# 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.
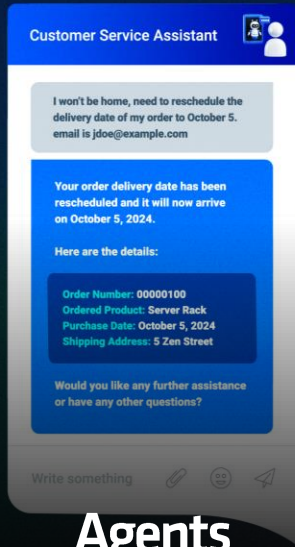
# 2024 AMD AI SOLUTIONS DAY
## AI 無限進化 . AMD 驅動未來

**AMD 台灣區商用市場業務處**
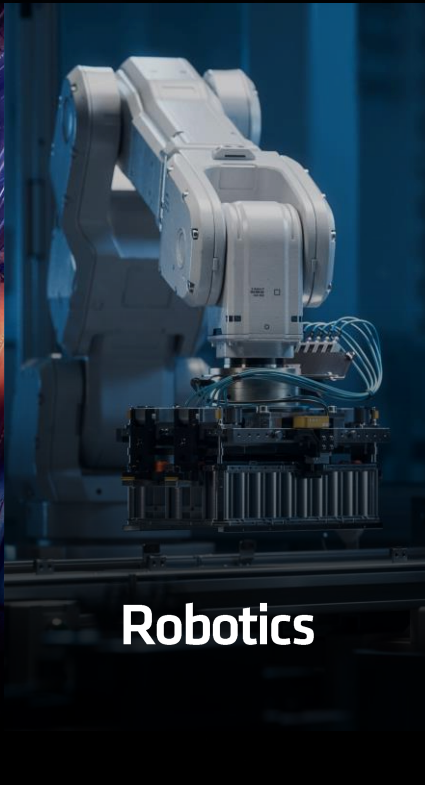**資深業務副總經理**

林建誠 Ken Lin

**AMD**
together we advance_

# AI

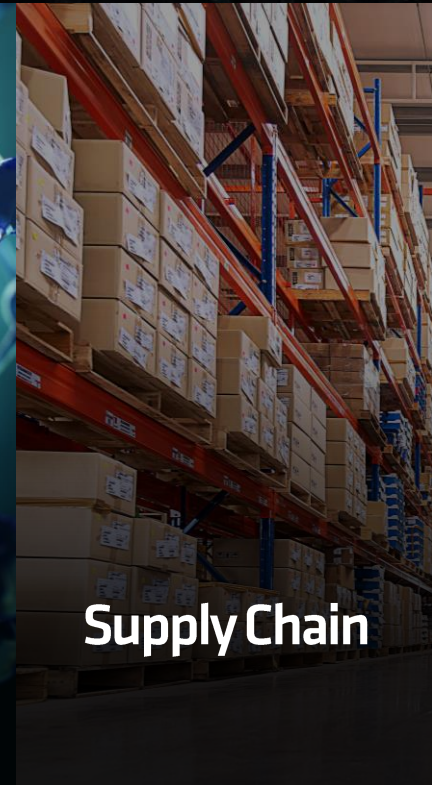## Most transformational technology in 50 years



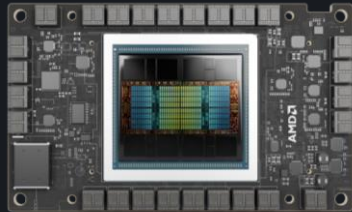**Agents** | **Smarter Cities** | **Robotics** | **Healthcare** | **Research** | **Supply Chain**

# Broad portfolio to address diverse spectrum of requirements

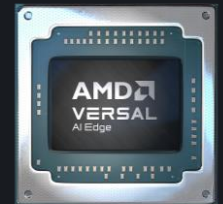| AMD EPYC™ Processors | AMD Instinct™ Accelerators | AMD Radeon™ Graphics | AMD Ryzen™ Mobile Processors | AMD Versal™ Adaptive SoCs |
|---|---|---|---|---|

‹ **From cloud to endpoint to embedded** ›

# AMD

## Advancing end-to-end AI infrastructure

| Cloud | Endpoint | Edge |

# Evolving Complexity: 250KW Racks, 100K node clusters

**Silicon** > **Server** > **Rack** > **Data Center**

**3x** chip power
**15x** rack power/density

**32x** GPU POD size
**100x** cluster size

**10-100x** datacenter power
**16x** networking bandwidth

AMD internal estimates

# Advancing Data Center Solutions

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

# AMD EPYC™ record market share…and growing

2018 — 2% — 1st Gen Processor Family
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**

# #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

📞    f    📷    NETFLIX    Office 365    ORACLE EXADATA    salesforce    SAP    Tji    Uber    💬    zoom

# Trusted by industry leaders on-prem

Adobe · amazon · Saudi Aramco أرامكو السعودية · ARISTA · BASF · BEST BUY · BNP PARIBAS · BNY MELLON · casa systems · CGG Passion for Geoscience

CITADEL | Securities · core42 A G42 company · DBS · DXC TECHNOLOGY · Emirates NBD · eni · f5 · GE · HERSHEY'S · Honeywell

HYUNDAI · IBM · Jane Street · mastercard · Medtronic · NETFLIX · NOKIA · NORTHROP GRUMMAN · PETRONAS · QTS

Qubit. · Raytheon · Reliance · RBC · Shell · SILICON LABS · ST · SUBARU · SYNOPSYS · TATA

TESLA · TOPCON Healthcare SEEING EYE HEALTH DIFFERENTLY · tsmc · TURTLE ROCK STUDIOS · UD TRUCKS · UNION PACIFIC · WELLS FARGO · weta DIGITAL

# 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

# Industry's Highest Performing Server CPU

**AMD EPYC™**
**5th Gen 9965**

192 cores — 2.7

**AMD EPYC™**
**4th Gen 9754**

128 cores — 1.7

**Intel™ Xeon®**
**5th Gen 8592+**

64 cores — 1.0

SPECrate®_2017_int_base

# 2.7x

vs. top-of-stack
"Emerald Rapids"

# Leadership Engines for Enterprise AI Workloads

| | Data Input | Data Cleaning | Pre-Processing | Model Training | Deployment | |
|---|---|---|---|---|---|---|
| **AMD EPYC** | | | | | | **AMD INSTINCT** |

**x86 CPU**

**GPU**

**Data from x86-Based Systems**

**Agentic AI**

**LLM Training and Inference**

From **analytics** to **generative AI** to **agentic AI**

# Easily Upgrade to 5ᵗʰ Gen AMD EPYC™ CPUs

Modernize your data center – Add more capacity for your compute needs

| 1000 Old Servers | 131 Modern Servers |
|---|---|
| 2P Intel® Xeon® Platinum 8280 servers | 2P AMD EPYC™ 9965 servers |

**7 to 1**

Easy to migrate to AMD
- X86 architecture
- Mature ecosystem
- Robust tools

Up to **68%**
Less power

Up to **87%**
Fewer Servers

Up to **67%**
Lower 3-yr TCO

Servers required to achieve a total of 391,000 SPECrate®2017_ int_base performance score

See endnotes 9xx5TCO-002A

15

# End-to-End AI and Inference Performance



**Machine Learning**
XGBoost (Higgs)

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

| 5th Gen Intel® | 4th Gen AMD | 5th Gen AMD |
|---|---|---|
| Xeon® 8592+ 64C | EPYC™ 9654 96C | EPYC™ 9965 192C |

up to **3.8x**

AI performance on CPU

# 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

# Advancing Data Center Solutions

**Data Center CPUs**

**Data Center GPUs**

**Networking**

AMD EPYC

AMD INSTINCT

AMD PENSANDO

# 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

Microsoft · ORACLE · DELL Technologies · Hewlett Packard Enterprise · Lenovo · SUPERMICRO

EVIDEN an atos business · CISCO · ASUS · GIGABYTE · TYAN · PENGUIN COMPUTING · AiVRES · ASRock Rack

Aligned · Elio · MiTAC · QCT · NSCALE · VULTR · Crusoe

Cirrascale CLOUD SERVICES · TENSORWAVE · COMPAL · Inventec · ingrasys · wiwynn · wistron

**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

**AMD Instinct™**
# MI325X Platform

See endnotes MI325-001A, MI325-002

*Dense flops

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

1x GPU

8x GPU

~1.1x

~1x

| H200 | MI325X |

| H200 HGX | MI325X Instinct™ Platform |

**Meta Llama-2**
**7B**

**Meta Llama-2**
**70B**

Nvidia
**H200 HGX**

AMD Instinct™
**MI325X Platform**

See endnotes MI325-013, MI325-012.

# AMD Instinct™ MI325X GPU

## Production starting in **Q4 2024**

DELL Technologies  EVIDEN an atos business  GIGABYTE™  Hewlett Packard Enterprise  Lenovo.  SUPERMICRO

AiVRES  ASRock Rack  ASUS®  Celestica™  ingrasys®  Inventec  MiTAC  QCT  wistron  wiwynn®

## Available from leading system and infrastructure solution partners starting Q1 2025

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

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

Planned availability **2H 2025**

# AMD Instinct™ MI300X Accelerator

## Performant out-of-box support on popular generative AI models

**1M+**
models supported
out of the box

🤗 Hugging Face

**Extended support
for leading models**

---

∞ Meta
Llama 3.1

∞ Meta
Llama 3.2

Stable
Diffusion
3

**Day 0 support
for AMD GPUs**

---

Llama 3 405B
latency
improvement

MI300X vs. H100

**Leadership performance
on popular models**

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

~**2.4x**
average performance improvement

~**1.9x** Mixtral 8x22B

~**2.1x** Mixtral 8x7B

~**2.6x** Qwen2 72B

~**2.6x** Llama3.1 70B

~**2.8x** Llama3.1 8B

Runtime Optimization | Kernel Fusion | Collective Communication | Subgraph

# Advancing Data Center Solutions

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

# Ultra Ethernet
### Consortium

## Steering Members

AMD

ARISTA

BROADCOM

CISCO

EVIDEN
an atos business

Hewlett Packard
Enterprise

intel

Meta

Microsoft

ORACLE

# Ethernet is always the preferred choice

**> 50%**
TCO Saving

1,000,000+ GPU

up to
**48,000** GPU

InfiniBand

Ethernet
RoCEv2

InfiniBand

Ethernet
RoCEv2

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

**Scalability**
Higher is better

Sources: 1) 650Group Datacenter AI Networking and Server SmartNIC Forecast Reports 2Q24 .

# AMD

## Advancing end-to-end AI infrastructure

| Cloud | Endpoint | Edge |

# AMD Ryzen™ AI Leads the AI PC Era

**Q1 2023**

AMD
RYZEN AI

1st Gen
"Phoenix Point"

**10**
TOPS NPU

**Q4 2023**

AMD
RYZEN AI

2nd Gen
"Hawk Point"

**16**
TOPS NPU

**Q2 2024**

AMD
RYZEN AI

3rd Gen
"Strix Point"

**50+**
TOPS NPU

See endnote GD-243.

AMD Ryzen™ AI PRO 300 Series

First Copilot+ laptops enabled for enterprise PCs

ZEN 5 | AMD XDNA 2 | AMD RDNA 3.5 | Copilot+PC

See endnote STXP-05

# Enterprise AI PC Application Ecosystem

Adobe          Microsoft          webex          zoom
                                   by CISCO

---

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

Blackmagicdesign   nero   BORIS FX   Avid   Rhinoceros   GoPro   Topaz Labs™   ARKRUNR   voicemy.ai
                                                         Be a HERO.

blender   RADICAL   CyberLink   AFFINITY   ACCA   convai   OBS   MAXON
          AI-POWERED 3D ANIMATION   Photo 2   ACCA SOFTWARE         Open Broadcaster Software

# AMD

## Advancing end-to-end AI infrastructure

| Cloud | Endpoint | Edge |

# General Disclaimer and Attribution Statement 2024

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 products are as set forth in a signed agreement between the parties or in AMD's Standard Terms and Conditions of Sale. GD-18u.

# Endnotes

9xx5-001: Based on AMD internal testing as of 9/10/2024, geomean performance improvement (IPC) at fixed-frequency. 5th Gen EPYC CPU Enterprise and Cloud Server Workloads generational IPC Uplift of 1.170x  (geomean) using a select set of 36 workloads and is the geomean of estimated scores for total and all subsets of SPECrate®2017_int_base (geomean ), estimated scores for total and all subsets of SPECrate®2017_fp_base (geomean), scores for Server Side Java multi instance max ops/sec,   representative Cloud Server workloads (geomean), and representative Enterprise server workloads (geomean). "Genoa" Config (all NPS1): EPYC 9654 BIOS TQZ1005D 12c12t (1c1t/CCD in 12+1), FF 3GHz, 12x DDR5-4800 (2Rx4 64GB), 32Gbps xGMI; "Turin" config (all NPS1): EPYC 9V45    BIOS RVOT1000F 12c12t (1c1t/CCD in 12+1), FF 3GHz, 12x DDR5-6000 (2Rx4 64GB), 32Gbps xGMI Utilizing Performance Determinism and the Performance governor on Ubuntu® 22.04 w/ 6.8.0-40-generic kernel OS for all workloads. 5th Gen EPYC generational ML/HPC Server Workloads IPC Uplift of 1.369x (geomean) using a select set of 24 workloads and is the geomean of representative ML Server Workloads (geomean), and representative HPC Server Workloads (geomean). "Genoa Config (all NPS1) "Genoa" config: EPYC 9654 BIOS TQZ1005D 12c12t (1c1t/CCD in 12+1), FF 3GHz, 12x DDR5-4800 (2Rx4 64GB), 32Gbps xGMI; "Turin" config (all NPS1):   EPYC 9V45 BIOS RVOT1000F 12c12t (1c1t/CCD in 12+1), FF 3GHz, 12x DDR5-6000 (2Rx4 64GB), 32Gbps xGMI. Utilizing Performance Turrance Determinism and the Performance governor on Ubuntu 22.04 w/ 6.8.0-40-generic kernel OS for all workloads except LAMMPS, HPCG, NAMD, OpenFOAM, Gromacs  which utilize 24.04 w/ 6.8.0-40-generic kernel. SPEC® and SPECrate® are registered trademarks for Standard Performance Evaluation Corporation. Learn more at spec.org.

9xx5-002C: SPECrate®2017_int_base comparison based on published scores from www.spec.org as of 10/10/2024. 2P AMD EPYC 9965 (3000 SPECrate®2017_int_base, 384 Total Cores, 500W TDP, $14,813 CPU $), 6.060 SPECrate®2017_int_base/CPU W, 0.205 SPECrate®2017_int_base/CPU $, https://www.spec.org/cpu2017/results/res2024q3/cpu2017-20240923-44833.html). 2P AMD EPYC 9755 (2720 SPECrate®2017_int_base, 256 Total Cores, 500W TDP, $12,984 CPU $), 5.440 SPECrate®2017_int_base/CPU W, 0.209 SPECrate®2017_int_base/CPU $, https://www.spec.org/cpu2017/results/res2024q4/cpu2017-20240923-44837.pdf). 2P AMD EPYC 9754 (1950 SPECrate®2017_int_base, 256 Total Cores, 360W TDP, $11,900 CPU $), 5.417 SPECrate®2017_int_base/CPU W, 0.164 SPECrate®2017_int_base/CPU $, https://www.spec.org/cpu2017/results/res2023q2/cpu2017-20230522-36617.html). 2P AMD EPYC 9654 (1810 SPECrate®2017_int_base, 192 Total Cores, 360W TDP, $11,805 CPU $), 5.028 SPECrate®2017_int_base/CPU W, 0.153 SPECrate®2017_int_base/CPU $, https://www.spec.org/cpu2017/results/res2024q1/cpu2017-20240129-40896.html). 2P Intel Xeon Platinum 8592+ (1130 SPECrate®2017_int_base, 128 Total Cores, 350W TDP, $11,600 CPU $) 3.229 SPECrate®2017_int_base/CPU W, 0.097 SPECrate®2017_int_base/CPU $, http://spec.org/cpu2017/results/res2023q4/cpu2017-20231127-40064.html). 2P Intel Xeon 6780E (1410 SPECrate®2017_int_base, 288 Total Cores, 330W TDP) 4.273 SPECrate®2017_int_base/CPU W, 0.124 SPECrate®2017_int_base/CPU $, https://spec.org/cpu2017/results/res2024q3/cpu2017-20240811-44406.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/.

9xx5-012: TPCxAI @SF30 Multi-Instance 32C Instance Size throughput results based on AMD internal testing as of 09/05/2024 running multiple VM instances. The aggregate end-to-end AI throughput test is derived from the TPCx-AI benchmark and as such is not comparable to published TPCx-AI results, as the end-to-end AI throughput test results do not comply with the TPCx-AI Specification. 2P AMD EPYC 9965 (384 Total Cores), 12 32C instances, NPS1, 1.5TB 24x64GB DDR5-6400 (at 6000 MT/s), 1DPC, 1.0 Gbps NetXtreme BCM5720 Gigabit Ethernet PCIe, 3.5 TB Samsung MZWLO3T8HCLS-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -l 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVOT1000C (SMT=off, Determinism=Power, Turbo Boost=Enabled). 2P AMD EPYC 9755 (256 Total Cores), 8 32C instances, NPS1, 1.5TB 24x64GB DDR5-6400 (at 6000 MT/s), 1DPC, 1.0 Gbps NetXtreme BCM5720 Gigabit Ethernet PCIe, 3.5 TB Samsung MZWLO3T8HCLS-00A07 NVMe®, Ubuntu 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -l 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVOT0090F (SMT=off, Determinism=Power, Turbo Boost=Enabled) 2P AMD EPYC 9654 (192 Total cores) 6 32C instances, NPS1, 1.5TB 24x64GB DDR5-4800, 1DPC, 2 x 1.92 TB Samsung MZQL21T9HCJR-00A07 NVMe, Ubuntu 22.04.3 LTS, BIOS 1006C (SMT=off, Determinism=Power) Versus 2P Xeon Platinum 8592+ (128 Total Cores), 4 32C instances, AMX On, 1TB 16x64GB DDR5-5600, 1DPC, 1.0 Gbps NetXtreme BCM5719 Gigabit Ethernet PCIe, 3.84 TB KIOXIA KCMYXRUG3T84 NVMe, , Ubuntu 22.04.4 LTS, 6.5.0-35 generic (tuned-adm profile throughput-performance, ulimit -l 132065548, ulimit -n 1024, ulimit -s 8192), BIOS ESE122V (SMT=off, Determinism=Power, Turbo Boost = Enabled) Results. CPU Median Relative Generational. Turin 192C, 12 Inst 6067.531 3.775 2.278. Turin 128C, 8 Inst 4091.85 2.546 1.536 Genoa 96C, 6 Inst 2663.14 1.657 1. EMR 64C, 4 Inst 1607.417 1 NA. Results may vary due to factors including system configurations, software versions and BIOS settings. TPC, TPC Benchmark and TPC-C are trademarks of the Transaction Processing Performance Council.

.

# 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

# Endnotes

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 may vary. 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/nb5zzsjdf/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 GPU memory 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*R0NMLjE3MTM5NjQ3NTAuQ2p3S0NBancyNkt4QmhCREVppd0F1NktYdDlweXY1dlUtaHNKNmhPdHM4UVdPSlM3dFdQaE40WkI4THZBaWFVajFyTGhYd3hLQmlZQ3pCb0NsVElRQXZEX0J3RQ..*_gcl_au*MTIwNjg4NjU0Ny4xNzExMDM1NTQ3

MI325-012: Overall GPU-normalized Training Throughput (processed tokens per second) for text generation using the Llama2-7b chat model running Megatron-LM v0.12 (BF16) when using a maximum sequence length of 4096 tokens comparison based on AMD internal testing as of 10/4/2024. Batch size according to largest micro-batch that fits in GPU memory for each system. AMD Instinct batch size 8, Nvidia batch size 2. Configurations: AMD Development system: 1P AMD Ryzen 9 7950X (16-core), 1x AMD Instinct™ MI325X (256GB, 1000W) GPU, 128 GiB memory, ROCm 6.3.0 (pre-release), Ubuntu 22.04.2 LTS with Linux kernel 5.15.0-72-generic, PyTorch 2.4.0.Vs. An Nvidia DGX H200 with 2x Intel Xeon Platinum 8468 Processors, 1x Nvidia H200 (141GB, 700W) GPUs, 2 TiB (32 DIMMs, 64 GiB/DIMM), CUDA 12.6.37-1, 560.35.03, Ubuntu 22.04.5, PyTorch 2.5.0a0+872d972e41.nv24.8. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations. MI325-012

MI325-013: Overall GPU-normalized Training Throughput (processed tokens per second) for text generation using the Llama2-70b chat model running Megatron-LM v0.12 (BF16) when using a maximum sequence length of 4096 tokens comparison based on AMD internal testing as of 10/7/2024. Batch size according to largest micro-batch that fits in GPU memory for each system. AMD Instinct global batch size 32, Nvidia global batch size 49. Configurations:2P Intel Xeon Platinum 8480+ CPU server with 8x AMD Instinct™ MI325X (256GB, 1000W) GPUs, 4 TiB memory, ROCm® 6.3 (Pre-release), Ubuntu® 22.04.2, PyTorch 2.4.0. An Nvidia DGX H200 with 2x Intel Xeon Platinum 8468 Processors, 8x Nvidia H200 (141GB, 700W) GPUs, 2 TiB memory, CUDA 12.6.37-1, 560.35.03, Ubuntu 22.04.5, PyTorch 2.5.0a0+872d972e41.nv24.8.. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations. MI325-013

# Endnotes

MI325-002: 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

AMD
together we advance_

# ENDNOTES

- MI325-002 -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/. MI325-002

- MI300-55: Inference performance projections as of May 31, 2024 using engineering estimates based on the design of a future AMD CDNA 4-based Instinct MI350 Series accelerator as proxy for projected AMD CDNA™ 4 performance. A 1.8T GPT MoE model was evaluated assuming a token-to-token latency = 70ms real time, first token latency = 5s, input sequence length = 8k, output sequence length = 256, assuming a 4x 8-mode MI350 series proxy (CDNA4) vs. 8x MI300X per GPU performance comparison. Actual performance will vary based on factors including but not limited to final specifications of production silicon, system configuration and inference model and size used.

- MI325-001A: Calculations conducted by AMD Performance Labs as of September 26th, 2024, based on current specifications and /or estimation. The AMD Instinct™ MI325X OAMaccelerator will have 256GB HBM3e memory capacity and 6 TB/s GPU peak theoretical memory bandwidth performance. Actual results based on production silicon may vary. 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 GPU memory 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*R0NMLjE3MTM5NjQ3NTAuQ2p3S0NBancyNkt4QmhCREVpd0F1NktYdDlweXY1dlUtaHNKNmhPdHM4UVdPSlM3dFdQaE40WkI4THZaBaWFVajFyTGhYd3hLQmlZQ3pCb0NsVElRQXZEX0J3RQ..*_gcl_au*MTIwNjg4NjU0Ny4xNzExMDM1NTQ3 . MI325-001A

- 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 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. Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations. MI325-004

- MI355-003. Calculations conducted by internal AMD Performance Labs as of September 26, 2024 on current specifications and/or internal engineering calculations. Comparison of the AMD Instinct™ MI355X platform (2.3 TB HBM3e) vs AMD Instinct™ MI300X platform (1.5 TB HBM3) on Large Language Model (LLM) to determine the maximum model size supported  Max Model Size Supported on a Platform

  FP16 Datatype
  MI300X Platform: 1.5 TB HBM3 (in FP16 datatype)
  [(Max Model Size in Billion Parameters) * 2.1 = 1.5 TB*1000
  Max Model Size = ~715 billion Parameters
  FP4 Datatype
  MI355X Platform: 2.3 TB HBM3 (in FP4 datatype)
  [(Max Model Size in Billion Parameters) * 0.55 = 2.3 TB*1000
  Max Model Size = ~4181 billion parameters/4.2 trillion parameters
  Assumptions:
  Batch size 1
  Memory needs for model = 2 Bytes per Parameter
  Memory size needs for activations and others = +10%

  Actual maximum LLM parameter size that can run on each platform may vary upon performance testing with physical servers. Calculations rely on published and sometimes preliminary model memory sizes. Model size results estimated on MI355X, MI325X and MI300X platforms due to system/part availability. Actual performance will vary based on final specifications and system configuration. MI355-003

- MI355-004: Calculations conducted by AMD Performance Labs as of September 26th, 2024 for the AMD Instinct™ MI300X GPU platform and accelerator and AMD Instinct™ MI300X GPU platform and accelerator performance comparing FP16 and FP4 datatypes.

  MI355X 8xGPU Platform
  Peak theoretical Four-bit Precision (FP4) Performance - 74 PFLOPs

  MI300X 8xGPU Platform
  Peak theoretical Half Precision (FP16) Performance - 10.4 PFLOPs

  Actual performance will vary based on final specifications and system configuration. MI355-004

  MI325-015: Based on testing completed on 10/08/2024 by AMD performance lab measuring overall latency for text generated using LLaMA 3.1-405B model using FP8 datatype. Test was performed using input length of 2048 tokens and an output length of 2048 tokens with a batch size of 32 for the following configurations of AMD Instinct™ MI325X 8xGPU platform and NVIDIA H200 HGX GPU platform. Configurations: AMD Instinct™ MI325X platform: Dell PowerEdge XE9680 with 2x Intel Xeon Platinum 8480+ Processors, 8x AMD Instinct MI325X (256GiB, 1000W) GPUs, Ubuntu 22.04, and a pre-release build of ROCm 6.3 NVIDIA H200 HGX platform: Supermicro SuperServer with 2x Intel Xeon Platinum 8468 Processors, 8x Nvidia H200 (140GB, 700W) GPUs, Ubuntu 22.04) 8x MI325X platform with vLLM performance: 23.033 seconds (latency in seconds) Vs. 8x H200 DGX platform with TensorRT-LLM: 27.743 seconds (latency in seconds) Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations. MI325-015

**AMD**
together we advance_