1 NUMA per socket), 1.5 TiB (24 DIMMs, 4800 mts memory, 64 GiB/DIMM), 4x 3.49TB Micron 7450 storage, BIOS version: 1.8, , ROCm 6.2.0-00, vLLM 0.5.5, PyTorch 2.4.0, Ubuntu® 22.04 LTS with Linux kernel 5.15.0-119-generic. Vs. 1P AMD EPYC 9534 CPU server with 8x AMD Instinct™ MI300X (192GB, 750W) GPUs, Supermicro AS-8125GS-TNMR2, NPS1 (1 NUMA per socket), 1.5TiB 24 DIMMS, 4800 mts memory, 64 GiB/DIMM), 4x 3.49TB Micron 7450 storage, BIOS version: 1.8, ROCm 6.0.0-00, vLLM 0.3.3, PyTorch 2.1.1, Ubuntu 22.04 LTS with Linux kernel 5.15.0-119-generic. Server manufacturers may vary configurations, yielding different results. Performance may vary based on factors including but not limited to different versions of configurations, vLLM, and drivers.

# Endnotes

MI300-63: Testing conducted by internal AMD Performance Labs as of September 29, 2024 training performance comparison between ROCm 6.2 software with compared to ROCm 6.0 software both with Megatron-LM on systems with 8 AMD Instinct™ MI300X GPUs running Llama 2-7B, Llama 2-70B (4K), Qwen1.5-14B models using custom docker container for each system. ROCm 6.2 with megatron-LM TFLOPs was measured against the TFLOPs with ROCm 6.0 with megatron-LM. Configurations:CPU: 1P AMD EPYC 9454 48-core processor,Host memory: 2x3.5 T GB GPU: AMD Instinct MI300X. 1P AMD EPYC™ 9454 CPU, 8x AMD Instinct™ MI300X (192GB, 750W) GPUs, American Megatrends International LLC BIOS version: 1.8, ROCm 6.2 internal release, Megatron-LM code branches hanl/disable_te_llama2 for Llama 2-7B, guihong_dev for LLama 2-70B, renwuli/disable_te_qwen1.5 for Qwen1.5-14B, PyTorch 2.4, Ubuntu 22.04 LTS with Linux kernel 5.15.0-117-generic. Vs. 1P AMD EPYC 9454 CPU 48-core processor, 8x AMD Instinct™ MI300X (192GB, 750W) GPUs, American Megatrends International LLC BIOS version: 1.8, ROCm 6.0.0, Megatron-LM code branches hanl/disable_te_llama2 for Llama 2-7B, guihong_dev for LLama 2-70B, renwuli/disable_te_qwen1.5 for Qwen1.5-14B, PyTorch 2.2, Ubuntu 22.04 LTS with Linux kernel 5.15.0-72-generic. Server manufacturers may vary configurations, yielding different results. Performance may vary based on factors including but not limited to different versions of configurations, megatron-LM, and drivers. Results: MI300X with ROCm 6.2 delivers average 1.83X the (83% higher) training throughput than ROCm 6.0.

MI300-64: Based on testing completed on 10/09/2024 by AMD performance lab measuring overall throughput for text generated using LLaMA 3.1-405B model using FP8 datatype. Test was performed using various input token length and an output token length for the following configurations of AMD Instinct™ MI325X 8xGPU platform and NVIDIA H100 platform. Configurations: AMD Instinct™ MI300X platform: Supermicro AS - 8125GS-TNMR2 server with 2x AMD EPYC 9654 Processors, 8x AMD MI300X (192GB, 750W) GPUs, Ubuntu 22.04).ROCm 6.2 NVIDIA H100 HGX platform: Supermicro AS - 8125GS-TNHR server with 2x AMD EPYC 9654 Processors, 8x Nvidia H100 (80GB, 700W) GPUs, Ubuntu 22.04) Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations. MI300-64

MI325-001A: Calculations conducted by AMD Performance Labs as of September 26th, 2024, based on current specifications and /or estimation . The AMD Instinct™ MI325X OAM accelerator will have 256GB HBM3E memory capacity and 6 TB/s GPU peak theoretical memory bandwidth performance. Actual results based on production silicon mayvary. The highest published results on the NVidia Hopper H200 (141GB) SXM GPU accelerator resulted in 141GB HBM3E memory capacity and 4.8 TB/s GPU memory bandwidth performance. https://nvdam.widen.net/s/nb5zzzsjdf/hpc-datasheet-sc23-h200-datasheet-3002446. The highest published results on the NVidia Blackwell HGX B100 (192GB) 700W GPU accelerator resulted in 192GB HBM3E memory capacity and 8 TB/s GPUmemory bandwidth performance.  The highest published results on the NVidia Blackwell HGX B200 (192GB) GPU accelerator resulted in 192GB HBM3E memory capacity and 8 TB/s GPU memory bandwidth performance.  Nvidia Blackwell specifications at https://resources.nvidia.com/en-us-blackwell-architecture?_gl=1*1r4pme7*_gcl_aw*R0NMLjE3MTM5NjQ3NTAuQ2p3S0NBancyNkt4QmhCREVppd0F1NktYdDlweXY1dlUtaHNKNmhPdHM4UVdPSlM3dFdQaE40WkI4THZBaWFVajFyTGhYd3hLQmlZQ3pCb0NsVElROXZEX0J3RQ..*_gcl_au*MTIwNjg4NjU0Ny4xNzExMDM1NTQ3

# Endnotes

MI325-02: Calculations conducted by AMD Performance Labs as of May 28th, 2024 for the AMD Instinct™ MI325X GPU resulted in 1307.4 TFLOPS peak theoretical half precision (FP16), 1307.4 TFLOPS peak theoretical Bfloat16 format precision (BF16), 2614.9 TFLOPS peak theoretical 8-bit precision (FP8), 2614.9 TOPs INT8 floating-point performance. Actual performance will vary based on final specifications and system configuration. Published results on Nvidia H200 SXM (141GB) GPU: 989.4 TFLOPS peak theoretical half precision tensor (FP16 Tensor), 989.4 TFLOPS peak theoretical Bfloat16 tensor format precision (BF16 Tensor), 1,978.9 TFLOPS peak theoretical 8-bit precision (FP8), 1,978.9 TOPs peak theoretical INT8 floating-point performance. BFLOAT16 Tensor Core, FP16 Tensor Core, FP8 Tensor Core and INT8 Tensor Core performance were published by Nvidia using sparsity; for the purposes of comparison, AMD converted these numbers to non-sparsity/dense by dividing by 2, and these numbers appear above. Nvidia H200 source: https://nvdam.widen.net/s/nb5zzzsjdf/hpc-datasheet-sc23-h200-datasheet-3002446 and https://www.anandtech.com/show/21136/nvidia-at-sc23-h200-accelerator-with-hbm3e-and-jupiter-supercomputer-for-2024 Note: Nvidia H200 GPUs have the same published FLOPs performance as H100 products https://resources.nvidia.com/en-us-tensor-core/.

MI325-004: Based on testing completed on 9/28/2024 by AMD performance lab measuring text generated throughput for Mixtral-8x7B model using FP16 datatype. Test was performed using input length of 128 tokens and an output length of 4096 tokens for the following configurations of AMD Instinct™ MI325X GPU accelerator and NVIDIA H200 SXM GPU accelerator. 1x MI325X at 1000W with vLLM performance Vs. 1x H200 at 700W with TensorRT-LLM v0.13 Configurations: AMD Instinct™ MI325X reference platform:
 1x AMD Ryzen™ 9 7950X CPU, 1x AMD Instinct MI325X (256GiB, 1000W) GPU, Ubuntu® 22.04, and ROCm™ 6.3 pre-release Vs NVIDIA H200 HGX platform: Supermicro SuperServer with 2x Intel Xeon® Platinum 8468 Processors, 8x Nvidia H200 (140GB, 700W) GPUs [only 1 GPU was used in this test], Ubuntu 22.04) CUDA® 12.6. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI325-005:: Based on testing completed on 9/28/2024 by AMD performance lab measuring overall latency for LLaMA 3.1-70B model using FP8 datatype. Test was performed using input length of 2048 tokens and an output length of 2048 tokens for the following configurations of AMD Instinct™ MI325X GPU accelerator and NVIDIA H200 SXM GPU accelerator. MI325X at 1000W with vLLM performance: 48.025 sec (latency in seconds) Vs. 1x H200 at 700W with TensorRT-LLM v 0.13: 56.310 sec (latency in seconds) Configurations: AMD Instinct™ MI325X reference platform: 1x AMD Ryzen™ 9 7950X 16-Core Processor CPU, 1x AMD Instinct MI325X (256GiB, 1000W) GPU, Ubuntu® 22.04, and ROCm™ 6.3 pre-release Vs NVIDIA H200 HGX platform: Supermicro SuperServer with 2x Intel Xeon® Platinum 8468 Processors, 8x Nvidia H200 (140GB, 700W) GPUs, Ubuntu 22.04), CUDA 12.6. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.

MI325-006: Based on testing completed on 9/28/2024 by AMD performance lab measuring overall latency for LLaMA 3.1-70B model using FP8 datatype. Test was performed using input length of 2048 tokens and an output length of 2048 tokens for the following configurations of AMD Instinct™ MI325X GPU accelerator and NVIDIA H200 SXM GPU accelerator. MI325X at 1000W with vLLM performance: 48.025 sec (latency in seconds) Vs. 1x H200 at 700W with TensorRT-LLM v 0.13: 56.310 sec (latency in seconds) Configurations: AMD Instinct™ MI325X reference platform: 1x AMD Ryzen™ 9 7950X 16-Core Processor CPU, 1x AMD Instinct MI325X (256GiB, 1000W) GPU, Ubuntu® 22.04, and ROCm™ 6.3 pre-release Vs NVIDIA H200 HGX platform: Supermicro SuperServer with 2x Intel Xeon® Platinum 8468 Processors, 8x Nvidia H200 (140GB, 700W) GPUs, Ubuntu 22.04), CUDA 12.6. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations

# Endnotes

MI355-004: Calculations conducted by AMD Performance Labs as of September 26th, 2024 for the AMD Instinct™ MI300X GPU platform and AMD Instinct™ MI300X GPU platform performance comparing FP16, FP8 and FP4 datatypes. MI355X 8xGPU Platform  Peak theoretical Half Precision (FP16) Performance - 18.5 PFLOPs. Peak theoretical Eight-bit Precision (FP8) Performance - 37 PFLOPs. Peak theoretical Four-bit Precision (FP4) Performance - 74 PFLOPs . MI325X 8xGPU Platform: Peak theoretical Half Precision (FP16) Performance - 10.4 PFLOPs Peak theoretical Eight-bit Precision (FP8) Performance - 20.88 PFLOPs. MI300X 8xGPform: Peak theoretical Half Precision (FP16) Performance - 10.4 PFLOPs. Actual performance will vary based on final specifications and system configuration.

MI355-005: Calculations conducted by AMD Performance Labs as of October 2nd, 2024 for the AMD Instinct™ MI300X GPU  accelerator, AMD Instinct™ MI325X GPU  accelerator and AMD Instinct™ MI350X GPU accelerator performance comparing FP16, FP8 and FP4 datatypes. MI300X GPU Accelerator Peak theoretical Half Precision (FP16)  Performance - 1.3 PFLOPs Peak theoretical Eight-bit Precision (FP8) Performance - 2.61 PFLOPs. MI325X GPU Accelerator. Peak theoretical Half Precision (FP16) Performance - 1.3 PFLOPs. Peak theoretical Eight-bit Precision (FP8) Performance - 2.61 PFLOPs. MI355X GPU Accelerator: Peak theoretical Half Precision (FP16) Performance - 2.3 PFLOPs. Peak theoretical Eight-bit Precision (FP8) Performance - 4.614 PFLOPs. Peak theoretical Six-bit Precision (FP6) Performance – 9.227 PFLOPs.Peak theoretical Four-bit Precision (FP4) Performance - 9.227 PFLOPs. Actual performance will vary based on final specifications and system configuration

GD-173a: AMD defines "All Day Battery Life" as at least 8 hours of continuous battery life and "Multi-Day battery Life" as continuous runtime above 8 hours. All battery life scores are approximate. Actual battery life will vary based on several factors, including, but not limited to: system configuration and software, settings, product use and age, and operating conditions.

GD-243: Trillions of Operations per Second (TOPS) for an AMD Ryzen processor is the maximum number of operations per second that can be executed in an optimal scenario and may not be typical. TOPS may vary based on several factors, including the specific system configuration, AI model, and software version.

STXP-04: Based on product specifications and competitive products announced as of Oct 2024 and testing as of Sept 2024 by AMD performance labs using the following systems: HP EliteBook X G1a with AMD Ryzen AI 9 HX PRO 375 processor @23W, Radeon 880M graphics, 32GB of RAM, 512GB SSD, VBS=ON, Windows 11 PRO; Dell Latitude 7450 with Intel Core Ultra 7 165U processor @15W (vPro enabled), Intel Iris Xe Graphics, VBS=ON, 32GB RAM, 512GB NVMe SSD, Microsoft Windows 11 Professional; Dell Latitude 7450 with Intel Core Ultra 7 165H processor @28W (vPro enabled), Intel Iris Xe Graphics, VBS=ON, 16GB RAM, 512GB NVMe SSD, Microsoft Windows 11 Pro.  All systems were tested in Best Performance Mode. AI PC is defined as a laptop PC with a processor that includes a neural processing unit (NPU).

STXP-05: Based on Microsoft Copilot+ requirements of minimum 40 TOPS using AMD product specifications and competitive products announced as of Oct 2024. Microsoft requirements found here - https://support.microsoft.com/en-us/topic/copilot-pc-hardware-requirements-35782169-6eab-4d63-a5c5-c498c3037364.

STXP-12: Testing as of Sept 2024 by AMD performance labs on an HP EliteBook X G1a (14in) (40W) with AMD Ryzen AI 9 HX PRO 375 processor, Radeon™ 890M graphics, 32GB of RAM, 512GB SSD, VBS=ON, Windows 11 Pro vs. a Dell Latitude 7450 with an Intel Core Ultra 7 165H processor (vPro enabled), Intel Arc Graphics, VBS=ON, 16GB RAM, 512GB NVMe SSD, Microsoft Windows 11 Pro in the application(s) (Best Performance Mode): Cinebench R24 nT. Laptop manufactures may vary configurations yielding different results. STXP-12.

# Endnotes

# DISCLAIMER AND ATTRIBUTIONS

DISCLAIMER

The information contained herein is for informational purposes only, and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions and typographical errors, and AMD is under no obligation to update or otherwise correct this information.  Advanced Micro Devices, Inc. makes no representations or warranties with respect to the accuracy or completeness of the contents of this document, and assumes no liability of any kind, including the implied warranties of noninfringement, merchantability or fitness for particular purposes, with respect to the operation or use of AMD hardware, software or other products described herein. No license, including implied or arising by estoppel, to any intellectual property rights is granted by this document. Terms and limitations applicable to the purchase or use of AMD's products are as set forth in a signed agreement between the parties or in AMD's Standard Terms and Conditions of Sale. GD-18

© 2024 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, EPYC, Instinct, Pensando, Radeon, ROCm, Ryzen  and combinations thereof are trademarks of Advanced Micro Devices, Inc.  National laboratory names and logos are registered trademarks of the U.S. Department of Energy. Use of these marks does not constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof or its contractors or subcontractors. Other product names used in this publication are for identification purposes only and may be trademarks of their respective owners. Certain AMD technologies may require third-party enablement or activation. Supported features may vary by operating system. Please confirm with the system manufacturer for specific features. No technology or product can be completely secure.