# The Journey to AI Pervasiveness

**Victor Peng**
**President, AMD**

**AMD**
together we advance_

# Evolution of AI

**Early History: Supervised Learning, Perceptrons, MNIST**

1960s — 1989

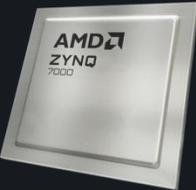| V11 M Chip | Early NPUs at ISSCC | Alpha 21264 |
| :---: | :---: | :---: |
| ~1984 | 1991 | ~1998 |

AMD
together we advance_

# Evolution of AI

**Early History: Supervised Learning, Perceptrons, MNIST**

**Deep Neural Networks enabled by ImageNet and GPUs** · **AlexNet**

**Transformers**

**Chat GPT**

**Agentic AI** Multi-Model MOEs SSMs

| 1960s | 1989 | 2012 | 2014 | 2016 | 2018 | 2020 | 2022 | 2024 |

First FPGA SOC
**AMD Zynq™ 7040**
~2013

First AIE Device
**VC1902**
~2021

Largest FPGA Ever
**VP1902**
~2022

Leadership GPU for AI Inference
**MI300X**
~2024

AMD
together we advance_

# The Journey to AI Pervasiveness in Supercycles

**Embedded**
**Automotive, Sensor Intelligence, Robotics, Medical**

**Endpoints**
**Personal Assistants, Games, Copilot**

**Cloud**
**Large-Scale Training and Inference**

Today's AI

**Virtual** ❯ **Personalized** ❯ **Physical**

**With Increasing Unit Volume**

AI Pervasiveness

AMD
together we advance_

4

# The Investment and Value Creation Race



- CAPEX Data Center AI
- Data Center AI Revenue

2023     2024 (est)     2025 (est)

Today's AI     AI Pervasiveness

AMD
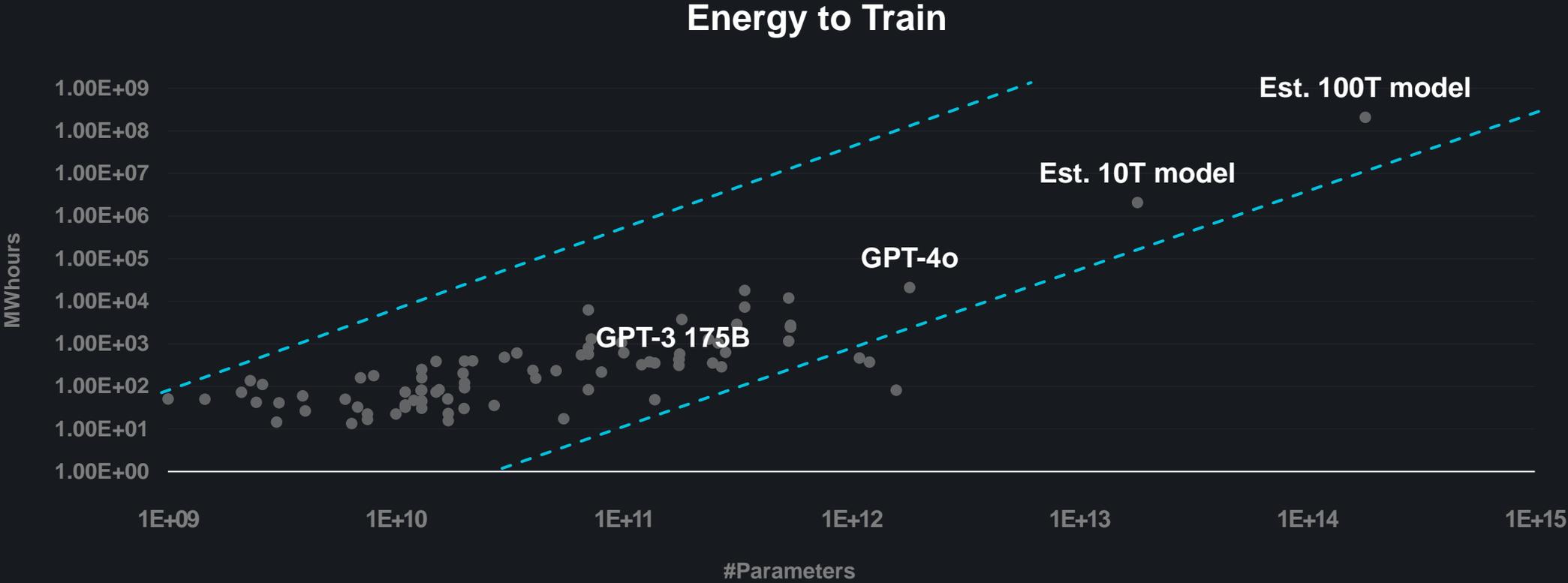together we advance_

# The Challenges Ahead

Today's AI

AI Pervasiveness

**Additional challenges:** High quality data, safety, reliability, confidentiality…

**On the technology side:** Foremost power source & distribution, and diverse requirements of diverse use cases

AMD
together we advance_

# Power to Train Frontier Models

**Energy to Train**



Exponential growth in model sizes drives massive increase in energy required for training

Total data center power capacity limits model size

AMD
together we advance_

# Energy Efficiency Concerns all Segments

## Cloud
- **Total data center power capacity limits model size**
- **Network power contribution is significant (~20%)**
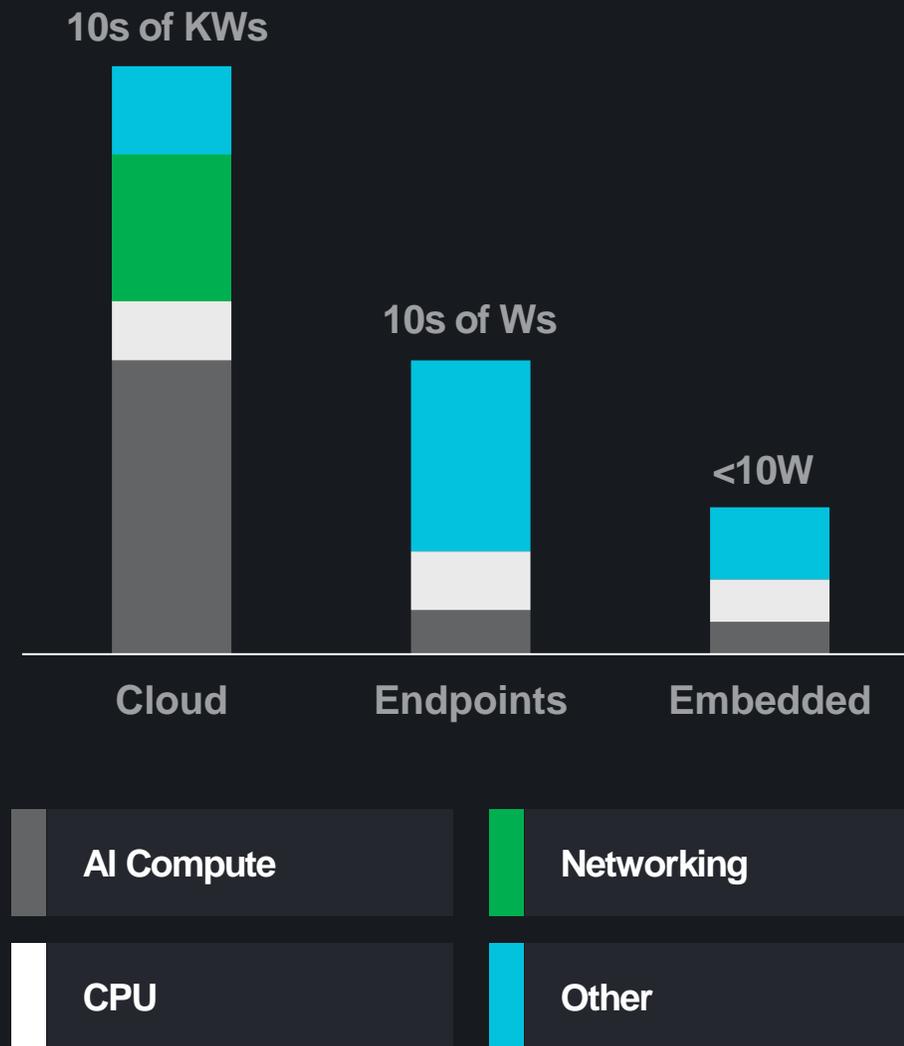
## Endpoints
- **Trade-off between user experience and battery life**

## Embedded
- **Total power envelop highly constrained**

At larger scale we can amortize system overheads for more efficiency

**System-Level Breakdown of Power**

10s of KWs

10s of Ws

<10W

Cloud     Endpoints     Embedded

- AI Compute
- Networking
- CPU
- Other

Based on AMD internal calculations

8

AMD together we advance_

# Different Requirements Across Segments
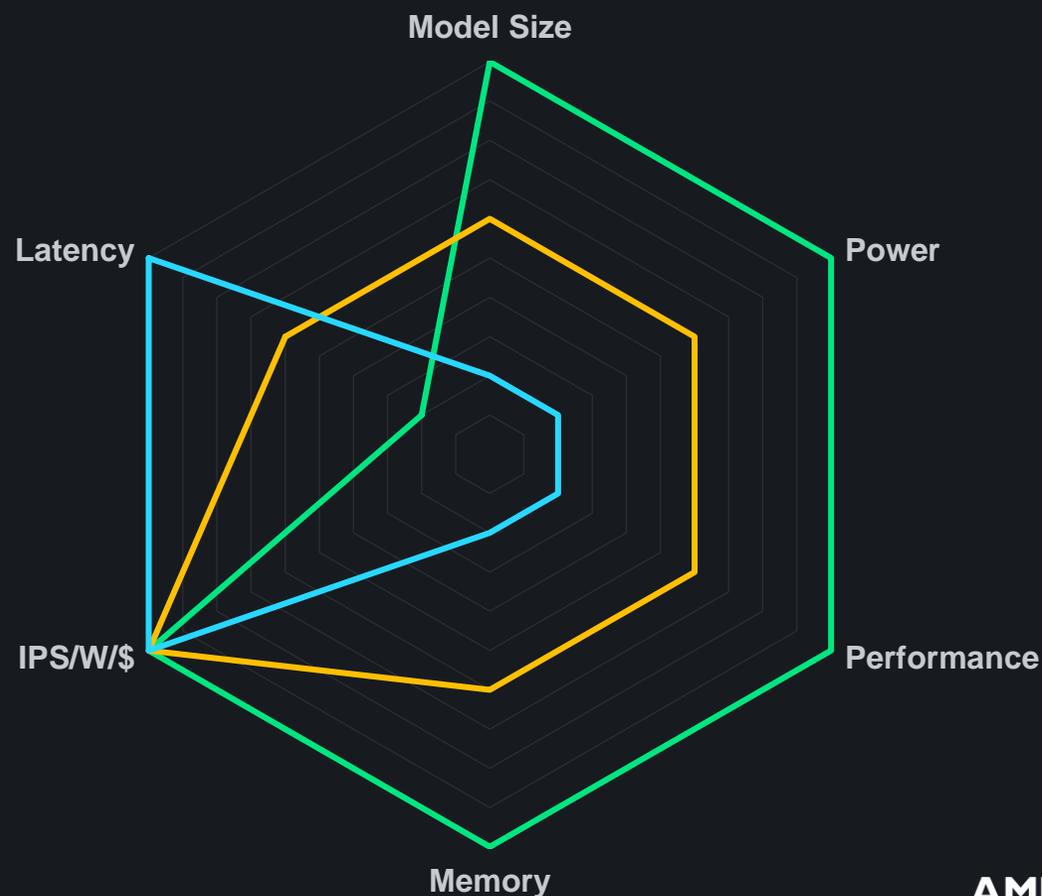
## Cloud AI

- 1000+ of AI TFLOPS per GPU
- TCO
- Interconnect to thousands of nodes
- Feature velocity critical

## Endpoint AI

- 100s of AI TOPS per APU
- Interactive latency
- Highly optimized vendor libraries
- Feature velocity for developers

## Embedded AI

- 10s of AI TOPS per SOC
- Real-time requirements driving low latency
- Form factor & functional safety focused workloads
- Customization, I/O

Radar chart axes: Model Size, Power, Performance, Memory, IPS/W/$, Latency

AMD
together we advance_

# AMD Solutions to Address Key Challenges

**1**

**Energy-Efficient Performance**

**2**

Diverse Requirements

Today's AI

AI Pervasiveness

AMD

together we advance_
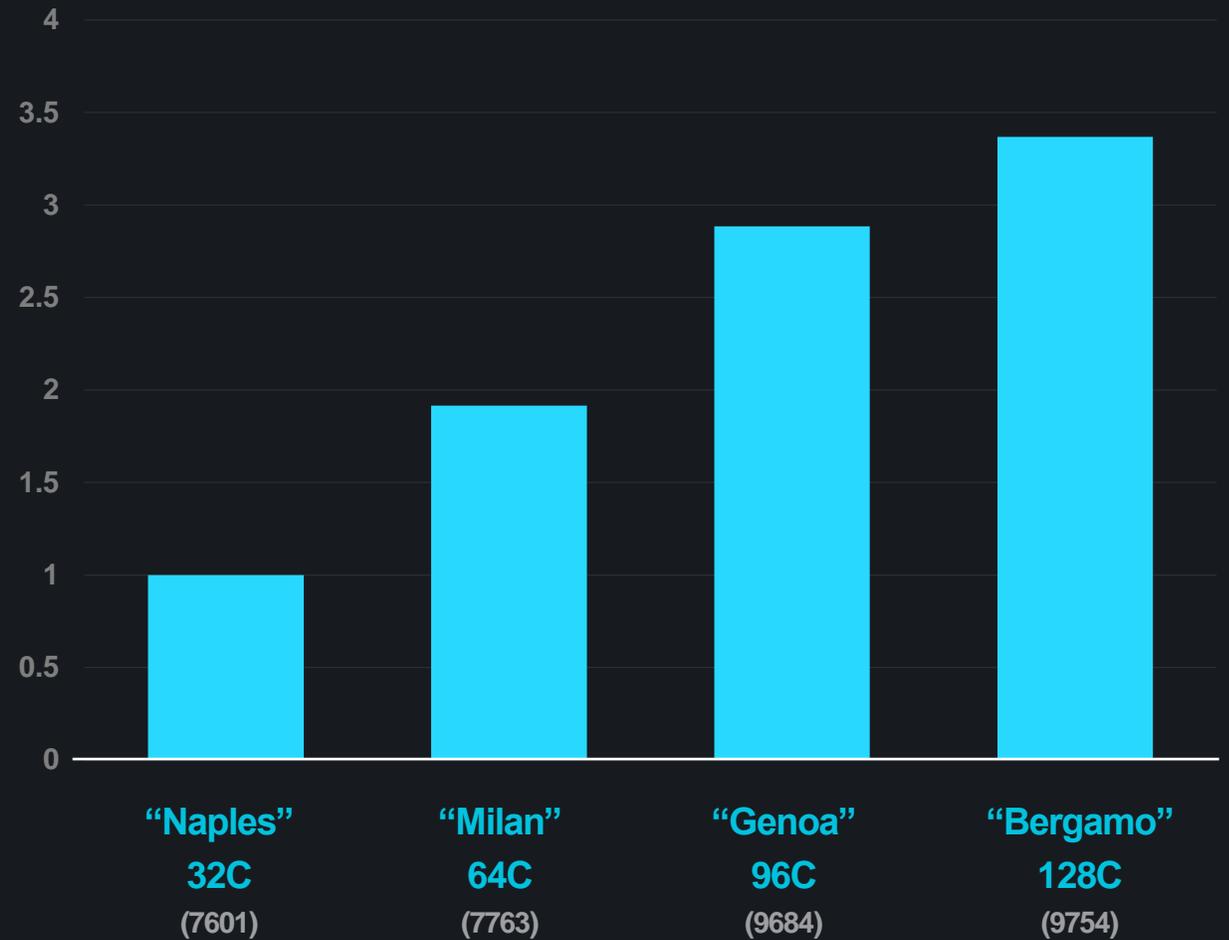
# AMD - Leadership Track Record in Energy Efficiency

**A decade of focus on energy efficiency**

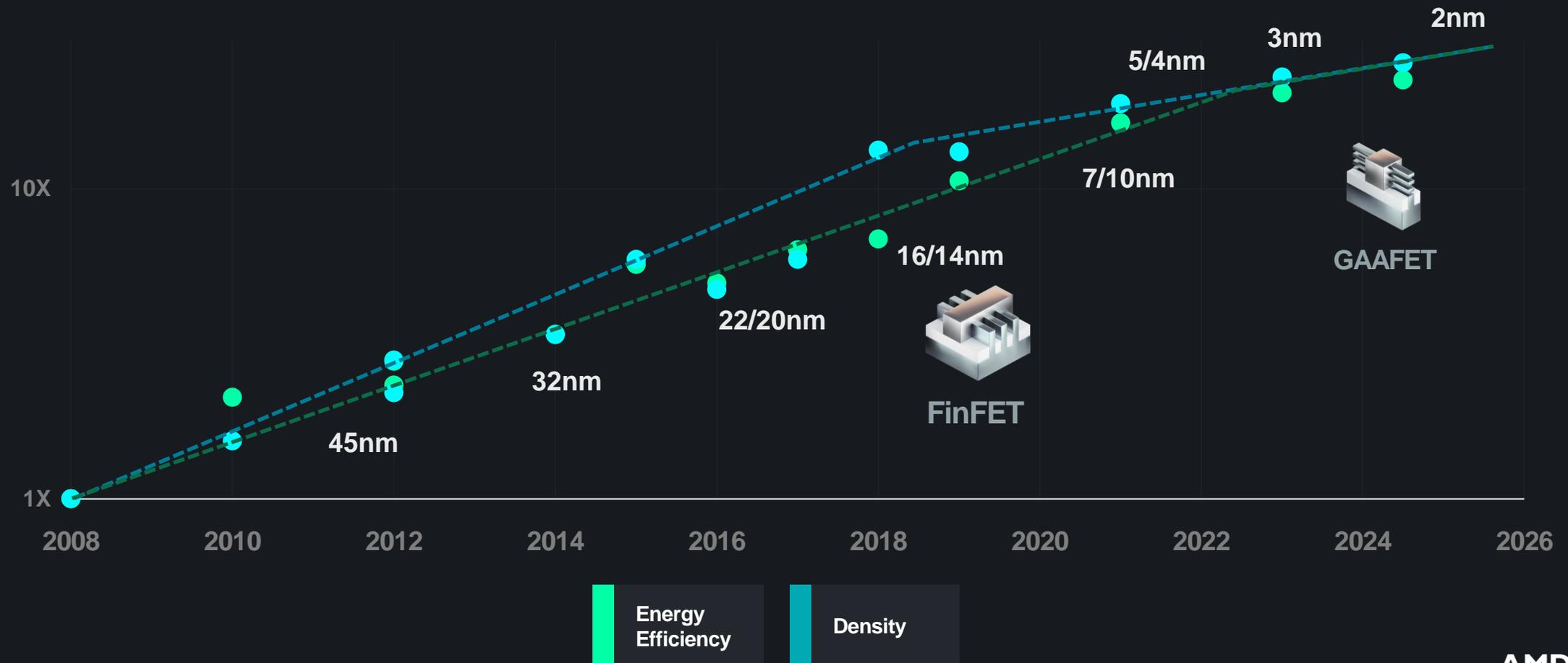**Leadership efficiency products in server, client, graphics and HPC**

## Performance/Watt



*Normalized to Naples 32C (7601), based on EPYC processor TDP and SpecIntRate2017
Naples: 141, 180W TDP; Milan: 424, 280W, Genoa: 904, 400W, Bergamo 948, 360W
https://www.spec.org/

AMD
together we advance_

# Relentless Focus on Performance and Energy Efficiency – All Levels



**Technology**

**Architecture**

**System**

**DC Infrastructure**

AMD together we advance_

# Energy Efficiency Gains at Silicon Level

## Technology gains slowing but essential

2nm

3nm

5/4nm

10X

7/10nm

22/20nm

16/14nm

GAAFET

32nm

FinFET

45nm

1X

2008   2010   2012   2014   2016   2018   2020   2022   2024   2026

Energy Efficiency

Density

Based on AMD internal estimates

13

AMD together we advance_

# Cost Scaling Challenges

## New node introduction rate is slowing

90nm
65nm
45nm
32nm
22nm
14nm
10/7nm
5nm
3nm
2nm

2004 2006 2008 2010 2012 2014 2016 2018 2020 2022 2024

## While costs continue to increase

**(Cost per yielded mm² for a 250mm² die)**

45nm 32nm 28nm 20nm 14/16nm 7nm 5nm 3nm 2nm

Source: Internal AMD Data

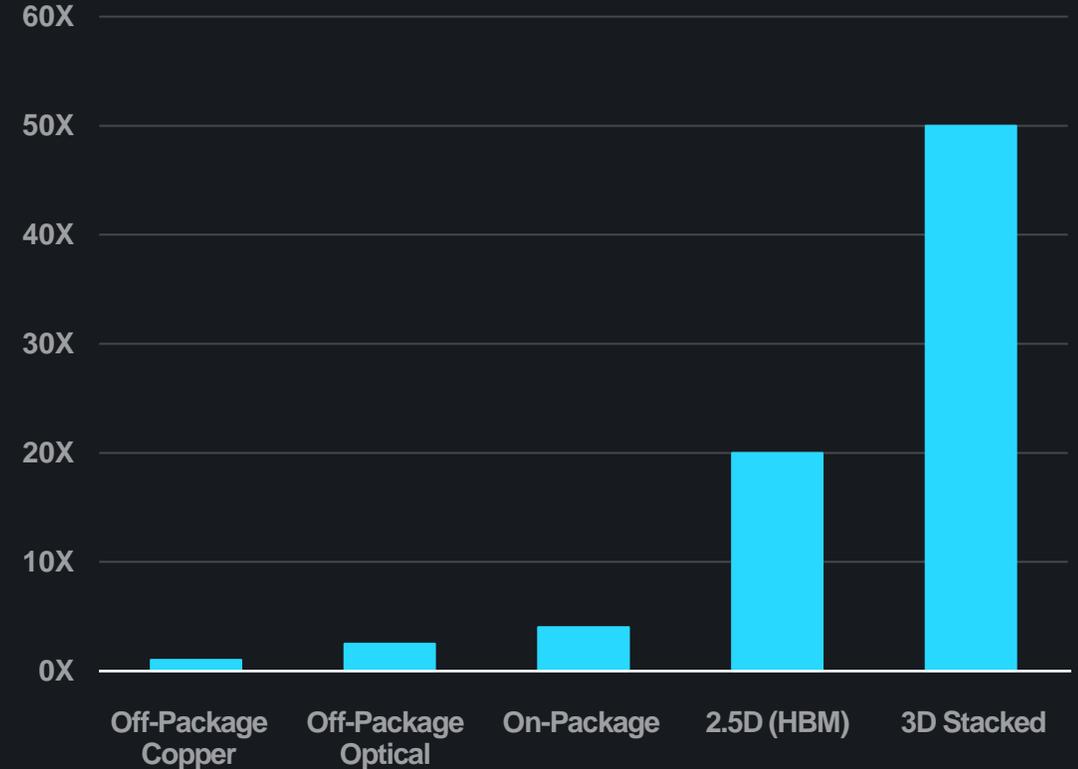# Advanced Packaging: Significant Gains in Energy-Efficient Performance

**Energy-efficient performance needs tight integration**

**2.5D enables co-packed compute with HBM**

**AMD 3D V-Cache™ technology drives energy efficiency leadership**

**Advanced 3D hybrid bonding provides by orders of magnitude the densest, most power efficient chiplet interconnect through higher bandwidth and lower latency**

L. Su and S. Naffziger, "1.1 Innovation For the Next Decade of Compute Efficiency," 2023 IEEE International Solid-State Circuits Conference (ISSCC), San Francisco, CA, USA, 2023, pp. 8-12, doi: 10.1109/ISSCC42615.2023.10067810

## Relative Bits/Joule

| | Value |
|---|---|
| Off-Package Copper | ~1X |
| Off-Package Optical | ~2.5X |
| On-Package | ~4X |
| 2.5D (HBM) | 20X |
| 3D Stacked | 50X |

Source: Internal AMD Data

**AMD** together we advance_

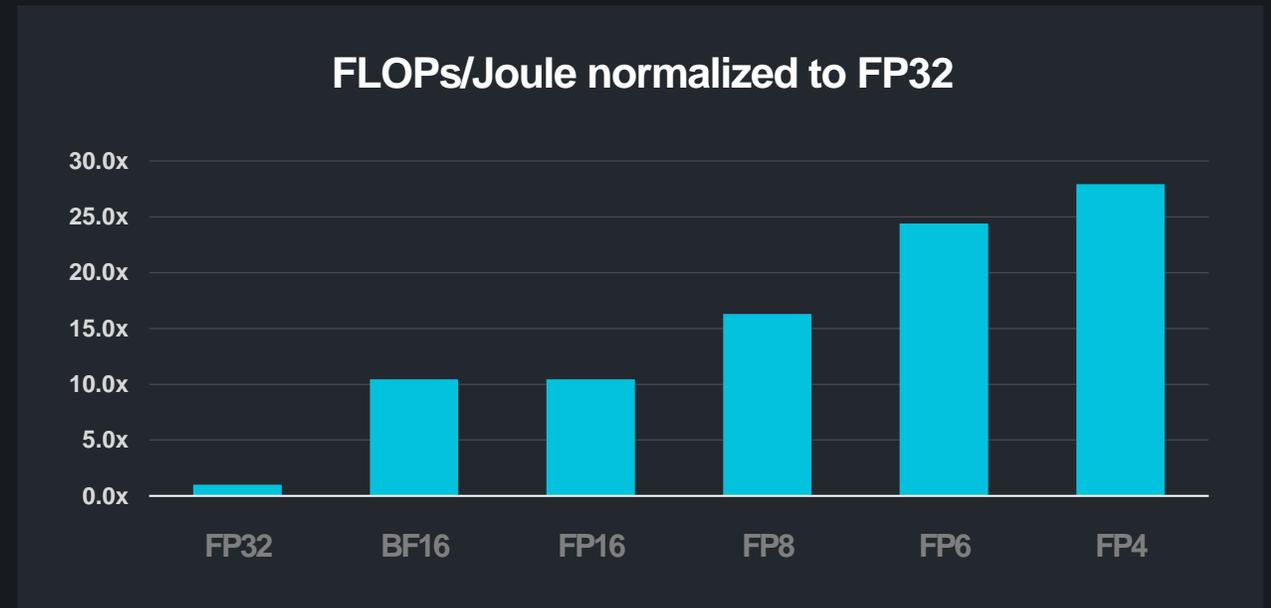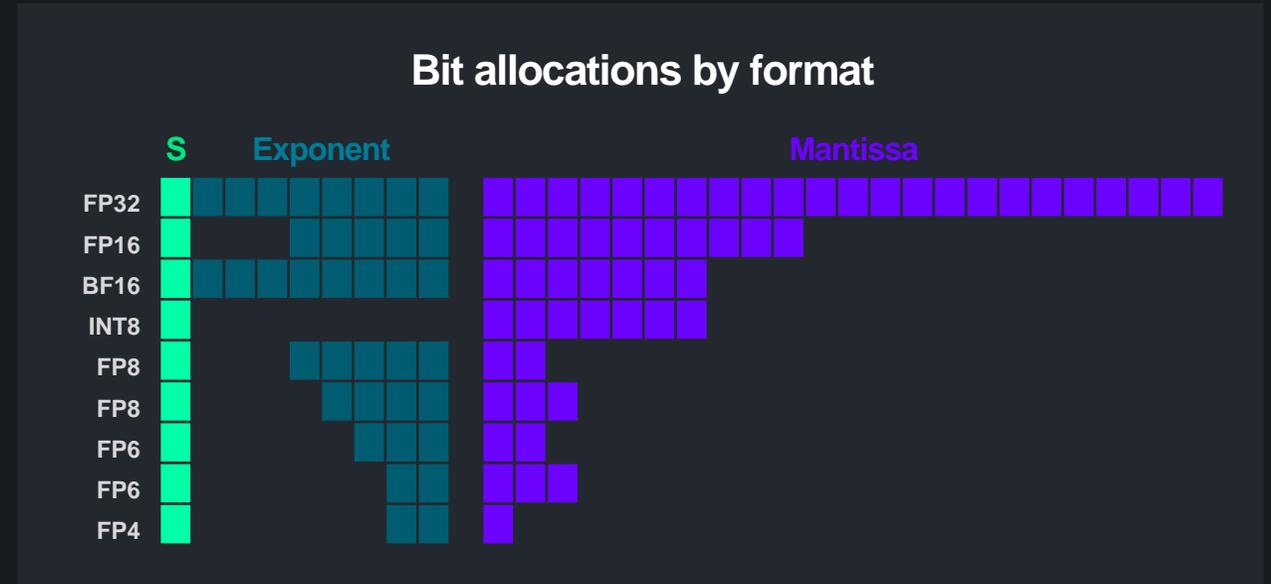# Enabling Further Efficiency through Advanced Quantization

Adapt algorithms to use lower precision math for significant improvements in energy efficiency

- FP8 in MI325X and FP4 and FP6 in MI350X

Advancing innovation in quantization

Novel research into accumulator-aware quantization[1]

Collaborating through open-source[2]

[1] https://arxiv.org/abs/2301.13376
[2] *brevitas/src/brevitas_examples/imagenet_classification/a2q

## Bit allocations by format



## FLOPs/Joule normalized to FP32



Source: Internal AMD Data

# Efficiency through Specialization with Dataflow



DNN Dataflow Execution Architecture

**AI model mapped as dataflow architecture**

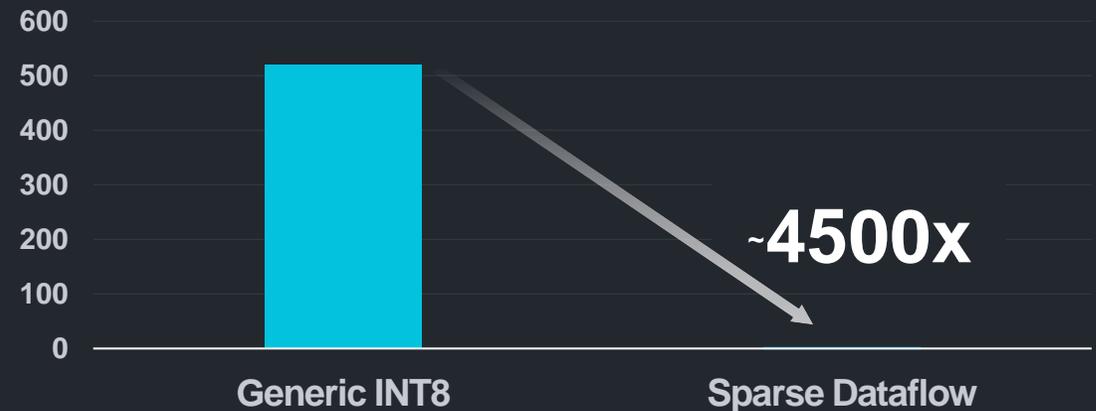**Eliminates intermediate buffering between layers**

**Additional opportunities/research**

- Fine-granular customization of datatypes, operations and connectivity
- Potentially replicating irregular sparse compute graphs
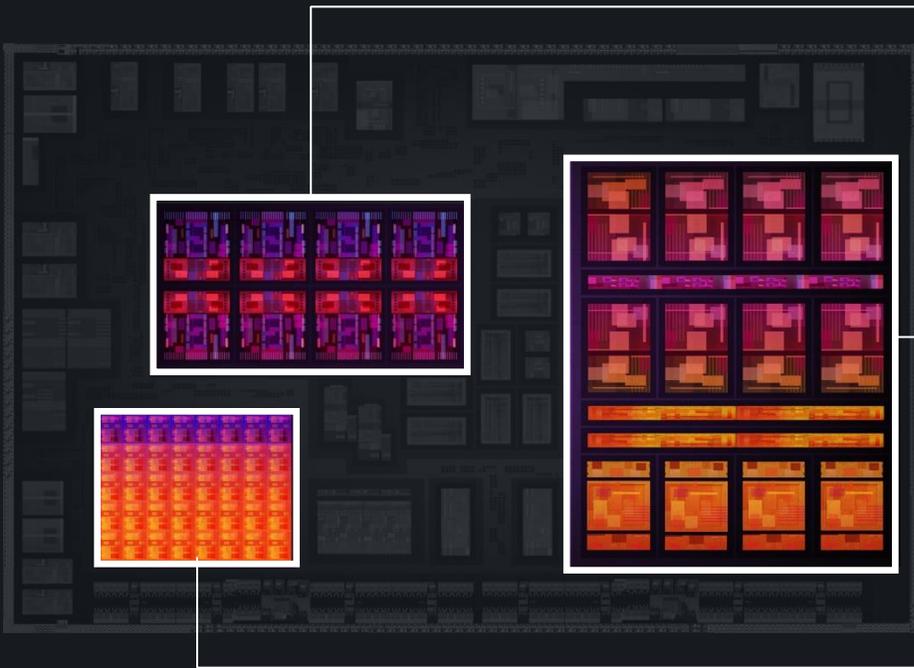- Demonstrates huge potential with ~4500x on FPGAs



**Energy per Inference [uJoules]**

~4500x

Generic INT8     Sparse Dataflow

AMD
together we advance_

# Efficiency through Heterogeneity
## 3rd Generation AMD Ryzen™ AI Technology

**Next-Gen GPU**
Up to 16 Compute Units

**AMD**
RDNA 3.5

**Next-Gen CPU**
Up to 12 Cores, 24 Threads

ZEN 5

**Next-Gen NPU**
Industry-leading 50 NPU TOPS

**AMD**
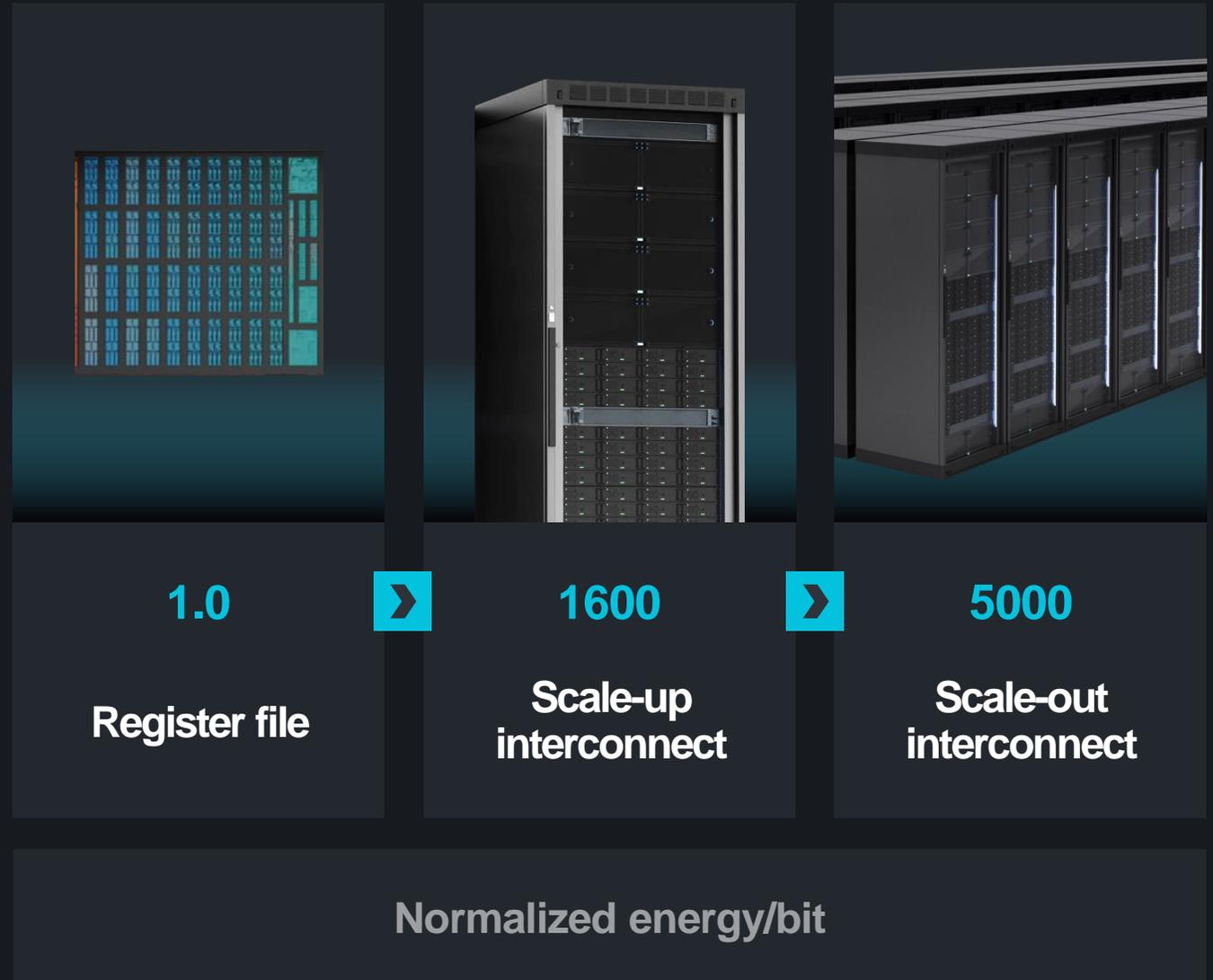XDNA 2

18

*See endnotes: STX-04

**AMD**
together we advance_

# Efficiency in Data Centers:
## Reducing Data Movement Energy

Networking is significant power contributor at data center scale

Maximizing locality at data center level becomes critical to efficiency
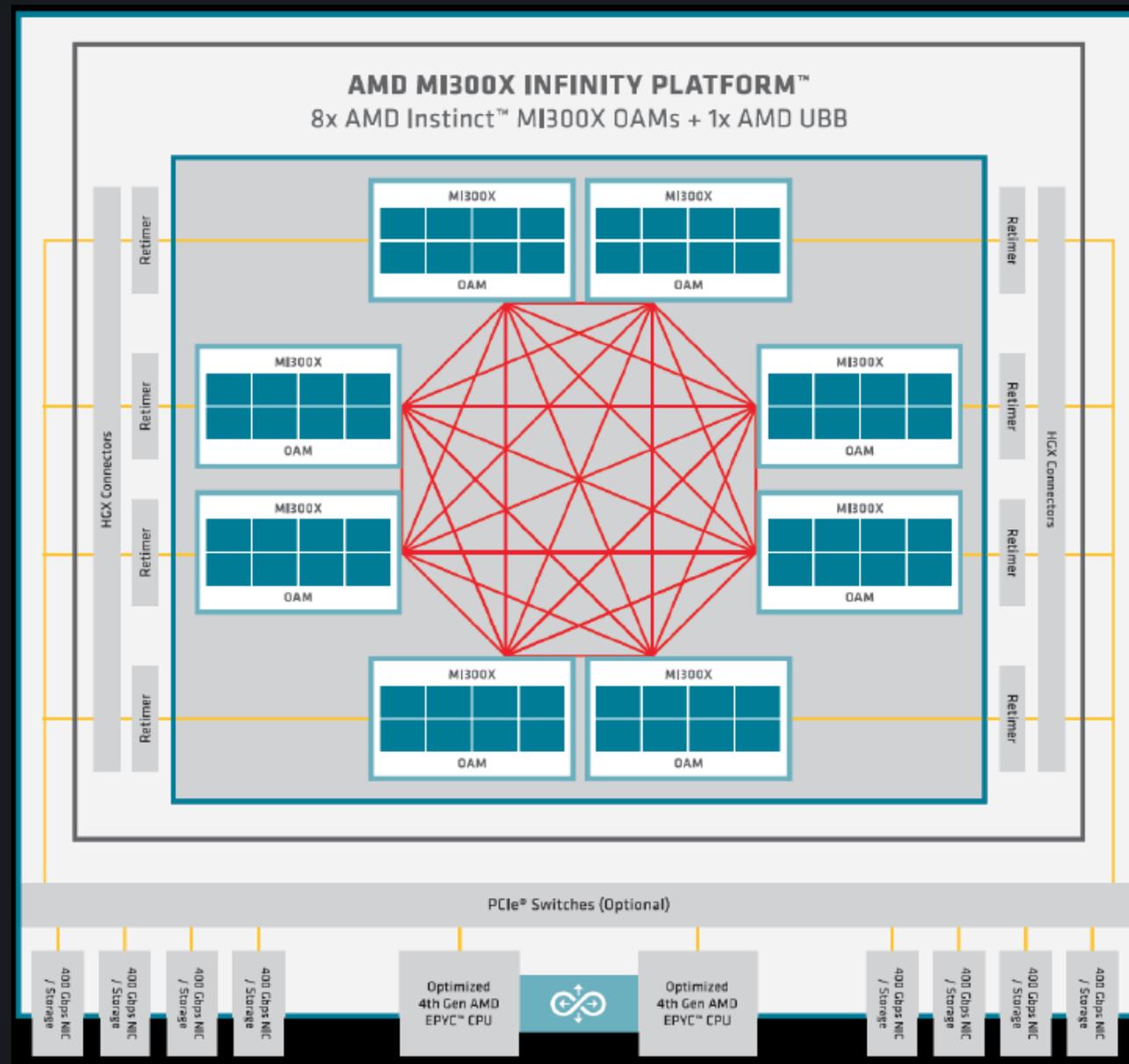
Hierarchical interconnects with scale-up and scale-out

| 1.0 | > | 1600 | > | 5000 |
|---|---|---|---|---|
| **Register file** | | **Scale-up interconnect** | | **Scale-out interconnect** |

**Normalized energy/bit**

Source: Internal AMD Data

**AMD**
together we advance_

# Efficiency Through Scale-up Interconnects
## MI300X Infinity Platform

**Direct connectivity for 8 OAMs via AMD Infinity Fabric™**

**896 GB/s AMD Infinity Fabric™ Bandwidth**

# AMD Solutions to Address Key Challenges
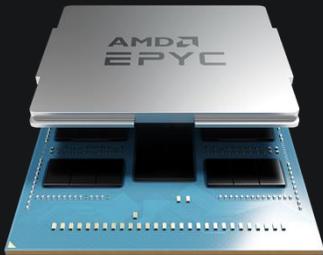
**1**

Energy-Efficient Performance

**2**

Diverse Requirements

Today's AI

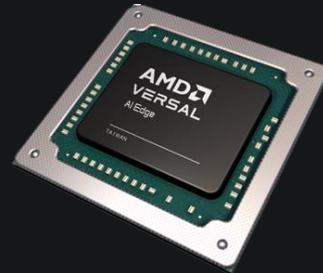AI Pervasiveness

AMD
together we advance_

# Broad Portfolio to Address Diverse Spectrum of Requirements
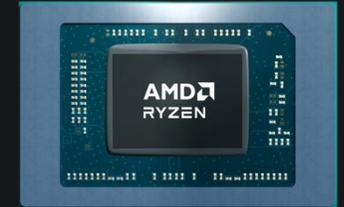


**4th Gen AMD EPYC™ Processors**

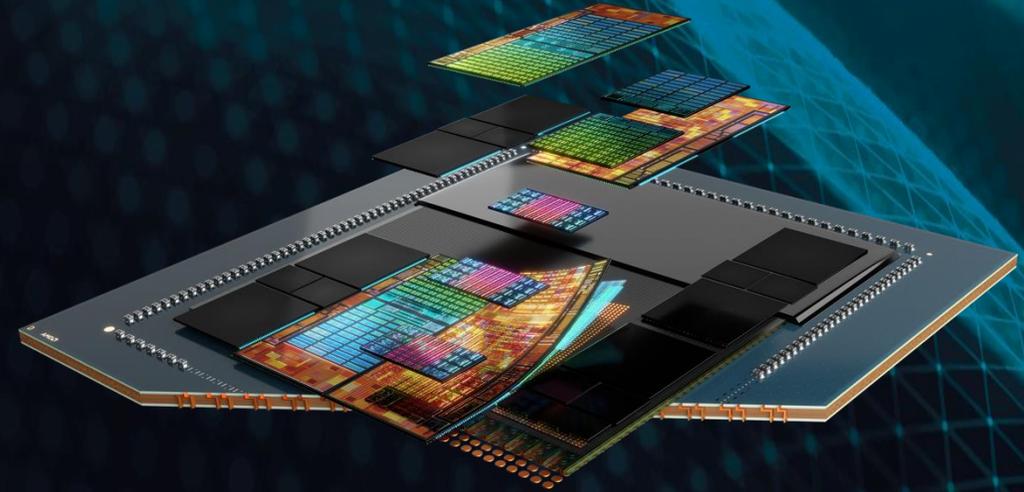**AMD Instinct™ AI Accelerators**

**AMD Versal™ FPGAs and SoCs**

**AMD Radeon™ GPUs**

**AMD Ryzen™ Mobile Processors**

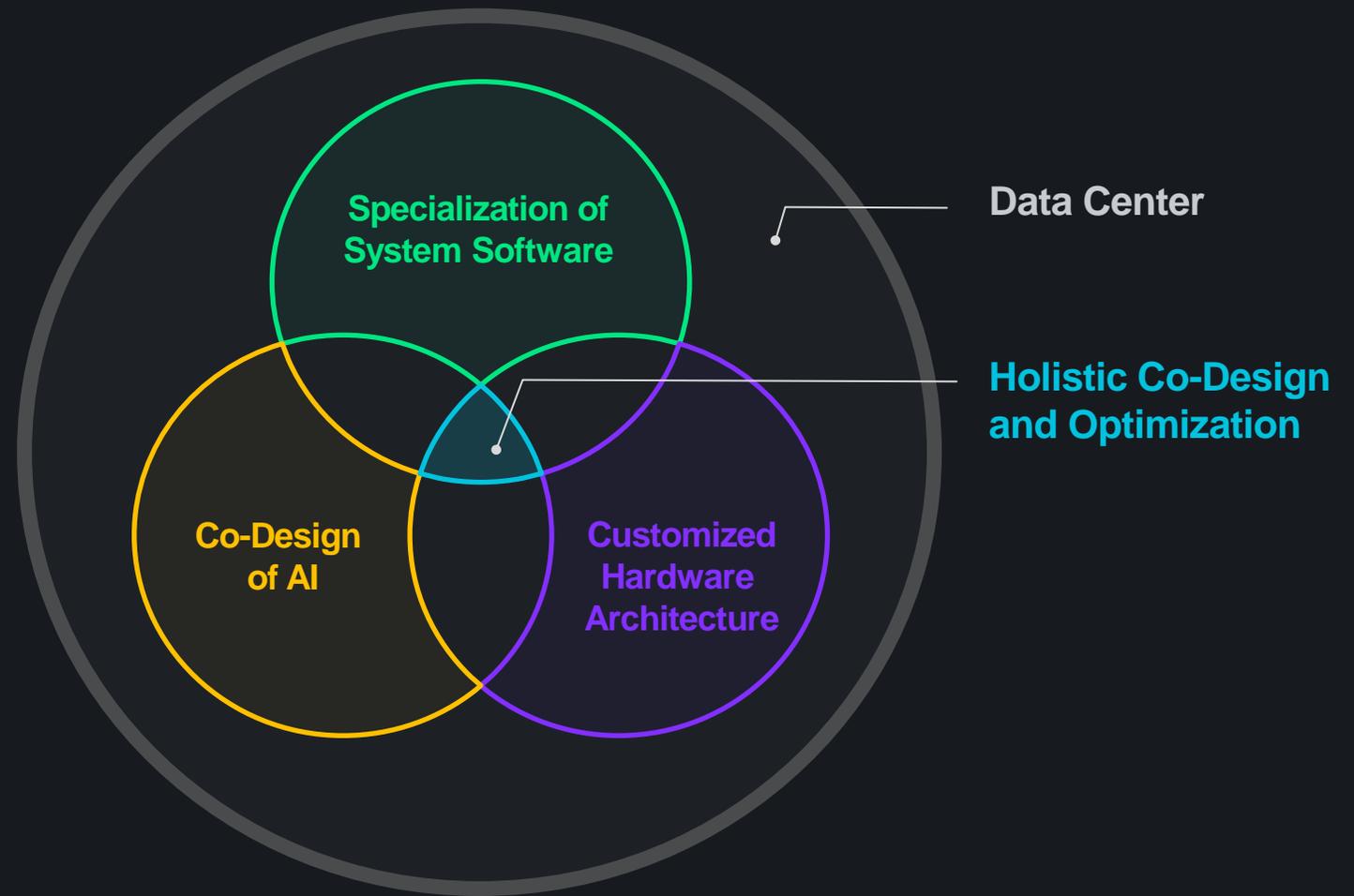## From Cloud to Embedded

AMD
together we advance_

# Device Diversity with Leadership in Chiplet and Packaging Technology

- **Long history in modularity through chiplet and advanced packaging**

- **Leveraging technology across markets**

- **Accelerated enabling of new device families to cater for emerging diverse AI requirements**

**AMD**
together we advance_

# Holistic Co-Design and Optimization to Unlock Full Efficiency and Performance

**Specialization of System Software**

**Co-Design of AI**

**Customized Hardware Architecture**

Data Center

Holistic Co-Design and Optimization

**From model parallelization, training libraries, compiler stacks, kernel libraries, runtimes, to hardware architectures**

AMD
together we advance_

# Open Ecosystem
## Driving Accelerated Innovation

**Open-source research, code sharing, standards and datasets enabled critical innovation in AI**

- **ImageNet**
- **PyTorch**
- **"Attention is all you need"**
- **Llama3**

**Committed to open collaboration and community innovation**

## 🤗 Hugging Face

**700,000+ models run out-of-box on AMD ROCm™ platform**

## ⚫ PyTorch

**Fully upstreamed AMD ROCm™ platform support**
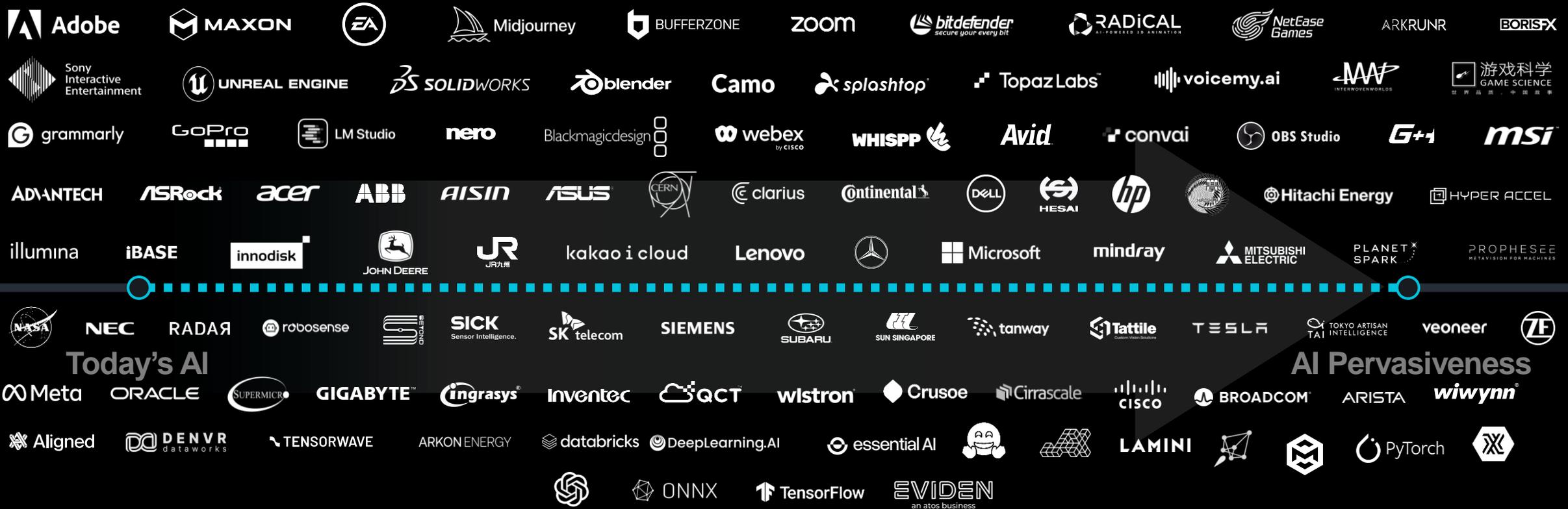
**Quantization Libraries such as Brevitas**

## OpenAI Triton

**Fully upstreamed AMD ROCm™ platform support**

**Used for key LLM kernel generation**

| | |
|---|---|
| JAX | vLLM |
| TensorFlow | MLIR IREE |
| ONNX Runtime | OpenXLA |

AMD
together we advance_

# On the Way to AI Pervasiveness
# Active Adoption from Cloud to Embedded



Today's AI

AI Pervasiveness

AMD
together we advance_

# Llama 3.1 enabled by MI300X GPUs

**Memory Requirements [GB]**

- 1536GB available
- 1215GB used

**8x AMD Instinct™ MI300X GPUs**
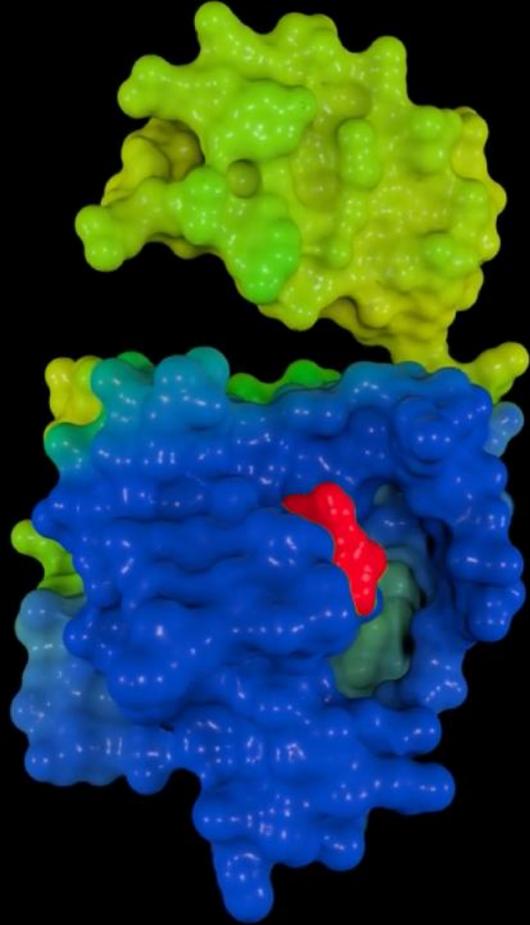**Dell PowerEdge™ XE9680**

## Llama 3 ∞

- Democratizing **Frontier-Level AI** Through Open Source
- Energy Efficiency Through Compact Parameter Footprint
- **405B** on a **single** 8-GPU server

**AMD**
together we advance_

**Program** a better asparaginase protein to starve cancer cells of asparagine

**Generative AI for Molecule Programming Accelerated by an AMD EPYC™ processor and AMD Instinct™ GPU based cloud infrastructure**

AMD◢

together we advance_
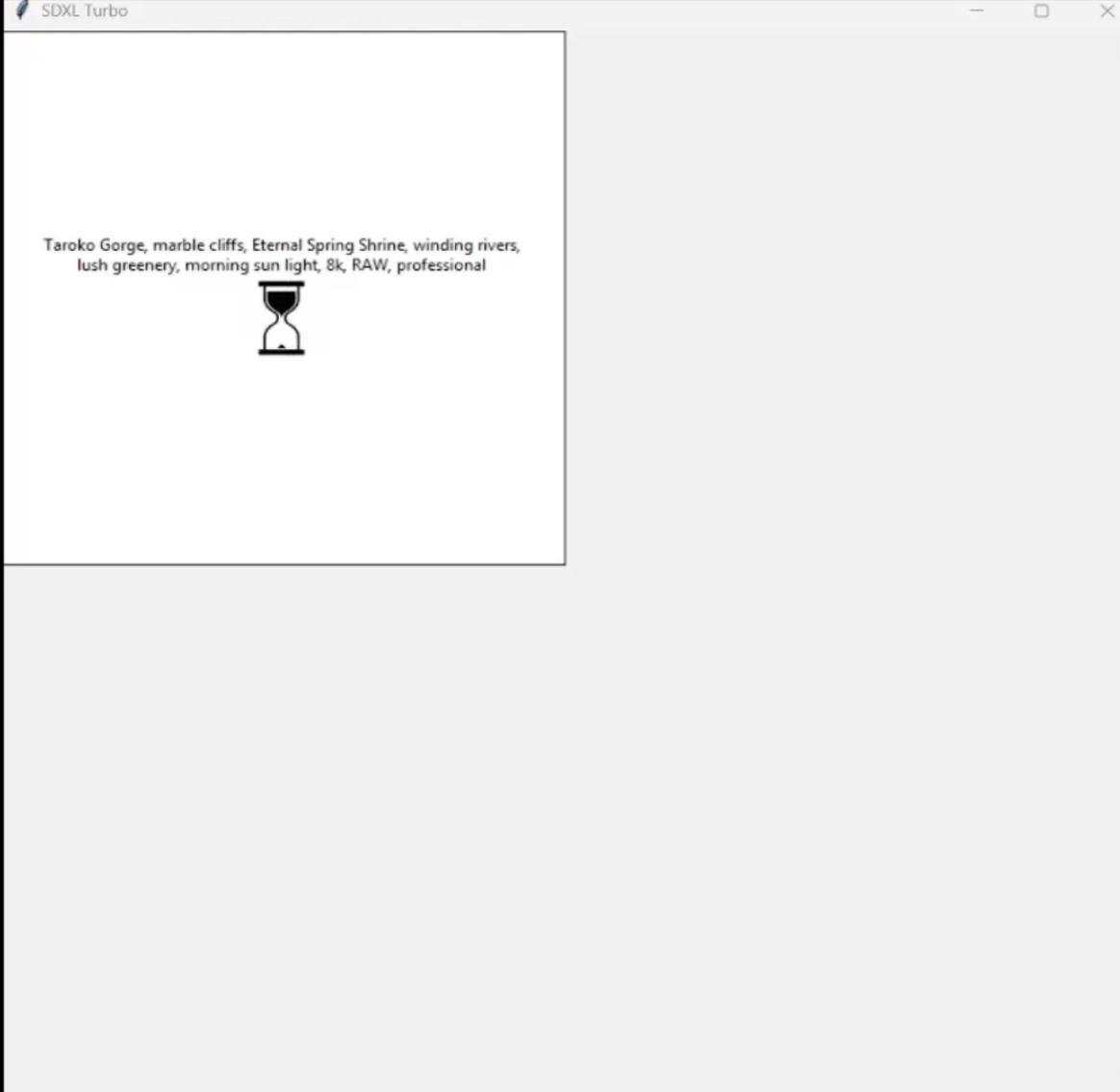
# 310 Foundation:
# AI revolutionizes biomachine Programmability.

**Generative AI for Molecule Programming Accelerated by an AMD EPYC™ processor and AMD Instinct™ GPU based cloud infrastructure**

AMD
together we advance_

**Text-to-image generation excels in photorealism, processing complex prompts, and generating outstanding results**

**Advanced by AMD Ryzen™ AI technology**

stability.ai

30

AMD
together we advance_

**radmantis**

Kria™ SOM - FPGA based SOC platforms: Accelerating AI in
real-time at low power in high constraint environments

AMD
together we advance_

# Search for Elusive Particles at Europe's Large Hadron Collider

**CERN use case: Needle-in-a-haystack search at extreme scale**

- ~ 2.4B collisions/sec producing Pb/s of raw data Real-time machine learning need data analysis

**10s thousands of FPGAs 100m underground,**

- Each inferencing Tb/s of data with nsec latency

**With highly customized dataflow AI accelerators in FPGAs at data-center scale**

**Extreme throughput and nanosecond latency with highly customized dataflow AI accelerators in FPGAs at data center scale.**

AMD
together we advance_

AMD is delivering full stack performant and power efficient solutions with partners and the ecosystem to customers from cloud to endpoints and embedded

# AMD

## together we advance_

# Endnotes

- GD-83: Use of third-party marks/logos/products is for informational purposes only and no endorsement of or by AMD is intended or implied.

- GD-220a: Ryzen™ AI is defined as the combination of a dedicated AI engine, AMD Radeon™ graphics engine, and Ryzen processor cores that enable AI capabilities. OEM and ISV enablement is required, and certain AI features may not yet be optimized for Ryzen AI processors. Ryzen AI is compatible with all AMD Ryzen 7040 series processors except the Ryzen 5 7540U and Ryzen 3 7440U. Please check with your system manufacturer for feature availability prior to purchase.

**AMD**
together we advance_

# Disclaimer

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions and typographical errors.

The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION.

AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY DIRECT, INDIRECT, SPECIAL OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

AMD
together we advance_