

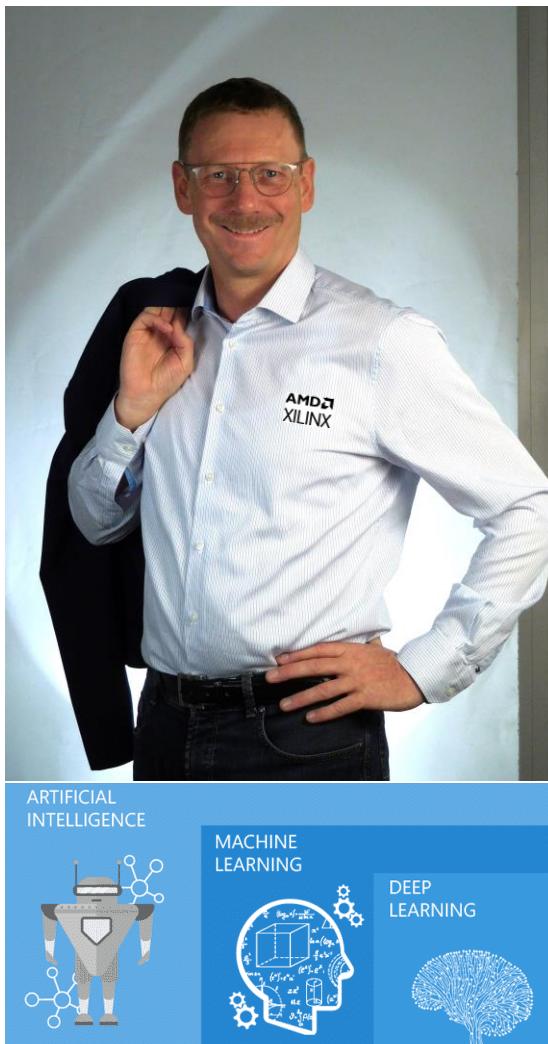


Adaptive Computing, a Disruptive Technology - AMD Unified AI Stack

Jens Stapelfeldt – Sr. DBM
Data Center AI & Compute Markets

Jens.Stapelfeldt@amd.com
[LinkedIn: www.linkedin.com/in/JensStapelfeldt](https://www.linkedin.com/in/JensStapelfeldt)

Jens Stapelfeldt TSL – EMEA – AMD / Xilinx



Jens technical background:

- 1996 /1998 - 3 Y ASIC designer  
- 1998 - 12 Y Business Manager CE and Trainer for Doulos in CE and ARM ATC WW
 - Design methodology (Verification, SystemC,..)
 - (V)HDL, FPGA, ARM Architecture (ARM7 - Cortex-M to Cortex-A)
- 2010 - 5 Y Sr. Embedded FAE at TI 
- Sitara, DaVinci, OMAP, Keystone II
- Industrial App. (IoT, Industrial Ethernet (EtherCAT, Profinet,..), Camera Vison ..)
- Since Oct. 2016 Tech Sales Lead for AMD/Xilinx in EMEA
 - Supporting DACH, EE, Russia, Israel, India
 - 2019 BDM Data Center
 - TSL DACH, Nordic UK & Ireland
 - 2022 Sr. BDM Data Center AI & Compute Markets
- 2016-2018 Part time MBA in Bremen with weeks in Hong Kong, Cambridge, Dublin Business school!
- Since Oct. 2018 guest lecturer for “Digital Signal Processing (VHDL / FPGA design)”.

MBA: In May 2018 I finished my MBA with Master Thesis in International Marketing looking into about 600 AI/ML Start-ups in EMEA!

[LinkedIn Article: https://www.linkedin.com/pulse/analysing-europe-s-aiml-start-up-landscape-using-jens-stapelfeldt/](https://www.linkedin.com/pulse/analysing-europe-s-aiml-start-up-landscape-using-jens-stapelfeldt/)



Jens.Stapelfeldt@amd.com
LinkedIn: www.linkedin.com/in/JensStapelfeldt



High Performance and Adaptive Computing



Cloud, Network,
Hyperscale &
Supercomputing



5G & Comms
Infrastructure



AI & Analytics
Everywhere



Adaptive
Intelligent Systems



Gaming, Simulation
and Visualization



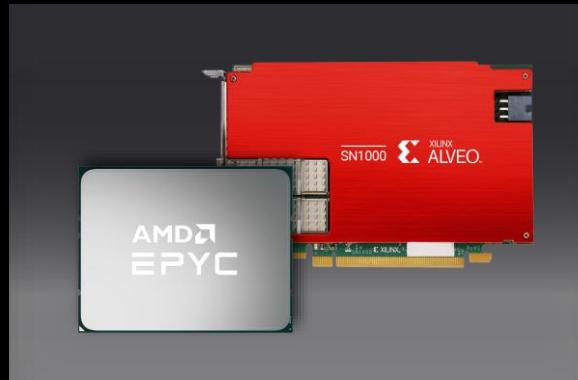
Smarter Client
Devices & Edge

At The Center of Today's Intelligent World

AMD + XILINX AI OPPORTUNITIES



AI Inference
and Training



Data Center and
Communications



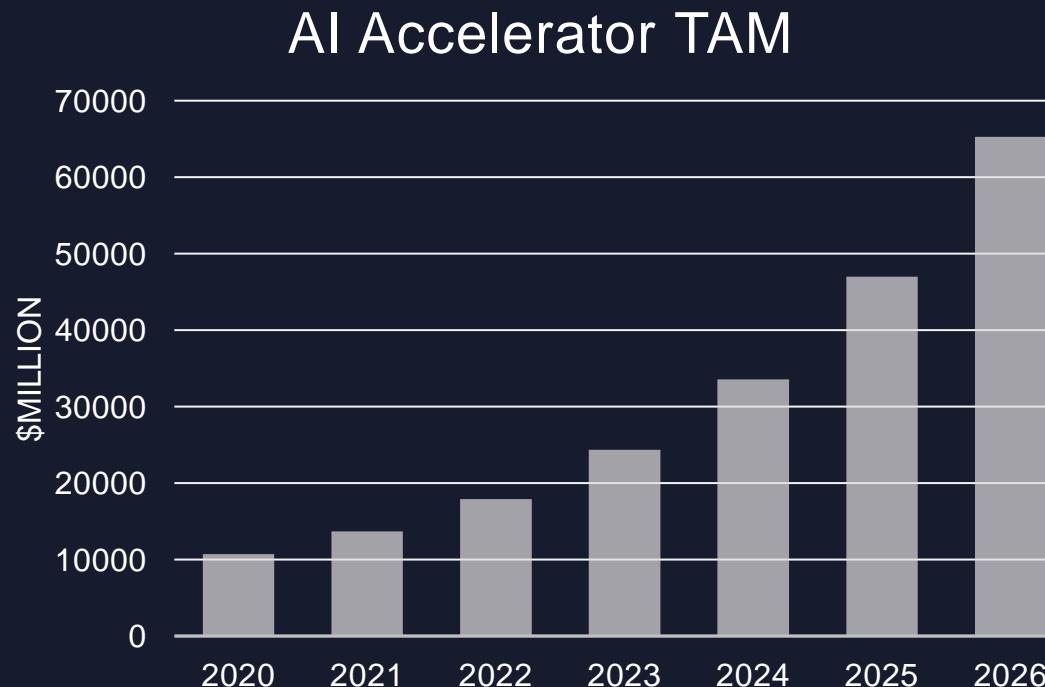
Automotive



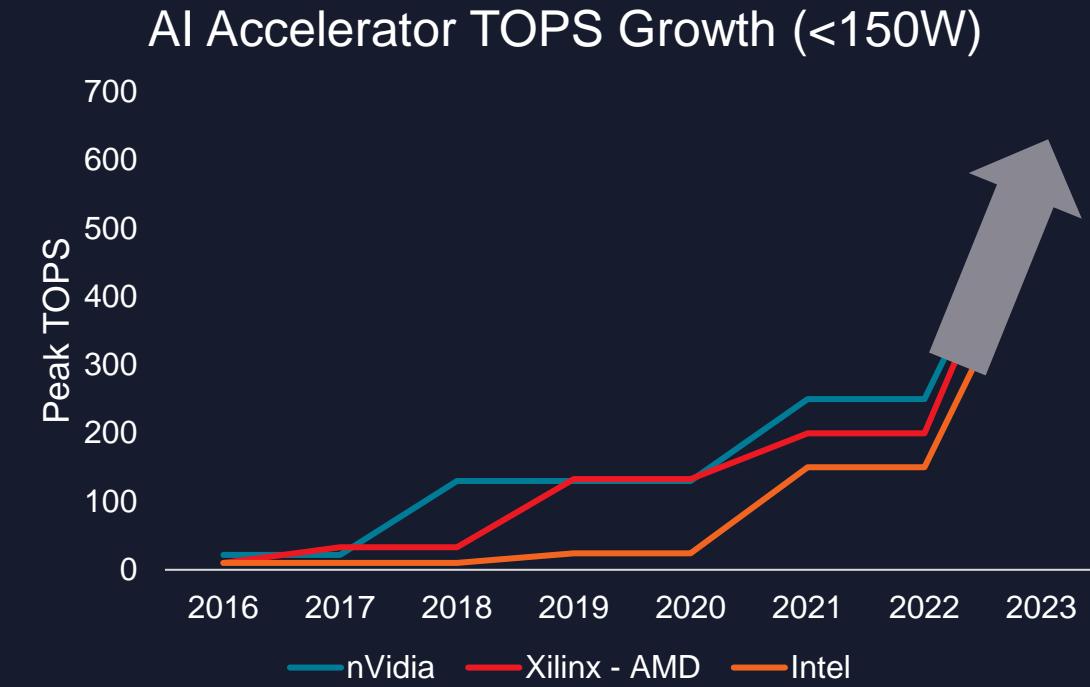
Embedded

- AI OPPORTUNITIES in all areas

AI Accelerator Demands Major Innovation



DC AI Capex Growing at 36.7% CAGR, projected to \$65B TAM by 2026



AI Accelerator Peak TOPS Growing Exponentially to Keep Up with the Model Innovation

5 Source: Markets and Markets (Data Center Accelerator Market Global Forecast to 2026)

*2022-2023 is projected based on history
** Only including AI cards with 150W or less power consumption



PERVASIVE AI



PERVASIVE AI



Commercial & Enterprise

EPYC™, Ryzen™ PRO CPUs

- Server
- PC
- Workstation



Cloud Data Center

EPYC™ CPUs, AMD Instinct™ GPUs

- AI Training
- AI Inference



Digital Home

Ryzen™, Radeon™ CPUs

- PCs and Consoles
- Metaverse

PERVASIVE AI



Cloud Data Center

EPYC™ CPUs, AMD Instinct™ GPUs,
Alveo™ Accelerators, Kintex® FPGAs

- AI Training
- AI Inference
- Security
- Automatic Speech Recognition

Healthcare & Life Sciences

Zynq® Adaptive SoCs

- Imaging and Diagnosis
- Surgical Robotics

Transportation

Zynq® Adaptive SoCs

- ADAS
- Autonomous Vehicles

Commercial & Enterprise

EPYC™, Ryzen™ PRO CPUs

- Server
- PC
- Workstation
- Edge

Smart Retail

Zynq® Adaptive SoCs

- Cashier-less payment
- Customer Reidentification
- Inventory Management

Digital Home

Ryzen™, Radeon™ GPUs,
Zynq® Adaptive SoCs

- PCs and Consoles
- Metaverse
- Smart Home

Intelligent Factory

Zynq® Adaptive SoCs

- Machine Vision
- Predictive Maintenance
- AI Robotics

Smart City

Zynq® Adaptive SoCs

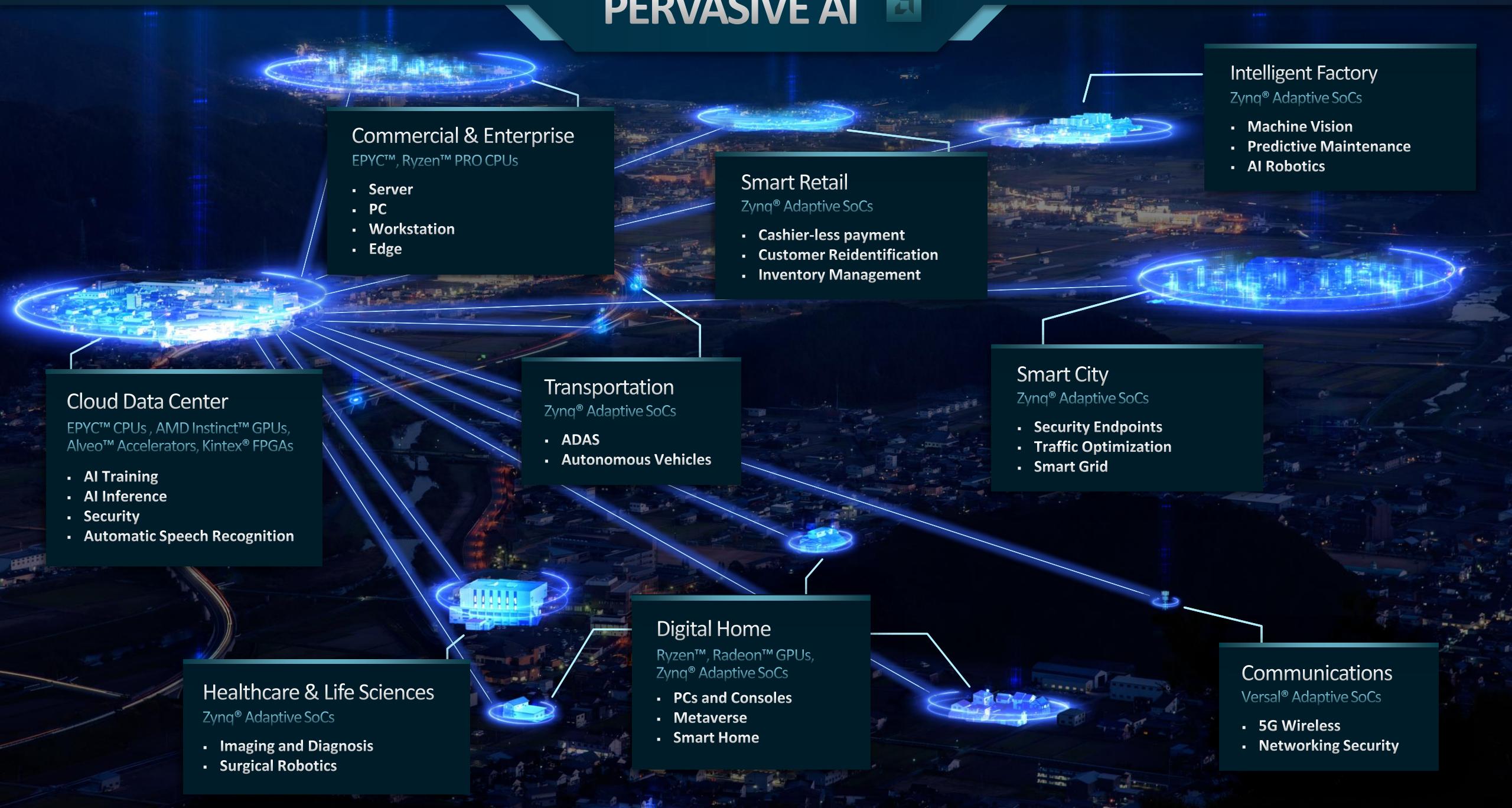
- Security Endpoints
- Traffic Optimization
- Smart Grid

Communications

Versal® Adaptive SoCs

- 5G Wireless
- Networking Security

PERVASIVE AI



PERVASIVE AI STRATEGY



— Leadership AI IP with AMD CDNA™ and scalable AI Engine architecture

— Broaden AMD AI product portfolio across cloud, edge, and endpoint applications

— Unified AI stack to empower developers across the AMD portfolio

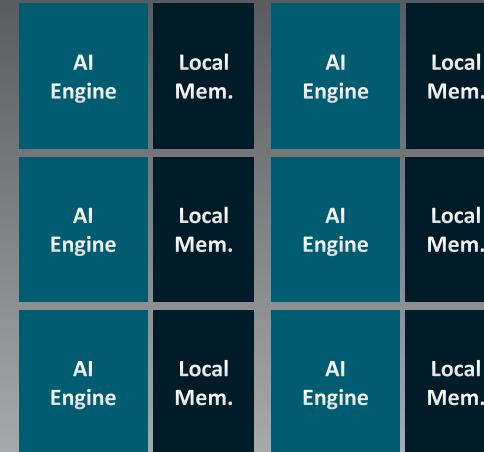
AMD XDNA: ADAPTIVE ARCHITECTURE IP



- Dataflow architecture optimal for AI and signal processing applications
- Highly-scalable array of engines with local memory and data movers
- Leverages deep expertise of compiling algorithms to FPGAs and adaptive SoCs

AI Engine
(AIE)

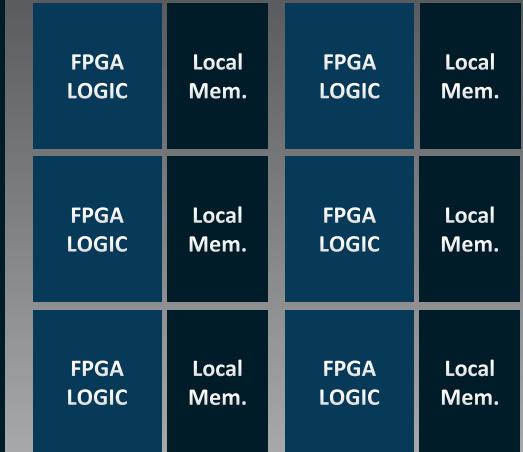
Adaptive Interconnect



High-Performance and Energy Efficiency
for AI and Signal Processing

FPGA Fabric
(FPGA)

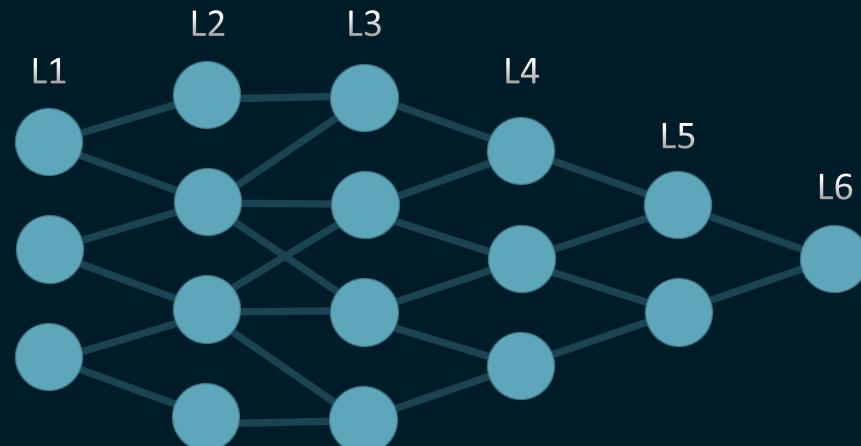
Adaptive Interconnect



Leading FPGA for broad set
of applications and AI

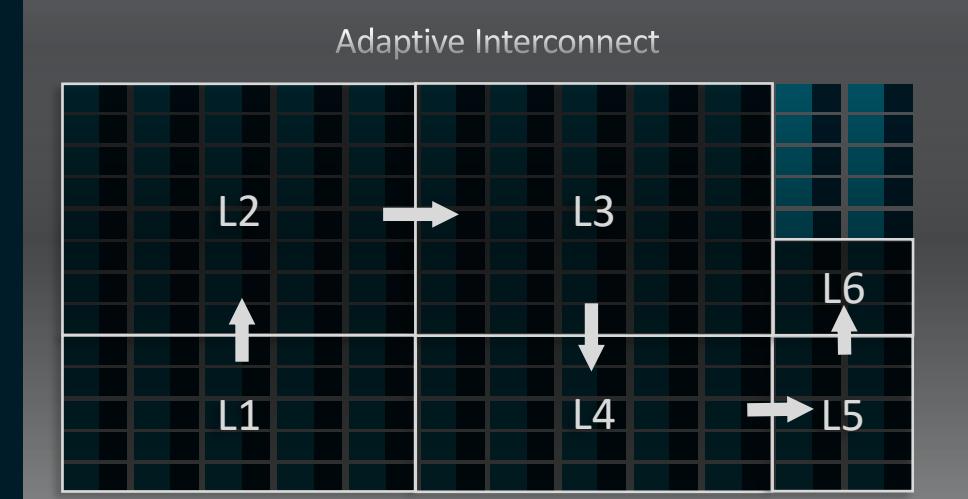
AMD AI ENGINE: ADAPTIVE DATAFLOW PROCESSOR

Deep Neural Network (DNN)



Data “flows” from layer to layer, connections between layers are often “sparse”

DNN Runs Optimally on AIE



Dataflow architecture, sparsity, efficient datatypes deliver high performance and low power

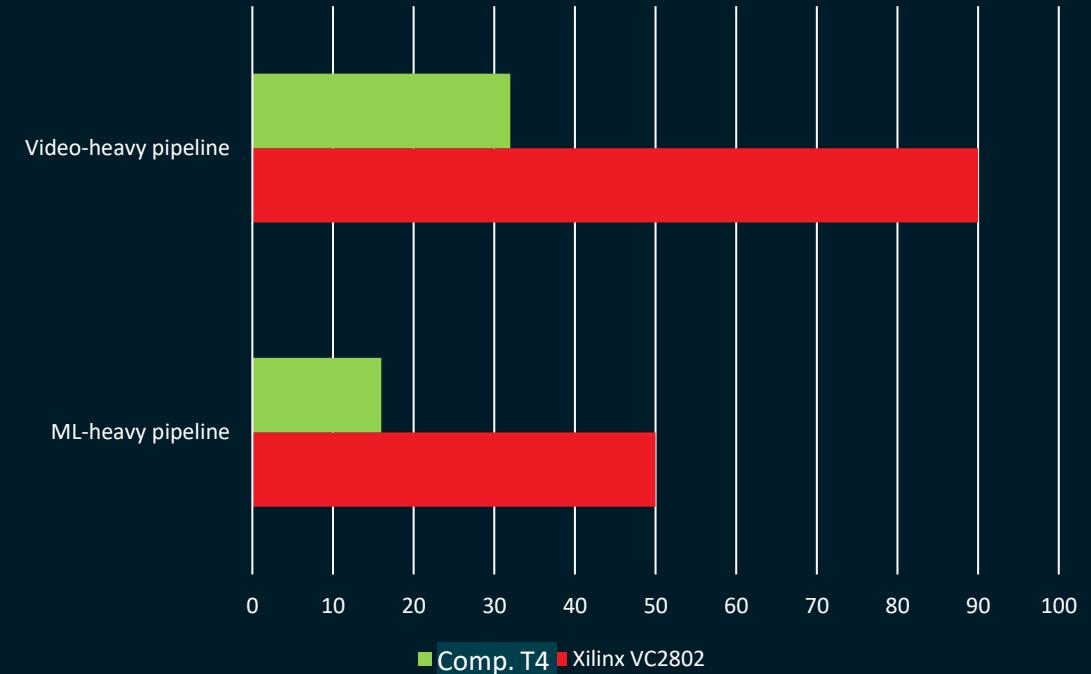
High-Performance, Energy Efficient, and Customizable for AI Workloads

Example 1: Video Analytics

ML-Heavy: H.264 Decode + CV + Yolov3 + 9 x RN18
Video-Heavy: H.264 Decode + CV + tinyYolov3 + 9 x RN50



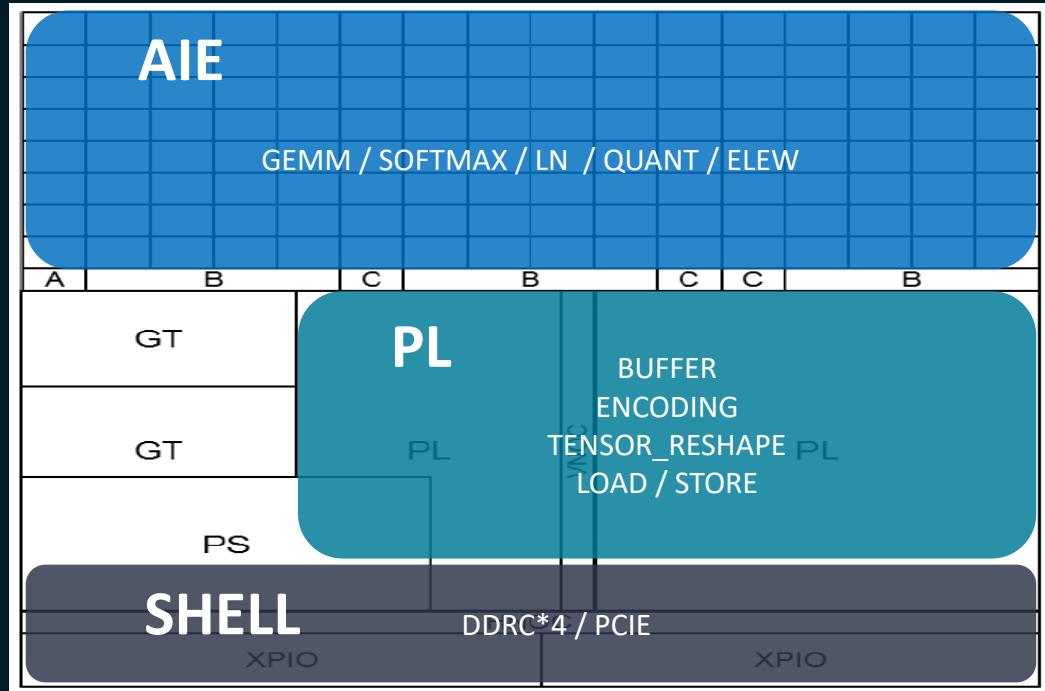
of FHD Streams @ 10 FPS



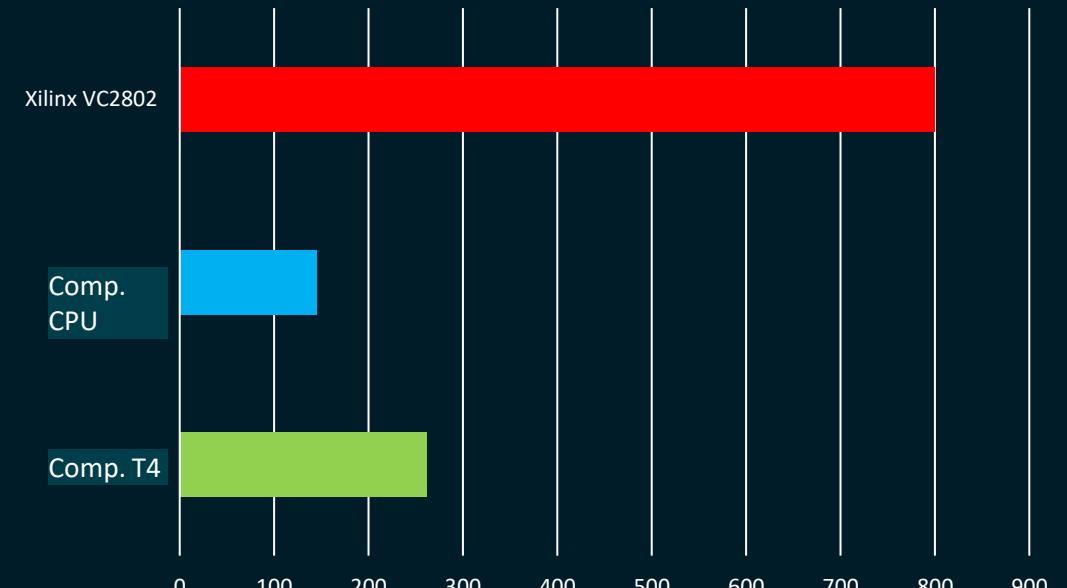
~ 3x competitor in End-to-End Throughput

Example 2: Natural Language

Transformer-based XDNA Design



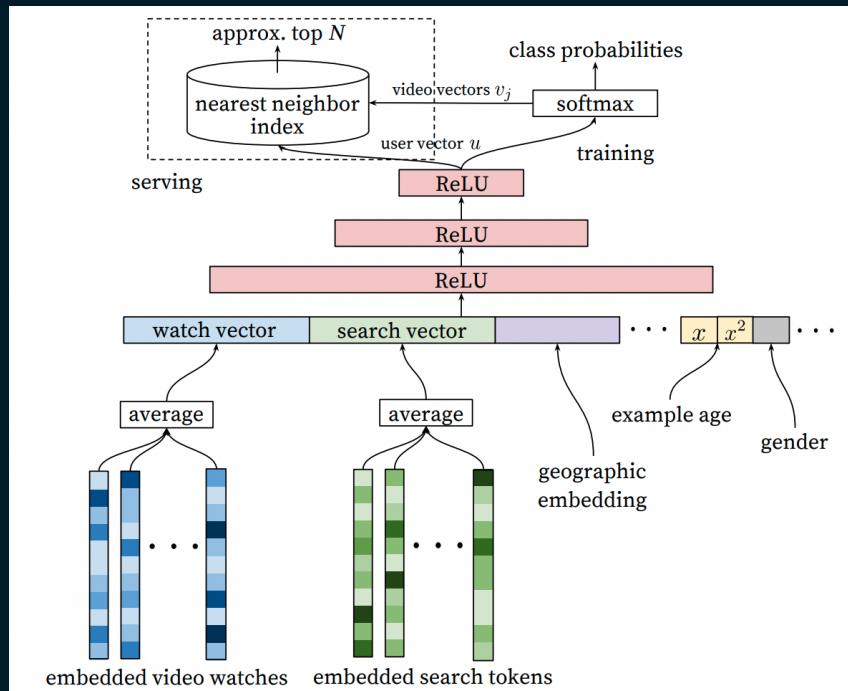
BERT-large (Sentences / sec)



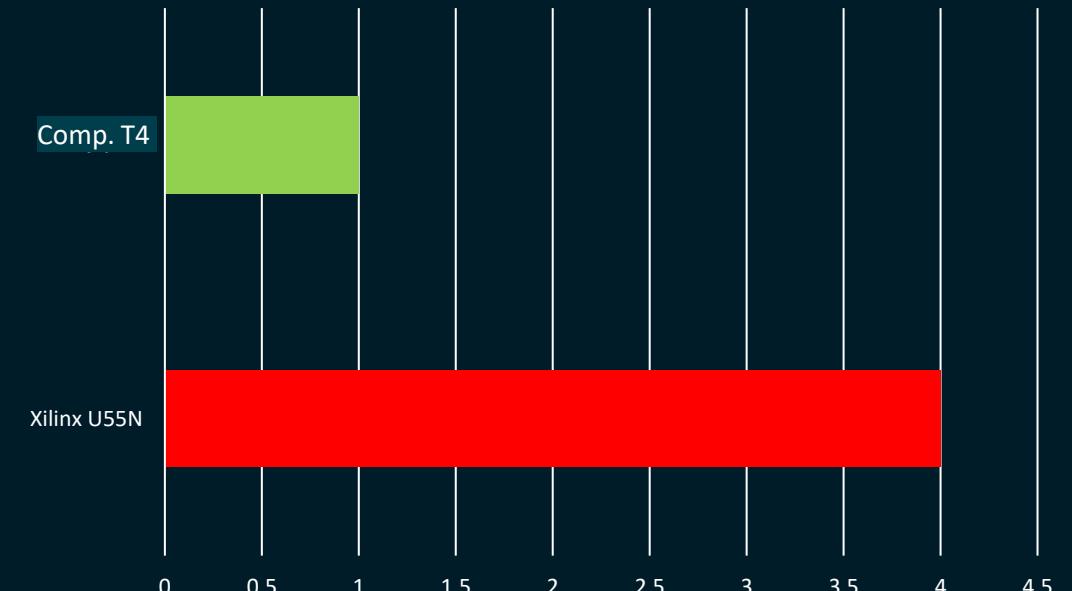
Flexible Architecture for GEMM and Non-Linear Operations

Example 3: Recommendation AI

HBM-based Pipeline to Hide Compute Latency

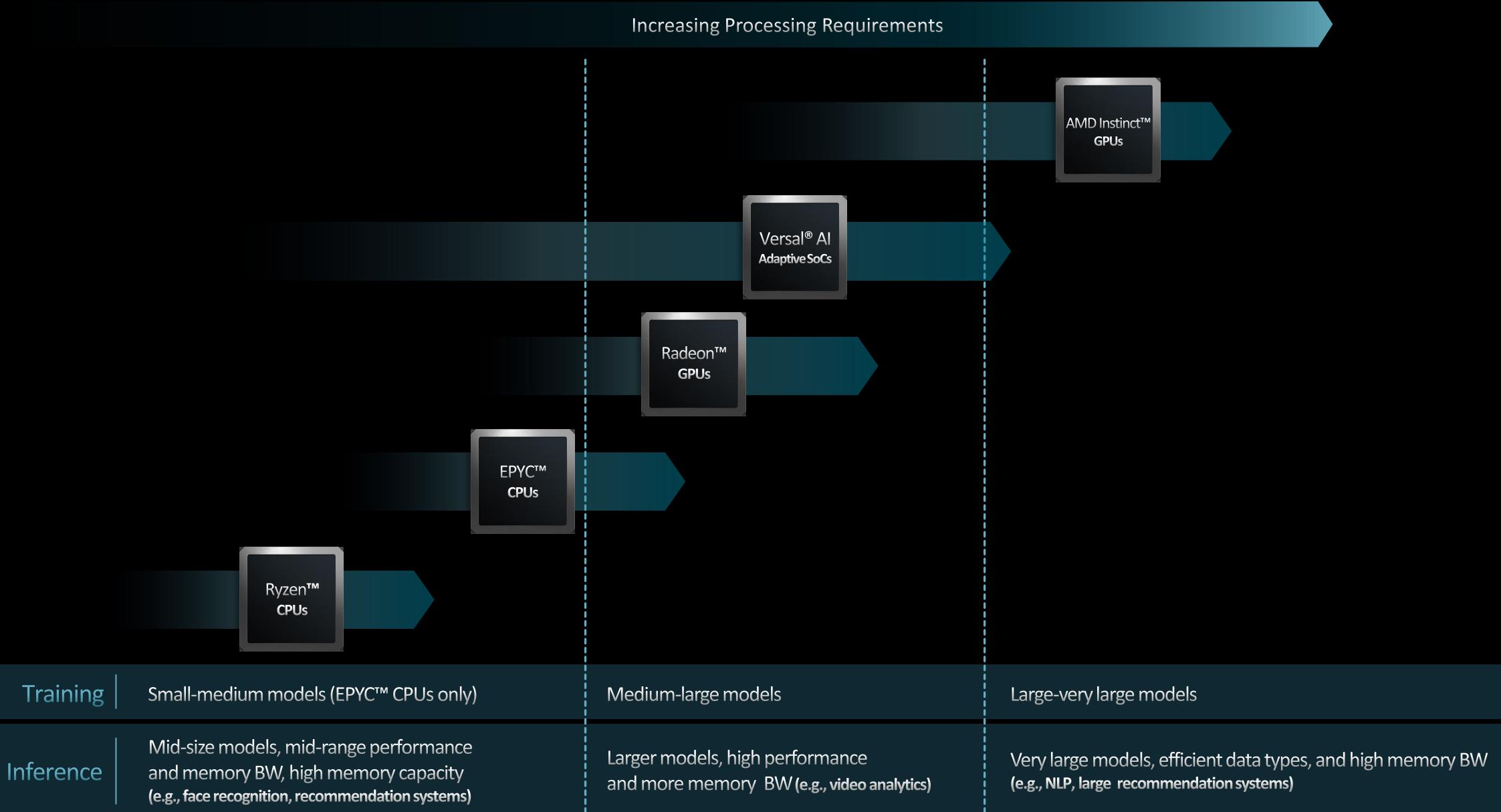


Vector Retrieval (Speed-up factor)

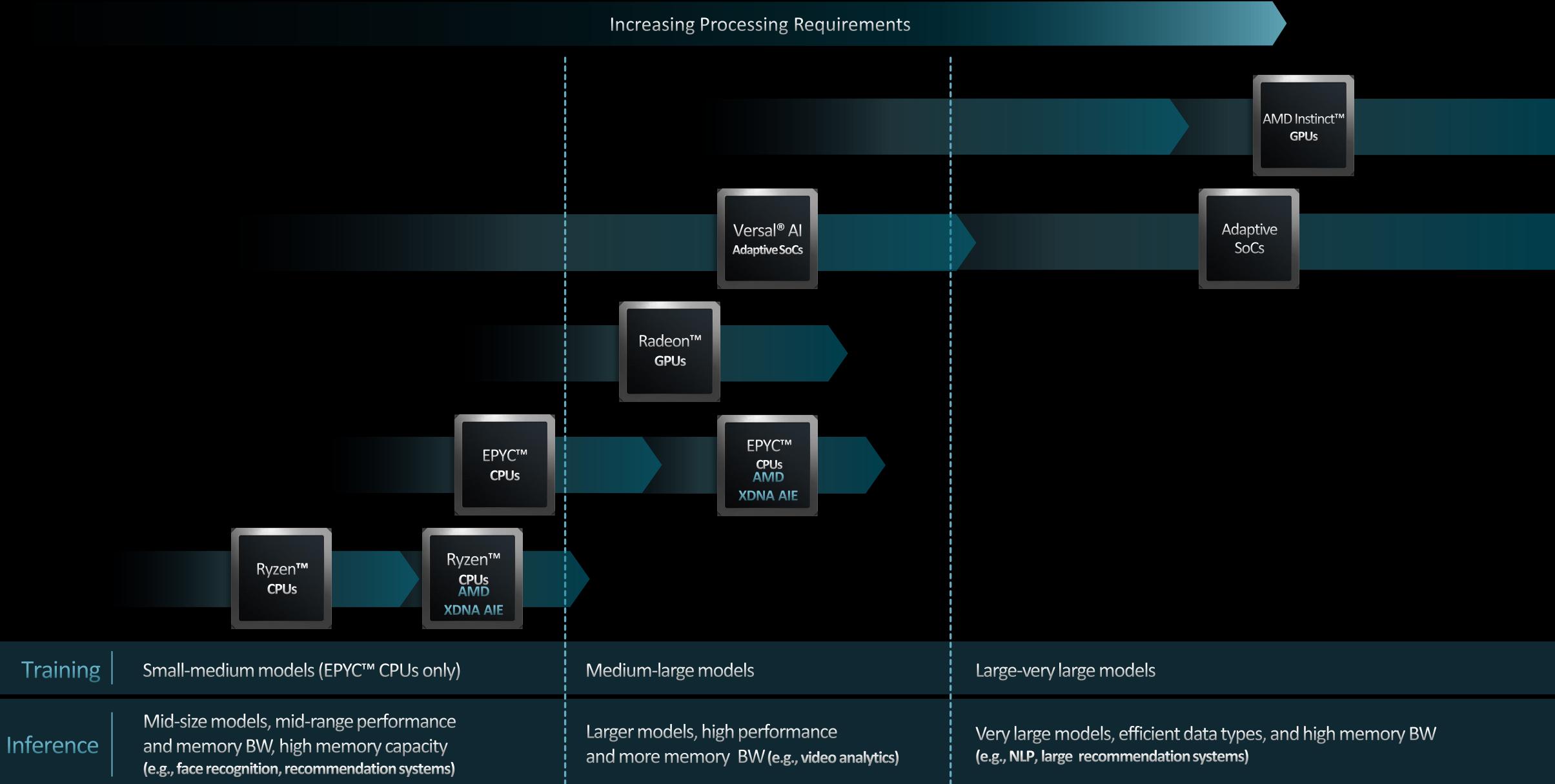


4x competitor in End-to-End Throughput at 10 ms Latency

AI APPLICATION COVERAGE



AI APPLICATION COVERAGE



AMD AI SOFTWARE TODAY

CPU Stack

Optimized Inference Models

WinML
ONNX Runtime

Pytorch

TensorFlow

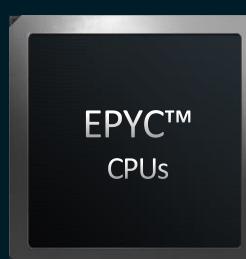
ML Graph Compiler

ZenDNN, AOCL Optimized Library

ZEN Studio (AOCC)

Windows Runtime

Linux Runtime



ROCM™ Platform

Optimized Inference Models

ONNX Runtime

Pytorch

TensorFlow

MIGraphX: AI Development Tools

MIGraphX: ML Graph Compiler

MIOpen, ROCBLAS

ROCM™ HIP Compiler and tools

Runtime



Vitis™ AI Platform

Optimized Inference Models

ONNX Runtime

Pytorch

TensorFlow

Vitis AI Development & Deployment Tools

Vitis AI ML Graph Compiler

Vitis ML Libraries

Vitis™ SW Platform (AIE Compiler, HLS and tools)

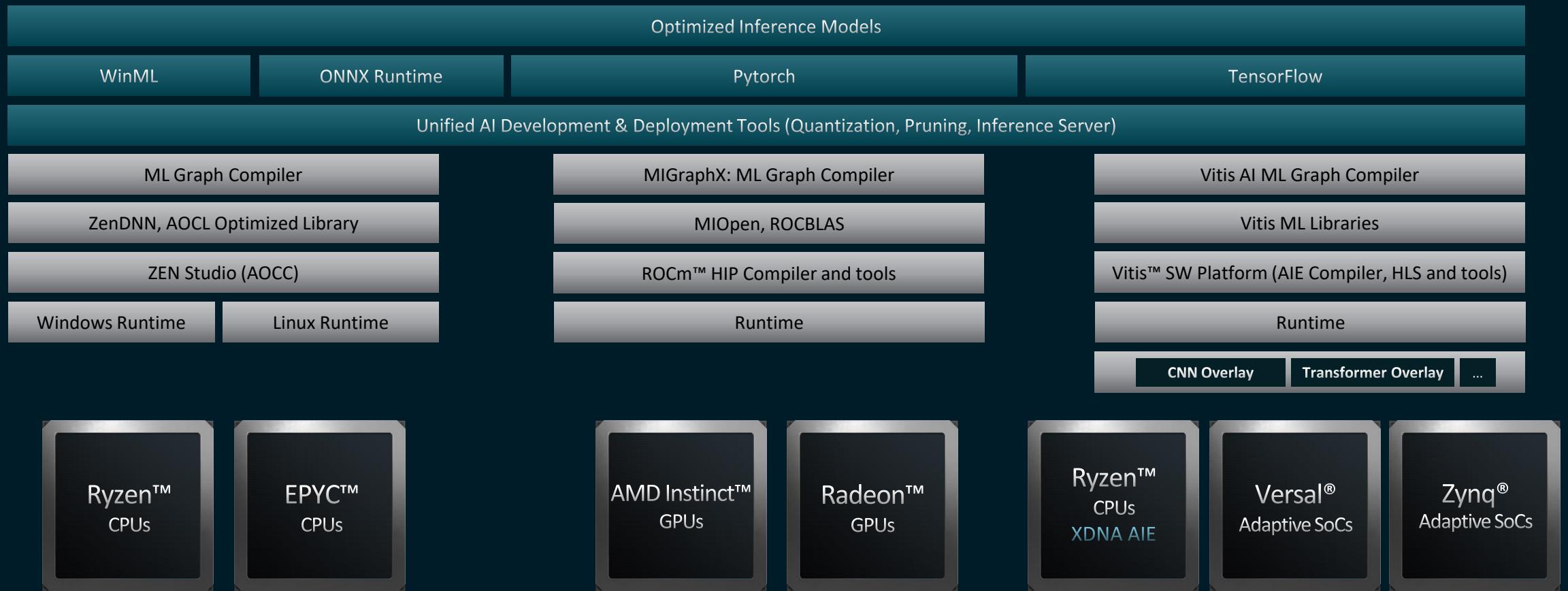
Runtime

CNN Overlay Transformer Overlay ...



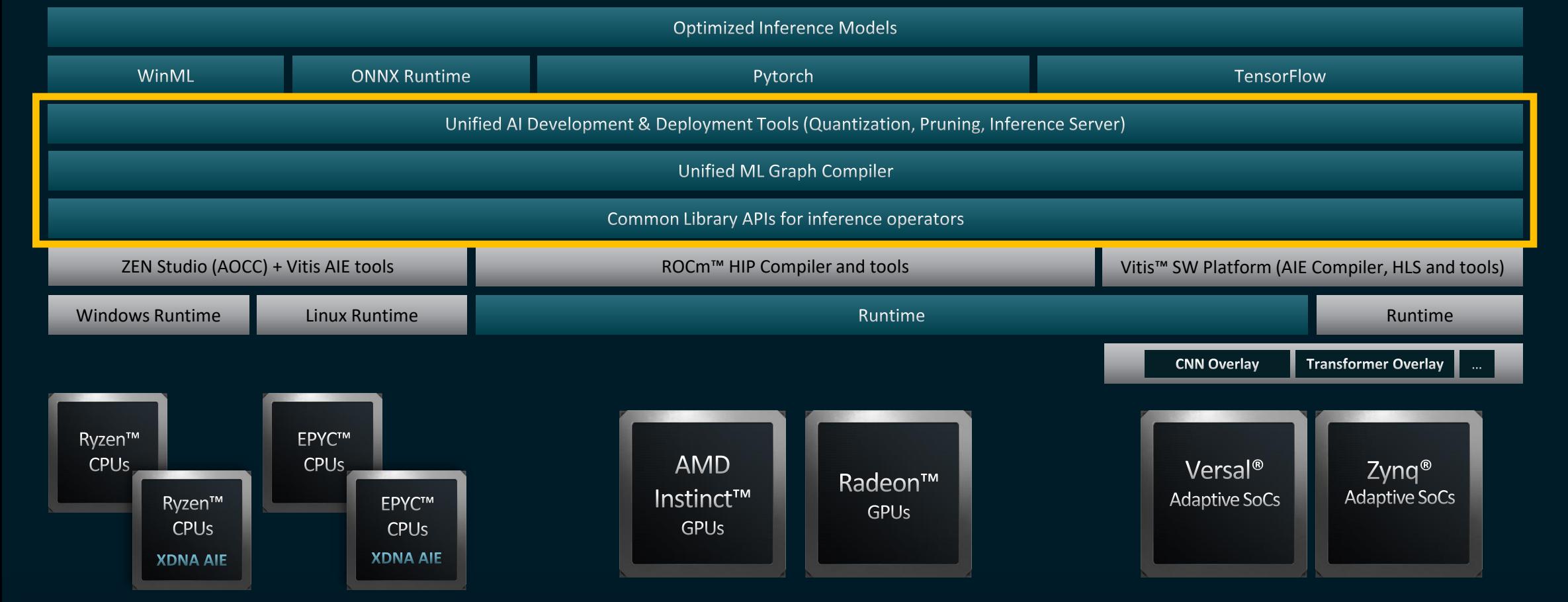
AMD UNIFIED AI STACK 1.0

Unified Inference Frontend (UIF) for AI Developers



AMD UNIFIED AI STACK 2.0

Seamless Workload Partitioning with UIF, Graph Compiler and Library APIs



AMD

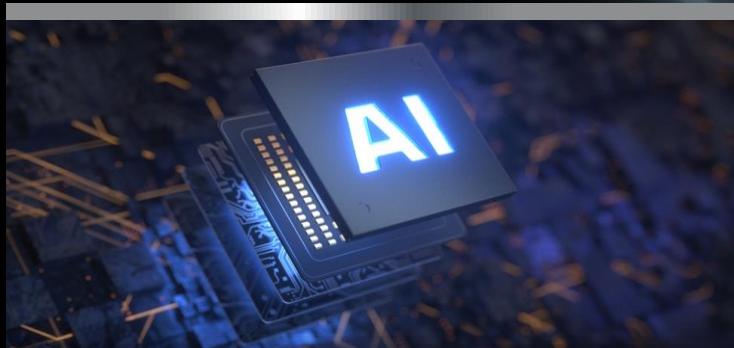
together we advance

AI INNOVATION



Large Model Training

Exponential growth in generative model size



Sustainable Inference

OPEX and CAPEX optimized deployment



Data Science First

Common environment for CPU, GPU and FPGA

SUMMARY



TAM Expansion

AI, Data Center, Automotive
and Embedded



Adaptive Computing

XDNA Architecture: Best of
ASICs and FPGAs



Pervasive AI

Unified AMD Hardware and
Software Roadmap