



# FINANCIAL ANALYST DAY 2022

together we advance\_

## Driving GPU Leadership

---

**David Wang**

Senior Vice President, Engineering, Radeon Technologies Group

# Cautionary Statement

This presentation contains forward-looking statements concerning Advanced Micro Devices, Inc. (AMD) including, but not limited to, AMD's GPU momentum and strategy; AMD's gaming GPU and compute GPU architecture roadmaps; and the timing, availability, features, functionality and expected benefits of AMD's products, 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.

# AMD GPU MOMENTUM



## Supercomputer

---

Frontier  
LUMI  
Adastra



## Data Center

---

AWS  
Google Cloud  
Microsoft Azure



## PC

---

Radeon™ RX 6000 Series  
Radeon™ W6000 Series  
Ryzen™ 6000 Series



## Console

---

PlayStation 5  
Xbox Series X | S  
Steam Deck



## Embedded

---

Magic Leap  
Tesla



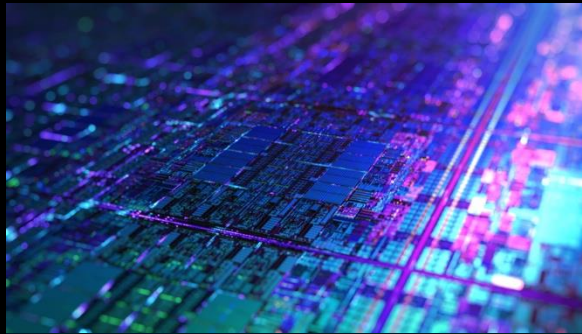
## Mobile

---

Samsung

# GPU TECHNOLOGY STRATEGY

DELIVERING PERFORMANCE AND EFFICIENCY LEADERSHIP



## Architecture

---

Domain-Specific  
Architecture Optimizations



## Technology

---

Advanced Process and  
Packaging Technologies



## Efficiency

---

Leadership Performance-  
Per-Watt Roadmap



## Ecosystem

---

Open-Source Software,  
Including AI



# AMD RDNA™ 2 ARCHITECTURE

Designed for enthusiast gaming

---

## — Innovation

High-Speed Compute Units, Raytracing Cores, AMD Infinity Cache™, DX12® Ultimate

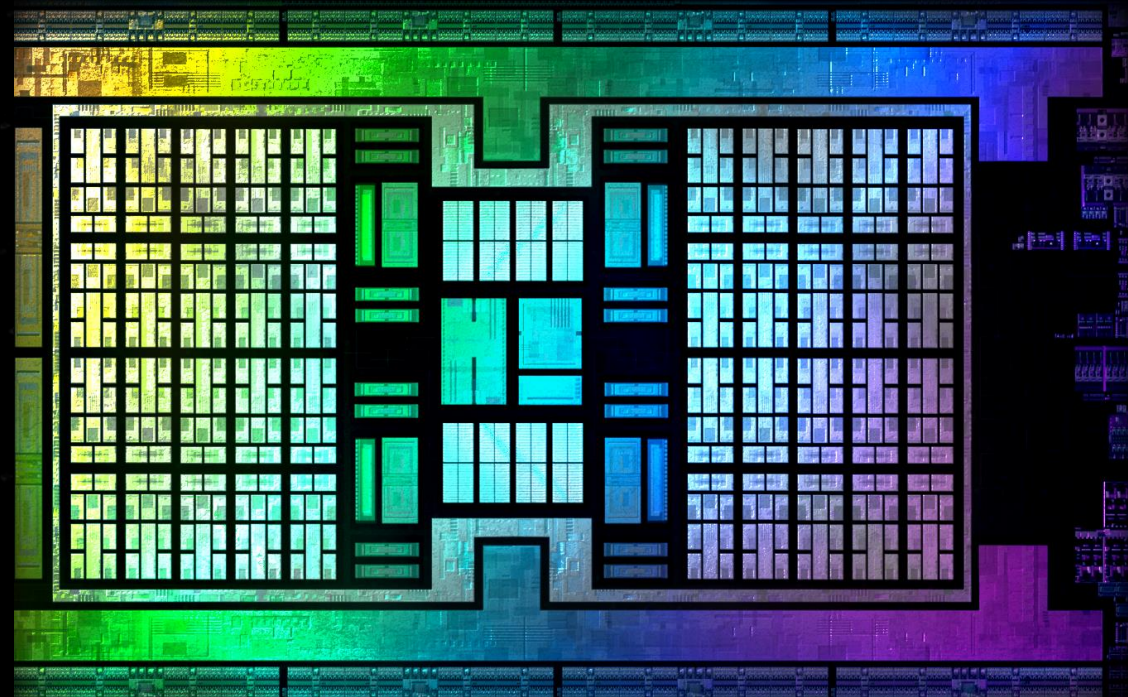
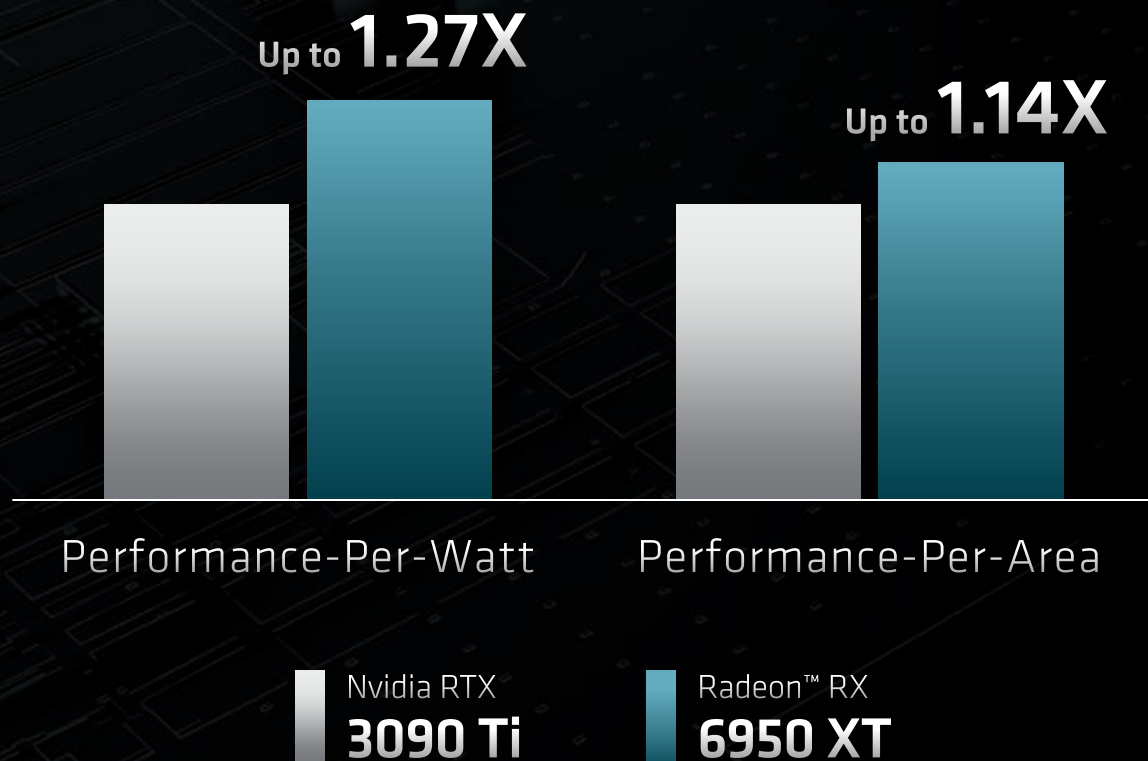
## — Performance

Delivered 2X Performance and More Than 50% Performance-Per-Watt Gains vs. AMD RDNA™

## — Scalability

Powers Mobile, Console, Ryzen™ APU, Radeon™ GPU, and Cloud

# AMD RDNA™ 2 LEADERSHIP GAMING PERFORMANCE AND EFFICIENCY



Average of 20 games, See endnote RX-785

**AMD**  
Software  
Adrenalin Edition

# GREAT HARDWARE NEEDS GREAT SOFTWARE



Community  
Engagement



Vigorous  
Testing



Regular  
Updates

**Focus on Quality  
and Stability**

**15%**  
Year-Over-Year Uplift

Day-0 Driver Support

**Continuous Performance  
Improvements**

AMD  
RADEON  
Anti-Lag

AMD  
RADEON  
Image Sharpening

AMD  
FreeSync

AMD  
FidelityFX

AMD  
Link

AMD  
RADEON  
Boost

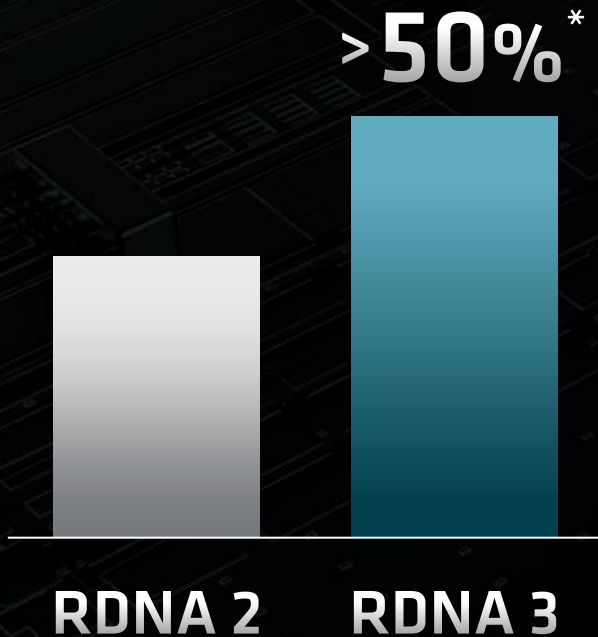
AMD  
RADEON  
Super Resolution

**Immersive  
Experiences**

# AMD RDNA 3

## THE JOURNEY CONTINUES

Projected Performance/Watt Uplift



Performance-per-watt uplift through:

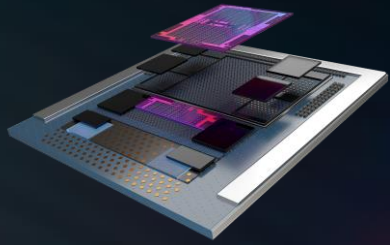
- 5nm Process
- Advanced Chiplet Packaging
- Rearchitected Compute Unit
- Optimized Graphics Pipeline
- Next-Gen AMD Infinity Cache™

Based on preliminary internal engineering estimates. Actual results subject to change.





# AT THE FOREFRONT OF GRAPHICS INNOVATION



## Advanced Chiplet Packaging

Leadership Performance and Scalability



## End-to-End Power Optimization

System Level Energy Efficiency



## Hybrid Rendering

Real-Time Immersive Experiences



## Next-Gen Multimedia

Enhanced Video and Display Capabilities



# AMD GAMING GPU ARCHITECTURE ROADMAP

## DRIVING PERFORMANCE AND EFFICIENCY LEADERSHIP



2019

2024

All roadmaps are subject to change.



# AMD CDNA™ 2 ARCHITECTURE

Exascale-class technology for HPC/AI

## — Innovation

High-Performance Dual Engines in MCM, 3<sup>rd</sup> Gen AMD Infinity Architecture, Ultra-Wide HBM Interface

## — Performance

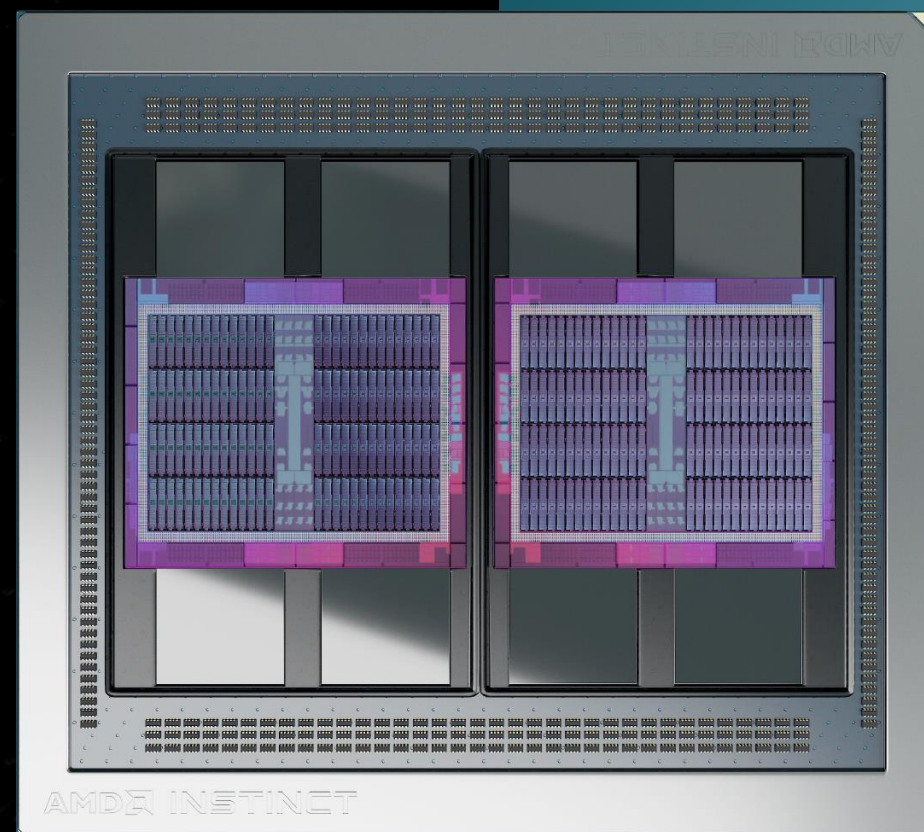
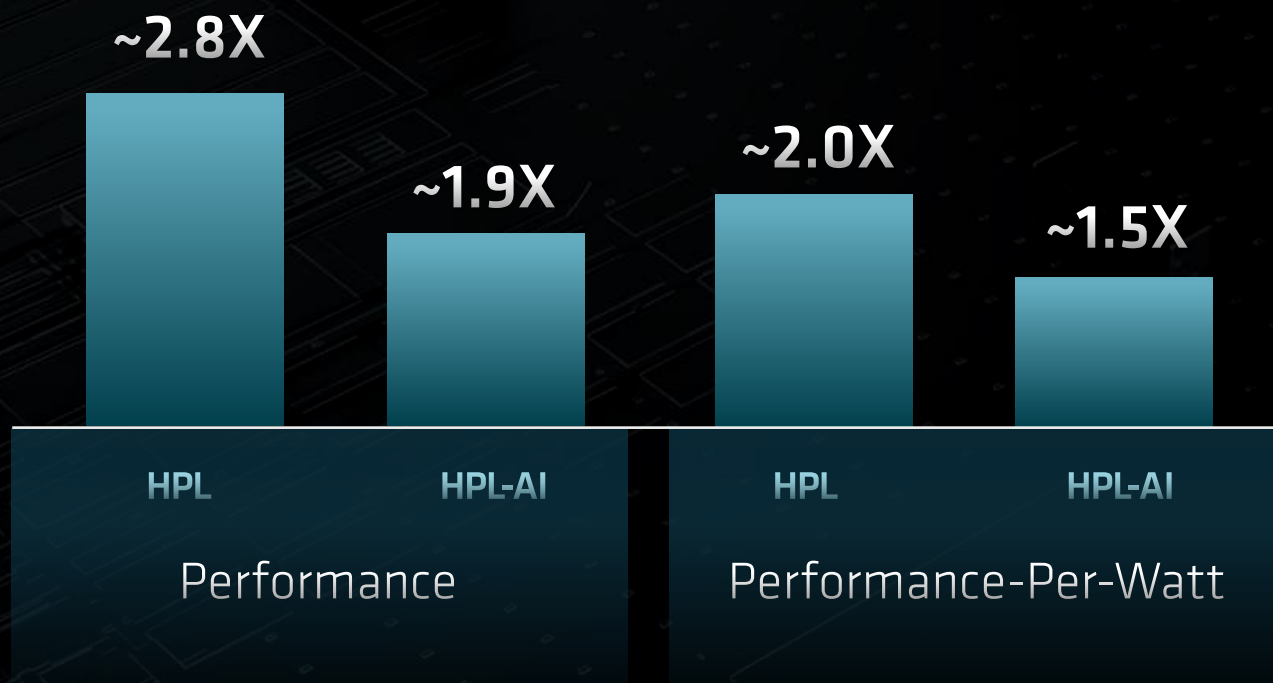
~4X Higher FP64 (HPC) and ~2X Mixed Precision (AI) Peak Performance than AMD CDNA™

## — Open Ecosystem

Broaden ROCm™ Open Software Platform, Including AI

# AMD CDNA™ 2 LEADERSHIP PERFORMANCE AND EFFICIENCY

AMD Instinct™ MI250X vs. Nvidia A100





# AMD ROCm

## OPEN SOFTWARE PLATFORM FOR GPU COMPUTE

- Unlocked GPU Performance to Accelerate Computational Tasks
- Optimized for HPC and AI Workloads at Scale
- Open Source Enabling Collaboration, Innovation, and Differentiation

AI Framework Support	ONNX Runtime						PyTorch			TensorFlow		
Programming Models	OpenMP API						HIP API			OpenCL™		
Libraries	BLAS	RAND	FFT	MIGraphX	MIVisionX	PRIM						
	SOLVER	ALUTION	SPARSE	THRUST	MIOpen	RCCL						
Compilers & Tools	Compiler	Profiler	Tracer	Debugger	hipify	TENSILE						
Drivers & Runtime	RedHat, CentOS, SLES & Ubuntu Device Drivers and Run-Time											
Deployment Tools	ROCm Validation Suite				ROCm Data Center Tool				ROCm SMI			



FROM HPC TO AI

# ROCm™ JOURNEY TO ECOSYSTEM ENABLEMENT

AMD   
**ROCm 4**  
HPC and Exascale

- Optimized** HPC Performance
- Enabled** AMD CDNA™ GPUs
- Hardened** for Scale



AMD   
**ROCm 5**  
Expanding to AI

- Optimized** Training & Inference Performance
- Enabling** AMD CDNA™ & AMD RDNA™ GPUs
- SDK for** Development & Deployment

2021

2022



# INNOVATING AI THROUGH DEEP PARTNERSHIPS



## Accelerating ROCm Performance on AI

“We’re also deepening our investment in the open-source PyTorch framework, working with the PyTorch Core team and AMD both to **optimize the performance and developer experience for customers running PyTorch on Azure, and to ensure that developers’ PyTorch projects work great on AMD hardware.**”

**Kevin Scott**

Executive Vice President and CTO, Microsoft



## Optimizing ROCm for PyTorch

“We are excited to partner with AMD to grow our PyTorch support on ROCm, enabling the vibrant PyTorch community to **adopt the latest generation of AMD Instinct GPUs quicker than ever with great performance for major AI use cases running on PyTorch.**”

**Soumith Chintala**

Co-Creator and Lead of PyTorch, Meta AI



## Developing Data Centric AI for Instinct™ GPUs

“We are excited to partner with AMD to leverage AMD MI200 GPUs and the ROCm stack to port and optimize LandingLens, our GPU-optimized computer vision software application. **The power of AMD GPUs and the maturity of ROCm will help Landing AI continue to deliver acceleration of high-resolution computer vision models with the purpose of providing better insights into manufacturing defects.**”

**Andrew Ng**

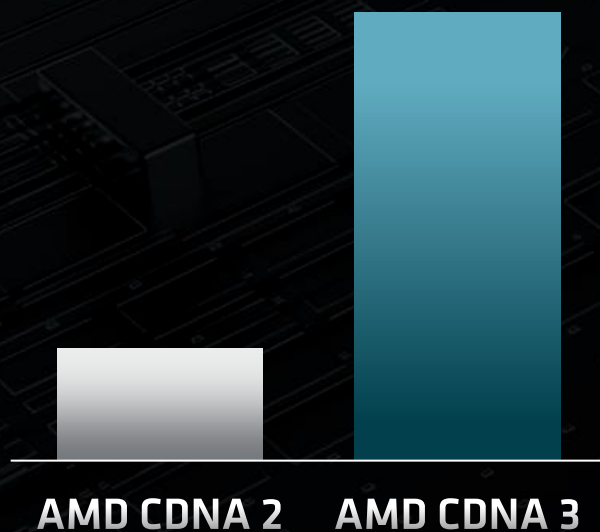
CEO of Landing AI and Adjunct Professor, Stanford University

# AMD CDNA 3

# THE JOURNEY CONTINUES

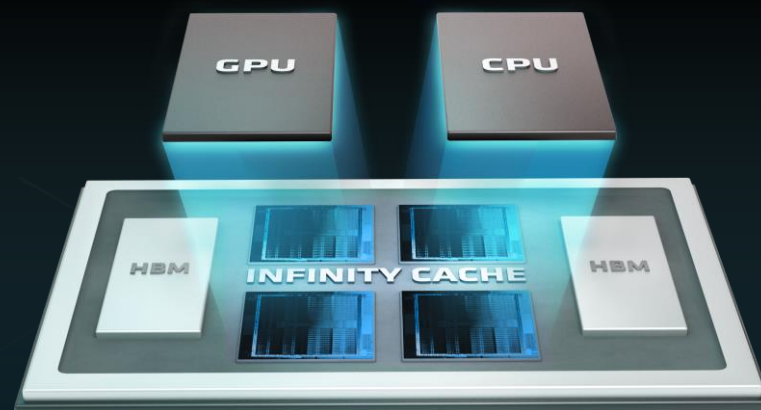
AI Performance/Watt Uplift

>5X



Expected performance-per-watt uplift through:

- 5nm Process and 3D Chiplet Packaging
- Next-Gen AMD Infinity Cache™
- 4<sup>th</sup> Gen Infinity Architecture
- Unified Memory APU Architecture
- New Math Formats







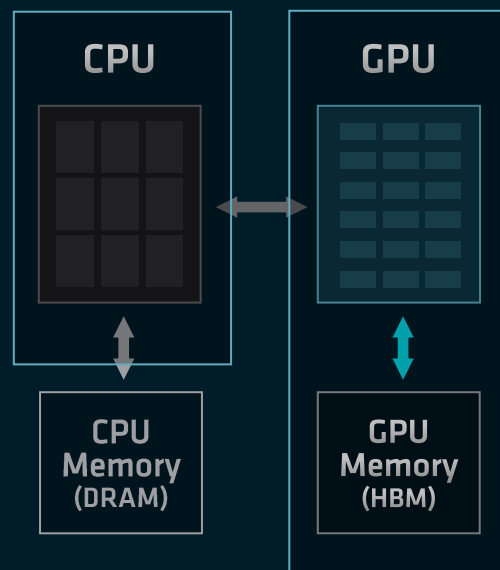
# UNIFIED MEMORY APU ARCHITECTURE BENEFITS

## AMD CDNA™ 2 Coherent Memory Architecture

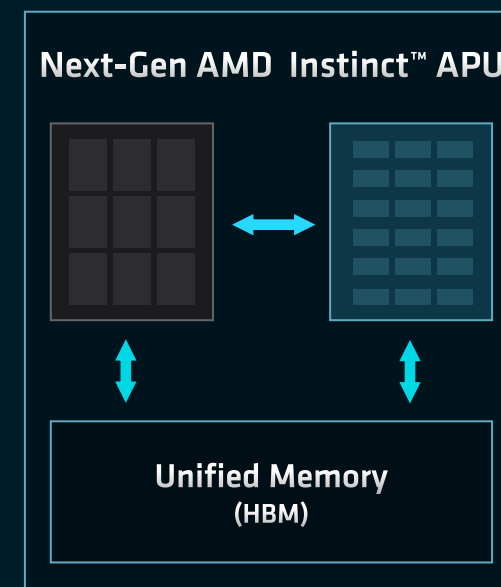


## AMD CDNA™ 3 Unified Memory APU Architecture

- Simplifies Programming
- Low Overhead 3<sup>rd</sup> Gen Infinity Interconnect
- Industry Standard Modular Design



- Eliminates Redundant Memory Copies
- High-Efficiency 4<sup>th</sup> Gen AMD Infinity Architecture
- Low TCO with Unified Memory APU Package





# AMD COMPUTE GPU ARCHITECTURE ROADMAP

## DRIVING PERFORMANCE AND EFFICIENCY LEADERSHIP



2020

2023

All roadmaps are subject to change.

OUR PATH FORWARD

# GPU TECHNOLOGY FOCUS

- Leadership AMD RDNA™ / AMD CDNA™ Architecture Roadmaps
- Advanced Process and Packaging Technologies
- Consistent Execution of Performance-Per-Watt Roadmap
- Expanding Open Software Ecosystems Including AI



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

© 2022 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, AMD CDNA, AMD Instinct, AMD RDNA, Infinity Cache, and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective owners.

# ENDNOTES

- RX-549 – Testing done by AMD performance labs 10/16/20, using Assassins Creed Odyssey (DX11, Ultra), Battlefield V (DX12, Ultra), Borderlands 3 (DX12, Ultra), Control (DX12, High), Death Stranding (DX12 Ultra), Division 2 (DX12, Ultra), F1 2020 (DX12, Ultra), Far Cry 5 (DX11, Ultra), Gears of War 5 (DX12, Ultra), Hitman 2 (DX12, Ultra), Horizon Zero Dawn (DX12, Ultra), Metro Exodus (DX12, Ultra), Resident Evil 3 (DX12, Ultra), Shadow of the Tomb Raider (DX12, Highest), Strange Brigade (DX12, Ultra), Total War Three Kingdoms (DX11, Ultra), Witcher 3 (DX11, Ultra no HairWorks) at 4K. System comprised of a Radeon RX 6800 XT with AMD Radeon Graphics driver 27.20.12031.1000 and an Radeon RX 5700 XT with AMD Radeon Graphics driver 26.20.13001.9005. Performance may vary.
- RX- 558 – Testing done by AMD performance labs October 20 2020 on RX 6900 XT and RX 5700 XT (20.45-201013n driver), AMD Ryzen 9 5900X (3.70GHz) CPU, 16GB DDR4-3200MHz, Engineering AM4 motherboard, Win10 Pro 64. The following games were tested at 4k at max settings: Battlefield V DX11, Doom Eternal Vulkan, Forza DX12, Resident Evil 3 DX11, Shadow of the Tomb Raider DX12. Performance may vary. RX-558RS-462 – Testing conducted by AMD as of February 16, 2022, on a test system configured with a Ryzen 5 5600X CPU, 16GB DDR4, Radeon RX 6800 XT GPU, and Windows 10 Pro, with AMD Software: Adrenalin 22.3.1 vs. 20.12.2 at 4K, Max settings. Performance may vary. Games tested: Age of Empires 4, Borderlands 3, Call of Duty: Vanguard, Cyberpunk 2077, F1 2021, Far Cry 6, Forza Horizon 4, Guardians of the Galaxy, Hitman 3, The Medium, Metro Exodus Enhanced Edition, Myst, RDR2, Resident Evil Village, and PUBG.
- RX – 785 Testing done by AMD performance labs May 31, 2022, on (11) AMD Radeon™ RX 6000 Series graphics cards, using systems configured with Ryzen™ 9 5900X and Ryzen™ 5 5600X CPUs, each with 16GB DDR4-3600MHz and AMD Smart Access Memory enabled, Win 10 Pro versus similarly configured systems with (11) Nvidia GeForce RTX 3000 Series, GeForce GTX 1650 and GTX 1050 Ti GPUs, each with ReBAR enabled. Performance tested across 20 games at 4K, 1440P and 1080P resolutions, at intended settings for each of the (22) AMD and NVIDIA GPUs. Performance per watt and per die size calculated using the total board power (TBP) and die sizes of the (22) individual AMD and NVIDIA GPUs over the average FPS scores. Performance may vary.
- RS-462 – Testing conducted by AMD as of February 16, 2022, on a test system configured with a Ryzen 5 5600X CPU, 16GB DDR4, Radeon RX 6800 XT GPU, and Windows 10 Pro, with AMD Software: Adrenalin 22.3.1 vs. 20.12.2 at 4K, Max settings. Performance may vary. Games tested: Age of Empires 4, Borderlands 3, Call of Duty: Vanguard, Cyberpunk 2077, F1 2021, Far Cry 6, Forza Horizon 4, Guardians of the Galaxy, Hitman 3, The Medium, Metro Exodus Enhanced Edition, Myst, RDR2, Resident Evil Village, and PUBG.
- GD-164 – Day-0 driver compatibility and feature availability depend on system manufacturer and/or packaged driver version. For the most up-to-date drivers, visit [AMD.com](https://www.amd.com)
- MI200-05 - Measurements conducted by AMD Performance Labs as of Sep 10, 2021 on the AMD Instinct™ MI250X accelerator designed with AMD CDNA™ 2 6nm FinFET process technology with 1,700 MHz engine clock resulted in 47.9 TFLOPS peak double precision (FP64) floating-point, 383.0 TFLOPS peak Bfloat16 format (BF16) floating-point performance. The results calculated for AMD Instinct™ MI100 GPU designed with AMD CDNA 7nm FinFET process technology with 1,502 MHz engine clock resulted in 11.54 TFLOPS peak double precision (FP64) floating-point, 92.28 TFLOPS peak Bfloat16 format (BF16) performance.

# ENDNOTES

- MI200-01 - World's fastest data center GPU is the AMD Instinct™ MI250X. Calculations conducted by AMD Performance Labs as of Sep 15, 2021, for the AMD Instinct™ MI250X (128GB HBM2e OAM module) accelerator at 1,700 MHz peak boost engine clock resulted in 95.7 TFLOPS peak theoretical double precision (FP64 Matrix), 47.9 TFLOPS peak theoretical double precision (FP64), 95.7 TFLOPS peak theoretical single precision matrix (FP32 Matrix), 47.9 TFLOPS peak theoretical single precision (FP32), 383.0 TFLOPS peak theoretical half precision (FP16), and 383.0 TFLOPS peak theoretical Bfloat16 format precision (BF16) floating-point performance.  
Calculations conducted by AMD Performance Labs as of Sep 18, 2020 for the AMD Instinct™ MI100 (32GB HBM2 PCIe® card) accelerator at 1,502 MHz peak boost engine clock resulted in 11.54 TFLOPS peak theoretical double precision (FP64), 46.1 TFLOPS peak theoretical single precision matrix (FP32), 23.1 TFLOPS peak theoretical single precision (FP32), 184.6 TFLOPS peak theoretical half precision (FP16) floating-point performance.  
Published results on the NVidia Ampere A100 (80GB) GPU accelerator, boost engine clock of 1410 MHz, resulted in 19.5 TFLOPS peak double precision tensor cores (FP64 Tensor Core), 9.7 TFLOPS peak double precision (FP64), 19.5 TFLOPS peak single precision (FP32), 78 TFLOPS peak half precision (FP16), 312 TFLOPS peak half precision (FP16 Tensor Flow), 39 TFLOPS peak Bfloat 16 (BF16), 312 TFLOPS peak Bfloat16 format precision (BF16 Tensor Flow), theoretical floating-point performance. The TF32 data format is not IEEE compliant and not included in this comparison.  
<https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/nvidia-ampere-architecture-whitepaper.pdf>, page 15, Table 1.
- MI200-26B – Testing Conducted by AMD performance lab as of 10/14/2021, on a single socket Optimized 3rd Gen AMD EPYC™ CPU (64) server, with 1x AMD Instinct™ MI250X OAM (128 GB HBM2e, 560W) GPU with AMD Infinity Fabric™ technology using benchmark HPL v2.3, plus AMD optimizations to HPL that are not yet upstream. vs. Nvidia DGX dual socket AMD EPYC 7742 (64C) @2.25GHz CPU server with 1x NVIDIA A100 SXM 80GB (400W) using benchmark HPL Nvidia container image 21.4-HPL Information on HPL:  
<https://www.netlib.org/benchmark/hpl/Nvidia>.  
HPL Container Detail: <https://ngc.nvidia.com/catalog/containers/nvidia:hpc-benchmarks>.  
  
Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.
- MI200-58 – Testing Conducted by AMD performance lab as of 5/25/2022, on a dual socket AMD EPYC™ 7200 Series CPUs (64C) server, with 8x AMD Instinct™ MI250X OAM (128 GB HBM2e, 500W) GPU with AMD Infinity Fabric™ technology using benchmark HPL-AI compiled with HIP version 5.1.20531.cacfa990, AMD clang version 14.0.0, OpenMPI 4.1.2. vs. dual socket AMD EPYC 7002 (64C) Series CPU (64) server with 8x NVIDIA A100 SXM 80GB (400W) using benchmark HPL-AI with CUDA 11.6. HPL-AI container 21.4-hpl Information on HPL-AI: <https://hpl-ai.org/>  
AMD HPL-AI container detail: <https://github.com/ROCmSoftwarePlatform/hpl-ai.git> rev bae3342  
Nvidia HPL-AI Container Detail: <https://catalog.ngc.nvidia.com/orgs/nvidia/containers/hpc-benchmarks> Nvidia container image 21.4-HPL  
Server manufacturers may vary configurations, yielding different results. Performance may vary based on use of latest drivers and optimizations.
- MI300-004 – Measurements by AMD Performance Labs June 4, 2022. MI250X (560W) FP16 (306.4 estimated delivered TFLOPS based on 80% of peak theoretical floating-point performance). MI300 FP8 performance based on preliminary estimates and expectations. MI300 TDP power based on preliminary projections. Actual results based on production silicon may vary.