

Modern Compute Architecture: From CPU to GPU



Norwegian research infrastructure services

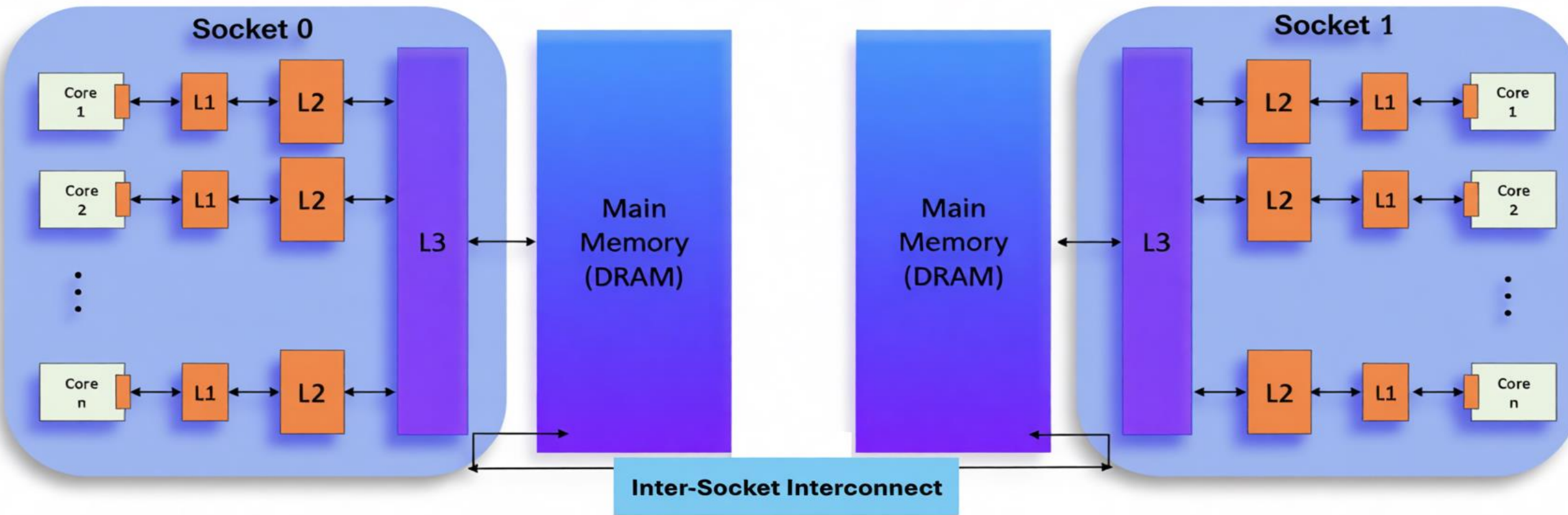
Hicham Agueny, PhD
Scientific Computing Group
Info. Tech. Department, UiB/NRIS



Outline

- **Why GPUs matter ?**
- **CPU & GPU Architecture Overview**
- **Evolution of GPUs**
- **Architecture of Modern GPUs: NVIDIA & AMD GPUs**
- **Software ecosystem of GPUs: CUDA & ROCm**

Traditional CPU architecture



Adapted from https://indico.cern.ch/event/814979/contributions/3401193/attachments/1831477/3105158/comp_arch_codas_2019.pdf

AMD EPYC 7742, Zen 2 architecture, 2 sockets, 128 cores total.

	L1 cache	L2 cache	L3 cache	
Size	32 KiB per core	256 KiB per core	256 MiB per socket	
Latency	~4 cycles	~12–15 cycles	~35–45 cycles	
Scope	Private per core	Private per core	Shared by 4 cores (per CCX)	



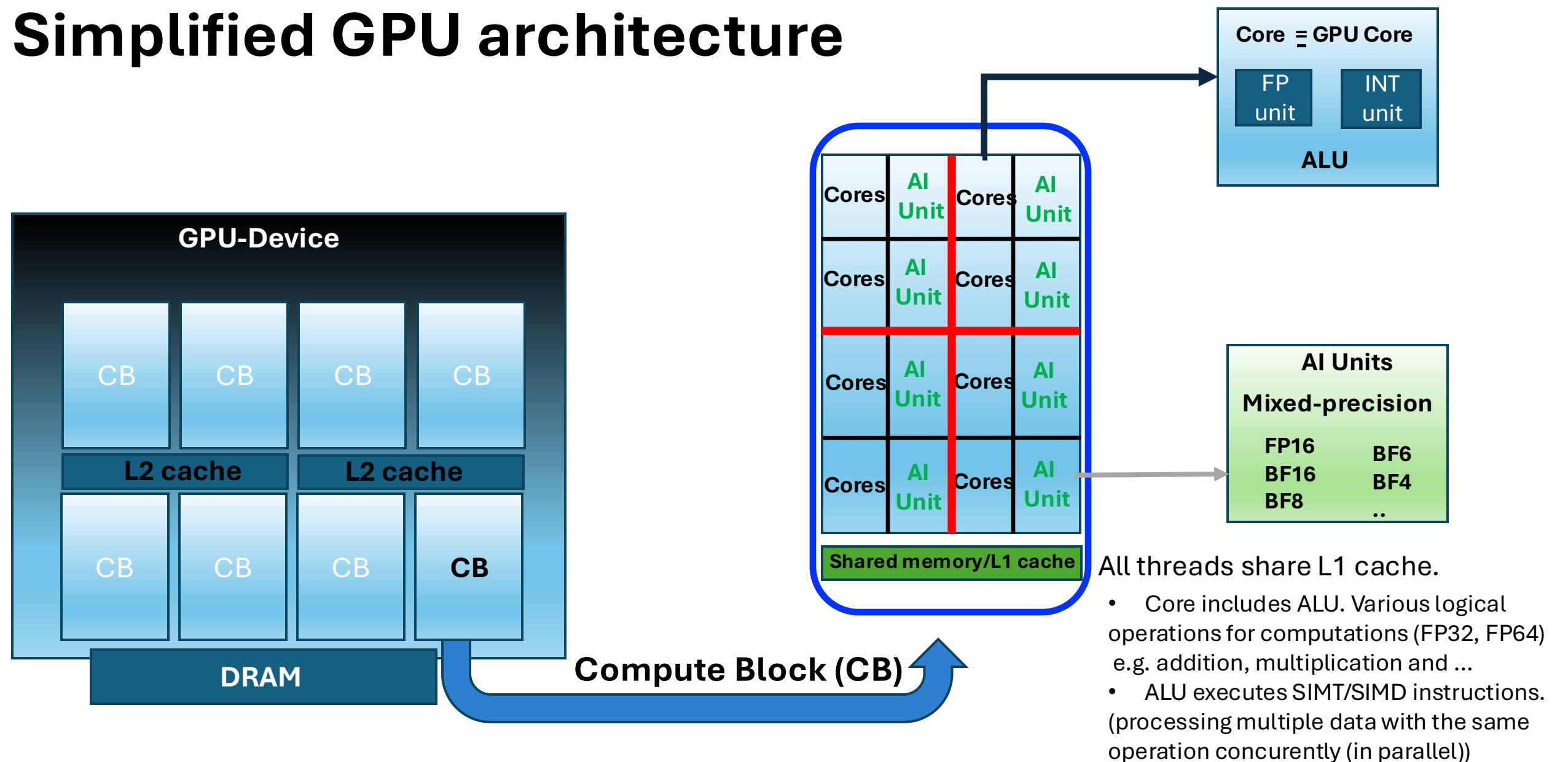
Traditional CPU architecture

CPU's are designed for

- Low latency (use caches L1/L2/L3)
- Responsive (fast change of instruction e.g. “if” block, loops with variable lengths)
- single-threaded performance (running one main task faster, even with irregular steps)

The CPU design becomes less efficient when processing millions of simple operations in parallel.

Simplified GPU architecture



CPU vs GPU

Traditional CPUs built for:

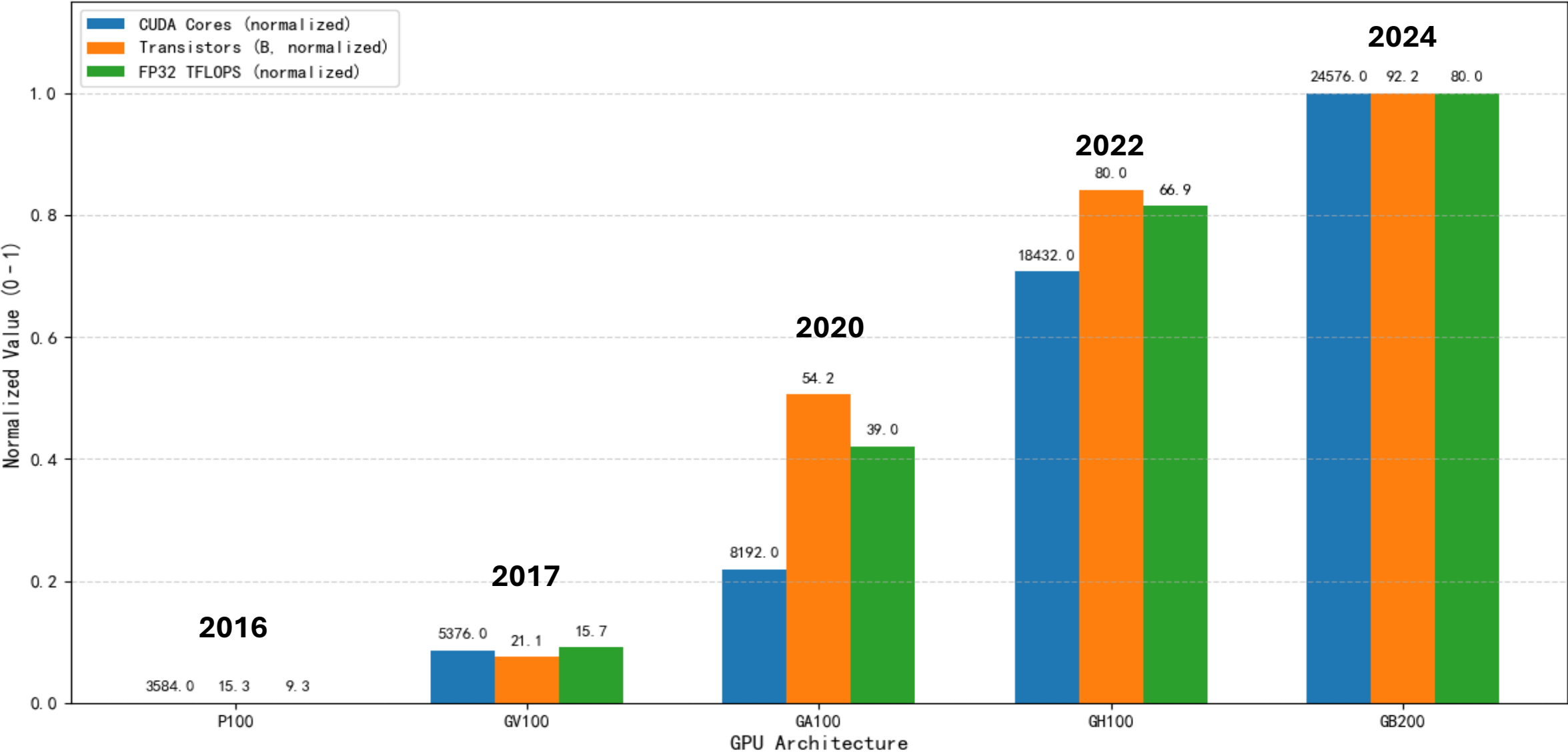
- Low latency
- Complex control flow (if, loops, function calls)
- Single-thread speed
- Not ideal for massive parallelism

Modern GPUs built for:

- High throughput
- Thousands of simple, synchronized threads
- Regular, data-parallel workloads (SIMT/SIMD)
- Not ideal for branching or irregular code

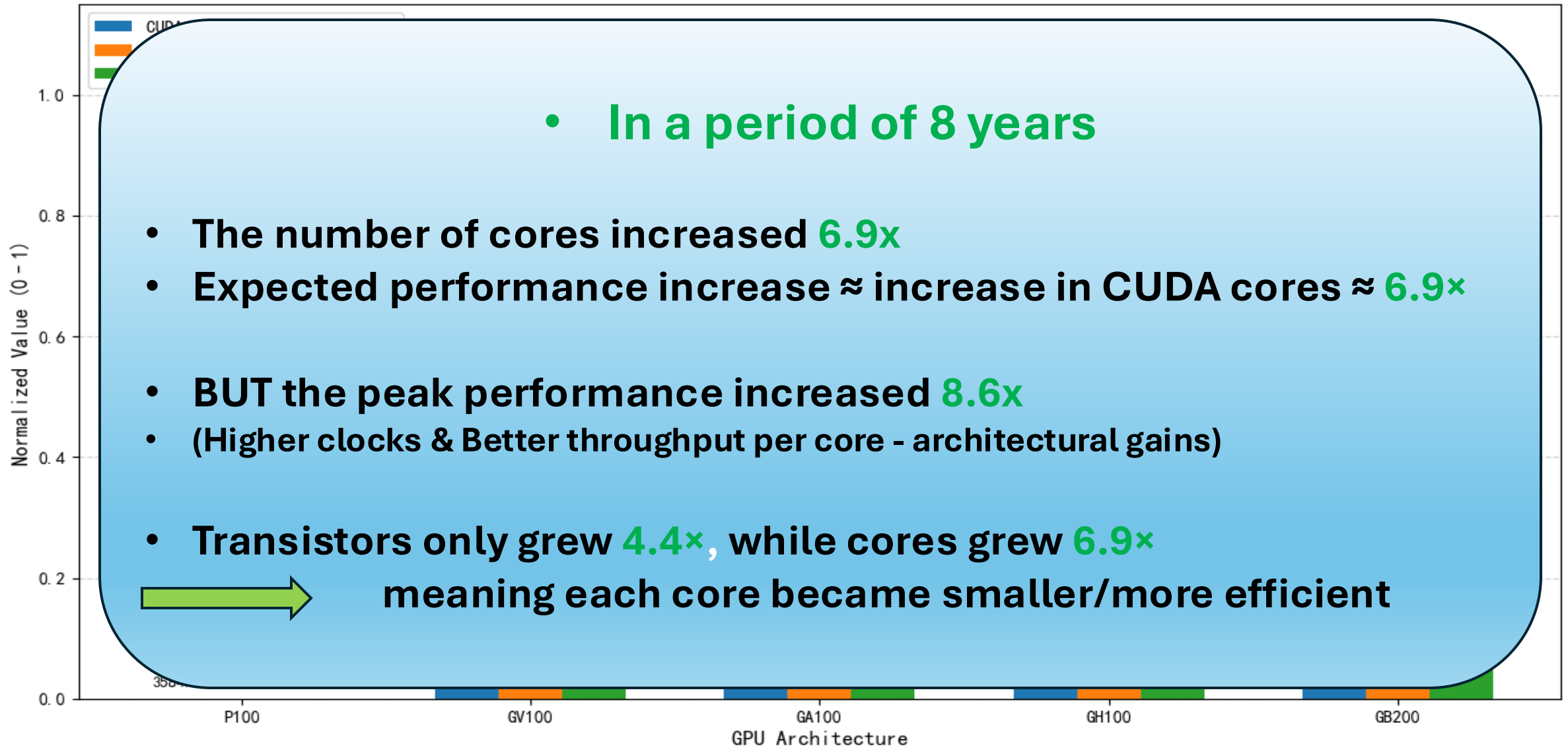
Evolution of GPUs

Normalized Comparison of CUDA Cores, Transistors (B), and FP32 TFLOPS
with Actual Values Displayed



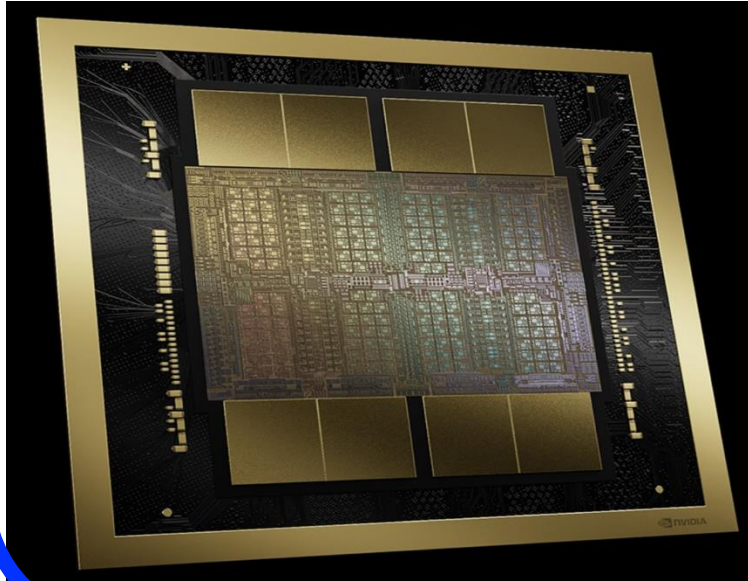
Evolution of GPU

Normalized Comparison of CUDA Cores, Transistors (B), and FP32 TFLOPS
with Actual Values Displayed

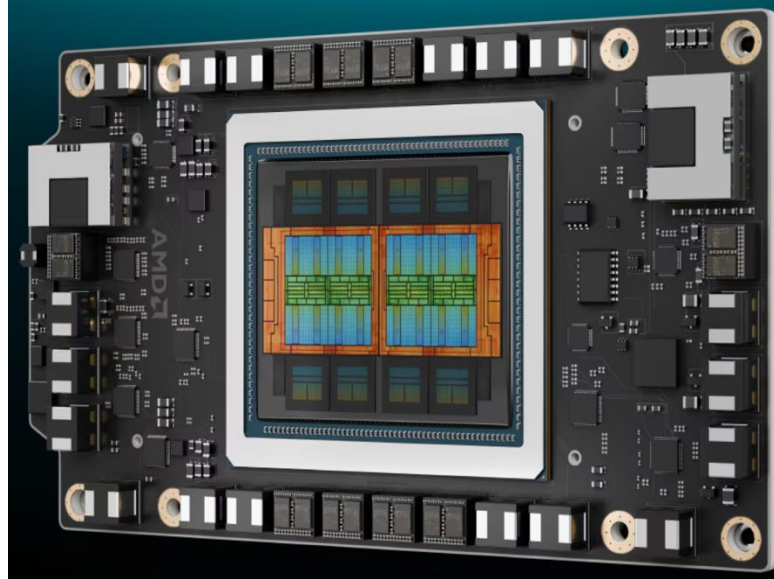


Modern GPU architectures

NVIDIA Blackwell GPU



AMD Instinct MI350 GPUs



Intel® Arc™ Pro B60 Graphics



NVIDIA Blackwell GPU <https://www.nvidia.com/en-us/data-center/technologies/blackwell-architecture/>

AMD MI350 GPU <https://www.amd.com/en/products/accelerators/instinct.html>

Intel Arc B60 GPU <https://www.intel.com/content/www/us/en/products/docs/discrete-gpus/arc/workstations/b-series/overview.html>

Modern GPU architectures

- **Chiplet-Based GPU Design**

Multi-die GPU with high-bandwidth die-to-die (5.5-10 TB/s)

- **Next-Generation Matrix/AI Units**

Enhanced AI Units optimized for mixed-precision: FP16/FP8/BF16

Higher throughput per compute block compared to previous GPUs

- **Scaled memory Bandwidth** (up to 8 TB/s)

- **Higher memory capacity** per GPU (up to 288 GB)

- **Advanced CPU-GPU Interconnects** with high-bandwidth (128-900 GB/s)

- **Advanced GPU-GPU Interconnects** with high-bandwidth (1.0-1.8 TB/s)

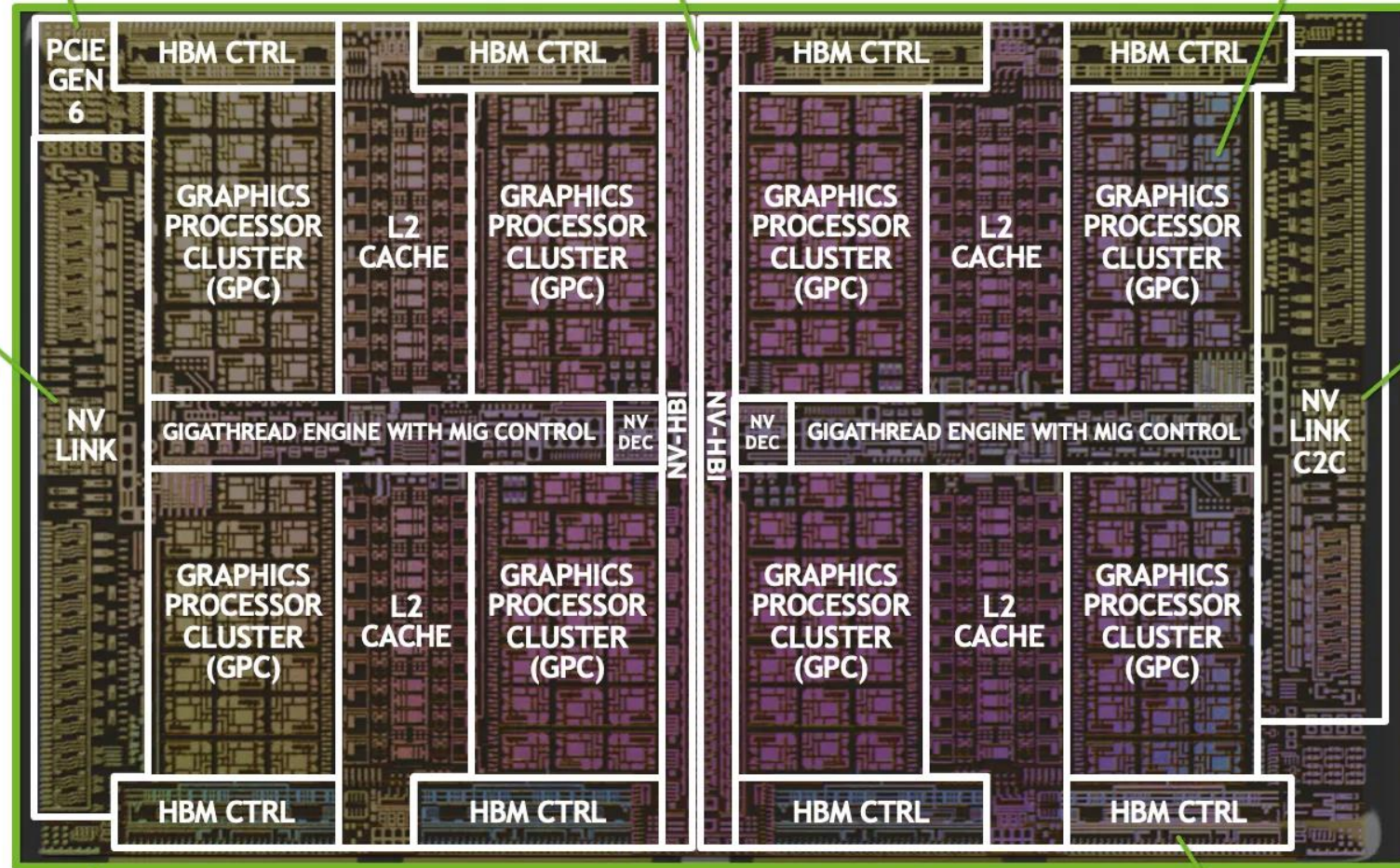
NVIDIA Blackwell Ultra GPU

x16 PCIe Gen 6
256 GB/s CPU Host Interface

High Bandwidth Interface
10TB/s Die-to-Die

160 SMs per GPU: 640 Tensor Cores
15 PetaFLOPS Dense NVFP4

NVLink v5
1,800 GB/s to
NVLink Switch



NVLink-C2C
900GB/s Coherent
CPU-GPU Interface

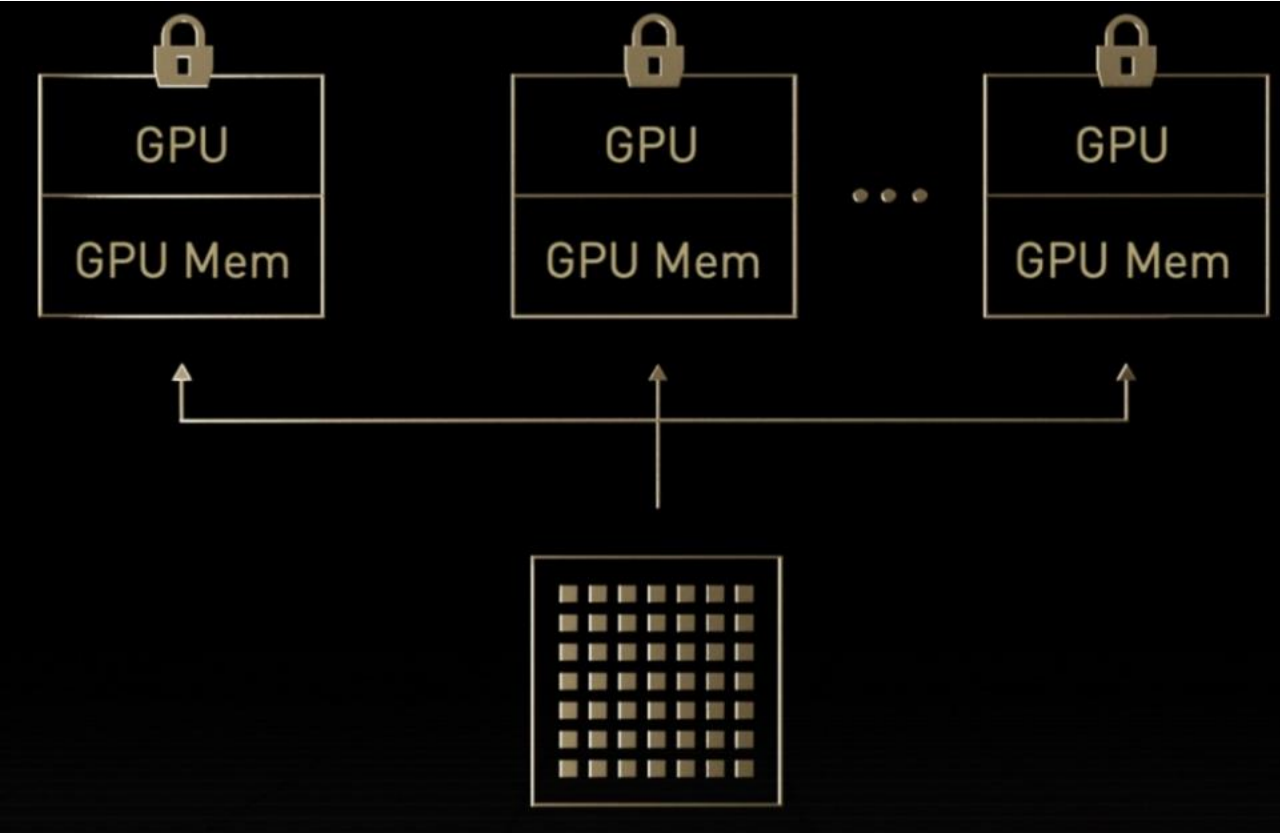
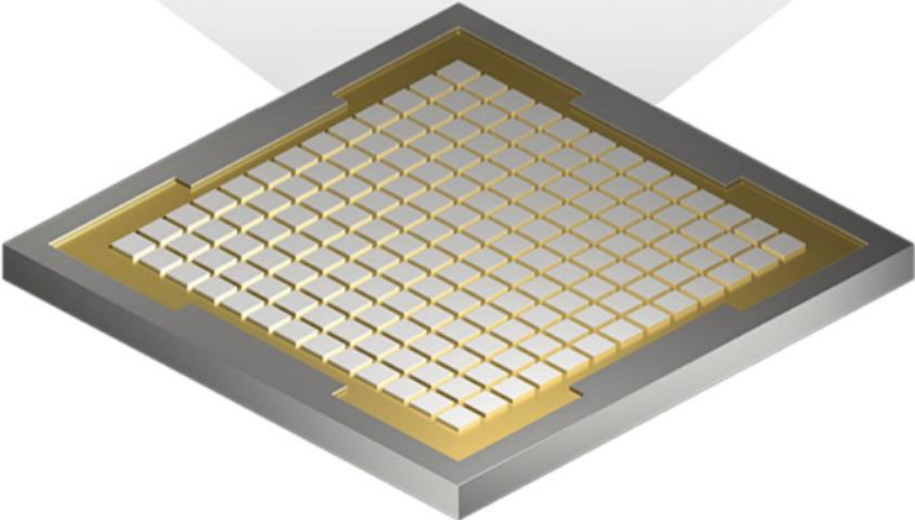
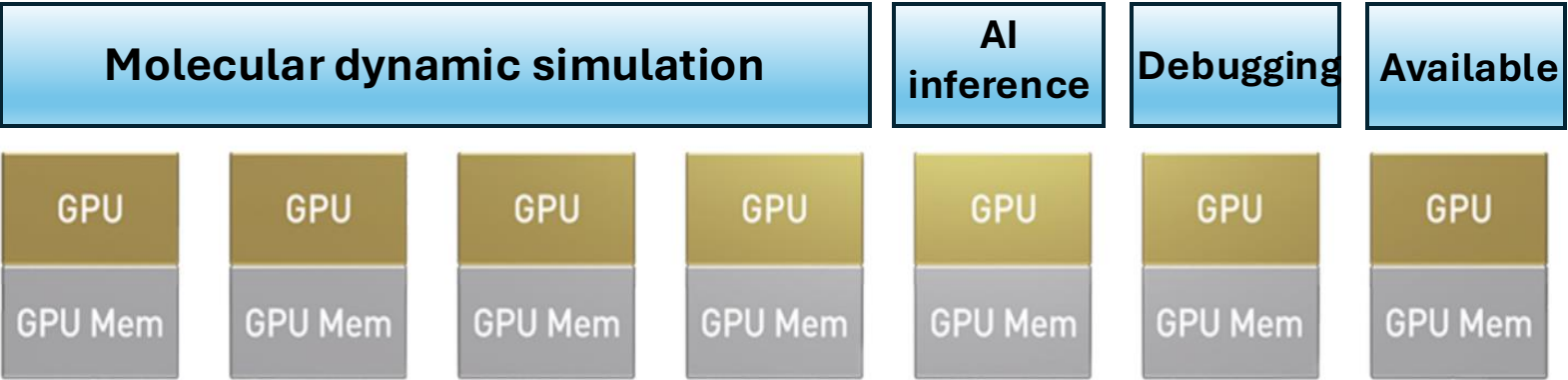
Confidential Computing
TEE-I/O Capable

288GB HBM3E Memory
(12 Stacks, Up to 8 TB/s)

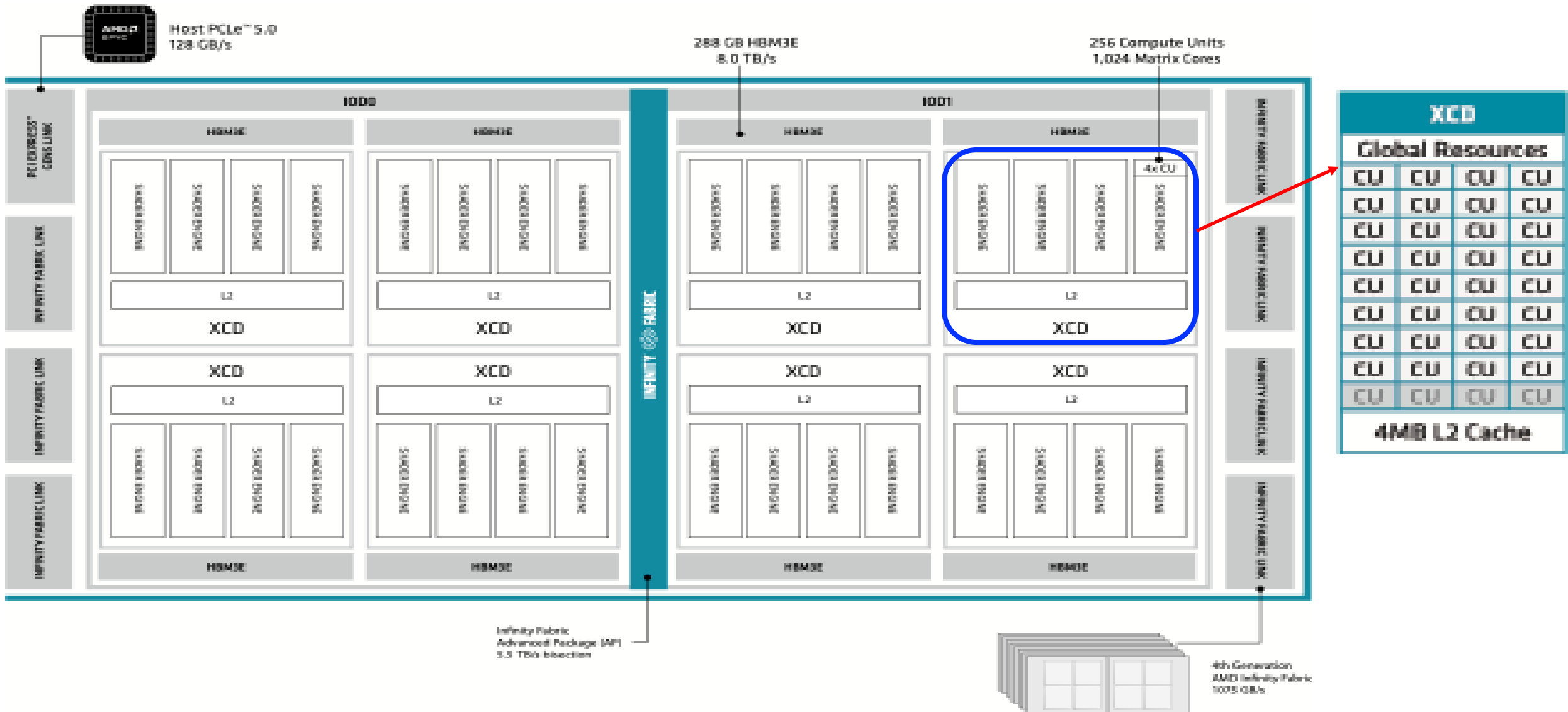
Streaming multiprocessor



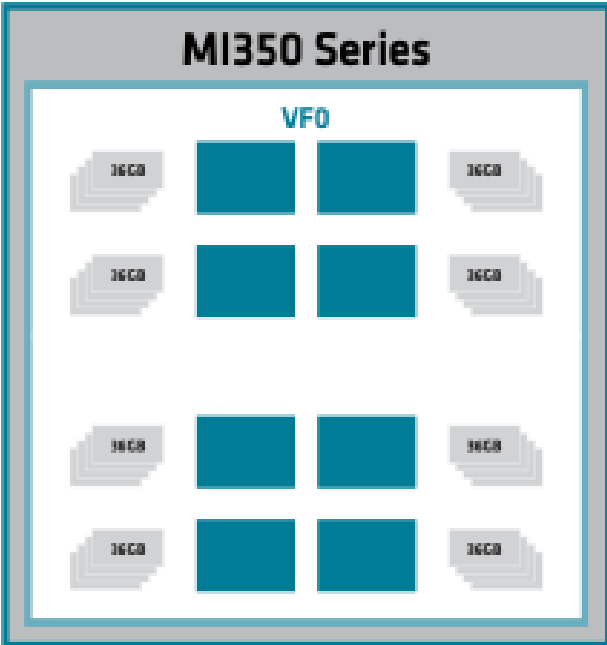
Multi-Instance GPU: NVIDIA GPU



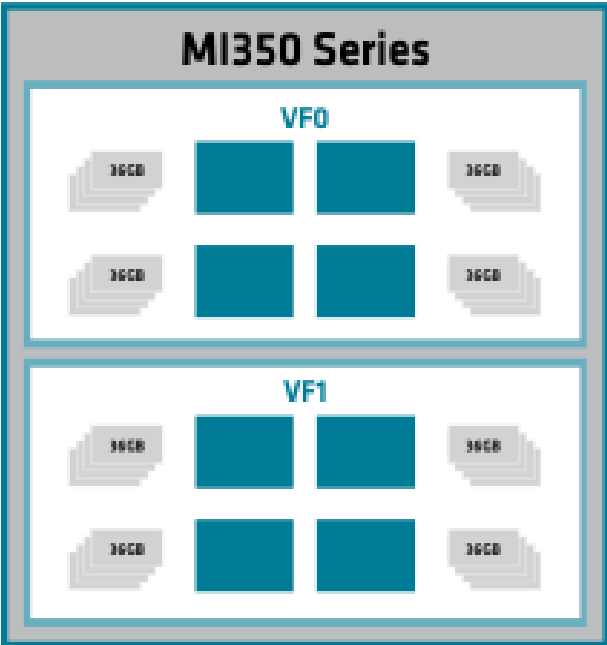
AMD Instinct MI350 Series GPU Multi-Die Chiplet (XCD)



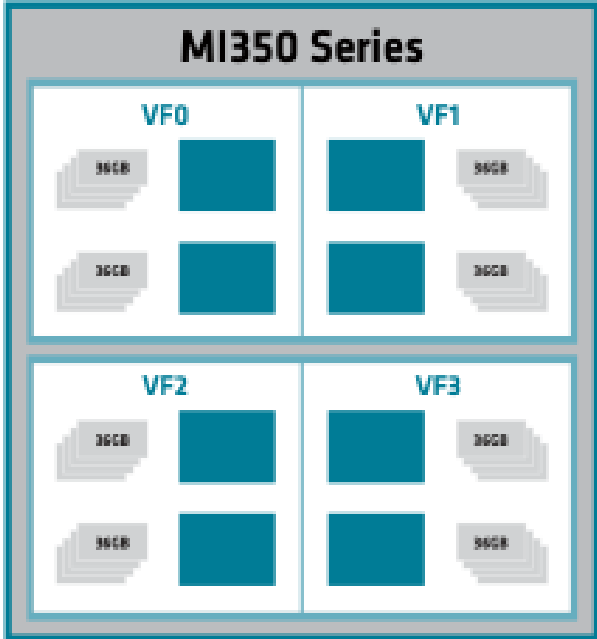
Compute Partitioning (or MIG): AMD GPU



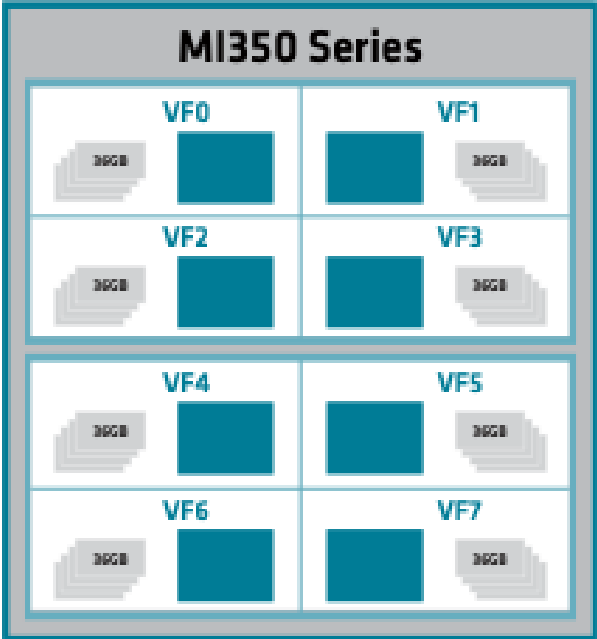
(a) Single Partition: SPX, NPS 1
8 XCDs per Partition / 8 XCDs per VM,
288 GB per Partition / 288 GB per VM



(b) Two Partitions: DPX, NPS 2
4 XCDs per Partition / 4 XCDs per VM,
144 GB per Partition / 144 GB per VM



(c) Four Partitions: QPX, NPS 2
2 XCDs per Partition / 2 XCDs per VM,
72 GB per Partition / 72 GB per VM



(a) Eight Partitions: CPX, NPS 2
1 XCDs per Partition / 1 XCDs per VM,
36 GB per Partition / 36 GB per VM

Key GPU Hardware Unit Terminology Comparison

Concept	NVIDIA	AMD		Explanation
Large GPU Compute Block	GPC (Graphics Processing Cluster)	XCD (Accelerator Compute Dies)		Handle parallel compute and graphics tasks via SMs/CUs.
Compute Block	SM (Streaming Multiprocessor)	CU (Compute Unit)		Contains ALUs, registers, schedulers.
GPU Core	CUDA Core	Stream Processor		Executes FP/INT operations.
AI / Matrix Acceleration Unit	Tensor Core	Matrix Engine		Accelerates matrix / AI computations. (mixed precision)
Thread	Thread	Work item/thread		Smallest independently scheduled unit of execution.
Thread Group Size	Warp = 32 threads	Wavefront = 32 RDNA: 32 or 64		Implement SIMT/SIMD, allowing one instruction to be executed across many threads/work-items at the same time.

Software ecosystem for GPU accelerators

- **CUDA** (created by NVIDIA)
- **ROCm** (created by AMD)

CUDA Software ecosystem

USE-CASES



Speech

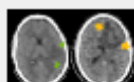


Translate

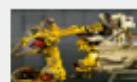


Recommender

CONSUMER INTERNET



Healthcare



Manufacturing

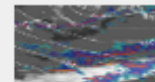


Finance

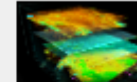
INDUSTRIAL APPLICATIONS



Molecular
Simulations



Weather
Forecasting



Seismic
Mapping

SUPERCOMPUTING

APPS & FRAMEWORKS



Amber
NAMD

+600
Applications



CUDA-X LIBRARIES

MACHINE LEARNING

cuDF

cuML

cuGRAPH

DL / HPC

cuDNN

CUTLASS

TENSORRT

CUDA Math Libraries

LANGUAGES

python OpenACC



LLVM Compiler
For CUDA

CUDA

CUDA TOOLKIT

CUDA
COMPILER

DEVELOPER TOOLS

DEBUGGERS

PROFILERS

CUDA C++
CORE

CUDA DRIVER

MEMORY
MANAGEMENT

WINDOWS &
GRAPHICS

COMMS
LIBRARIES

OS PLATFORMS



CentOS



Windows Server

What is new in CUDA Toolkit: 13.0

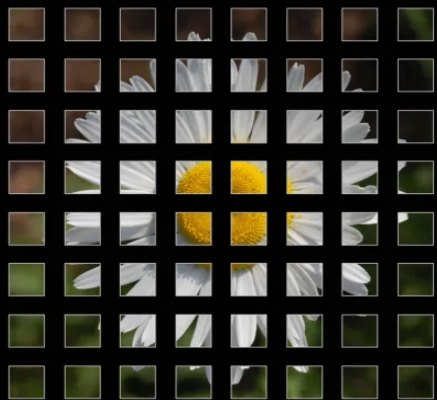
- **Unified support for Arm CPUs**, enabling CUDA on broader platforms (e.g., Grace CPU systems).
- **Traditional CUDA**: Based on thread-parallel **SIMT**, developers manage low-level thread blocks and warps.
- **CUDA 13+**: Introduces **tile-based programming**, a higher-level abstraction for GPU acceleration.
- **Compiler & runtime** now handle thread scheduling and hardware optimization automatically.
- **Developers focus on logic**, not threads, write less repetitive code, achieve full performance.
- **Tiles supported in Python, C++, and more**, accessible across languages.

What is new in CUDA Toolkit: 13.0

Tile-based programming: define tiles (array) of data and specify operations over those tiles.
Compiler handles the low-level execution.

Block-level models

Productivity & performance
for array-based programs



Application maps
data onto blocks

Compiler maps
blocks onto threads

Tile Programming

Block-wide cooperative execution on regular Tiles of data

Data & operation granularity is array / tensor

$$\underline{A} + \underline{B} = \underline{C}$$

Tile tensor addition
(array granularity)

Well-suited to regular / data-parallel problems

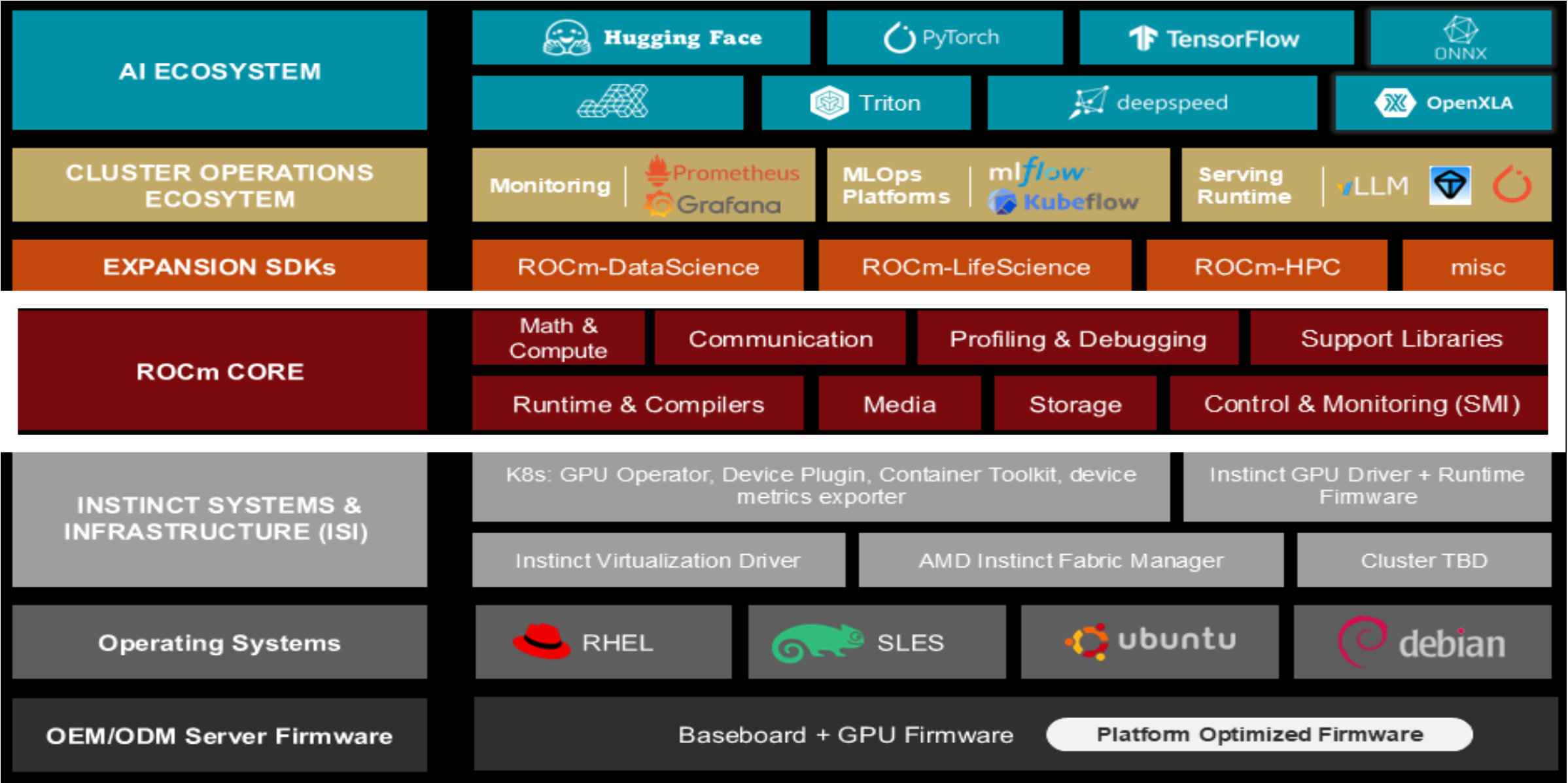
Allows for easy & performant kernel fusion

Simplifies array operations with implicit use of threads

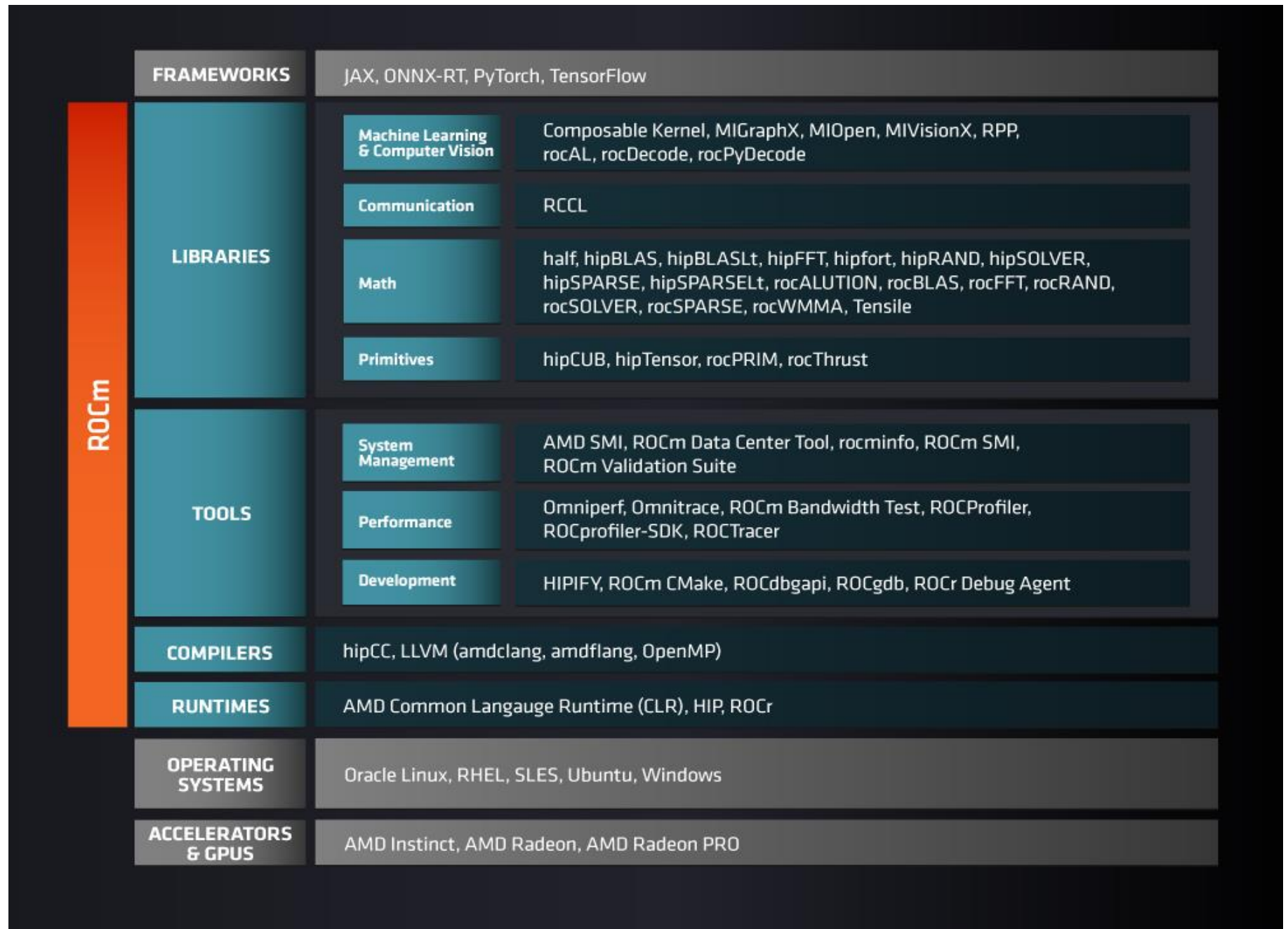
$$\begin{matrix} A_{00}A_{01} \\ A_{10}A_{11} \\ A_{20}A_{21} \\ A_{30}A_{31} \\ A_{40}A_{41} \\ A_{50}A_{51} \\ A_{60}A_{61} \\ A_{70}A_{71} \end{matrix} + \begin{matrix} B_{00}B_{01} \\ B_{10}B_{11} \\ B_{20}B_{21} \\ B_{30}B_{31} \\ B_{40}B_{41} \\ B_{50}B_{51} \\ B_{60}B_{61} \\ B_{70}B_{71} \end{matrix} = \begin{matrix} C_{00}C_{01} \\ C_{10}C_{11} \\ C_{20}C_{21} \\ C_{30}C_{31} \\ C_{40}C_{41} \\ C_{50}C_{51} \\ C_{60}C_{61} \\ C_{70}C_{71} \end{matrix}$$

SIMT tensor addition
(element granularity)

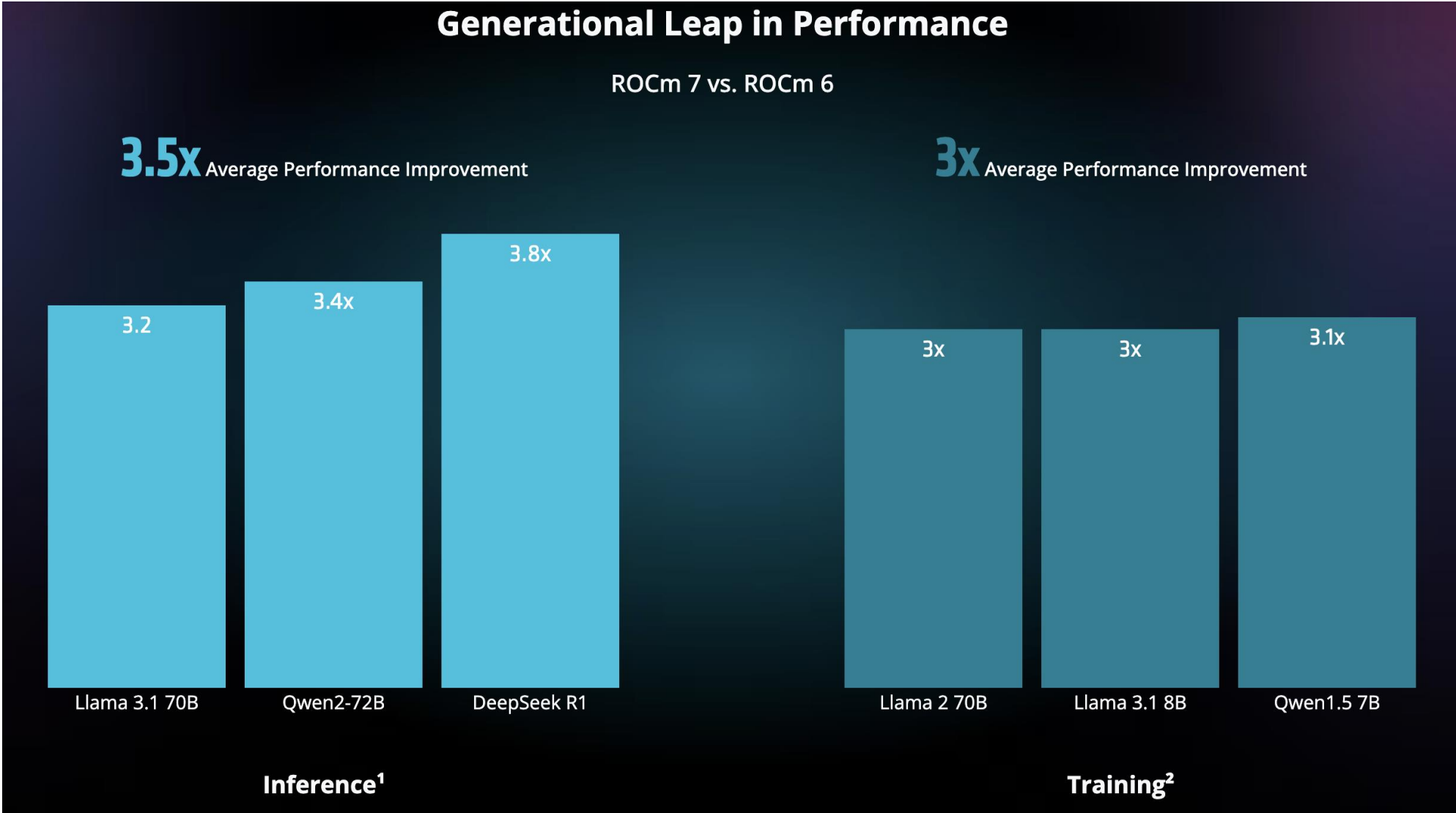
ROCM Software ecosystem



ROCm CORE



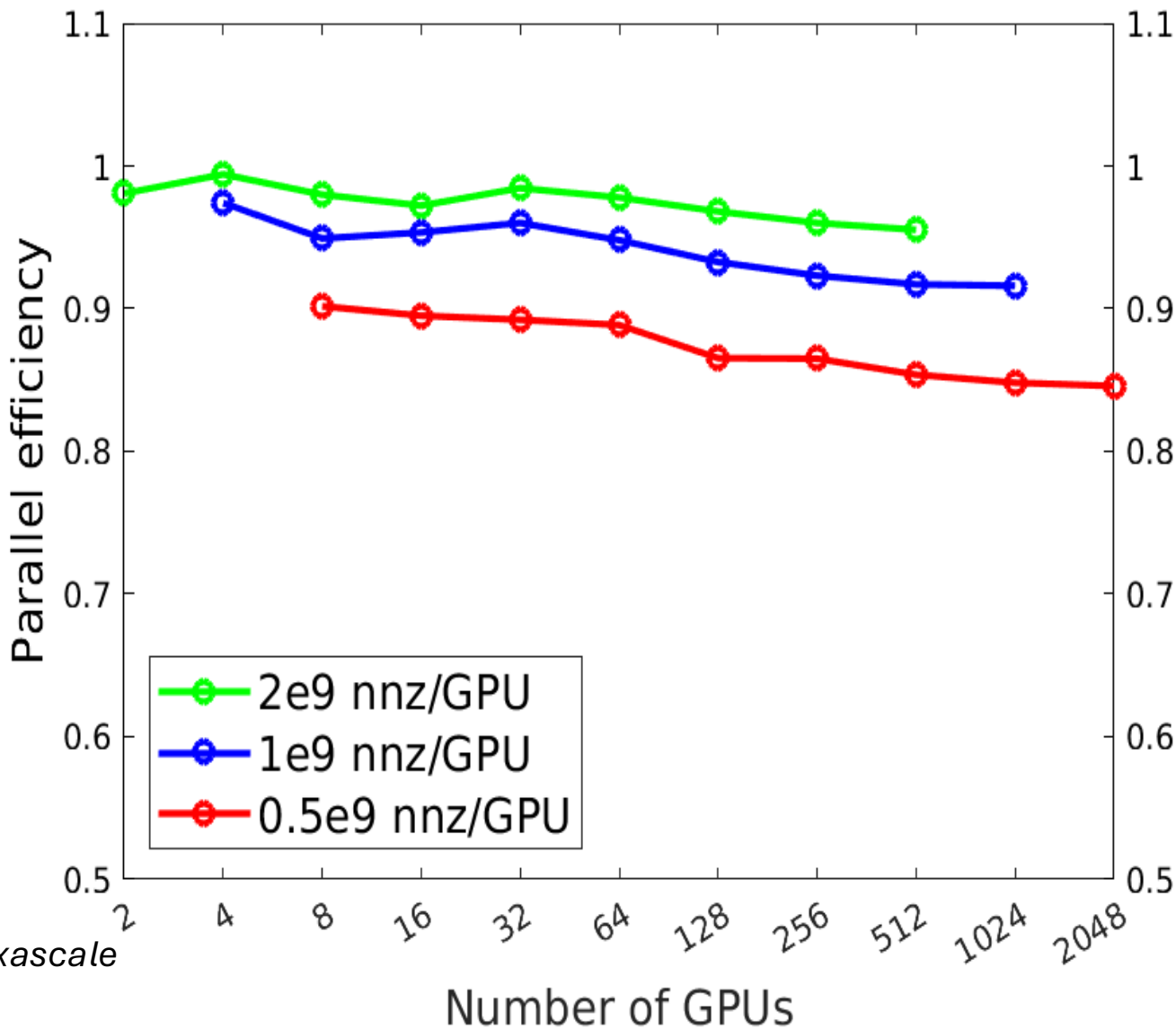
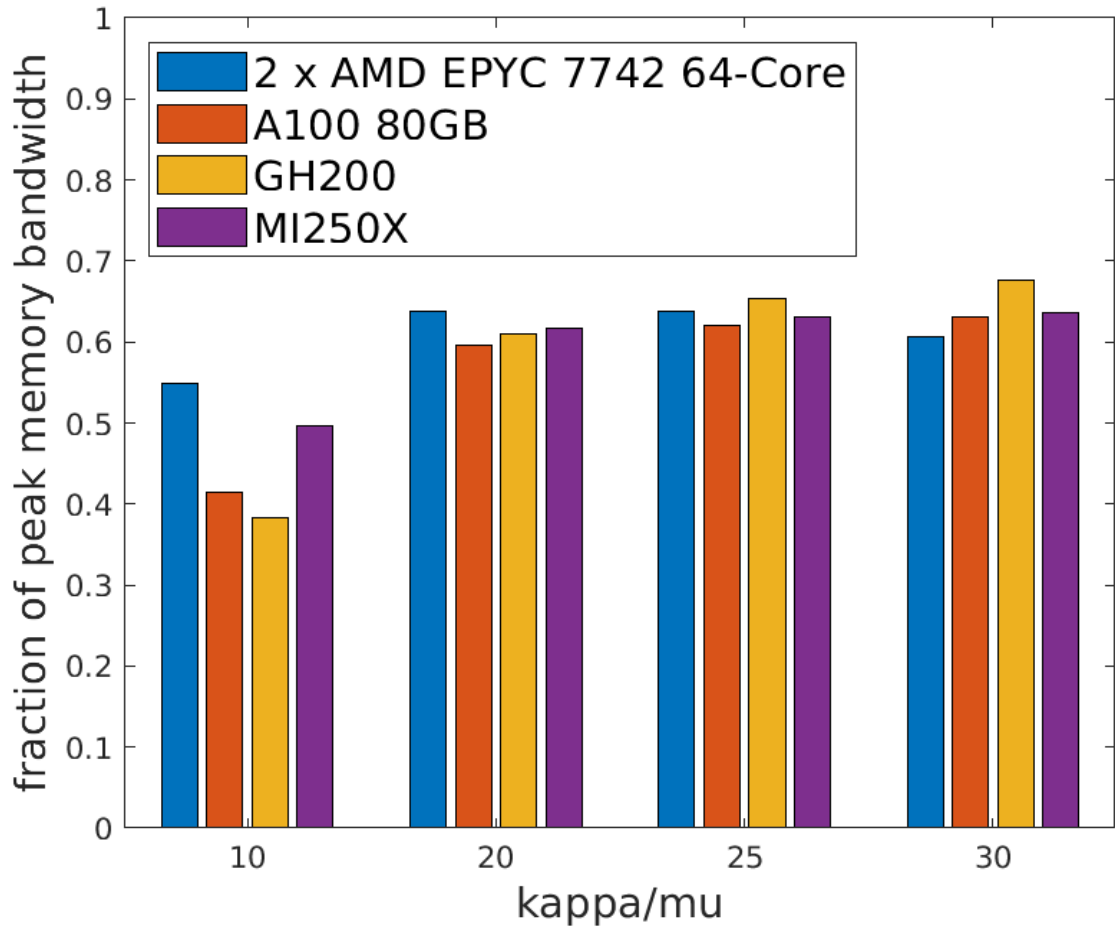
Latest ROCm Software: 7.



Benchmark: GPU-accelerated solver

Time-dependent Dirac Eq.

LUMI-G



Final Takeaway

Hardware & Architecture

- Chiptlets, high-bandwidth interconnects, and AI engines power modern GPUs (GB300, MI350).
- Scaling driven by multi-die designs & massive memory bandwidth.

Software Ecosystems

- CUDA and ROCm deliver full-stack acceleration for AI & HPC.
- New tooling makes GPU programming faster and simpler.

Big Picture

- GPUs are the engines of modern AI and scientific computing.

In short: Modern GPUs deliver scalable, efficient performance, from hardware to software.

I stop here 😊

Access Slides



https://github.com/HichamAgueny/multigpu_mpi_course

**Get in touch:
My LinkedIn profile**



<https://www.linkedin.com/in/hicham-agueny-956a1368/>