



Technological Update

November 2024

NVIDIA CONFIDENTIAL PRESENTATION

- All information in the following presentation is NVIDIA confidential, including codenames, future products, and performance projections

- No information in this presentation is allowed to be revealed or published without NVIDIA consent

- Sharing or distributing copies of this presentation to anyone is strictly prohibited

- Use of cameras to capture information is strictly prohibited

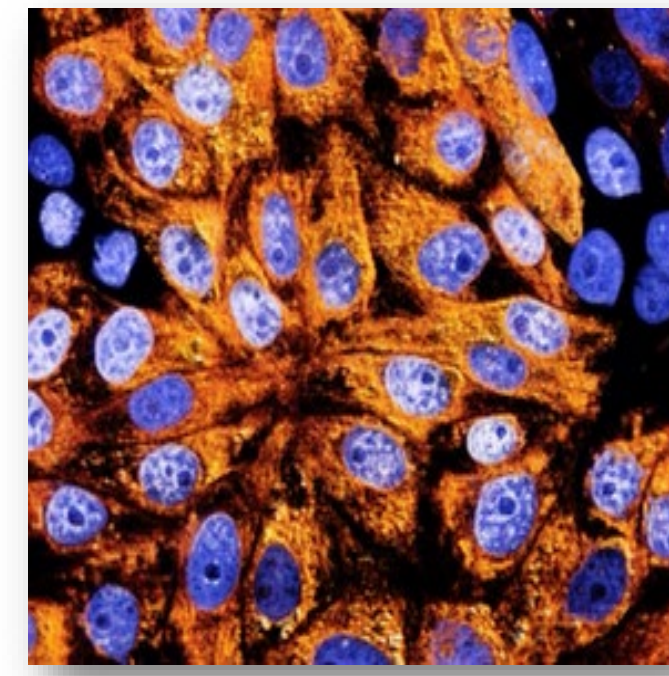
Customer Guidelines

Agenda

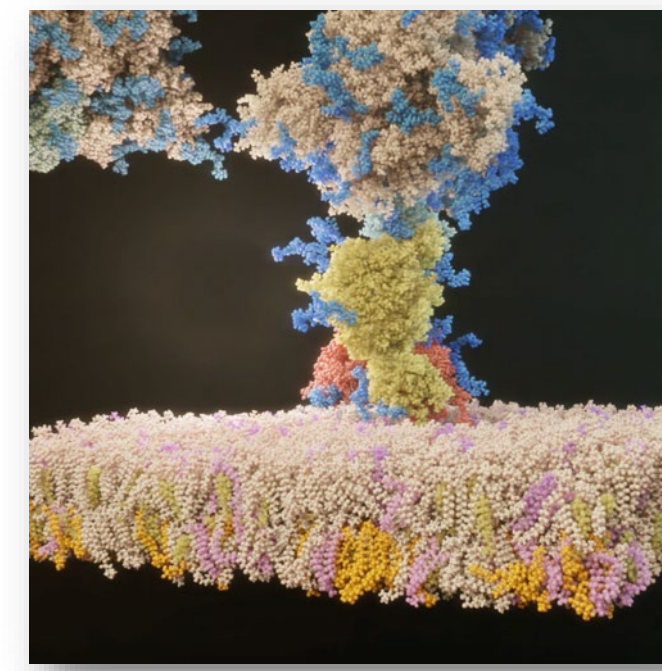
- NVIDIA hardware architecture overview;
- Main differences between the different NVIDIA GPU architectures;
- Overview of the main drivers / software libraries needed for the different hardware components;
- Programming languages available for GPUs;
- Tools for debugging and performance analysis (e.g. Visual/Compute Profiler);
- Main techniques for optimizing code performance;
- Main techniques for optimizing intra-node and inter-node multi-GPU communications with infiniband (e.g. GPUDirect P2P, RDMA); Roadmap of NVIDIA technological developments (for what concerns CPU/GPU);
- References to in-depth courses on the different topics.

Pre-Exascale Supercomputing

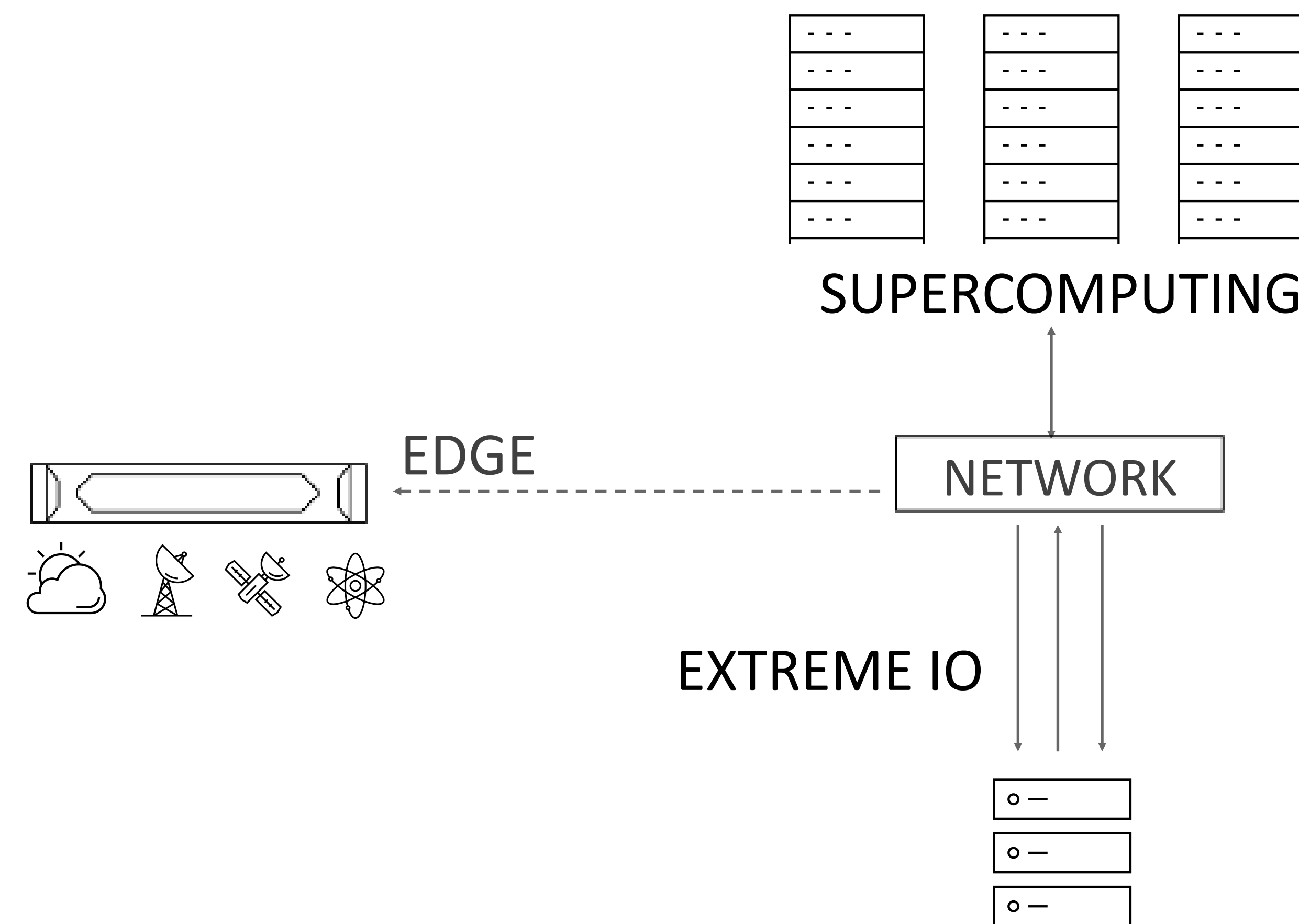
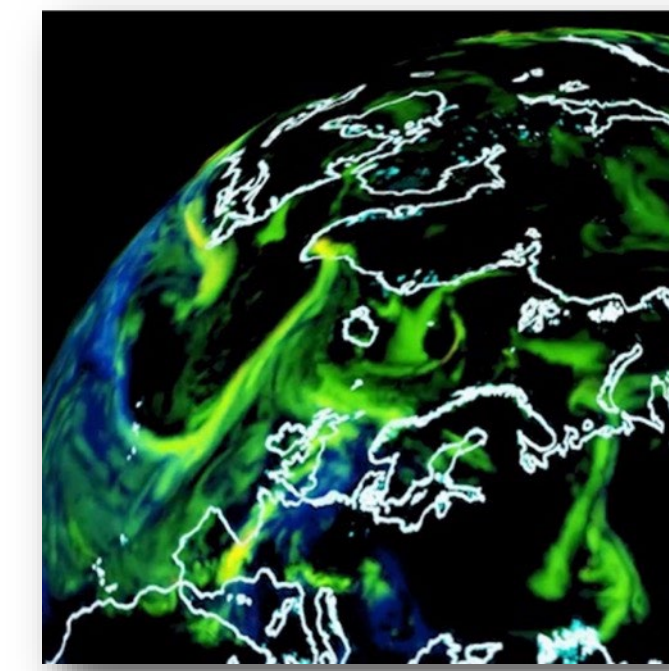
EDGE



SIMULATION

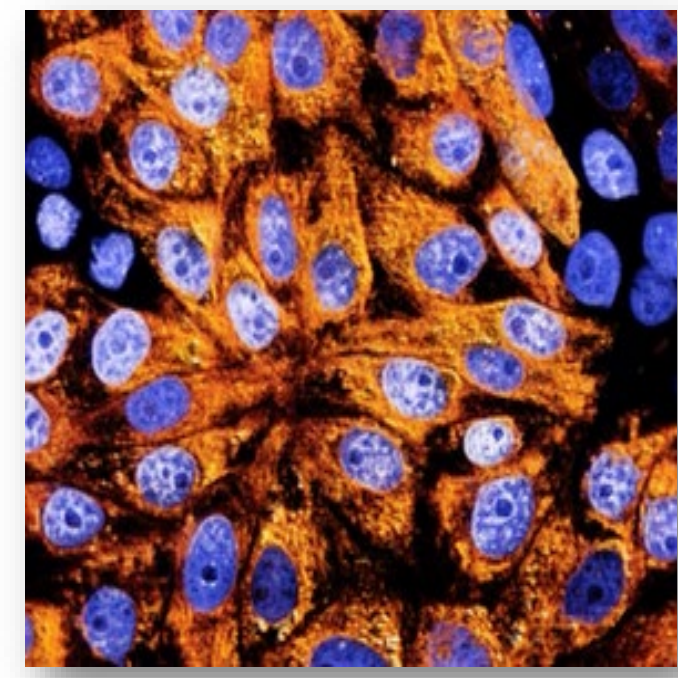


Viz

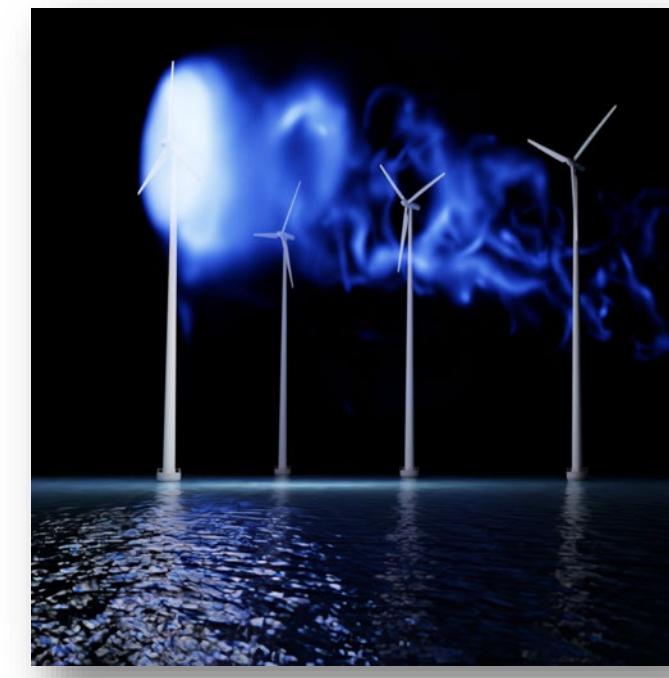


Exascale Supercomputing

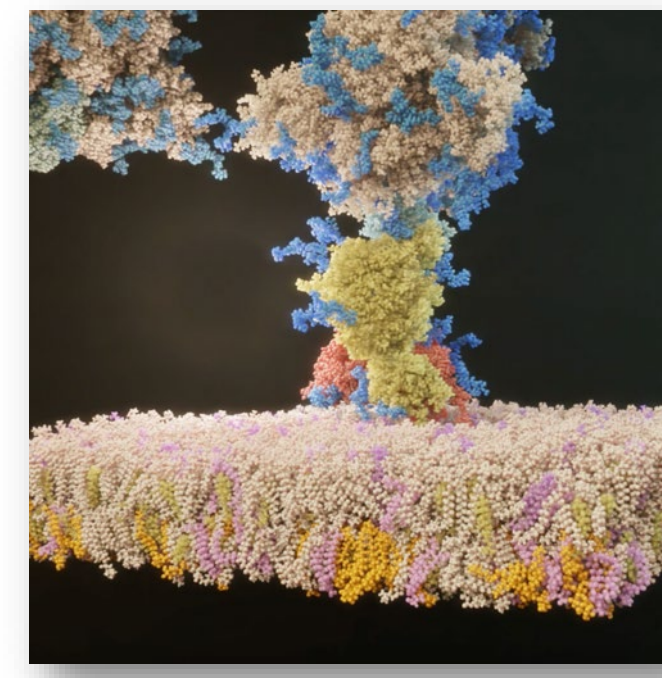
EDGE



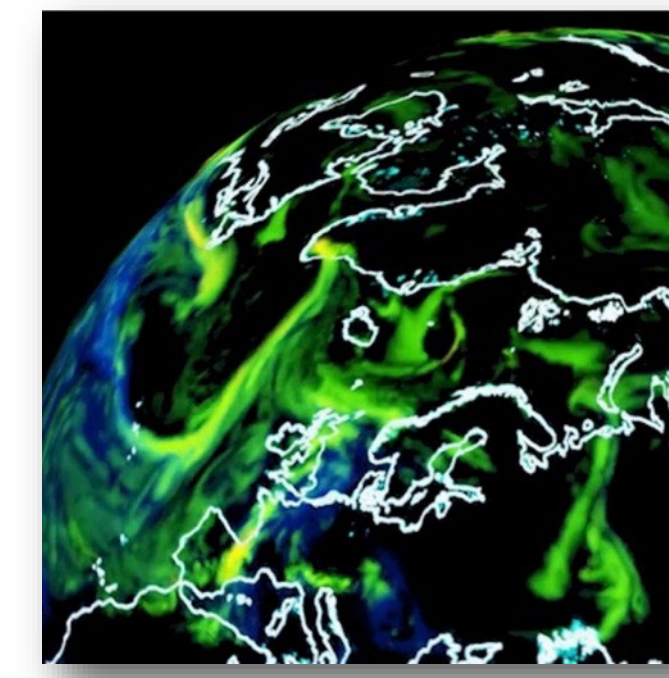
SIM + AI



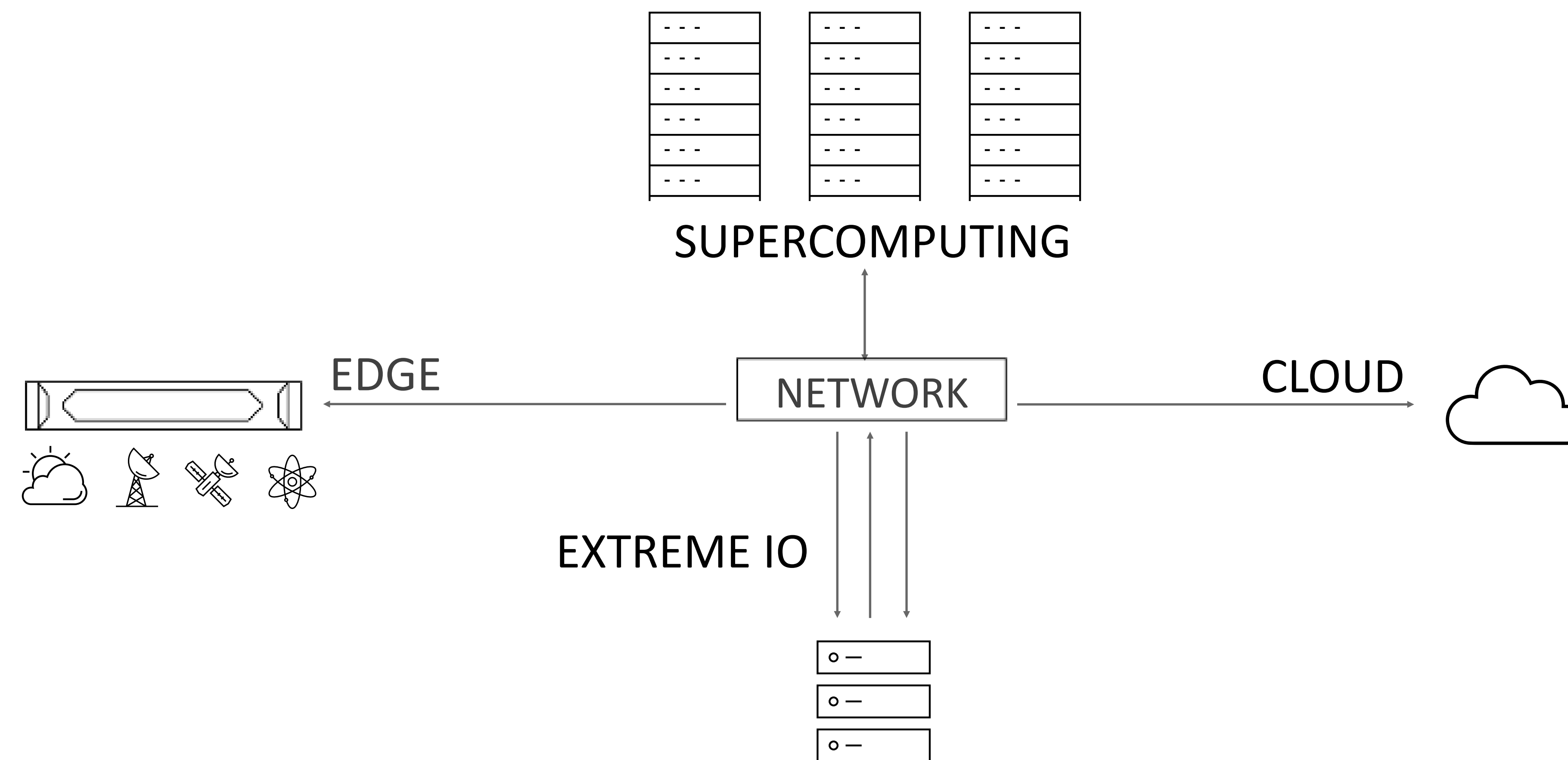
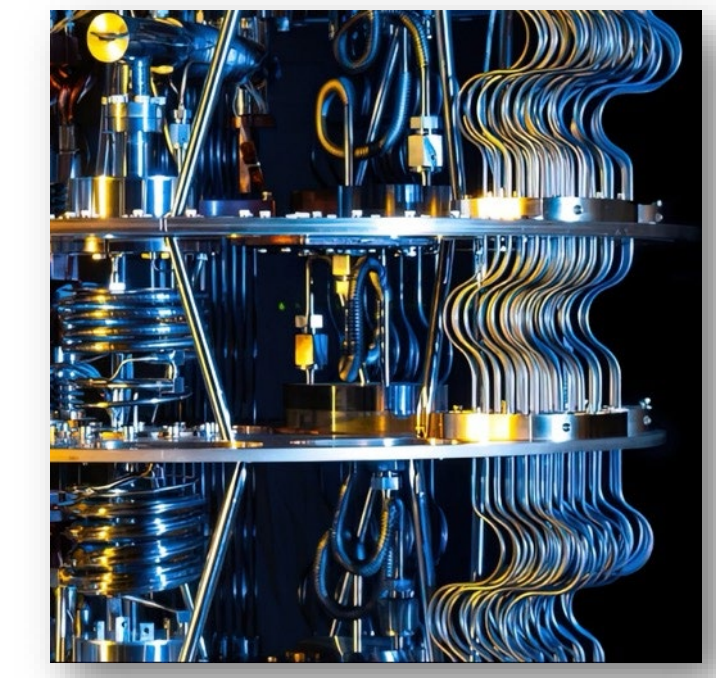
SIMULATION



DIGITAL TWIN

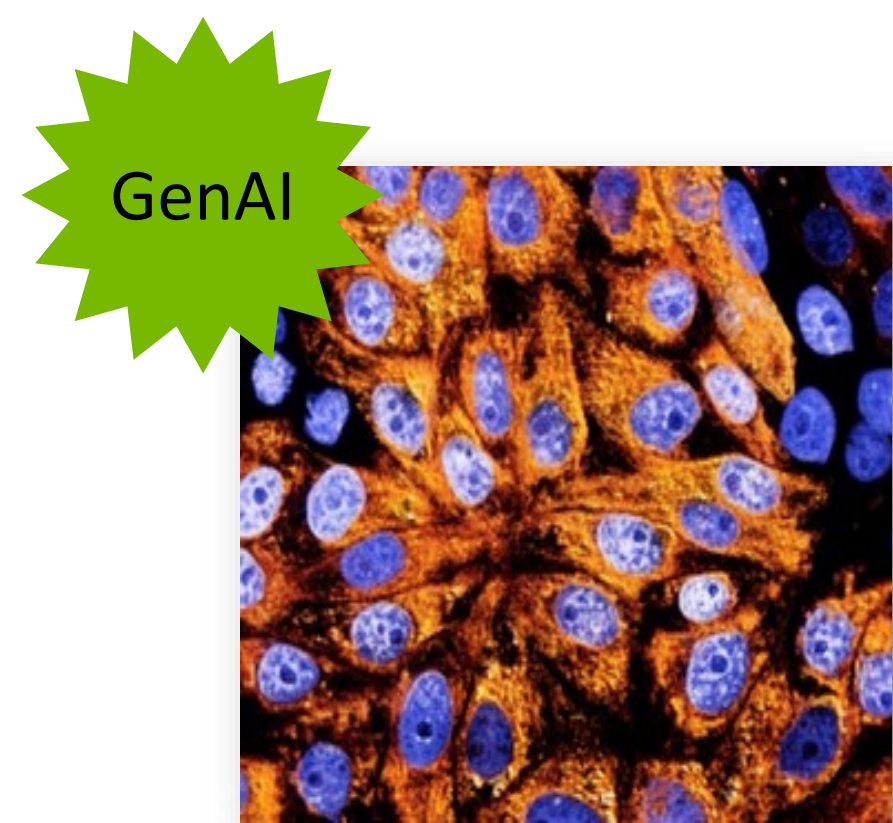


QUANTUM COMPUTING



AI Factory for Research

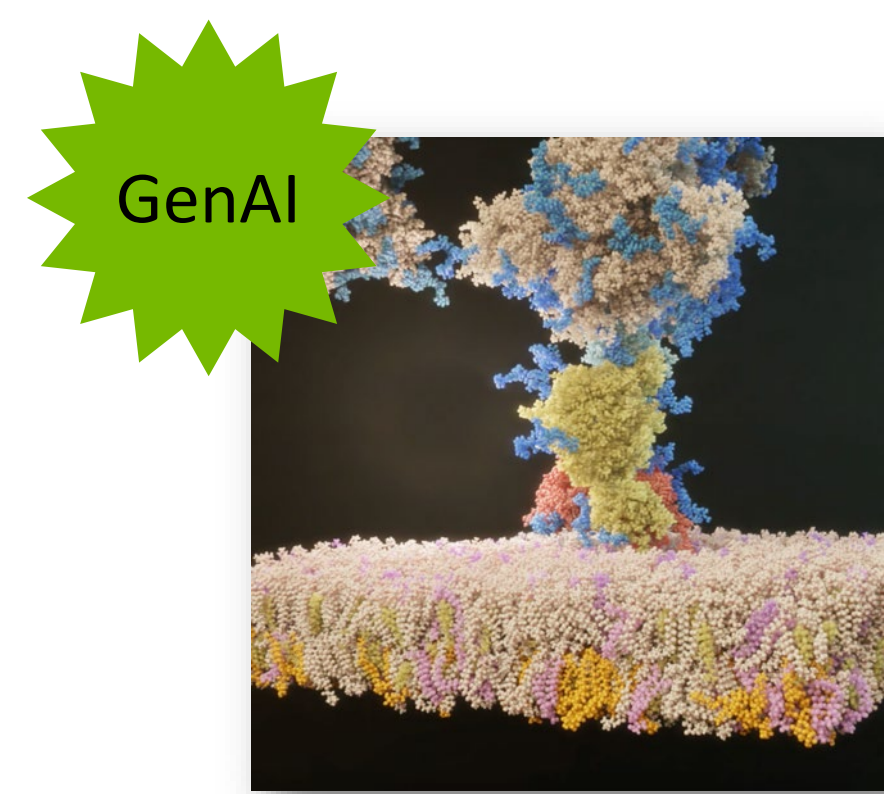
EDGE



SIM + AI



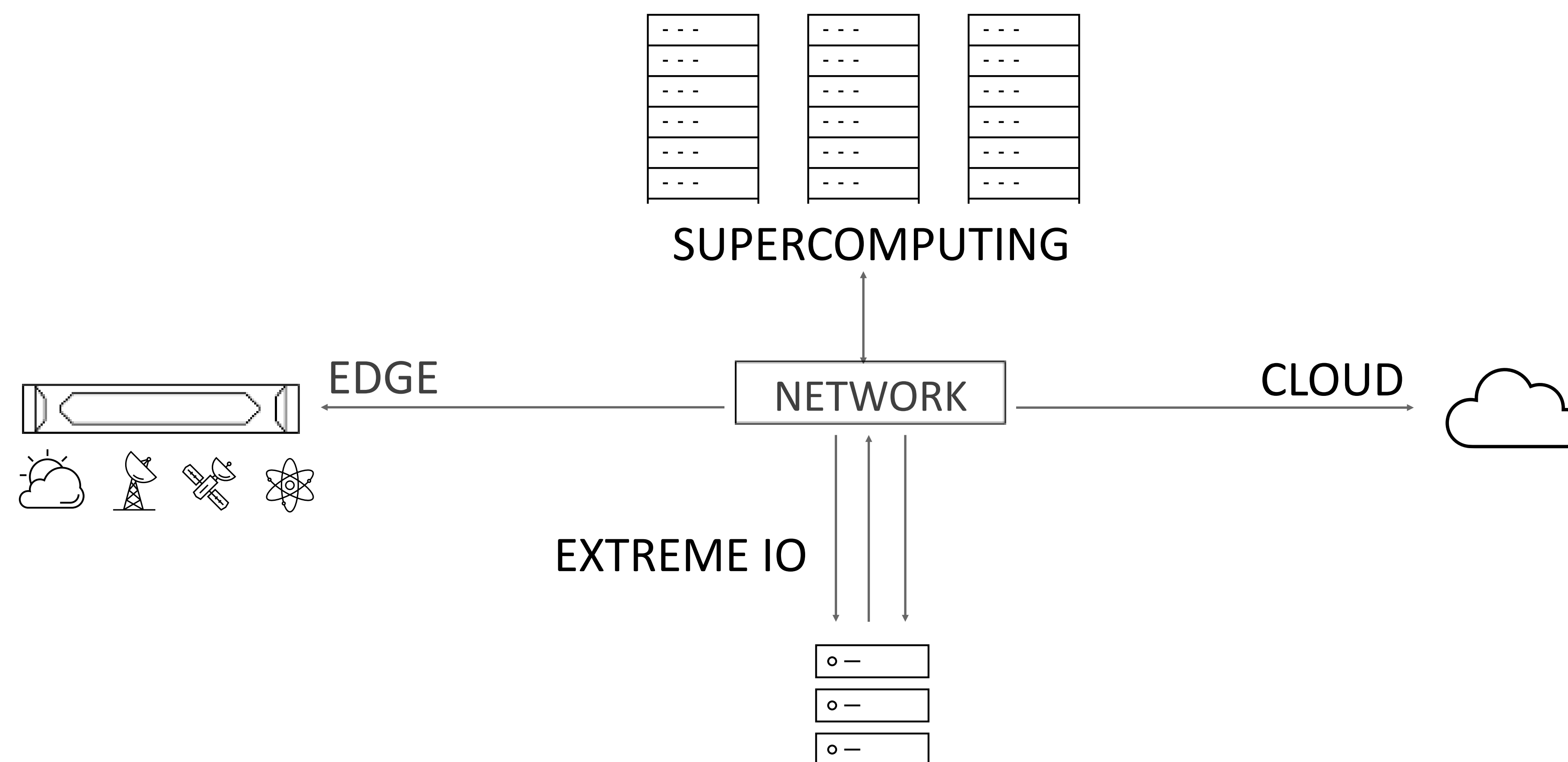
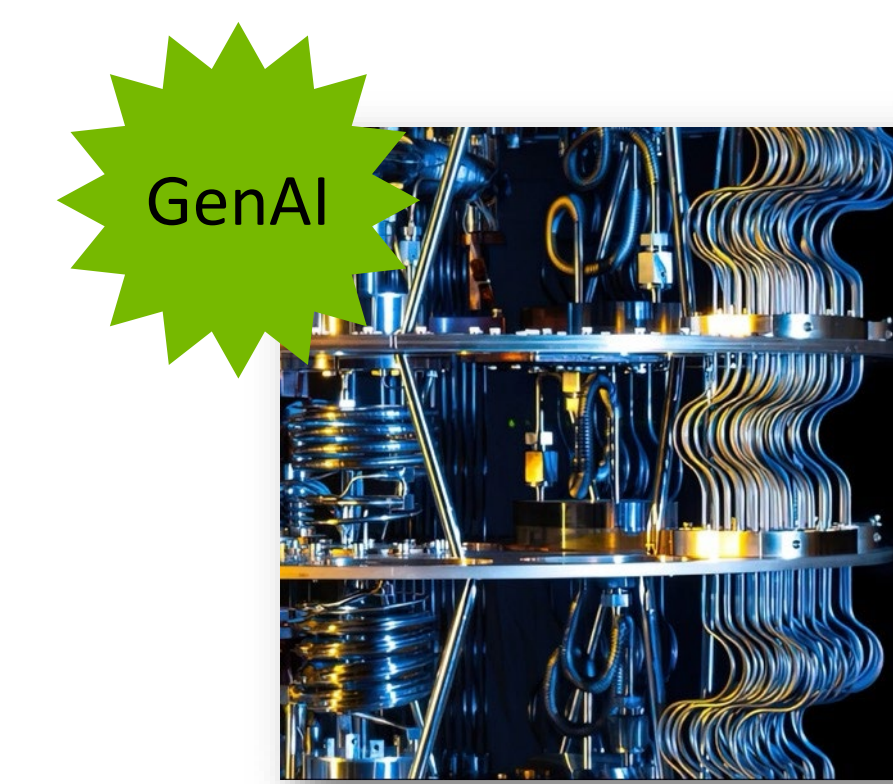
SIMULATION



DIGITAL TWIN

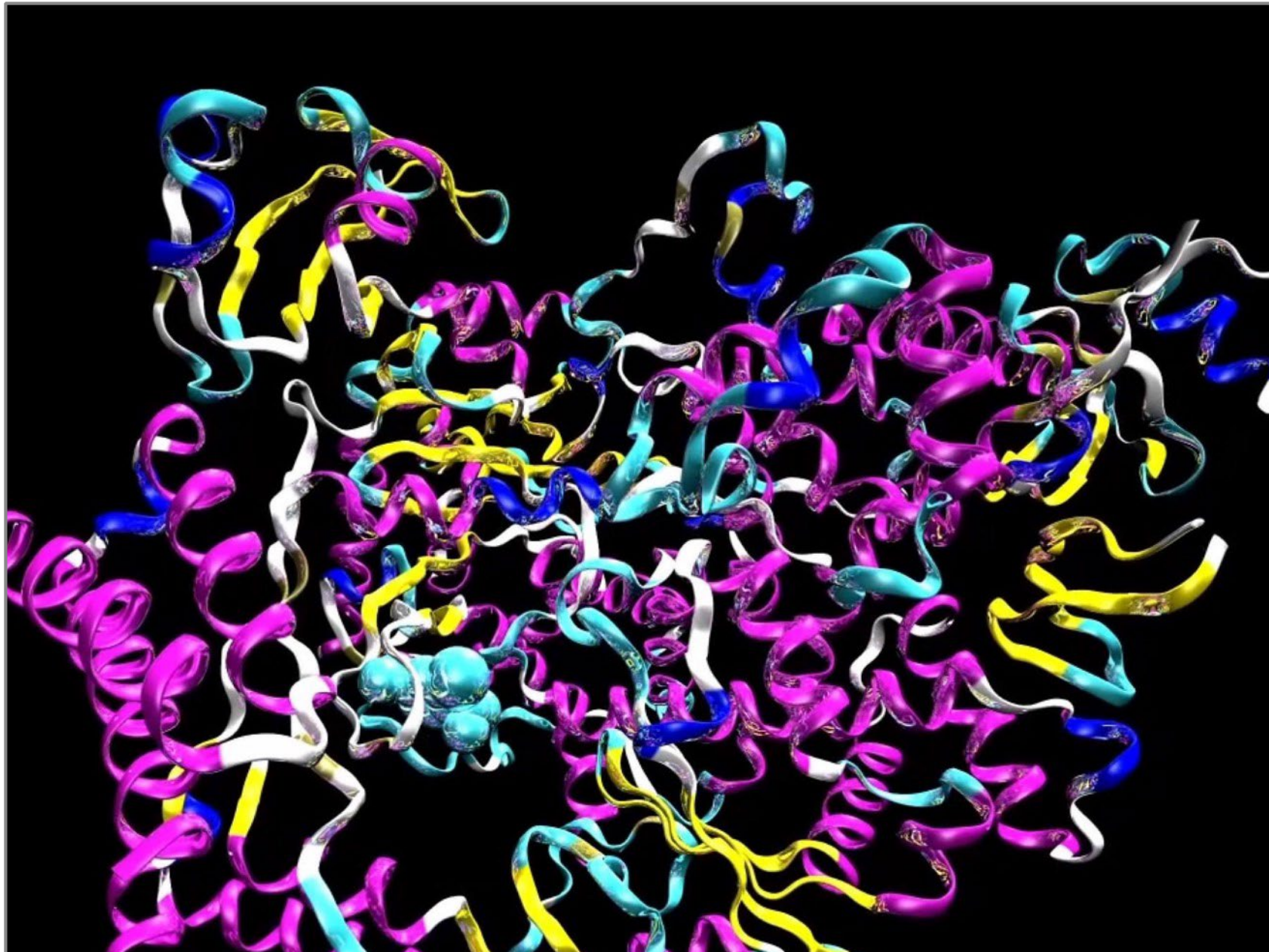


QUANTUM COMPUTING

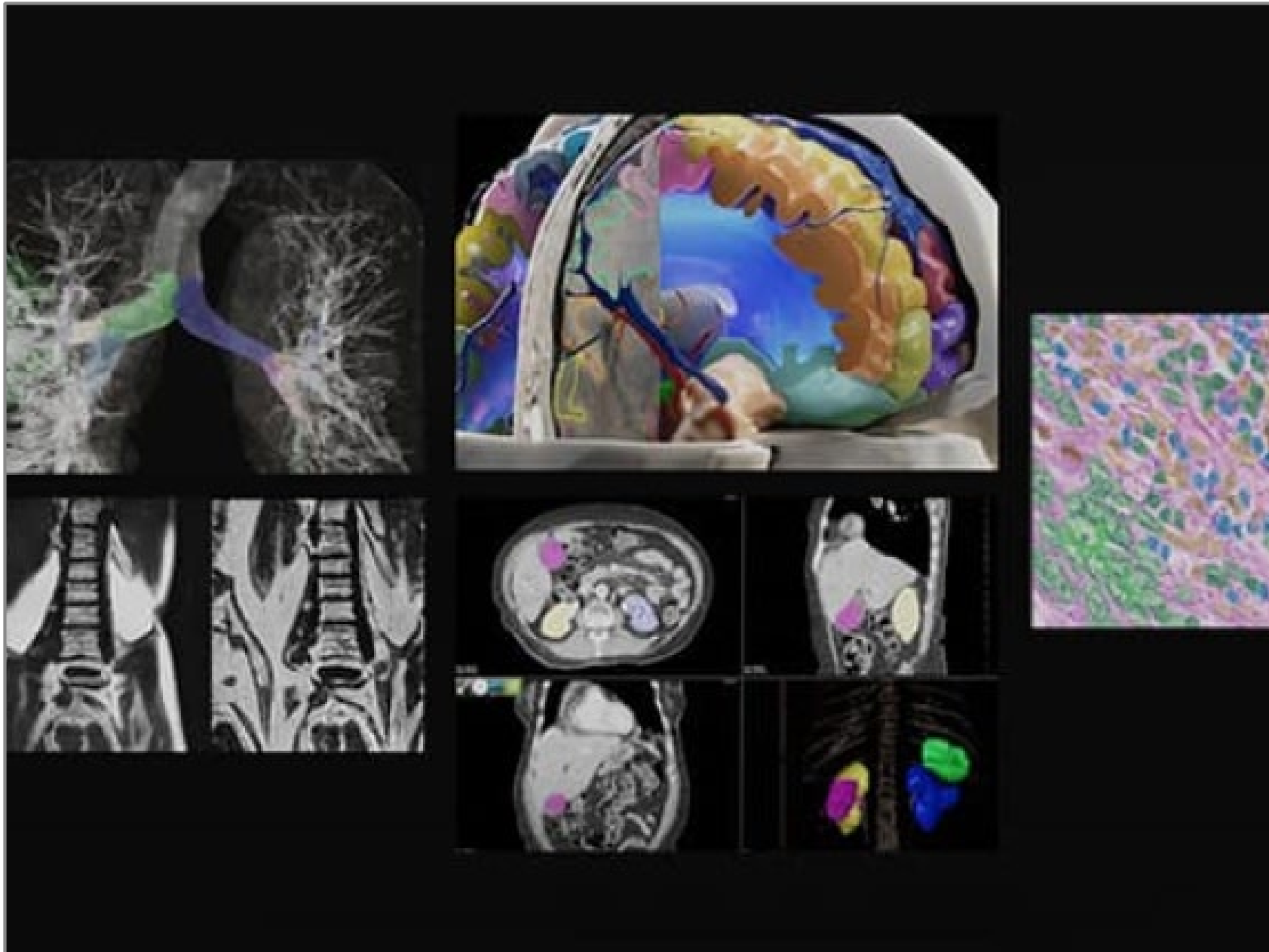


AI: The new tool for Science

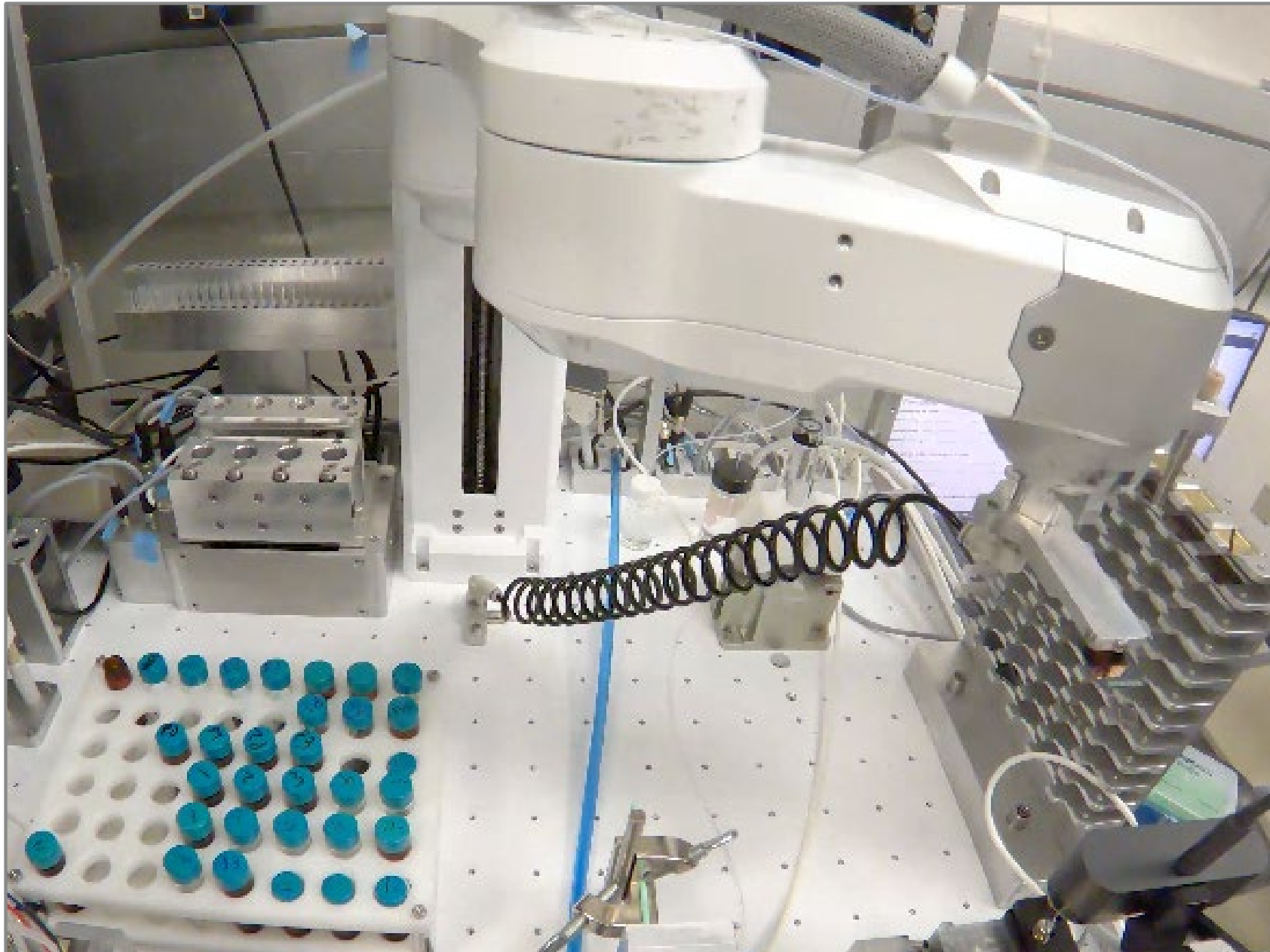
DRUG DISCOVERY
EvolutionaryScale



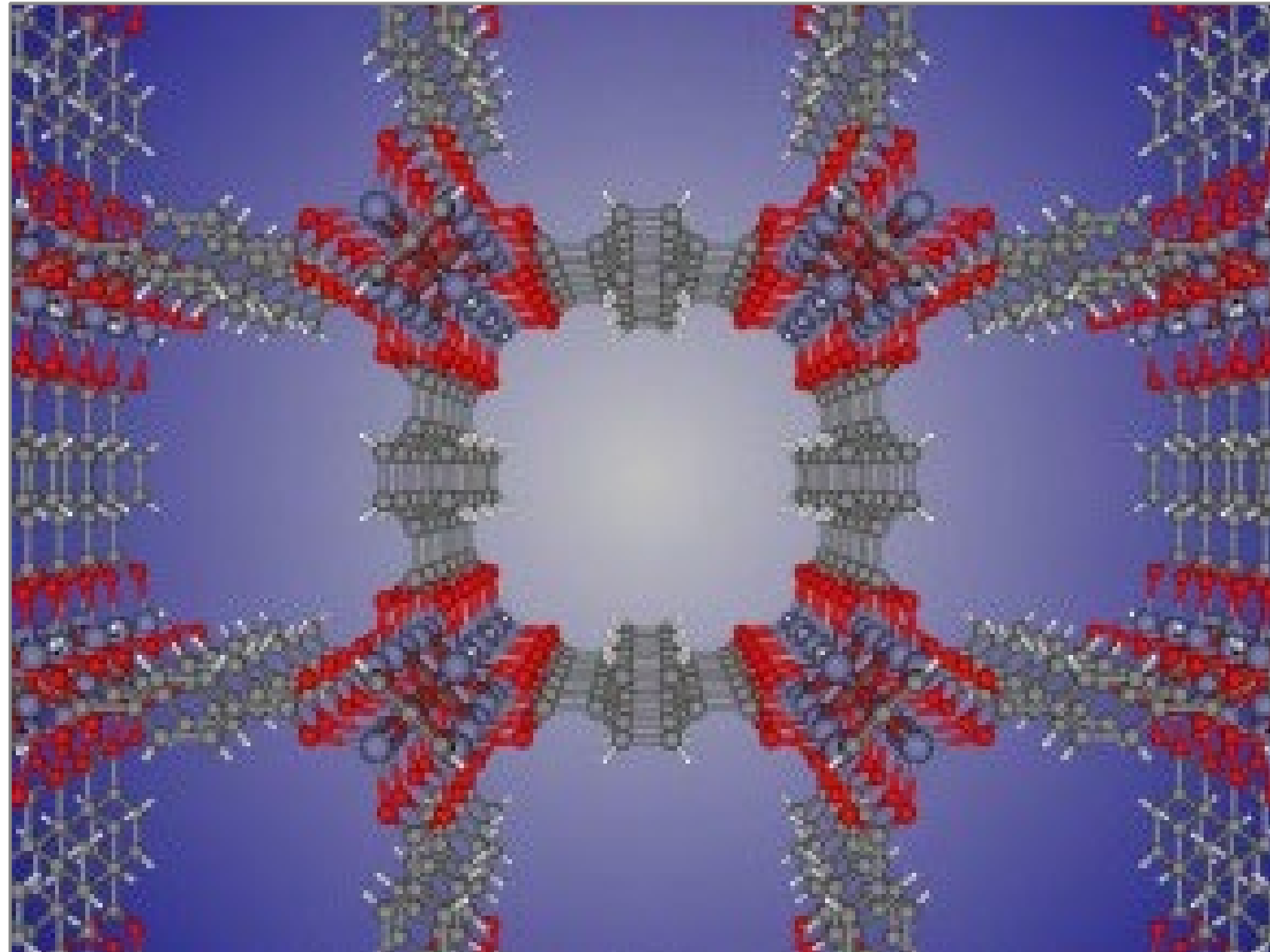
CANCER RESEARCH
Wellcome Sanger Institute



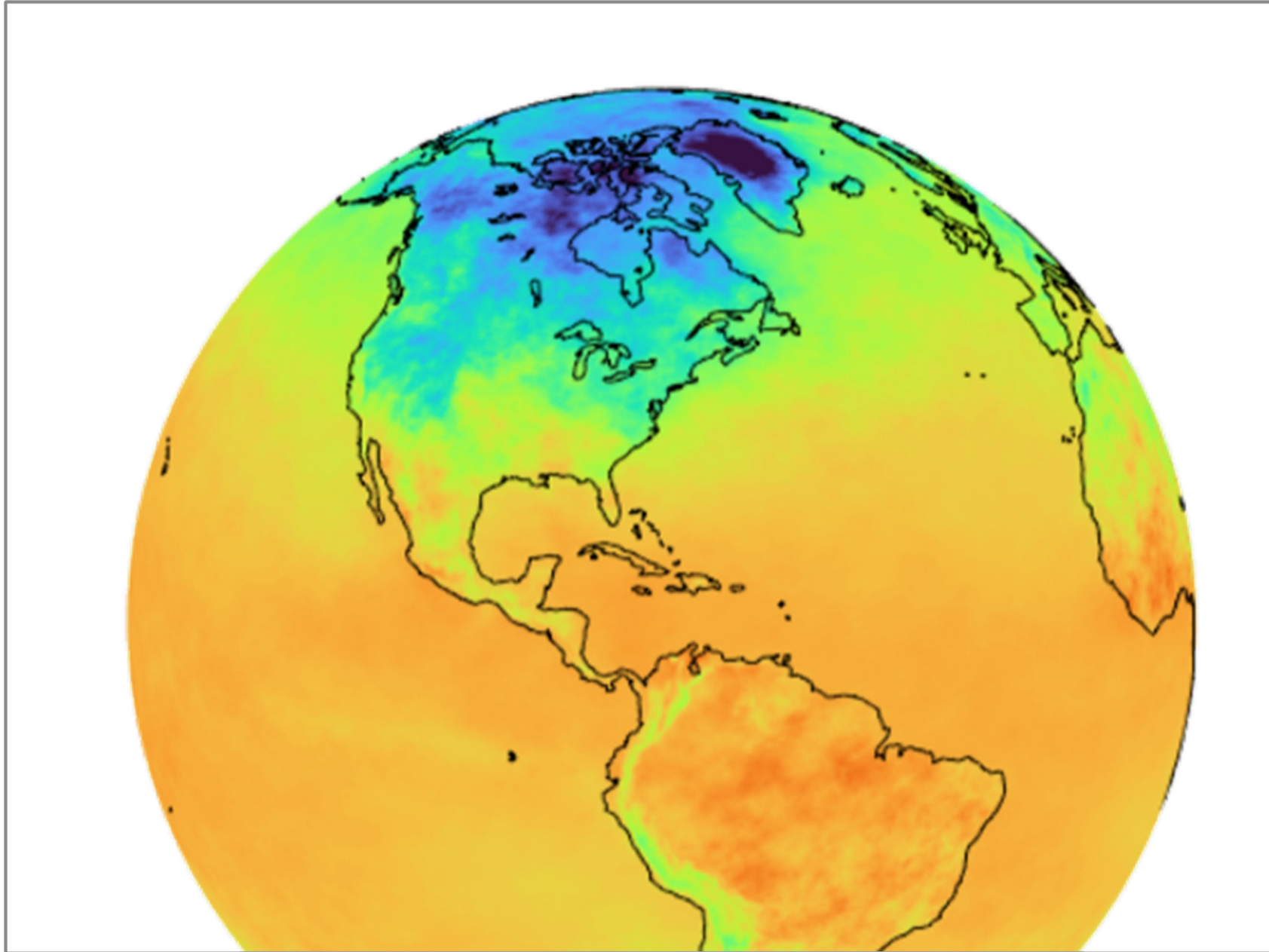
LABORATORY ROBOTICS
Argonne National Lab



MATERIALS DISCOVERY
Microsoft Research



CLIMATE MODELING
KAUST



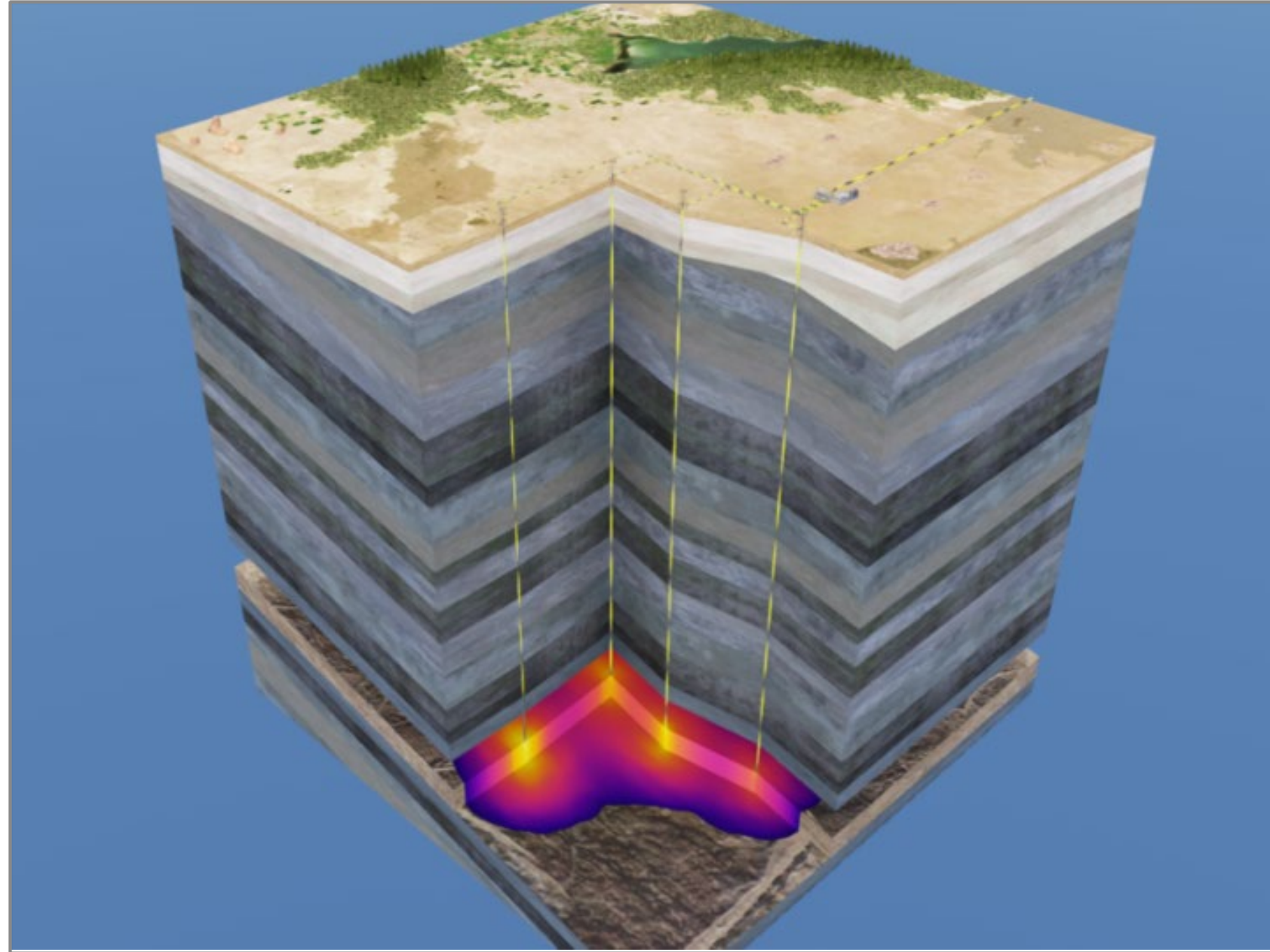
ASTRO FOUNDATION MODELS
The Flatiron Institute



AGRICULTURAL HEALTH
Fermata

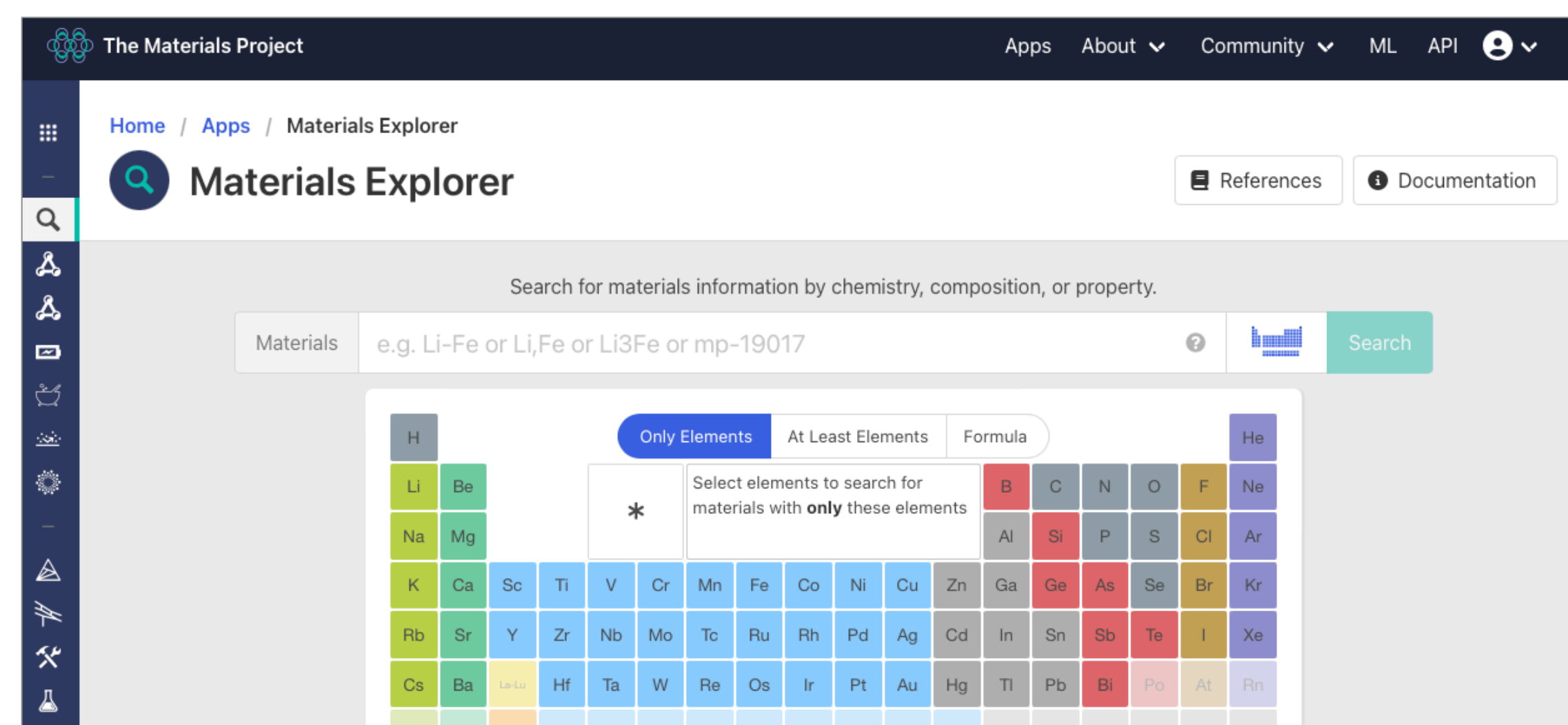


SEISMIC GEOSCIENCE
California Institute of Technology



AI is Transforming Material Science and Chemistry

Broad industry innovation and achievement



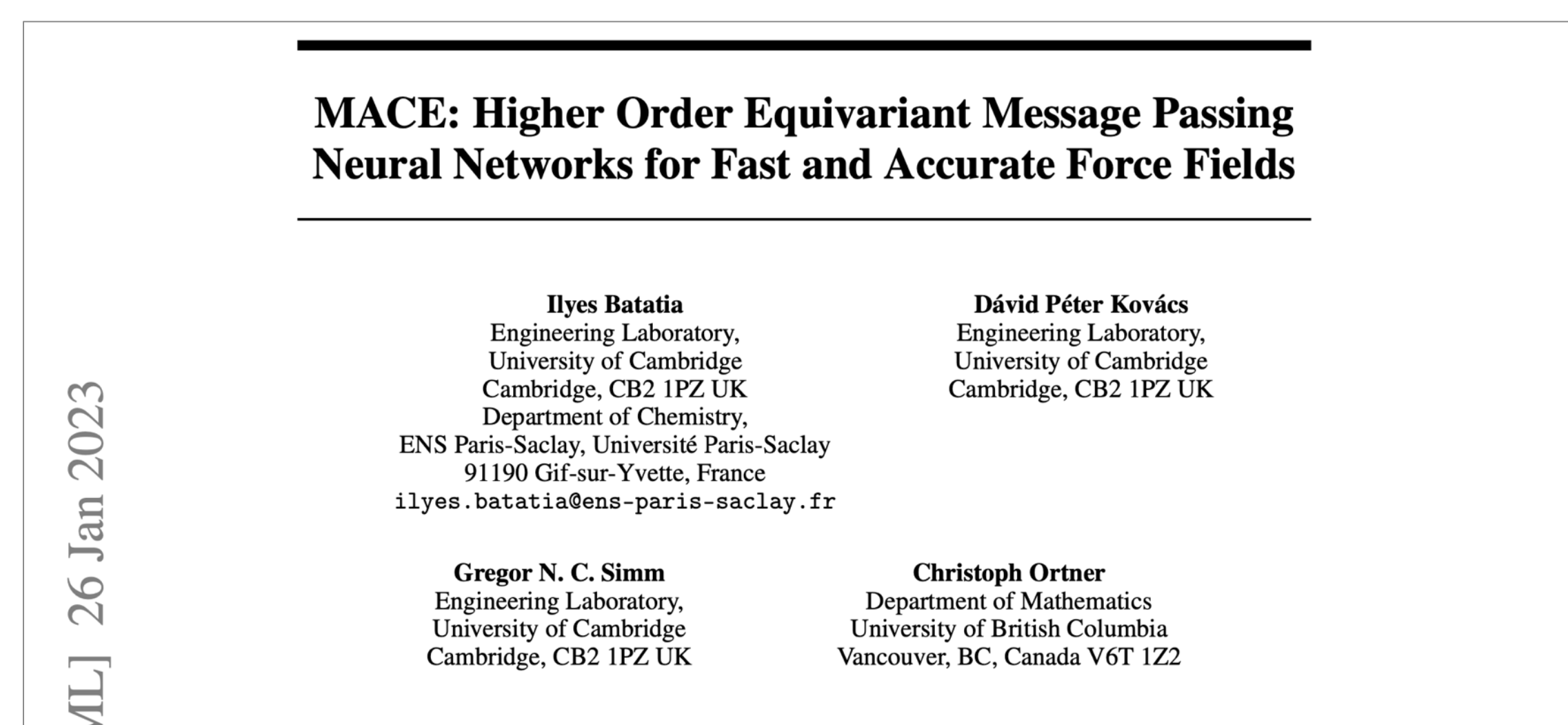
Materials Project
Multi-National Effort



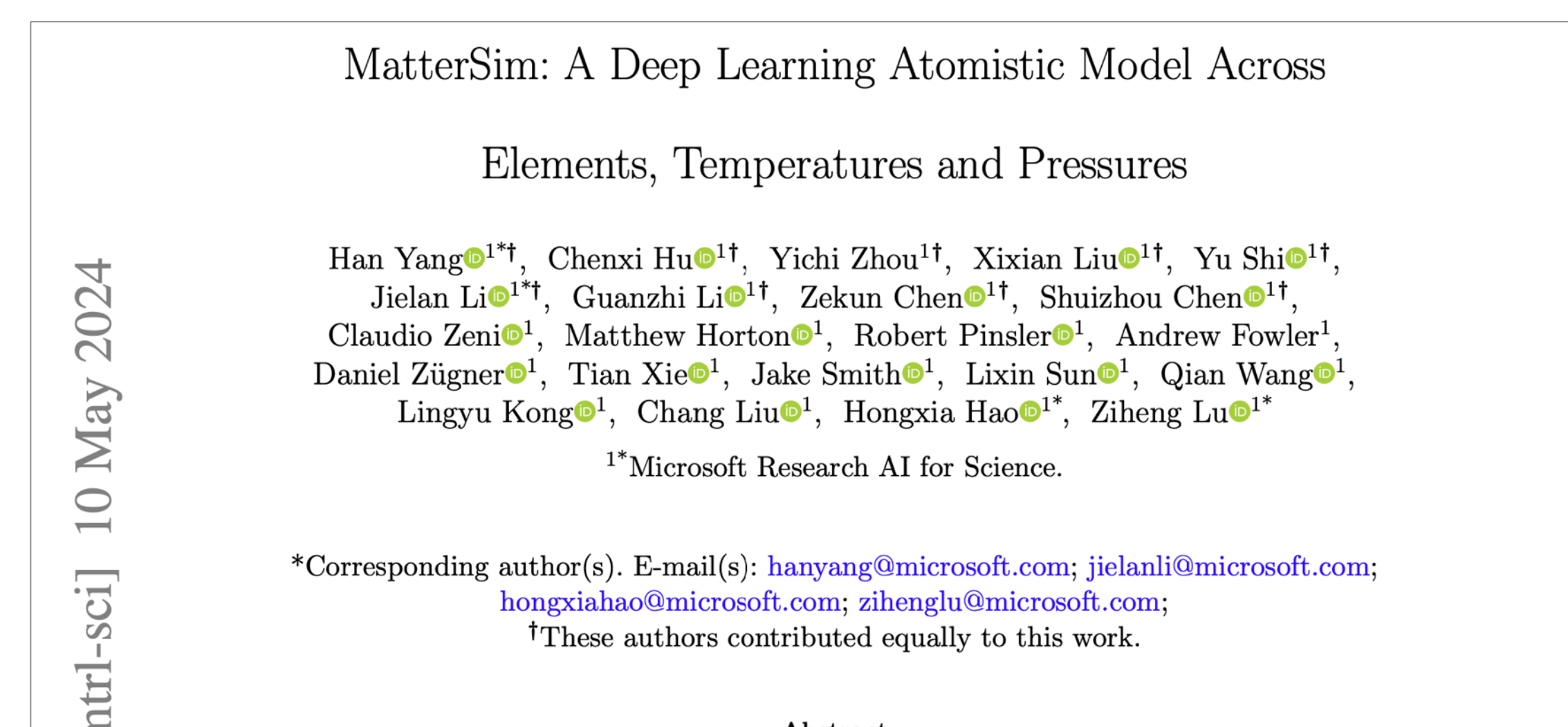
GNoME
Google DeepMind



MatterGen
Microsoft Research



MACE
OpenSource



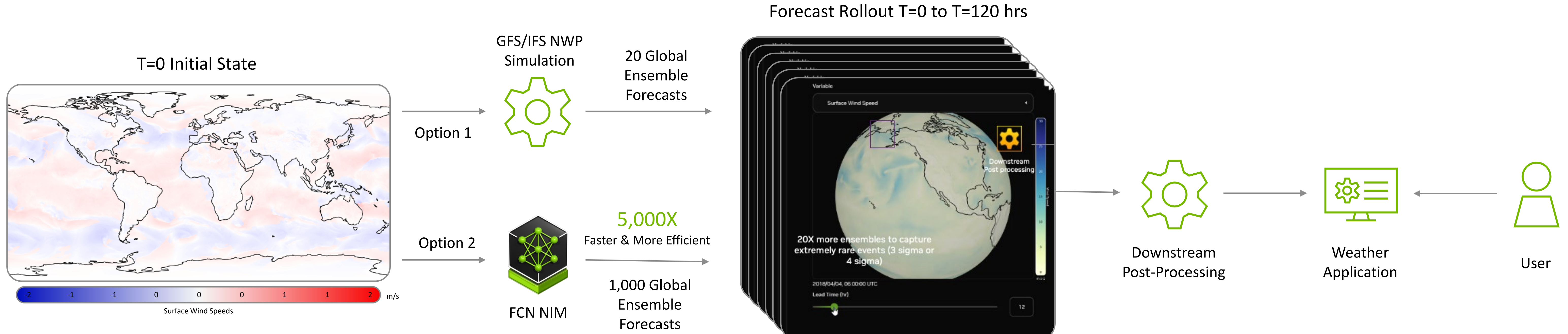
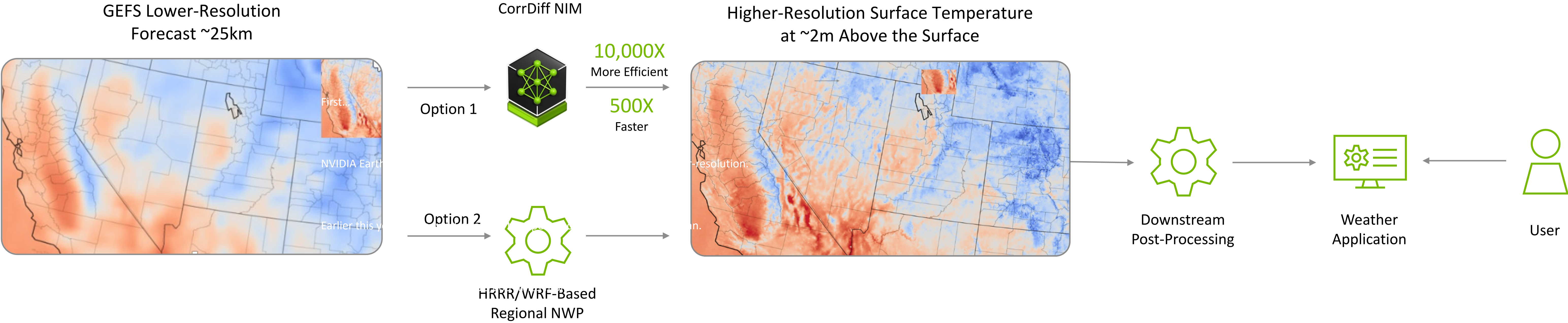
MatterSIM
Microsoft Research



Open Materials 2024
Meta FAIR

Announcing Earth-2 NIMs for CorrDiff & FourCastNet

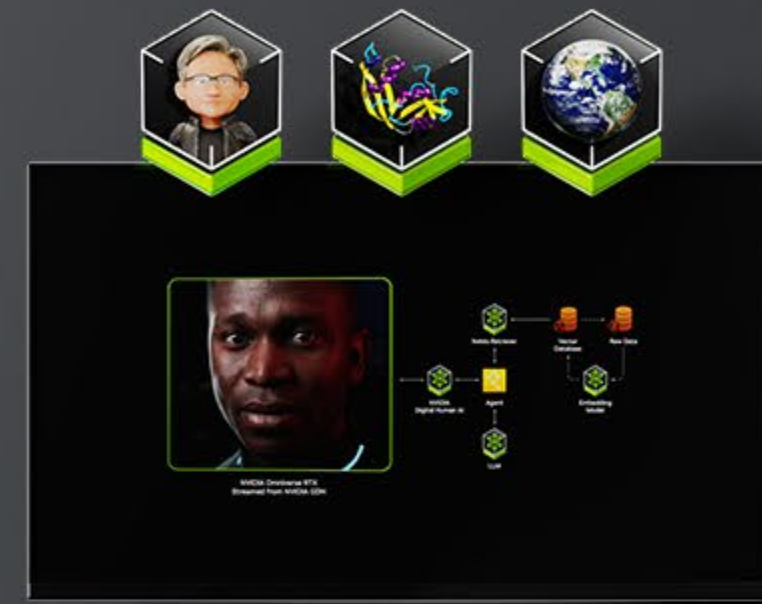
CorrDiff NIM for Generative AI Powered Downscaling | FourCastNet NIM for Global Weather Forecasting



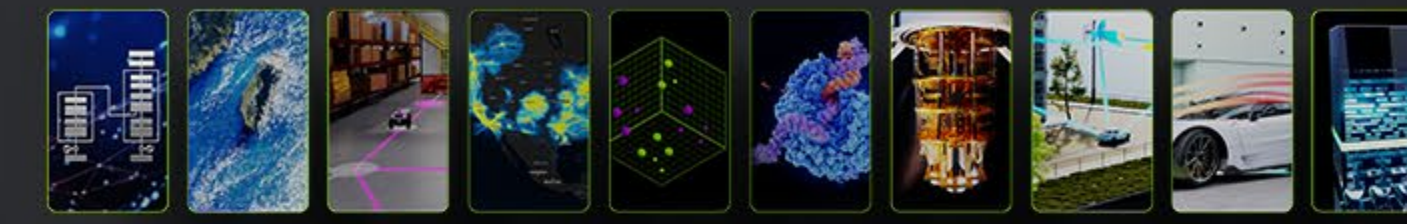
NVIDIA BUILDS AI SUPERCOMPUTING INFRASTRUCTURE

One Year Rhythm | Supercluster Scale | Full-Stack | CUDA Everywhere

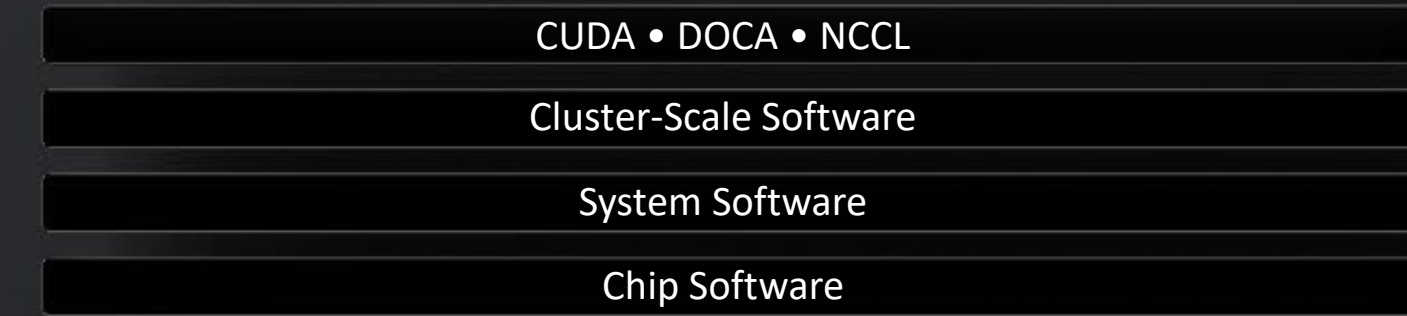
NIM
CUDA-Accelerated
Agentic AI Libraries



Omniverse
CUDA-Accelerated
Physical AI Libraries

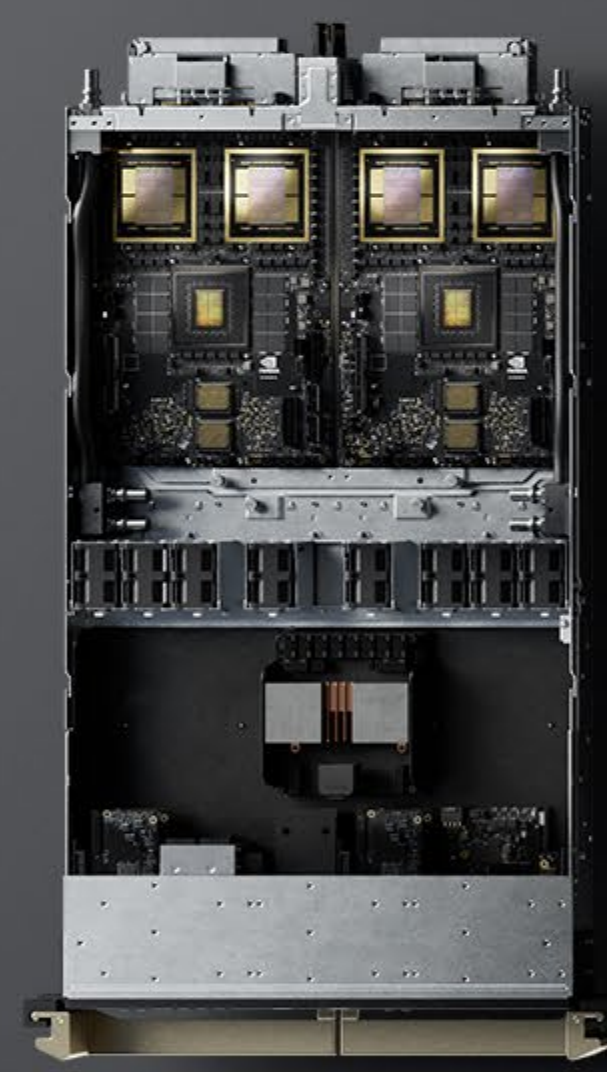


CUDA-X Libraries

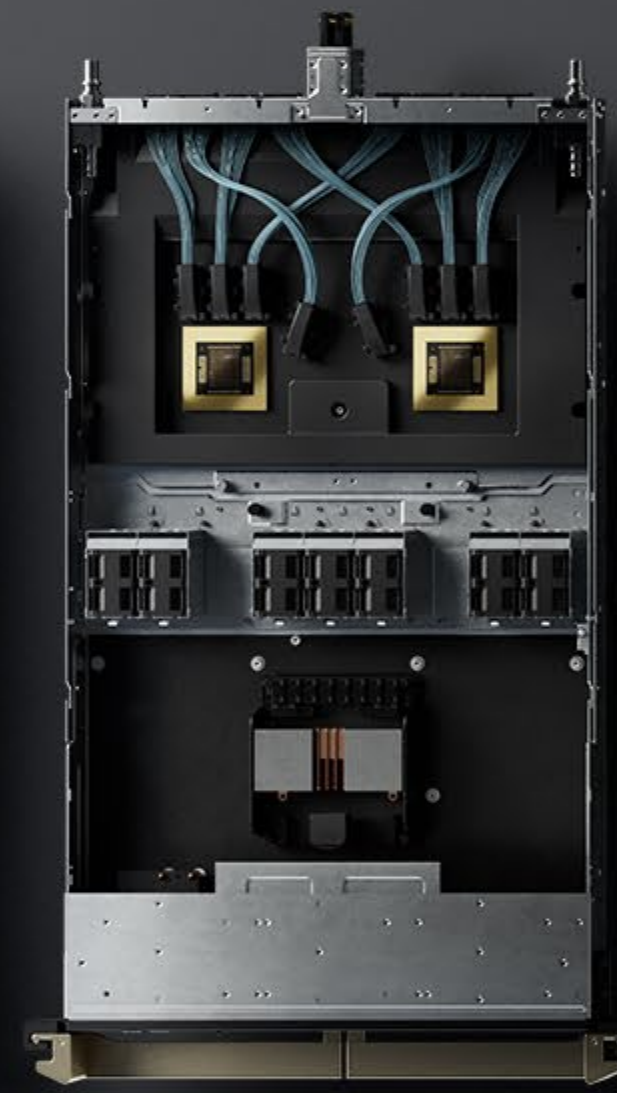


Accelerated
Software Stack

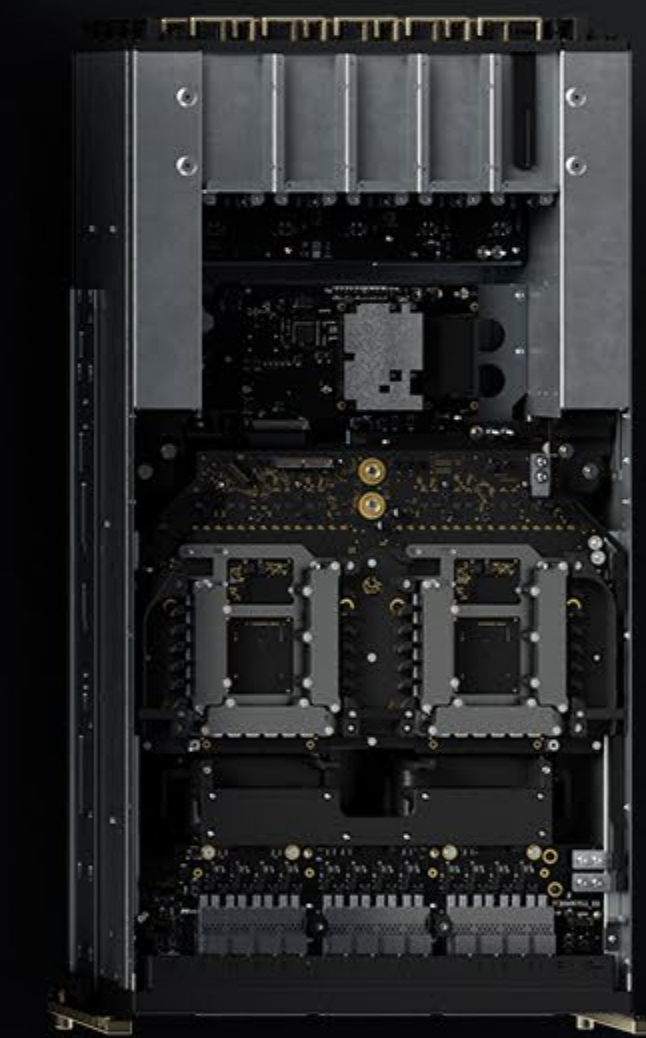
GB200 NVL72 SuperPOD



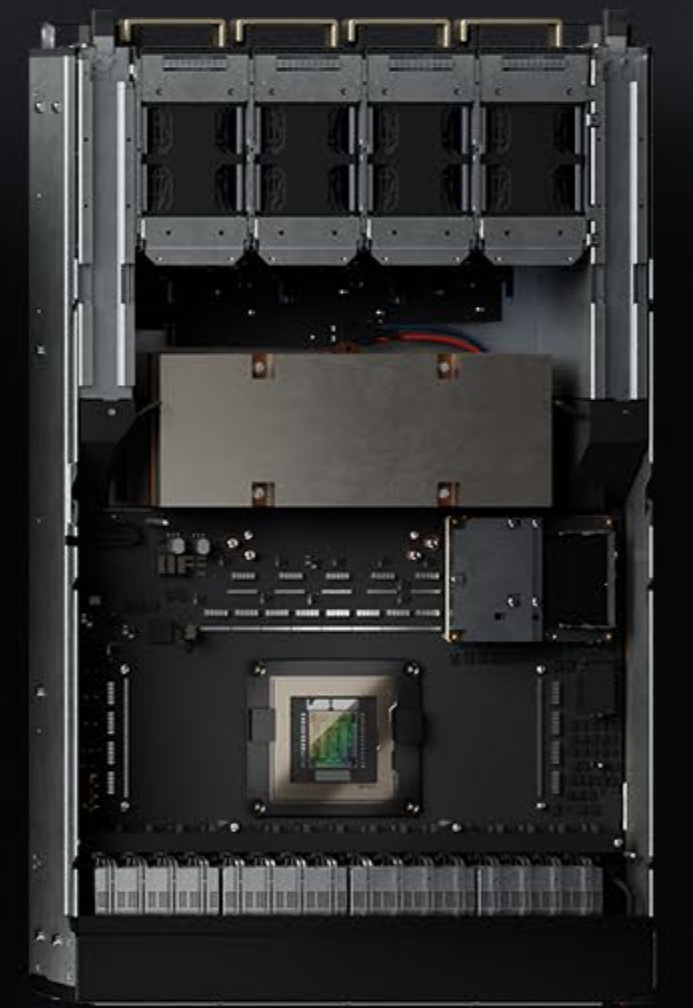
Grace Blackwell
MGX Node



NVLink Switch



Quantum Switch

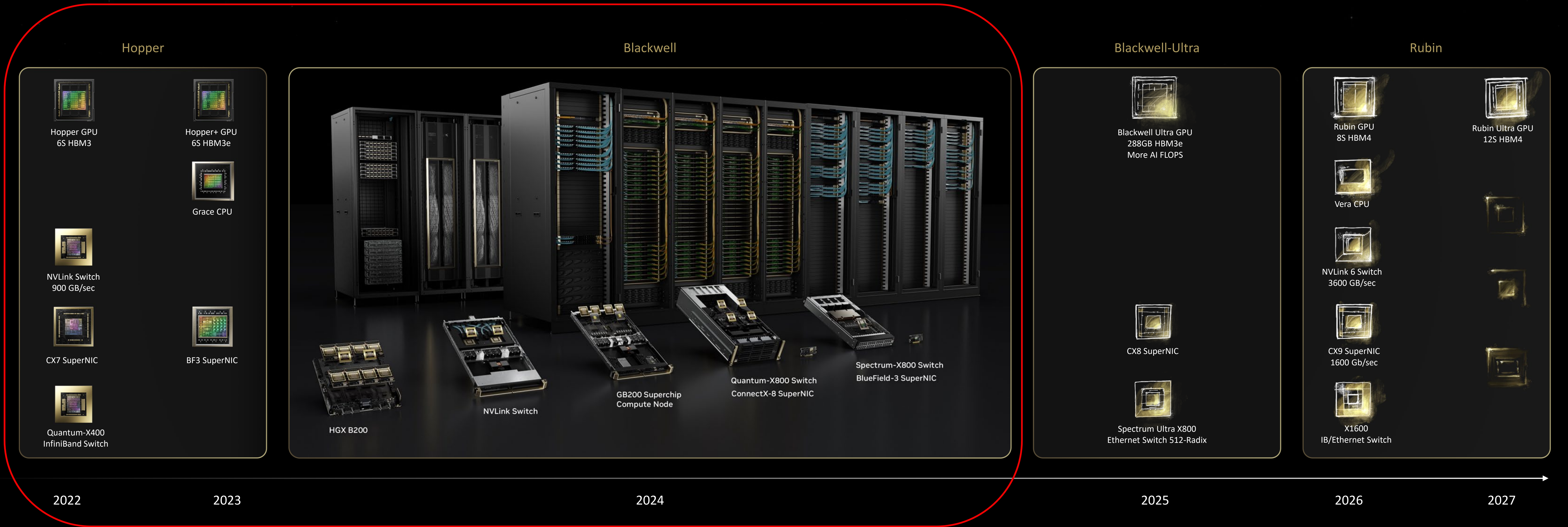


Spectrum-X Switch



Chips Purpose-Built for AI Supercomputing
GPU | CPU | DPU | NIC | NVLink Switch | IB Switch | ENET Switch

ONE YEAR RHYTHM | SUPERCLUSTER SCALE | FULL-STACK | CUDA EVERYWHERE



Ecosystem Powering the Next Wave of AI Supercomputing Systems

Partners Supercharge HPC and AI

AIVRES

ASRock
Rack

ASUS®

DELL Technologies

EVIDEN

GIGABYTE™

Hewlett Packard
Enterprise

Ingrasys®

Inventec

Lenovo

msi®

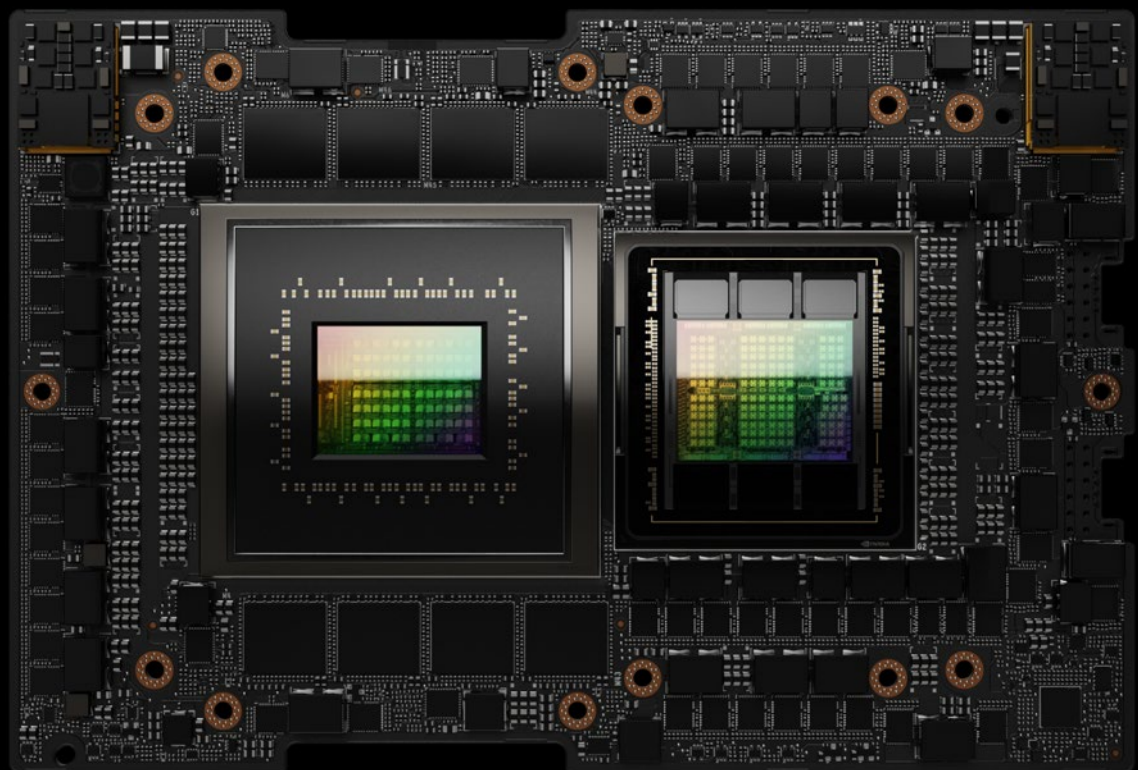
PEGATRON

QCT

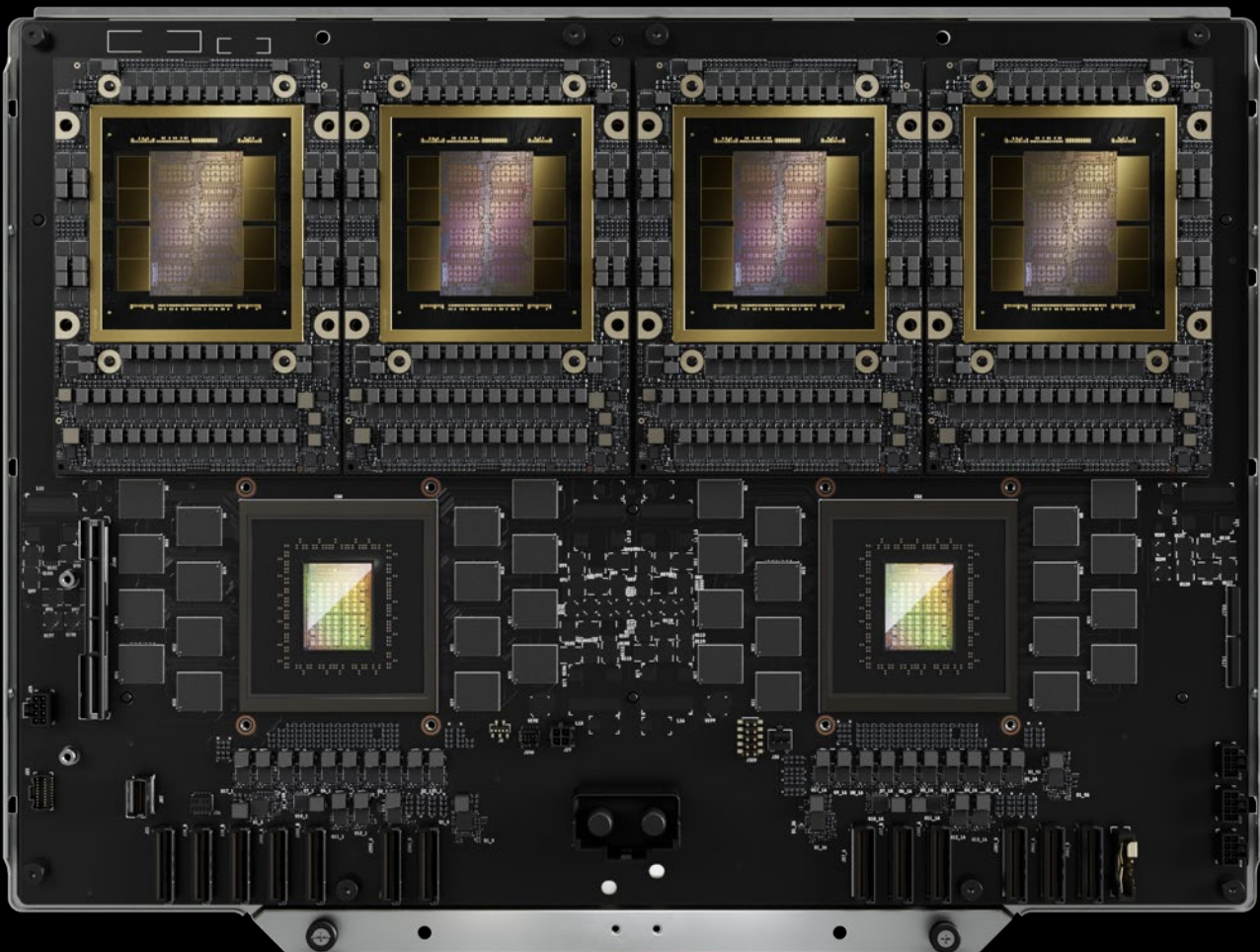
SUPERMICR®

wlstron®

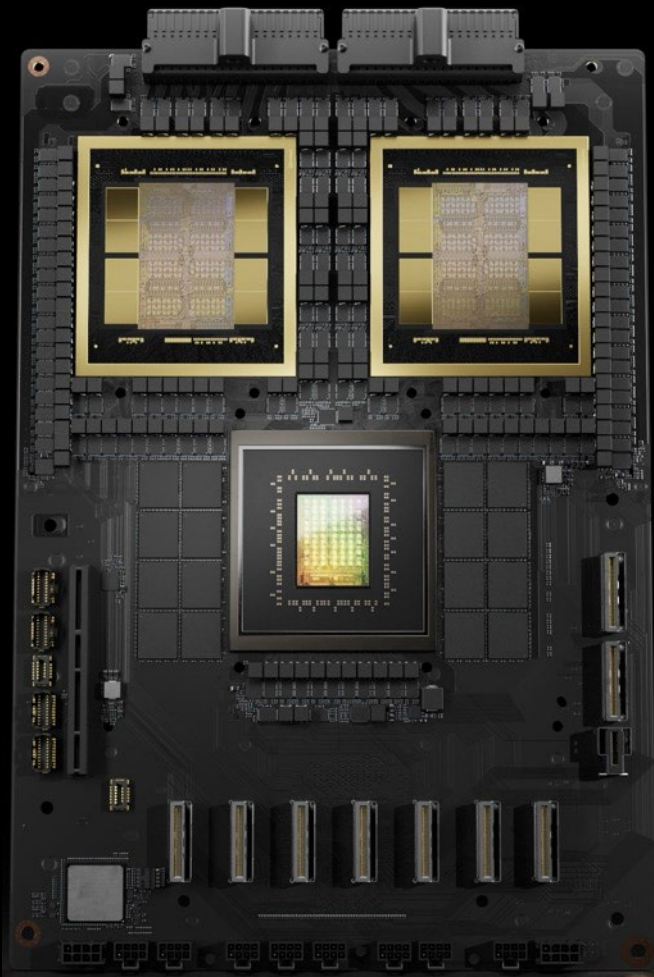
wiwynn®



Grace Hopper
Superchip



NVIDIA GB200 Grace Blackwell
NVL4 Superchip



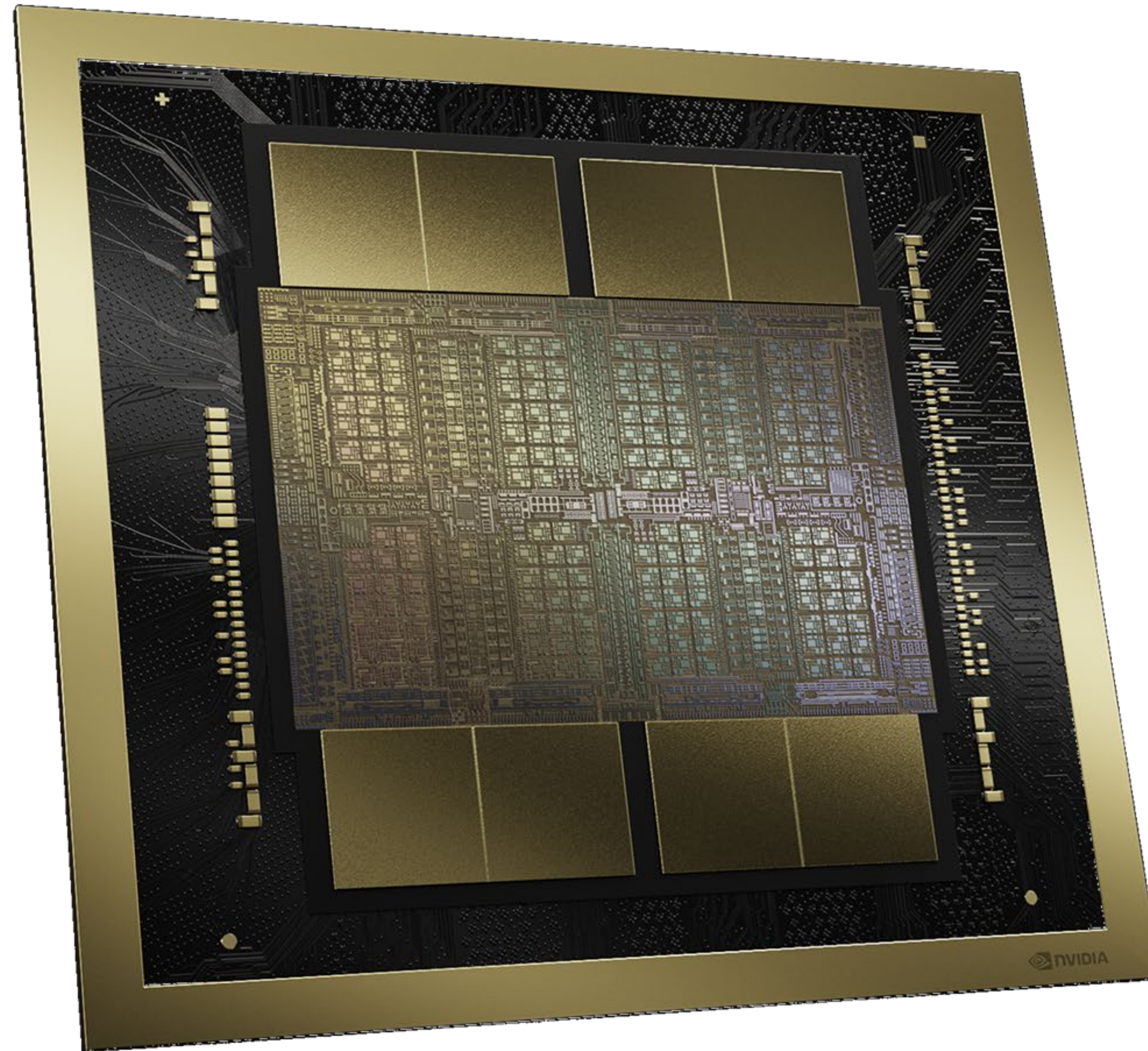
NVIDIA GB200 Grace Blackwell
Superchip



H200 NVL

NVIDIA Blackwell

The Engine of the New Industrial Revolution

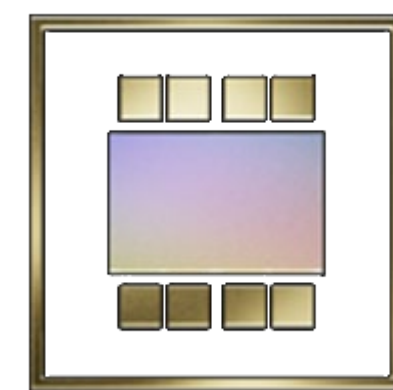


Built to Democratize Trillion-Parameter AI

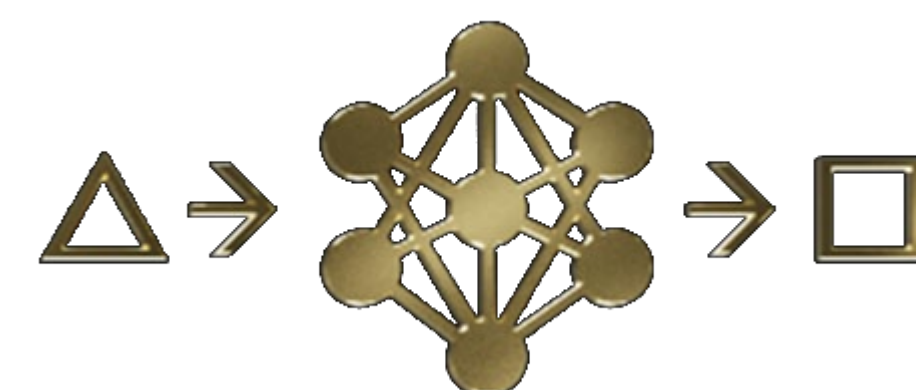
20 PetaFLOPS of AI performance on a single GPU

4X Training | 30X Inference | 25X Energy Efficiency & TCO

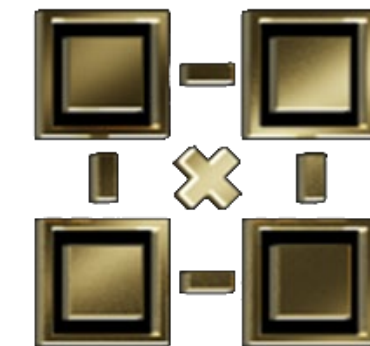
Expanding AI Datacenter Scale to beyond 100K GPUs



AI SUPERCHIP
208B Transistors



2nd GEN TRANSFORMER ENGINE
FP4/FP6 Tensor Core



5th GENERATION NVLINK
Scales to 576 GPUs



RAS ENGINE
100% In-System
Self-Test



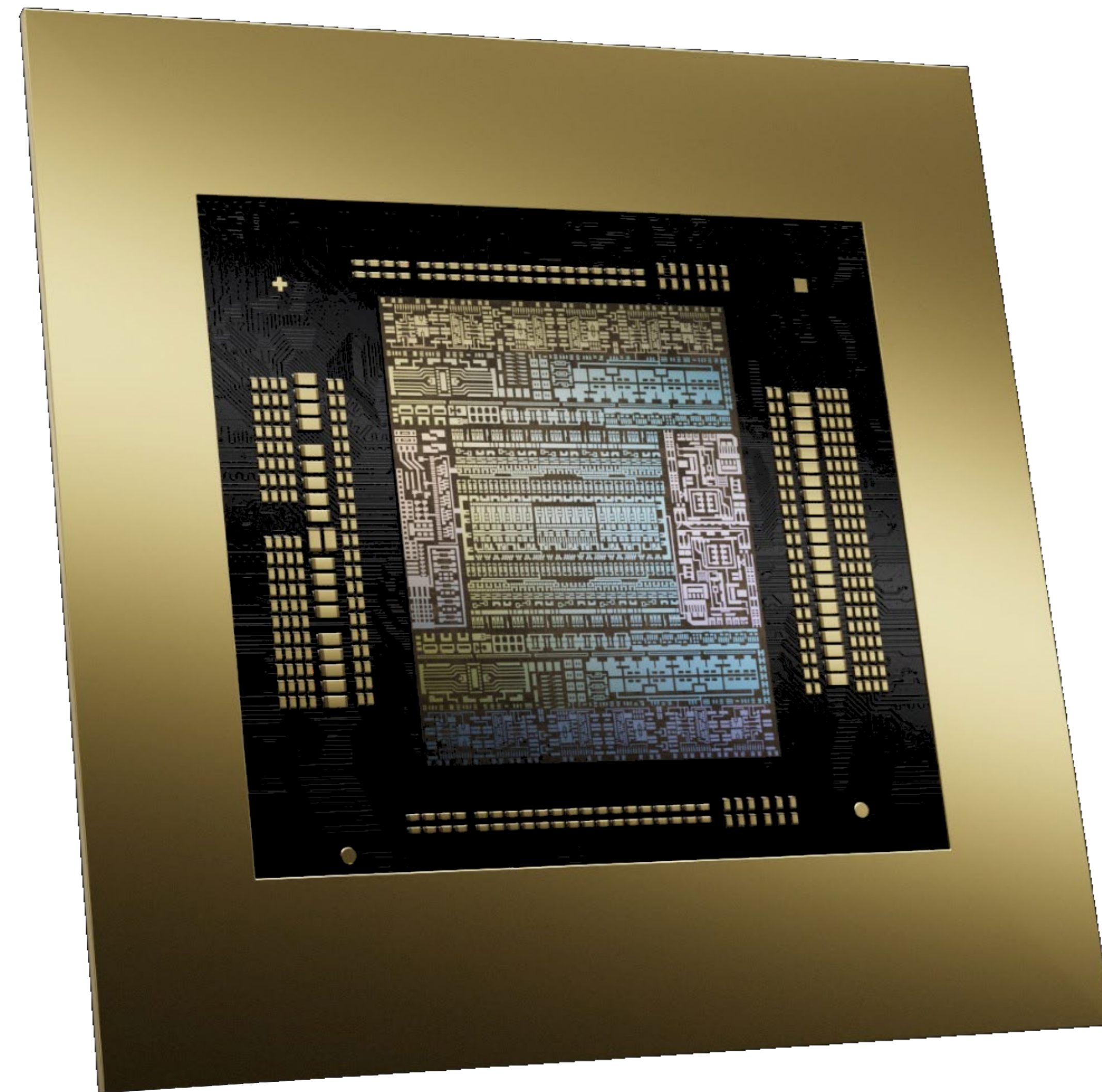
SECURE AI
Full Performance
Encryption & TEE



DECOMPRESSION ENGINE
800 GB/s

Announcing Fifth Generation NVLink and NVLink Switch Chip

Efficient Scaling for Trillion Parameter Models



7.2 TB/s Full all-to-all Bidirectional Bandwidth

Sharp v4 plus FP8

3.6 TF In-Network Compute

Expanding NVLink up to 576 GPU NVLink Domain

18X Faster than Today's Multi-Node Interconnect

GB200 NVL72

Delivers New Unit of Compute



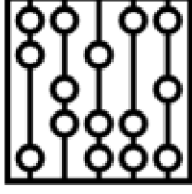

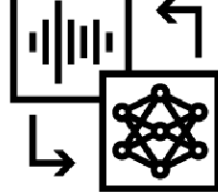
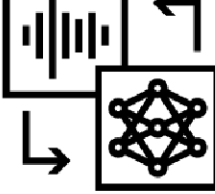


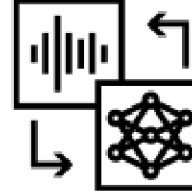


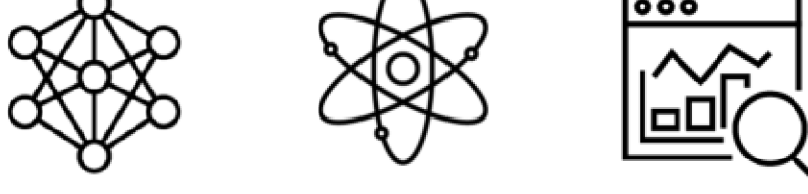


GB200 NVL72

- 36 GRACE CPUs
- 72 BLACKWELL GPUs
- Fully Connected NVLink Switch Rack

Training	720 PFLOPs
Inference	1,440 PFLOPs
NVL Model Size	27T params
Multi-Node All-to-All	130 TB/s
Multi-Node All-Reduce	260 TB/s

OEM and DGX
options

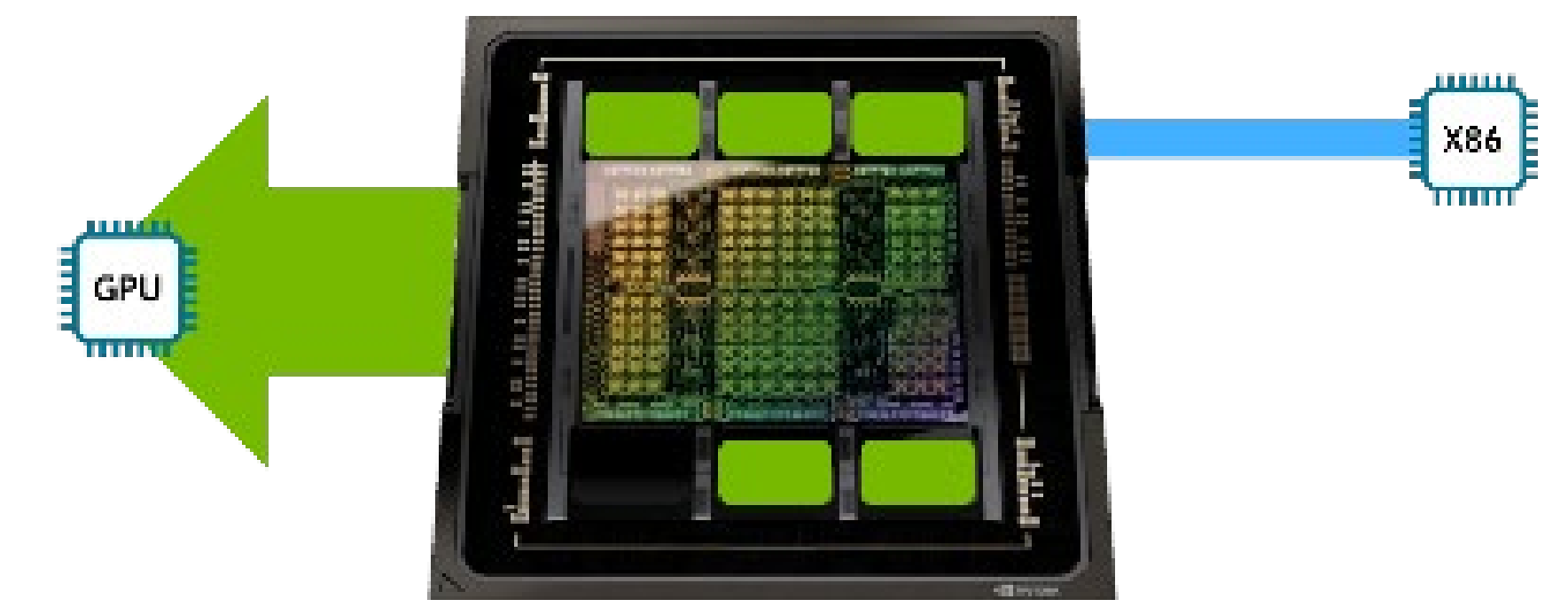
Choose The Right Solution For The Job

Scale-up, CPU+GPU & HGX Products		Scale-out, CPU+GPU & PCIe Products			
Real-Time Trillion-Parameter Models LLM & MoE	Highest Compute Performance HPC & AI	Most Versatile Platform Diverse Workloads	Most Efficient Compute LLM Inference, AI + HPC	Fast Universal AI + Graphics Text to Image/Video	AI Video Entry-tier Universal & Edge AI
 Graph Neural Networks	 AI Training, Inference & Scientific Research	 LLM Inference & Retrieval-Augmented Generation	 LLM Inference & Retrieval-Augmented Generation	 Text to Image/Video AI Multi-modal Generative AI	 Edge AI, Inference + Video and AI
 Massive Scale Model Training & Inference	 Data Analytics	 Data Analytics, Vector Database & HPC	 AI, HPC & Data Analytics	 Fine Tune Training/Inference, GenAI + Omniverse	 Mobile Cloud Gaming + vWS
405B - 1T+	70B-405B	70B-175B	70B-175B	13B-70B	Up to 7B
Rack Power: ~120 / 70kW	Node Power: ~11 / 5.5kW	Node Power: ~3kW	Node Power: ~8kW	Node Power: ~3kW	Node Power: ≤2kW
NVLink Domain 72	NVLink Domain 8 & 4	NVLink Domain 2	NVLink Domain 4 & 2	NA	NA
Max GPUs per NVLink Domain Max Performance & Capability	4-8 GPUs per Baseboard Highest Compute Performance AI, HPC & Data Analytics	1 – 2 GPUs per node Best Inference TCO Adv. Arch for AI & HPC	4-8 GPUs per node 2 nd Generation MIG 5-year NVIDIA AI Enterprise	Fastest RT Graphics Largest Render Models	1-16 GPUs per node Video & Graphics Compact & Versatile
GB200 NVL72 1000W, 480GB, 144GB	HGX H200 H100 700W, 141GB 80GB	GH200 NVL2 1000W, 480GB, 144GB	H200 NVL H100 NVL 600W, 141GB 400W, 94GB NVL2 & 4 NVL2	L40S 350W 48GB 2-Slot FHFL	L4 72W 24GB 1-Slot HHHL
Quantum-2/CX7	Quantum-2/ConnectX-7				
N/A	Spectrum-X / B3140H SuperNIC				
BlueField-3 DPU					
					Ethernet/Wi-Fi

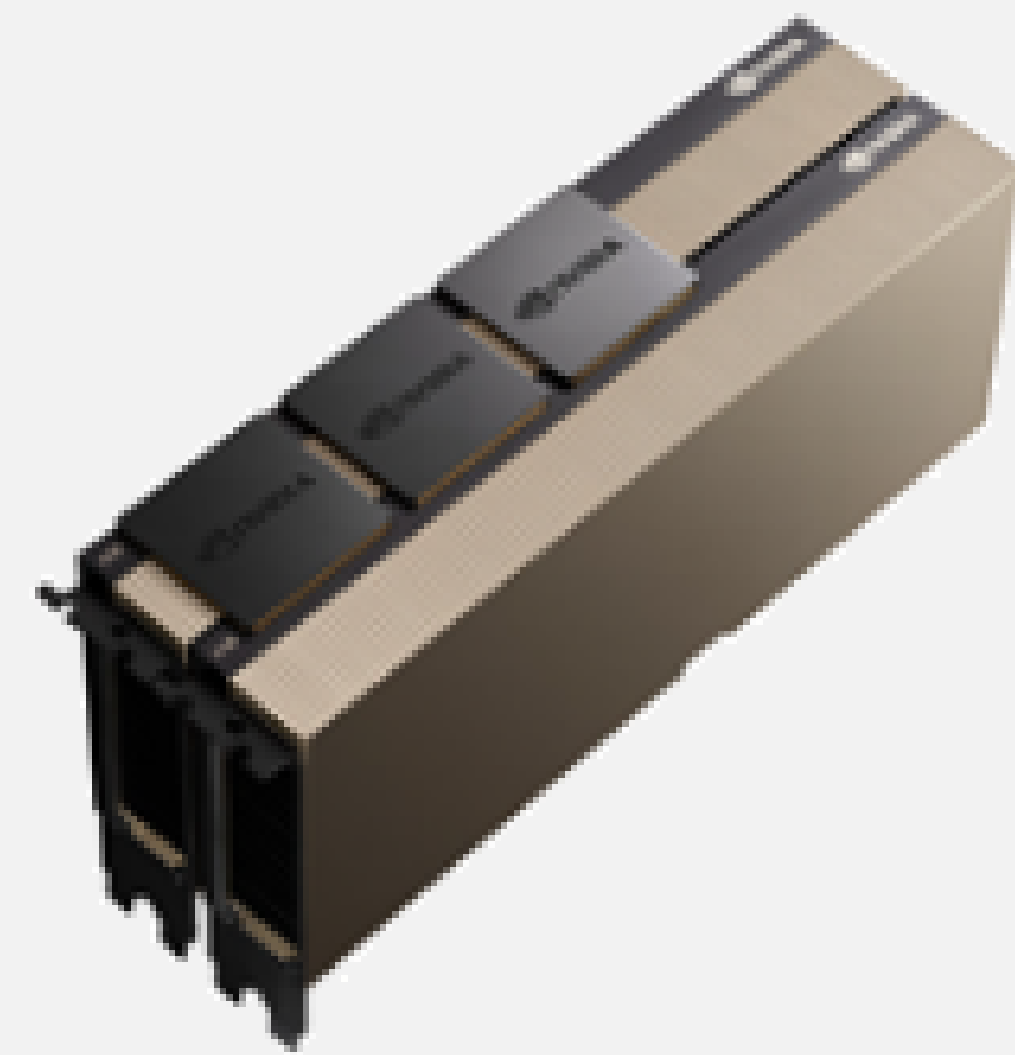
EW
InfiniBand
EW
Ethernet
NS
MGMT, IO

X86 + HOPPER

Architectures & Connectivity



PCIe



H100 up to NVL2
H200 up to NVL4

94GB HBM3
141GB HBM3e

NEW

HGX 4-way

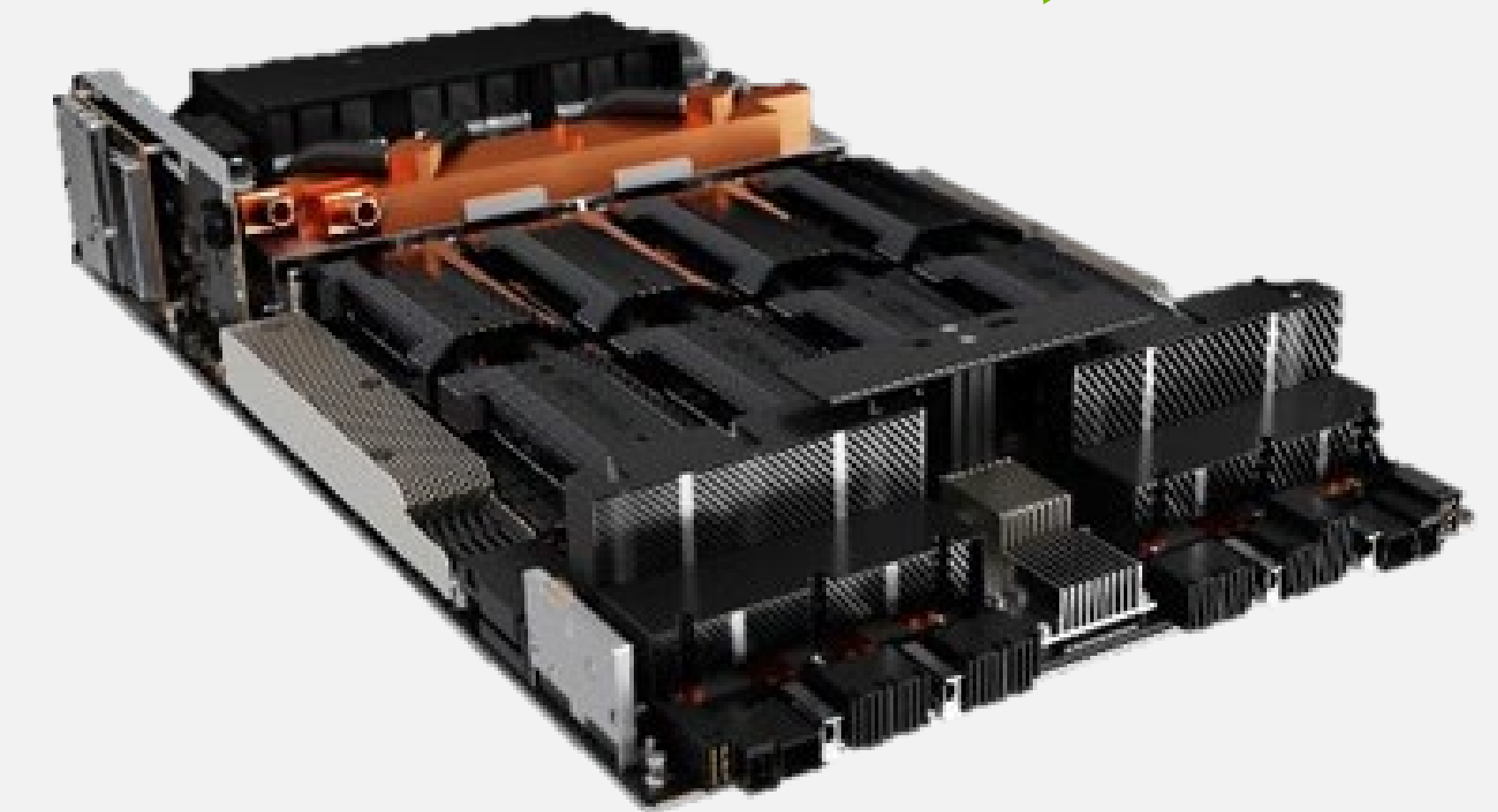


HGX H100 4-GPU
HGX H200 4-GPU

80GB HBM3 3.4TB/s
141GB HBM3e 4.8TB/s

NEW

HGX 8-way



OEM and DGX options

HGX H100 8-GPU
HGX H200 8-GPU

80GB HBM3 3.4TB/s
141GB HBM3e 4.8TB/s

NEW

NVIDIA GH200 Grace Hopper Superchip

Built for the New Era of Accelerated Computing and Generative AI

Most versatile compute

Best performance across CPU, GPU or memory intensive applications

Easy to deploy and scale out

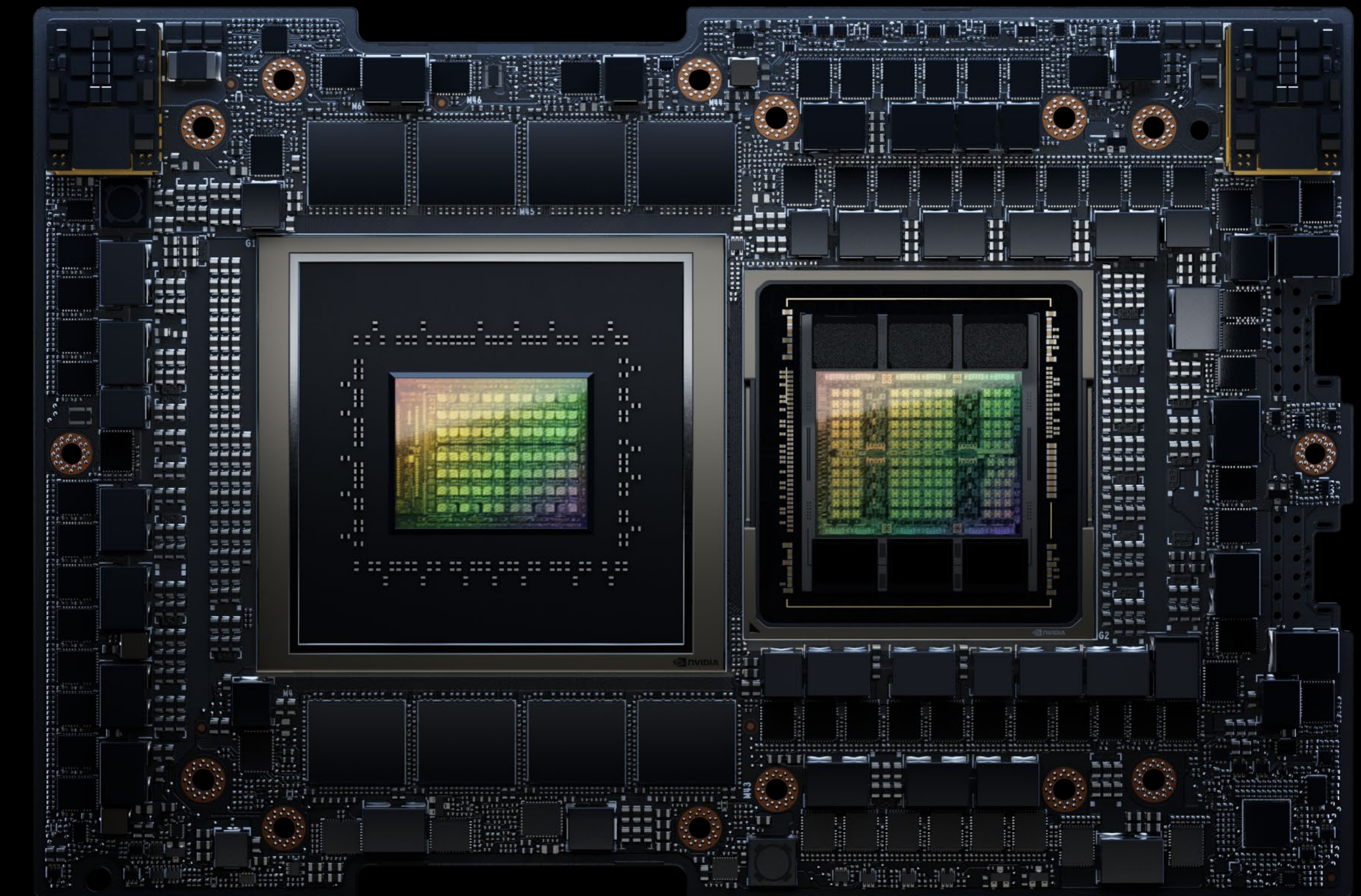
1 CPU:1 GPU node simple to manage and schedule for for HPC, enterprise, and cloud

Best Perf/TCO for diverse workloads

Maximize data center utilization and power efficiency

Continued Innovation

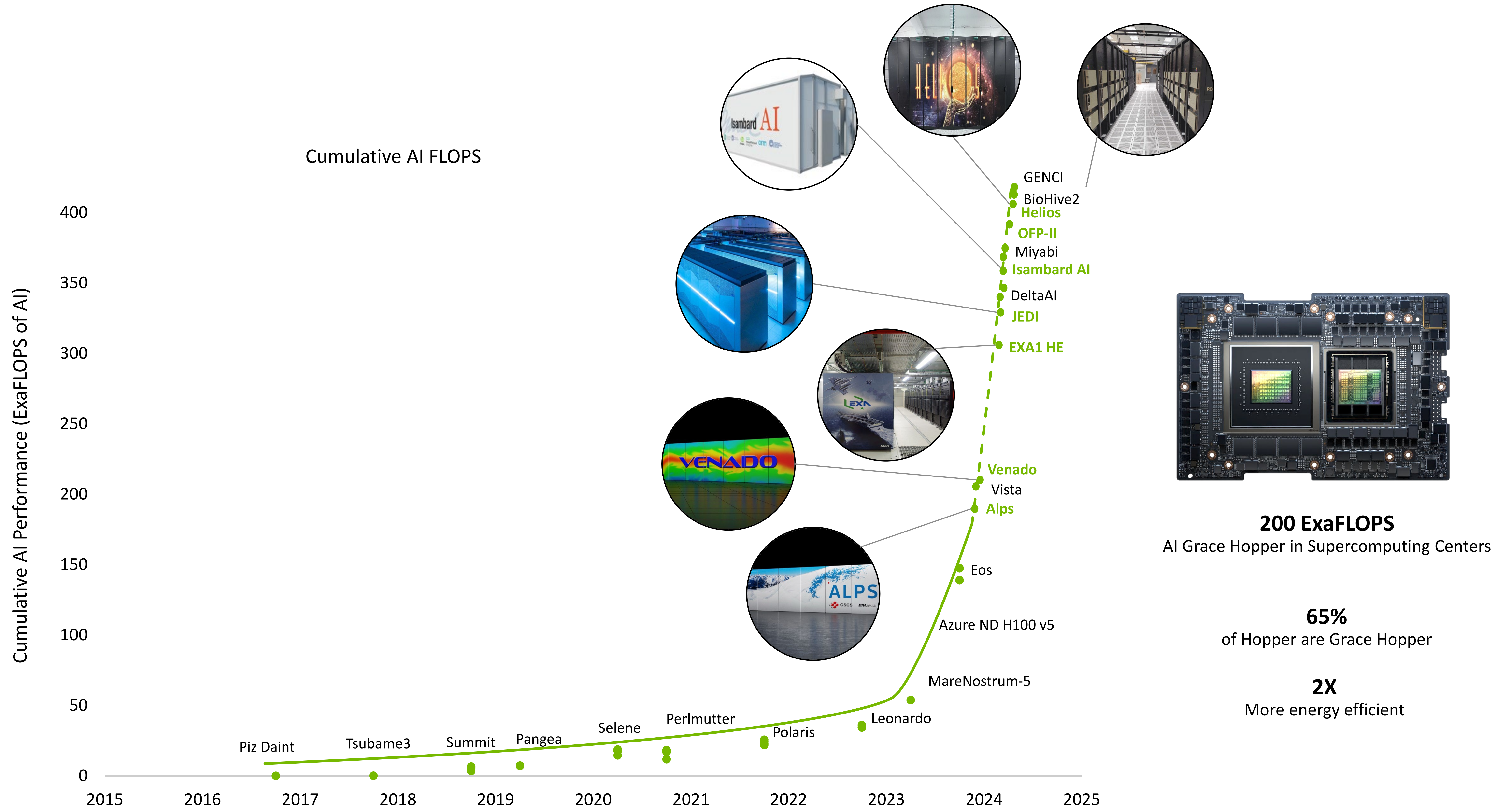
Grace and Blackwell in 2025



900GB/s NVLink-C2C | 624GB High-Speed Memory
4 PF AI Perf | 72 Arm Cores

Grace Hopper Powers AI Supercomputing Datacenters

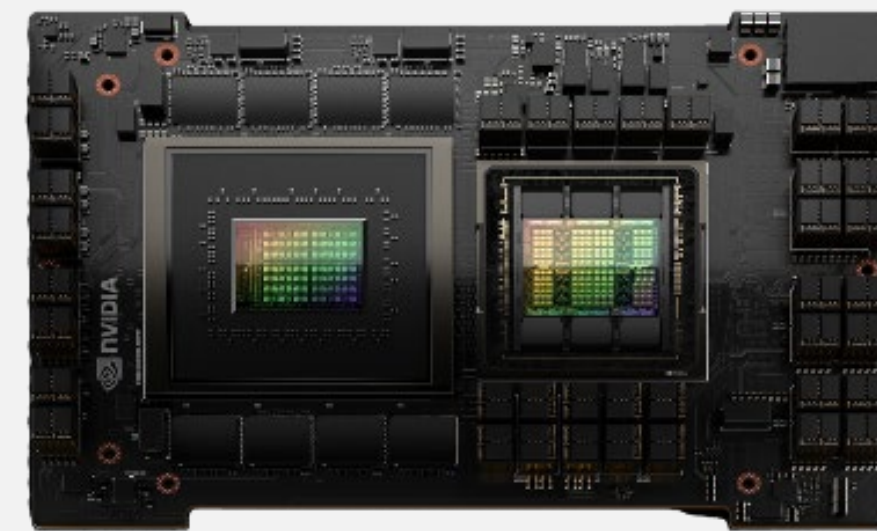
Grace Hopper Will Deliver 200 Exaflops of AI performance for Groundbreaking Research



GRACE GPU-GPU NVLINK

Architectures & Cost of Connectivity

1-Way
1:1

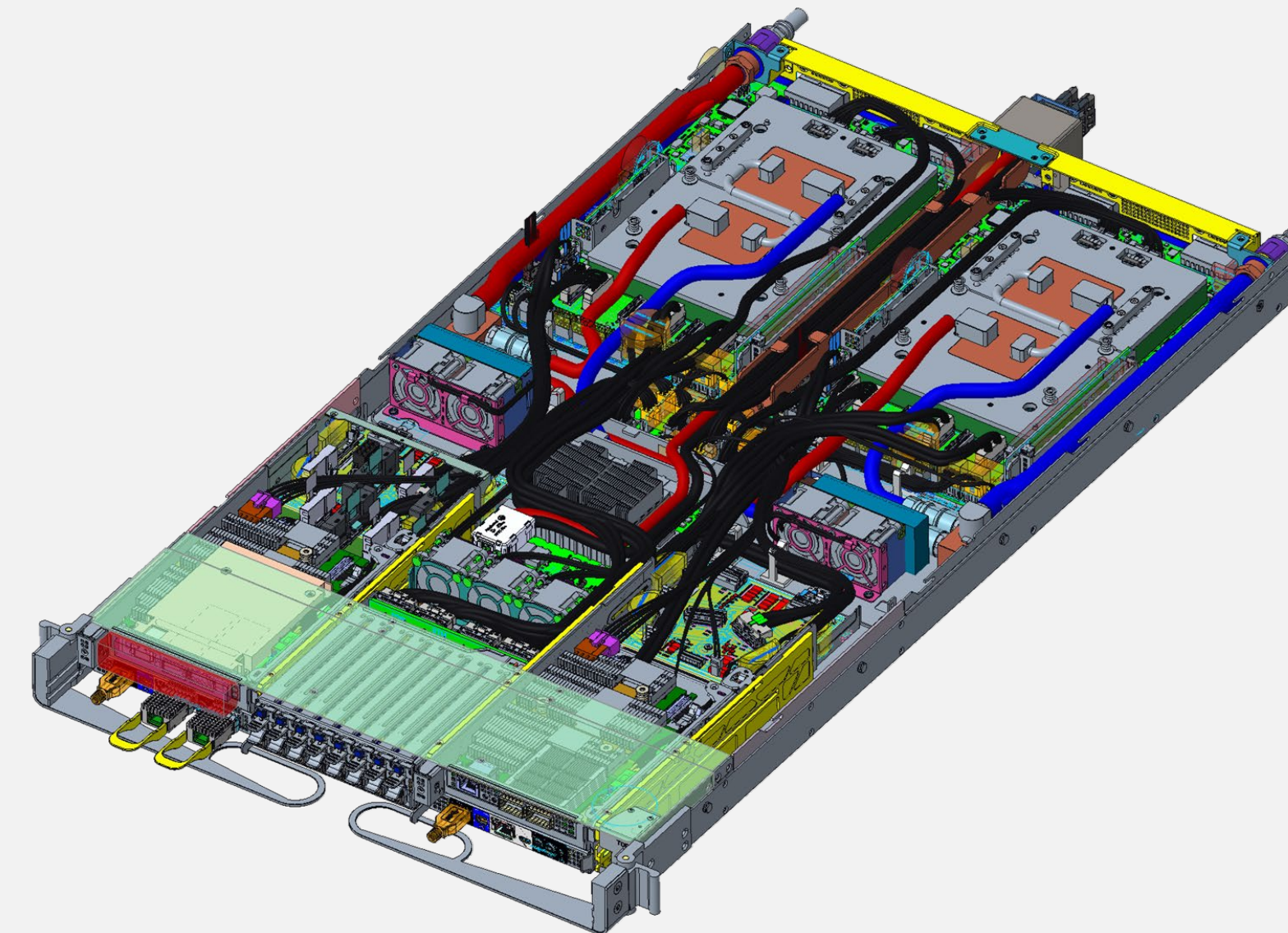


GH200 96GB

480 GB LP5x
96GB HBM3

Scale Out
AI Inference

2-Way
2:2

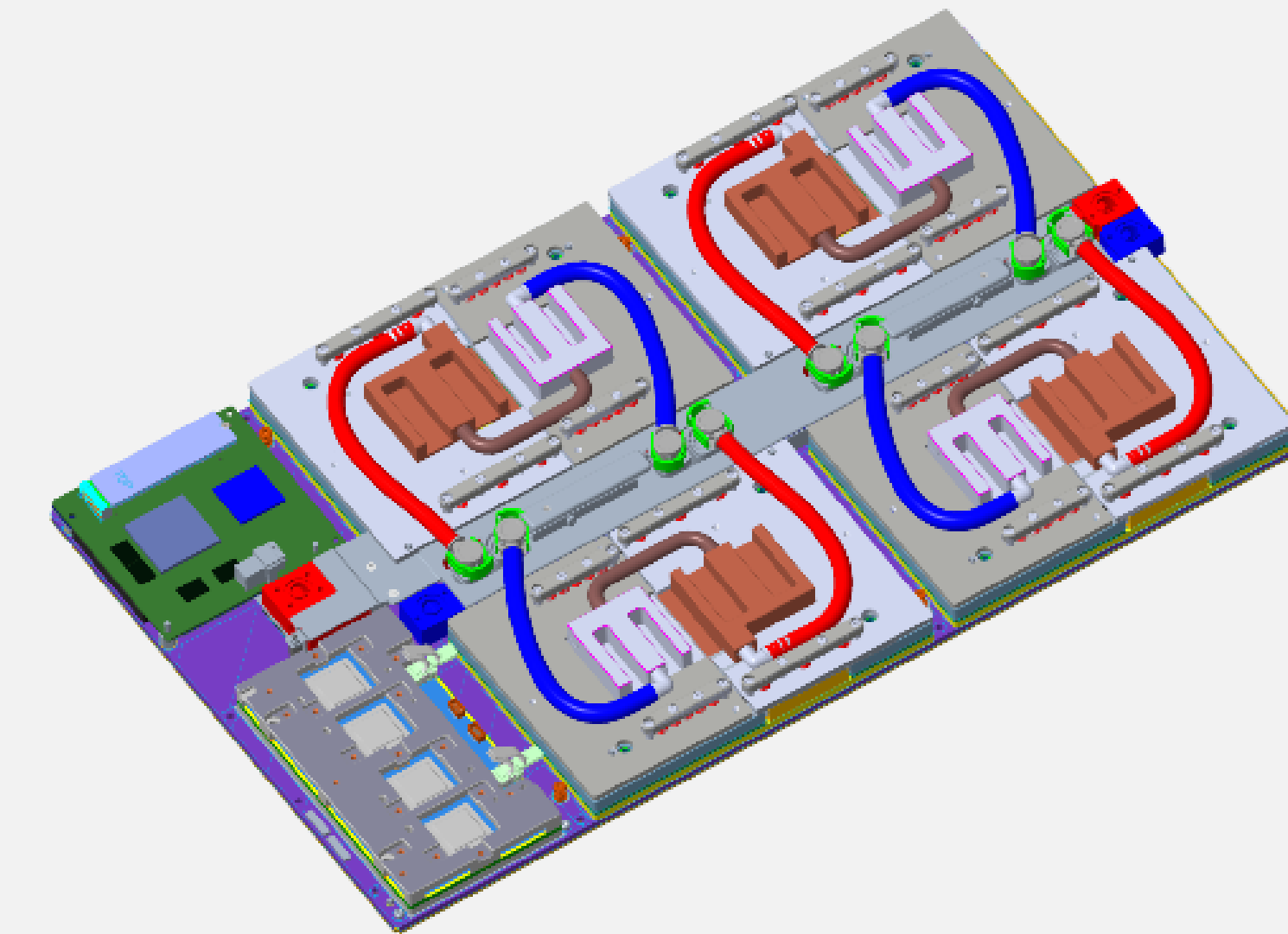


GH200 NVL2

480 GB LP5x
144GB HBM3e

Scale Out
AI Inference

4-Way
4:4



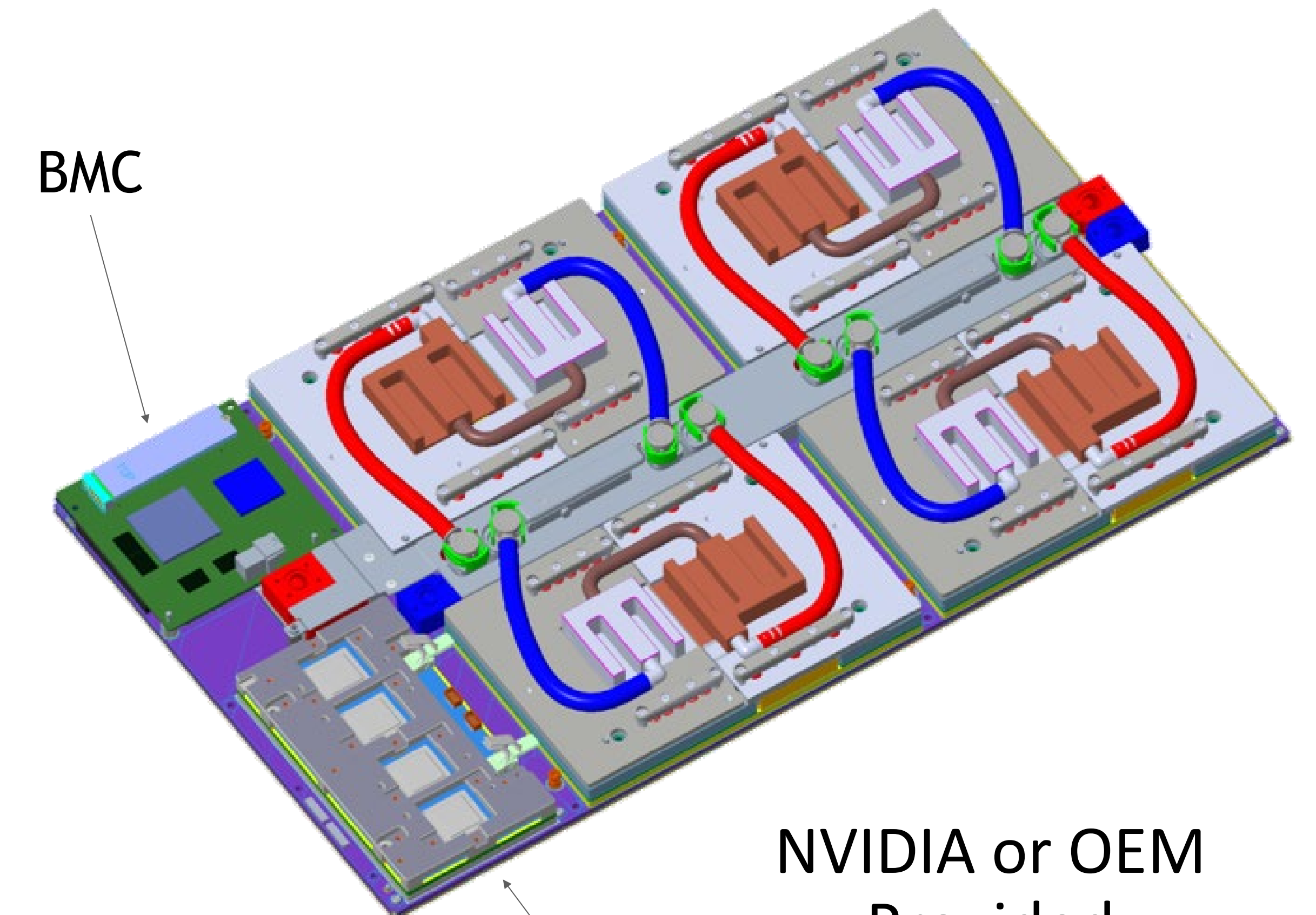
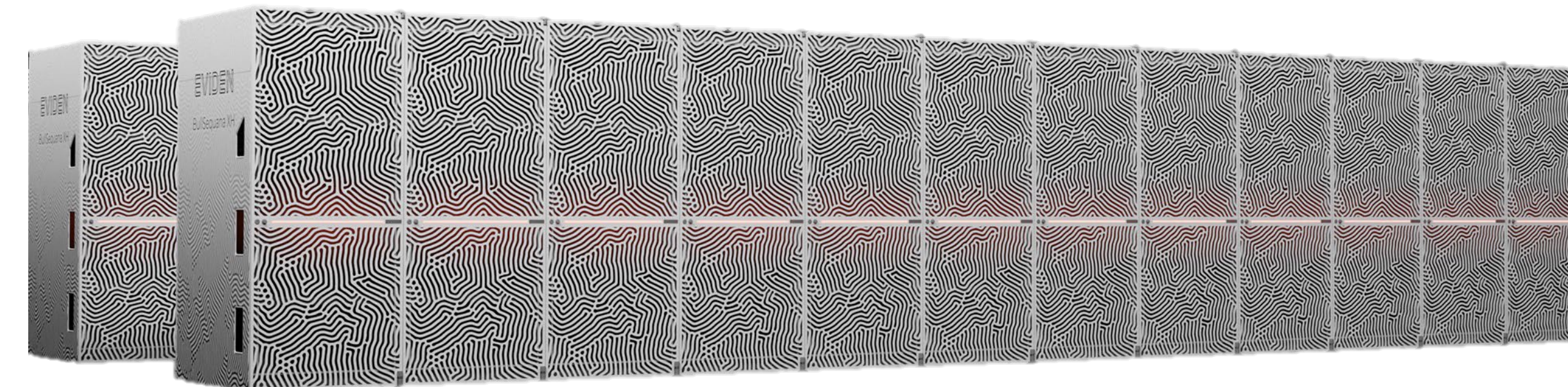
GH200 96GB

120 GB LP5x
96GB HBM3

Scale Out
Exascale HPC & AI Training

Grace Hopper Superchip 4-Way Design

The choice for the world's fastest supercomputers



BMC

NVIDIA or OEM
Provided
Baseboard

4X CX7


**Hewlett Packard
Enterprise**

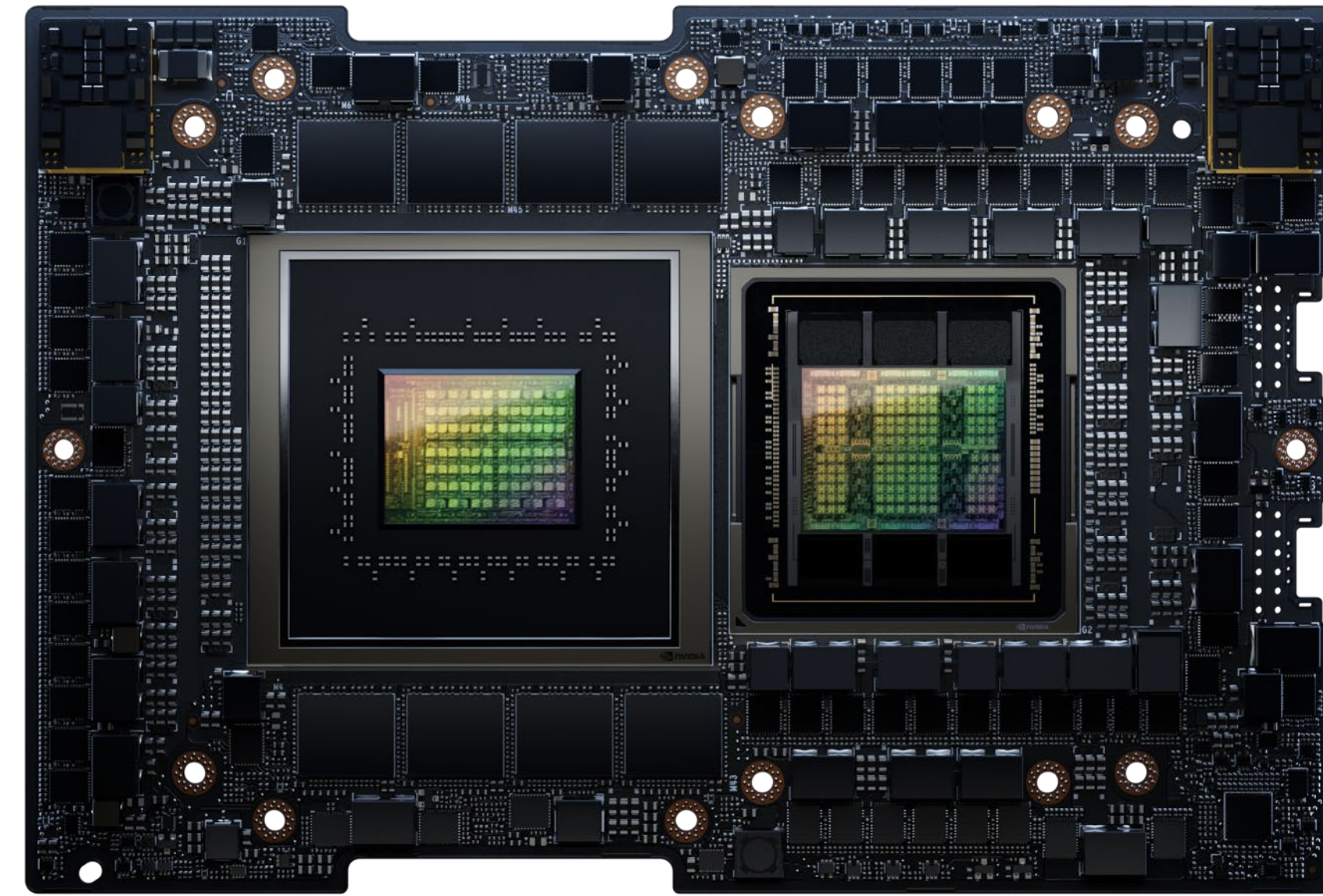
EVIDEN
an atos business

GH200 Grace Hopper HPC Platform

Unified Memory and Cache Coherence for Next Gen HPC Performance

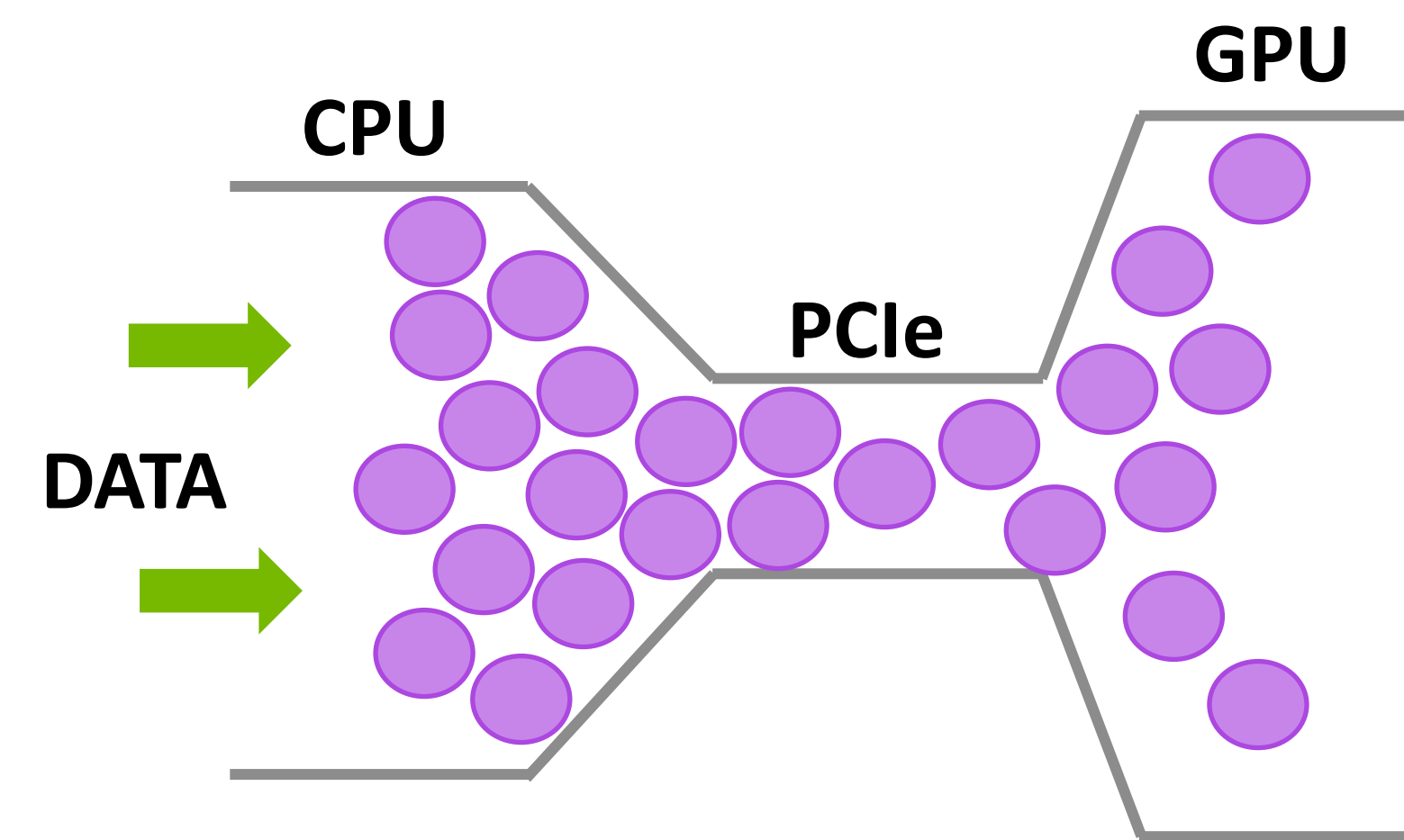
Partially GPU Accelerated Apps

Big performance gains with no code changes



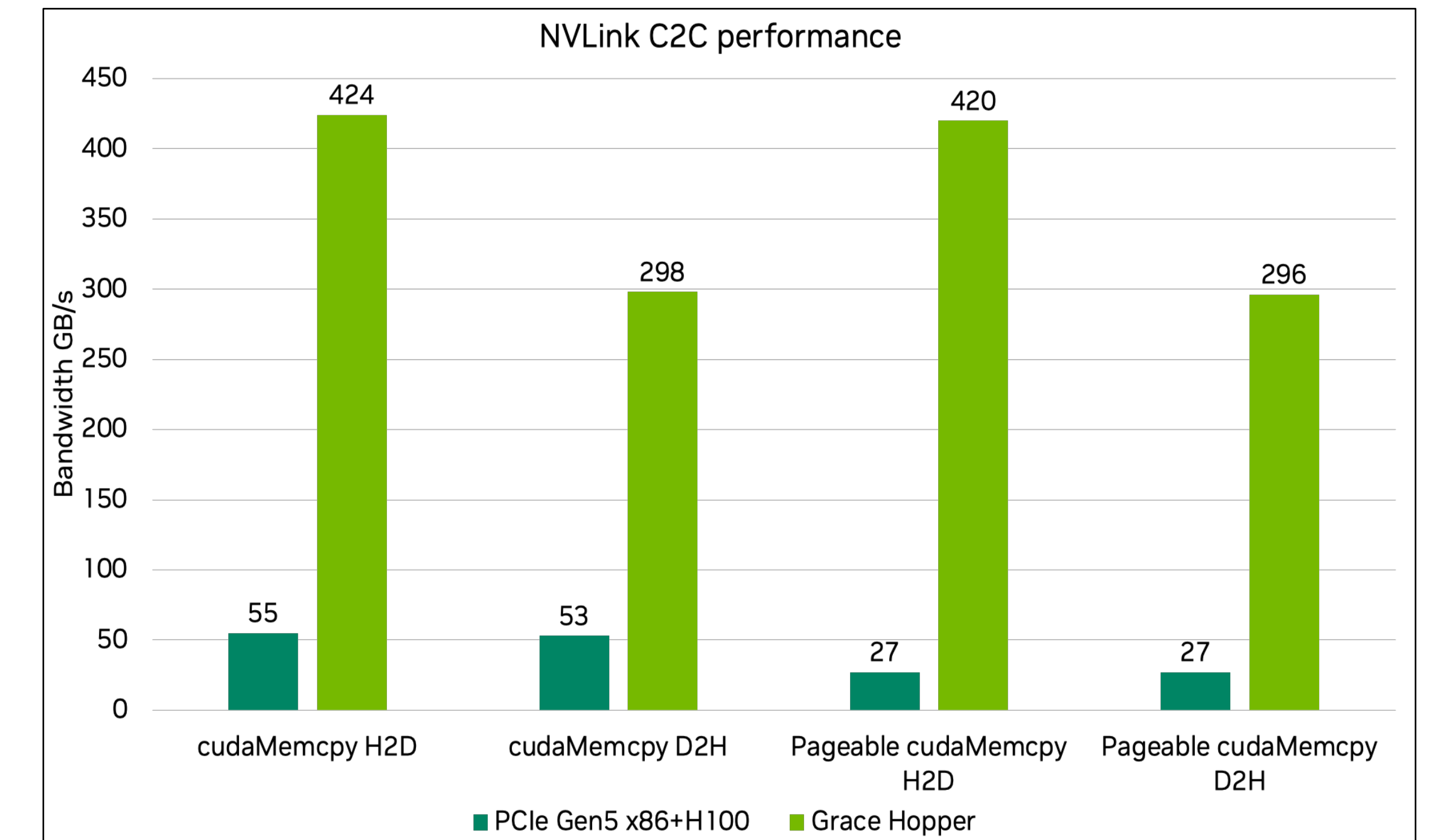
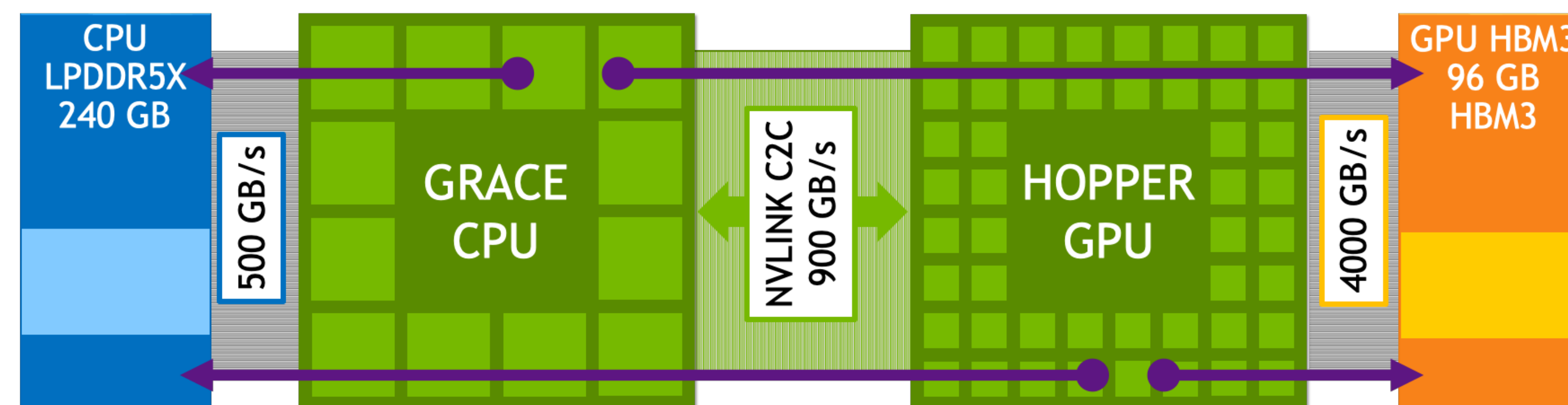
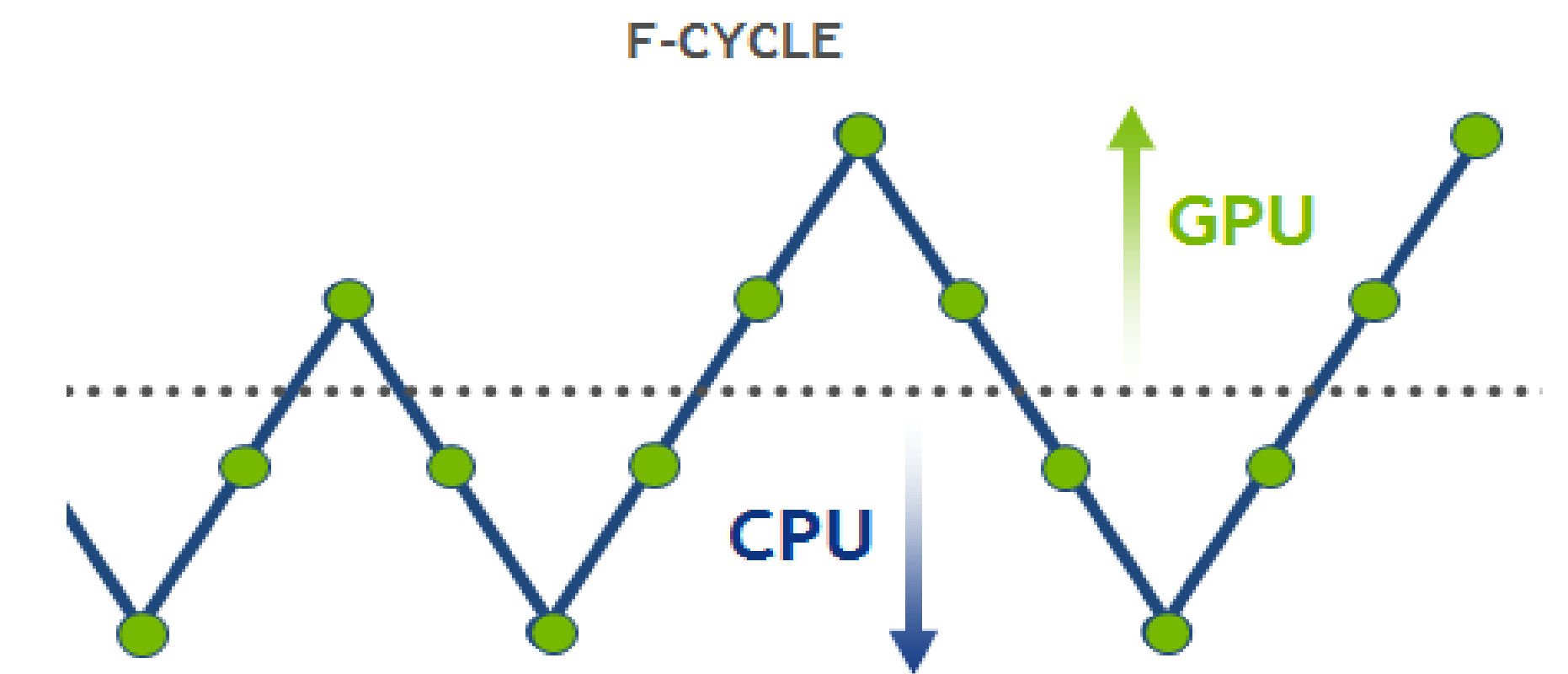
No More PCIe Bottleneck

NVLink-C2C is 7X PCIe BW

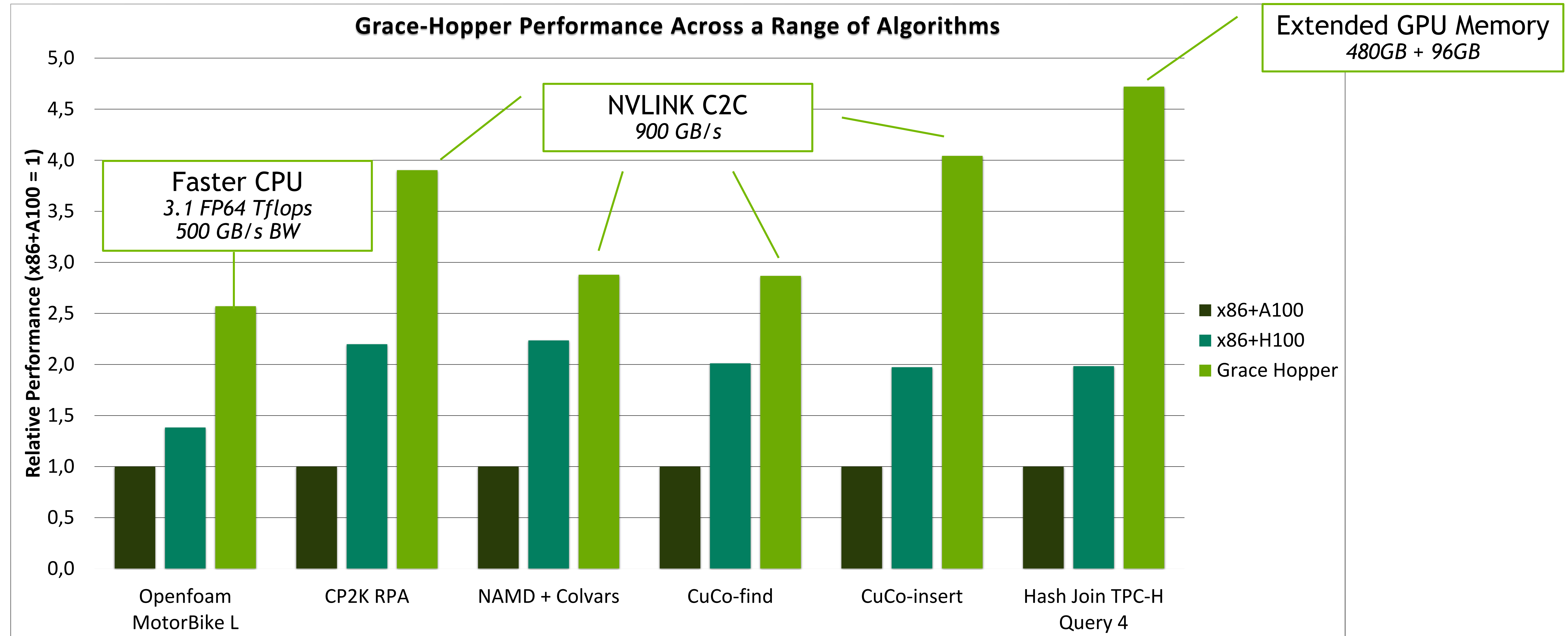


CPU & GPU Cache Coherence

Incremental code changes yield big gains



Grace-Hopper Superchip Workload Performance



The background features a series of parallel, wavy lines in various shades of green, creating a sense of depth and movement. A solid green vertical bar is positioned on the far left side of the image.

Hopper architecture

HOPPER H100 TENSOR CORE GPU

80B Transistors, TSMC 4N

2nd Gen Multi-Instance GPU
Confidential Computing
PCIe Gen5

New Memory System
World's First HBM3 DRAM
Larger 50MB L2

132 SMs 2x Performance per Clock
4th Gen Tensor Core
Thread Block Clusters

4th Gen NVLink 900GB/s total BW
New SHARP support
NVLink Network



NEW HOPPER SM ARCHITECTURE

- 2x faster FP32 & FP64 FMA
- 256 KB L1\$ / Shared Memory
- New 4th Gen Tensor Core
- New DPX instruction set
- New Tensor Memory Accelerator
 - Fully asynchronous data movement
- New Thread Block Clusters
 - Turn locality into efficiency

SM



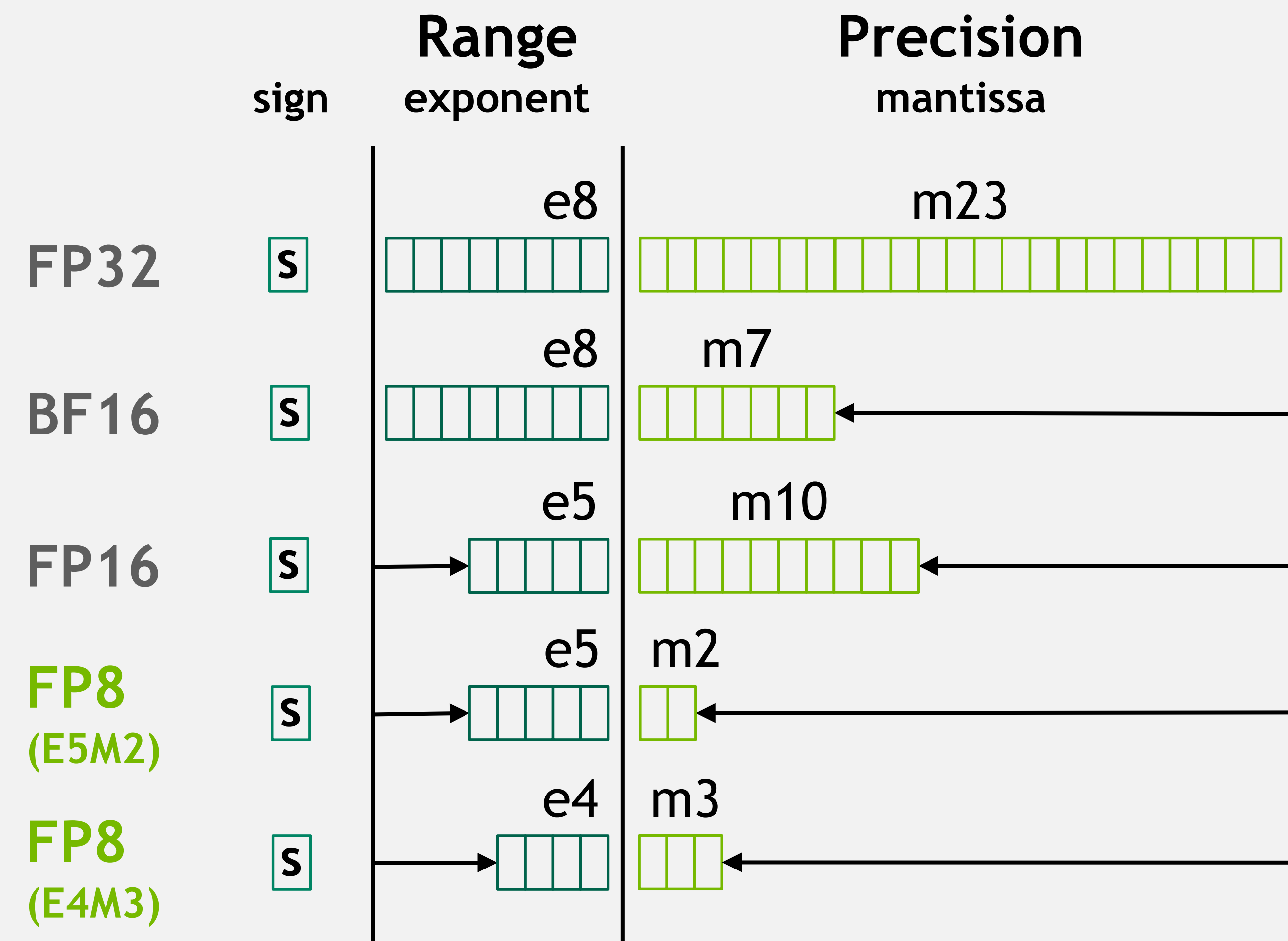
Tensor Cores

- Specialized high-performance compute cores for matrix multiply and accumulate (MMA) math operations for AI and HPC applications.
- Operating in parallel across SMs in one NVIDIA GPU deliver massive increases in throughput and efficiency compared to standard floating-point (FP), integer (INT), and fused multiply-accumulate (FMA) operations.
- Support for a wide range of data types (fp64, fp32, tf32, fp16, bfloat16, fp8, int8) and mixed precision
- New Transformer Engine designed specifically to accelerate Transformer model training and inference (chooses dynamically between FP8 and 16-bit calculations)
- Tensor Memory Accelerator feeds the H100 Tensor Cores with transfers large blocks of data and multi-dimensional tensors from global memory to shared memory and vice-versa.

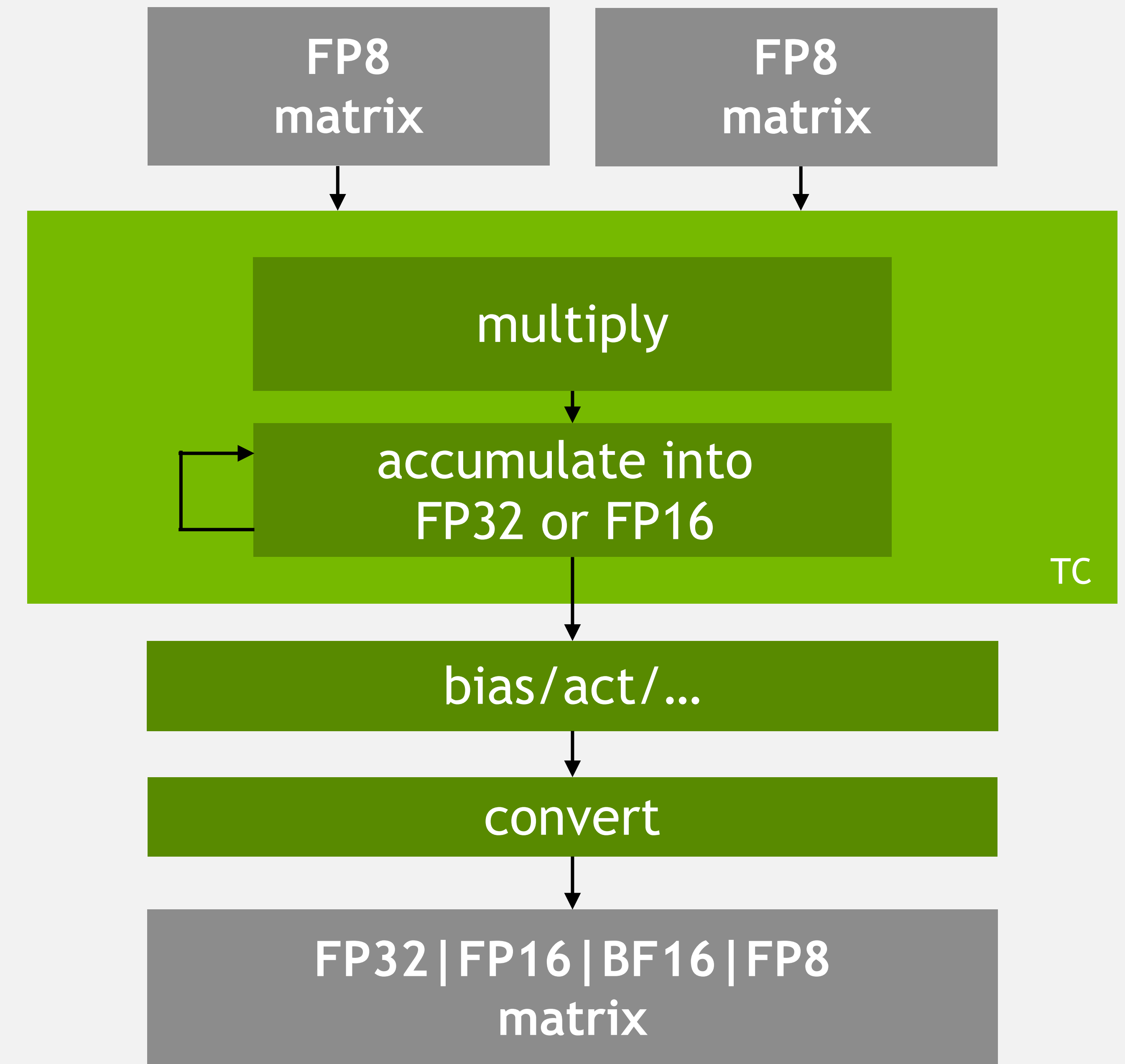
FP8 TENSOR CORE

Allocate 1 bit to either range or precision

Support for multiple accumulator and output types

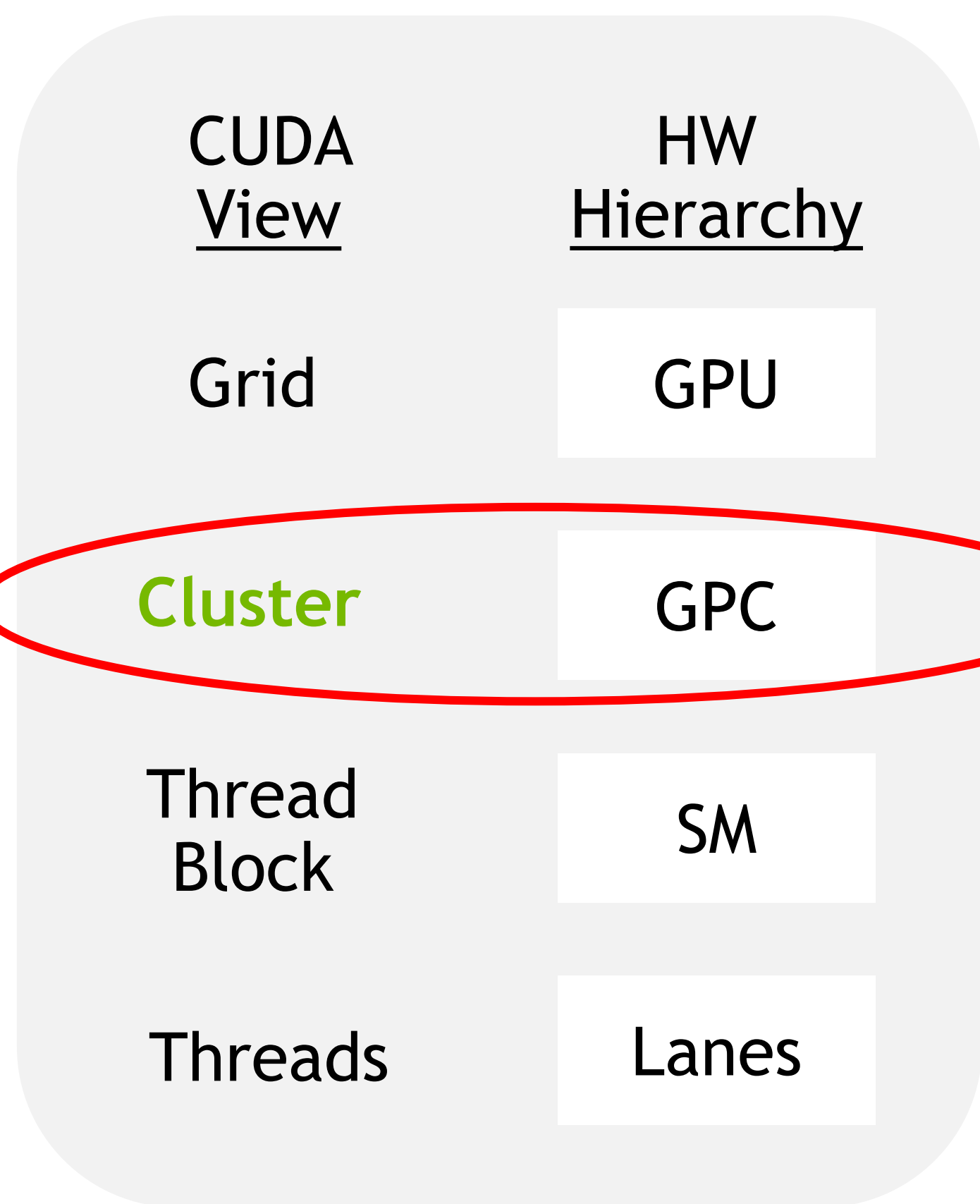


2x throughput & half footprint of FP16/BF16

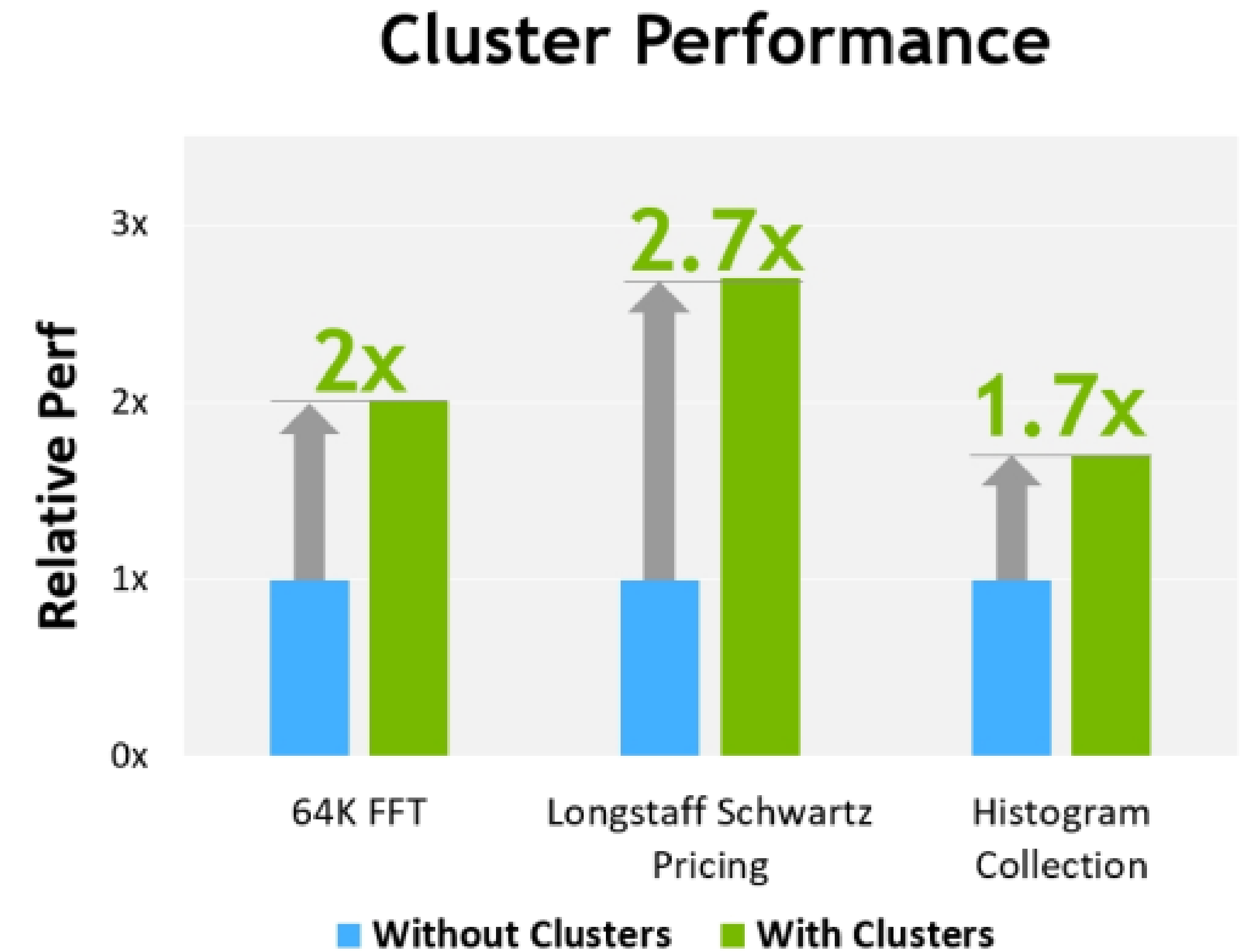
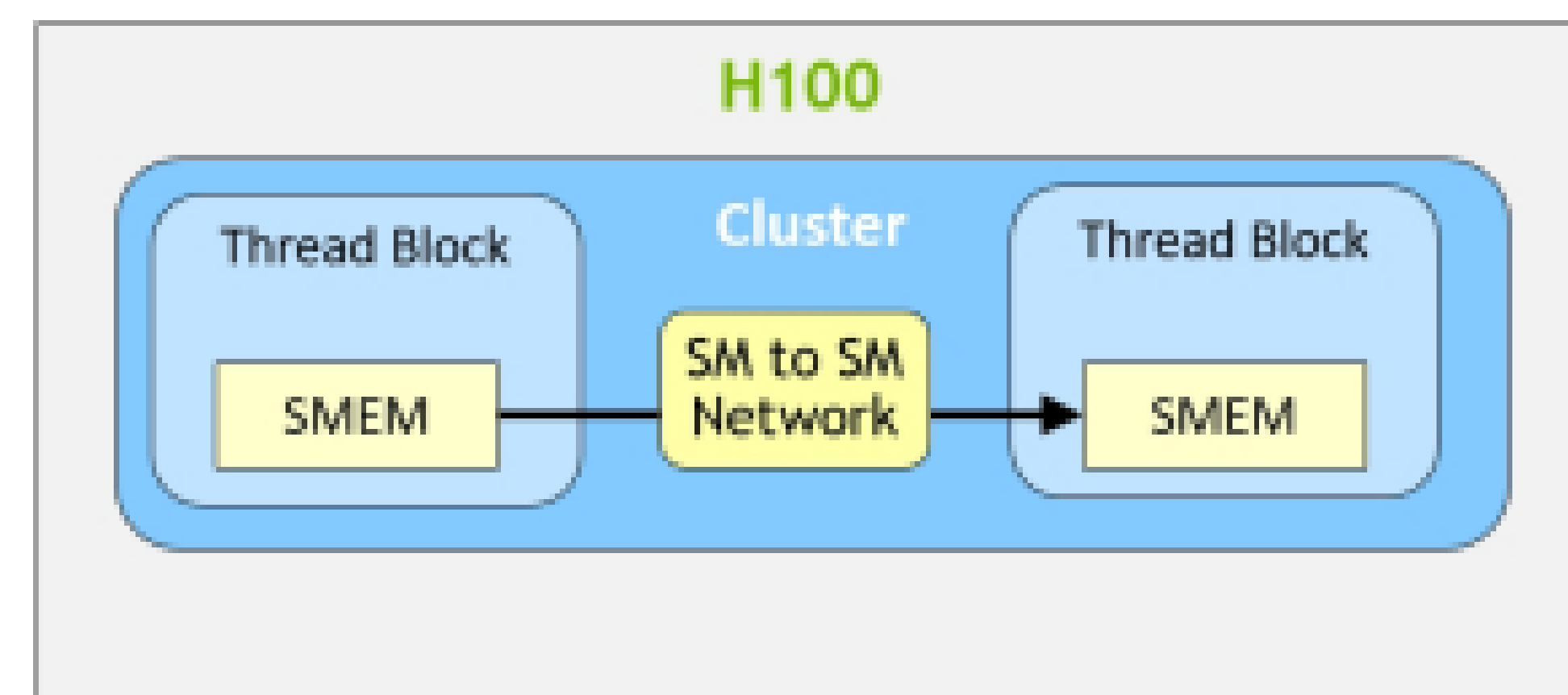
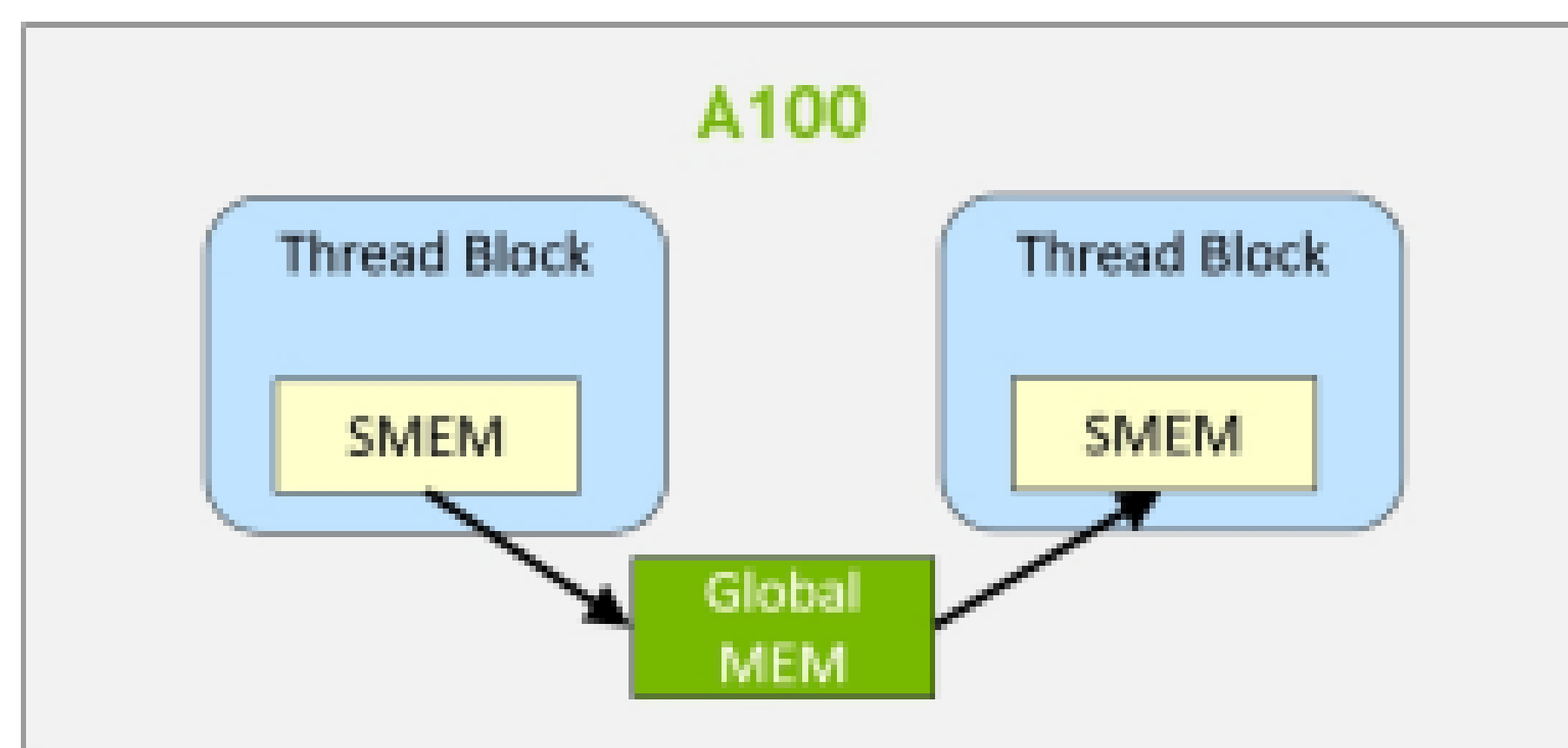


SM

Thread Block Clusters



- New feature introduces programming locality within clusters of SMs
- Shared memory blocks of SMs within a GPU Processing Cluster (GPC) can communicate directly (w/o going to HBM)
- Leveraged with CUDA cooperative groups API



For details, see “NVIDIA H100 Tensor Core GPU Architecture” white paper available for download

TENSOR MEMORY ACCELERATOR UNIT

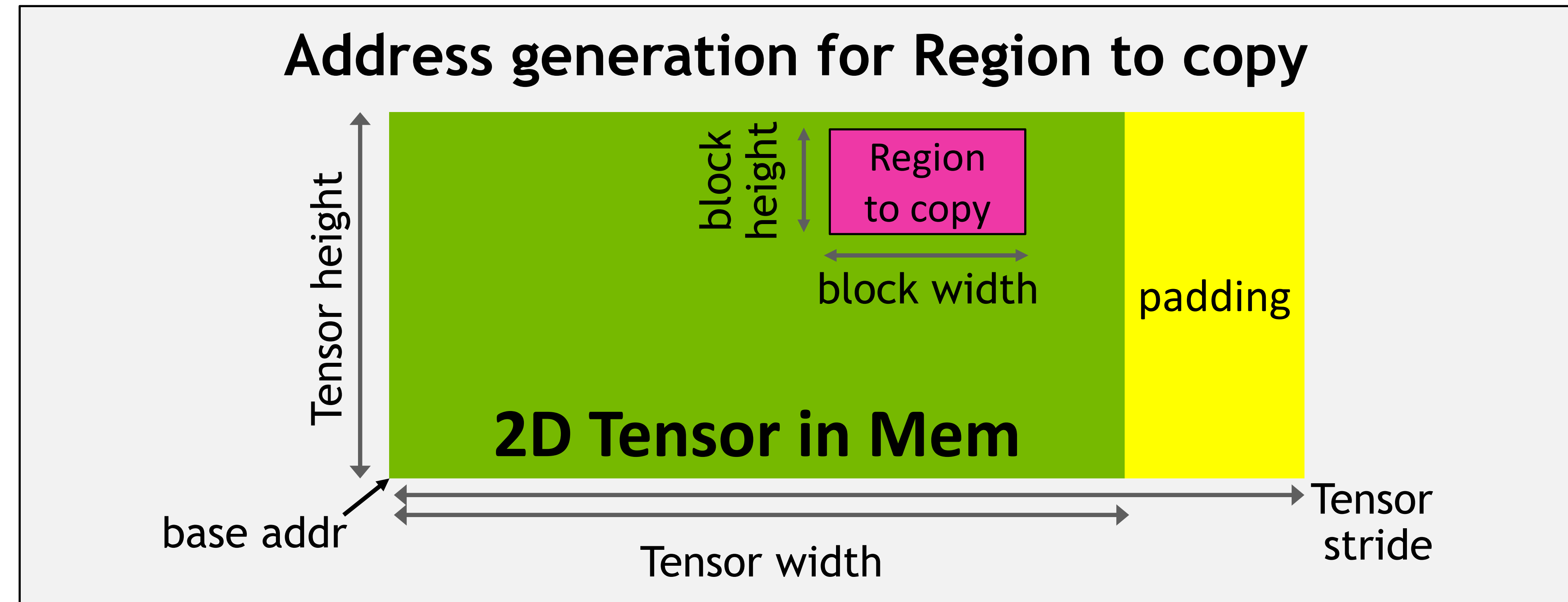
ASYNC MEM COPY USING TMA

HW-accelerated mem_copies

Global \Leftrightarrow Shared Mem

Shared Mem \Leftrightarrow Shared Mem for Clusters

Address generation for 1D to 5D Tensors

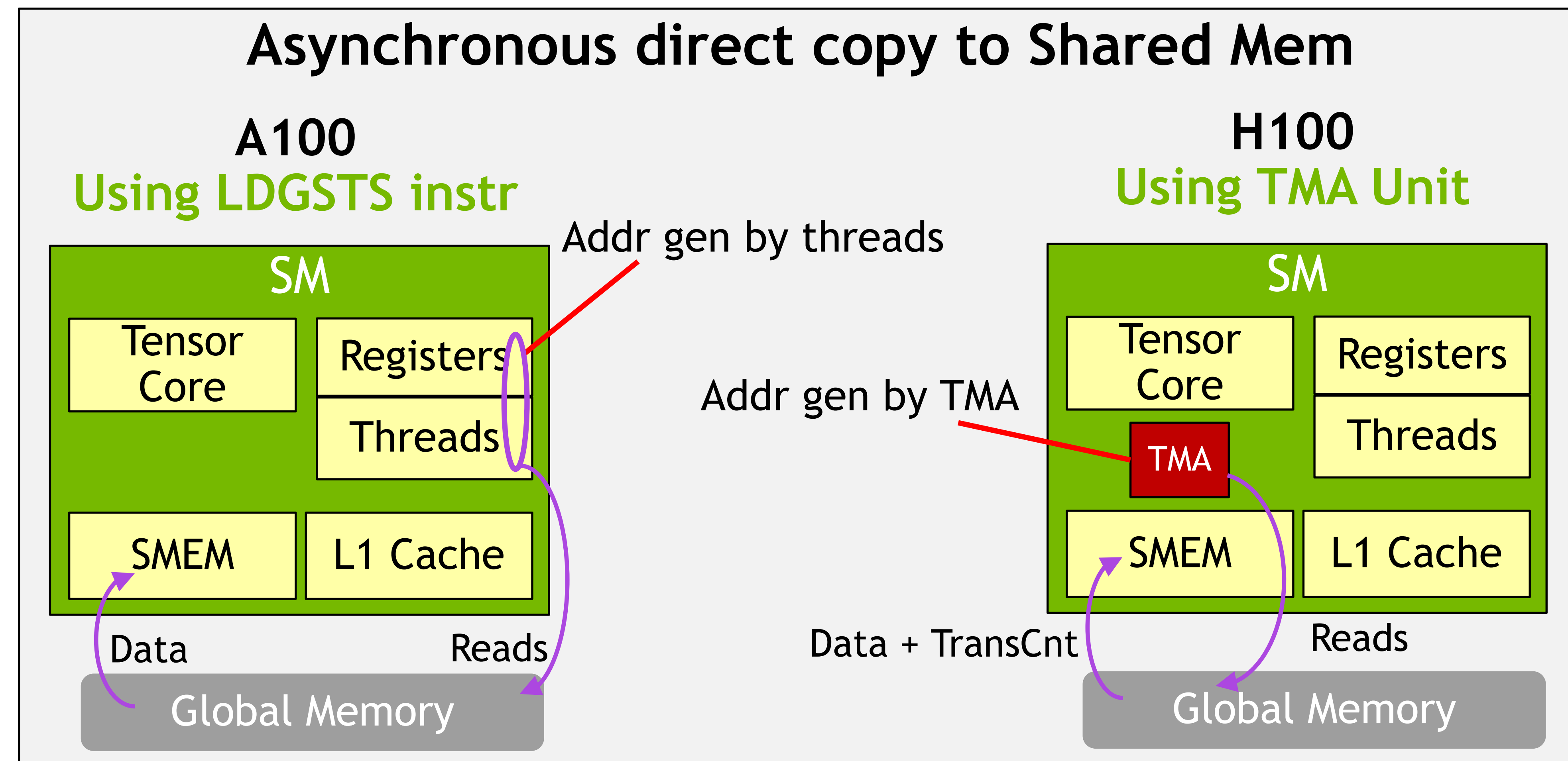


Simplified programming model

Fully asynchronous with threads

No addr gen or data movement overhead

Synchronize with transaction barrier





Software



CUDA

CUDA: NVIDIA's Computing Platform

Used Everywhere

Speech



Visual Search



Robotics



VALOSSA
Video Analysis

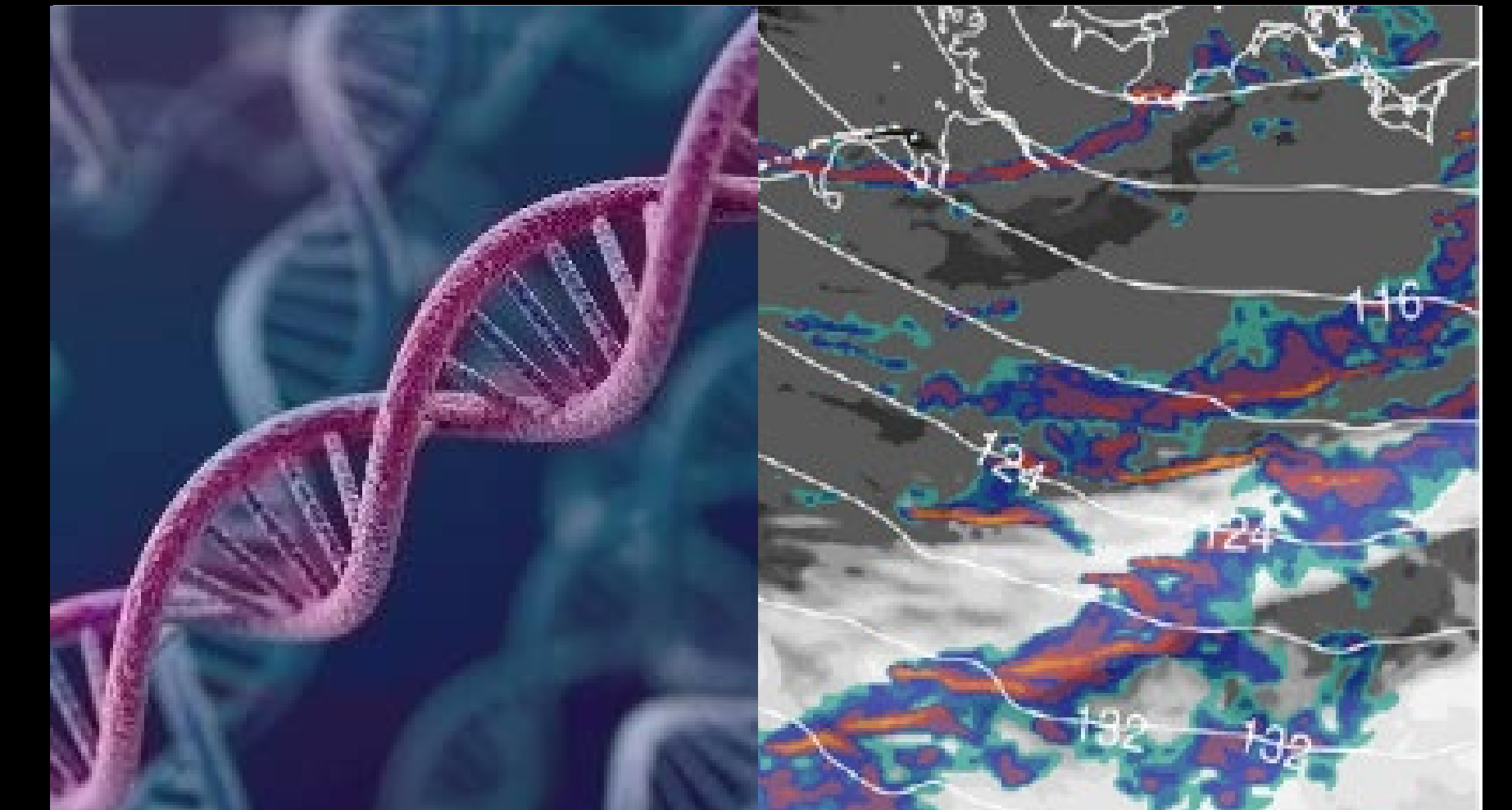
AI Training & Inference



Gaming



Media & Entertainment



Science

<http://developer.nvidia.com/cuda-downloads>

CUDA TOOLKIT

Libraries, Languages and Development Tools for GPU Computing

Programming Approaches

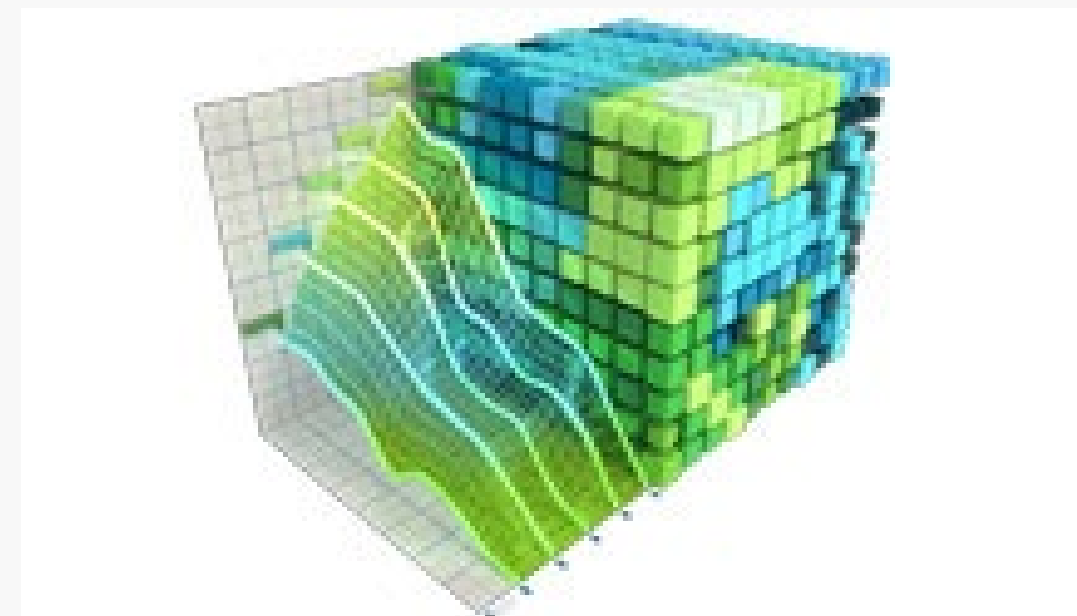
Libraries

“Drop-in” Acceleration

Programming Languages

Maximum Flexibility

Development Environment



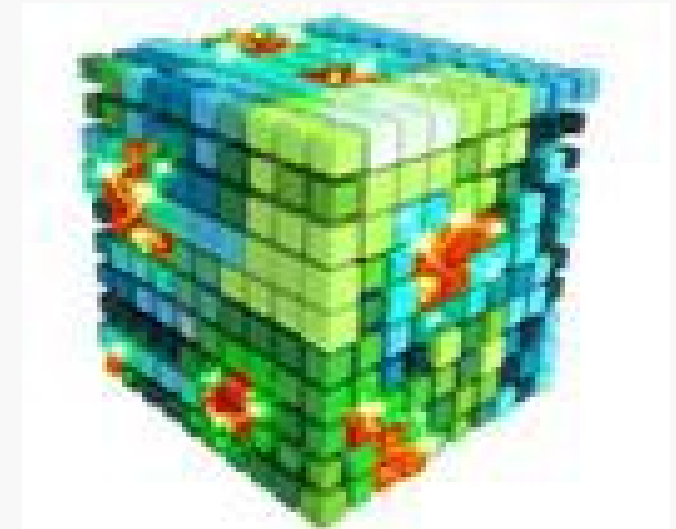
NVIDIA
Visual Profiler



CUDA Profiling
Tools Interface



CUDA-GDB
Debugger



CUDA
MEMCHECK

Language Support

C

C++

Fortran

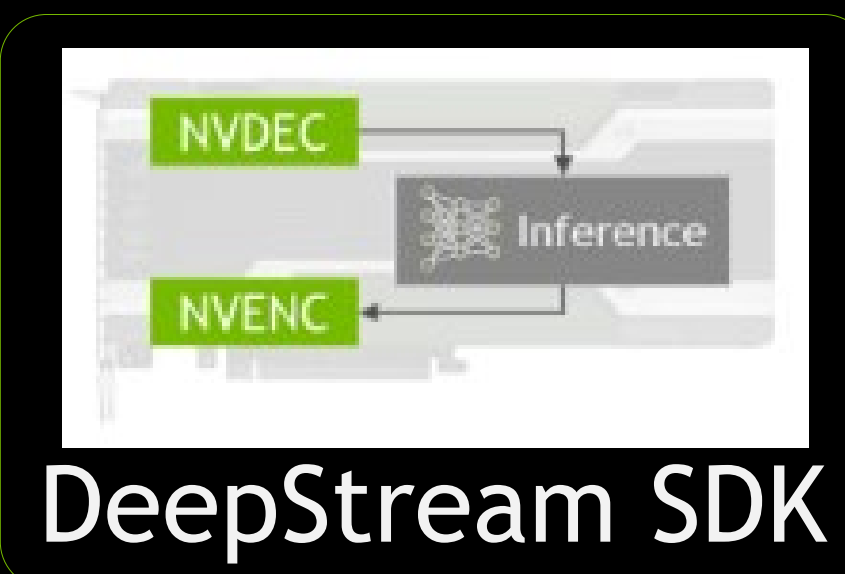
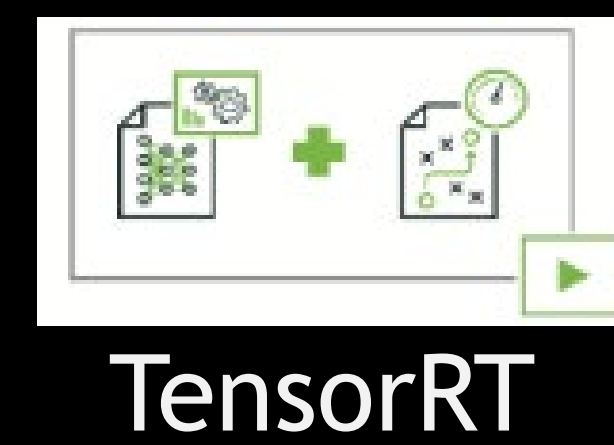
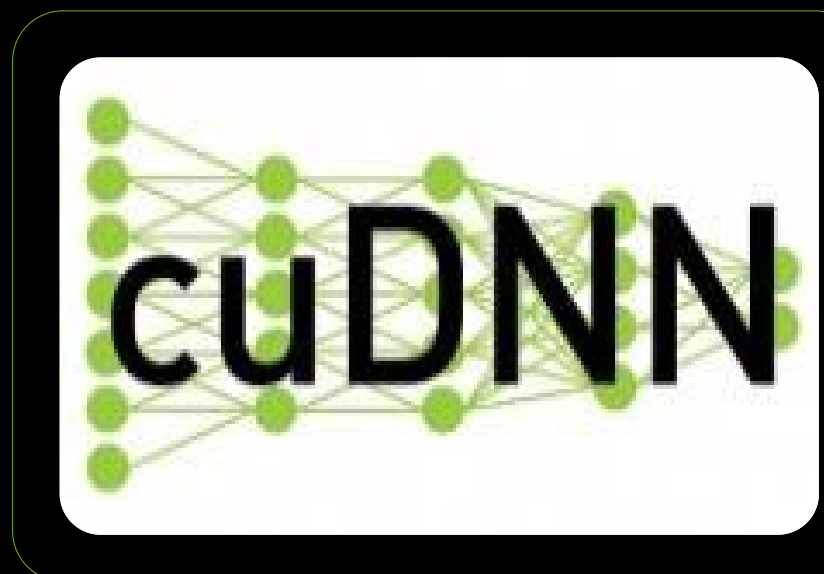


Compile new languages to CUDA

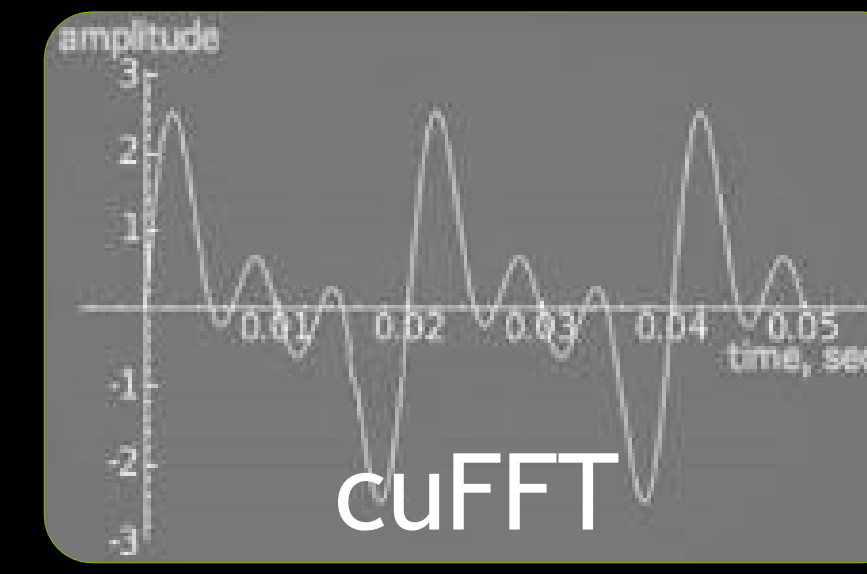
GPU Accelerated Libraries

“Drop-In” Acceleration For Your Applications

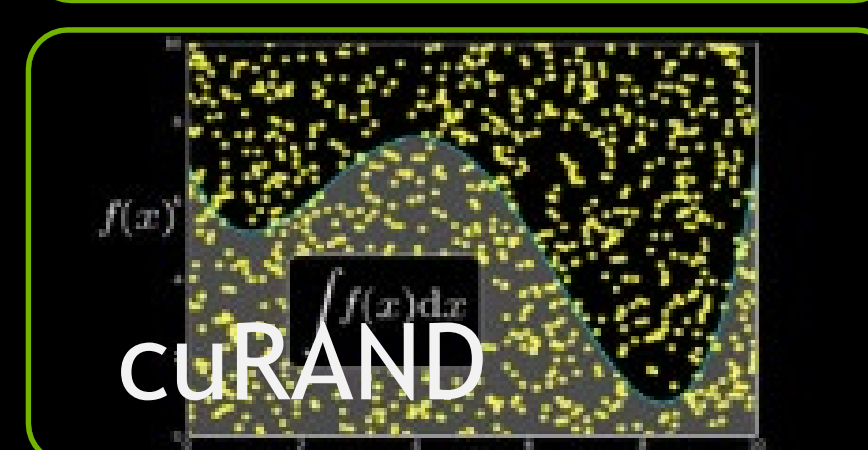
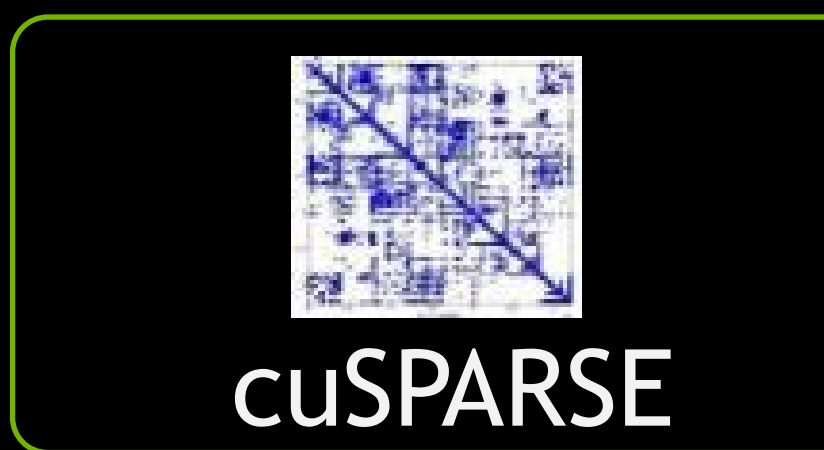
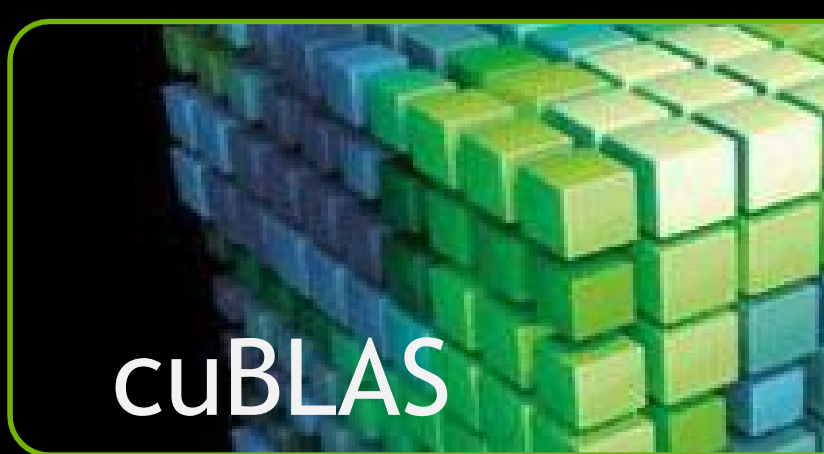
DEEP LEARNING



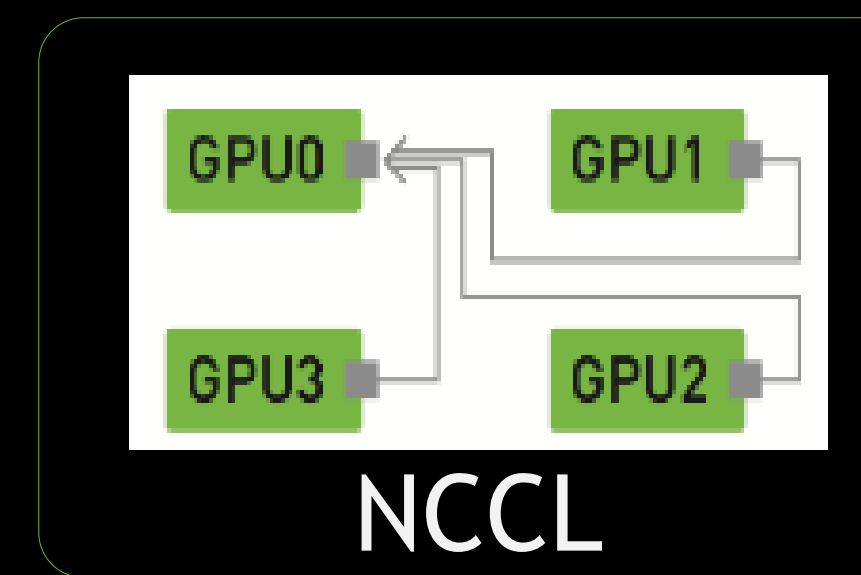
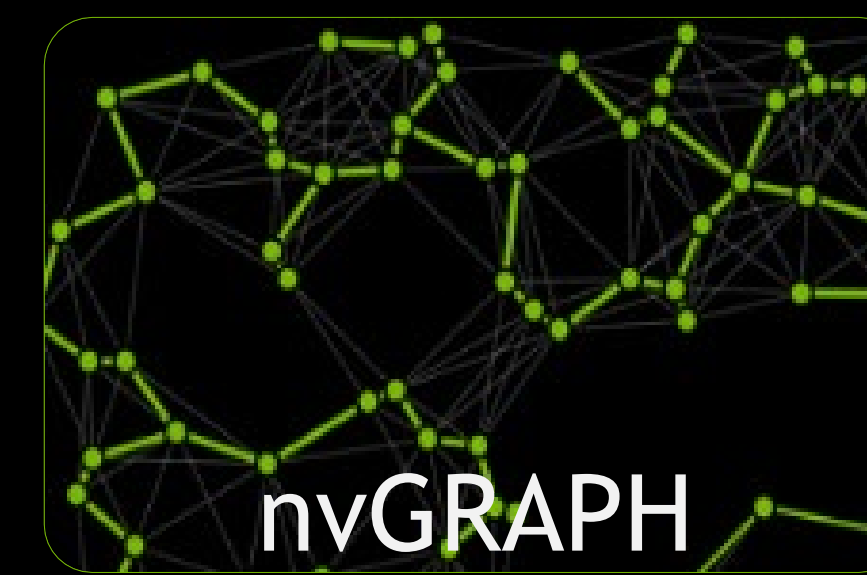
SIGNAL, IMAGE & VIDEO



LINEAR ALGEBRA

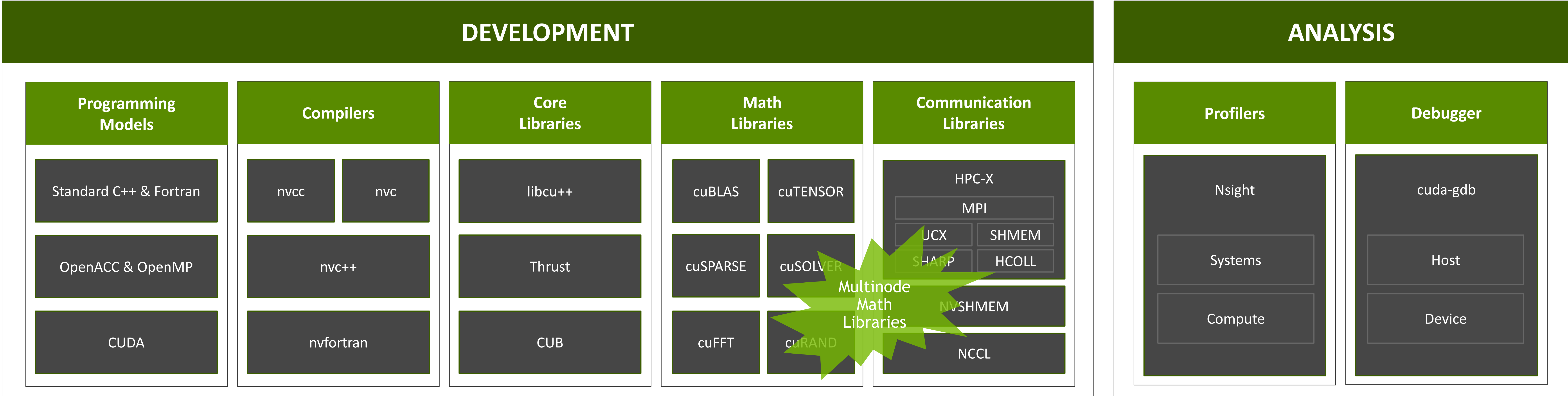


PARALLEL ALGORITHMS



NVIDIA HPC SDK

Available at developer.nvidia.com/hpc-sdk, on NGC, via Spack, and in the Cloud

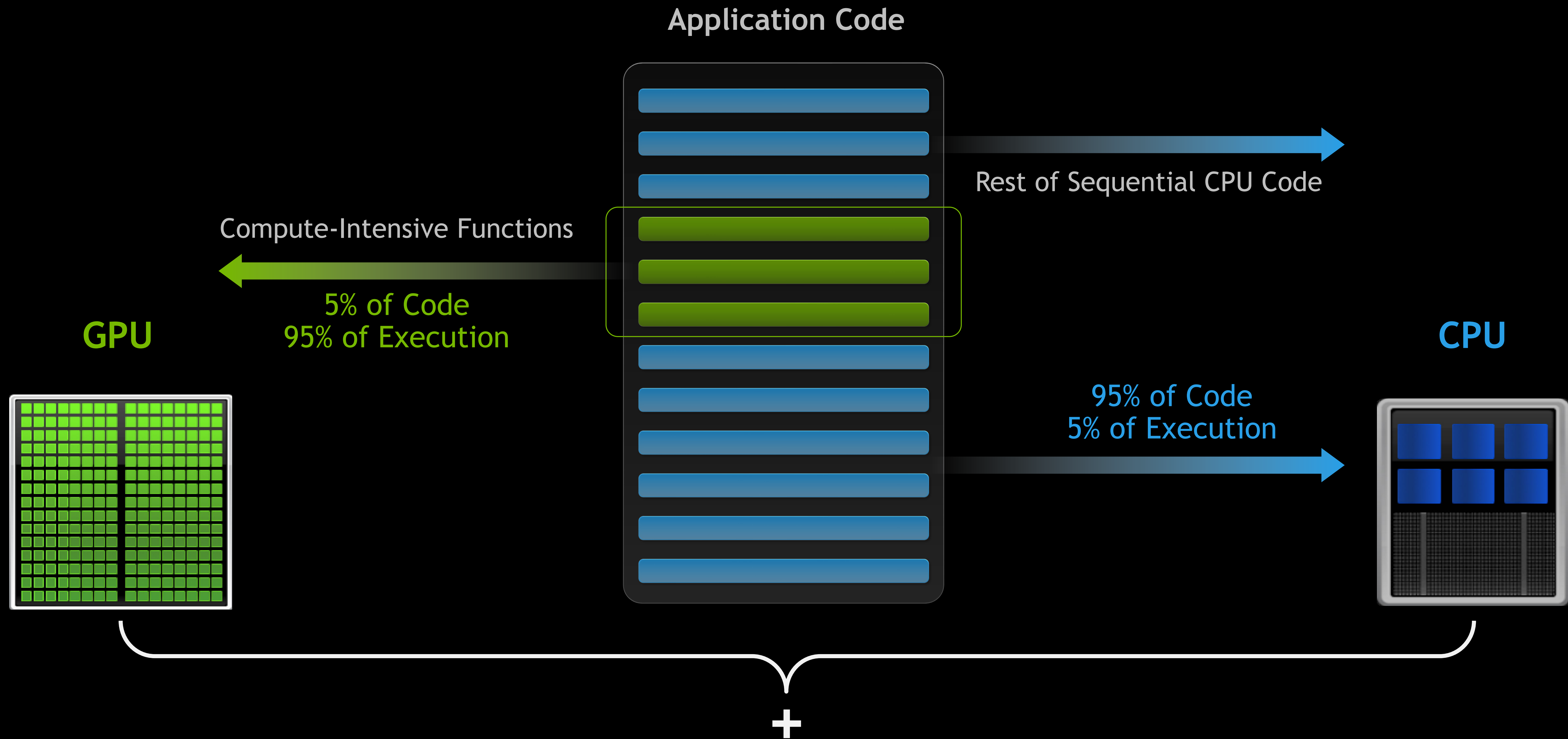


Develop for the NVIDIA Platform: GPU, CPU and Interconnect
Libraries | Accelerated C++ and Fortran | Directives | CUDA
x86_64 | Arm | OpenPOWER
7-8 Releases Per Year | Freely Available



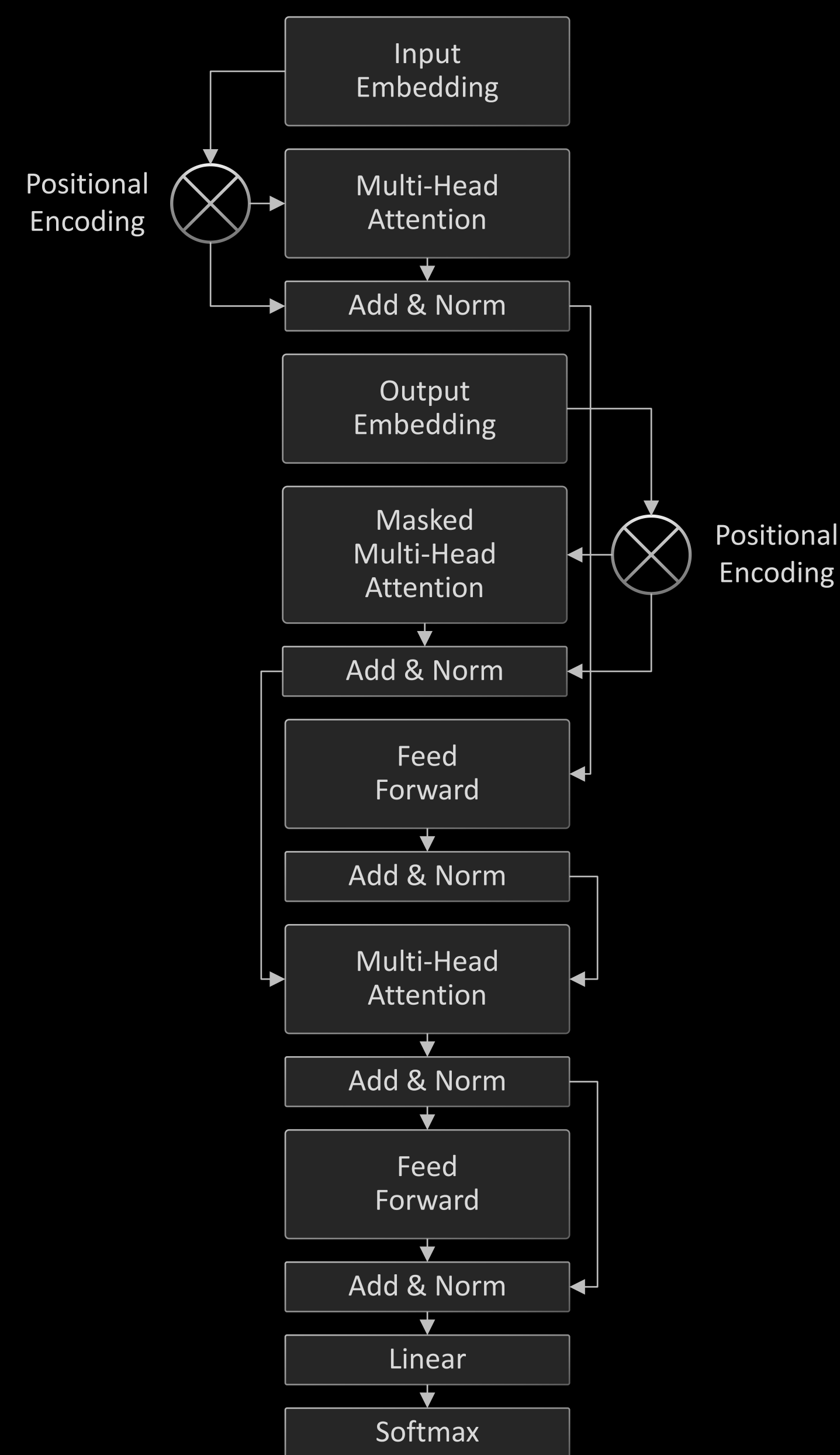
OK, but, What Does It Mean?

How GPU Acceleration Works



Piece-by-Piece, not All-or-Nothing

Incrementally accelerate key components of an application



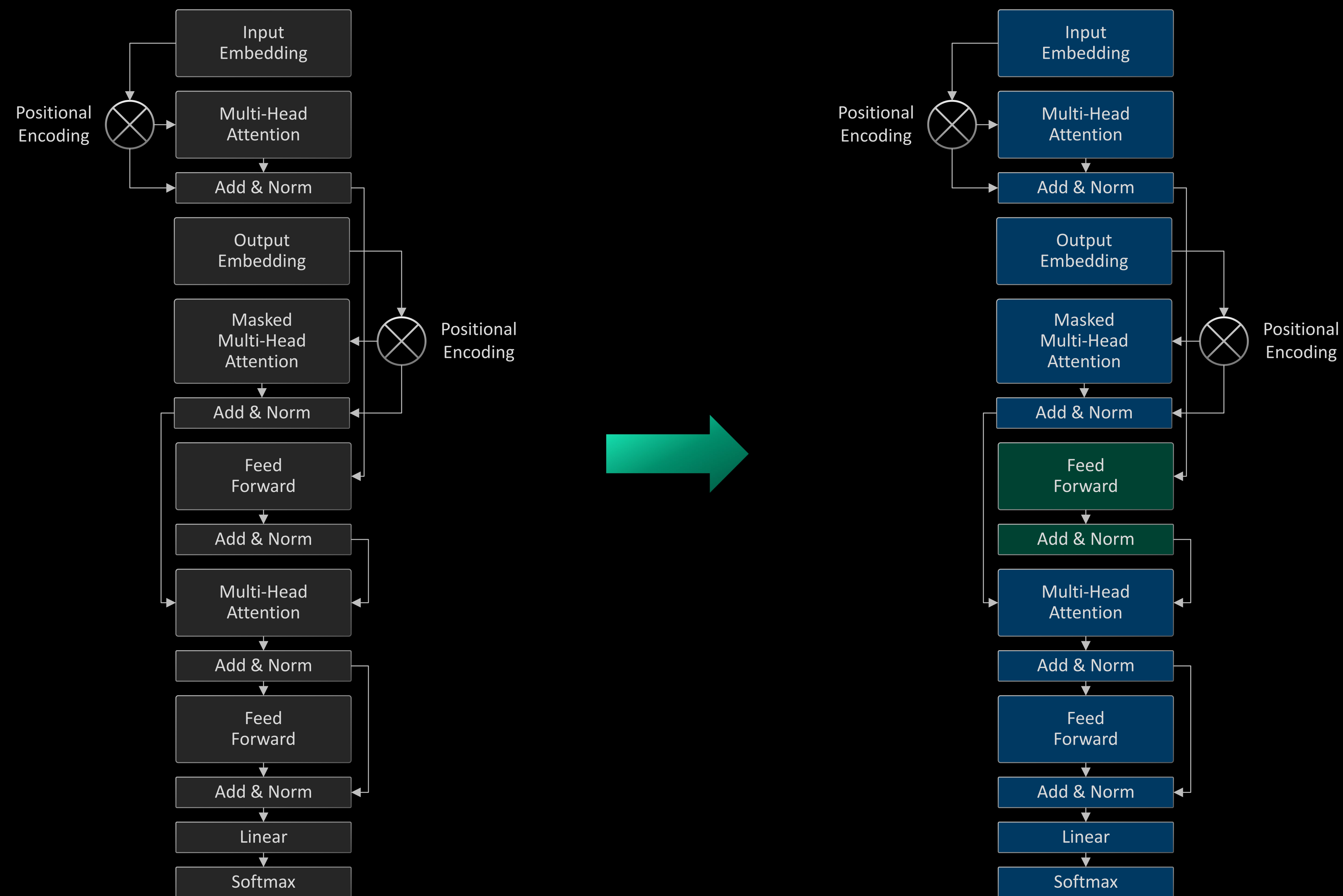
Real applications are complex

No need to port the whole thing in one go

Example: "Transformer" deep neural network

Piece-by-Piece, not All-or-Nothing

Incrementally accelerate key components of an application



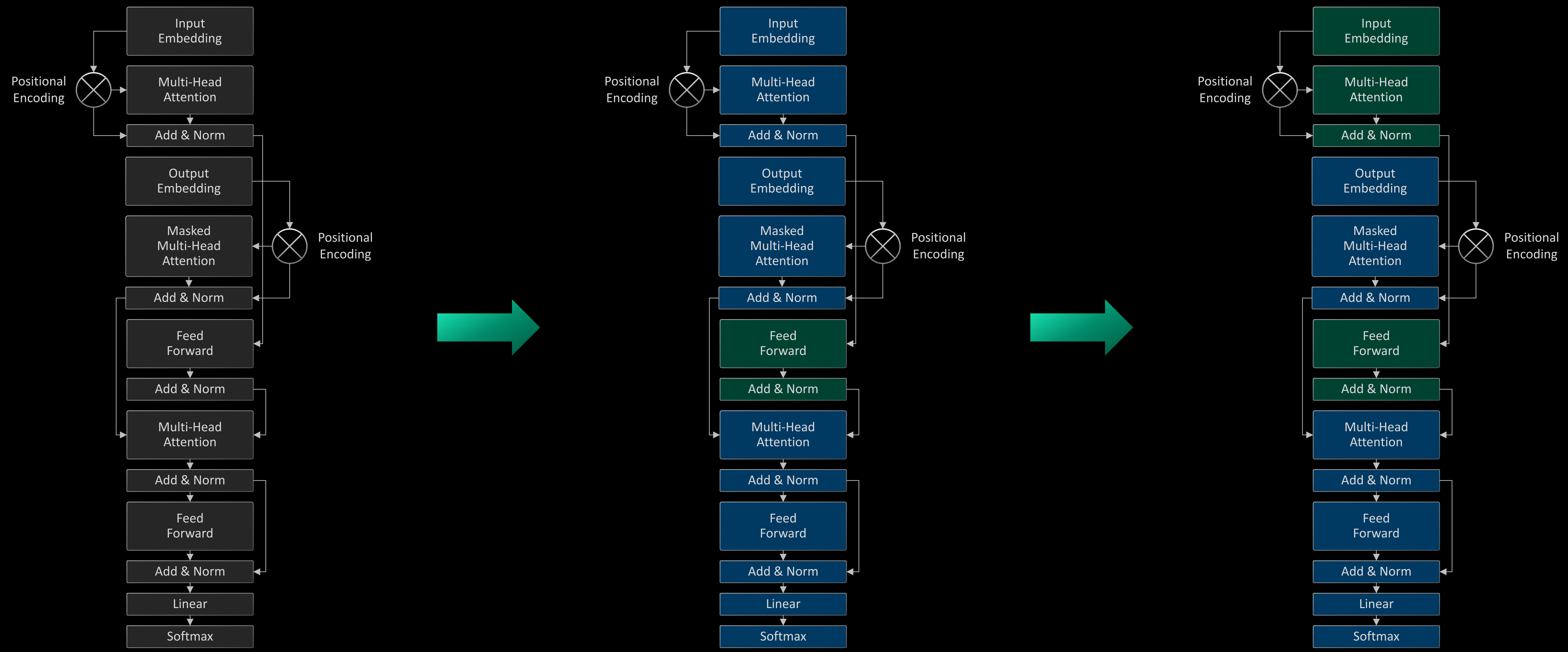
CUDA includes heterogeneous profiling tools to help evaluate which components to port next

Example: "Transformer" deep neural network

■ Runs on GPU ■ Runs on CPU

Piece-by-Piece, not All-or-Nothing

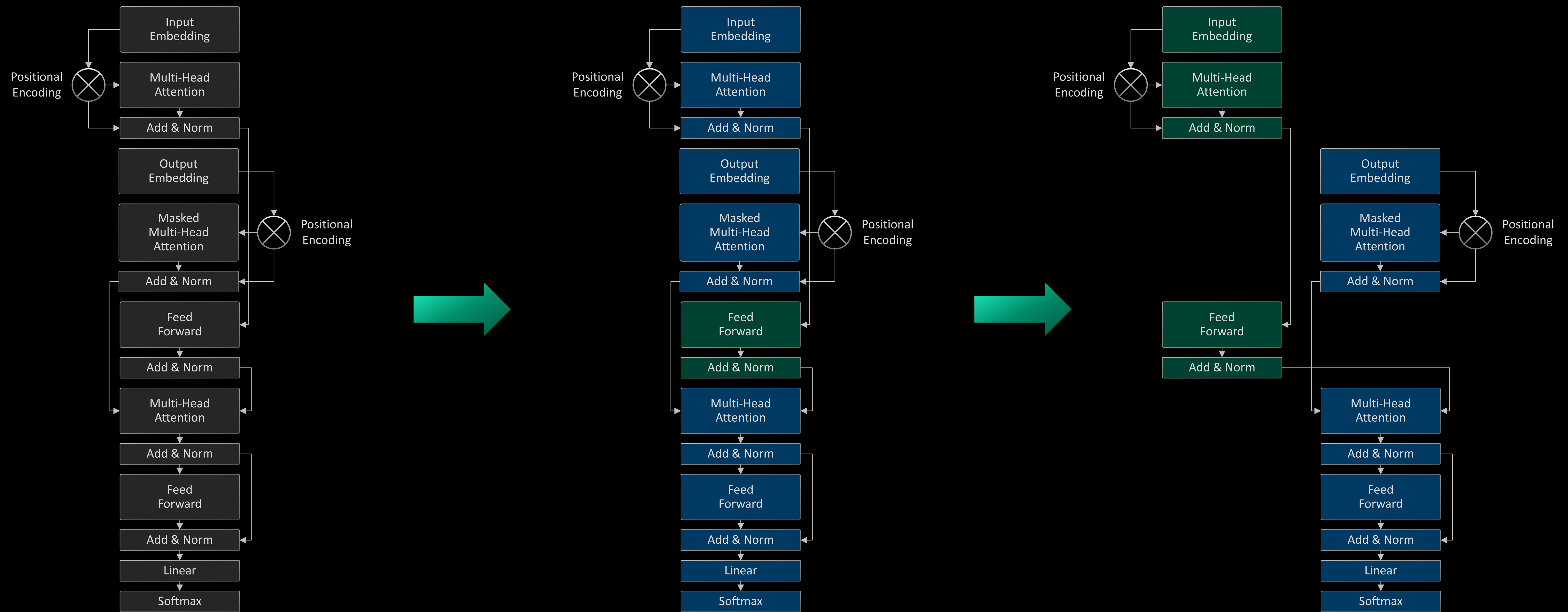
Incrementally accelerate key components of an application



Example: "Transformer" deep neural network

Piece-by-Piece, not All-or-Nothing

Incrementally accelerate key components of an application



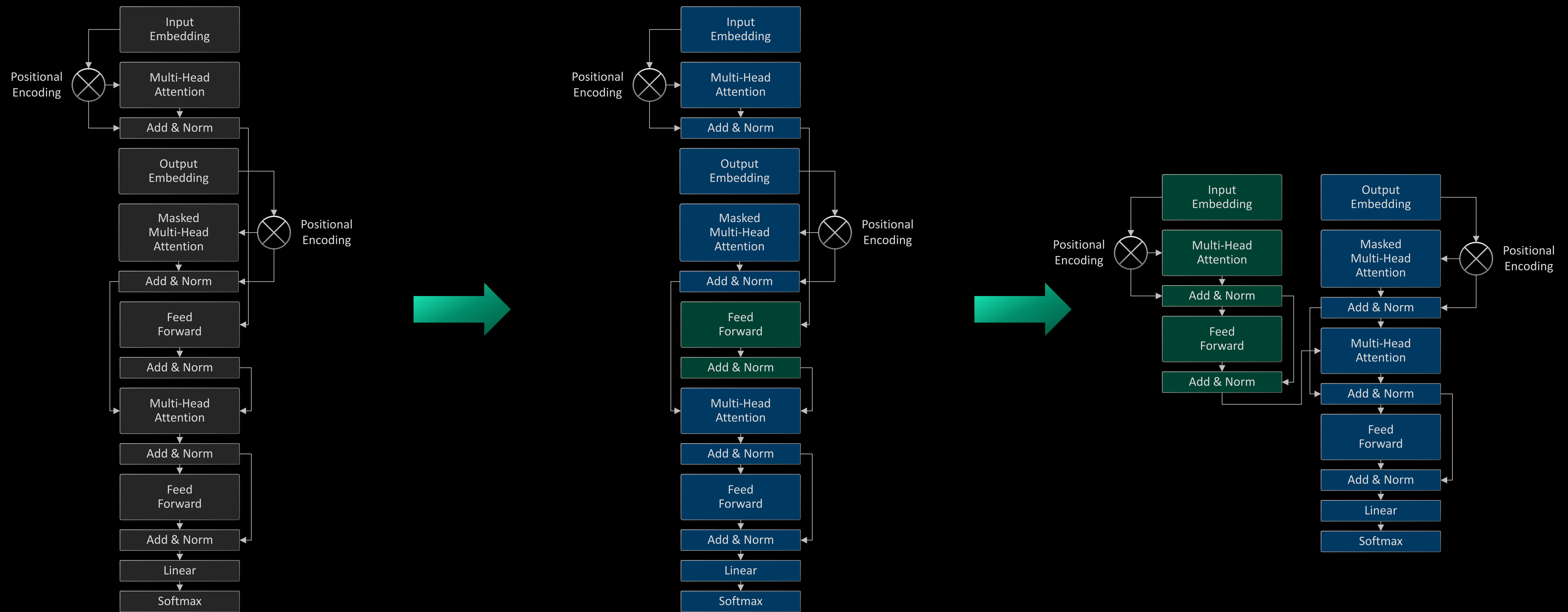
Example: "Transformer" deep neural network

Runs on GPU

Runs on CPU

Piece-by-Piece, not All-or-Nothing

Incrementally accelerate key components of an application



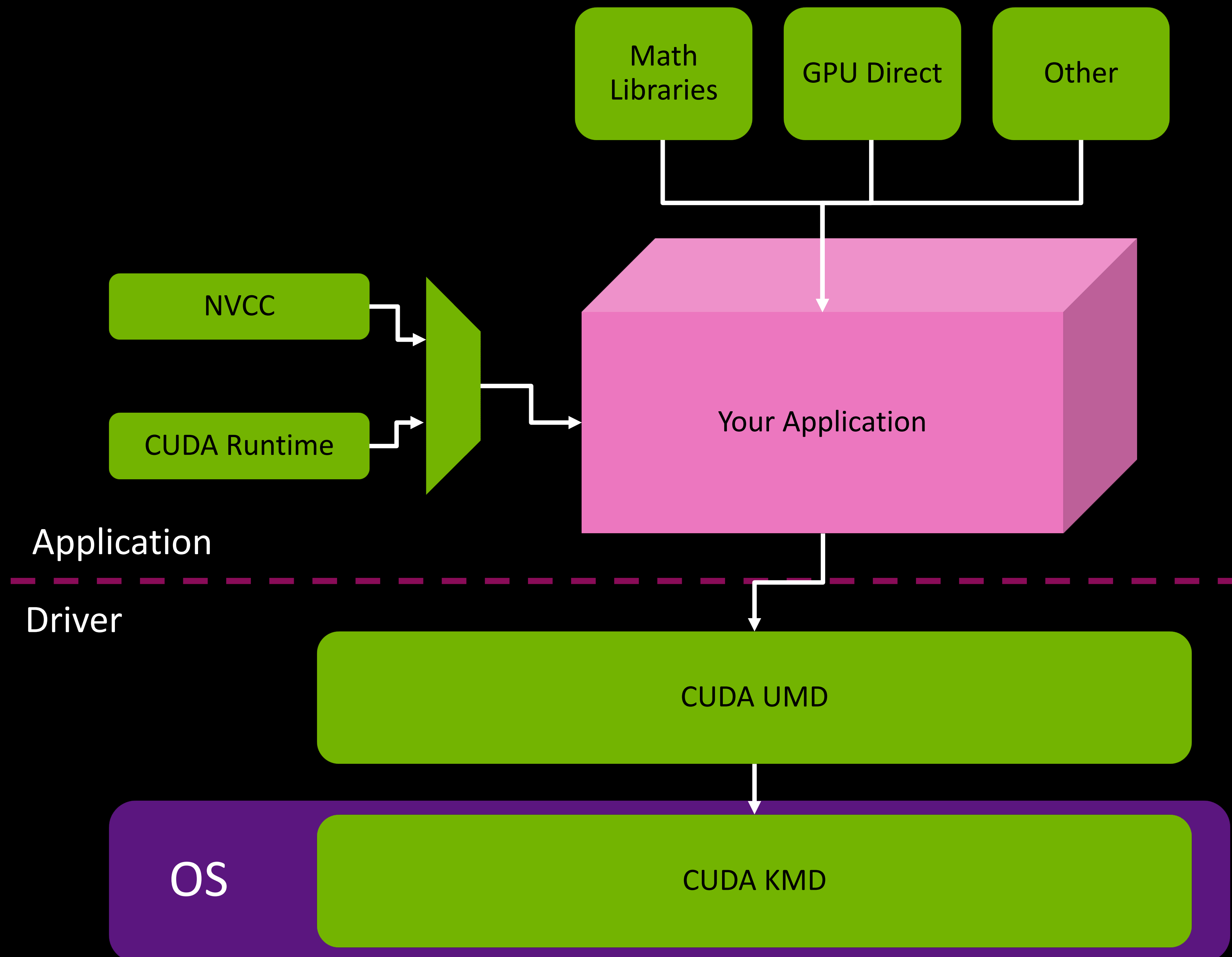
Example: "Transformer" deep neural network

■ Runs on GPU

■ Runs on CPU

What is CUDA

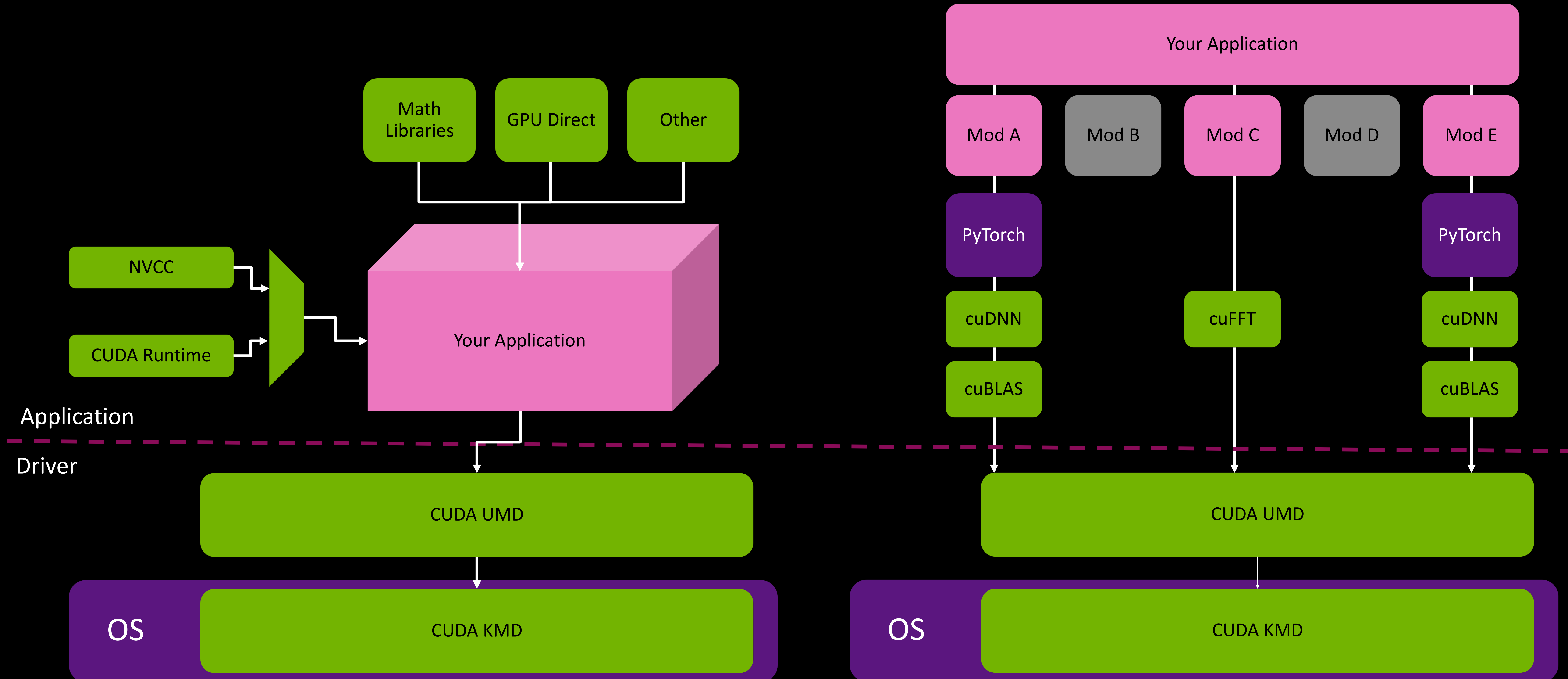
A Simplified View



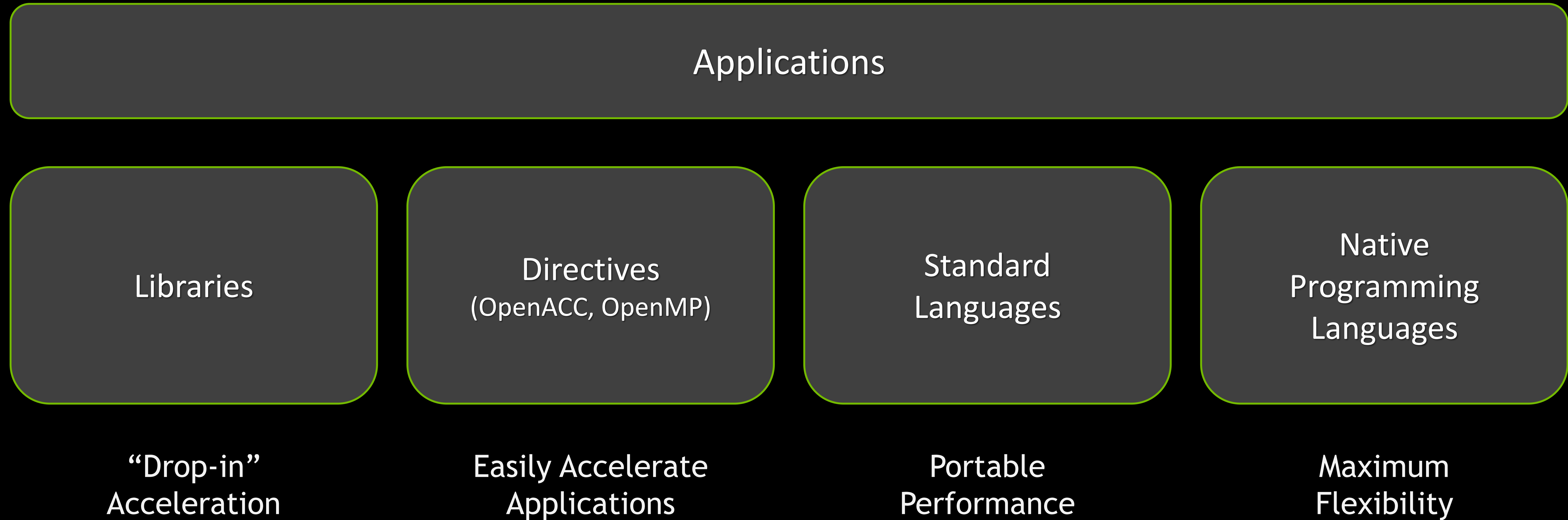
- Domain-Specific Libraries
- CUDA Programming Environment
 - Compiler
 - CUDA programming model (will talk more about this)
- CUDA Runtime Libraries
- Driver:
 - **Kernel Mode Driver** – Lives in the OS, handles low-level hardware interaction
 - **User Mode Driver** – Integrates with your application, maps low-level CUDA API calls to your specific HW

But, Of Course, a Real Application is Complex

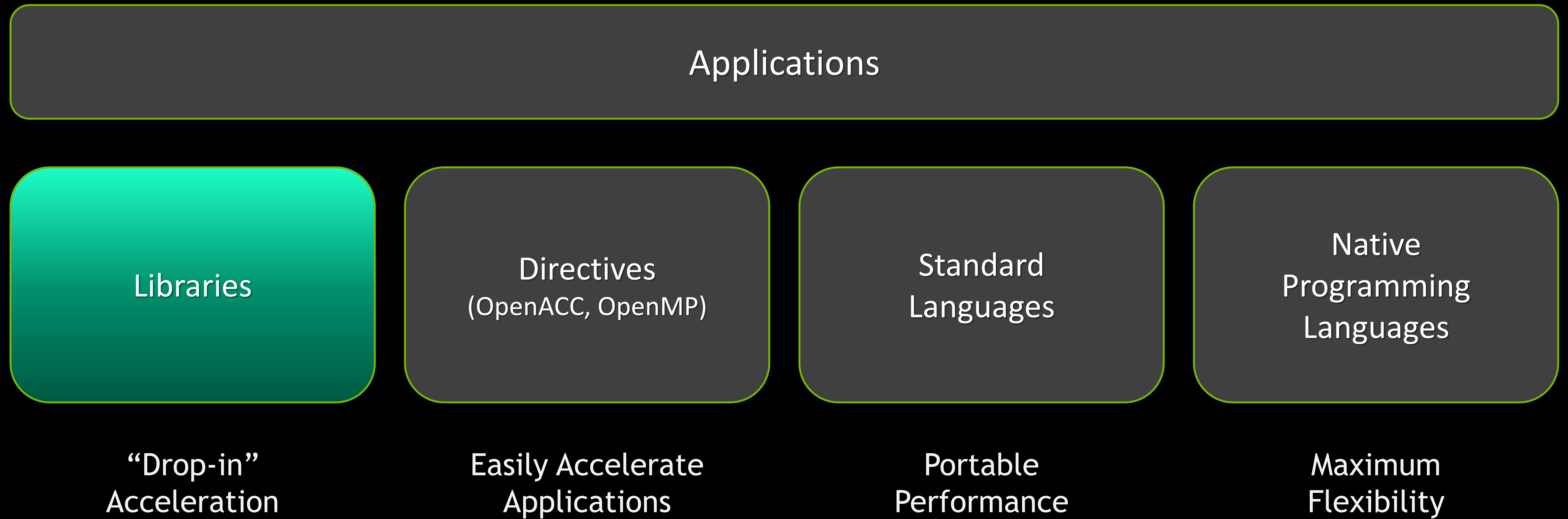
Many Components, Many Dependencies



4 Ways To Accelerate Applications



4 Ways To Accelerate Applications



Libraries: Easy, High-Quality Acceleration

EASE OF USE

Using libraries enables GPU acceleration without in-depth knowledge of GPU programming

“DROP-IN”

Many GPU-accelerated libraries follow standard APIs, thus enabling acceleration with minimal code changes

QUALITY

Libraries offer high-quality implementations of functions encountered in a broad range of applications

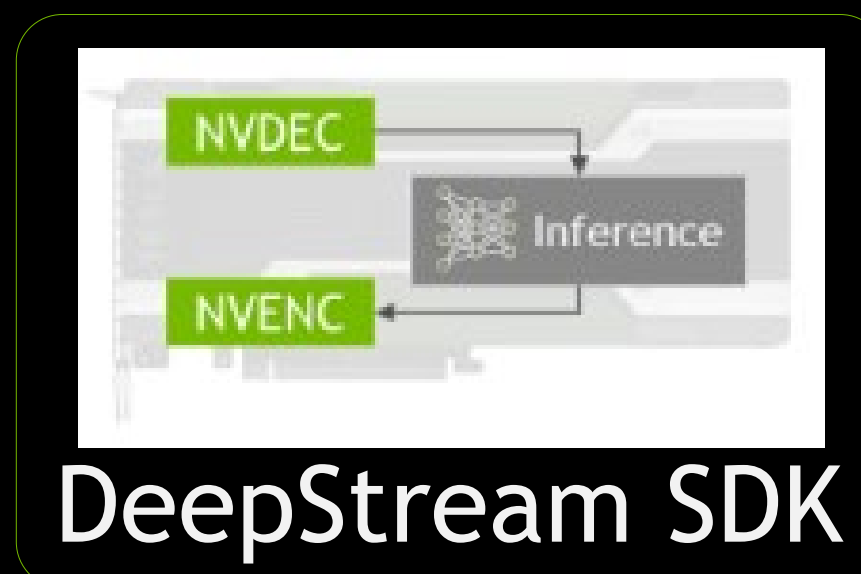
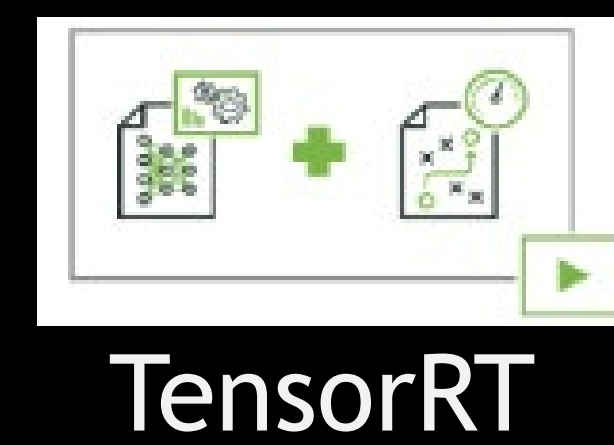
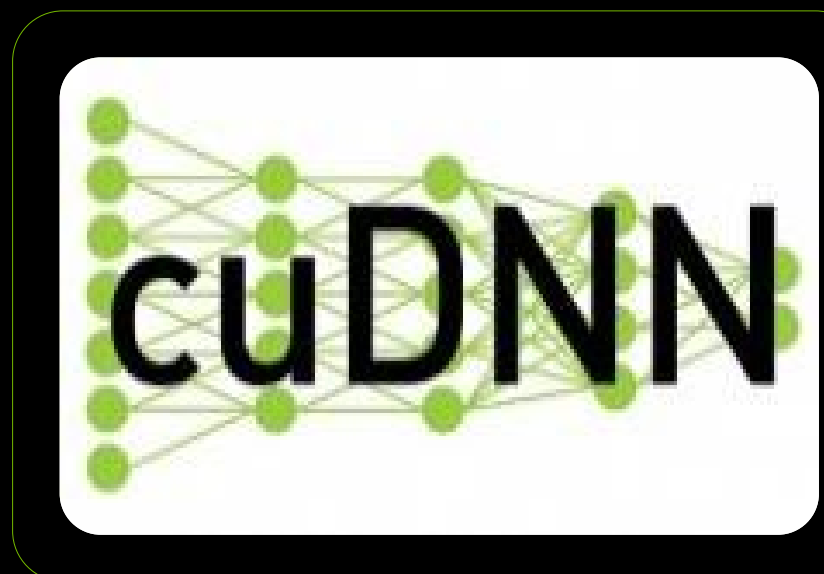
PERFORMANCE

NVIDIA libraries are tuned by experts

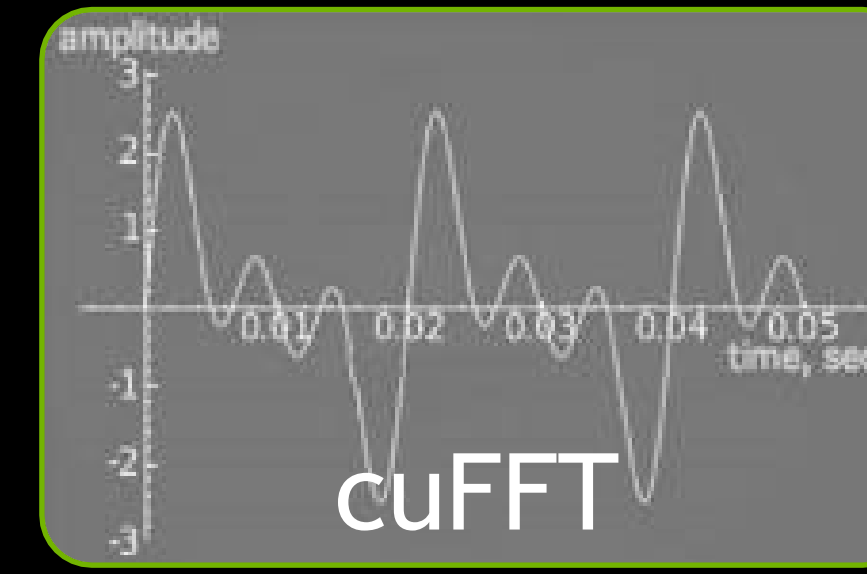
GPU Accelerated Libraries

“Drop-In” Acceleration For Your Applications

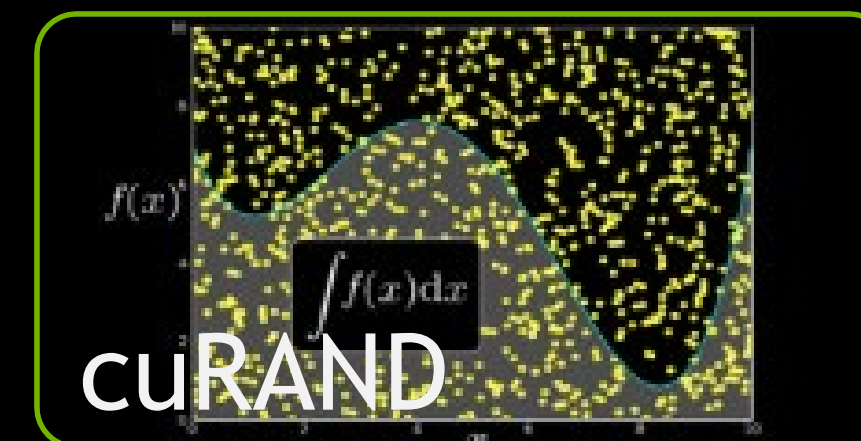
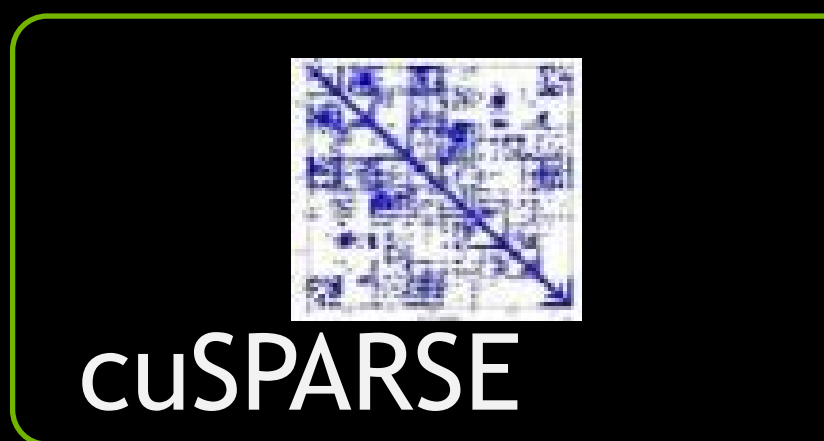
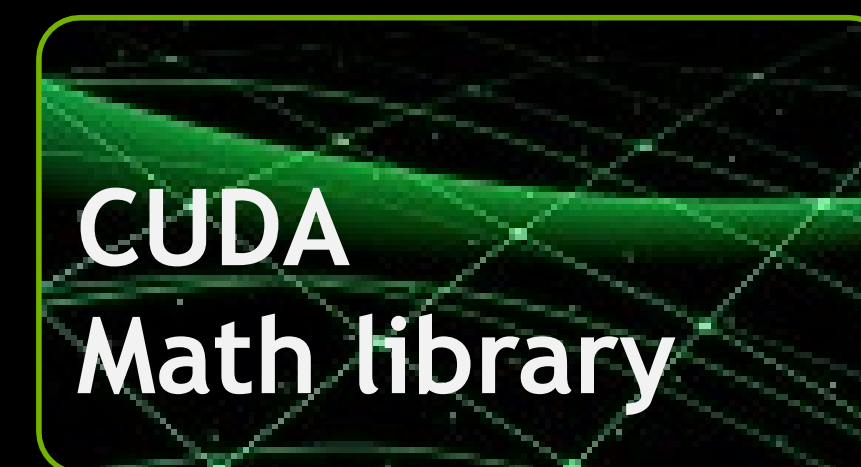
DEEP LEARNING



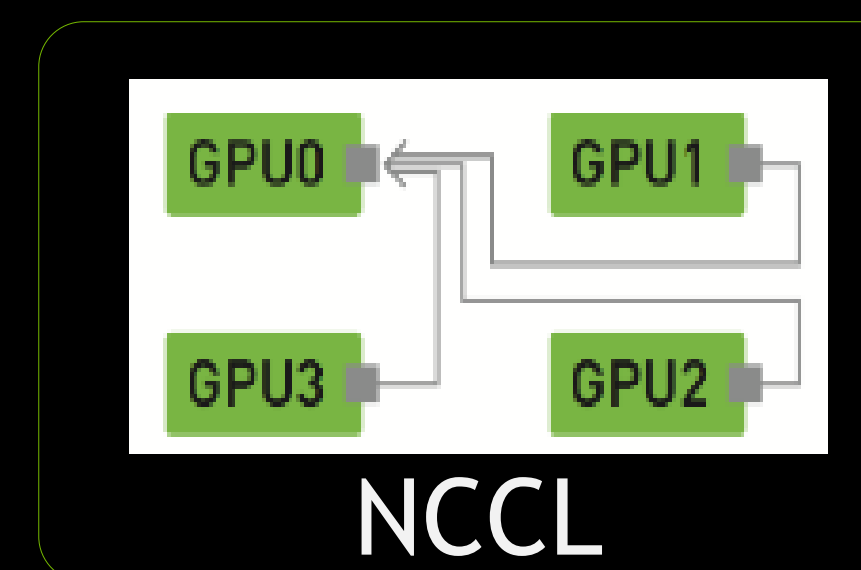
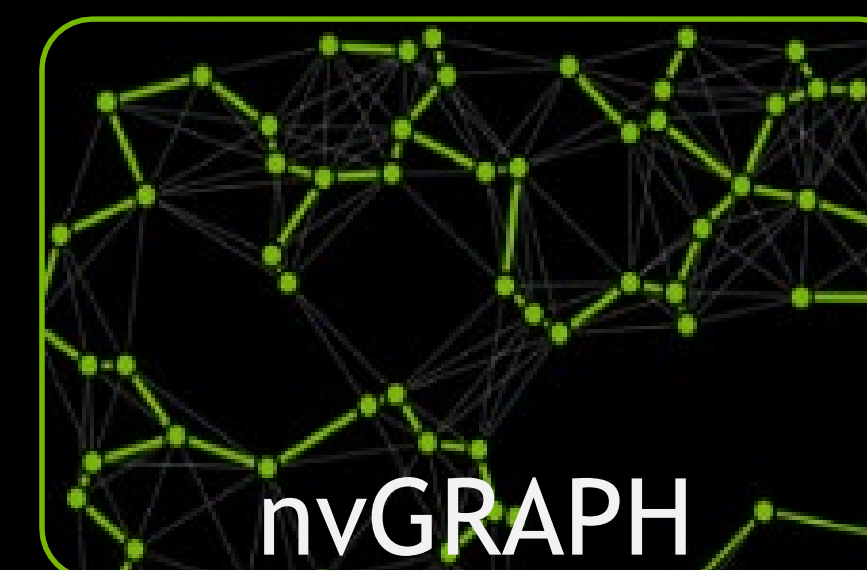
SIGNAL, IMAGE & VIDEO



LINEAR ALGEBRA



PARALLEL ALGORITHMS



3 Steps To A CUDA-Accelerated Application

STEP 1

Substitute library calls with equivalent CUDA library calls

`saxpy (...)` → `cublasSaxpy (...)`

STEP 2

Manage data locality

with CUDA: `cudaMallocManaged(), cudaMemcpy()`
with CUBLAS: `cublasAlloc(), cublasSetVector()`

STEP 3

Rebuild and link the application with the CUDA-accelerated library

`$ gcc myobj.o -l cublas`

Single Precision Alpha X Plus Y (SAXPY)

Part of Basic Linear Algebra Subroutines (BLAS) library

$$z = \alpha x + y$$

x, y, z : vector

α : scalar

Drop-In Acceleration With CUDA Maths Libraries

In two easy steps

```
int N = 1 << 20;           // 1M elements

x = (float *)malloc(N * sizeof(float));
y = (float *)malloc(N * sizeof(float));
initData(x, y);

// Perform SAXPY on 1M elements: y[]=a*x[]+y[]
saxpy(N, 2.0, x, 1, y, 1);

useResult(y);
```

Original Code

```
int N = 1 << 20;           // 1M elements

x = (float *)malloc(N * sizeof(float));
y = (float *)malloc(N * sizeof(float));
initData(x, y);

// Perform SAXPY on 1M elements: y[]=a*x[]+y[]
saxpy(N, 2.0, x, 1, y, 1);

useResult(y);
```

GPU-Accelerated Code

Drop-In Acceleration With CUDA Maths Libraries

Step 1: Update memory allocation to be CUDA-aware

```
int N = 1 << 20;           // 1M elements

x = (float *)malloc(N * sizeof(float));
y = (float *)malloc(N * sizeof(float));
initData(x, y);

// Perform SAXPY on 1M elements: y[]=a*x[]+y[]
saxpy(N, 2.0, x, 1, y, 1);

useResult(y);
```

Original Code

```
int N = 1 << 20;           // 1M elements

cudaMallocManaged(&x, N * sizeof(float));
cudaMallocManaged(&y, N * sizeof(float));
initData(x, y);

// Perform SAXPY on 1M elements: y[]=a*x[]+y[]
saxpy(N, 2.0, x, 1, y, 1);

useResult(y);
```

GPU-Accelerated Code

Here, we use Unified Memory which automatically migrates between host (CPU) and device (GPU) as needed by the program

Drop-In Acceleration With CUDA Maths Libraries

Step 2: Call CUDA library version of API

```
int N = 1 << 20;           // 1M elements

x = (float *)malloc(N * sizeof(float));
y = (float *)malloc(N * sizeof(float));
initData(x, y);

// Perform SAXPY on 1M elements: y[]=a*x[]+y[]
saxpy(N, 2.0, x, 1, y, 1);

useResult(y);
```

Original Code

```
int N = 1 << 20;           // 1M elements

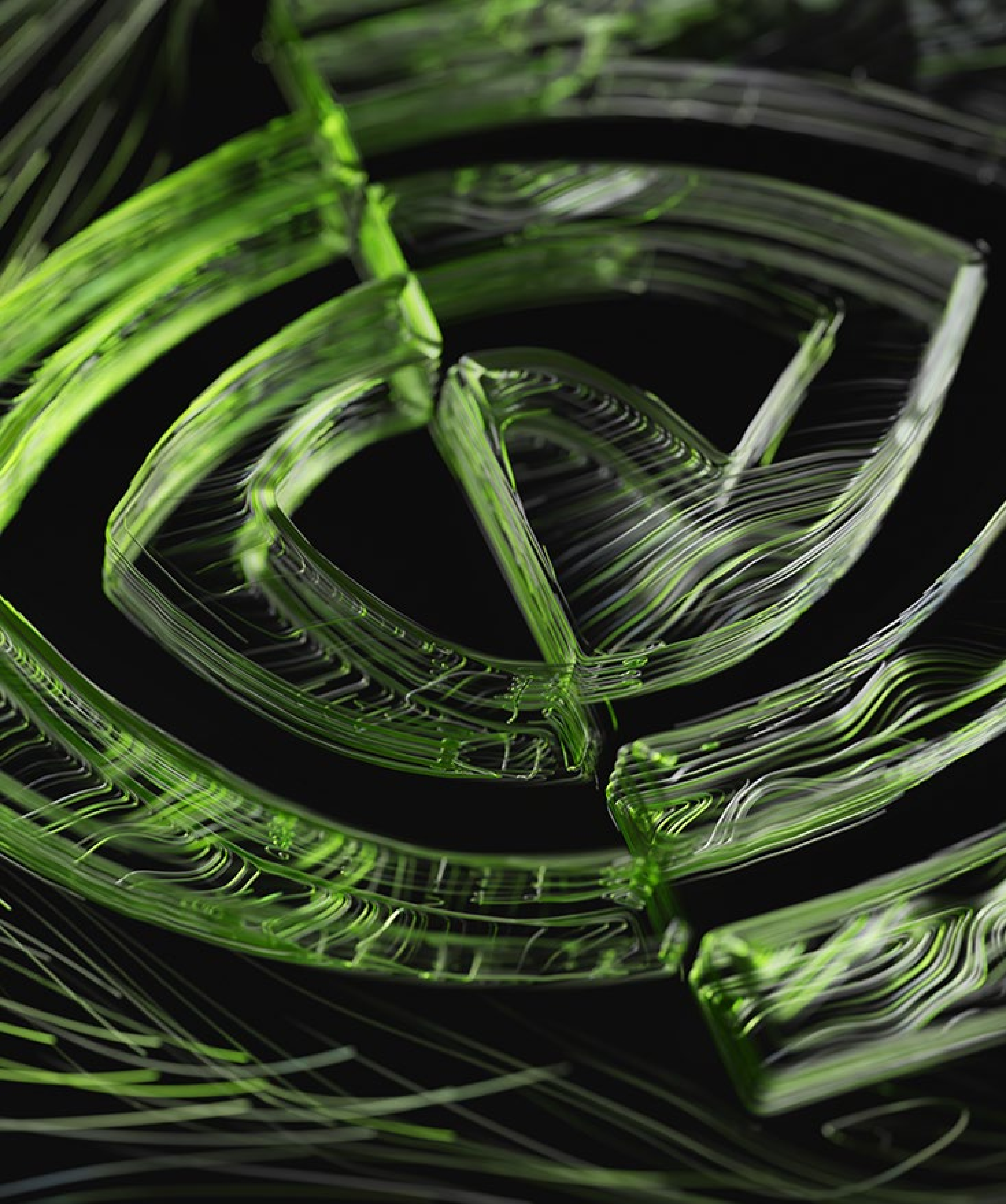
cudaMallocManaged(&x, N * sizeof(float));
cudaMallocManaged(&y, N * sizeof(float));
initData(x, y);

// Perform SAXPY on 1M elements: y[]=a*x[]+y[]
cublasSaxpy(N, 2.0, x, 1, y, 1);

useResult(y);
```

GPU-Accelerated Code

Here, we use Unified Memory which automatically migrates between host (CPU) and device (GPU) as needed by the program



SIX WAYS TO SAXPY

Programming Languages for GPU Computing

Single Precision **Alpha X Plus Y (SAXPY)**

Part of Basic Linear Algebra Subroutines (BLAS) library

$$z = \alpha x + y$$

x, y, z : vector

α : scalar

GPU SAXPY in multiple languages and libraries

A selection of possibilities, not a tutorial

1

OpenACC Compiler Directives

Serial C code

```
void saxpy(int n,
           float a,
           float *x,
           float *y)
{
    for (int i = 0; i < n; ++i)
        y[i] = a*x[i] + y[i];
}

// Perform SAXPY on 1M elements
int N = 1<<20;
saxpy(N, 2.0, x, y);
```

Parallel C code with OpenACC

```
void saxpy(int n,
           float a,
           float *x,
           float *y)
{
    #pragma acc kernels
        for (int i = 0; i < n; ++i)
            y[i] = a*x[i] + y[i];
}

// Perform SAXPY on 1M elements
int N = 1<<20;
saxpy(N, 2.0, x, y);
```

cuBLAS Library

Serial BLAS code

```
int N = 1<<20;

...

// Use your choice of blas library

// Perform SAXPY on 1M elements
blas_saxpy(N, 2.0, x, 1, y, 1);
```

Parallel cuBLAS code

```
int N = 1<<20;

cublasInit();
cublasSetVector(N, sizeof(x[0]), x, 1, d_x, 1);
cublasSetVector(N, sizeof(y[0]), y, 1, d_y, 1);

// Perform SAXPY on 1M elements
cublasSaxpy(N, 2.0, d_x, 1, d_y, 1);

cublasGetVector(N, sizeof(y[0]), d_y, 1, y, 1);

cublasShutdown();
```

You can also call cuBLAS from Fortran,
C++, Python and other languages

Serial C code

```
void saxpy(int n, float a,
           float *x, float *y)
{
    for (int i = 0; i < n; ++i)
        y[i] = a*x[i] + y[i];
}

// Perform SAXPY on 1M elements
int N = 1<<20;
saxpy(N, 2.0, x, y);
```

CUDA C++ code

```
__global__
void saxpy(int n, float a,
           float *x, float *y)
{
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n) y[i] = a*x[i] + y[i];
}

// Perform SAXPY on 1M elements
int N = 1<<20;
saxpy<<< 4096, 256 >>>(N, 2.0, d_x, d_y);
```

CUDA C++ Core Libraries (CCCL)

Serial Standard C++ code

```
int N = 1<<20;
std::vector<float> x(N), y(N);

...

// Perform SAXPY on 1M elements
std::transform(x.begin(), x.end(),
               y.begin(), y.end(),
               2.0f * _1 + _2);
```

CUDA C++ code

```
int N = 1<<20;
thrust::host_vector<float> x(N), y(N);
thrust::device_vector<float> d_x = x;
thrust::device_vector<float> d_y = y;

...

// Perform SAXPY on 1M elements
thrust::transform(d_x.begin(), d_x.end(),
                  d_y.begin(), d_y.begin(),
                  2.0f * _1 + _2)
```

Standard C++ Parallel Algorithms (**stdpar**)

Serial Standard C++ code

```
int N = 1<<20;
std::vector<float> x(N), y(N);

...

// Perform SAXPY on 1M elements
std::transform(x.begin(), x.end(),
               y.begin(), y.end(),
               2.0f * _1 + _2);
```

CUDA C++ code

```
int N = 1<<20;
std::vector<float> x(N), y(N), out;
out.reserve(N);

...

// Perform SAXPY on 1M elements
std::transform(
    std::execution::par_unseq,
    x.begin(), x.end(), y.begin(), y.end(),
    std::back_inserter(out),
    [](int a, int b) {
        return 2.0f * a + b;
    });
```

Standard Python

```
import numpy as np

def saxpy(a, x, y):
    return [a * xi + yi
            for xi, yi in zip(x, y)]

x = np.arange(2**20, dtype=np.float32)
y = np.arange(2**20, dtype=np.float32)

cpu_result = saxpy(2.0, x, y)
```

Numba Parallel Python

```
import numpy as np
from numba import vectorize

@vectorize(['float32(float32, float32,
                    float32)'], target='cuda')
def saxpy(a, x, y):
    return a * x + y

N = 1048576

# Initialize arrays
A = np.ones(N, dtype=np.float32)
B = np.ones(A.shape, dtype=A.dtype)
C = np.empty_like(A, dtype=A.dtype)

# Add arrays on GPU
C = saxpy(2.0, X, Y)
```

Anatomy of a CUDA binary

Hello world example

```
__global__ void kernel() {  
    printf("Hello, CUDA\n");  
}  
  
void main() {  
    kernel<<< 1, 1 >>>();  
    cudaDeviceSynchronize();  
}
```

```
nvcc example.cu -o example
```

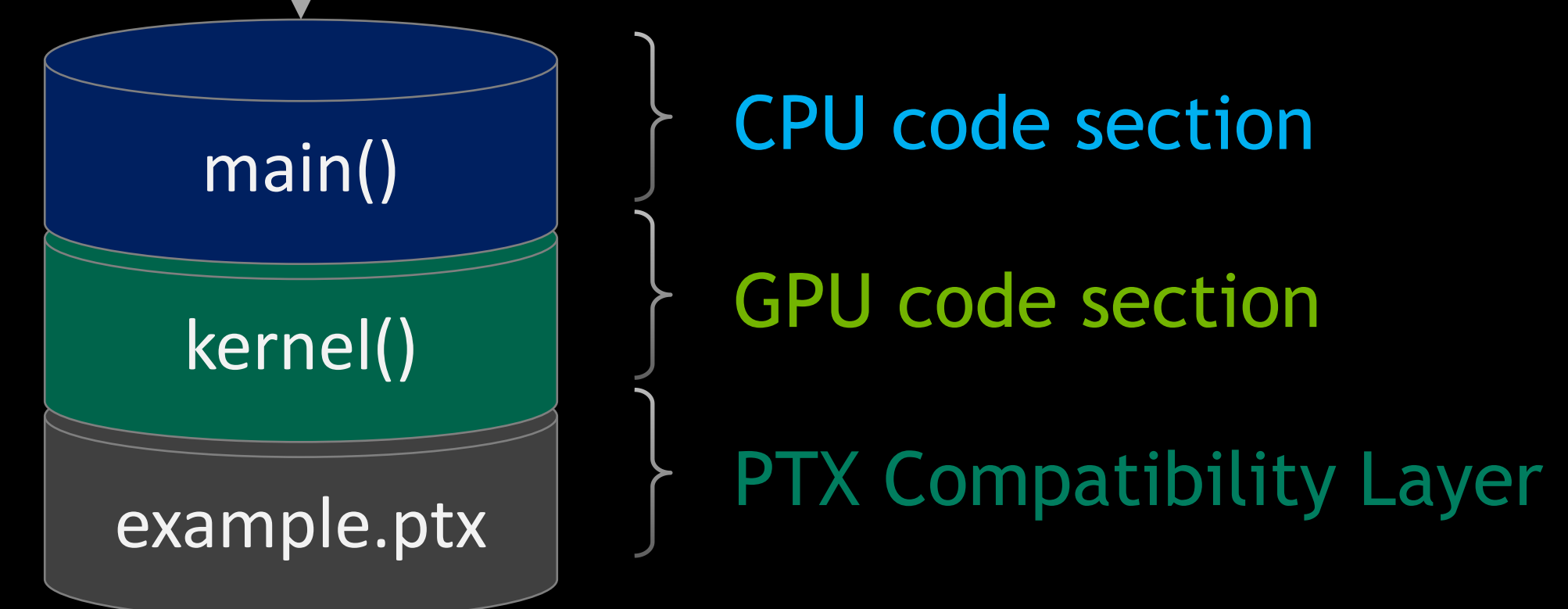
"example"
executable

Anatomy of a CUDA binary

Hello world example

```
__global__ void kernel() {  
    printf("Hello, CUDA\n");  
}  
  
void main() {  
    kernel<<< 1, 1 >>>();  
    cudaDeviceSynchronize();  
}
```

```
nvcc example.cu -o example
```



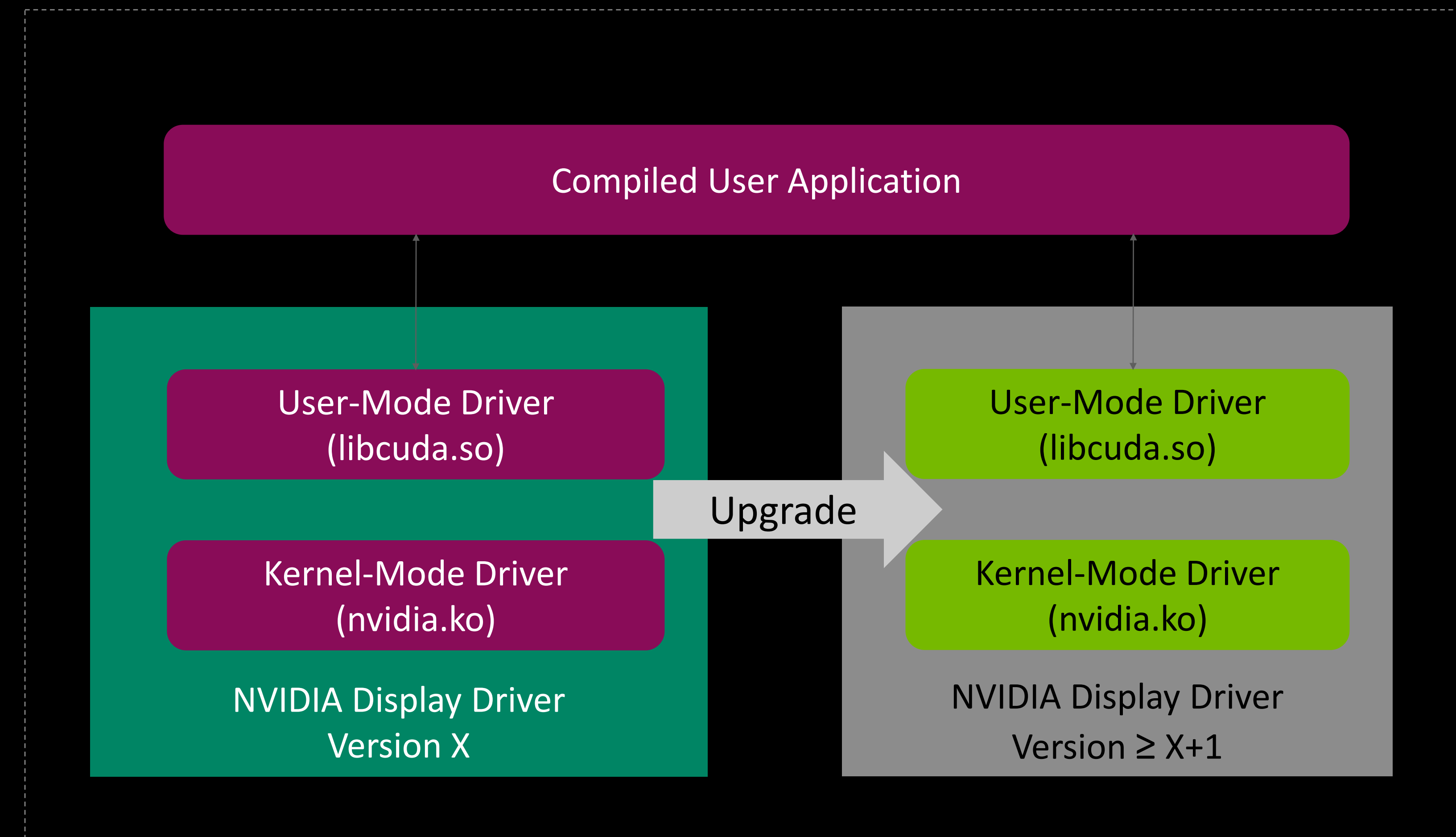
How Do We Keep Things Working Together?



CUDA Compatibility

“Backward Compatibility” Software Considerations

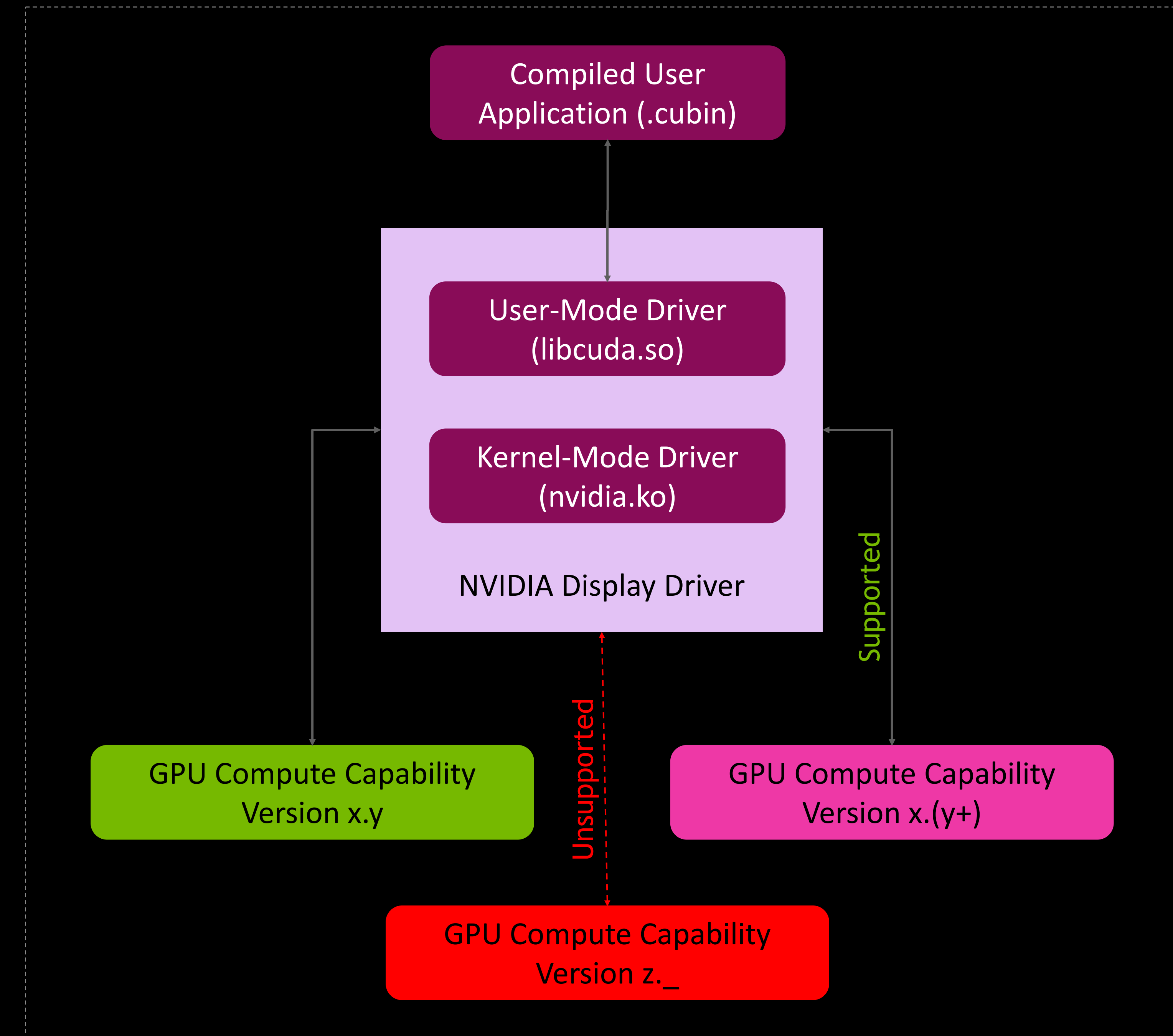
- ▶ The simplest use case:
 - ▶ “Your compiled application will work forever on NVIDIA GPUs, regardless of installed driver”
- ▶ All newer GPU drivers will be binary-compatible with older binaries
 - ▶ Requires statically linking libraries like the CUDA runtime
- ▶ Recompiling from *source* may require API changes
 - ▶ Only binary compatibility is guaranteed



CUDA Compatibility

“SM/Compute Compatibility” Hardware Considerations (Binaries)

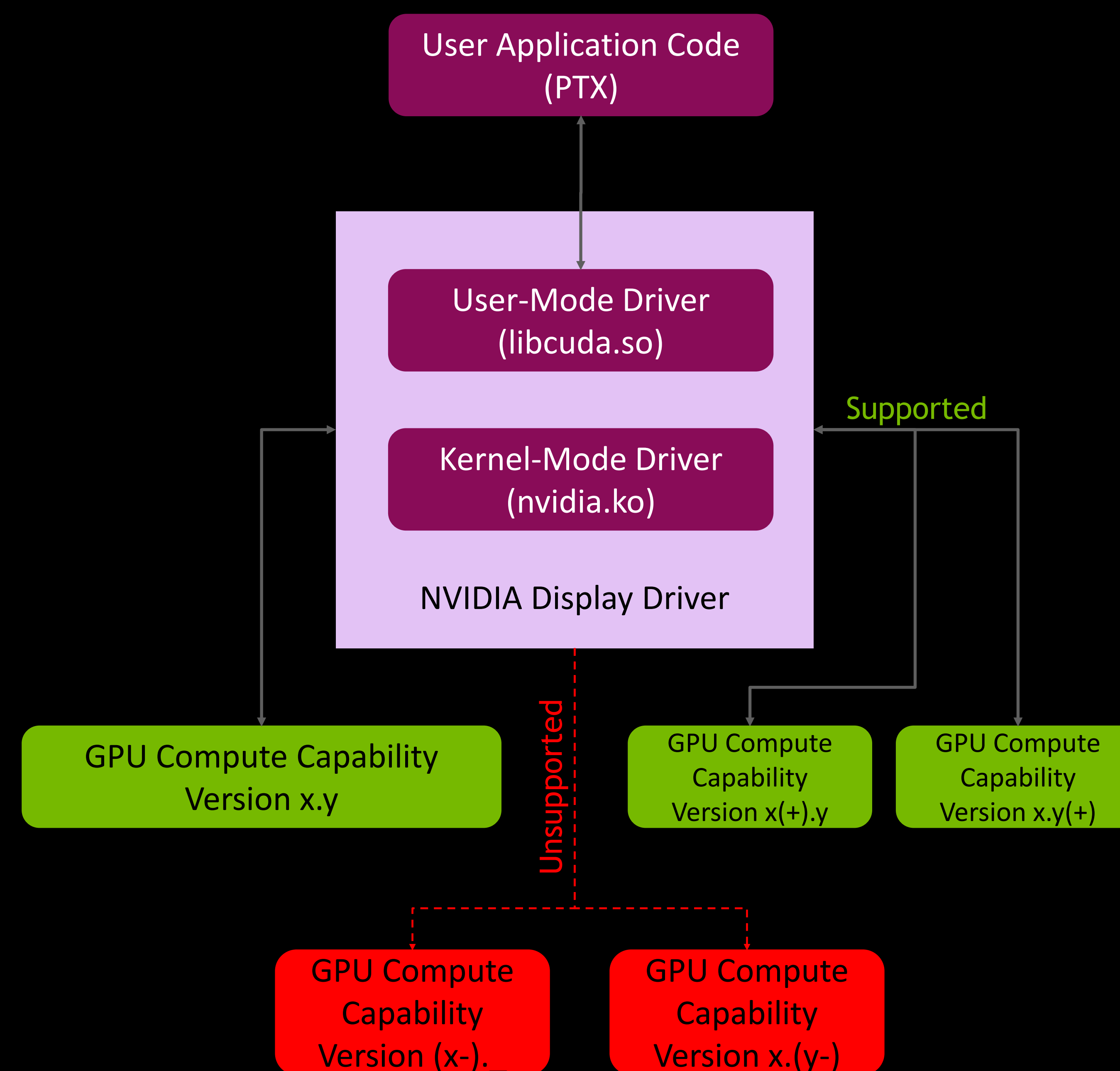
- ▶ Binaries are built for a specific GPU family, PTX is used to target additional families
 - ▶ Each architecture supports a given ISA, or **compute capability**
 - ▶ PTX enables compatibility between architectures
- ▶ Compiled applications target a specific CC, with *some* compatibility within a family (newer but not older)
- ▶ Supported:
 - ▶ **CC 8.0** cubin runs on **CC 8.6** (A100 → A40)
- ▶ Unsupported:
 - ▶ **CC 8.6** cubin cannot run on **CC 8.0** (A40 → A100)
 - ▶ **CC 8.0** cubin cannot run on **CC 7.0** (A100 → V100)
 - ▶ **CC 7.0** cubin cannot run on **CC 8.0** (V100 → A100)



CUDA compatibility

“PTX Compatibility” Hardware Considerations

- ▶ **PTX Code** is compatible with future versions, both Major and Minor
- ▶ Supported PTX Migration:
 - ▶ CC 8.0 PTX runs on CC 8.6 (A100 PTX → A40)
 - ▶ CC 7.0 PTX runs on CC 8._ (V100 PTX → A100)
- ▶ Unsupported PTX Migration:
 - ▶ CC 8.6 PTX cannot run on CC 8.0 (A40 PTX → A100)
 - ▶ CC 8.0 PTX cannot run on CC 7.0 (A100 PTX → V100)



CUDA COMPATIBILITY

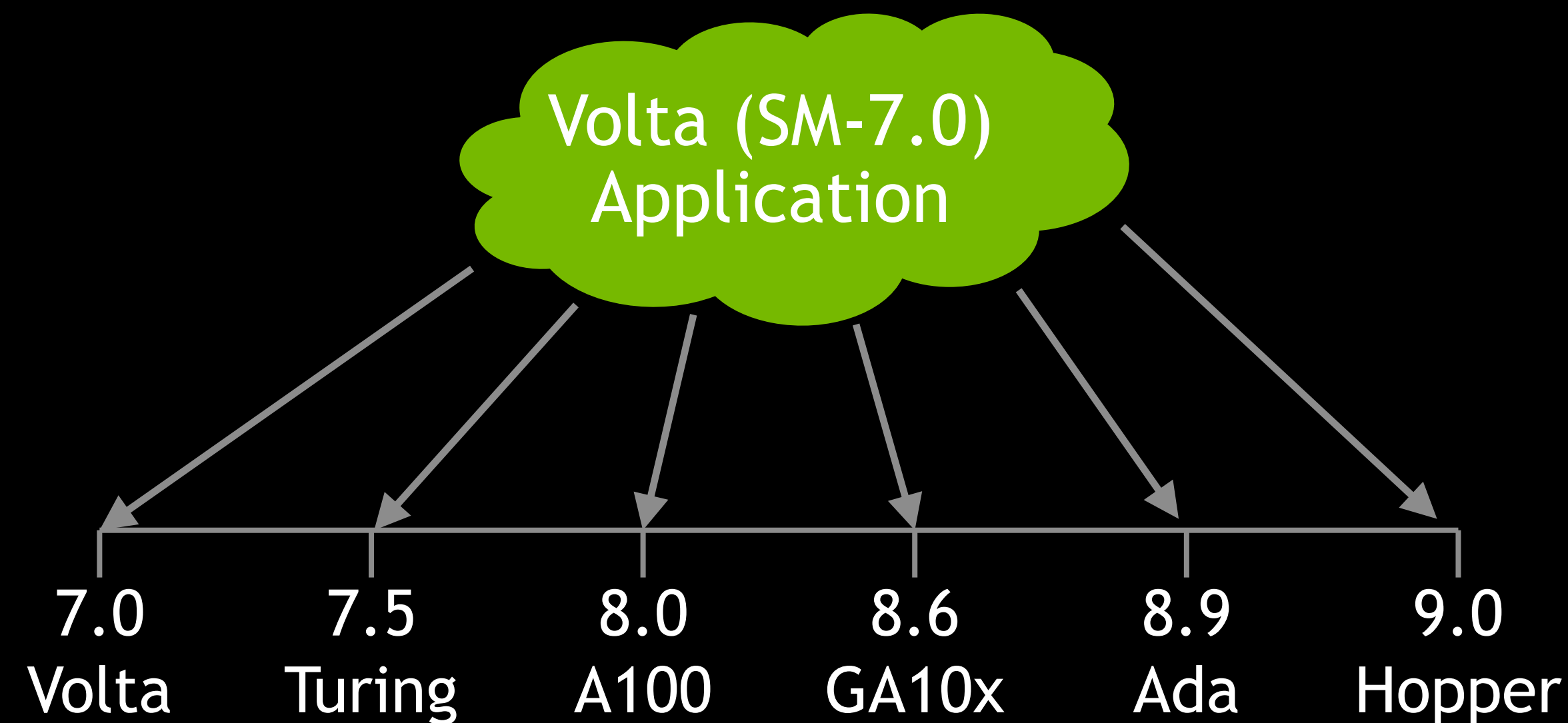
“PTX Compatibility” Hardware Considerations

“CUDA Everywhere”

Code for one GPU runs on all GPUs with newer SM version

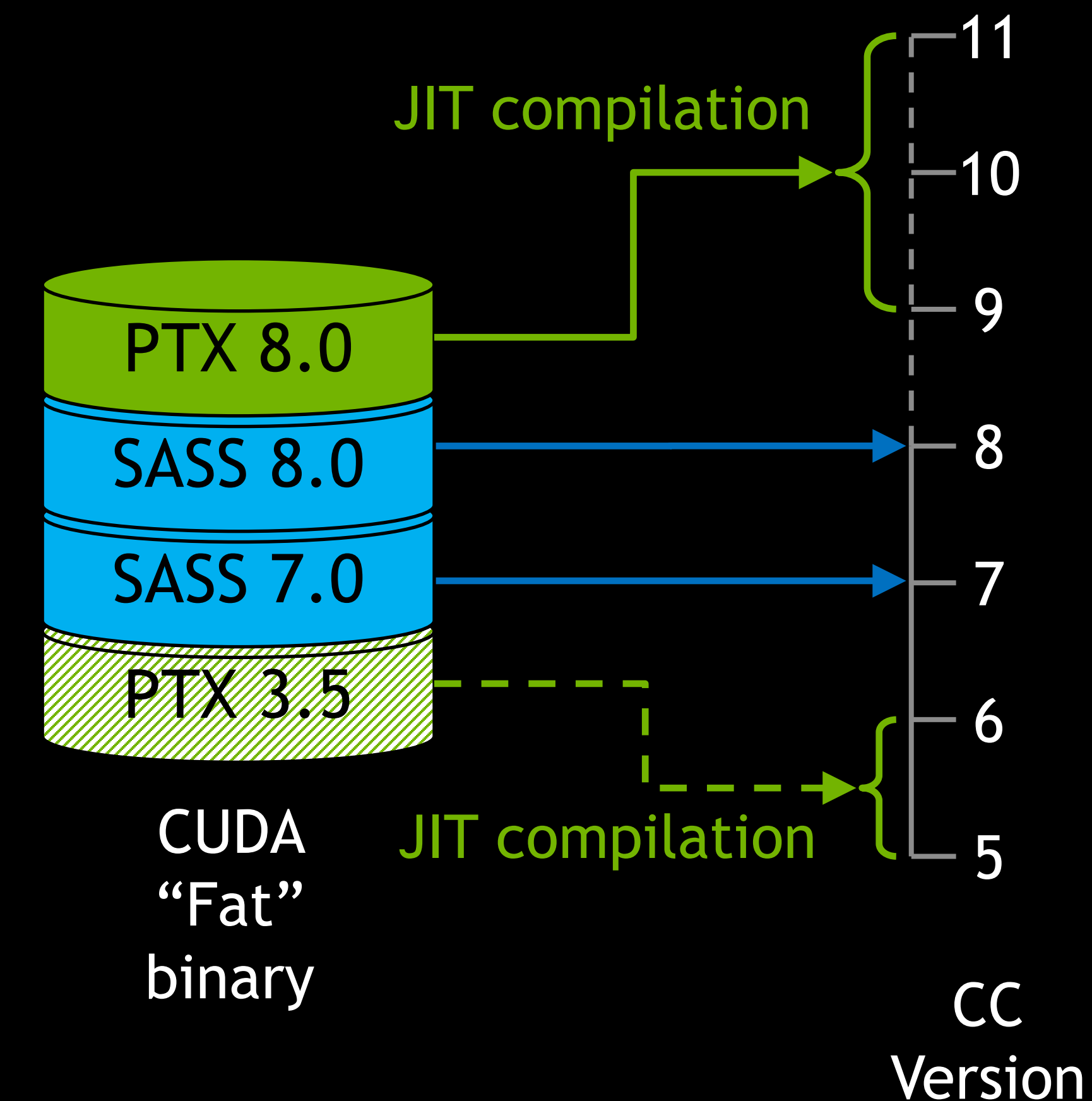
CUDA applications are portable between all chip classes (100, 10x, 20x, 21x, 10b, etc.)

- All current features supported on all future architectures
- Performance & capacities vary (e.g. SM count)
- A few features much slower but still functional (e.g. FP64)



Volta applications “just work” on Turing/Ampere/Hopper
Datacenter libraries “just work” for Quadro, GeForce, etc.

Portability depends on PTX Just-In-Time Compilation



Forward compatibility guarantee: PTX 8.0 runs on CC 9,10,11, ...

Exact match of SASS runs natively (many may exist)

PTX 8.0 won’t run on an older CC. Applications occasionally include older PTX to avoid shipping lots of SASS.

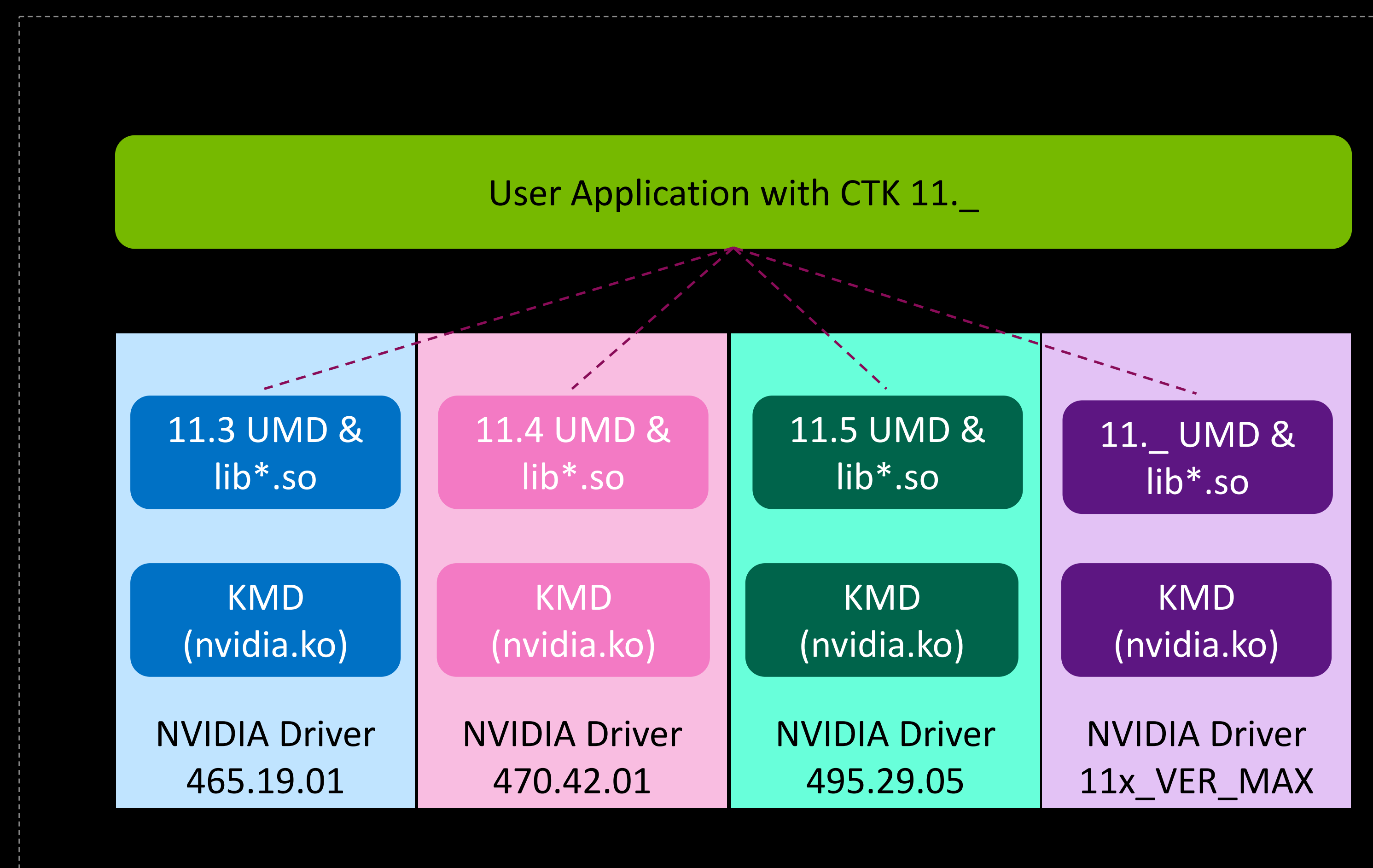
SASS is pre-compiled binary code native to a specific GPU architecture - multiple versions may be packaged together

PTX is assembly code JIT compiled by CUDA when an application is run on a new GPU for which there is no SASS

CUDA compatibility

“Minor-Version Compatibility” (Previously “Enhanced Compatibility”)

- ▶ Applications created **within a major-release** of CUDA may run on a system with the minimum driver version
 - ▶ E.g., 11.x CTK requires **450.80.02**
- ▶ Works with:
 - ▶ Newer driver than CTK
 - ▶ Newer CTK than driver
- ▶ New CTK features that require a new driver will return errors
 - ▶ Programmers must write code to check if features exist and if libraries are supported (e.g., cublas must match cudnn)*
- ▶ PTX JIT unsupported (matching driver required)

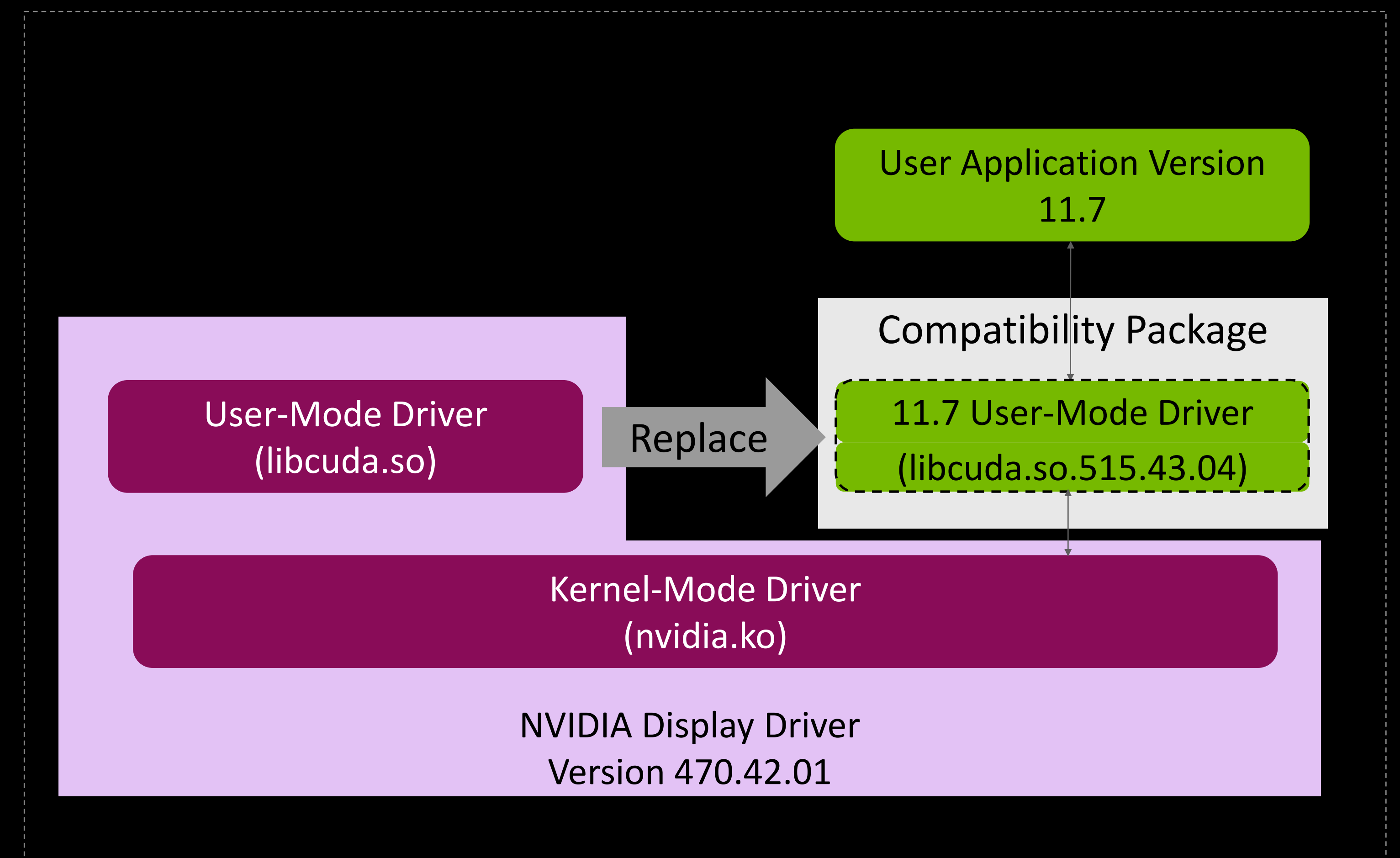


* See NVIDIA “Best Practices” Documentation

CUDA Compatibility

“Forward Compatibility”

- ▶ Using a CUDA toolkit with **higher-versioned** UMD with a **lower-versioned** KMD
- ▶ Deployment & upgrade of Drivers may be very disruptive, especially in CSP and enterprise datacenters
- ▶ Can be used across **major** and minor versions of CTK
- ▶ Compatibility Package to be installed, includes user-mode driver (among other files)
 - ▶ Via symbolic links, multiple compatibility versions can be installed together in a single system
- ▶ Programmers must check for supported features & supported hardware
- ▶ Supports PTX JIT compilation



Key Takeaways

- CUDA applications are compatible – forever
- CUDA programs within a major version – *generally* are compatible
- CUDA applications run against older drivers – with compatibility shims
 - Matters in e.g. containers, data center environments

Multi GPU Multi Node programming

Example: Jacobi solver

Solves the 2D-Laplace Equation on a rectangle

$$\Delta \mathbf{u}(x, \mathbf{y}) = \mathbf{0} \quad \forall (x, \mathbf{y}) \in \Omega \setminus \delta\Omega$$

Dirichlet boundary conditions (constant values on boundaries) on left and right boundary

Periodic boundary conditions on top and bottom boundary

Example: Jacobi Solver

Single GPU

While not converged

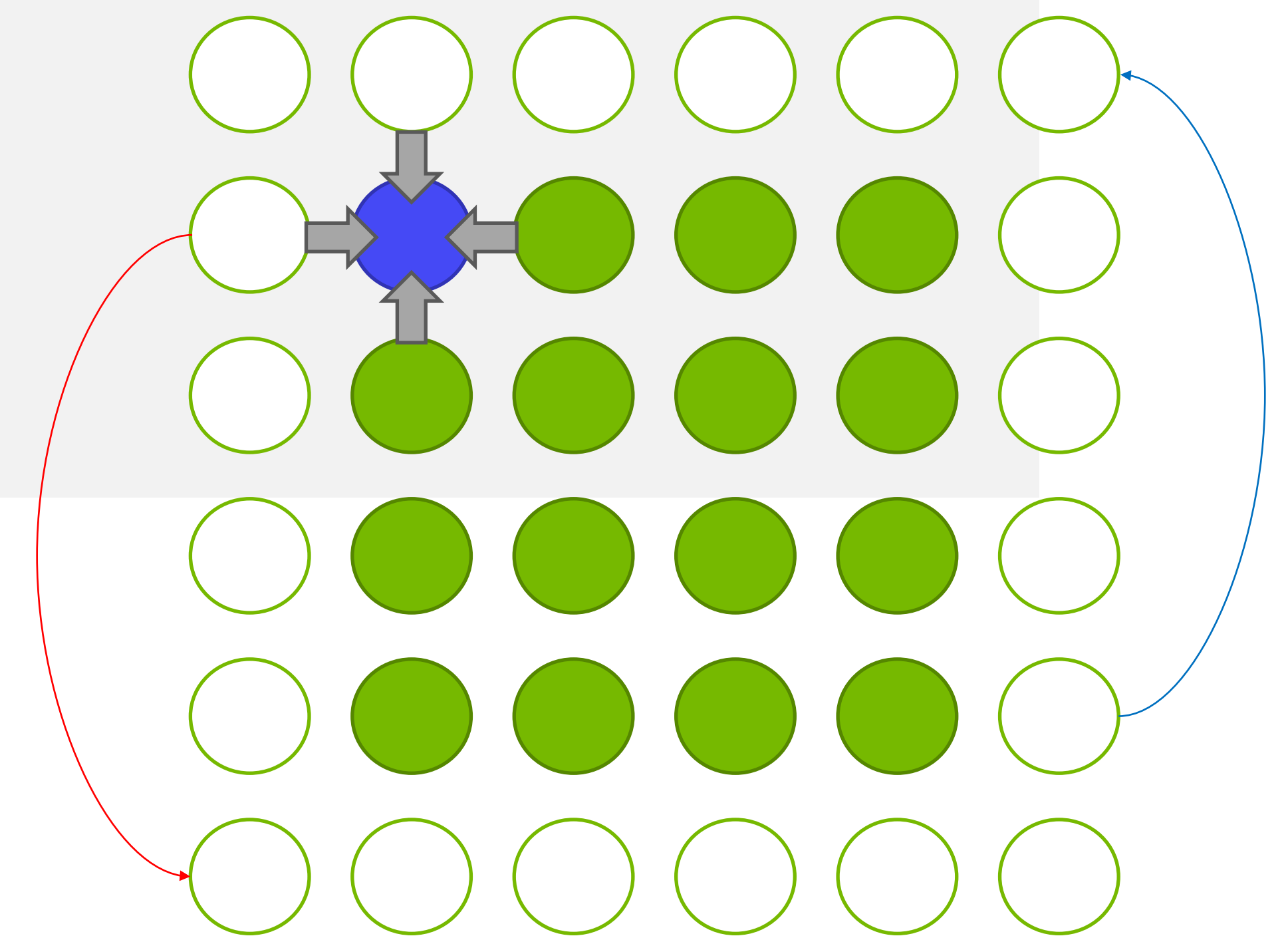
Do Jacobi step:

```
for( int iy = 1      ; iy < ny-1  ; iy++ )
for( int ix = 1      ; ix < nx-1  ; ix++ )
    a_new[iy*nx+ix] = -0.25 *
        -( a[ iy      *nx+(ix+1)] + a[ iy      *nx+ix-1]
          + a[(iy-1)*nx+ ix      ] + a[(iy+1)*nx+ix      ] );
```

Apply periodic boundary conditions

Swap a_new and a

Next iteration



Example: Jacobi Solver

Multi GPU

While not converged

Do Jacobi step:

```
for( int iy = iy_start; iy < iy_end; iy++ )
for( int ix = 1 ; ix < nx-1 ; ix++ )
  a_new[iy*nx+ix] = -0.25 *
    -( a[ iy *nx+(ix+1)] + a[ iy *nx+ix-1]
      + a[(iy-1)*nx+ ix ] + a[(iy+1)*nx+ix ] );
```

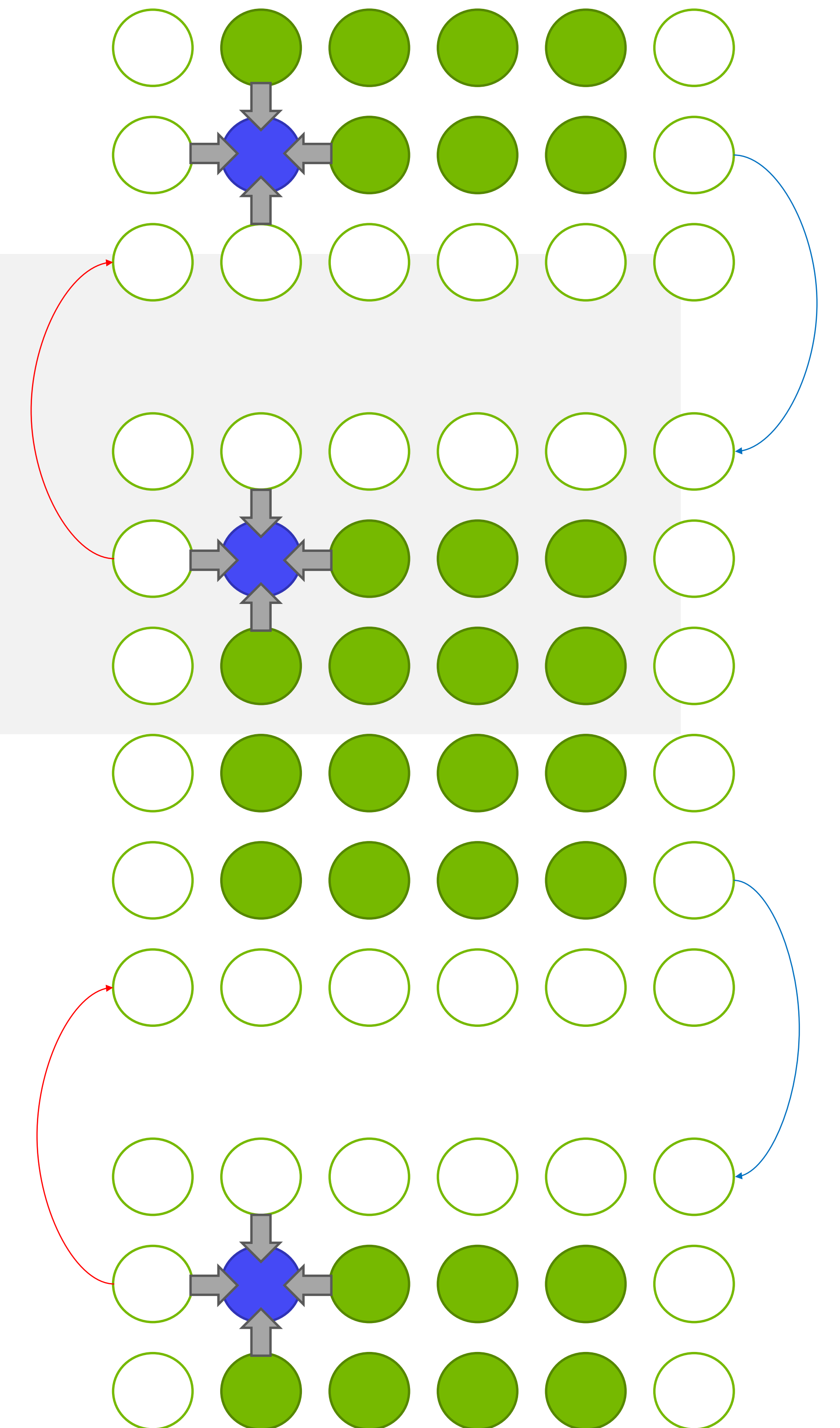
Apply periodic boundary conditions

Halo exchange

Swap a_new and a

Next iteration

One-step with ring exchange



Message Passing Interface - MPI

- Standard to exchange data between processes via messages
 - Defines API to exchanges messages
 - Point to Point: e.g. MPI_Send, MPI_Recv
 - Collectives: e.g. MPI_Reduce
- Multiple implementations (open source and commercial)
 - Bindings for C/C++, Fortran, Python, ...
 - E.g. MPICH, OpenMPI, MVAPICH, IBM Platform MPI, Cray MPT, ...

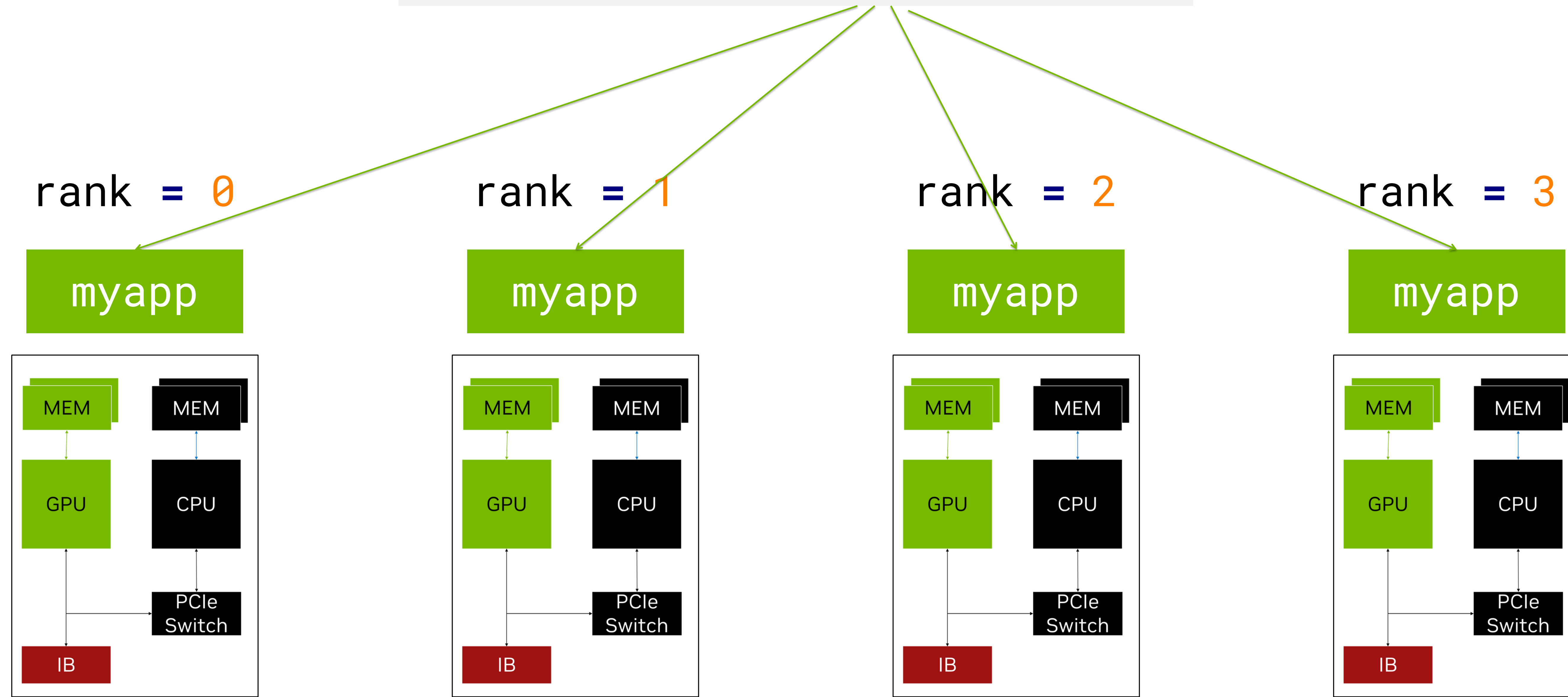
MPI - Skeleton

```
#include <mpi.h>
int main(int argc, char *argv[]) {
    int rank, size;
    /* Initialize the MPI library */
    MPI_Init(&argc, &argv);
    /* Determine the calling process rank and total number of ranks */
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    /* Call MPI routines like MPI_Send, MPI_Recv, ... */
    ...
    /* Shutdown MPI library */
    MPI_Finalize();
    return 0;
}
```

MPI

Launching

```
$ mpirun -np 4 ./myapp <args>
```



Multi Process Multi GPU Programming

Using CUDA-aware MPI

Handle GPU affinity on multi-GPU nodes:

```
int local_rank = -1;
MPI_Comm_rank(local_comm, &local_rank);
int num_devices = 0;
cudaGetDeviceCount(&num_devices);
cudaSetDevice(local_rank % num_devices);
```

(Use `MPI_Comm_split_type(MPI_COMM_WORLD, MPI_COMM_TYPE_SHARED, rank, info, &local_comm)`; to get `local_comm`.)

Multi Process Multi GPU Programming

Using CUDA-aware MPI

```
while (l2_norm > tol && iter < iter_max) {
    cudaMemsetAsync(l2_norm_d, 0, sizeof(real), compute_stream);
    launch_jacobi_kernel(a_new, a, l2_norm_d, iy_start, iy_end, nx, compute_stream);
    cudaEventRecord(compute_done, compute_stream);
    cudaMemcpyAsync(l2_norm_h, l2_norm_d, sizeof(real), cudaMemcpyDeviceToHost, compute_stream);

    cudaEventSynchronize(compute_done);
    const int top = rank > 0 ? rank - 1 : (size - 1);
    const int bottom = (rank + 1) % size;
    // Top/Bottom Halo exchange -> next slide

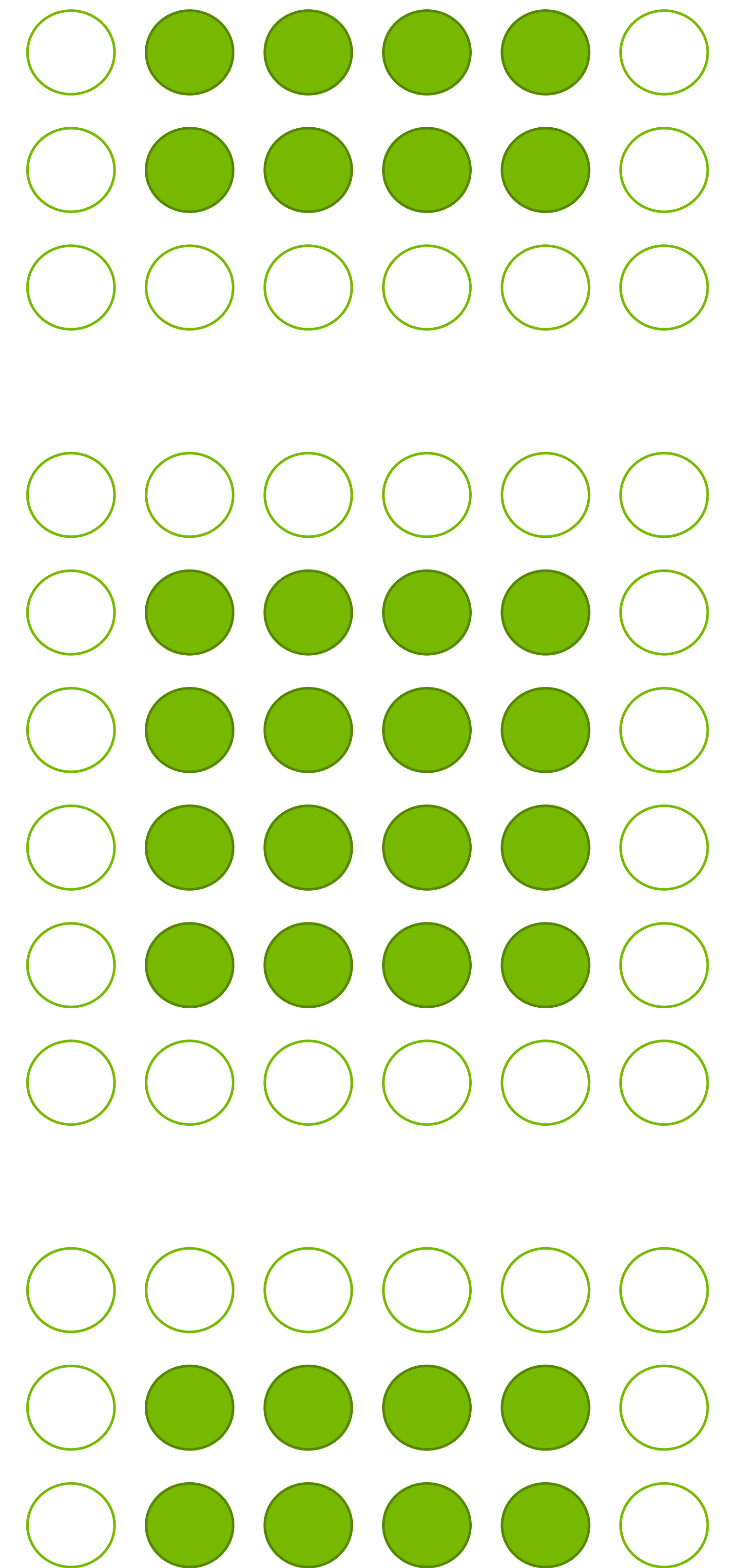
    cudaStreamSynchronize(compute_stream);
    MPI_CALL(MPI_Allreduce(l2_norm_h, &l2_norm, 1, MPI_REAL_TYPE, MPI_SUM, MPI_COMM_WORLD));
    l2_norm = std::sqrt(l2_norm);

    std::swap(a_new, a); iter++;
}
```

Example Jacobi

Top/Bottom Halo

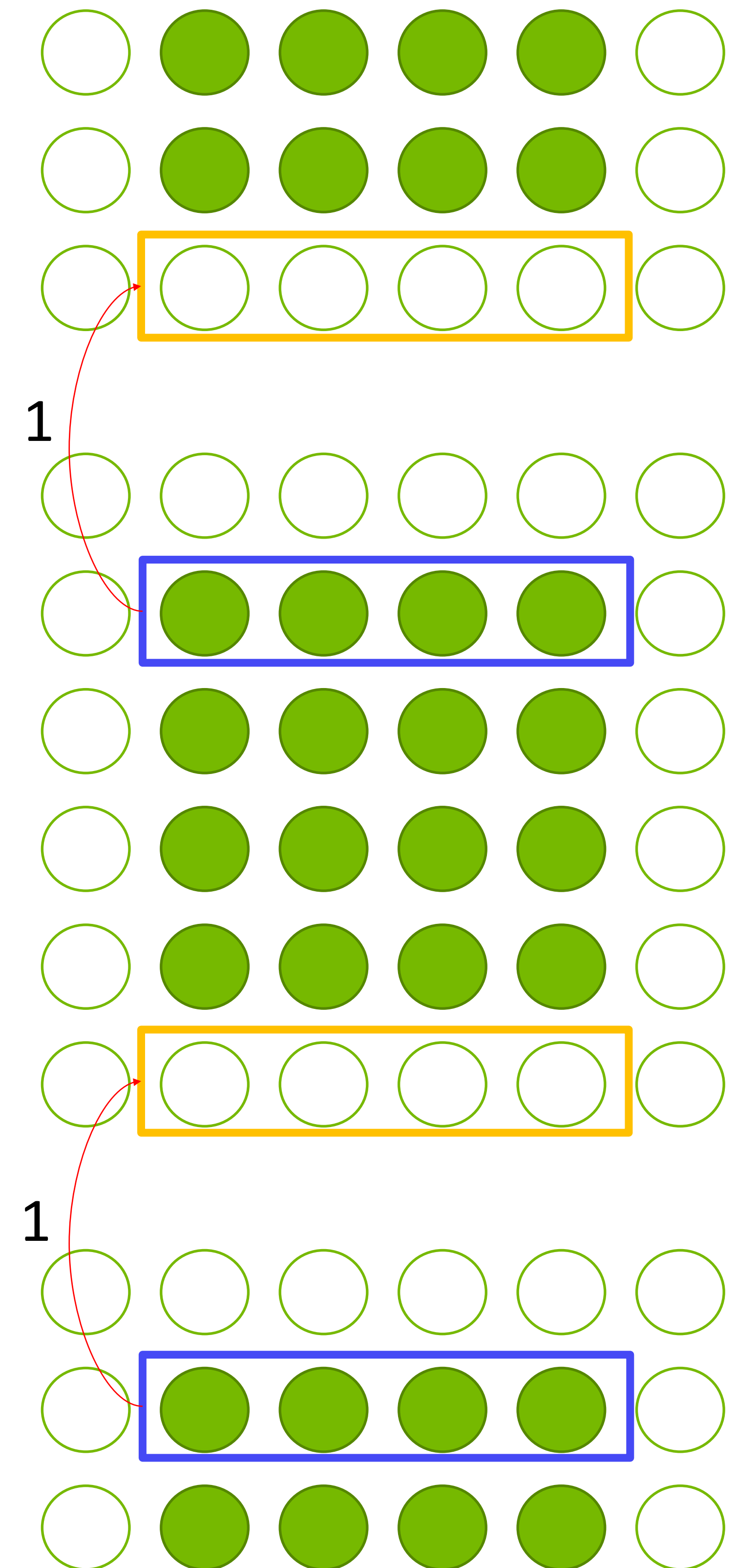
```
MPI_Sendrecv(a_new+iy_start*nx, nx, MPI_FLOAT, top, 0,  
             a_new+(iy_end*nx), nx, MPI_FLOAT, bottom, 0,  
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
```



Example Jacobi

Top/Bottom Halo

```
MPI_Sendrecv(a_new+iy_start*nx, nx, MPI_FLOAT, top, 0,  
a_new+(iy_end*nx), nx, MPI_FLOAT, bottom, 0,  
MPI_COMM_WORLD, MPI_STATUS_IGNORE);
```

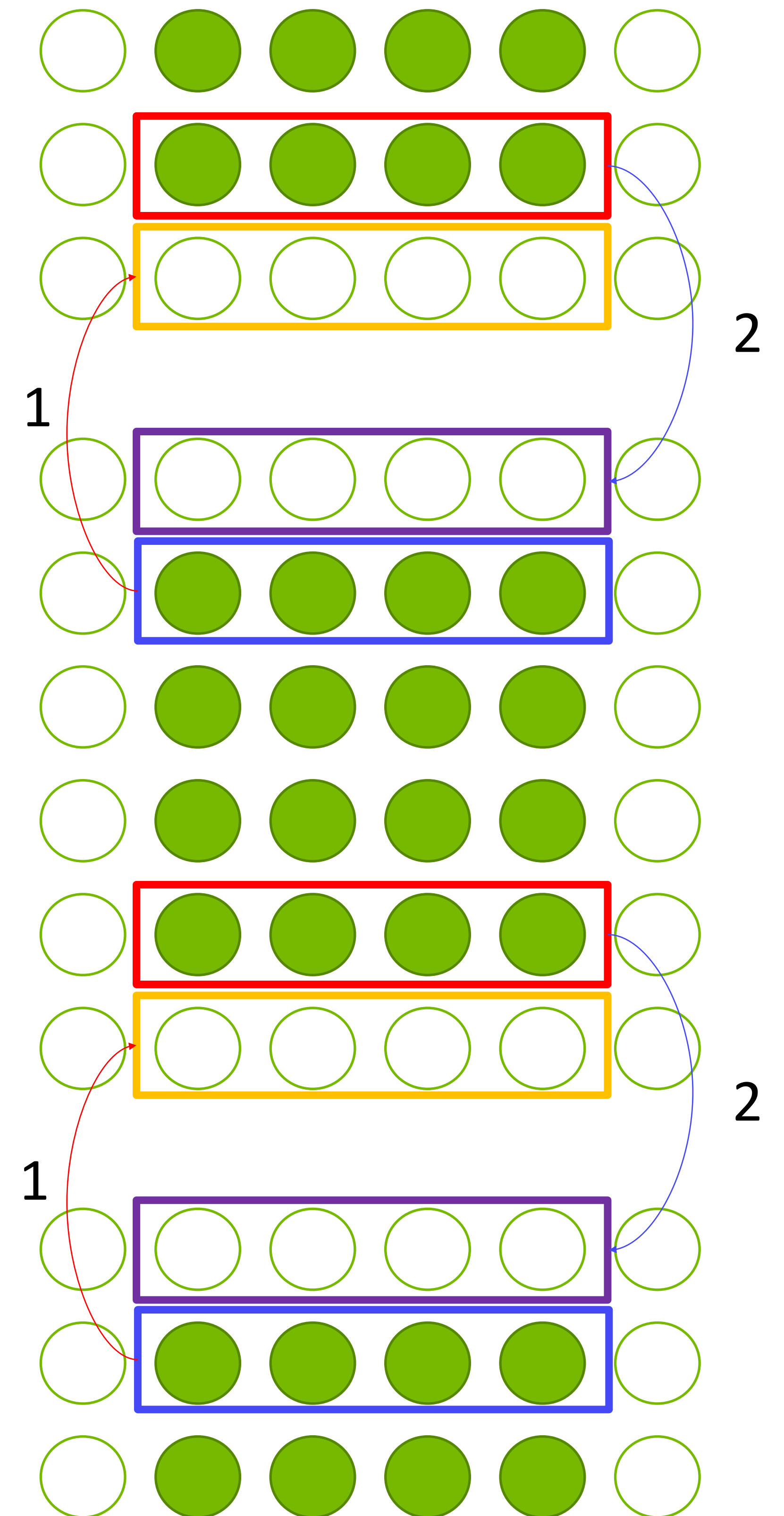


Example Jacobi

Top/Bottom Halo

```
MPI_Sendrecv(a_new+iy_start*nx, nx, MPI_FLOAT, top, 0,  
a_new+(iy_end*nx), nx, MPI_FLOAT, bottom, 0,  
MPI_COMM_WORLD, MPI_STATUS_IGNORE);
```

```
MPI_Sendrecv(a_new+(iy_end-1)*nx, nx, MPI_FLOAT, bottom, 0,  
a_new, nx, MPI_FLOAT, top, 0,  
MPI_COMM_WORLD, MPI_STATUS_IGNORE);
```



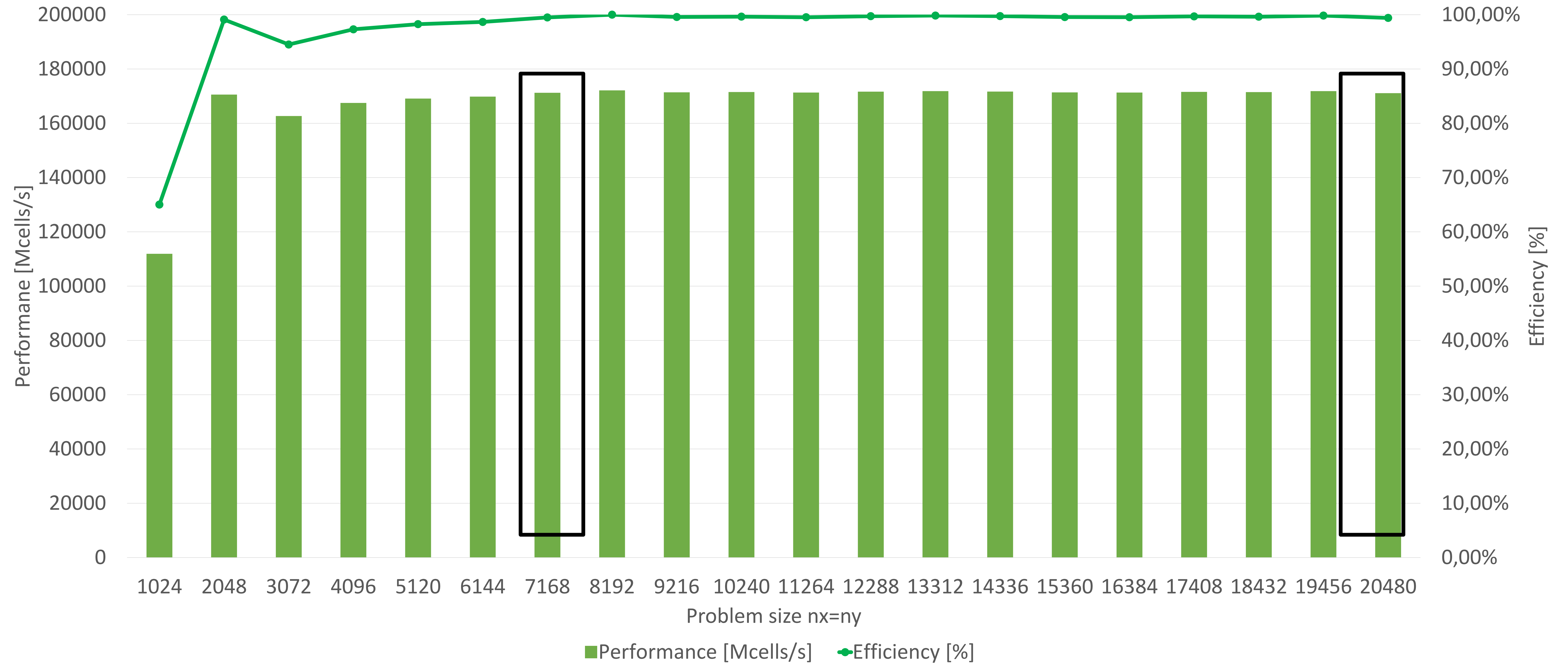
Benchmark Setup

DGX H100

- CUDA Driver 535.129.03
- [NVIDIA HPC SDK container](#): Tag `nvcr.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04`
- GPUs@1980 Mhz
- Reported Runtime is the minimum of 5 repetitions
- For all runs CPU and GPU affinities have been tuned: See `bench.sh` in <https://github.com/NVIDIA/multi-gpu-programming-models>

Example: Jacobi Solver

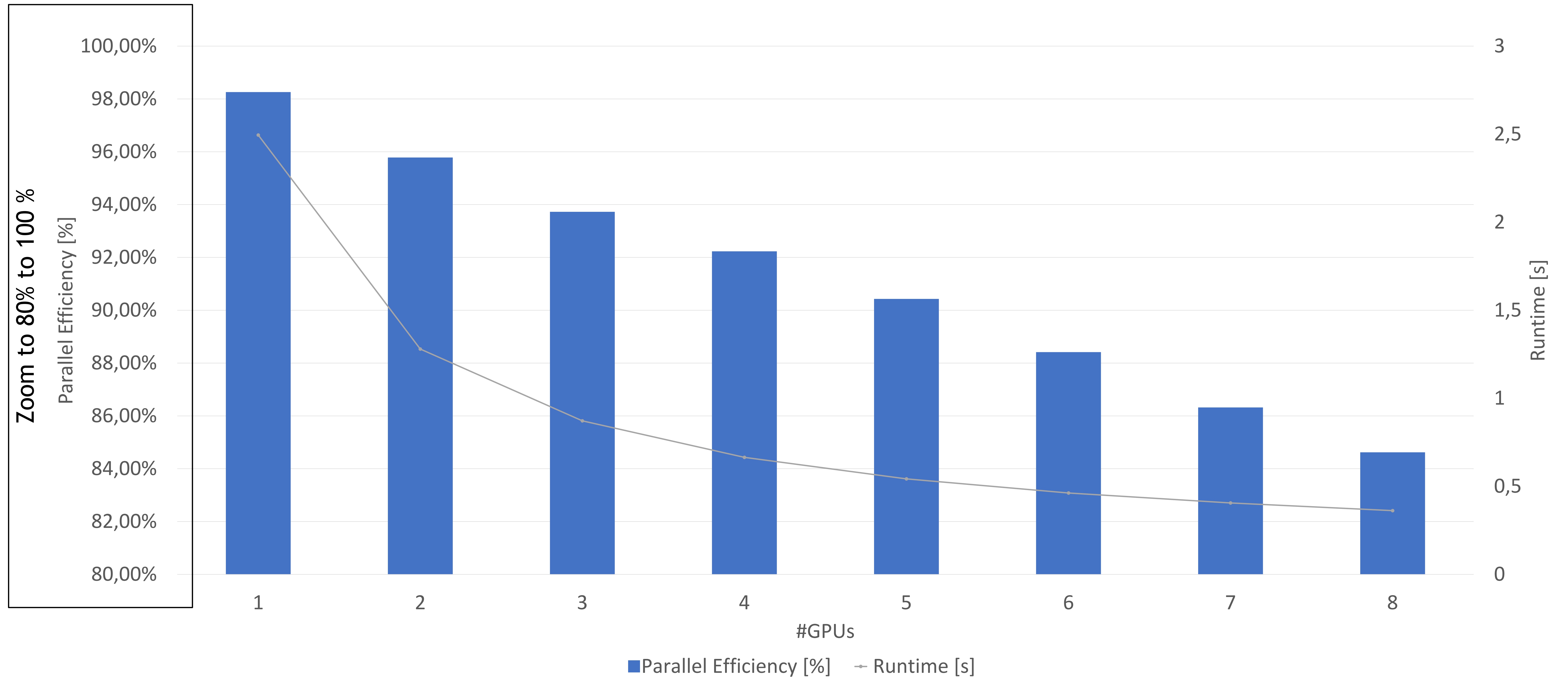
Single GPU performance vs. problem size – NVIDIA H100 80GB HBM3



Benchmarksetup: DGX H100, CUDA Driver 535.129.03, NVIDIA HPC SDK container: `nvcr.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04`, GPUs@1980Mhz AC, Reported Runtime is the minimum of 5 repetitions

Multi GPU Jacobi Runtime And Parallel Efficiency

MPI on DGX H100 – 20480 x 20480, 1000 iterations



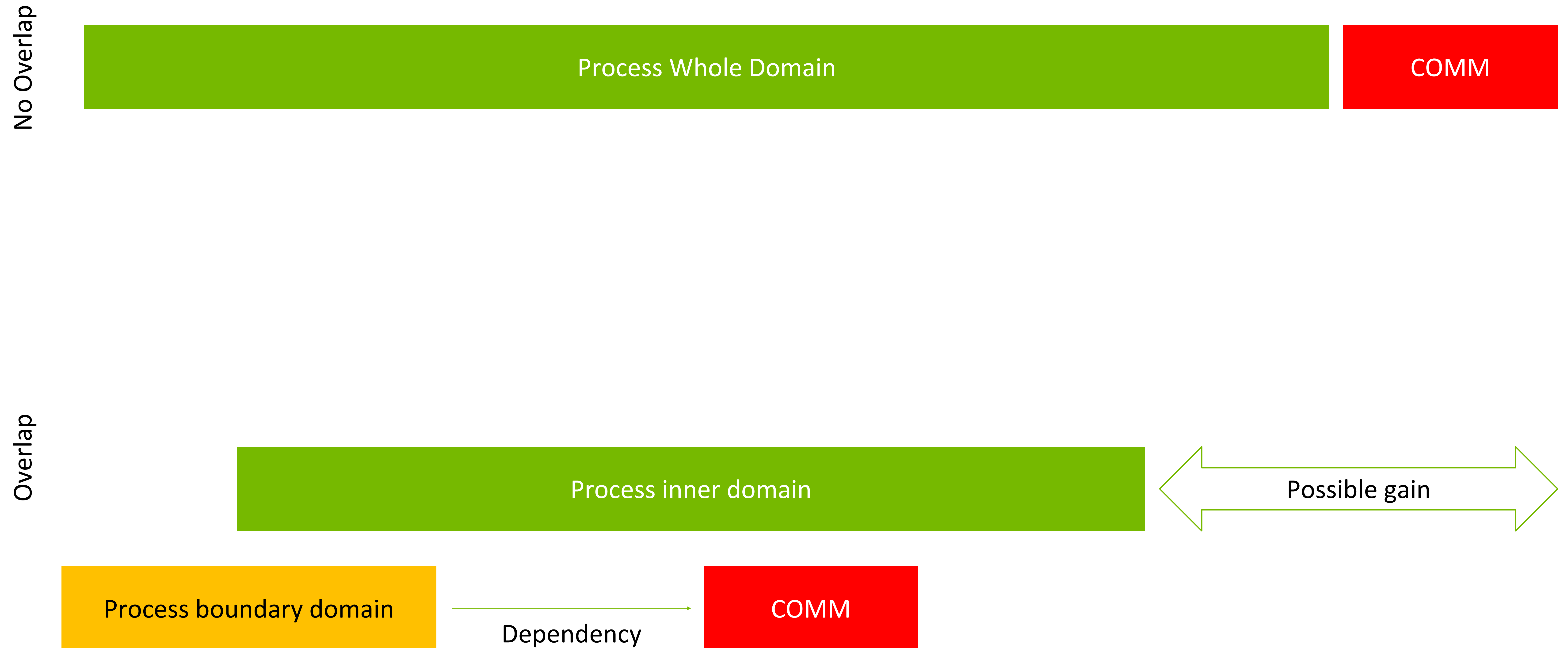
Benchmark setup: DGX H100, CUDA Driver 535.129.03, NVIDIA HPC SDK container: `nvcr.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04`, GPUs@1980Mhz AC, Reported Runtime is the minimum of 5 repetitions

Multi GPU Jacobi Nsight Systems Timeline

MPI 8 NVIDIA H100 80GB HBM3 on DGX H100



Overlapping Communication and Computation



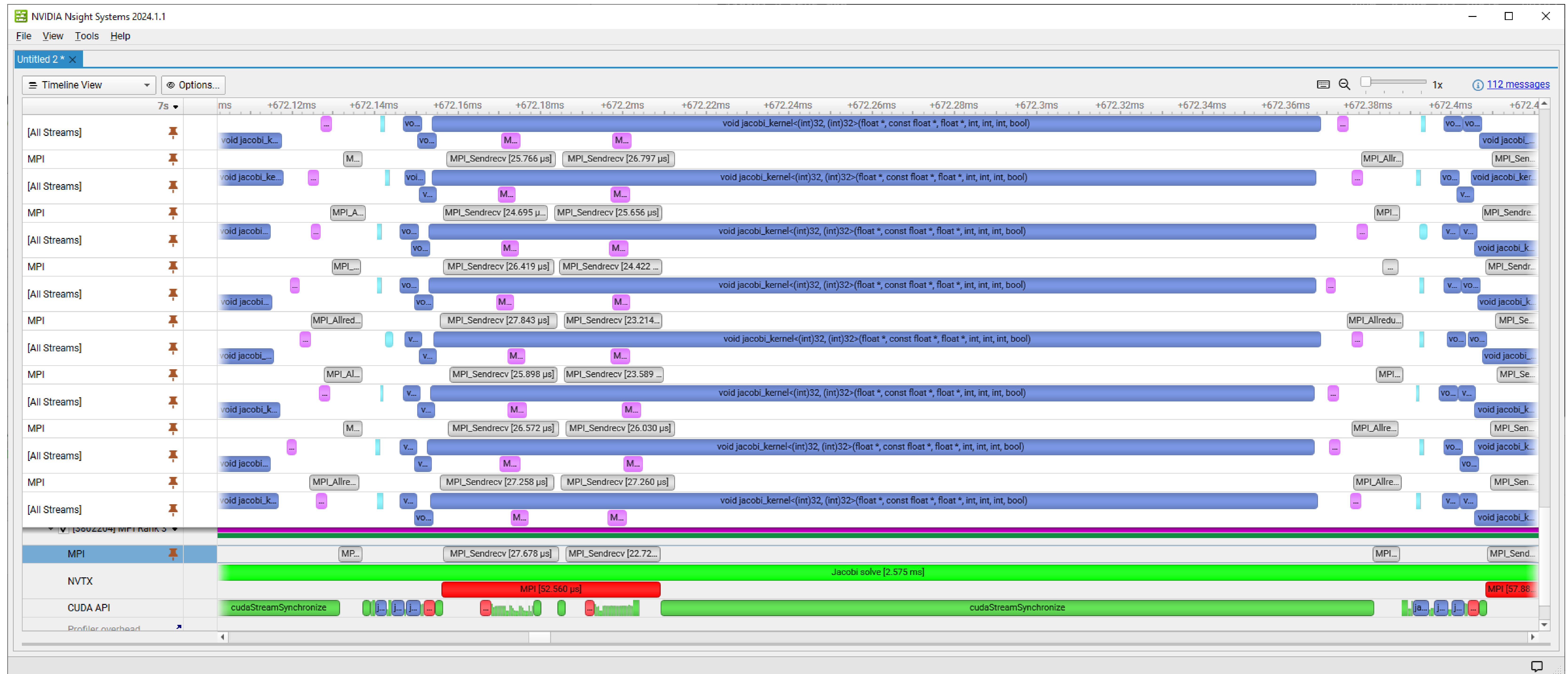
MPI

Overlapping Communication and Computation

```
launch_jacobi_kernel( a_new, a, l2_norm_d, iy_start, (iy_start+1), nx, push_top_stream );
launch_jacobi_kernel( a_new, a, l2_norm_d, (iy_end-1), iy_end, nx, push_bottom_stream );
launch_jacobi_kernel( a_new, a, l2_norm_d, (iy_start+1), (iy_end-1), nx, compute_stream );
const int top = rank > 0 ? rank - 1 : (size-1);
const int bottom = (rank+1)%size;
cudaStreamSynchronize( push_top_stream );
MPI_Sendrecv( a_new+iy_start*nx, nx, MPI_REAL_TYPE, top, 0,
              a_new+(iy_end*nx), nx, MPI_REAL_TYPE, bottom, 0,
              MPI_COMM_WORLD, MPI_STATUS_IGNORE );
cudaStreamSynchronize( push_bottom_stream );
MPI_Sendrecv( a_new+(iy_end-1)*nx, nx, MPI_REAL_TYPE, bottom, 0,
              a_new, nx, MPI_REAL_TYPE, top, 0, MPI_COMM_WORLD,
              MPI_STATUS_IGNORE );
```

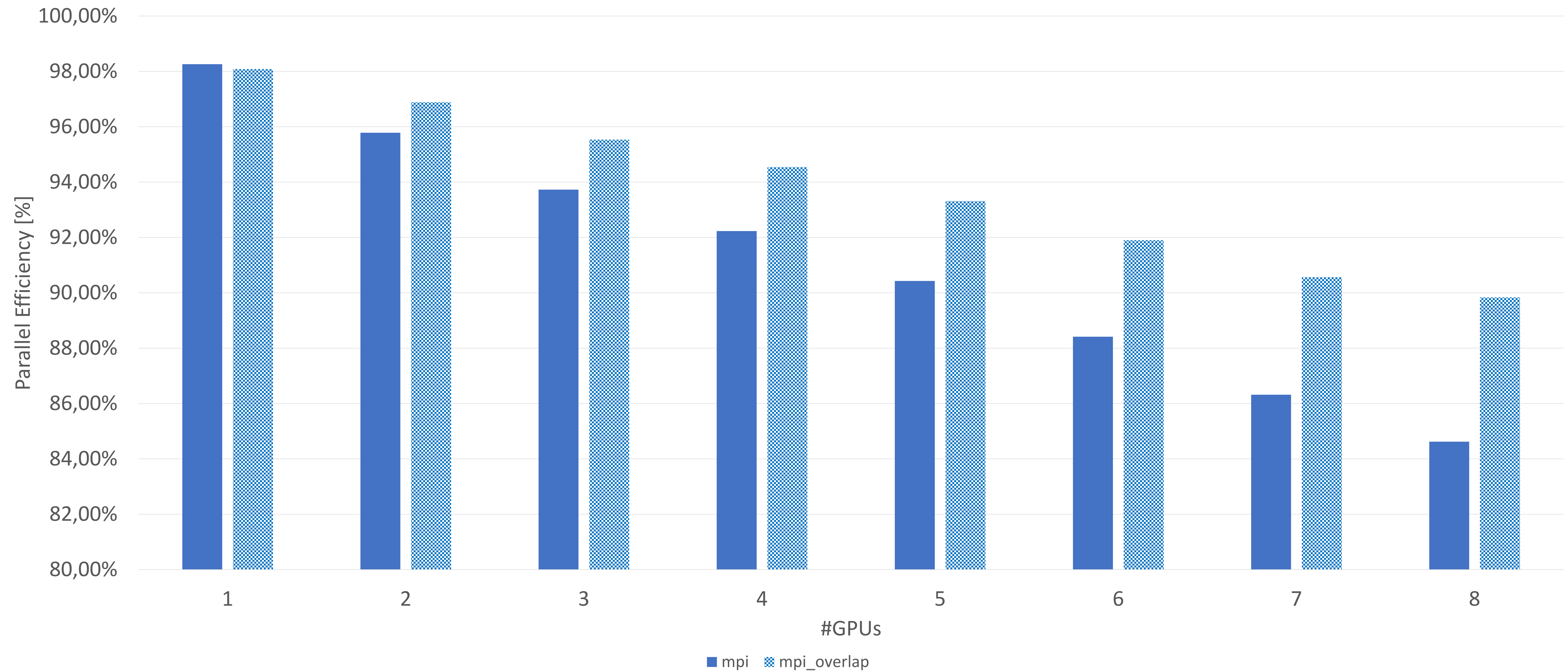
Multi GPU Jacobi Nsight Systems Timeline

MPI Overlap 8 NVIDIA H100 80GB HBM3 on DGX H100



Multi GPU Jacobi Parallel Efficiency

DGX H100 – 20480 x 20480, 1000 iterations



Benchmarksetup: DGX H100, CUDA Driver 535.129.03, NVIDIA HPC SDK container: `nvcr.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04`, GPUs@1980Mhz AC, Reported Runtime is the minimum of 5 repetitions

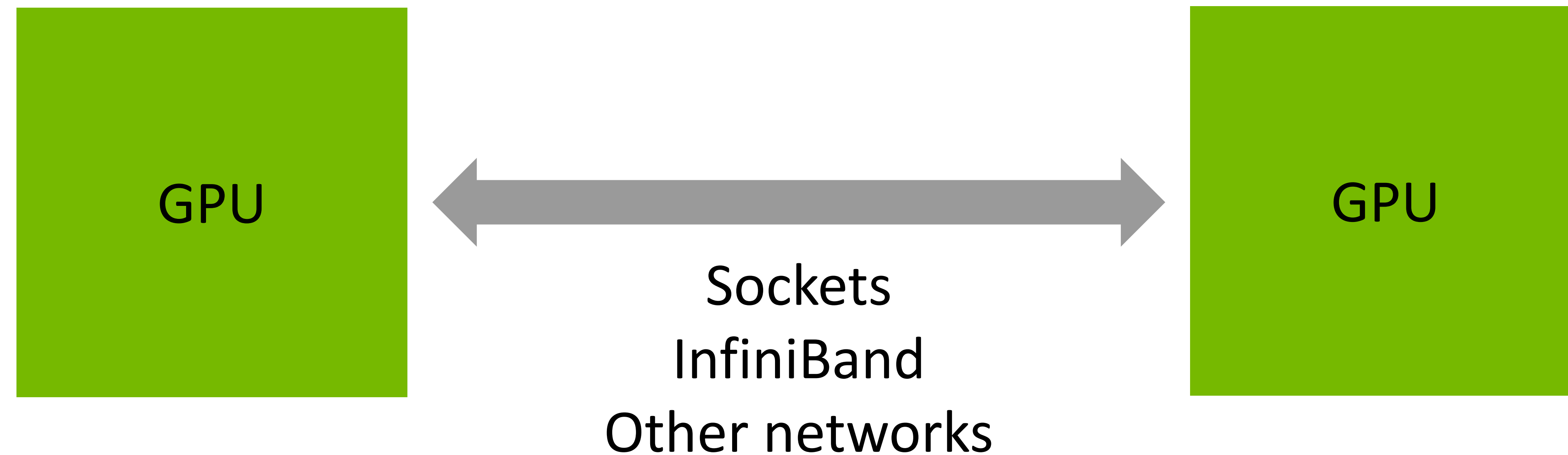
NCCL

Optimized inter-GPU communication

NCCL : NVIDIA Collective Communication Library

Communication library running on GPUs, for GPU buffers.

- Library for efficient communication with GPUs
- First: Collective Operations (e.g. Allreduce), as they are required for Deep Learning
- Since 2.8: Support for Send/Recv between GPUs
- Library running on GPU: Communication calls are translated to a GPU kernel (running on a stream)



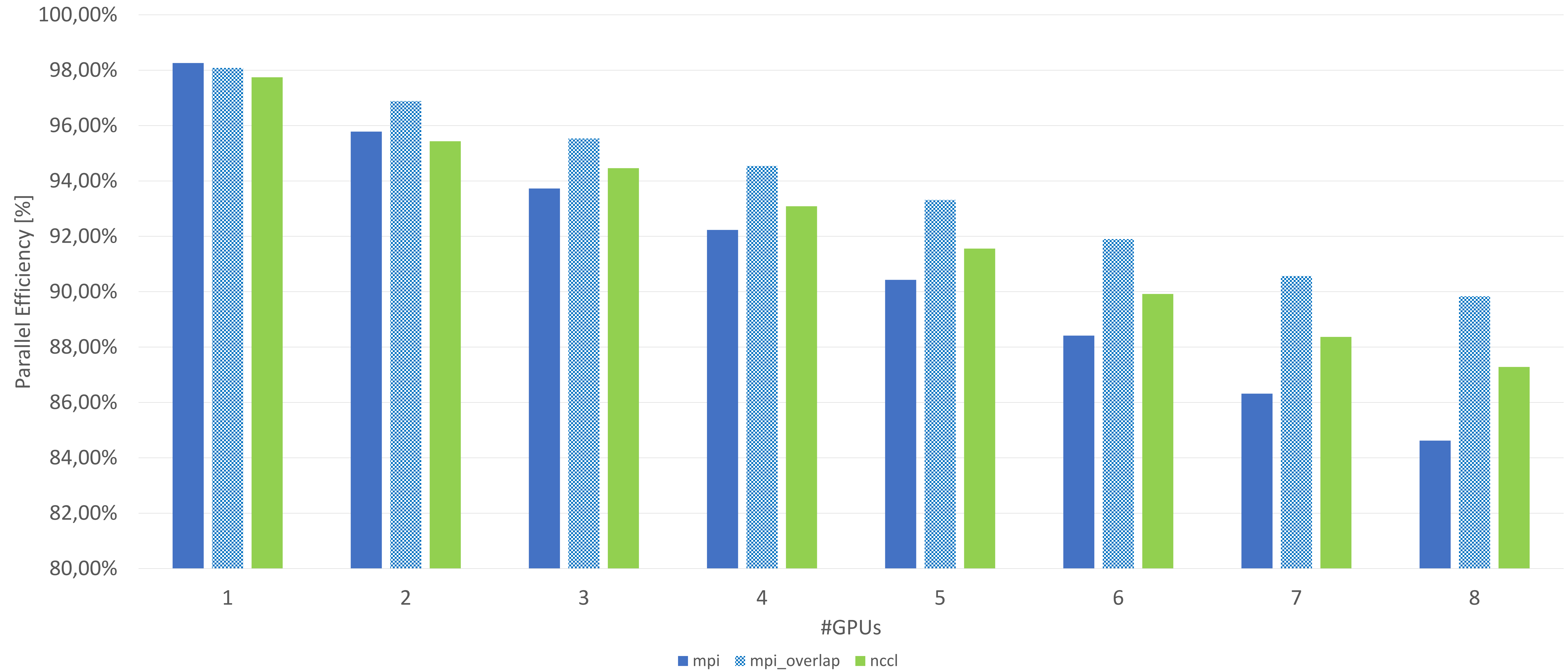
Binaries : <https://developer.nvidia.com/nccl> and in NGC containers

Source code : <https://github.com/nvidia/nccl>

Perf tests : <https://github.com/nvidia/nccl-tests>

Multi GPU Jacobi Parallel Efficiency

DGX H100 – 20480 x 20480, 1000 iterations



Benchmarksetup: DGX H100, CUDA Driver 535.129.03, NVIDIA HPC SDK container: `nvcr.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04`, GPUs@1980Mhz AC, Reported Runtime is the minimum of 5 repetitions

NCCL

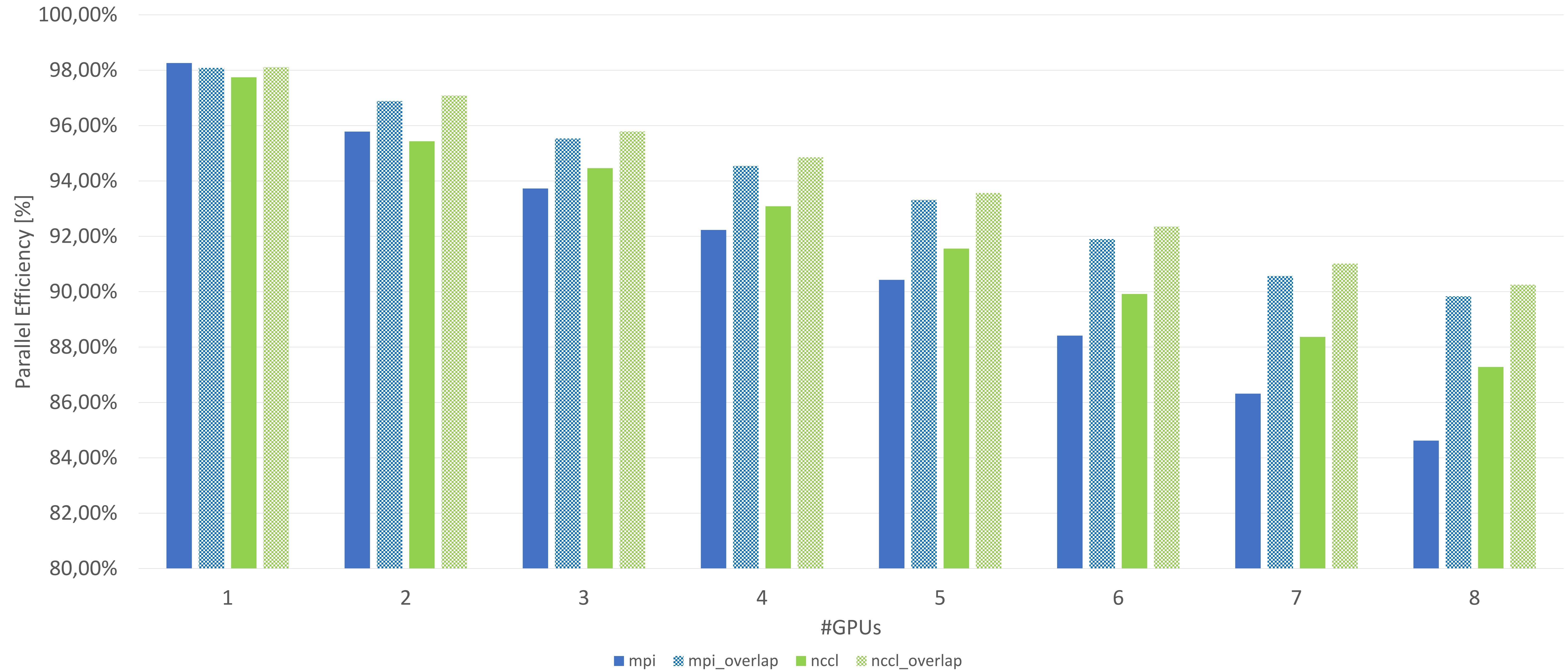
Overlapping Communication and Computation

```
int leastPriority = 0;
int greatestPriority = leastPriority;
cudaDeviceGetStreamPriorityRange(&leastPriority, &greatestPriority);
cudaStream_t compute_stream, push_stream;
cudaStreamCreateWithPriority(&compute_stream, cudaStreamDefault, leastPriority);
cudaStreamCreateWithPriority(&push_stream, cudaStreamDefault, greatestPriority);
...
launch_jacobi_kernel(a_new, a, l2_norm_d, iy_start,      (iy_start + 1), nx, push_stream);
launch_jacobi_kernel(a_new, a, l2_norm_d, (iy_end - 1), iy_end,      nx, push_stream);
launch_jacobi_kernel(a_new, a, l2_norm_d, (iy_start + 1), (iy_end - 1), nx, compute_stream);
ncclGroupStart();
ncclRecv(a_new,      nx, NCCL_REAL_TYPE, top, nccl_comm, push_stream)
ncclSend(a_new + (iy_end - 1) * nx, nx, NCCL_REAL_TYPE, btm, nccl_comm, push_stream);
ncclRecv(a_new + (iy_end * nx), nx, NCCL_REAL_TYPE, btm, nccl_comm, push_stream);
ncclSend(a_new + iy_start * nx, nx, NCCL_REAL_TYPE, top, nccl_comm, push_stream);
ncclGroupEnd();
```

Need to use CUDA high priority streams to avoid NCCL comms getting stuck behind compute.

Multi GPU Jacobi Parallel Efficiency

DGX H100 – 20480 x 20480, 1000 iterations

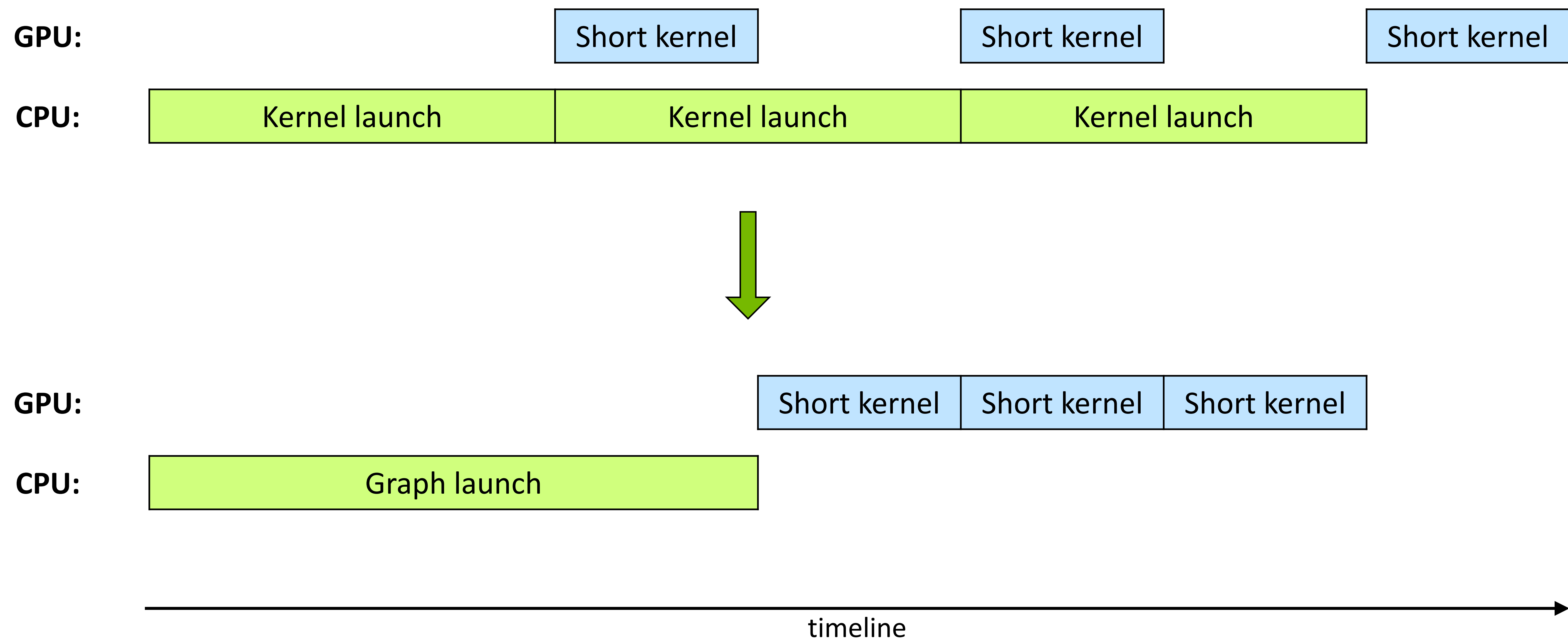


Benchmarksetup: DGX H100, CUDA Driver 535.129.03, NVIDIA HPC SDK container: `nvcr.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04`, GPUs@1980Mhz AC, Reported Runtime is the minimum of 5 repetitions

CUDA Graphs

Reducing launch overhead

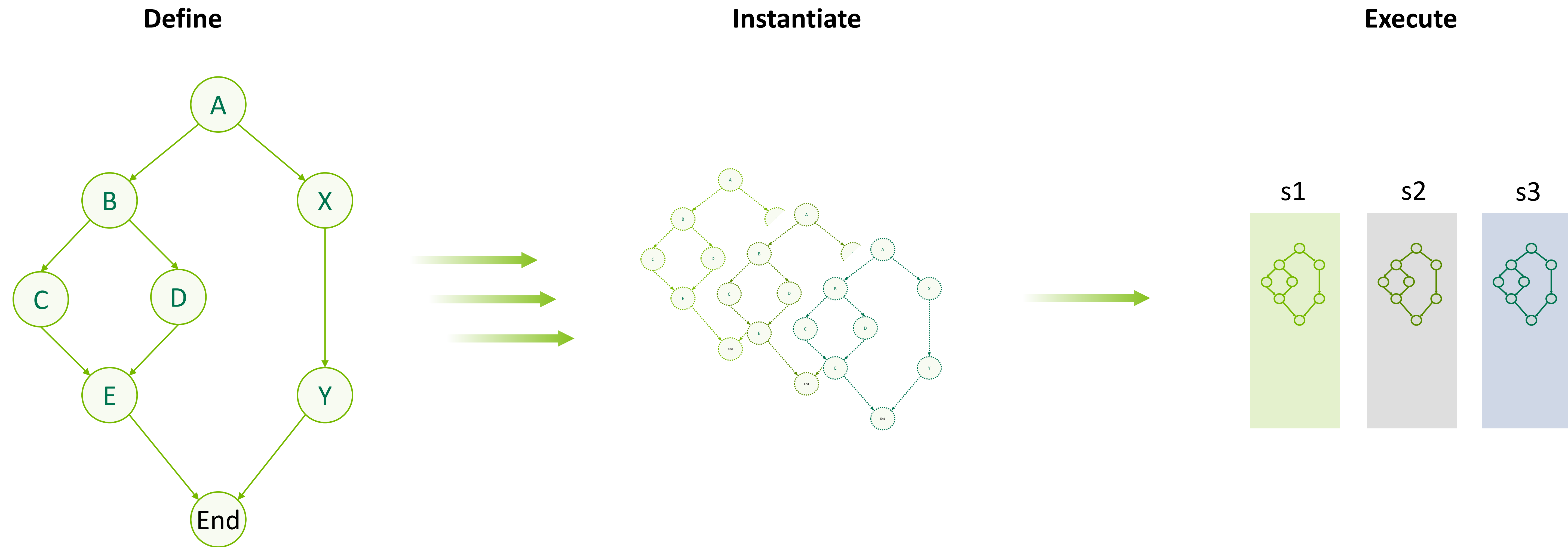
CUDA graphs reduce kernel launch latencies:



From [Advanced Performance Optimization in CUDA \[S62192\]](#) by Igor Terentyev more details on Graphs there.

Three-Stage Execution Model

Minimizes Execution Overheads – Pre-Initialize As Much As Possible



Single Graph “Template”

Created in host code
or built up from libraries

Multiple “Executable Graphs”

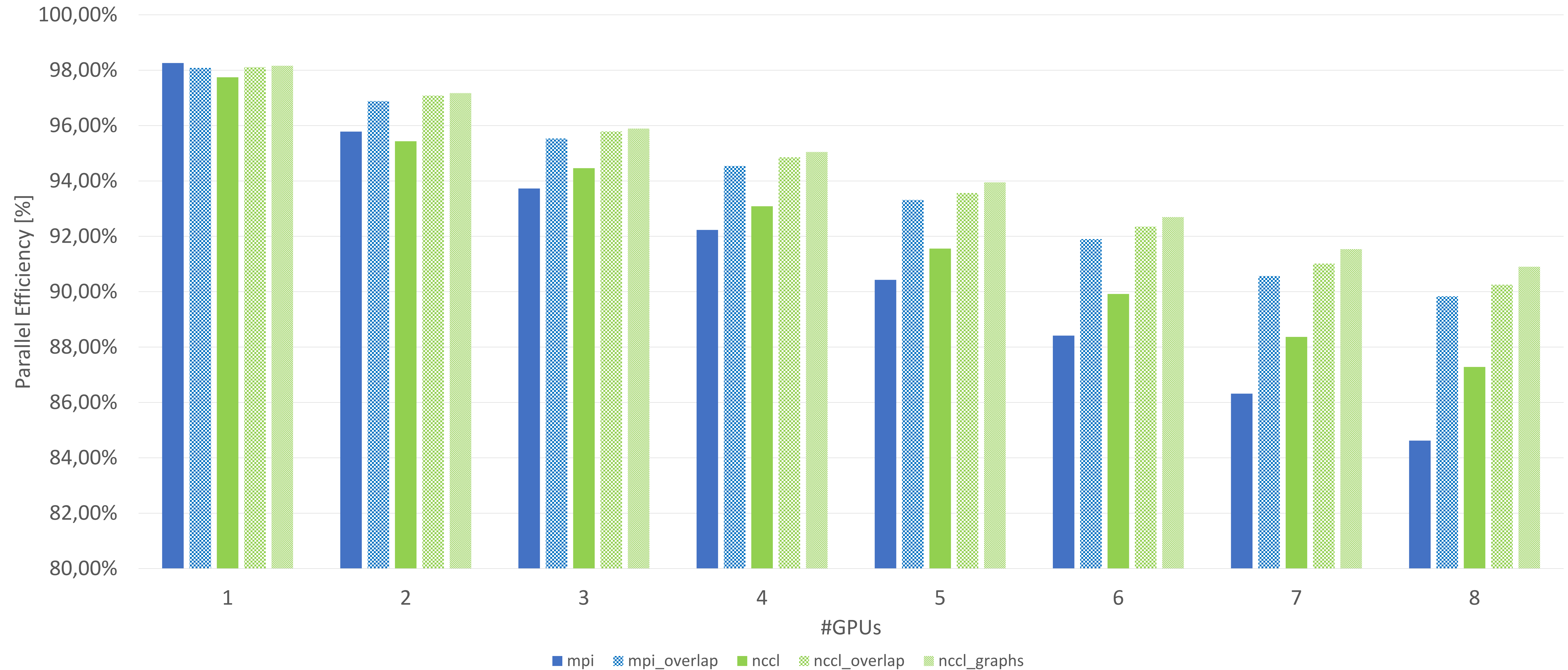
Snapshot of templates
Sets up & initializes GPU execution structures
(create once, run many times)

Executable Graphs Running in CUDA Streams

Concurrency in graph is **not** limited by stream

Multi GPU Jacobi Parallel Efficiency

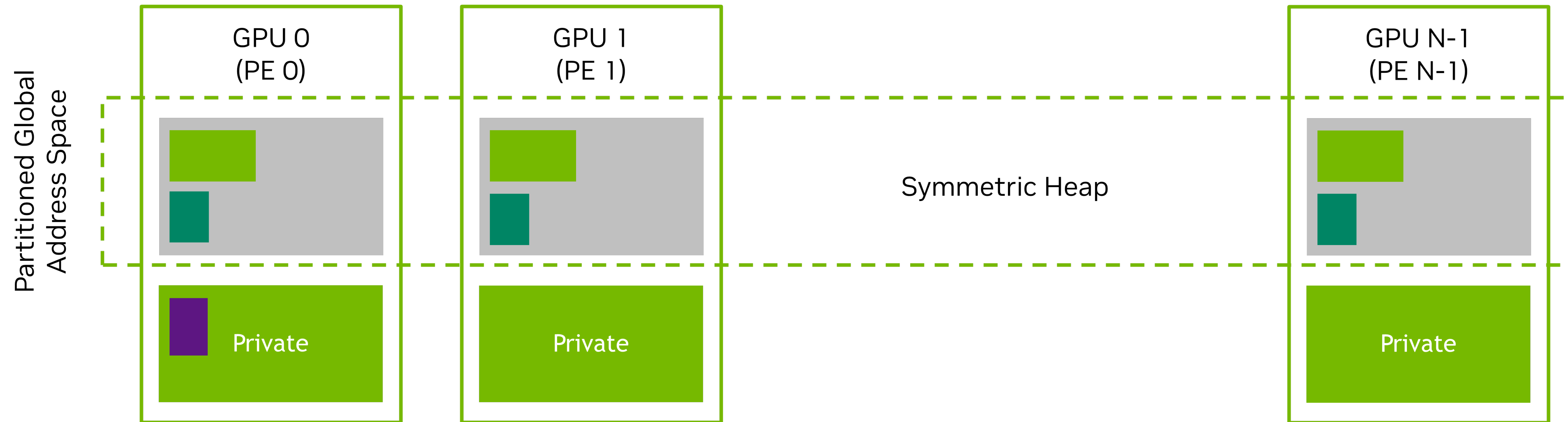
DGX H100 – 20480 x 20480, 1000 iterations



Benchmarksetup: DGX H100, CUDA Driver 535.129.03, NVIDIA HPC SDK container: `nvcr.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04`, GPUs@1980Mhz AC, Reported Runtime is the minimum of 5 repetitions

NVSHMEM

Implementation of OpenSHMEM, a Partitioned Global Address Space (PGAS) library



Symmetric objects are allocated collectively with the same size on every PE

Symmetric memory: `nvshmem_malloc(...)`; Private memory: `cudaMalloc(...)`;

CPU (blocking and stream-ordered) and CUDA Kernel interfaces

Read: `nvshmem_get(...)`; Write: `nvshmem_put(...)`; Atomic: `nvshmem_atomic_add(...)`;

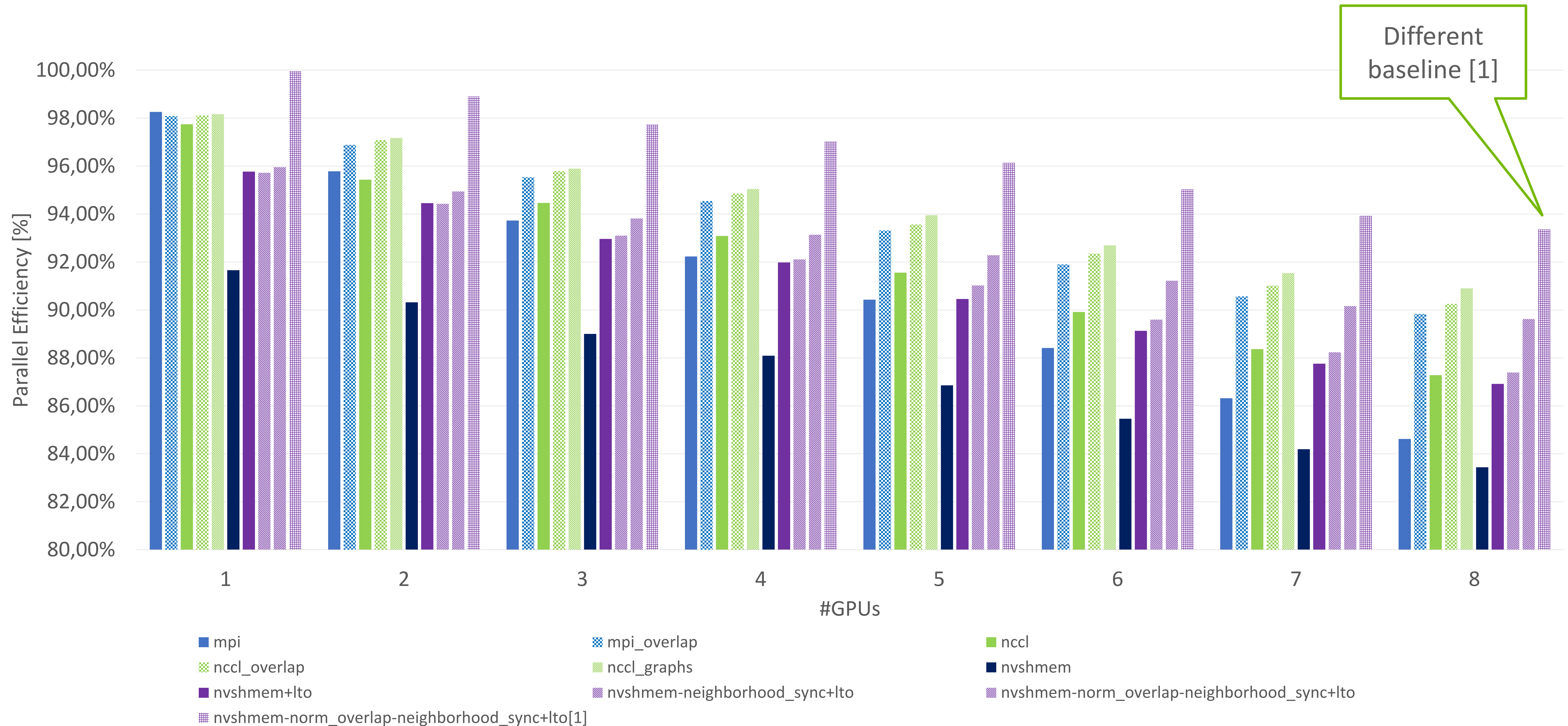
Flush writes: `nvshmem_quiet()`; Order writes: `nvshmem_fence()`;

Synchronize: `nvshmem_barrier()`; Poll: `nvshmem_wait_until(...)`;

Interoperable with MPI

Multi GPU Jacobi Parallel Efficiency

DGX H100 – 20480 x 20480, 1000 iterations



[1] Serial execution baseline uses same kernel as parallel version

Benchmarksetup: DGX H100, CUDA Driver 535.129.03, NVIDIA HPC SDK container: nvcv.io/nvidia/nvhpc:24.1-devel-cuda12.3-ubuntu22.04, GPUs@1980Mhz AC, Reported Runtime is the minimum of 5 repetitions

Conclusion

Thank you for your attention

	GPUDirect P2P/RDMA	CUDA stream/graph-aware	Kernel Initiated Communication
MPI	Improves Perf.	No	No
NCCL	Improves Perf.	Yes	No
NVSHMEM	Required	Yes	Yes

Source is on GitHub: <https://github.com/NVIDIA/multi-gpu-programming-models>

