# AMD INSTINCT™ Accelerator

Compute & Accelerator Forum June 2024

**AMD**
together we advance_

# Agenda

AMD
together we advance_

**Powers the daily lives of billions**
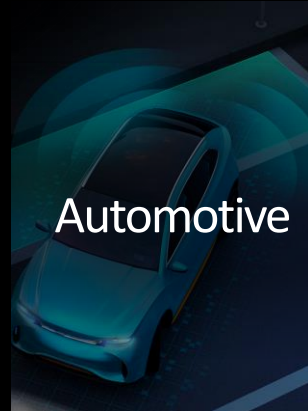
Connectivity   Healthcare   Industrial   Automotive   Cloud   PCs   Gaming   AI

| | APRIL 2024

AMD
together we advance_

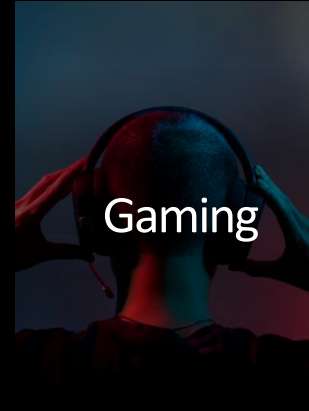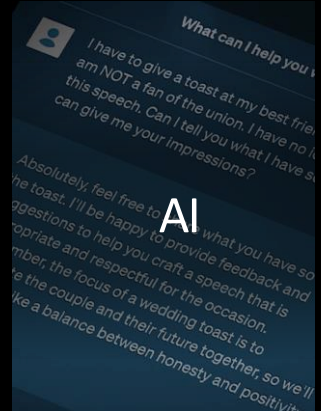**AMD**

# Advancing end-to-end AI infrastructure

Cloud          HPC          Enterprise          Embedded          PC

**AMD**
AI Platforms

Broad portfolio of training and inference compute engines

Open and proven software capabilities

AI ecosystem with deep co-innovation

**AMD**
together we advance_

# Only AMD powers the full range of data center workloads

Traditional IT workloads

Max performance and efficiency

**AMD EPYC**

CPU

Mix of AI and traditional workloads

Smaller model sizes

**AMD INSTINCT**

GPU

AI at scale and dedicated infrastructure

Larger model sizes

General Purpose

AI Inference

AI Training

# AMD Product Portfolio
# from cloud to client



**AMD Instinct GPU Accelerators**

Data center HPC and AI solutions

**4th Gen AMD EPYC**

Industry-leading x86 CPU server solution

**Embedded Versal and Alveo**

AI + sensor fusion for embedded, FPGA

**Radeon GPU**

GPU for AI in gaming and AI developers

**Ryzen Mobile Processors**

x86 with integrated GPU and Ryzen AI accelerator

AMD
together we advance_

# AMD
## Instinct™ Accelerators

Data center GPU for the most demanding AI and HPC workloads

# AMD Instinct Strategic Pillars

## Enabling customer success

**Ease of Migration**

Drop-in compatible with existing infrastructure for hardware and software

**Performance Leadership**

Leading performance without compromise

**Commitment to Openness**

Investment and participation in open standards across the entire ecosystem

**Customer Focused**

Roadmap and support structure geared towards customer success

AMD
together we advance_

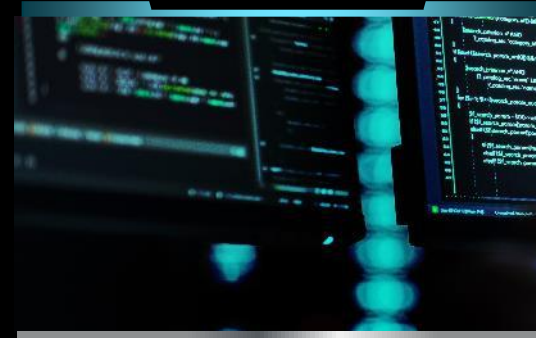# The AMD Instinct™ Accelerator Journey

## Multiple generations of architecture focused advancing HPC & AI compute

### MI100
AMD CDNA™

**ECOSYSTEM GROWTH**

First purpose-built GPU architecture to accelerate FP64 and FP32 HPC workloads

### MI200
AMD CDNA™ 2

**DRIVING HPC AND AI TO A NEW FRONTIER**

Denser compute architecture with leading memory capacity/bandwidth

### MI300
AMD CDNA™ 3

**DATA CENTER APU & DISCRETE GPU**

Focused improvements on Unified memory, AI data format performance and in-node networking

**2020** ──────────────────────────── **2023**

# Model Evolution Accelerating Rapidly

## AI performance needs driving GPU demand &cluster growth

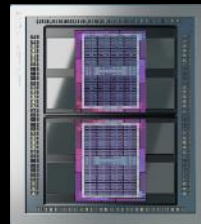- The pace of compute intensive model releases is accelerating with frontier models advancing rapidly
- Majority of the compute intensive models are LLM but newer multi modal and other domain models are emerging
- In 2020, only two models were trained with more than $10^{23}$ FLOP. This increased exponentially over the subsequent three years, and over 40 models trained at this scale were released in 2023

AMD
together we advance_

# AMD Instinct™ MI300X GPU vs. Competition

| | | MI300X (Up to) | H100 SXM | AMD Instinct™ Advantage (Up to) |
|---|---|---|---|---|
| Hardware Specifications | TBP | 750W | 700W | - |
| | Memory Capacity | 192 GB HBM3 | 80GB HBM3 | 2.4x |
| | Memory Bandwidth (Peak Theoretical) | 5.3 TB/s | 3.3TB/s | 1.6x |
| HPC Performance (Peak Theoretical) | FP64 Matrix \| Vector (TFLOPS) | 163.4 \| 81.7 | 66.9 \| 33.5 | 2.4x \| 2.4x |
| | FP32 Matrix \| Vector (TFLOPS)* | 163.4 \| 163.4 | N/A \| 66.9 | N/A \| 2.4x |
| AI Performance (Peak Theoretical) | TF32 (TFLOPS) | 653.7 | 494.7 | 1.3x |
| | FP16 (TFLOPS) | 1307.4 | 989.4 | 1.3x |
| | BFLOAT16 (TFLOPS) | 1307.4 | 989.4 | 1.3x |
| | FP8 (TFLOPS) | 2614.9 | 1978.9 | 1.3x |
| | INT8 (TFLOPS) | 2614.9 | 1978.9 | 1.3x |

See endnotes: MI300-05A, MI300-17, MI300-18
- Nvidia H100 GPUs don't support FP32 Tensor.
- Nvidia H100 source: https://resources.nvidia.com/en-us-tensor-core/

|

AMD
together we advance_

# AMD Instinct™ Platform

## 8x MI300X in a ready to deploy OCP form factor

| | | |
|---|---|---|
| **8x**<br>MI300X | **21** PF<br>BF16 \| FP16 | **1.5** TB/s<br>HBM3 |
| **896** GB/s<br>Infinity Fabric™ Bandwidth | Industry-Standard<br>OCP Design | |

**AMD**
together we advance_

# AMD Instinct™ MI300X Platform

## Infrastructure performance

| AMD Instinct™ **MI300X Platform** | Nvidia **H100 HGX** | AMD Instinct™ **MI300X Advantage** |
|---|---|---|
| **1.5** TB<br>**HBM3 memory** | **640** GB<br>HBM3 memory | **2.4X**<br>**More memory** |
| **~10.4** PF<br>**FP16 / BF16 FLOPS** | **7.9** PF<br>FP16 / BF16 FLOPS | **~1.3X**<br>**More Compute** |
| **~896** GB/s<br>**Aggregate bi-directional bandwidth** | **900** GB/s<br>Aggregate bi-directional bandwidth | **Comparable** |
| **448** GB/s<br>**Single node ring bandwidth** | **450** GB/s<br>Single node ring bandwidth | **Comparable** |
| Up to **400** GbE<br>**NIC / GPU** | Up to **400** GbE<br>NIC / GPU | **Equivalent** |
| **PCIe® Gen 5**<br>**128 GB/s** | **PCIe® Gen 5**<br>128 GB/s | **Equivalent** |

See endnotes:MI300-25
Nvidia H100 source: https://resources.nvidia.com/en-us-tensor-core/

**AMD**
together we advance_

# AMD Instinct™ Platform: Performance Advantage



**Nvidia**
**1** H100 HGX

640 **GB** HBM3 | 26.4 **TB/s**

**AMD Instinct™**
**1** MI300X Platform

1.5 **TB** HBM3 | 42.4 **TB/s**

| Training & Inference | | Training | Inference |
|---|---|---|---|
| **1x** | Performance per system | **1x**<br>MPT-30B | **2.1x**<br>Llama 70B |
| **1x** | Models per system | **2x**<br>~30B | **2x**<br>~70B |
| **1x** | Max LLM model size per system | **2x**<br>~70B vs.~30B | **2x**<br>~680B vs 290B |

Results may vary. See endnotes:MI300-34, MI300-40, MI300-39, MI300-42

# Delivering Exceptional Value to AI leaders

**Microsoft Azure**

**MI300X enables to serve larger AI models with fewer GPUs**

"With MI300X's larger memory capacity and bandwidth, we can serve larger models with fewer GPUs. We have already got GPT-4 up and running on MI300X"

**Satya Nadella**
**CEO, Microsoft**
*November 2023*

**Meta**

**Ecosystem growth over the years has made ROCm a highly competitive software platform**

"We have had a great experience with ROCm and the performance it has been able to deliver with MI300X. The optimizations and the ecosystem growth over the years have made ROCm a highly competitive software platform. We see great performance numbers which we believe will benefit the industry"

**Ajit Mathews**
**Sr. Director, Meta**
*December 2023*

**databricks**

**ROCm runs out of the box from day one**

ROCm runs out of the box from day one. It was was very easy to run and include ROCm in our stack . Many of the generative AI tools today are open source like PyTorch, Triton, Huggingface and these tools can run today on AMD ROCm software stack and this makes ROCm another key component of the open source ecosystem

**Ion Stoica**
**Co-Founder and Executive Chairman, Databricks**
*December 2023*

Open | Proven | Ready

# Frameworks Support Status

## Key frameworks fully upstreamed and optimized for AMD Instinct™ Accelerators

**⊙ PyTorch**

- Full Feature Support on Day 0 since Pytorch 2.0

**TensorFlow**

- Upstream Tensorflow Version Optimized For AMD Instinct (2.13, 2.14)

**JAX**

- Upstreamed JAX version optimized for AMD Instinct
- JAX supported w/ OpenXLA & Triton backends

**OpenAI Triton**

- AMD is the "top" 3rd party hardware contributor to OpenAI Triton
- Upstreamed support for AMD Instinct
- FP8 datatype supported on MI300X
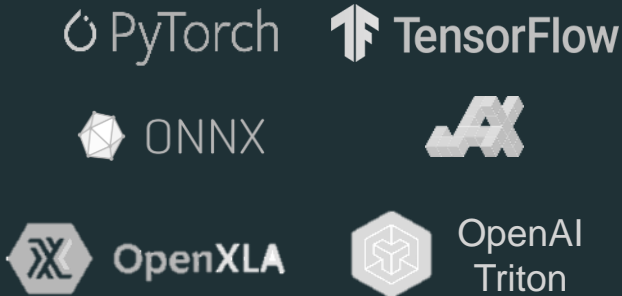- Available Now:  Docker pull rocm/oai-triton

**OpenXLA**

- OpenXLA project "founding member"
- AMD support functional (and upstreamed)
  - Focused on maintaining current AMD support for Tensorflow while code bases are being refactored
- Available Now:  https://github.Com/openxla/xla

**AMD**
together we advance_

# Transitioning Workloads to Instinct & ROCm
## Low friction softward porting for existing Nvidia users to AMD

| | |
|---|---|
| **DROP IN OUT-OF-THE-BOX SUPPORT**<br><br>For Existing Code | PyTorch  TensorFlow<br><br>ONNX  JAX<br><br>OpenXLA  OpenAI Triton |
| **PORT & OPTIMIZE**<br><br>For Custom Kernels | Leverage AMD HIPIFY tool for large custom kernels or code re-write if smaller number of lines of code (typical) |
| **EQUIVALENT LIBRARIES**<br><br>For New Code Dev | ROCm Libraries Developed to Mirror CUDA-based libraries<br>• rocBLAS, rocSparse, rocFFT, RCCL, MIOpen… |

**PURPOSELY DESIGNED TO LEVERAGE EXISTING CUSTOMER CODE WITH MINIMAL CHANGES**

- Vast majority of AI end users engaged by AMD are programing at the framework level and their code functions out of the box with no edits

- Performance optimizations for common models and customer driven asks underway to ensure out of the box performance

- Foundational model builders with custom CUDA kernels have the option to use AMD HIPIFY to convert CUDA code, but often find it to be a low lift to re-write that small portion of code for AMD GPUs

AMD
together we advance_

# Training: Case Studies

## Published AMD Instinct™ training runs

**Oak Ridge National Laboratory**
- 1T GPT model
- 3072 MI250s
- 87% strong scaling eff

**kisaled**
- 221B T5 based model
- 1200 MI250s
- Pre-tests outperformed A100

**UNIVERSITY OF TURKU**
- 13B Finnish model
- 768 MI250s
- Utilized Megatron DeepSpeed

**AI2**
- Olmo 7B (65B in progress)
- 1024 MI250s
- Utilized PyTorch FSDP

**SILOGEN**
- Poro 34B model
- 512 MI250s

**Microsoft**
- 6.7B RetNet
- 512 MI250s
- Reported "decent throughput"

**MW**
- MPT-1B,3B,7B
- 32 MI250
- Proved interoperability between AMD and NV GPUs

**LAMINI**
- Fine-tuning of open-source models
- Utilizes MI210
- Able to host 200B model in single server

AMD together we advance_

Announced last week

# Ultra Accelerator Link

Partner group of innovators for scale up AI infrastructure

AMD    BROADCOM    CISCO    Google    Hewlett Packard Enterprise    intel.    Meta    Microsoft

High Performance | Open | Scalable