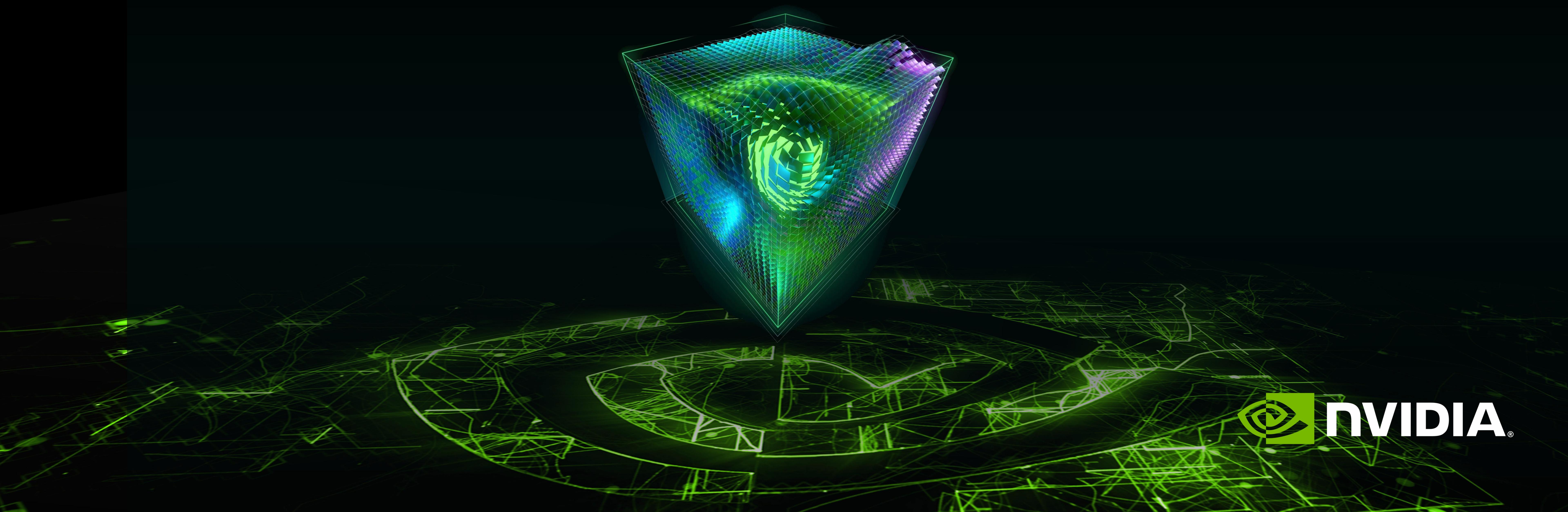


CUDA: New Features and Beyond

Stephen Jones, GTC 2023



Moore's Law is dead

- NVIDIA CEO Jensen Huang

Moore's Law is dead

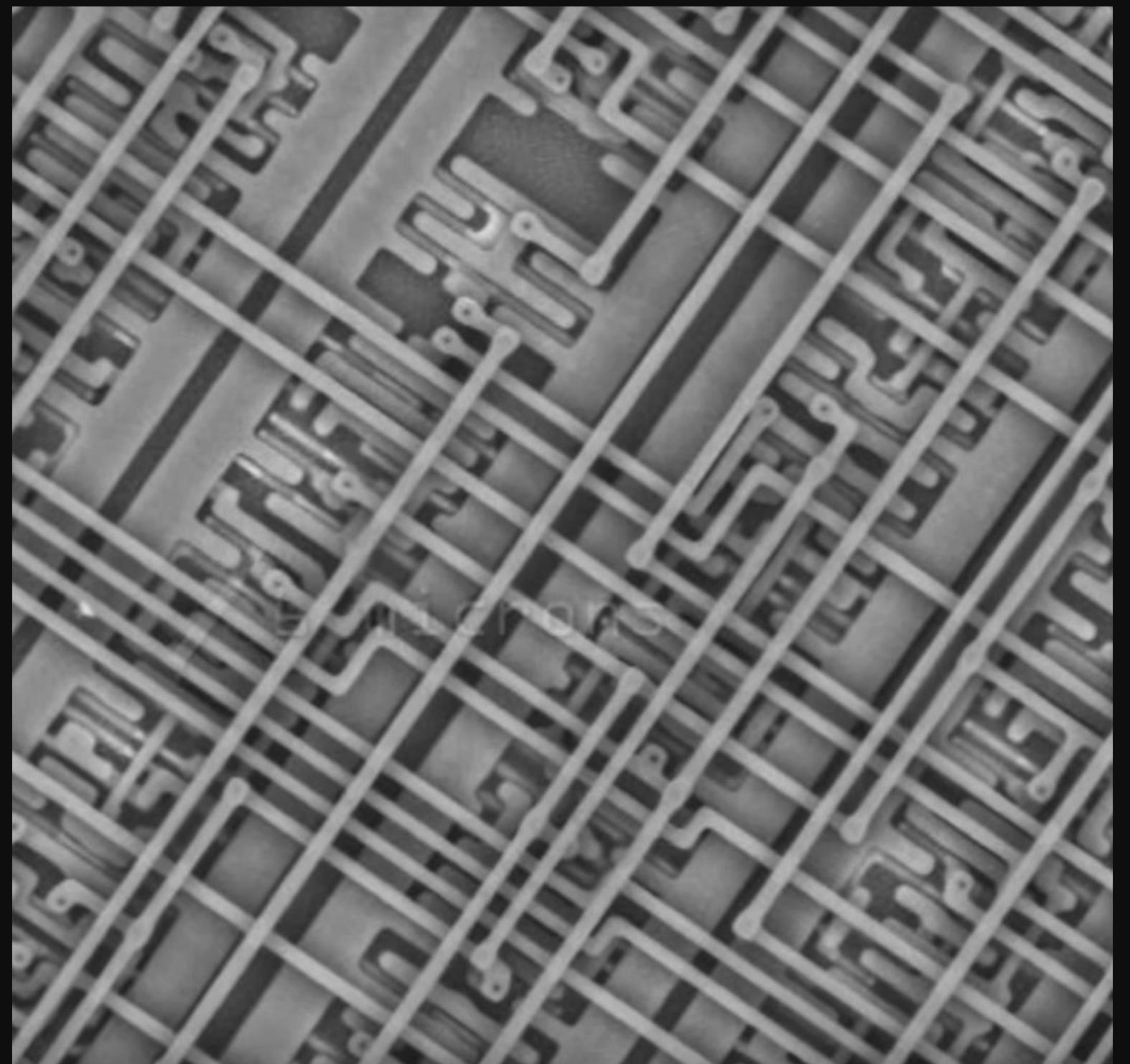
*... computing is not a chip problem
It's a software and chip problem*

- NVIDIA CEO Jensen Huang

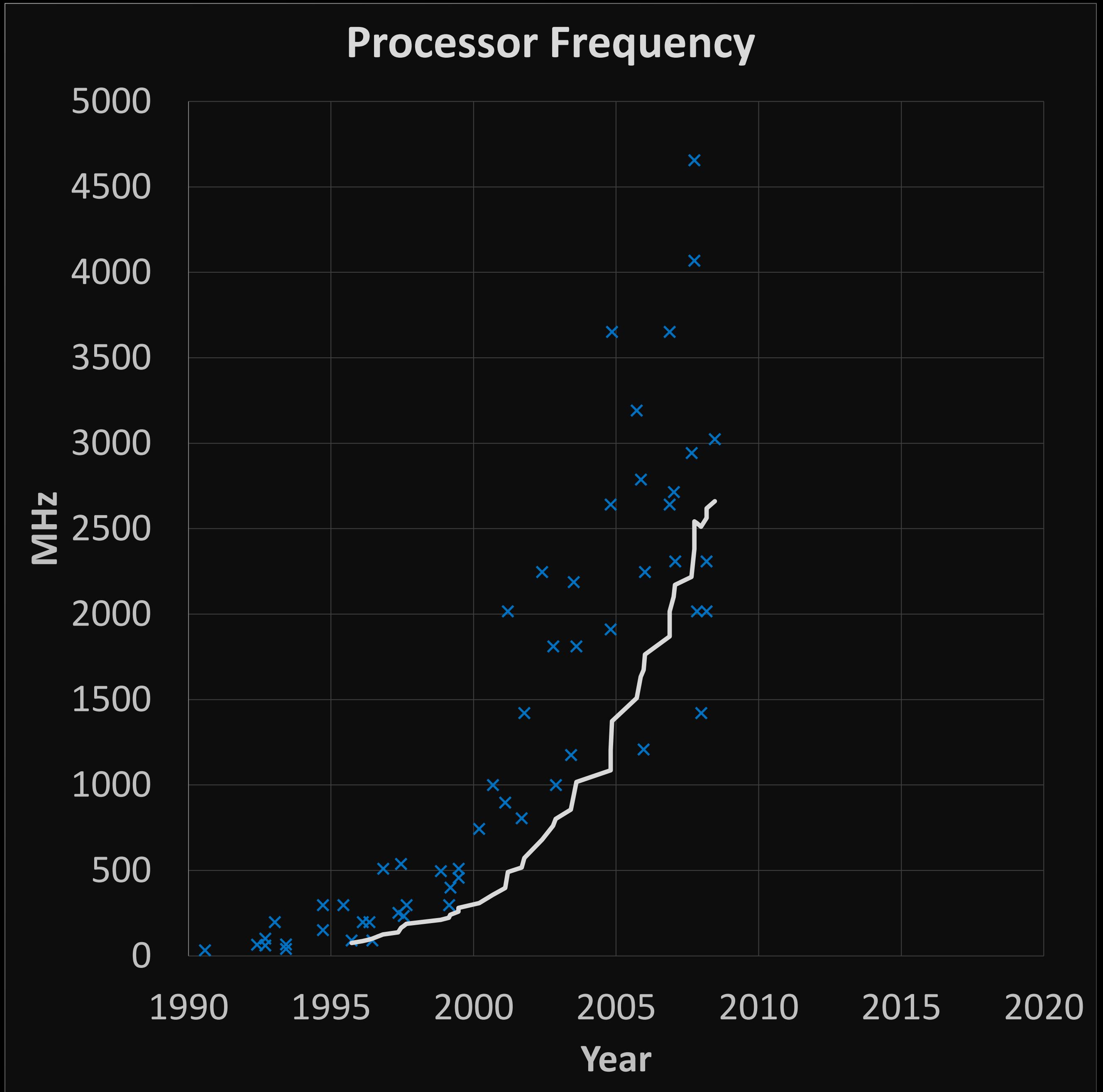
First Phase of Moore's Law

Dennard Scaling (until ~2007)

Doubling through smaller line size

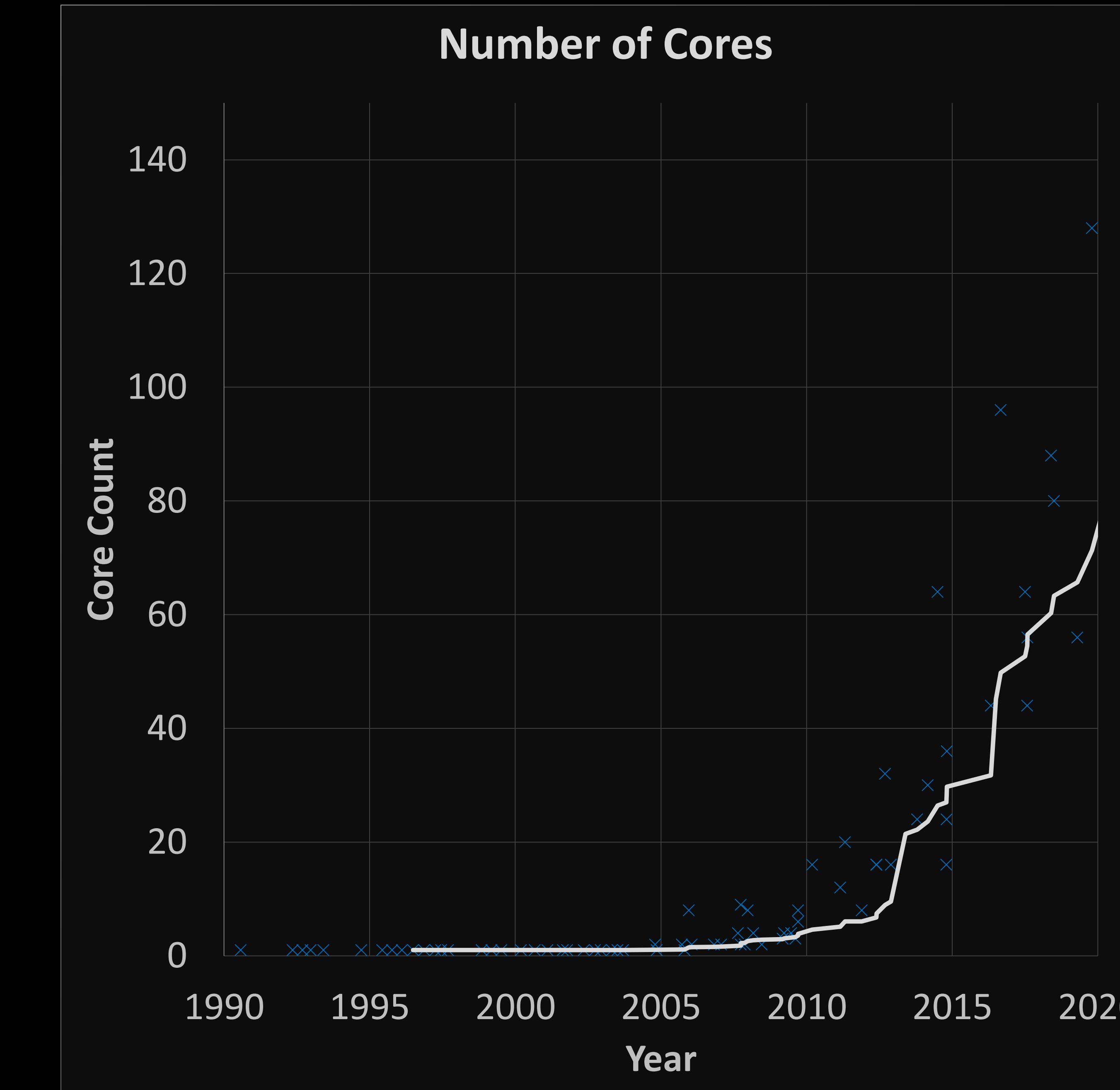
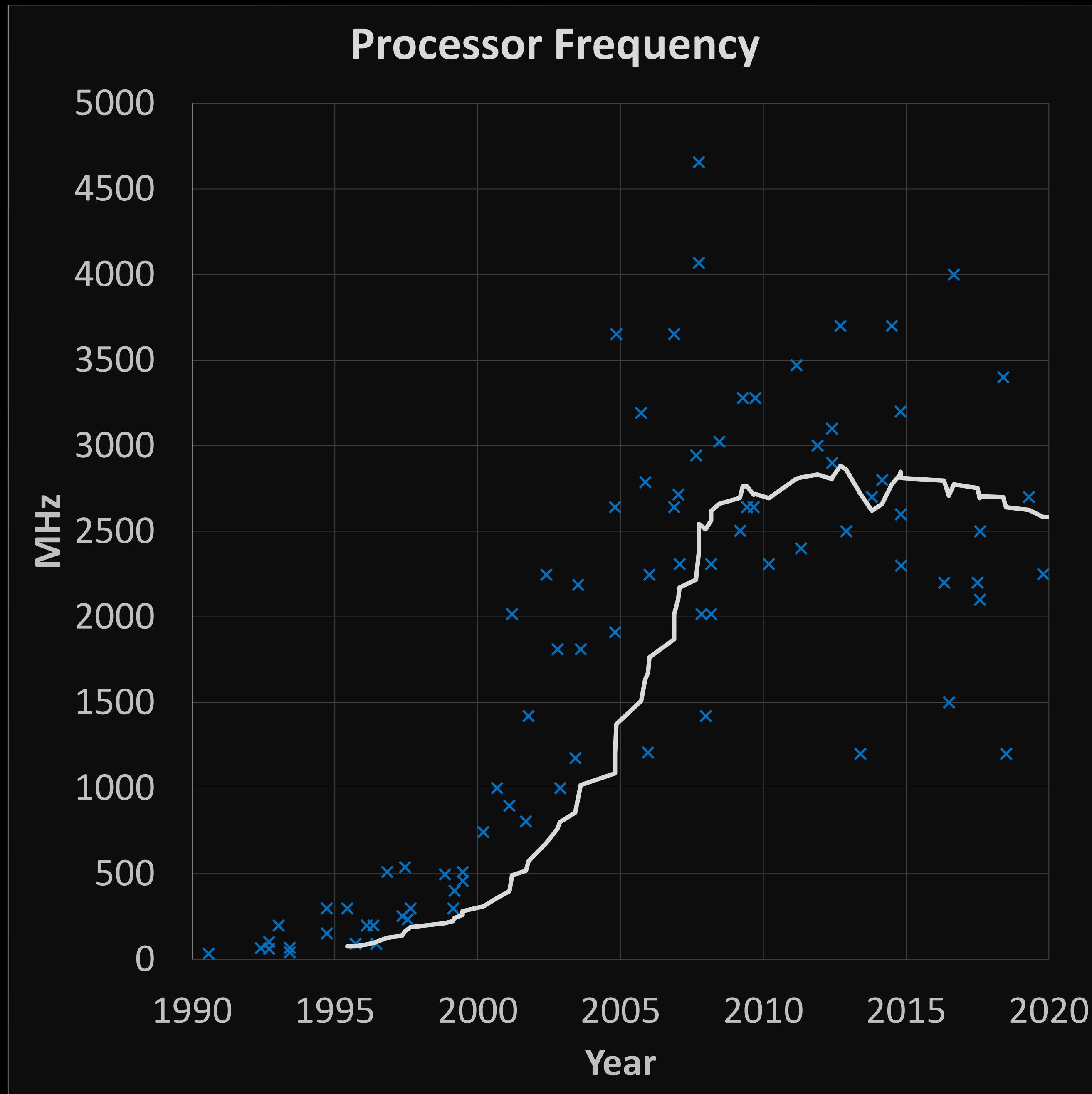


The End of Dennard Scaling



Source: Karl Rupp "40 years of microprocessor trend data"
<https://github.com/karlrupp/microprocessor-trend-data>

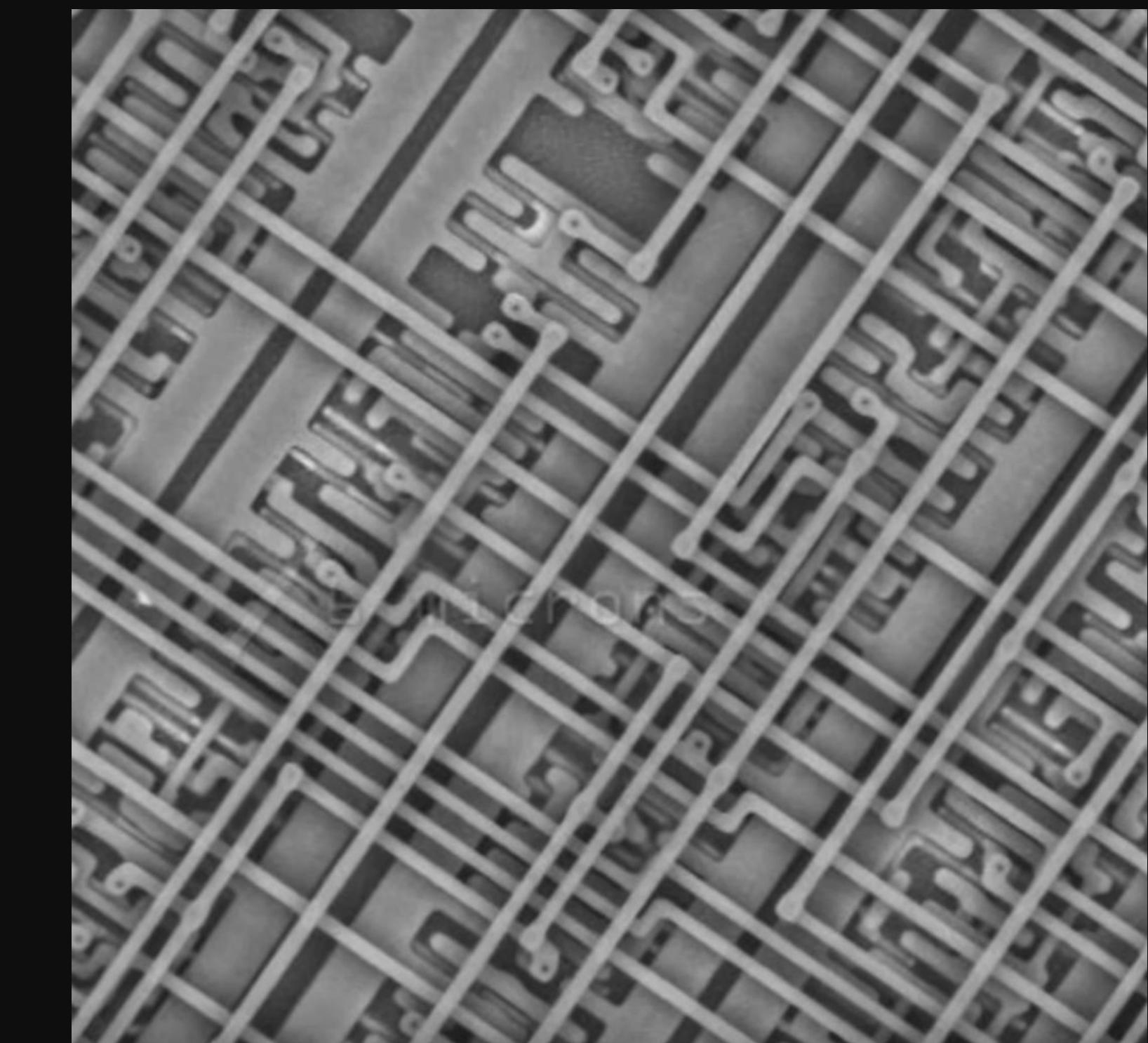
The End of Dennard Scaling



Source: Karl Rupp "40 years of microprocessor trend data"
<https://github.com/karlrupp/microprocessor-trend-data>

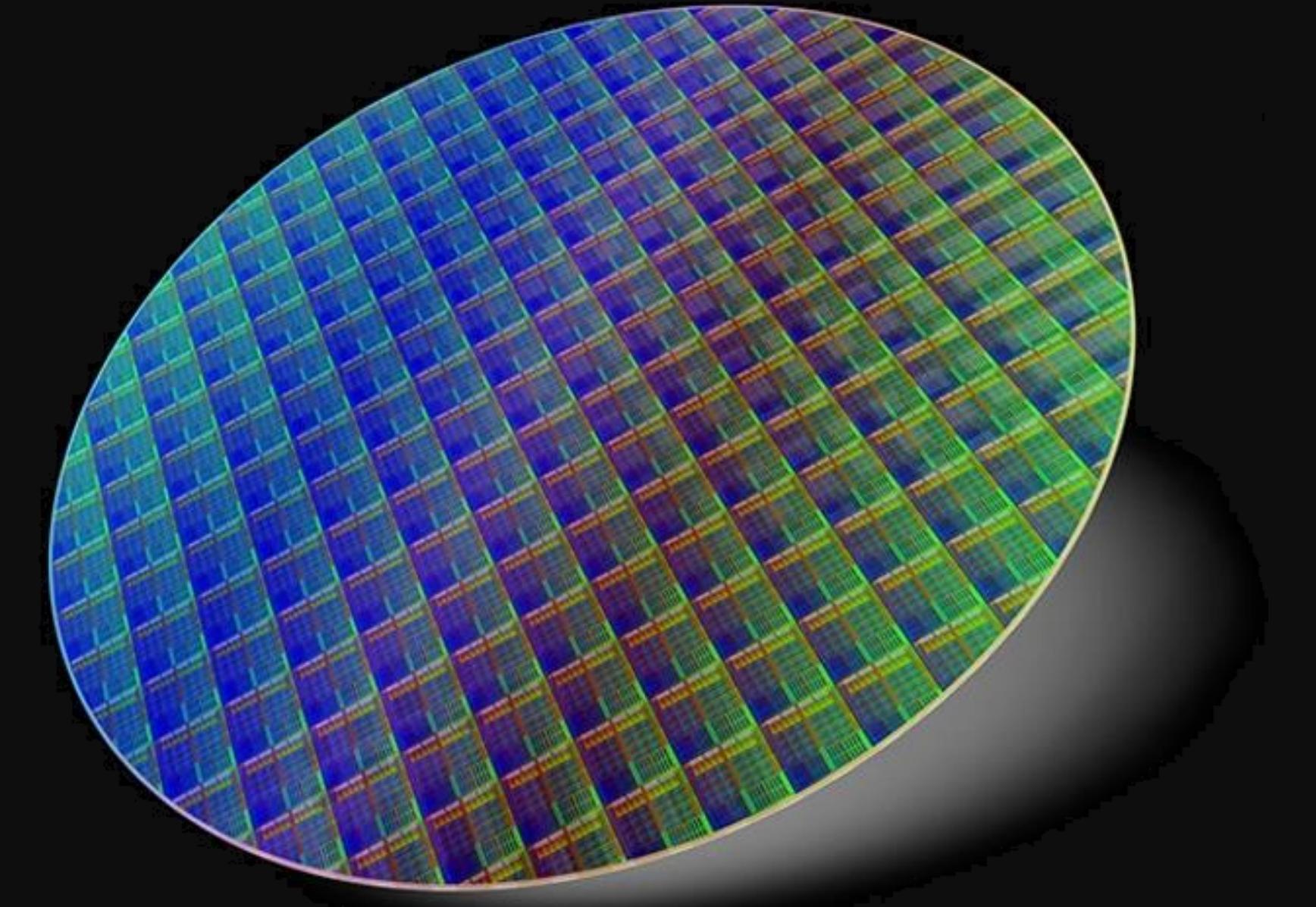
Second Phase of Moore's Law

Dennard Scaling (until ~2007)

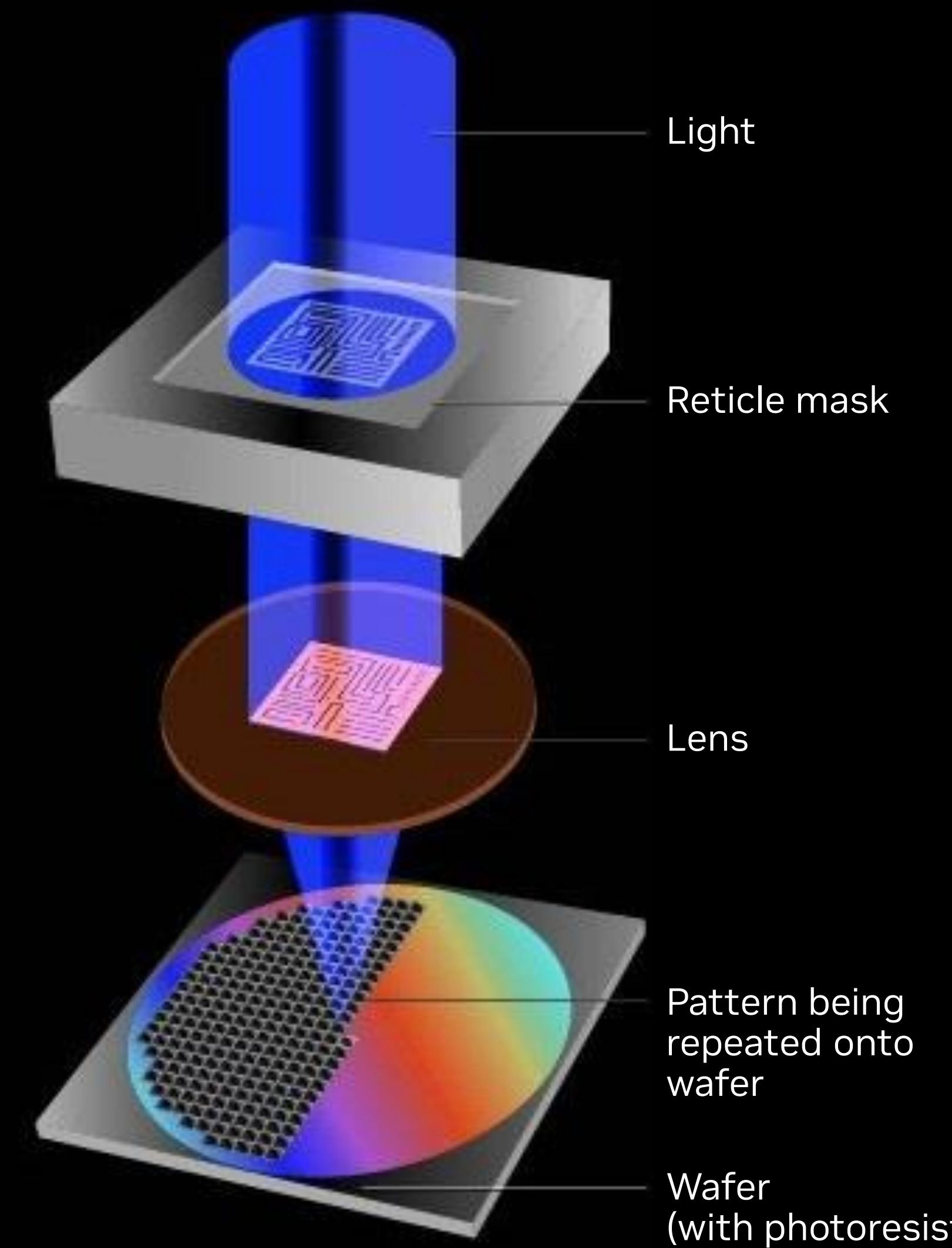
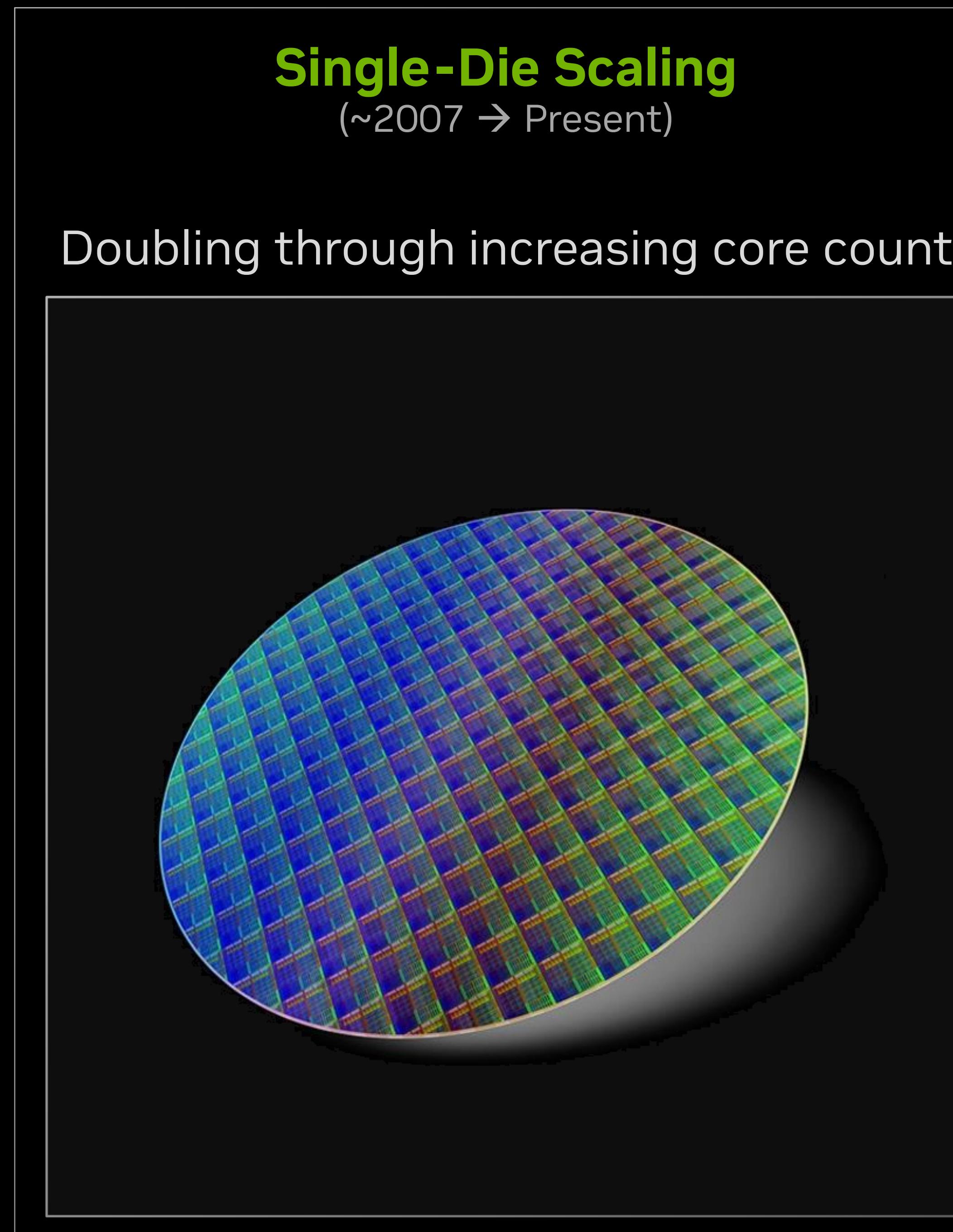


Single-Die Scaling (~2007 → Present)

Doubling through increasing core count



The End of Single-Die Scaling



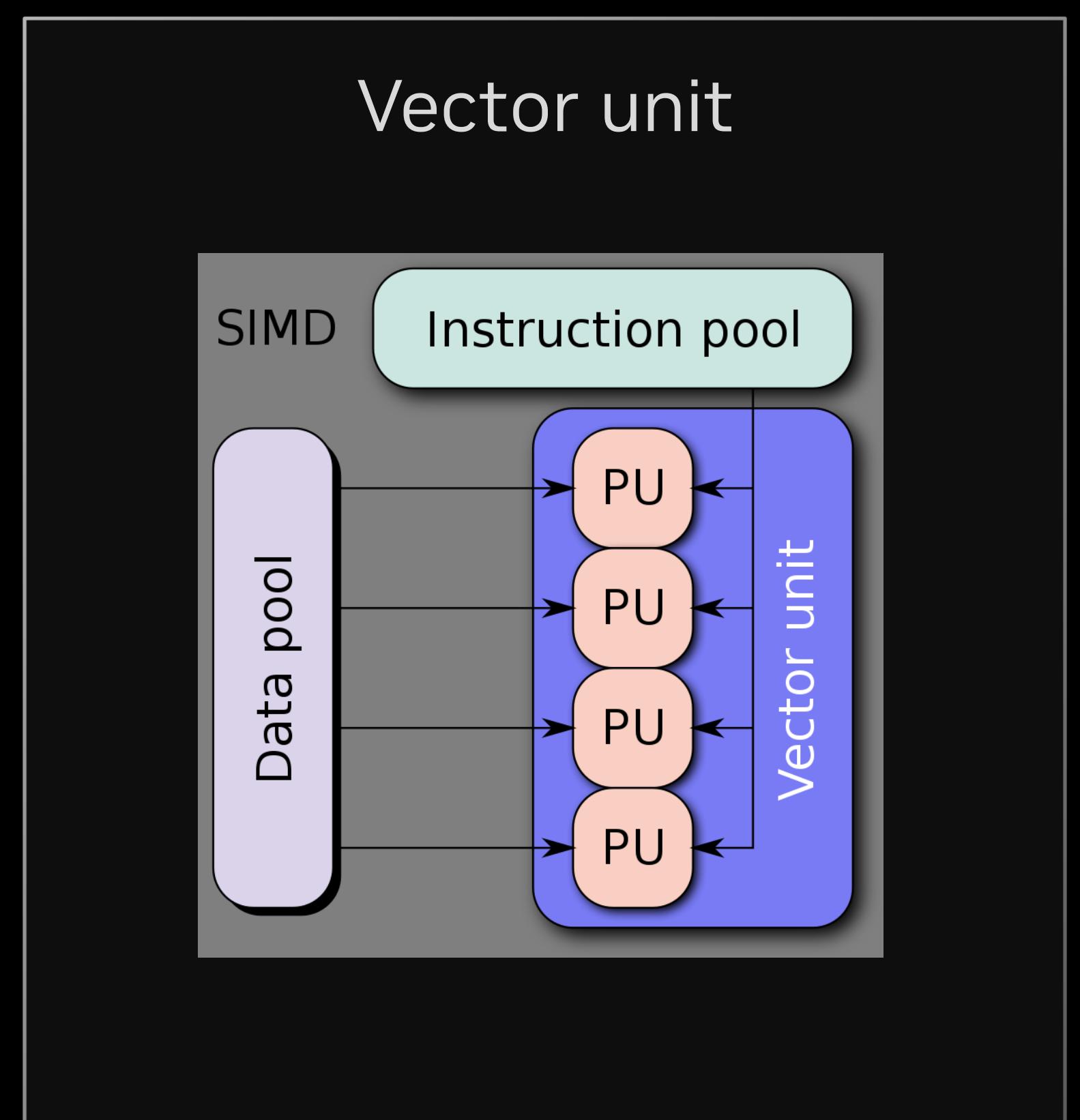
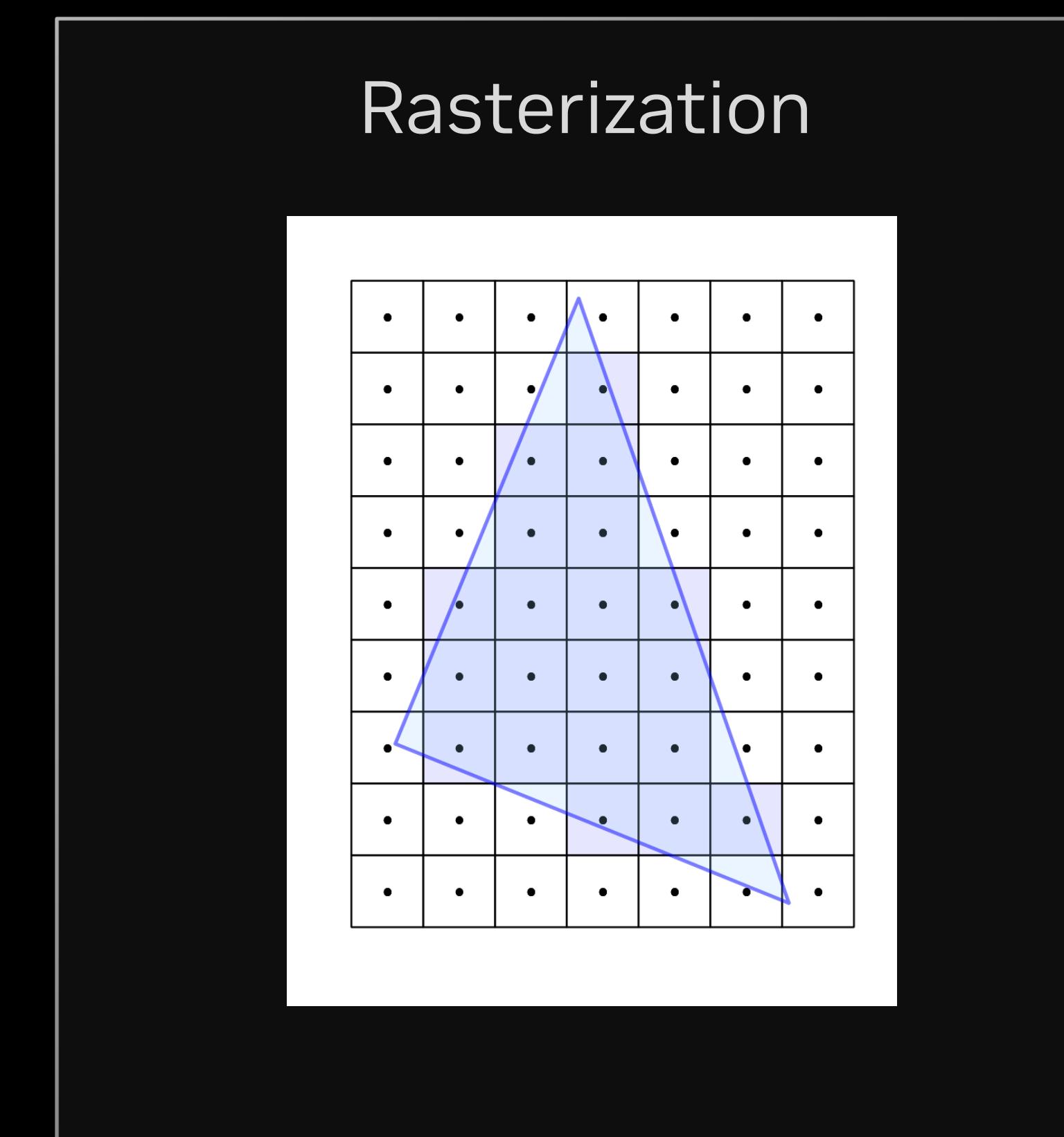
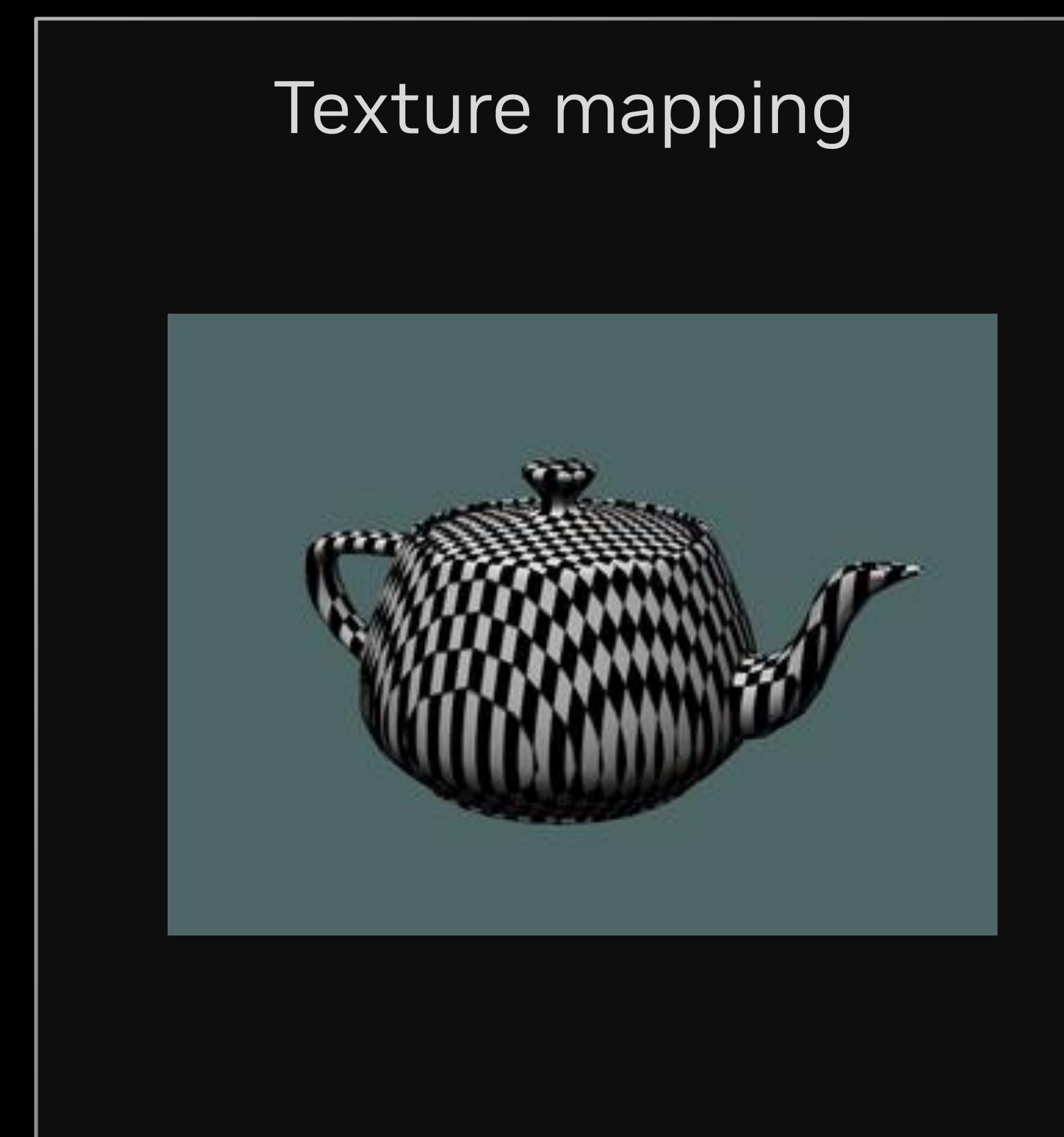
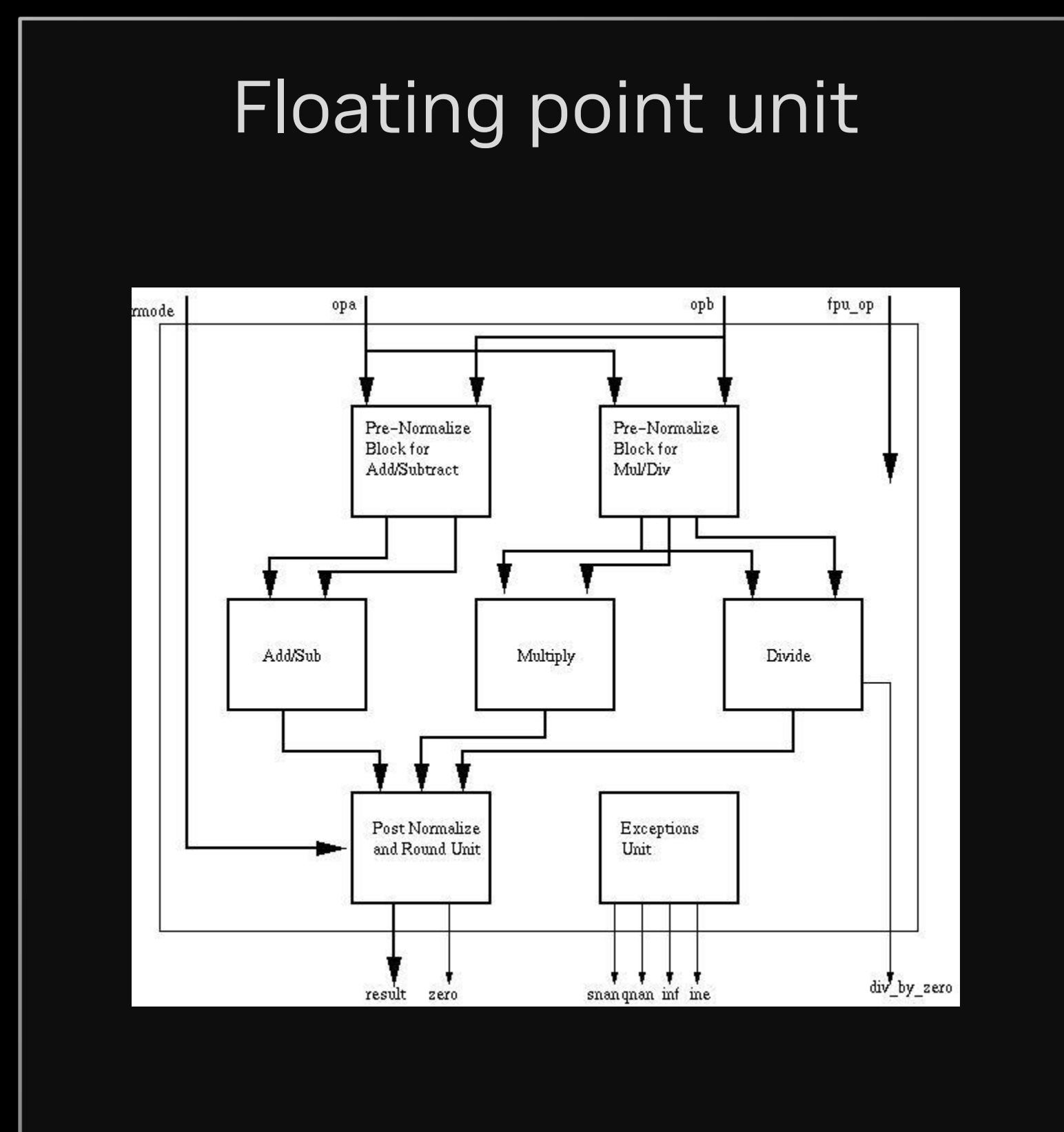
Overcoming the Reticle Limit: Heterogeneous Hardware

Specialization of hardware components

heterogeneous

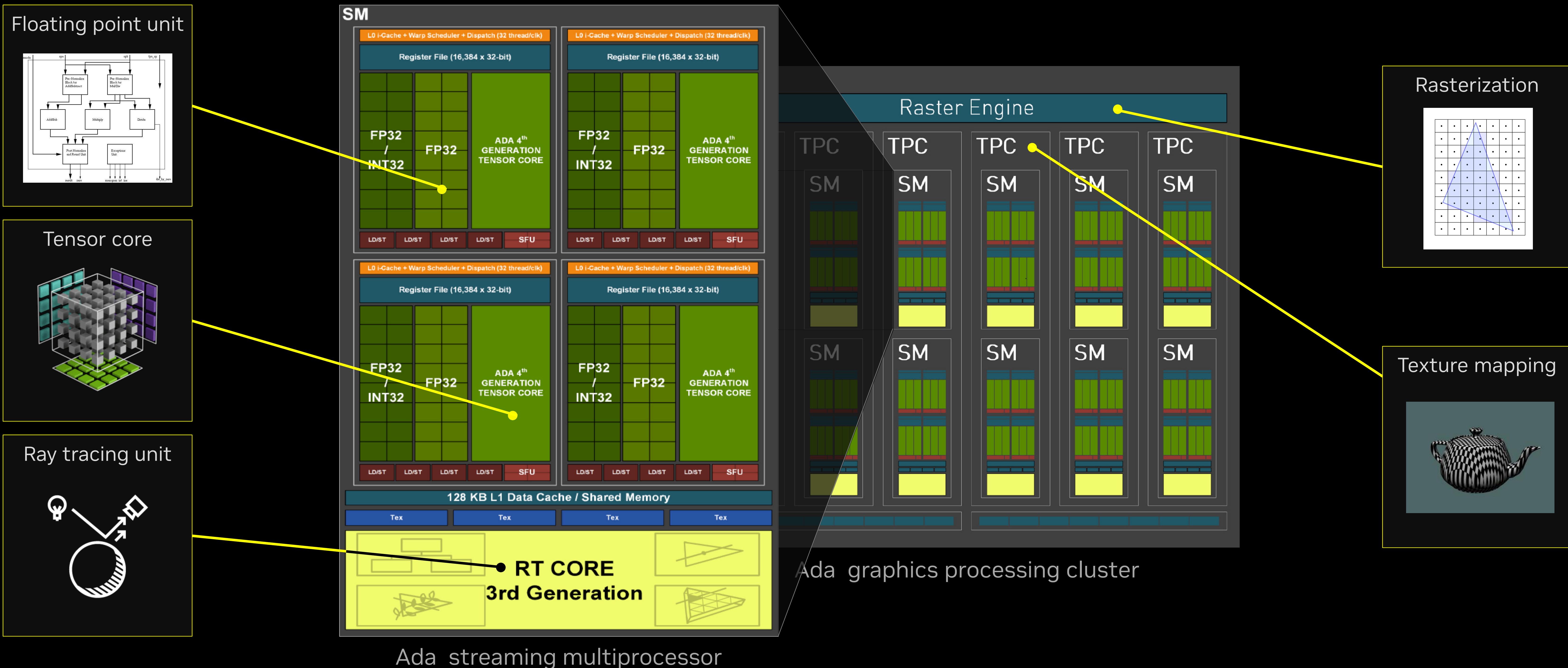
adjective US: /ˌhetəˈrəʊdʒiːniəs/ UK: /ˌhetərəˈdʒiːniəs/

consisting of parts or things that are very different from each other



Overcoming the Reticle Limit: Heterogeneous Hardware

NVIDIA Ada microarchitecture



Heterogeneous Software

Software is also made up of heterogeneous components written in many different languages

heterogeneous

adjective US: /,hetərə'rous/ UK: /,het.ər.ə'rous/

consisting of parts or things that are very different from each other

Example: Heterogeneous Programming in Python

```
import numpy as np

# Array of 1M random numbers
data = np.random.randint(1, 100, (1024 * 1024))

total = np.sum(data)      # Reduction across array
max = npamax(data)       # Max in array
average = np.mean(data)   # Average across array
```

← User application is written in Python for productivity

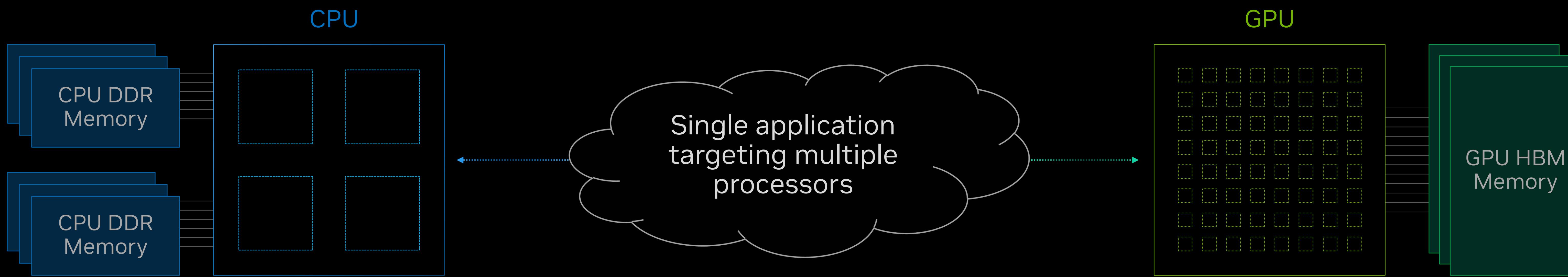
NumPy Library

NumPy library implements fast computation in C & C++

Applications often combine components & libraries
written in many different languages

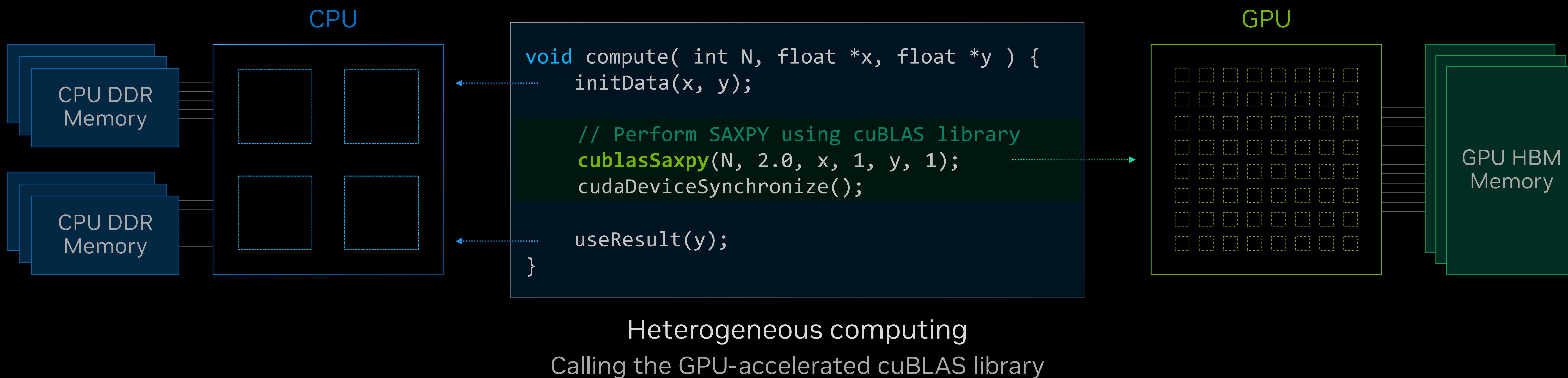
Heterogeneous Computing

Combining processors of different types, each specializing in different types of execution



Heterogeneous Computing

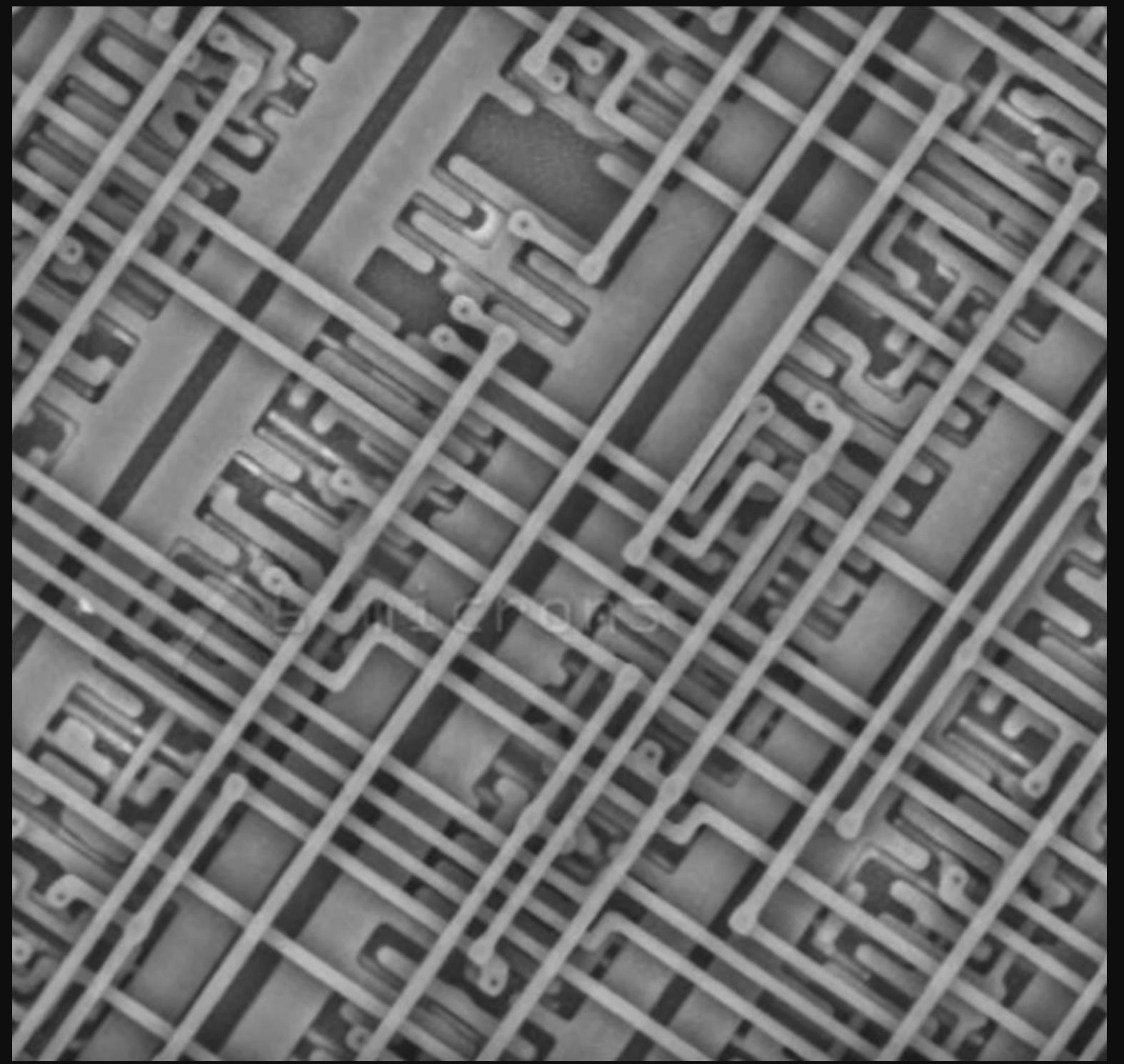
Combining processors of different types, each specializing in different types of execution



Third Phase of Moore's Law

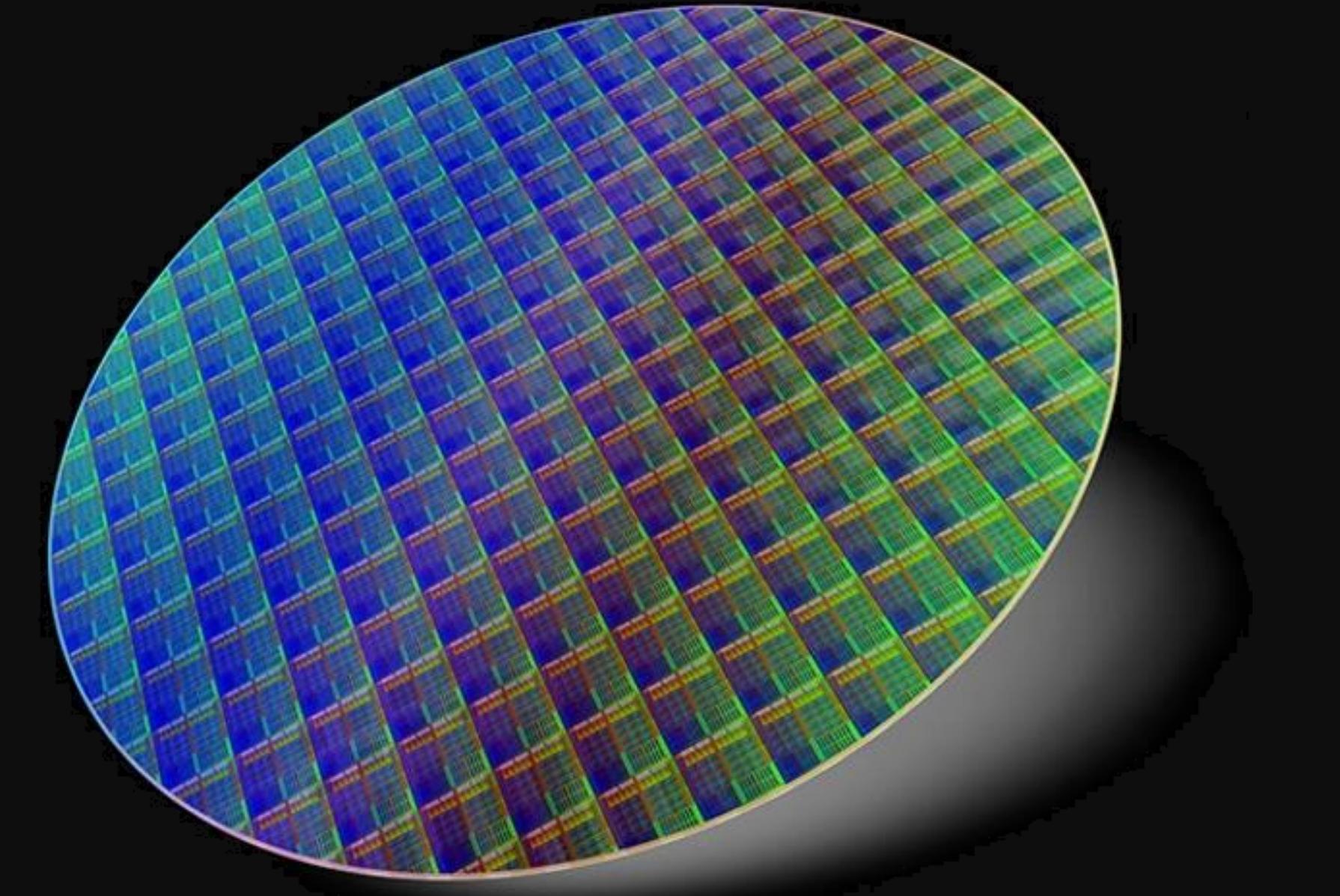
Dennard Scaling (until ~2007)

Doubling through smaller line size



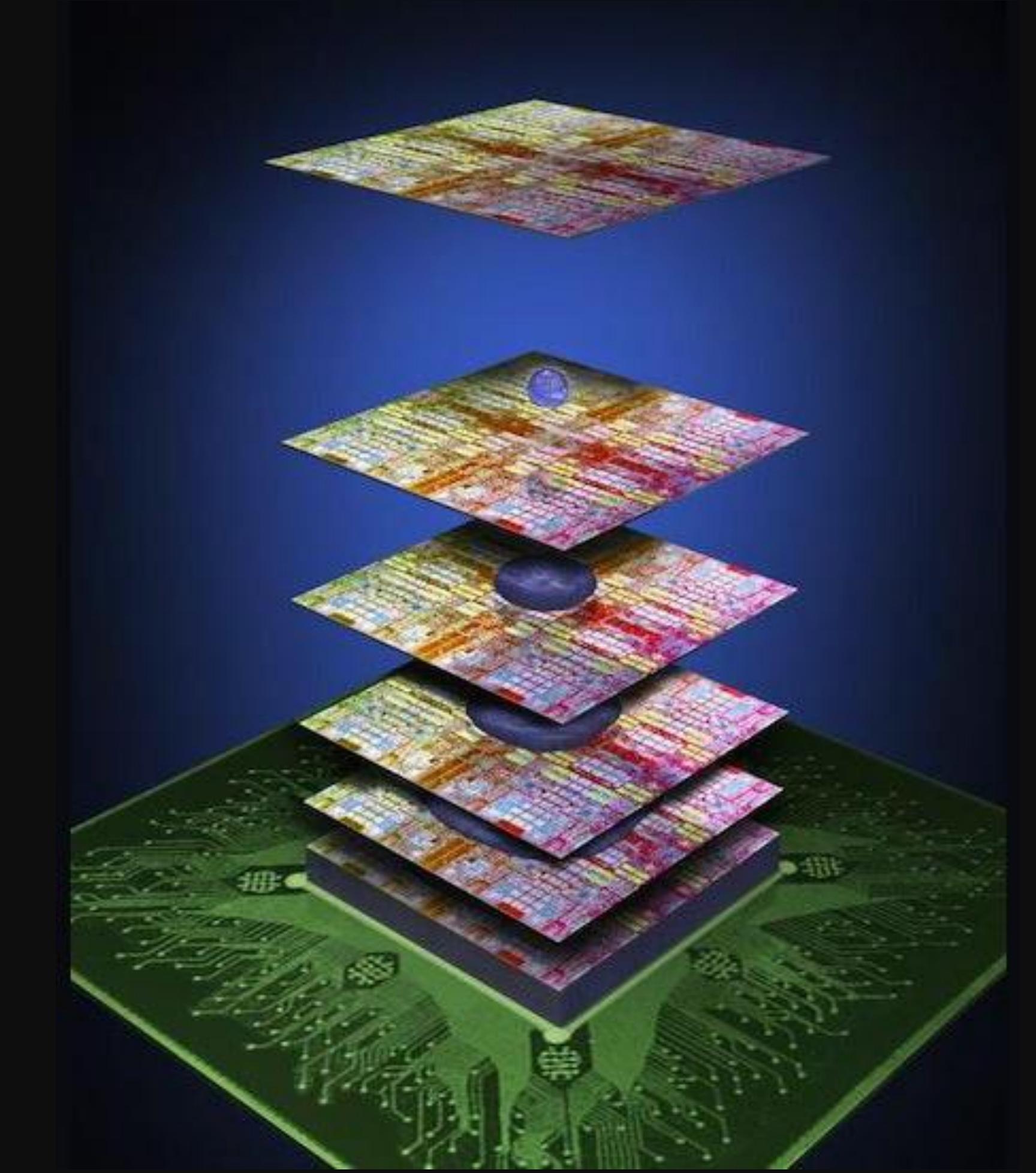
Single-Die Scaling (~2007 → Present)

Doubling through increasing core count

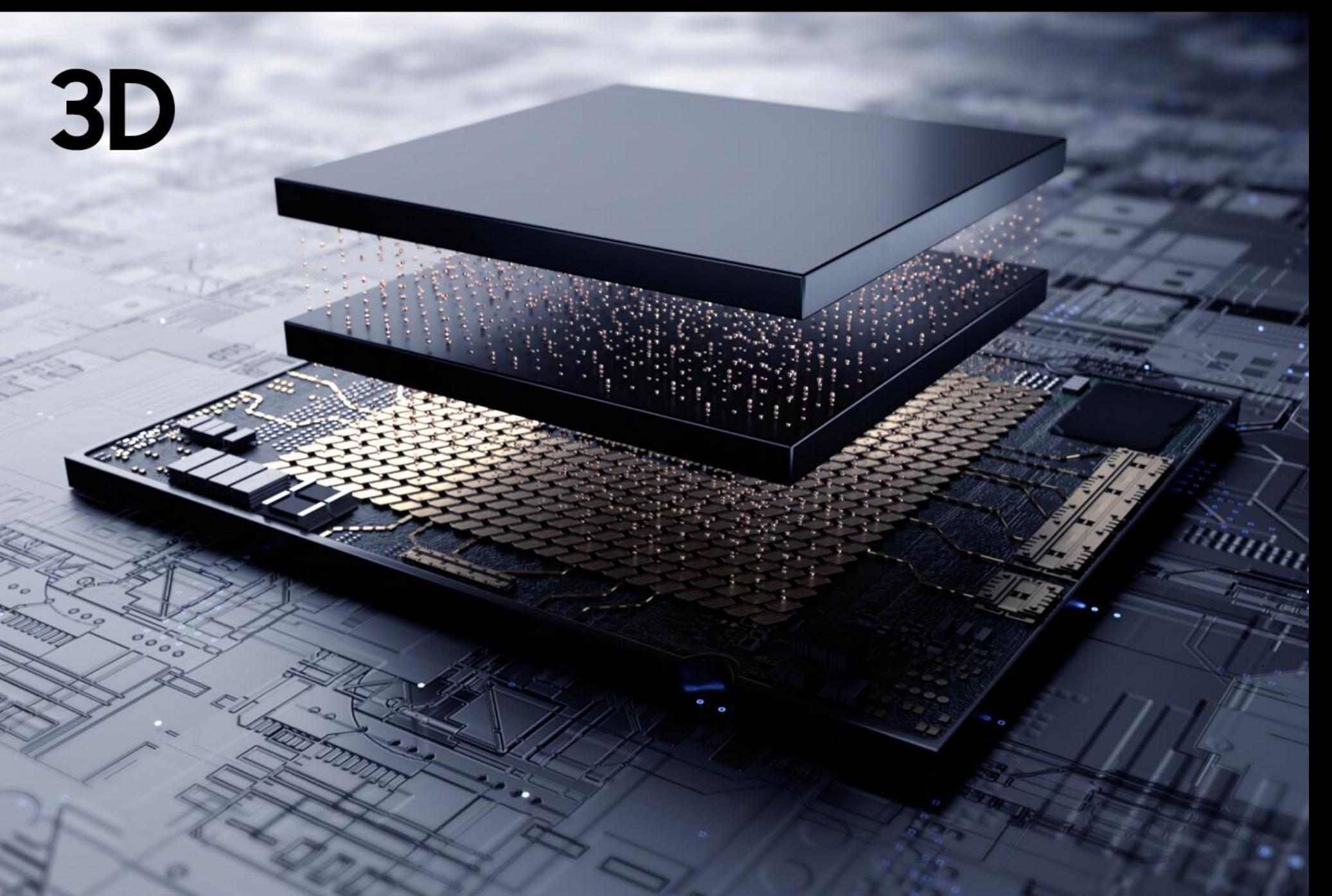
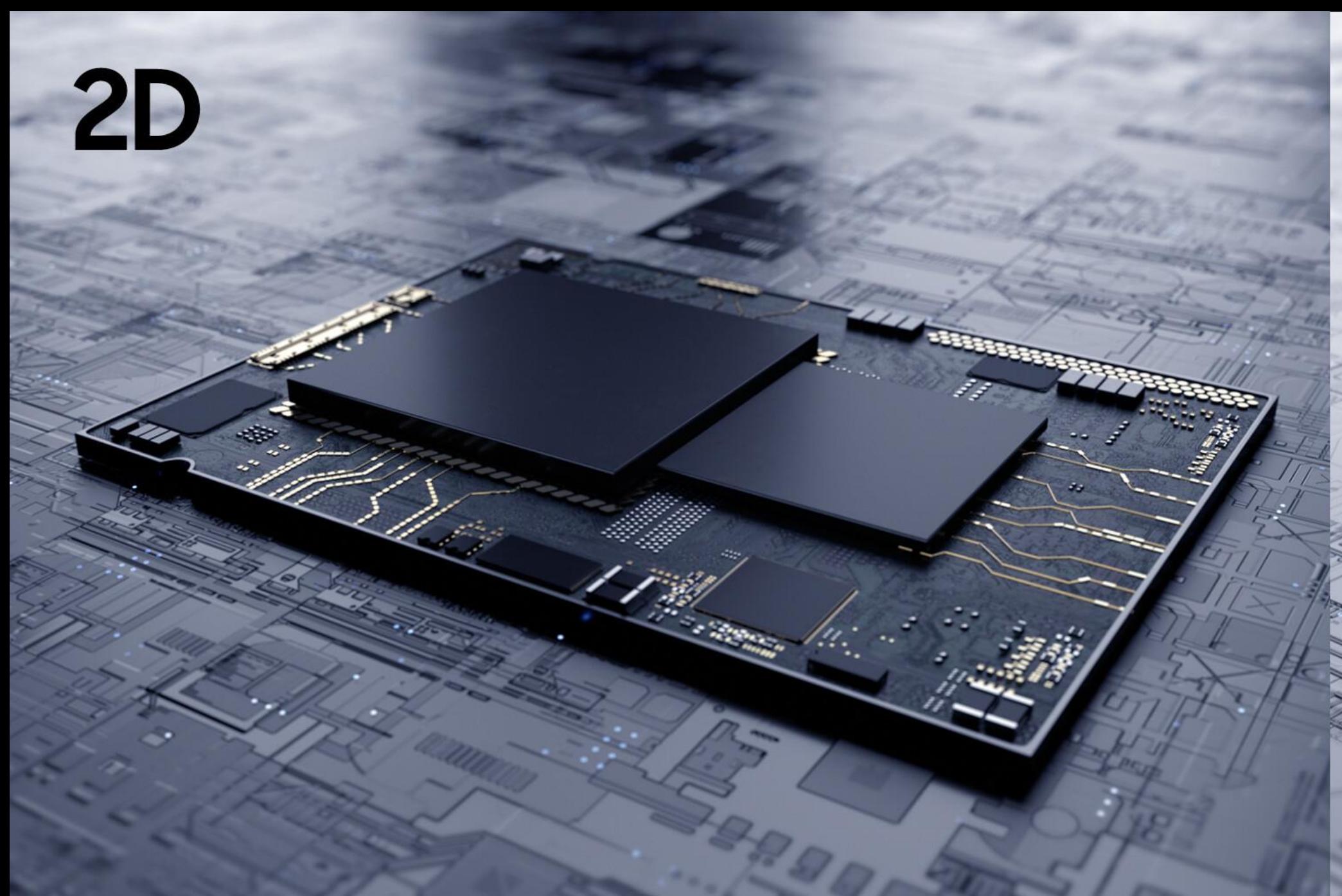


Multi-Die Scaling (Present → Future)

Doubling through combining dies



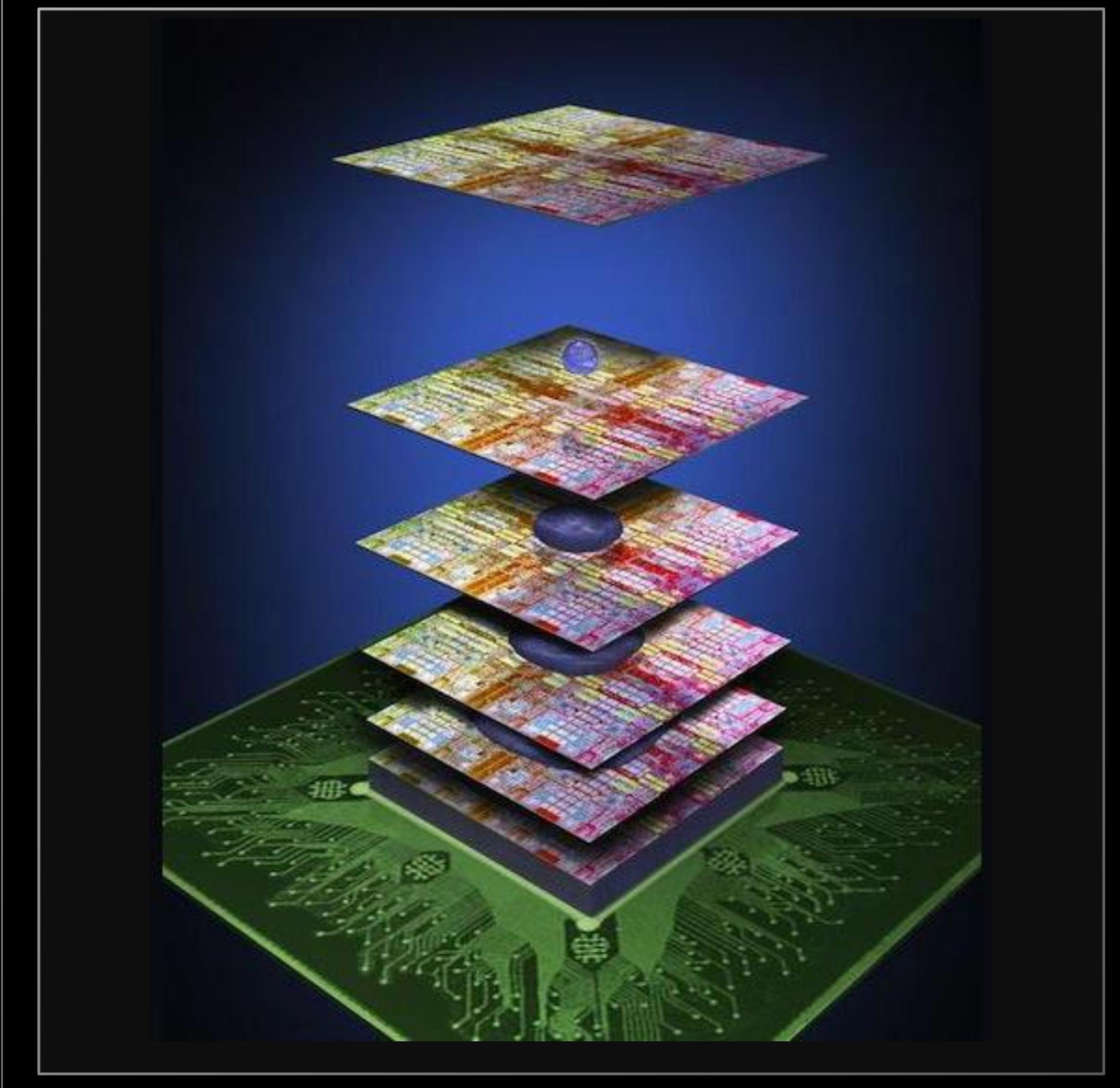
Third Phase of Moore's Law



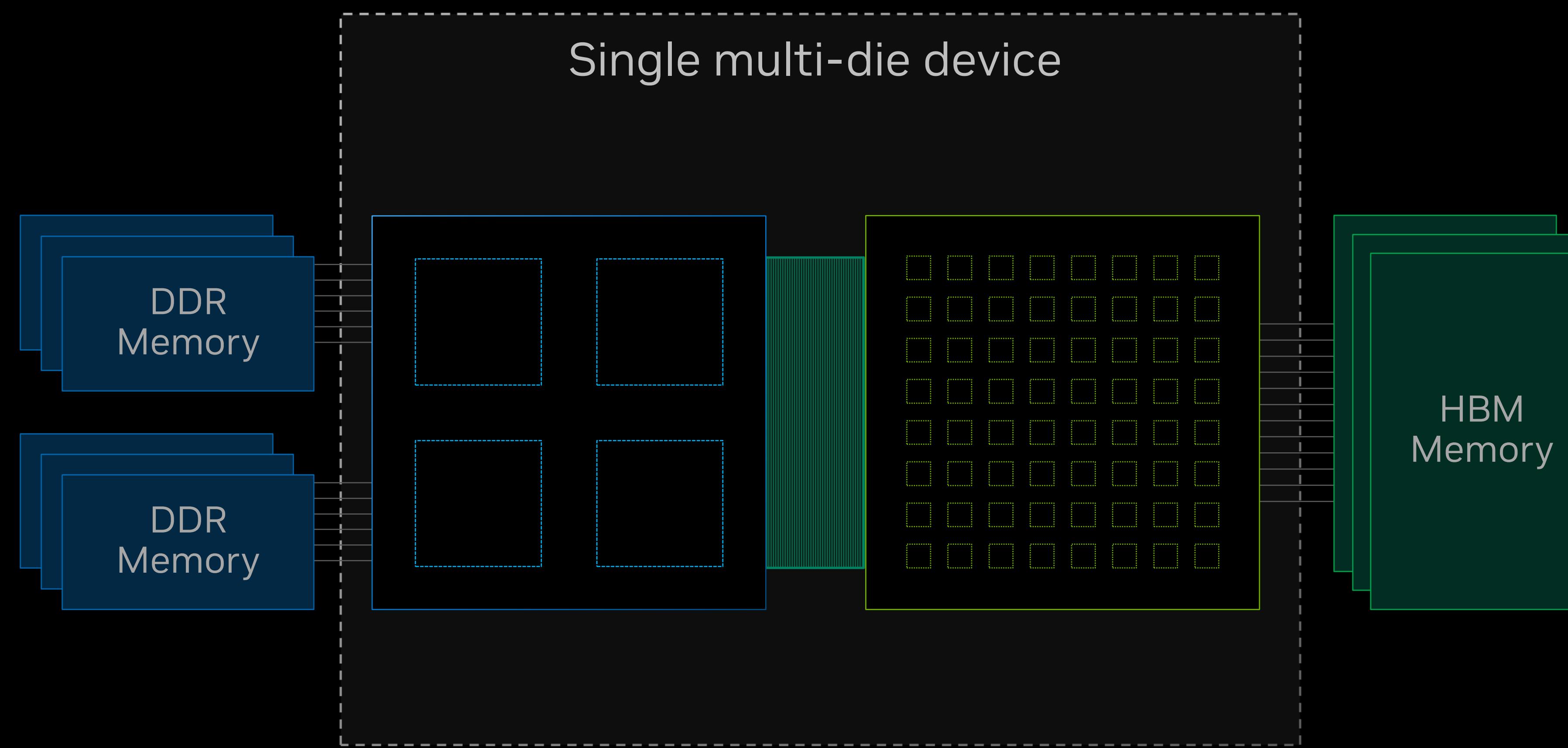
Multi-Die Scaling

(Present → Future)

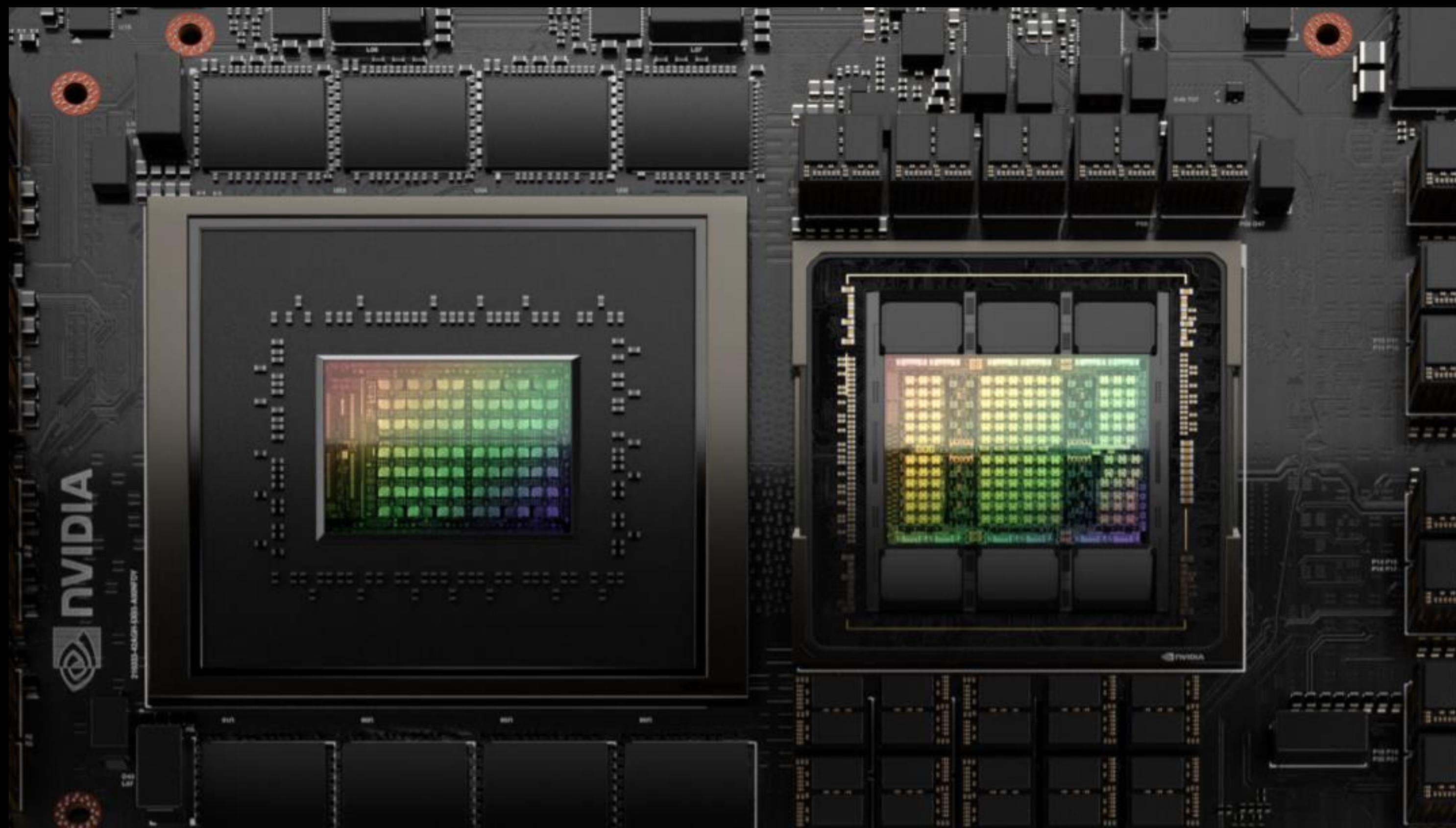
Doubling through combining dies



Multi-Die Heterogeneous Devices

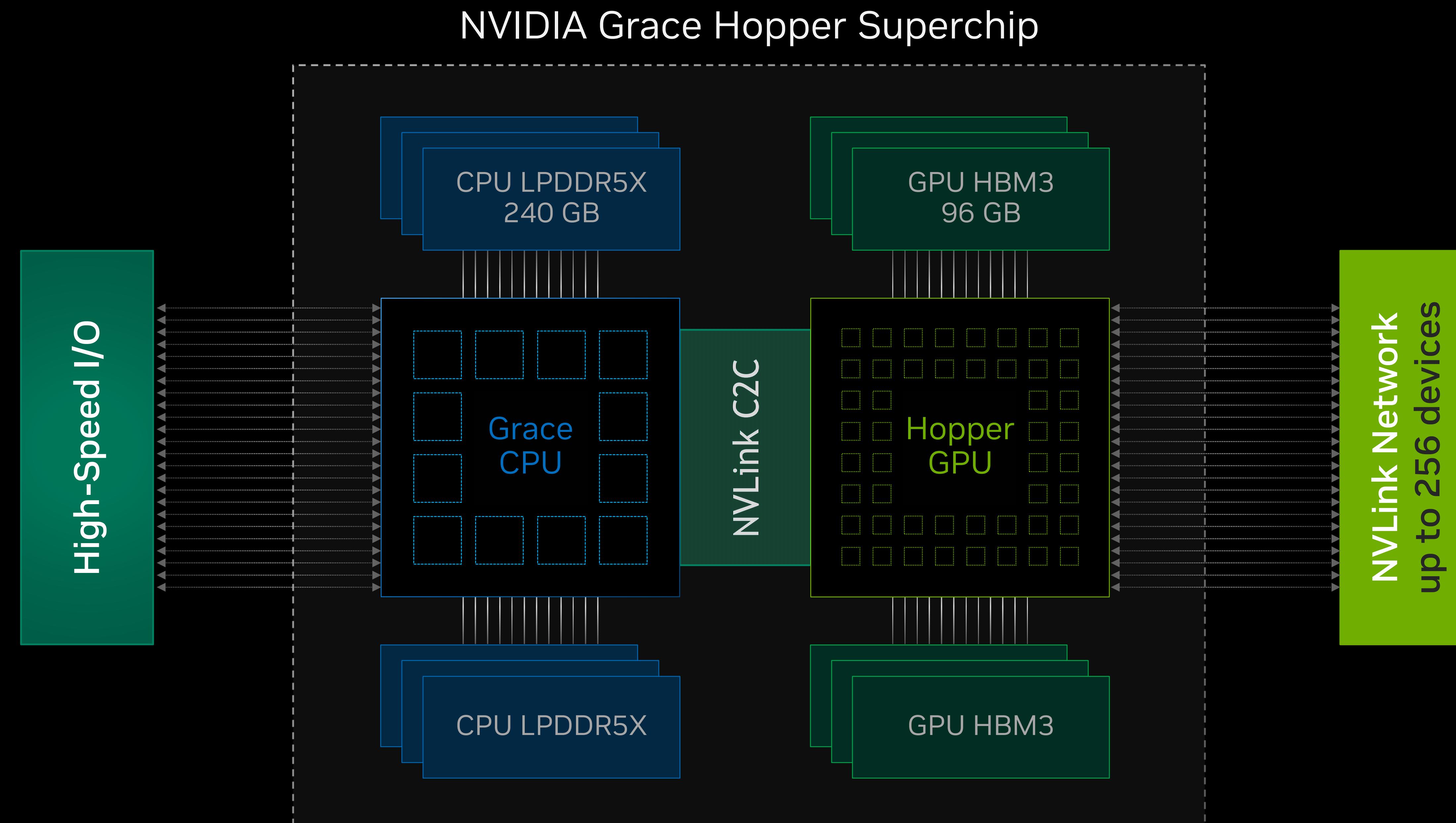


Grace/Hopper Superchip

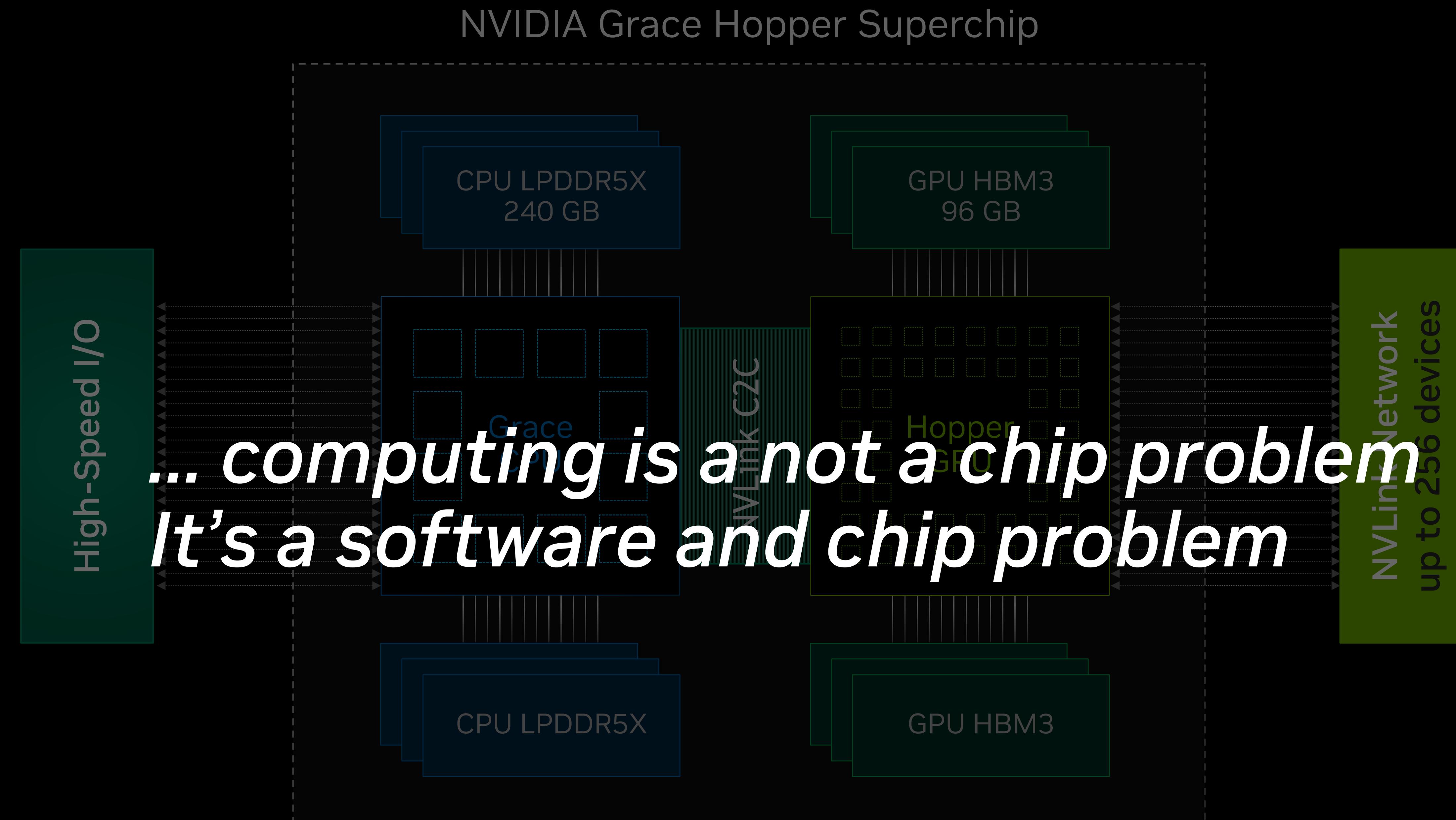


[NVIDIA Grace Hopper Superchip Architecture \[Paper\]](#)

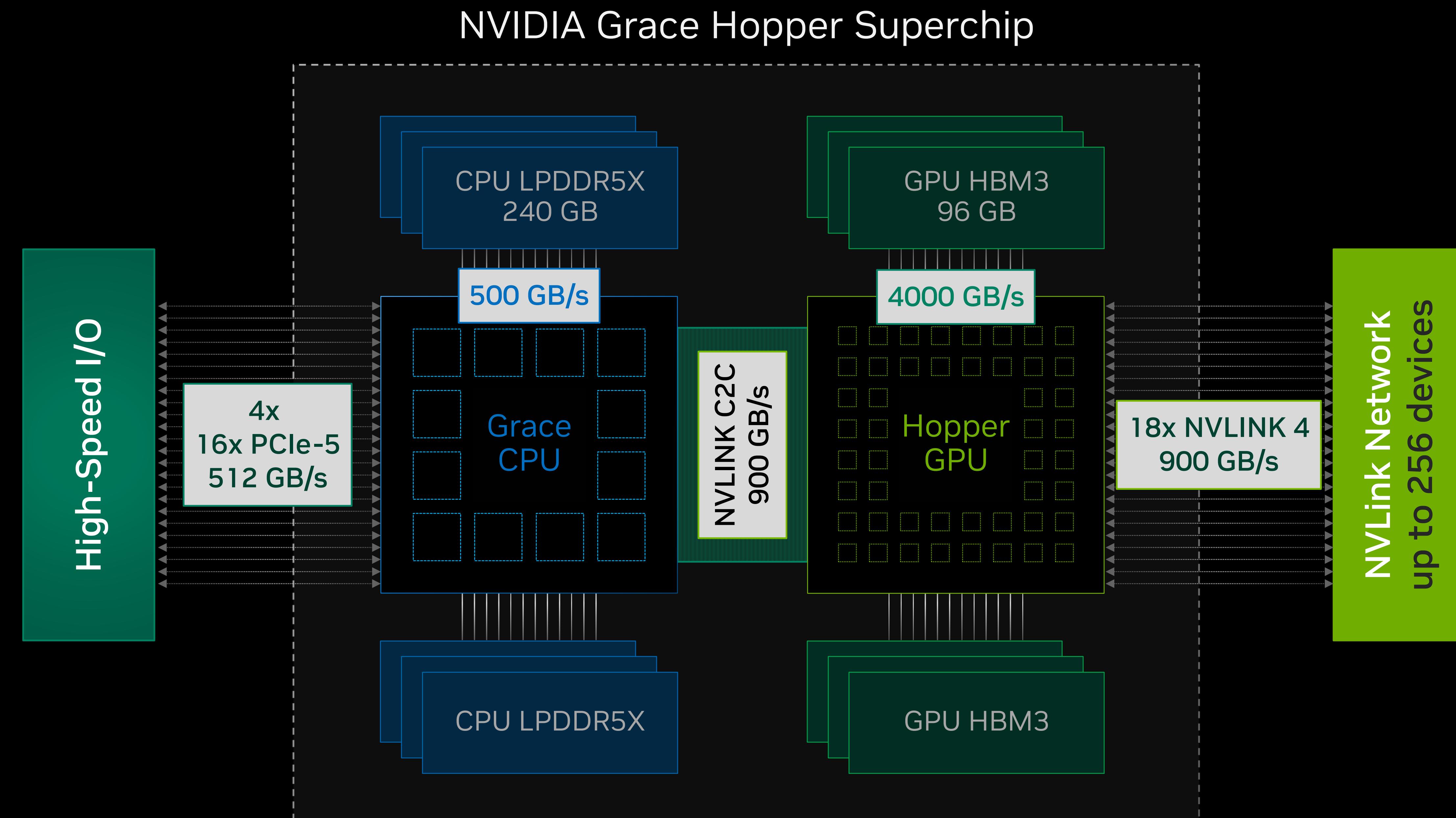
Grace/Hopper Superchip



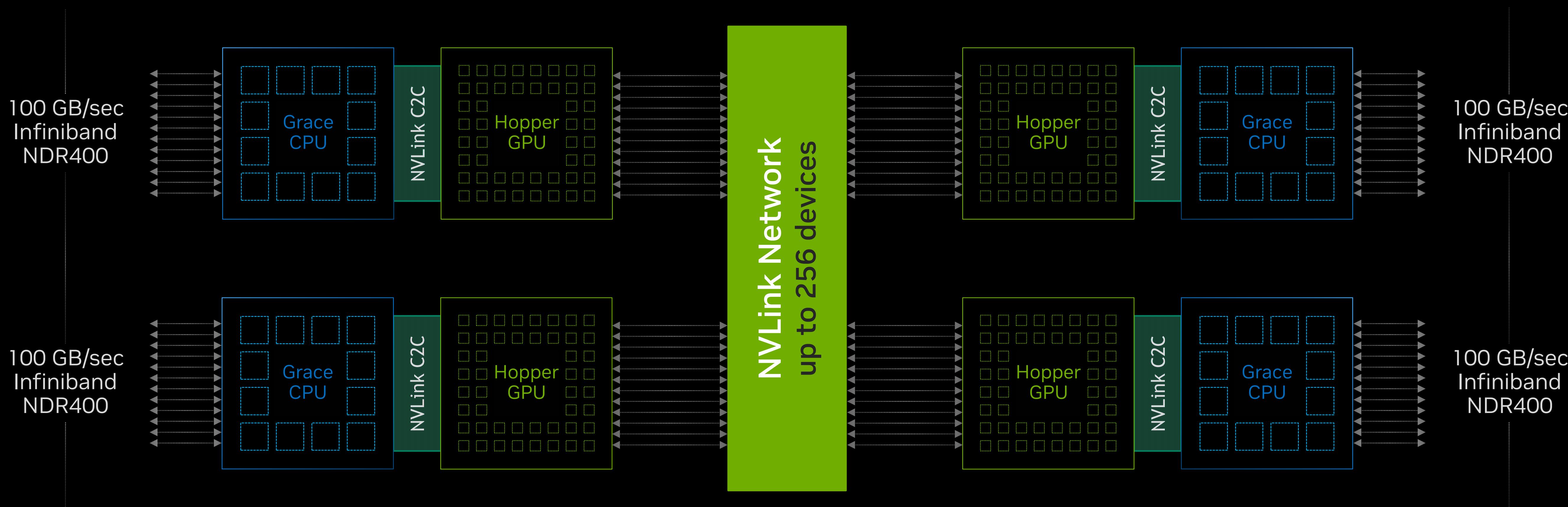
Grace/Hopper Superchip



Grace/Hopper Superchip

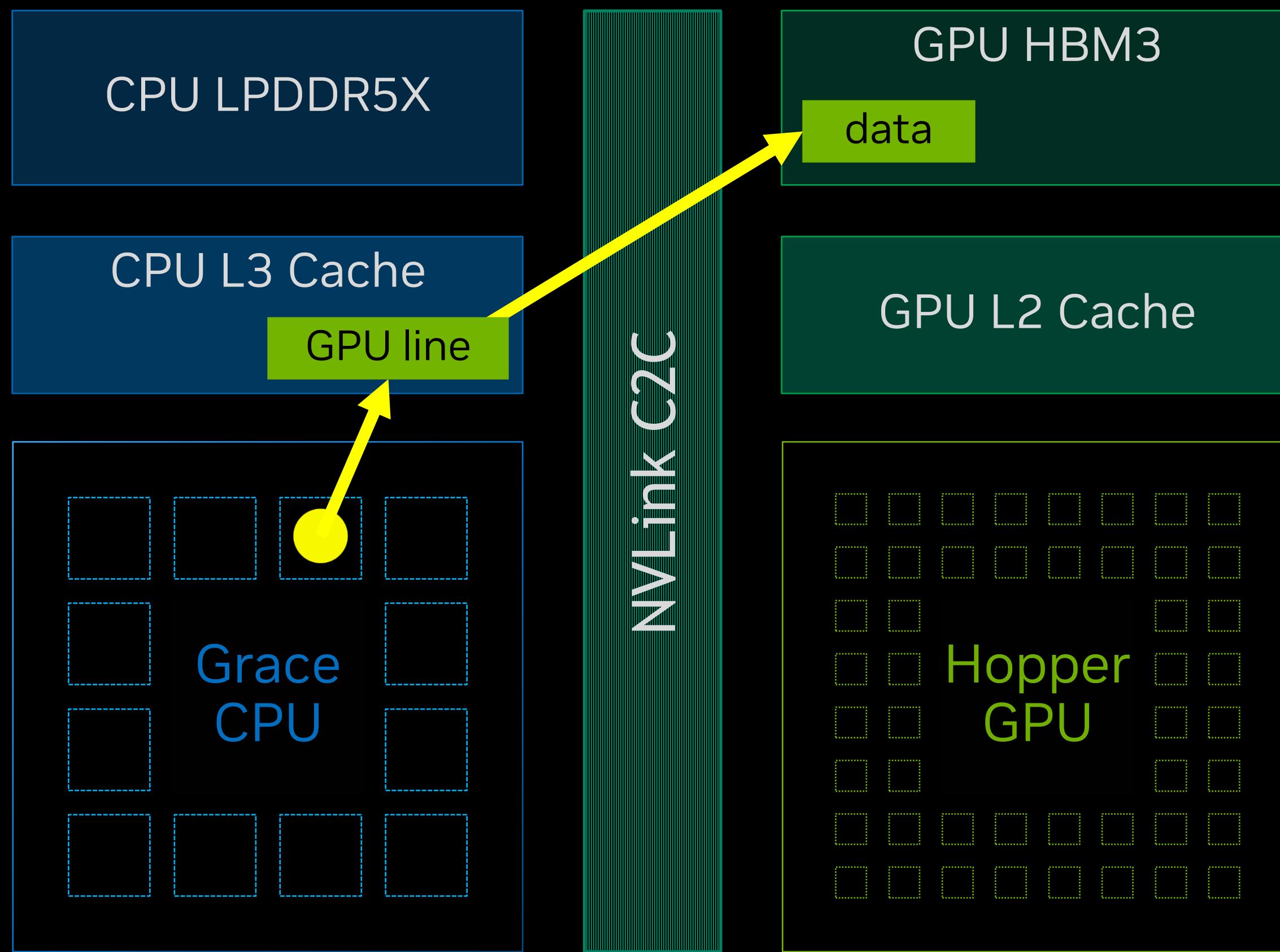


NVLink Connects Up To 256 Superchips



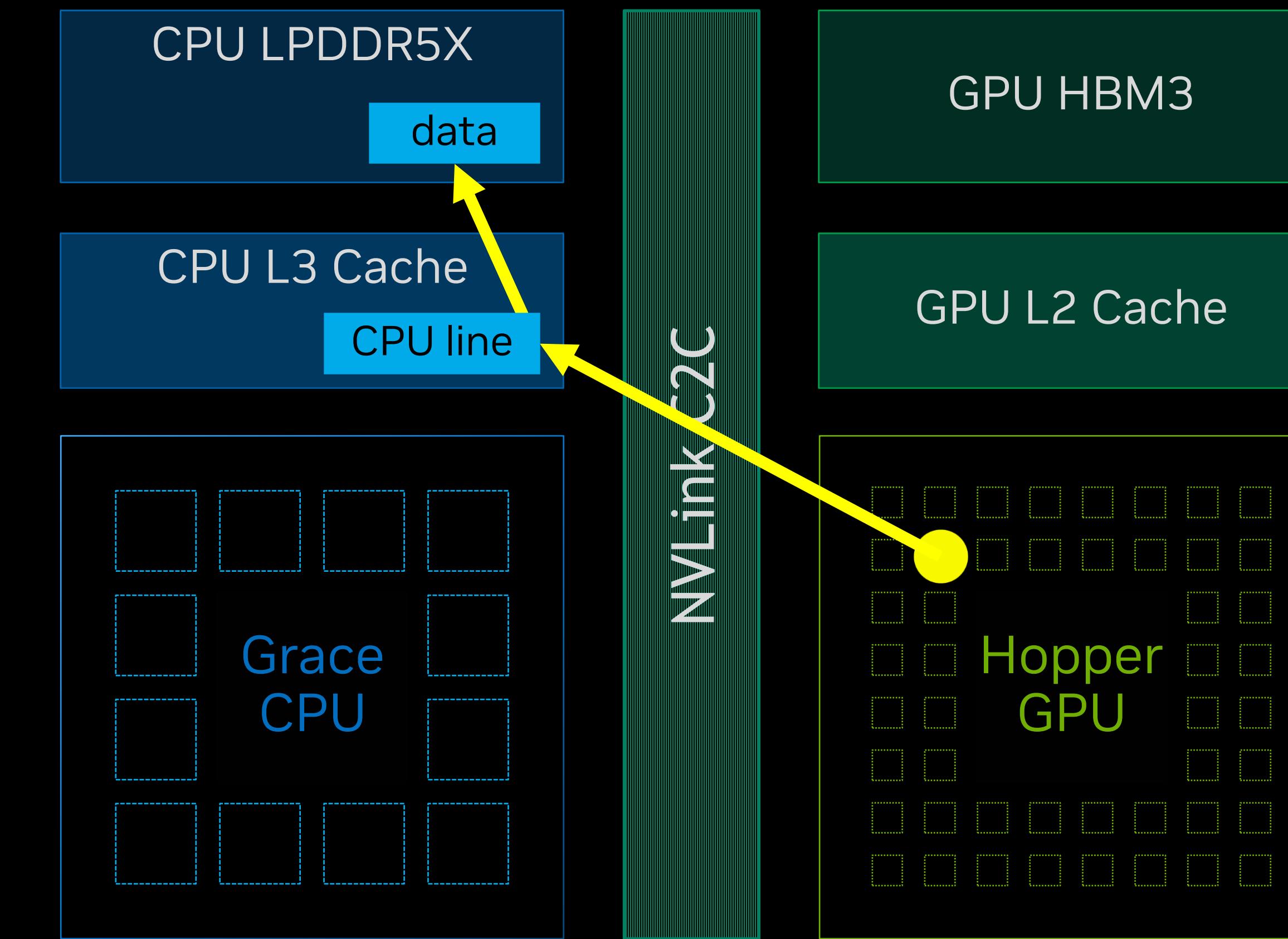
Global Access to All Data

Cache-coherent access via NVLink C2C from either processor to either physical memory



Grace directly reading Hopper's memory

CPU fetches GPU data into CPU L3 cache
Cache remains **coherent** with GPU memory
Changes to GPU memory **evict** cache line

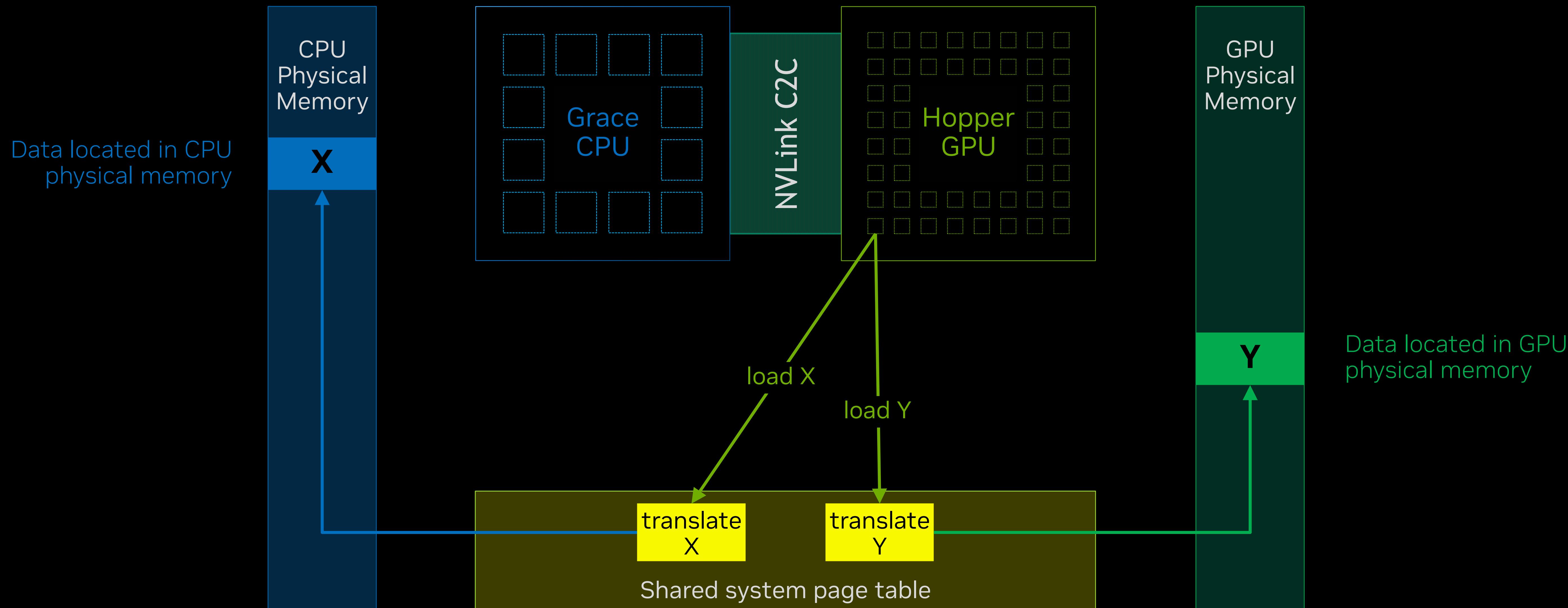


Hopper directly reading Grace's memory

GPU loads CPU data via CPU L3 cache
CPU and GPU **can both hit** on cached data
Changes to CPU memory **update** cache line

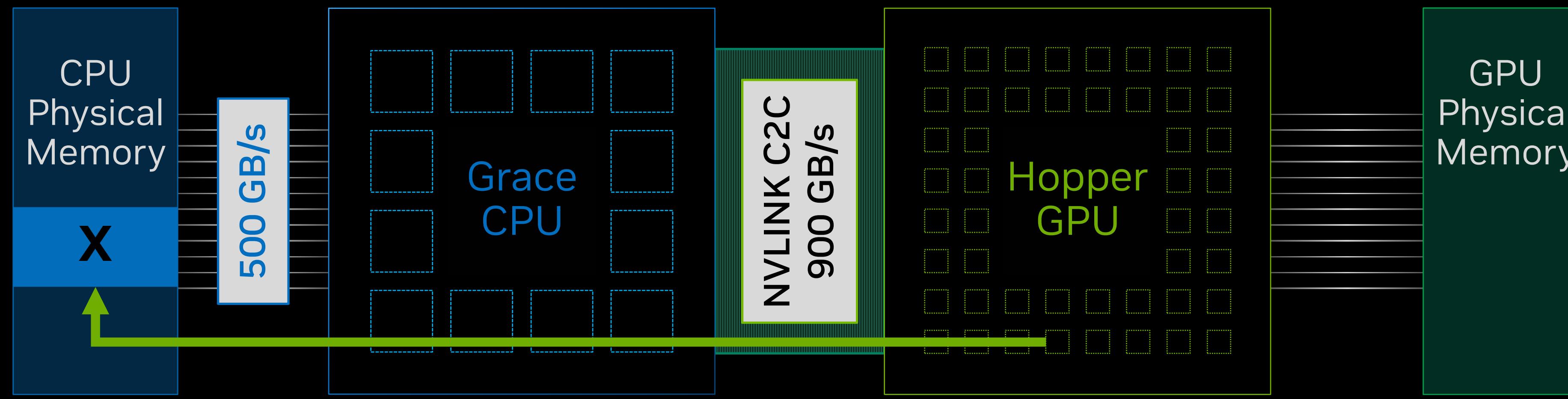
Grace/Hopper Unified Memory

Address Translation Service (ATS) allows full access to all CPU & GPU allocations

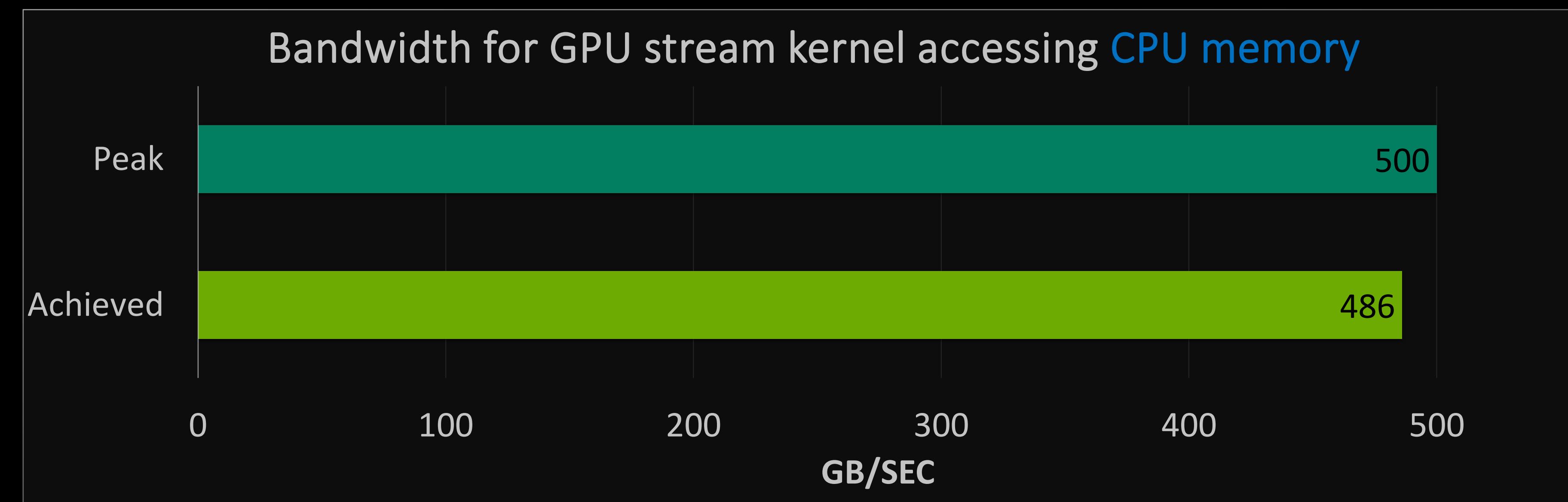


ATS creates a single page table for the whole system
NVLink C2C allows access to all physical memory **without migration**

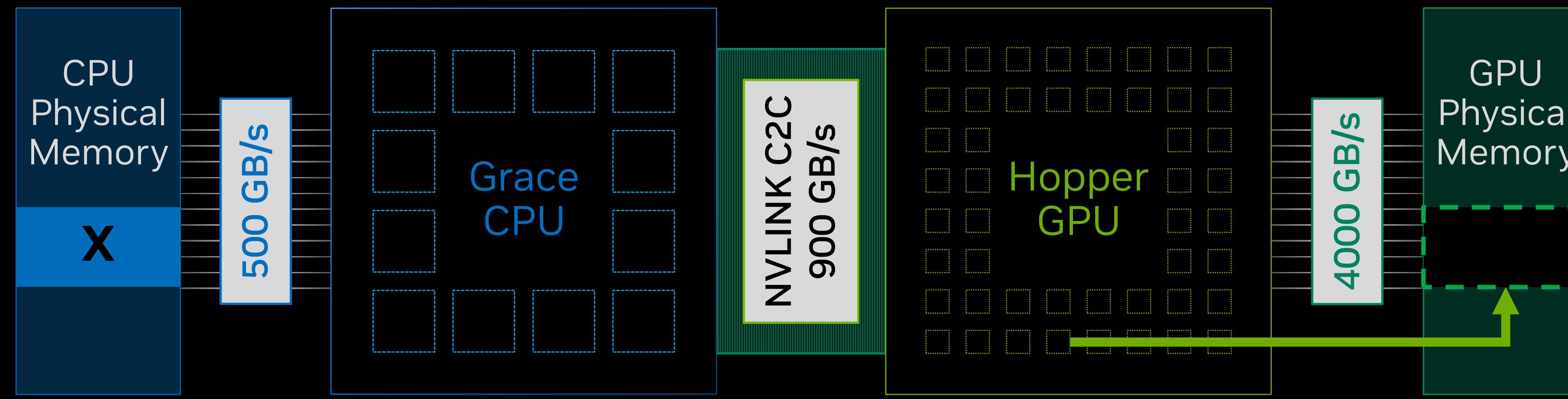
High Bandwidth Memory Access & Automatic Data Migration



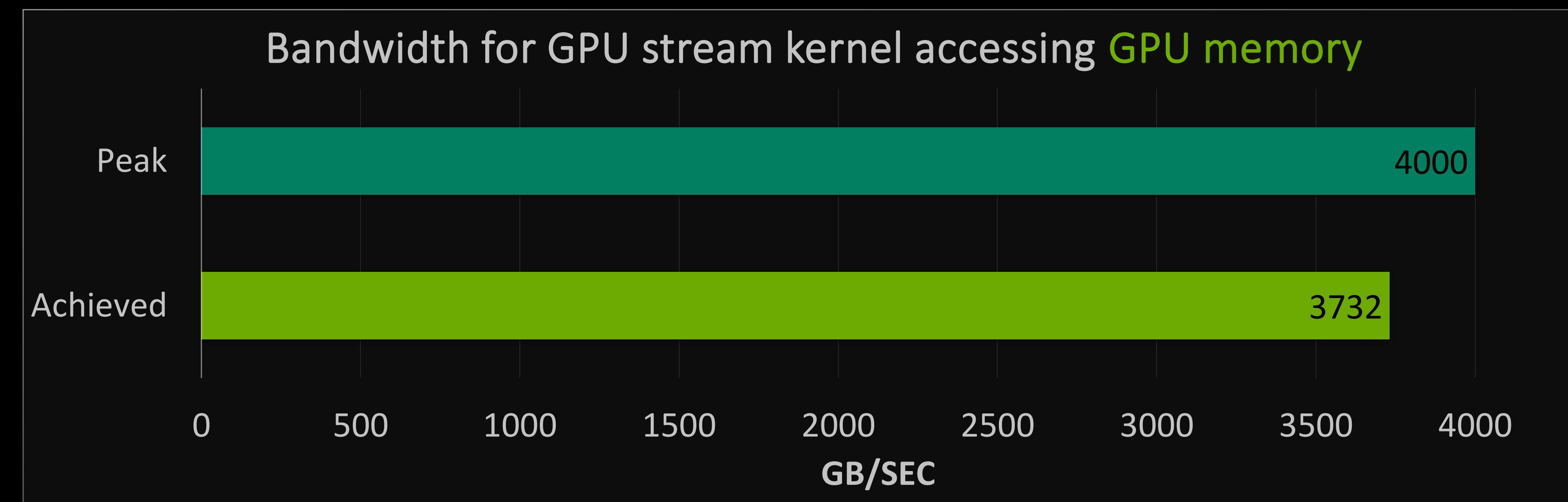
Hopper can access Grace memory
at **full CPU memory speed** of 500 GB/sec



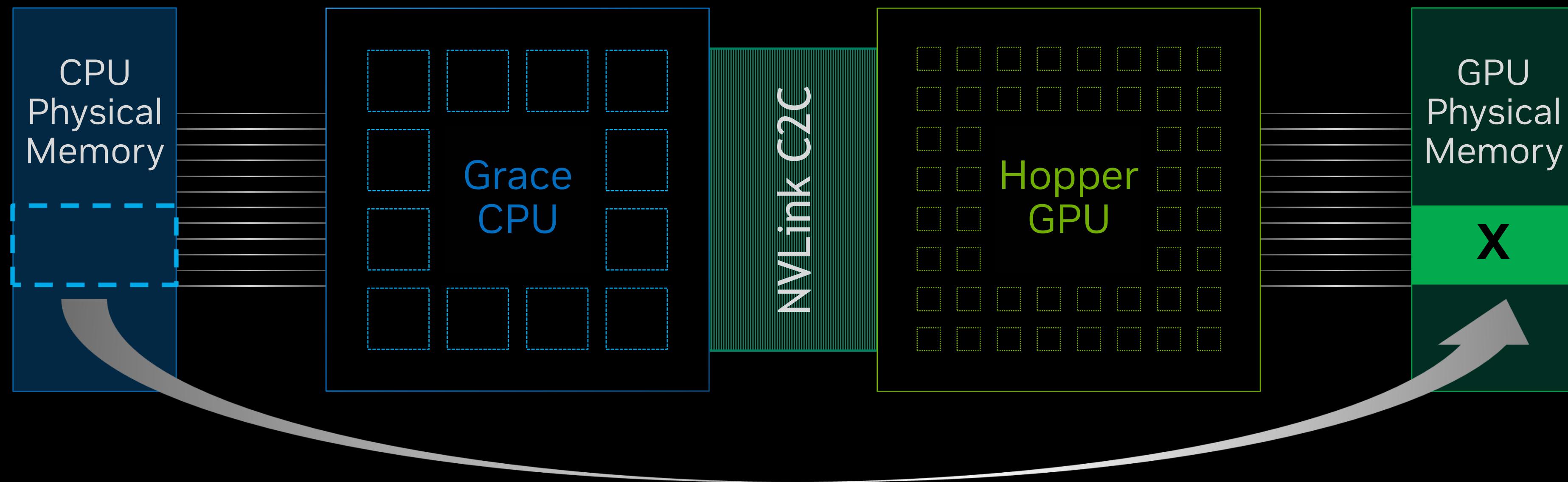
High Bandwidth Memory Access & Automatic Data Migration



But Hopper can access its own memory
at **full HBM speed** of 4000 GB/sec

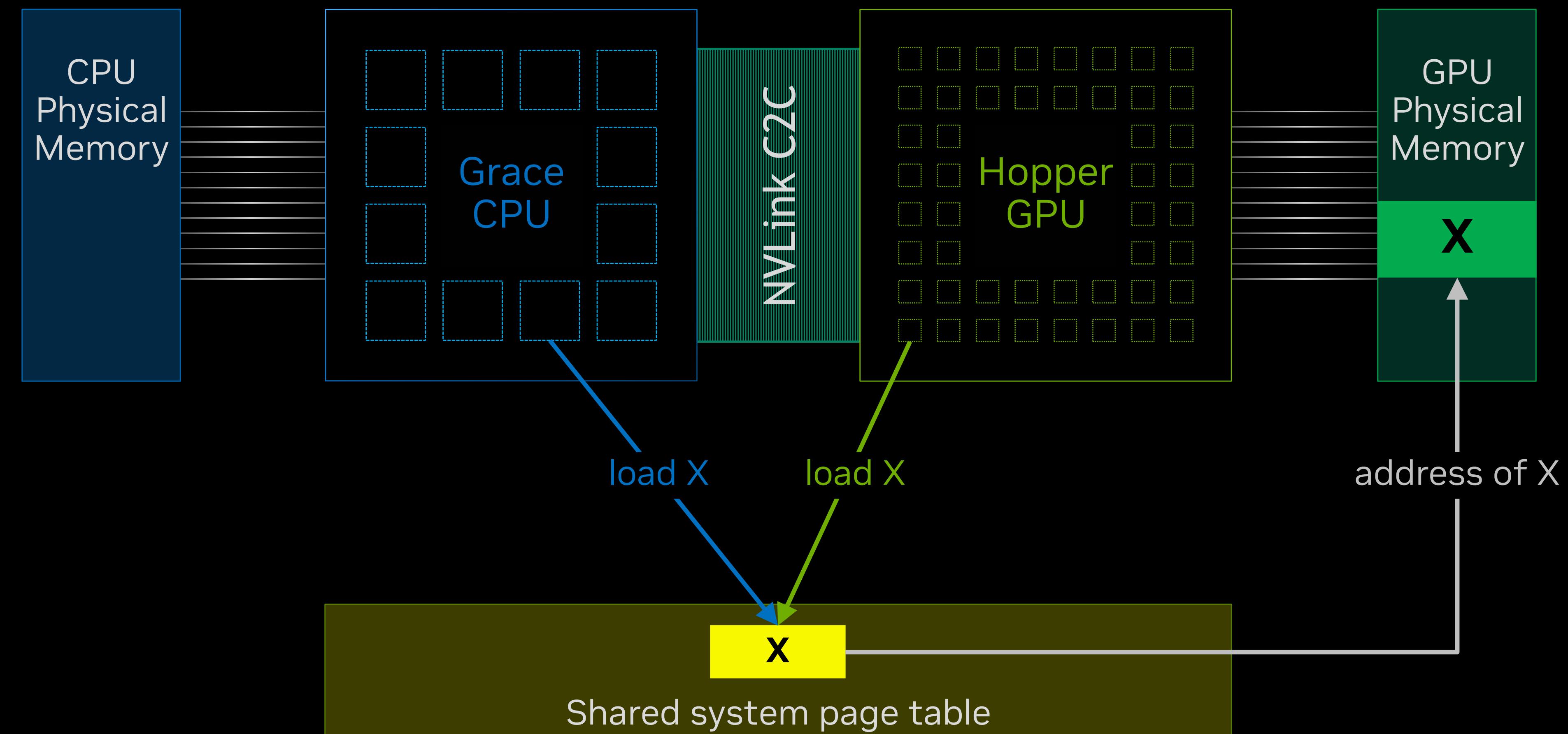


High Bandwidth Memory Access & Automatic Data Migration



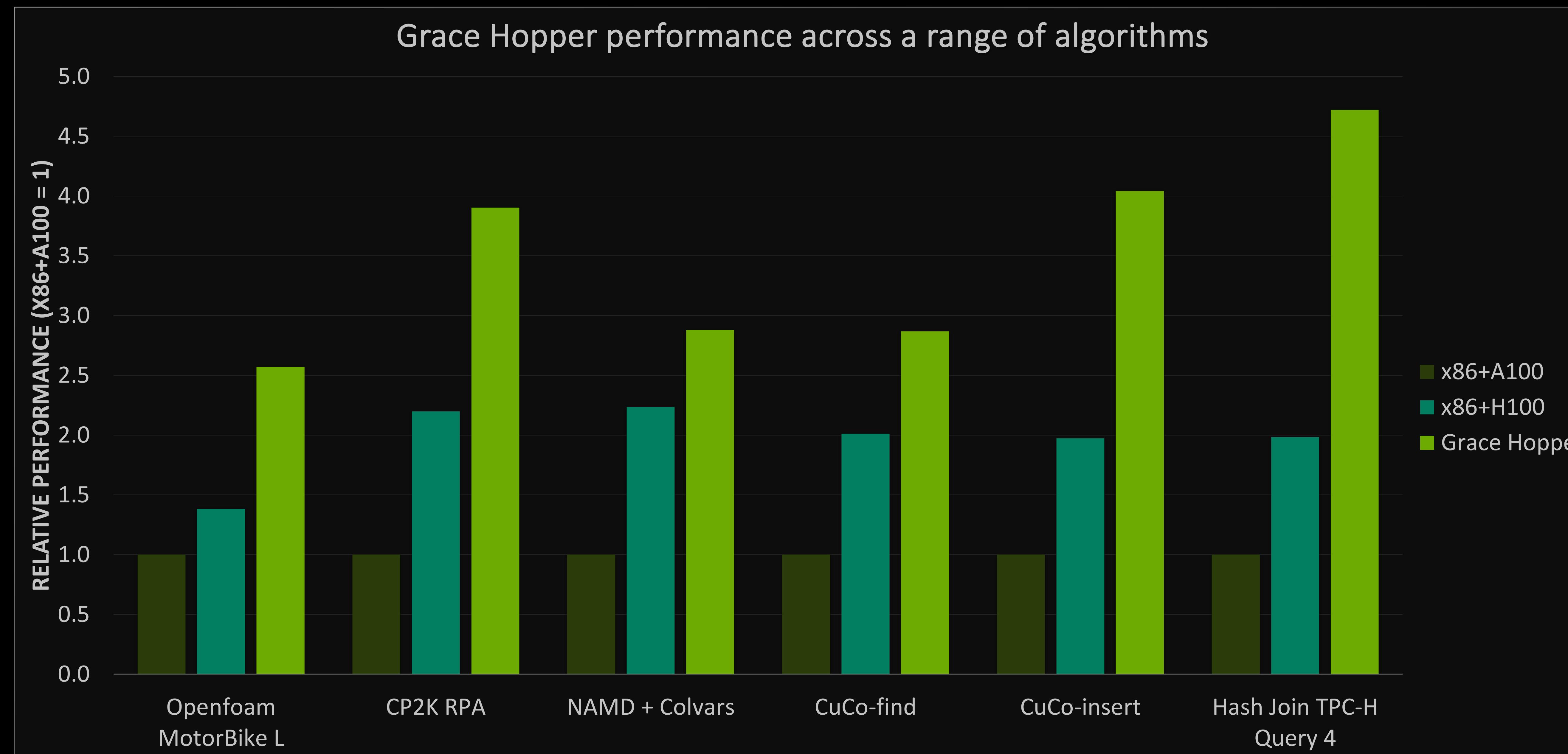
The system can **automatically** migrate both managed and CPU-allocated memory in order to optimize access speed

High Bandwidth Memory Access & Automatic Data Migration

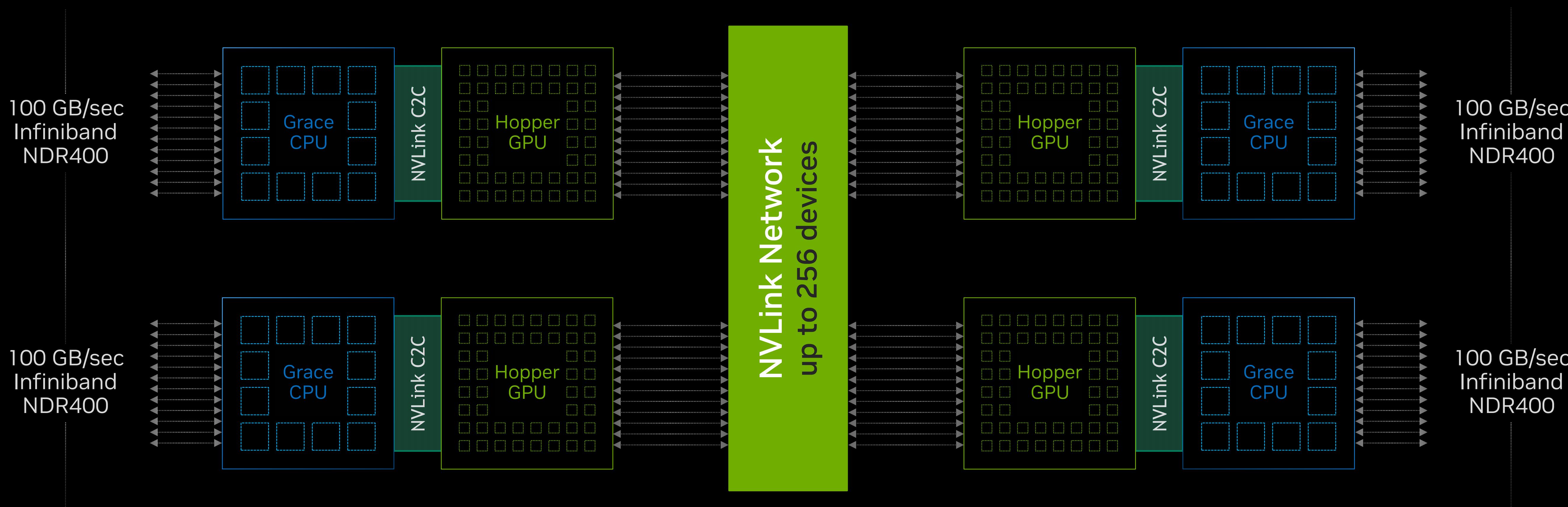


ATS shared page table means that both CPU and GPU automatically access X in its new location after migration

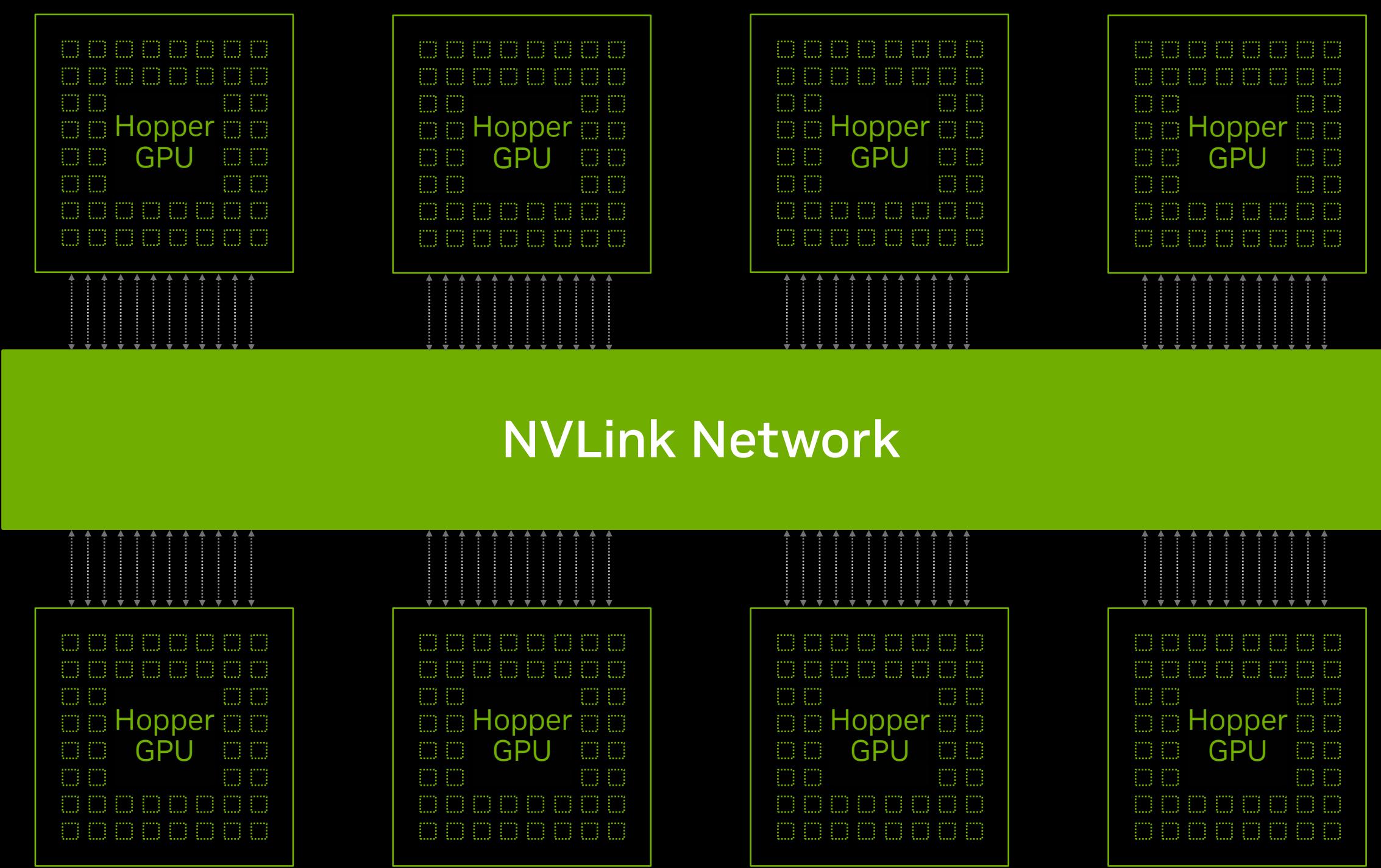
Grace/Hopper SuperChip Performance On Different Workloads

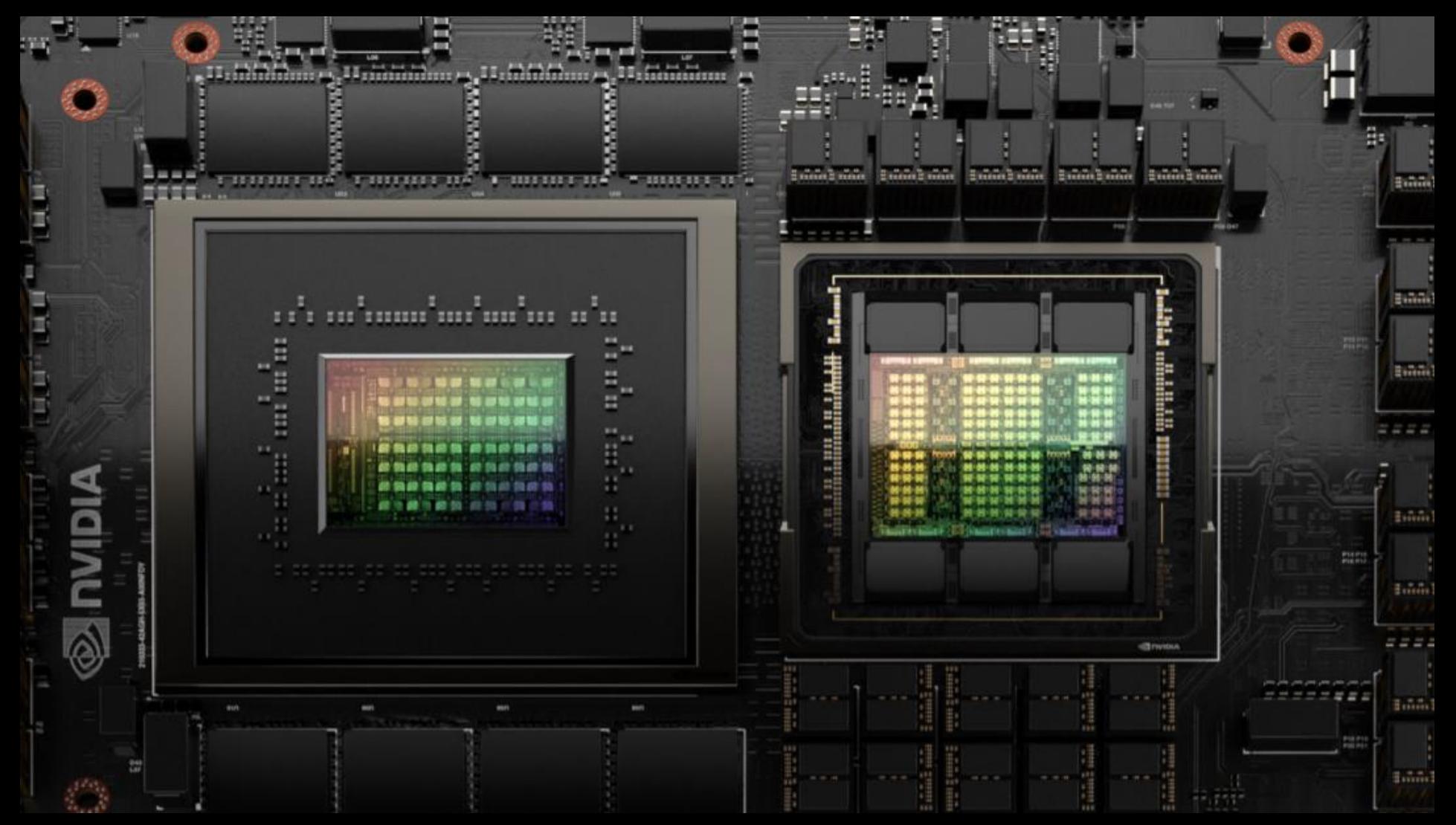


NVLink Connects Up To 256 Superchips

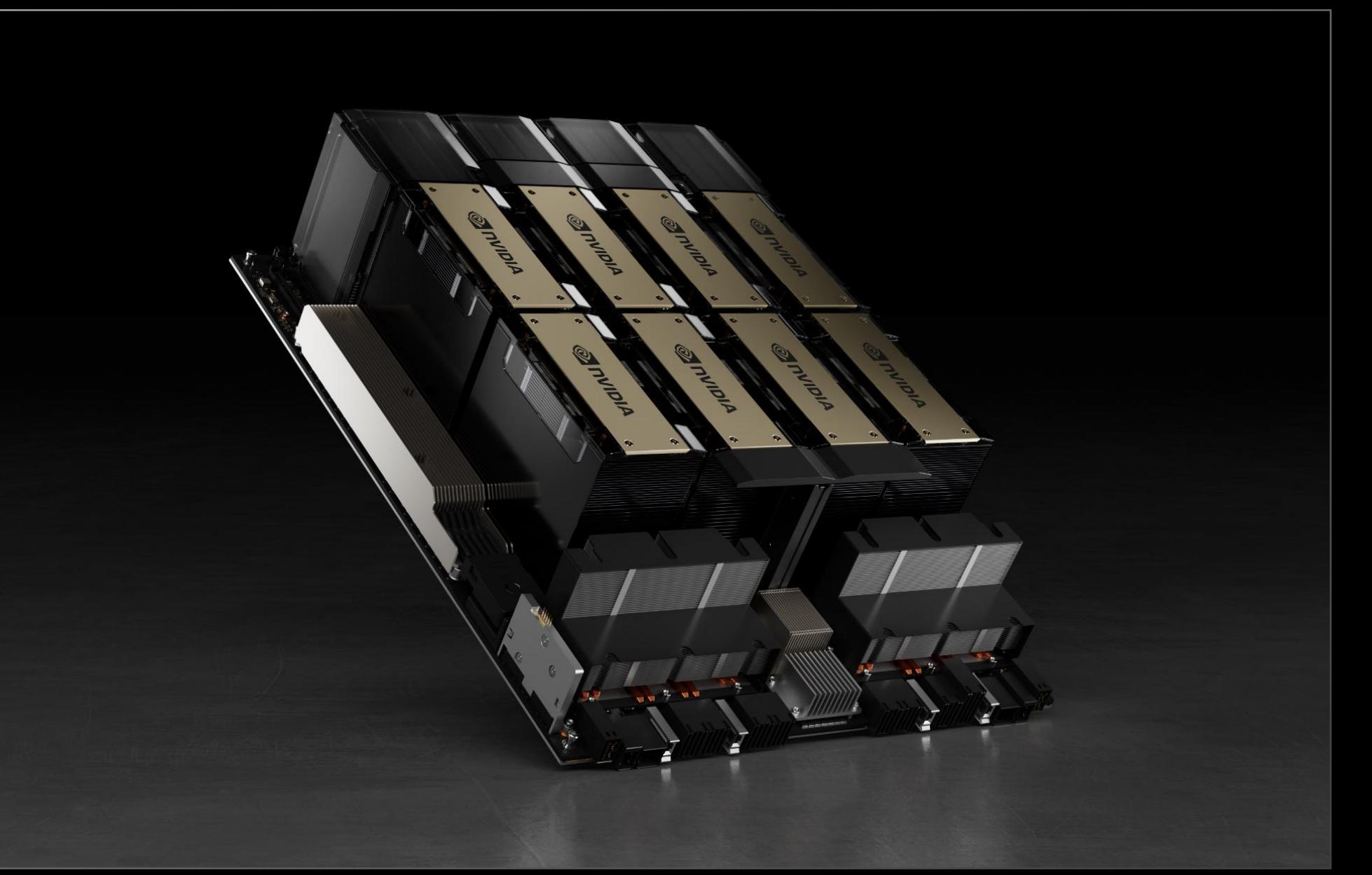


Multi-Chip Systems





Multi-die



Multi-chip



Multi-node

There Is No Future That Is Not Multi-Node



What is CUDA?

CUDA C++

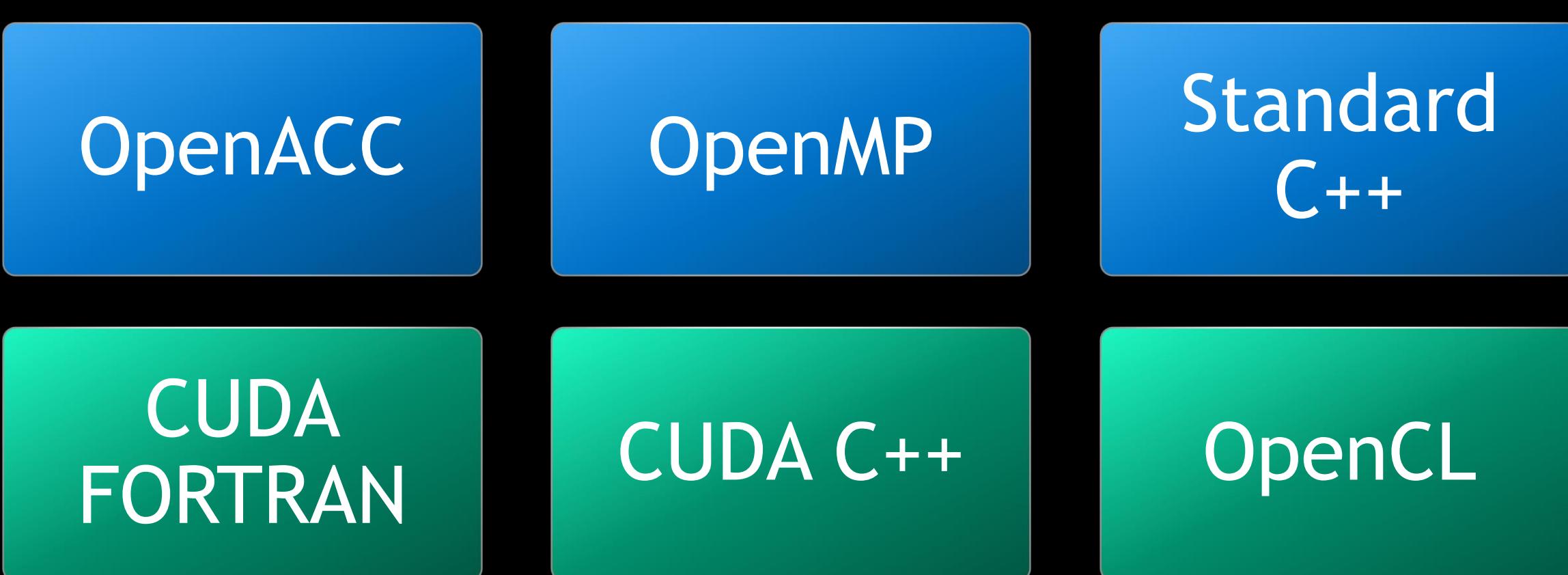
CUDA is More Than CUDA C++

CUDA
FORTRAN

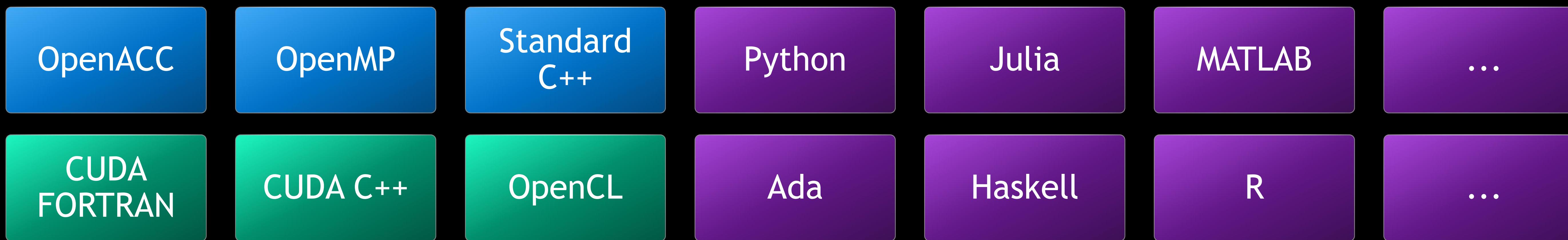
CUDA C++

OpenCL

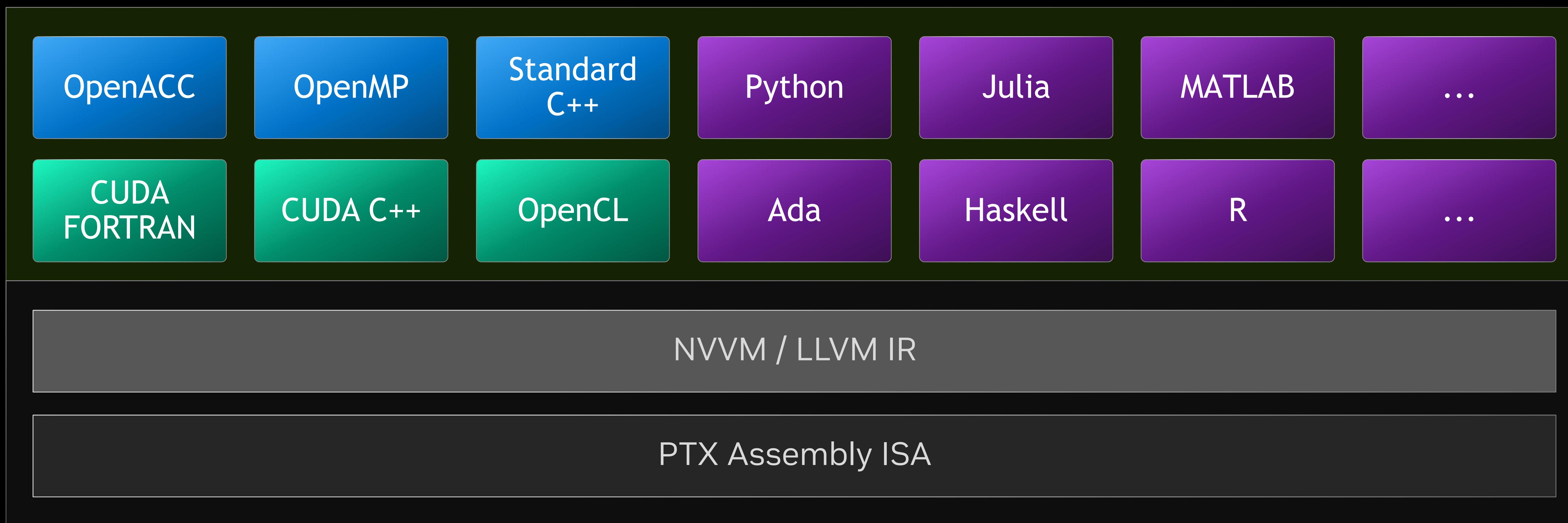
Many Standard Languages Support GPU Computing



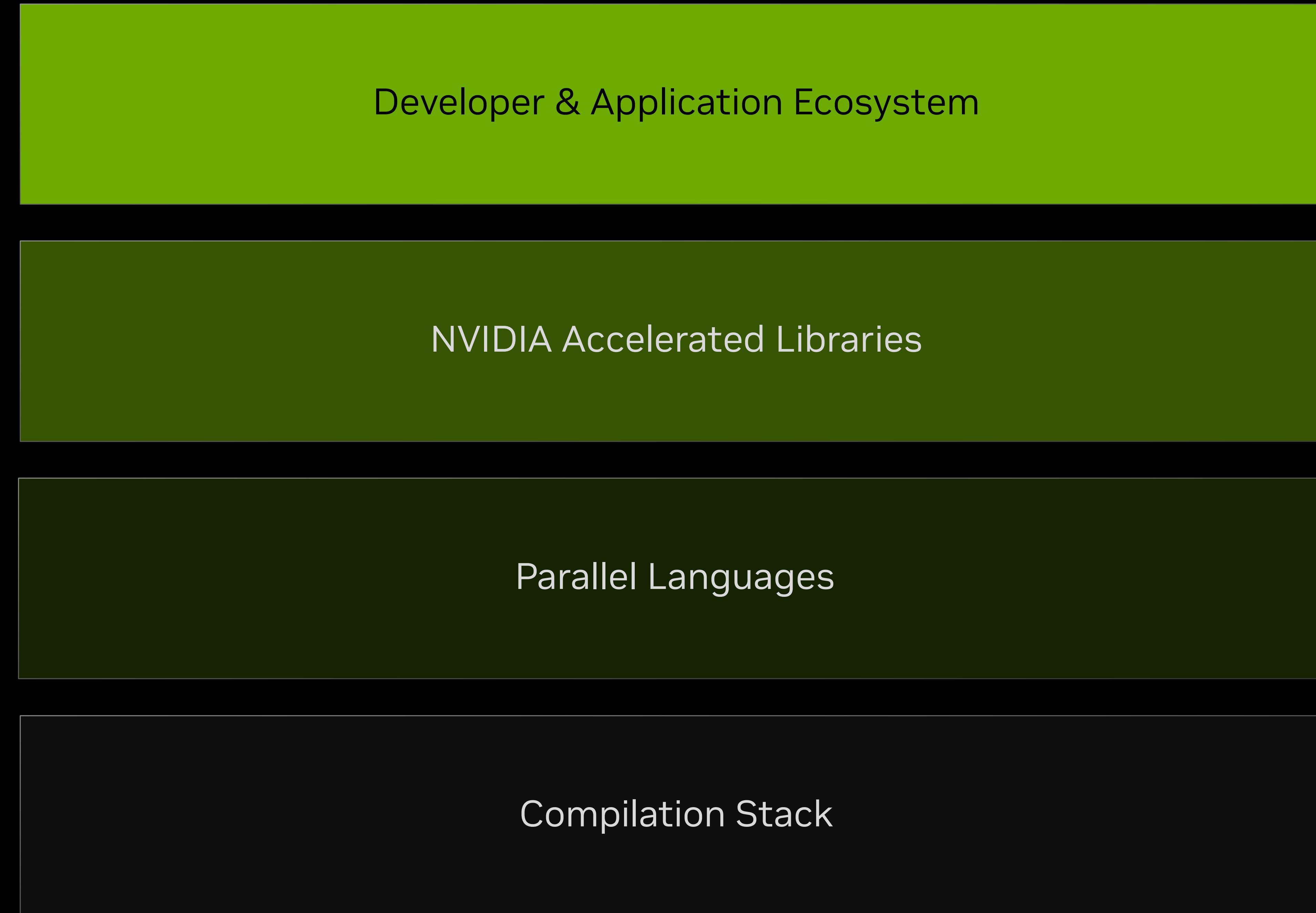
Along With Many Languages Not Supported Directly by NVIDIA



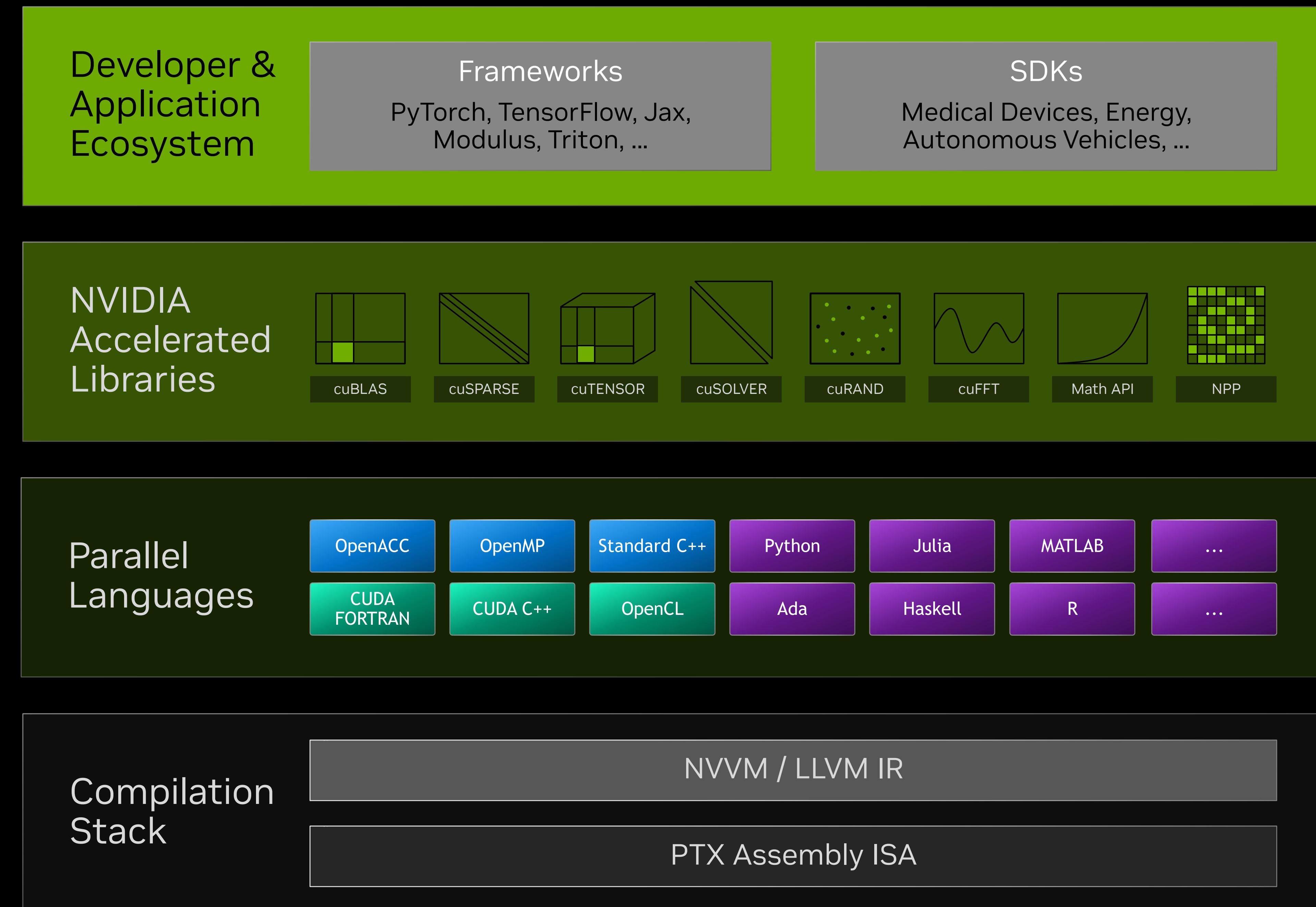
CUDA is a Platform for GPU Computing



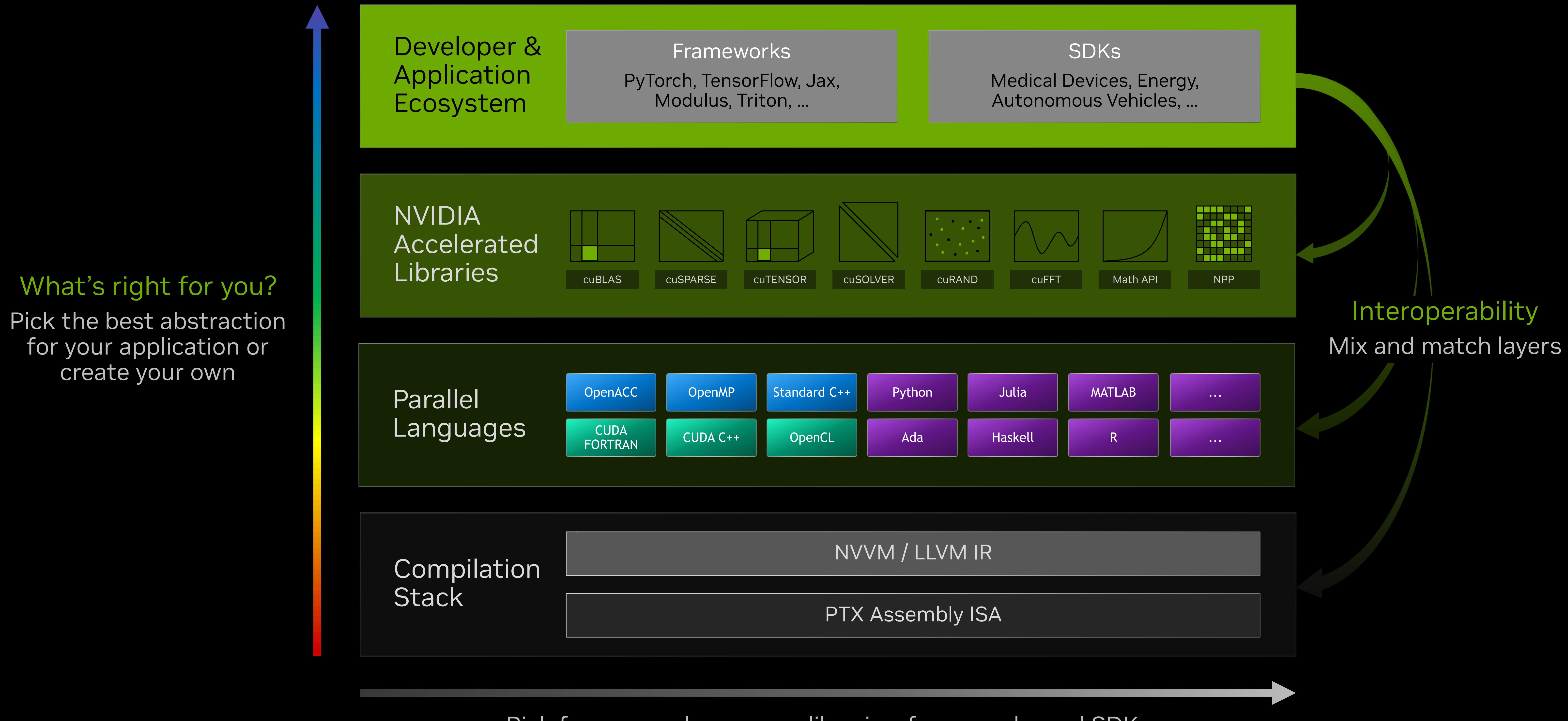
CUDA's Second Dimension: A Platform of Abstraction Layers



Target the Abstraction Layer That Works Best for Your Application



An Ecosystem With a Spectrum of Choices



[How to Write a CUDA Program \[S51210\]](#)





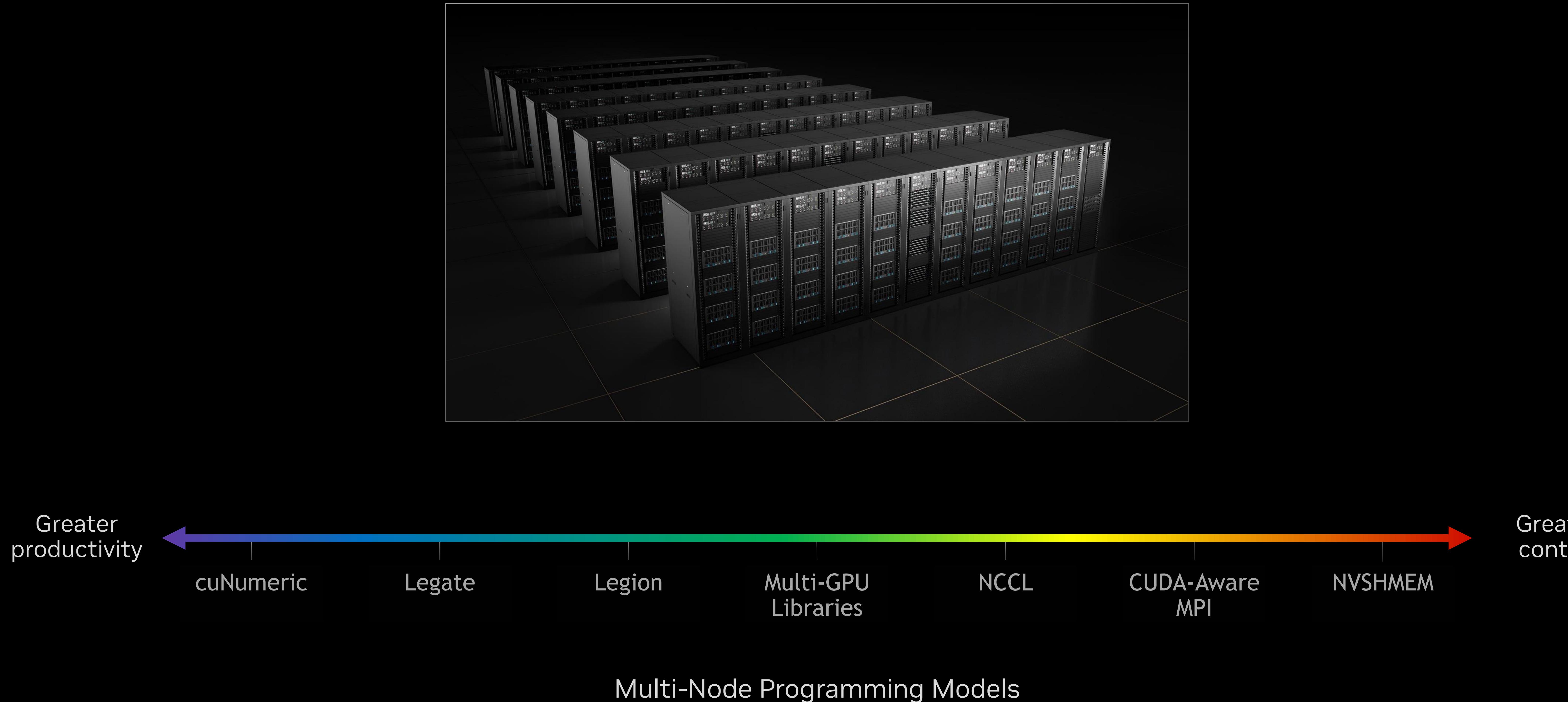
Greater
productivity



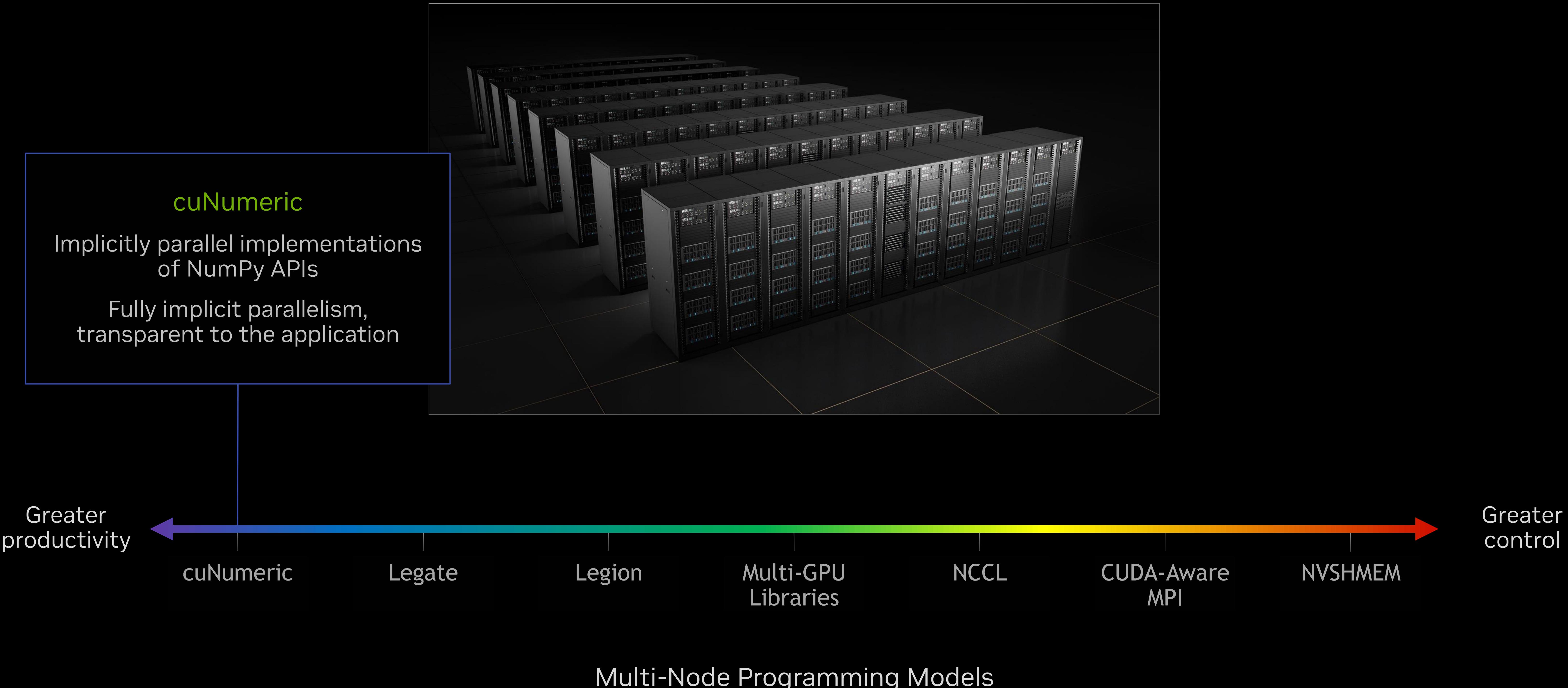
Greater
control



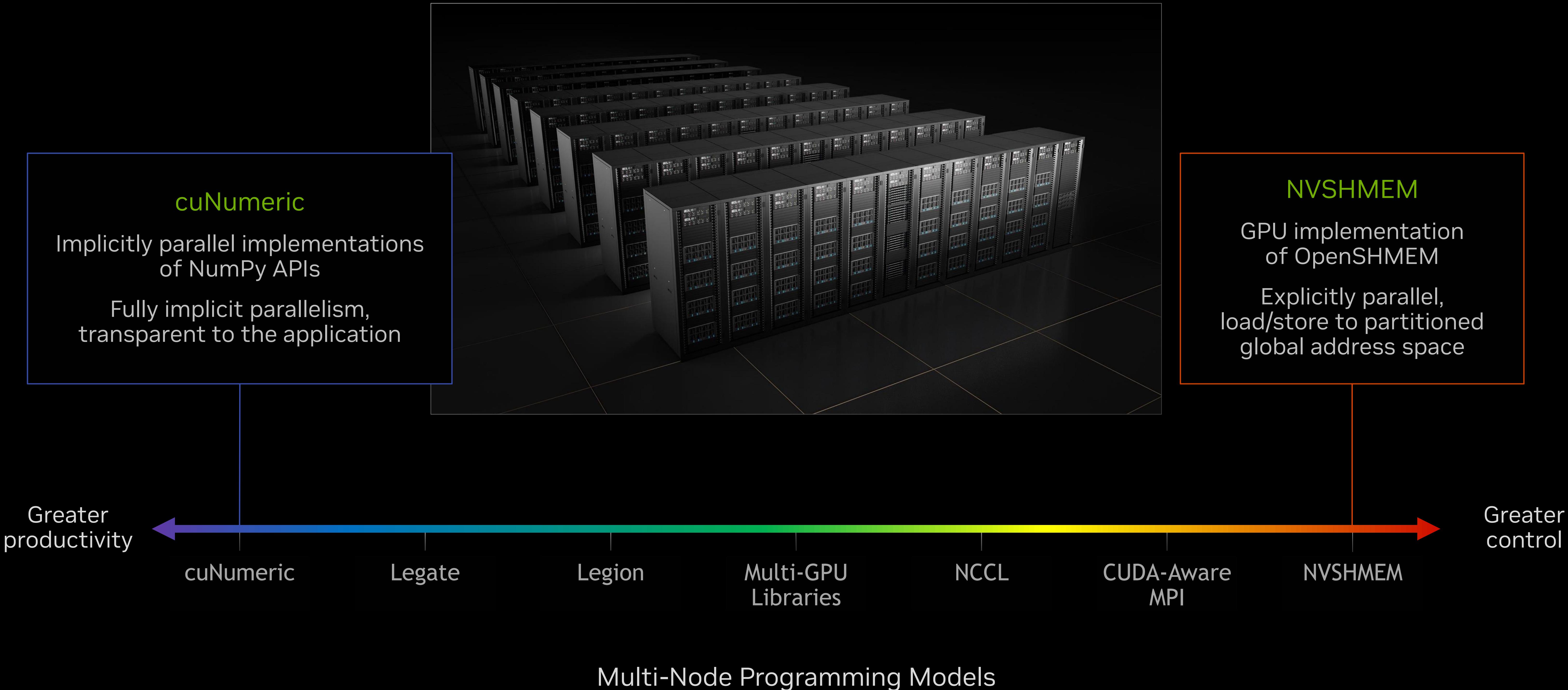
The Third Dimension is Scale



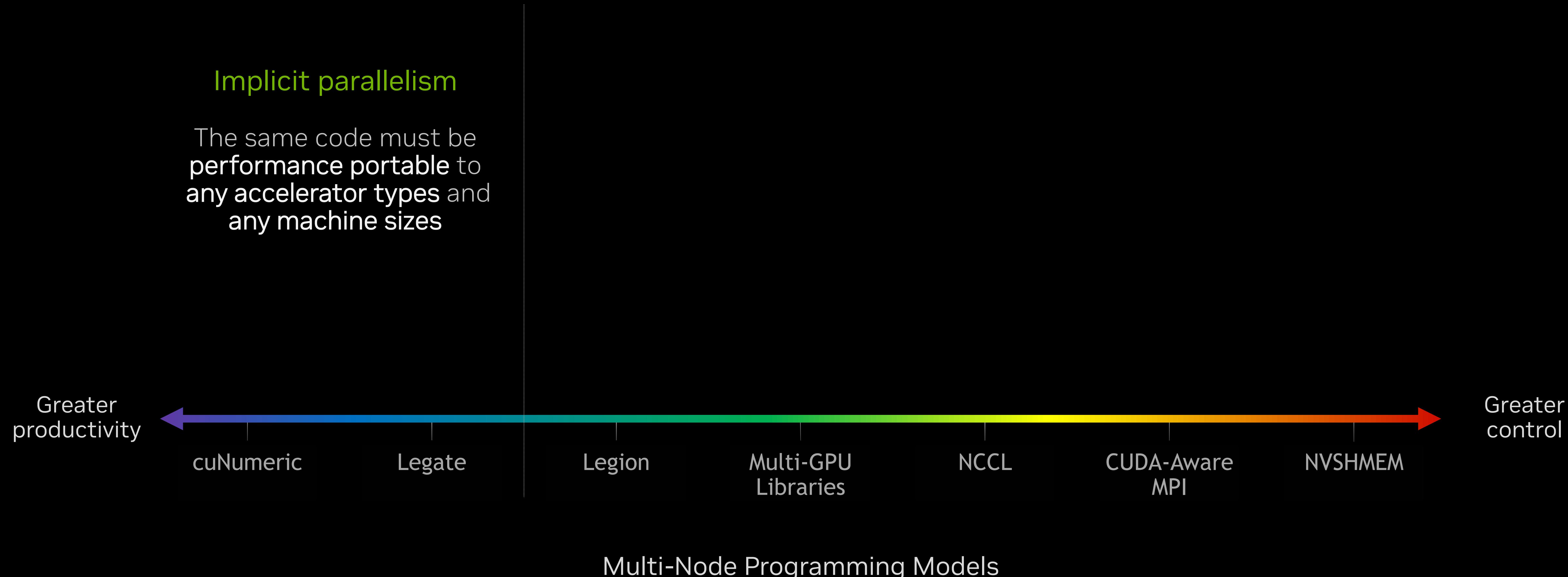
Intersect Where it Makes the Most Sense for Your Application



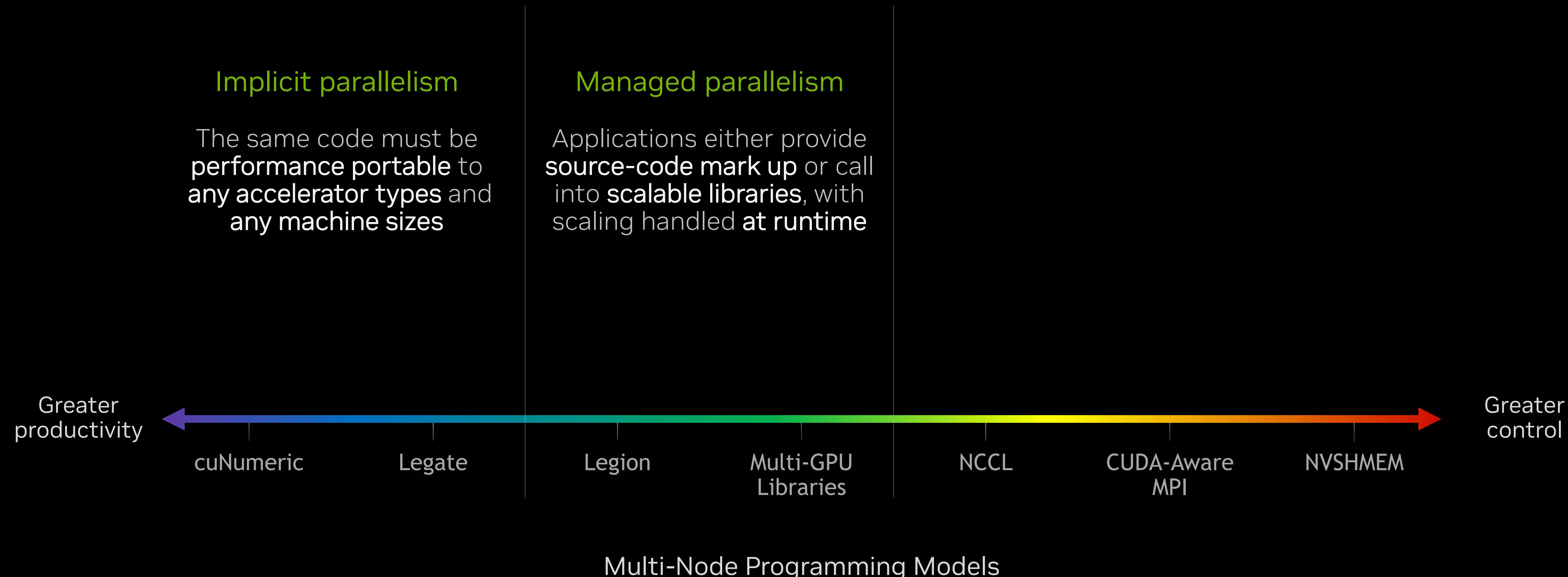
Intersect Where it Makes the Most Sense for Your Application



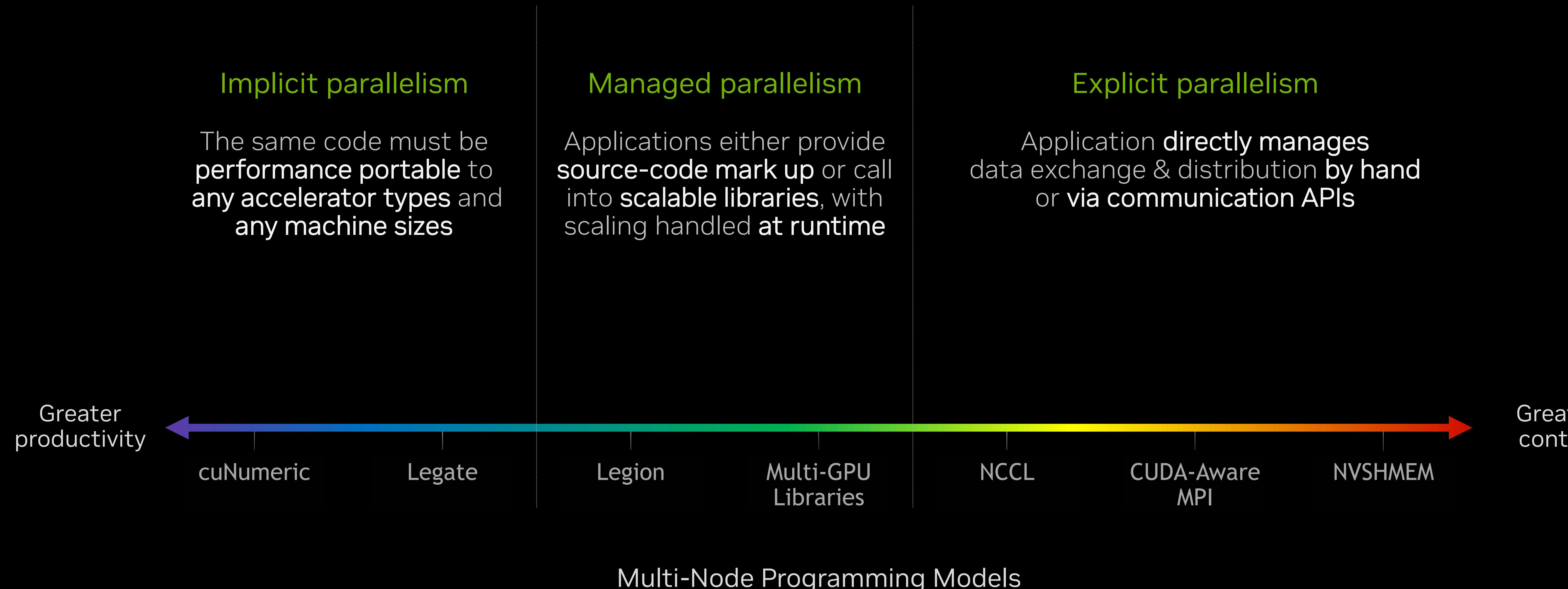
Developing For Multi-Node Systems



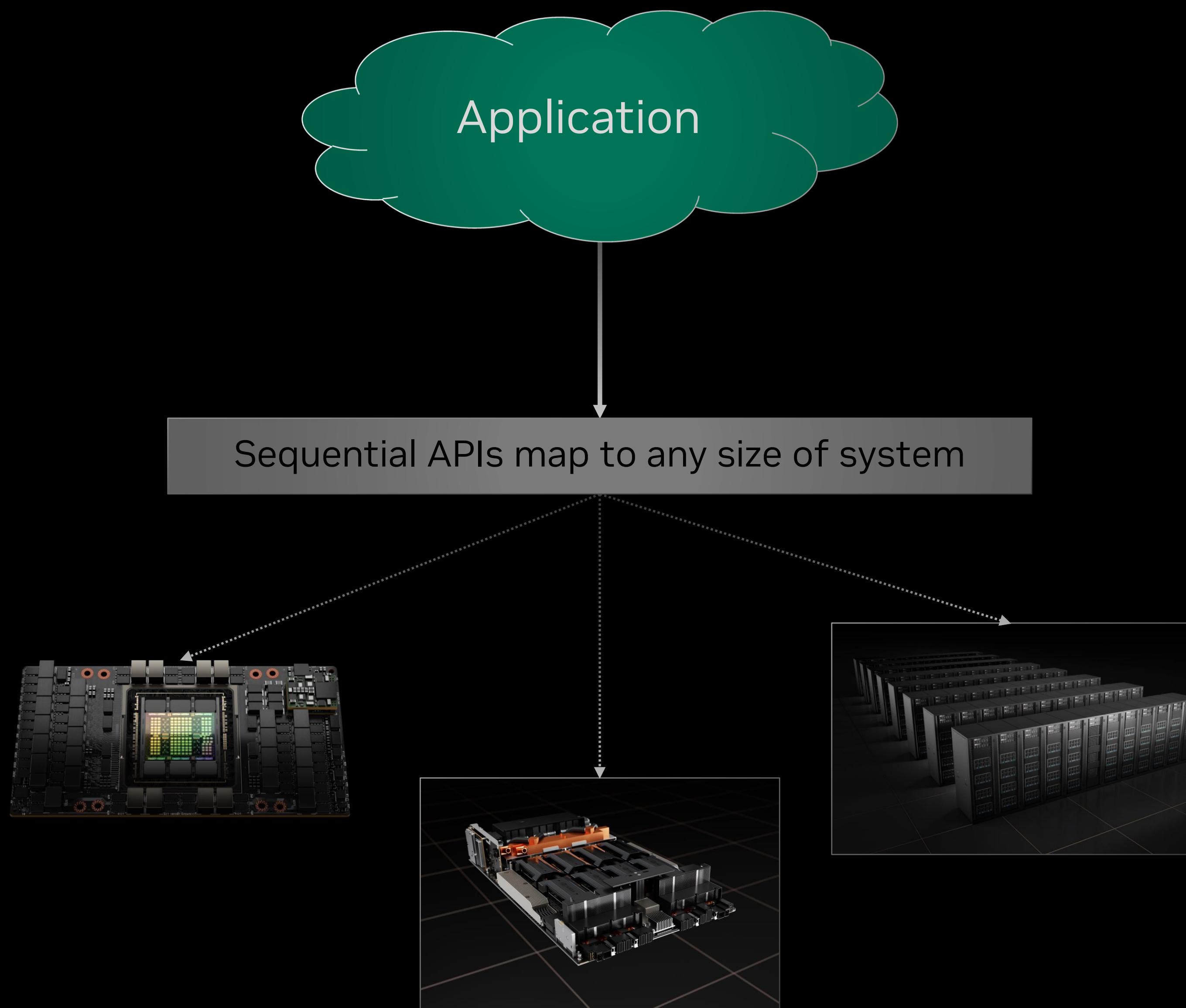
Developing For Multi-Node Systems



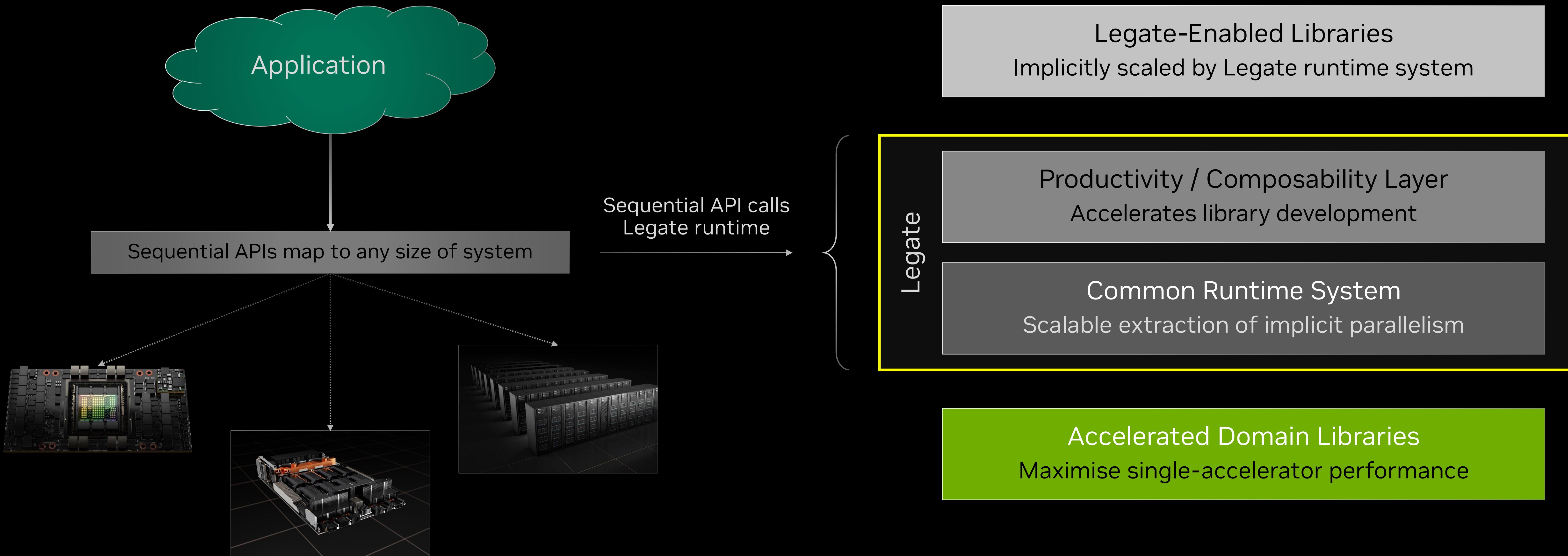
Developing For Multi-Node Systems



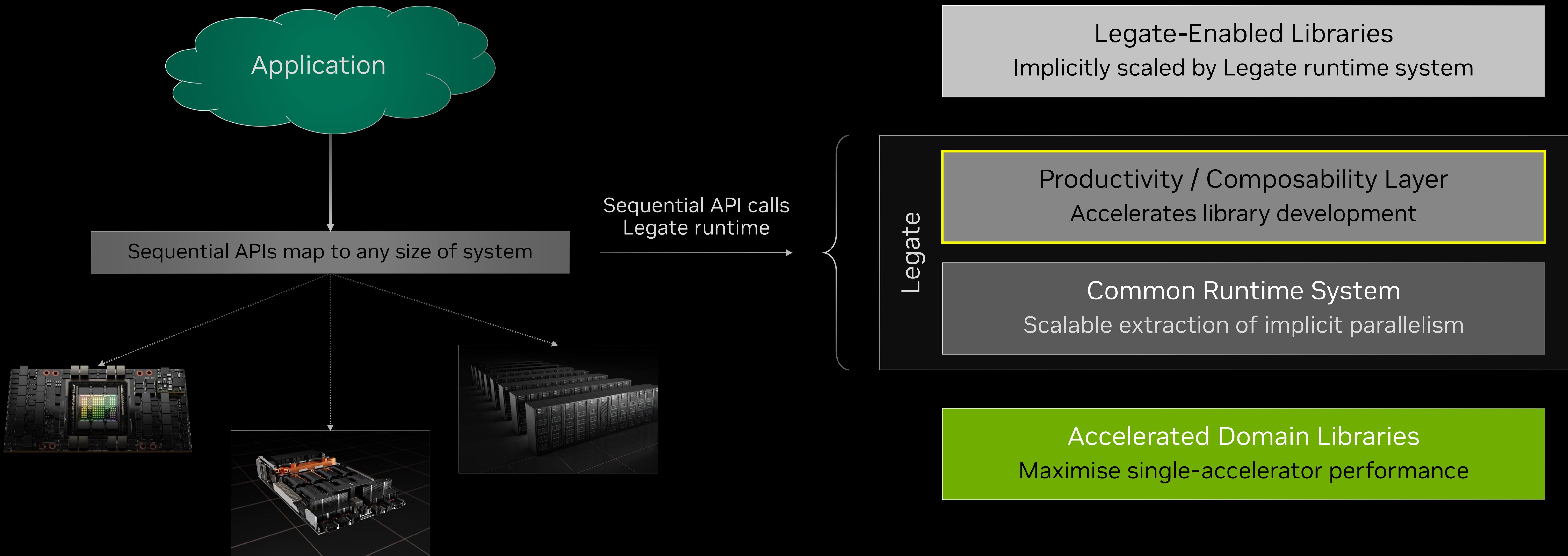
What is “Implicit Parallelism”?



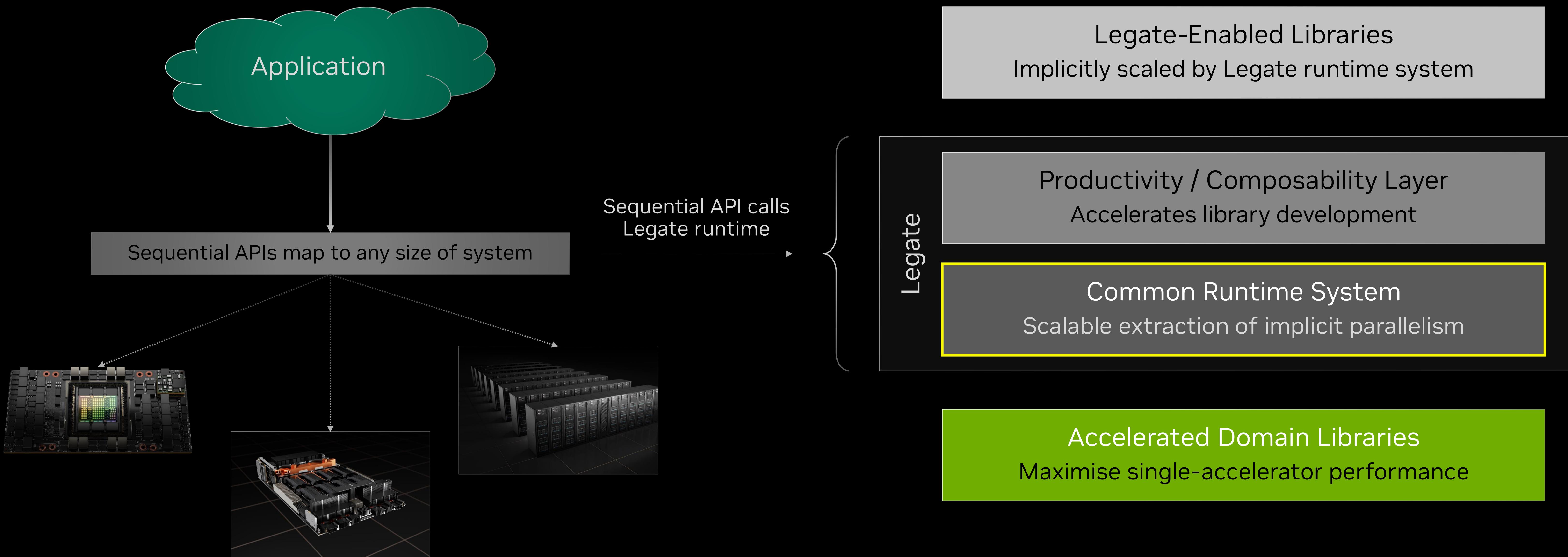
Legate: A Framework for Implicit Parallelism



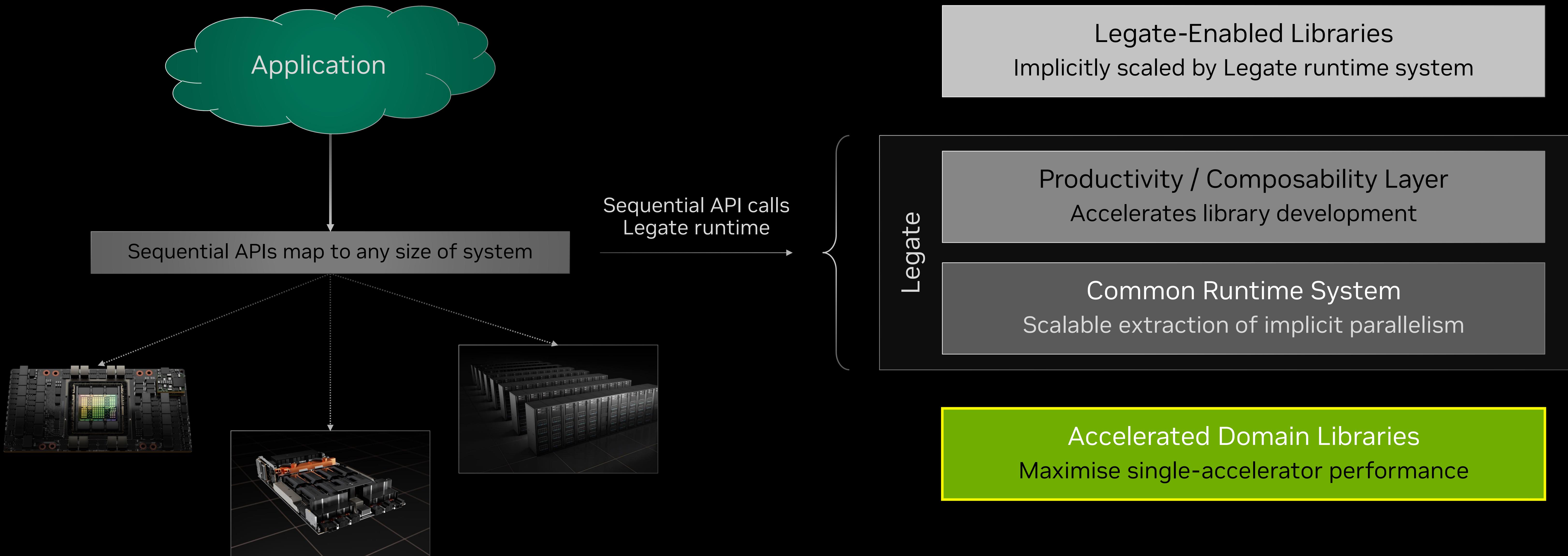
Legate: A Framework for Implicit Parallelism



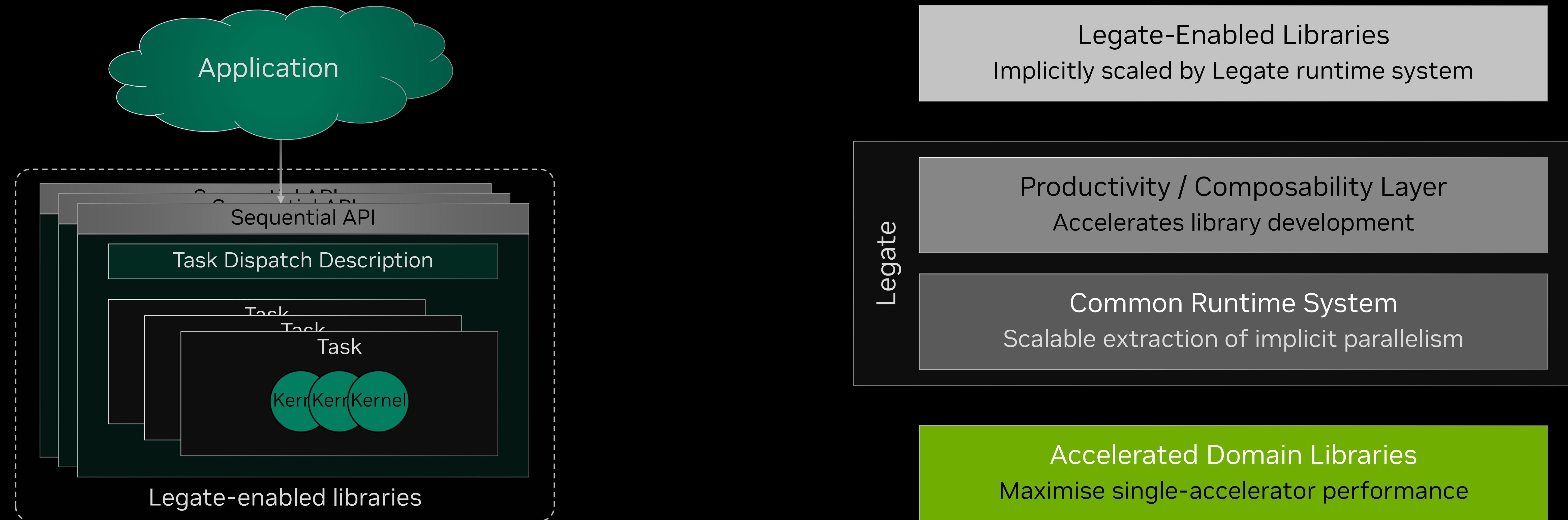
Legate: A Framework for Implicit Parallelism



Legate: A Framework for Implicit Parallelism



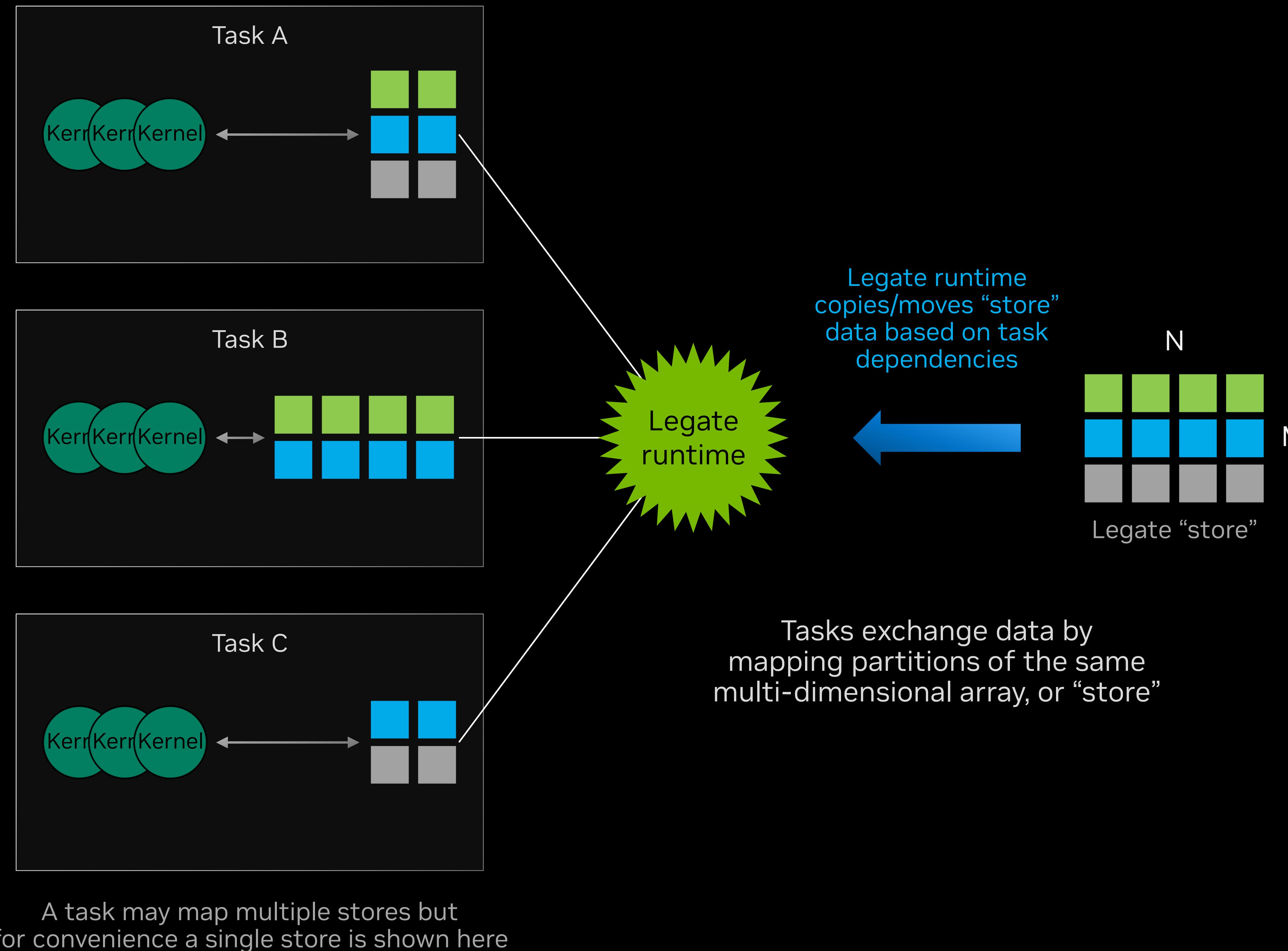
Legate: A Task-Based Model



Legate-enabled libraries

- Are **free** from explicit synchronization
- Are **free** from explicit data movement
- Are thereby **composable** with each other

Unified Data Abstraction for True Composability

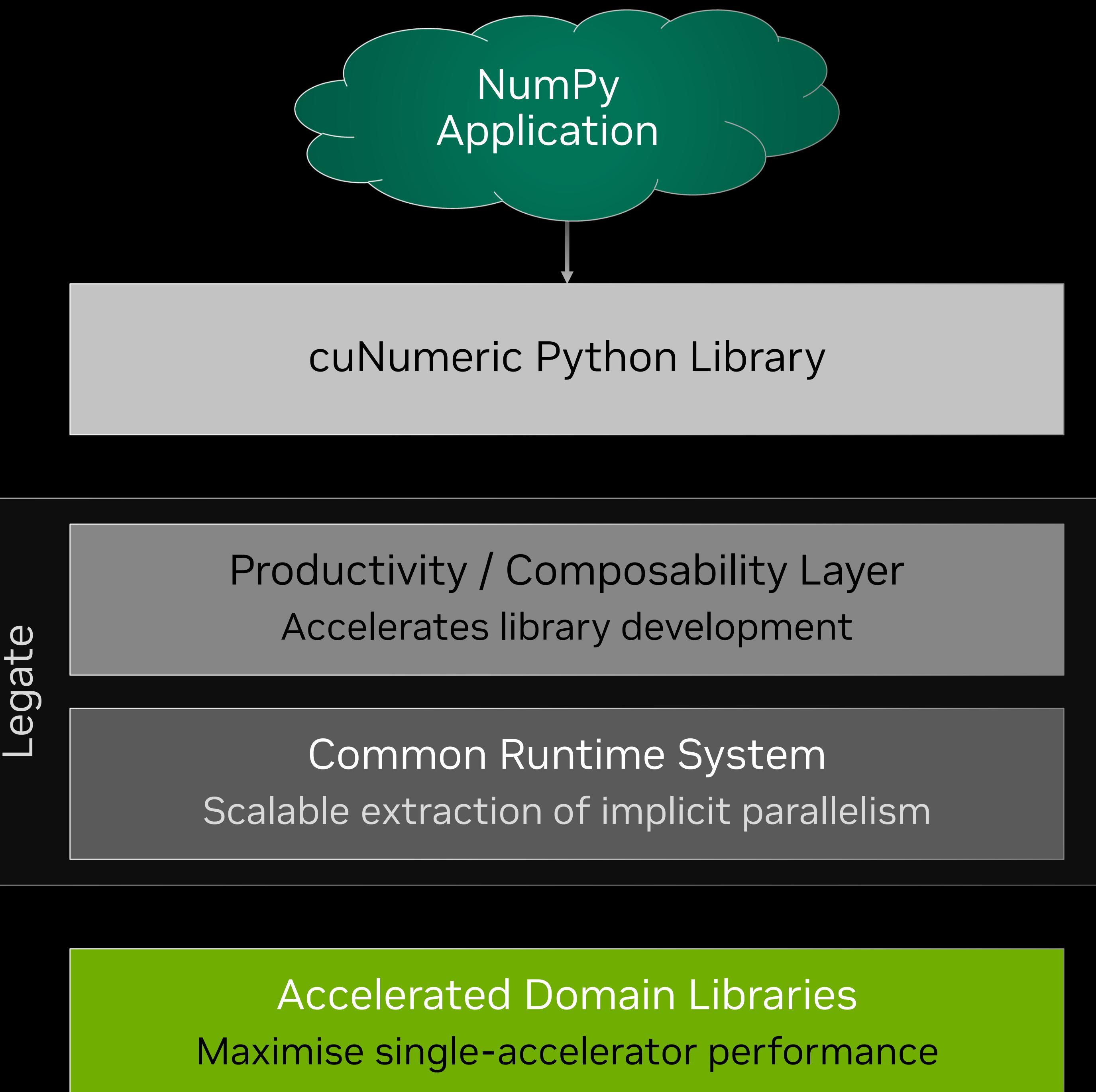


cuNumeric - Implicitly Parallel Implementations of NumPy APIs

Stencil Benchmark

No modifications required to scale to a thousand GPUs

```
 32
 33 def run_stencil(N, I, warmup, timing): # noqa: E741
 34     grid = initialize(N)
 35
 36     print("Running Jacobi stencil...")
 37     center = grid[1:-1, 1:-1]
 38     north = grid[0:-2, 1:-1]
 39     east = grid[1:-1, 2:]
 40     west = grid[1:-1, 0:-2]
 41     south = grid[2:, 1:-1]
 42
 43     timer.start()
 44     for i in range(I + warmup):
 45         if i == warmup:
 46             timer.start()
 47             average = center + north + east + west + south
 48             work = 0.2 * average
 49             center[:] = work
 50     total = timer.stop()
 51
 52     if timing:
 53         print(f"Elapsed Time: {total} ms")
 54     return total
```

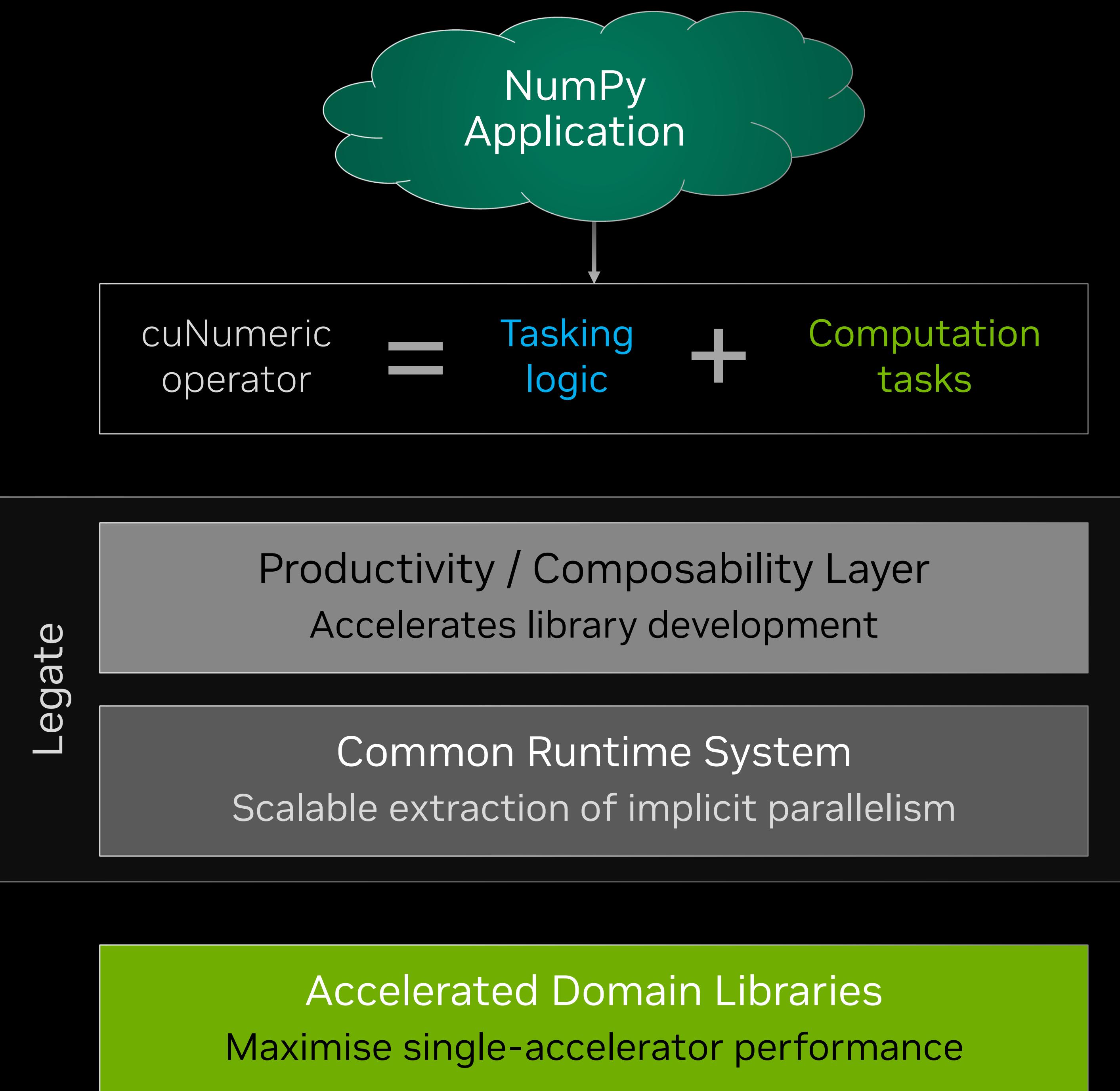


cuNumeric - Implicitly Parallel Implementations of NumPy APIs

Stencil Benchmark

No modifications required to scale to a thousand GPUs

```
 32
 33 def run_stencil(N, I, warmup, timing): # noqa: E741
 34     grid = initialize(N)
 35
 36     print("Running Jacobi stencil...")
 37     center = grid[1:-1, 1:-1]
 38     north = grid[0:-2, 1:-1]
 39     east = grid[1:-1, 2:]
 40     west = grid[1:-1, 0:-2]
 41     south = grid[2:, 1:-1]
 42
 43     timer.start()
 44     for i in range(I + warmup):
 45         if i == warmup:
 46             timer.start()
 47             average = center + north + east + west + south
 48             work = 0.2 * average
 49             center[:] = work
 50     total = timer.stop()
 51
 52     if timing:
 53         print(f"Elapsed Time: {total} ms")
 54     return total
```



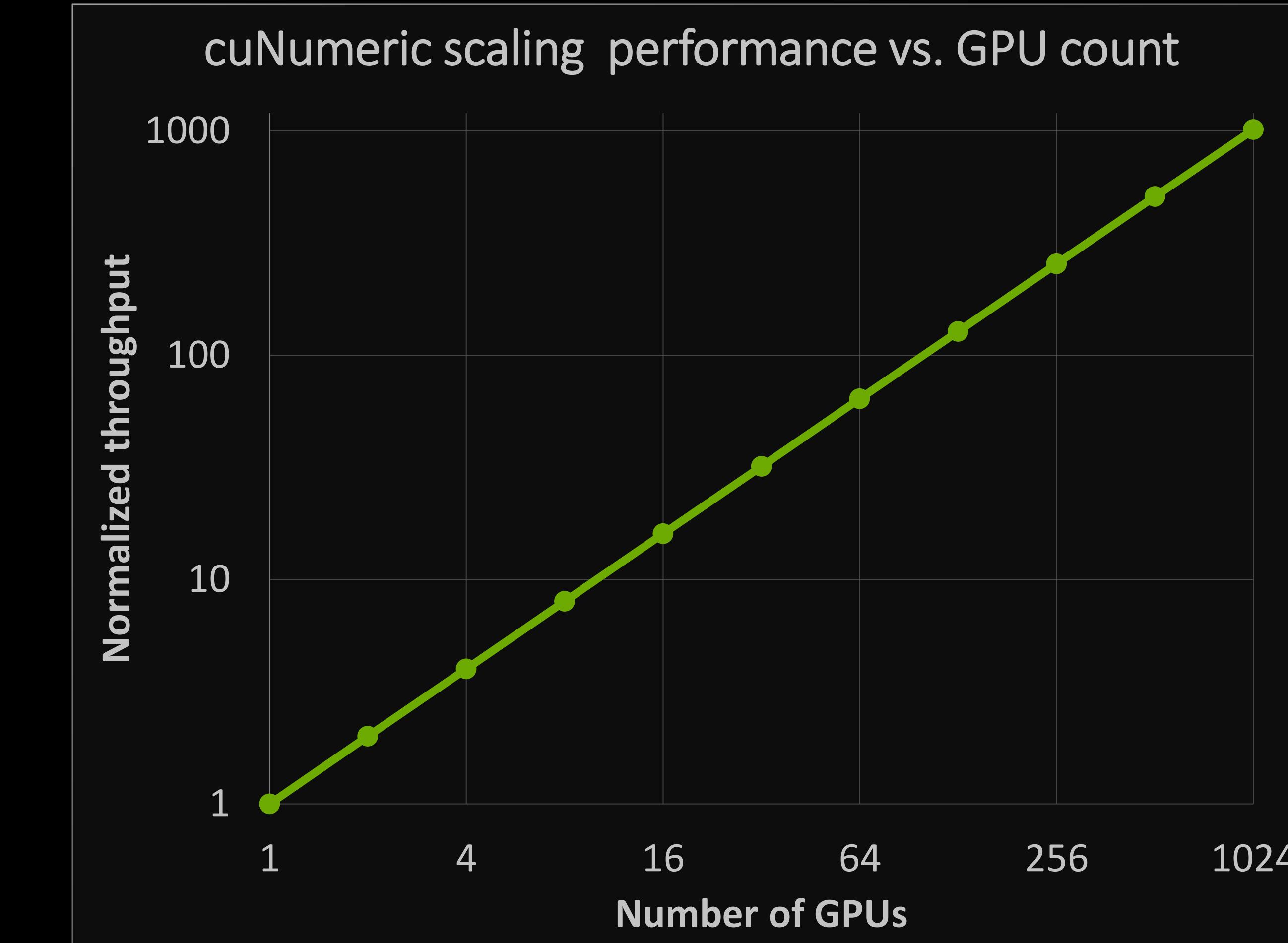
cuNumeric - Implicitly Parallel Implementations of NumPy APIs

cuNumeric Beta: out March 2023

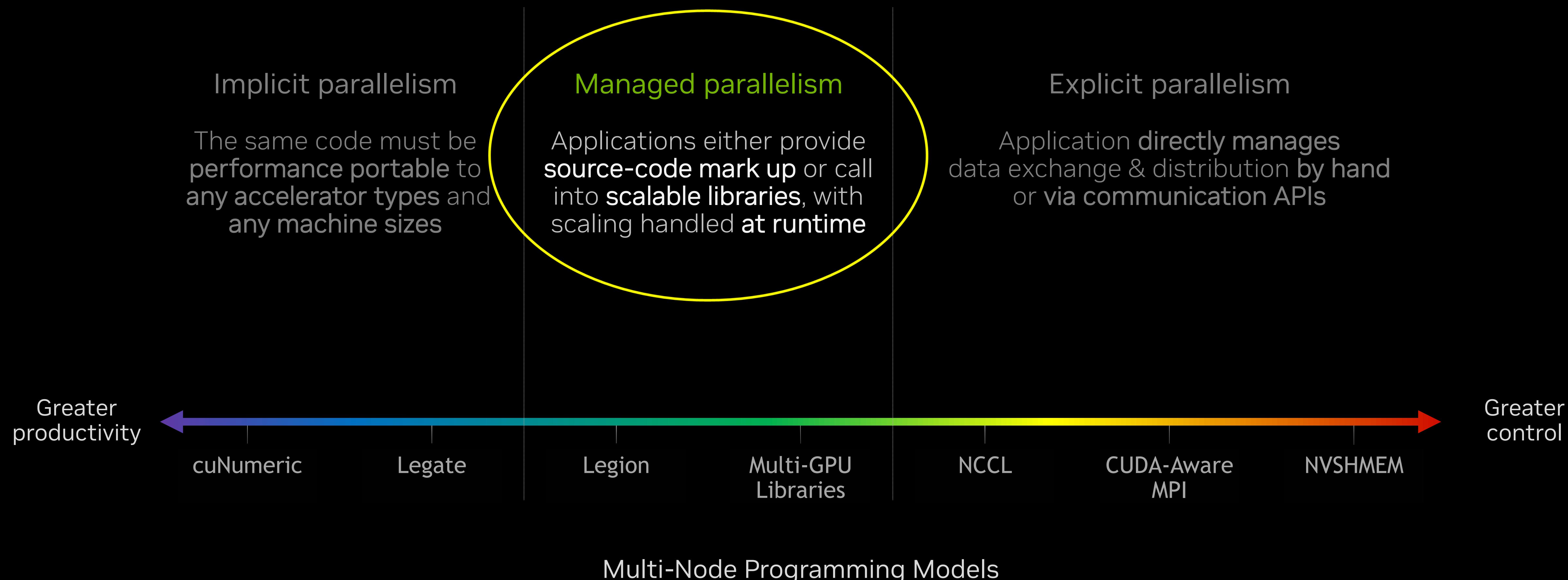
Stencil Benchmark

No modifications required to scale to a thousand GPUs

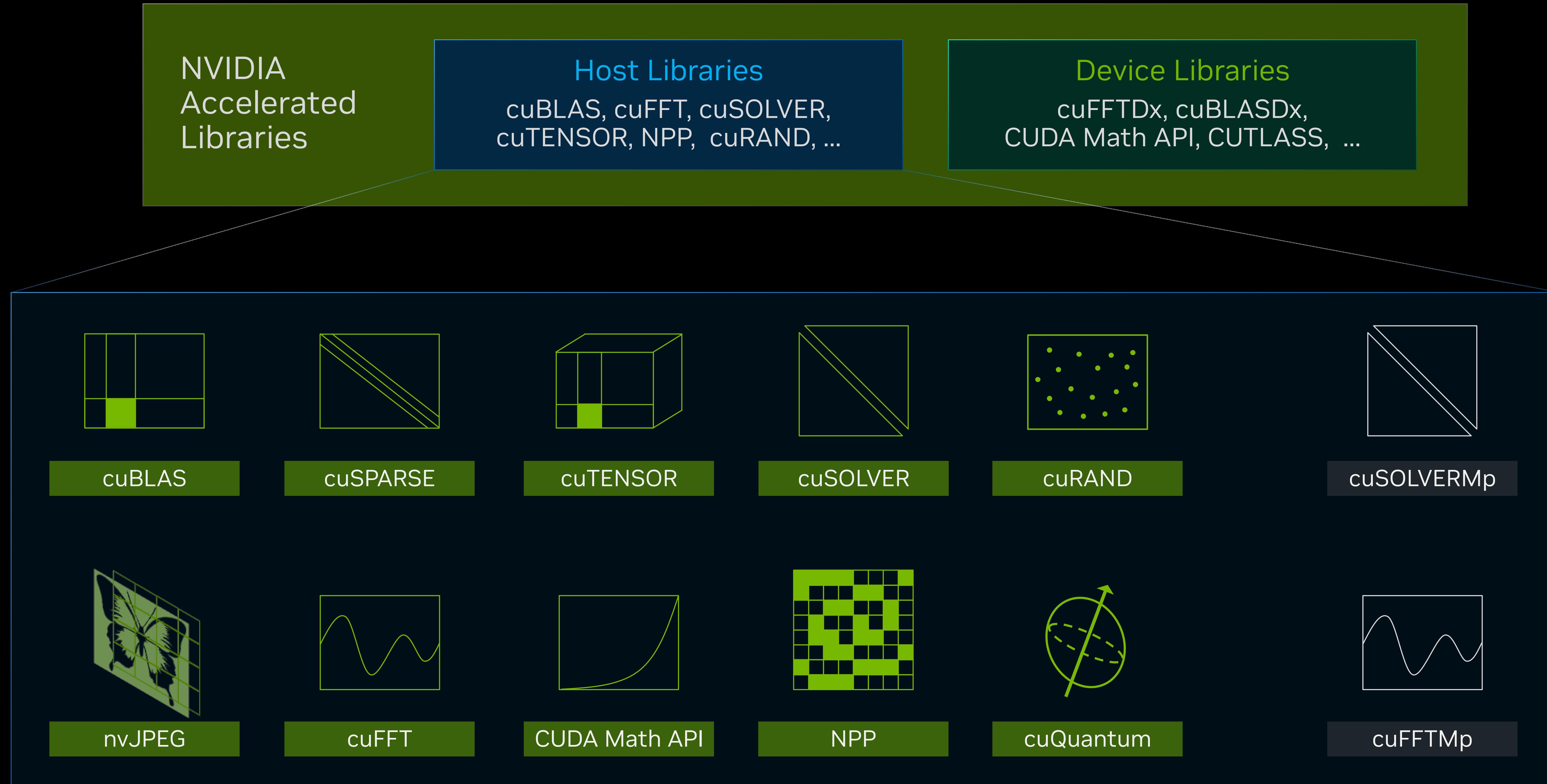
```
 32
33     def run_stencil(N, I, warmup, timing): # noqa: E741
34         grid = initialize(N)
35
36         print("Running Jacobi stencil...")
37         center = grid[1:-1, 1:-1]
38         north = grid[0:-2, 1:-1]
39         east = grid[1:-1, 2:]
40         west = grid[1:-1, 0:-2]
41         south = grid[2:, 1:-1]
42
43         timer.start()
44         for i in range(I + warmup):
45             if i == warmup:
46                 timer.start()
47                 average = center + north + east + west + south
48                 work = 0.2 * average
49                 center[:] = work
50         total = timer.stop()
51
52         if timing:
53             print(f"Elapsed Time: {total} ms")
54         return total
```



Managed Parallelism Through Accelerated Libraries



Multi-GPU Multi-Node (MGMN) Libraries



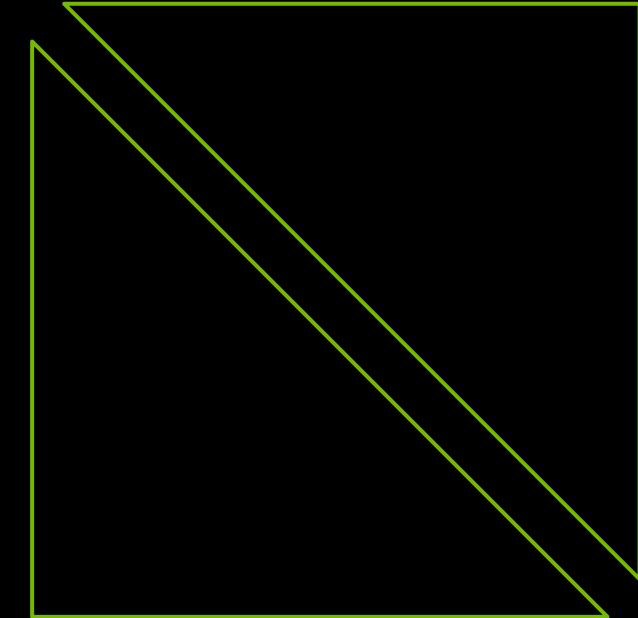
cuSOLVERMp: Distributed Symmetric Eigenvalue Solver

Local matrix size per GPU (32768^2)

Double precision real data

Speedup of $\sim 1.6x$ over ELPA on
130K, 262K, 524K sizes

Speedup of 1.35x on $1M^2$ problem



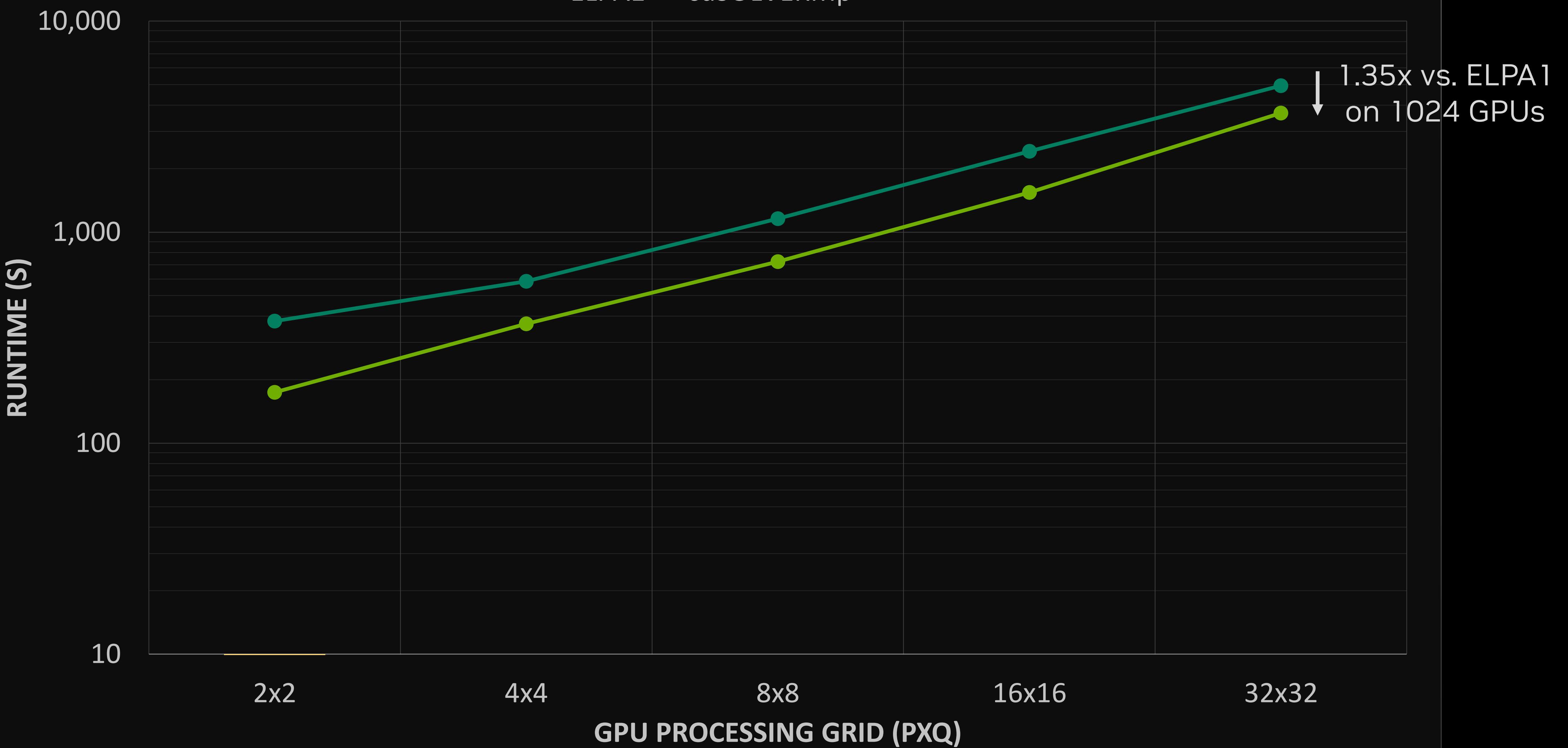
cuSOLVERMp

Weak scaling of PDSYEVD on Selene SuperPOD

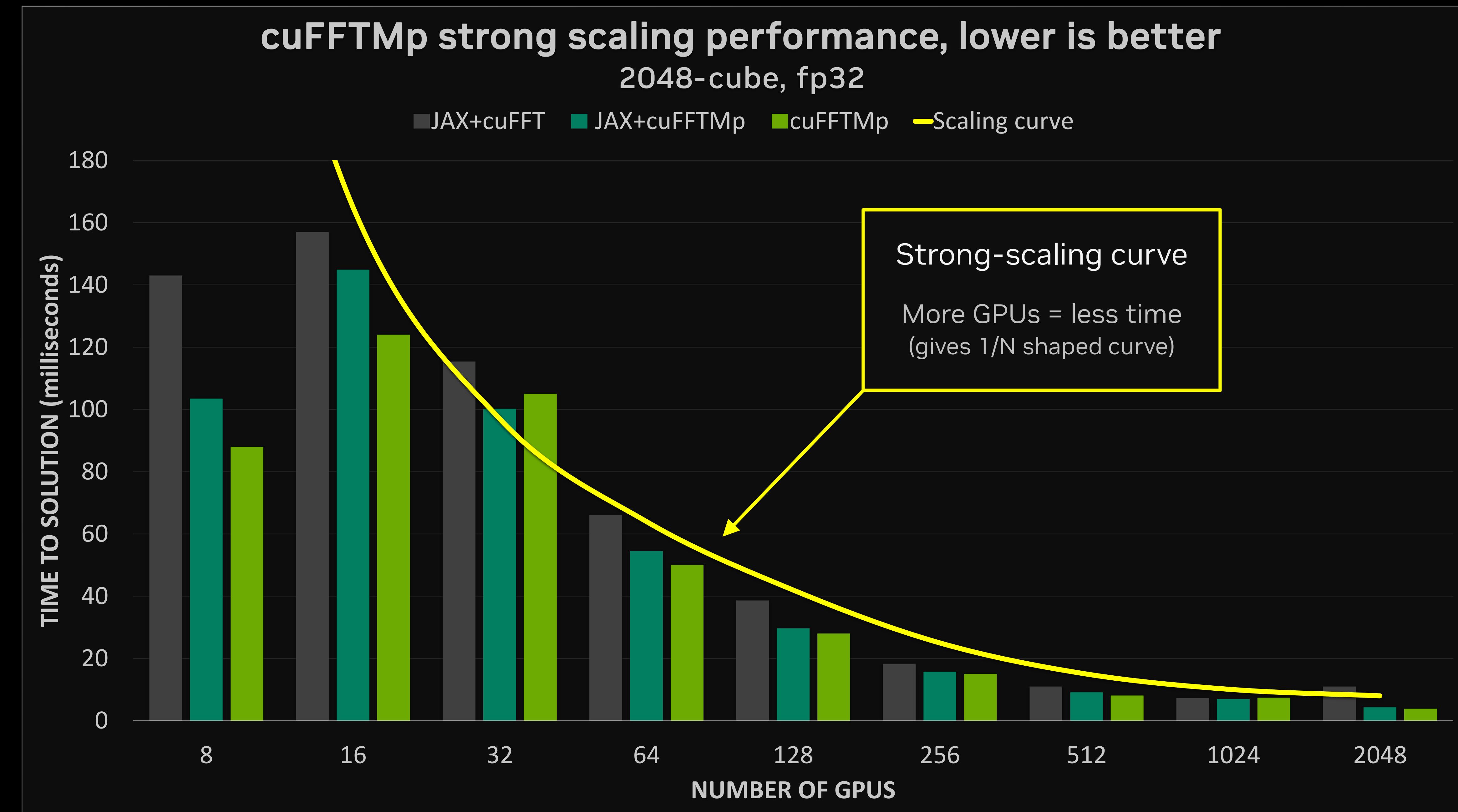
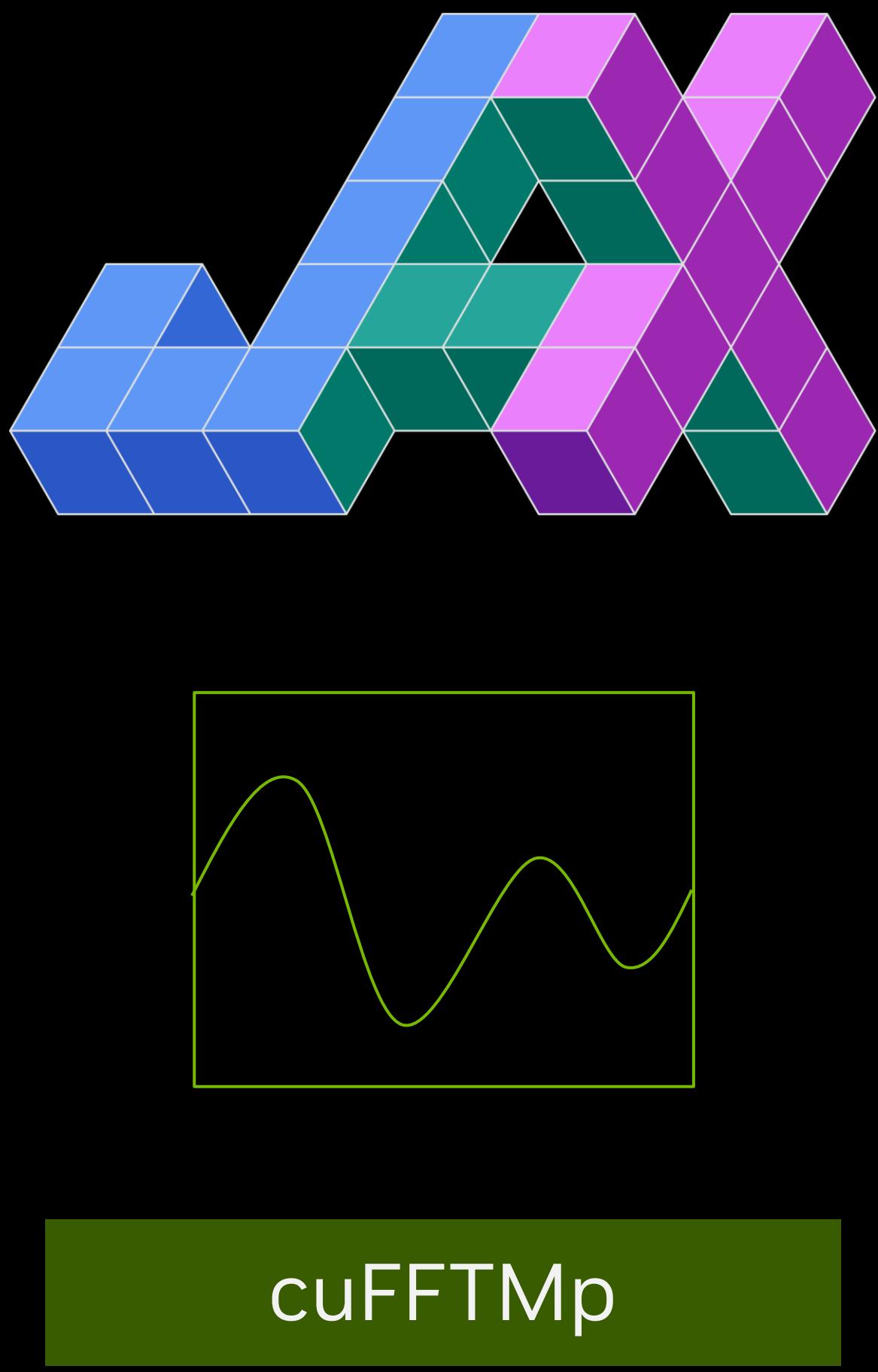
Comparison of ELPA1 vs. cuSOLVERMp

ELPA1 cuSOLVERMp

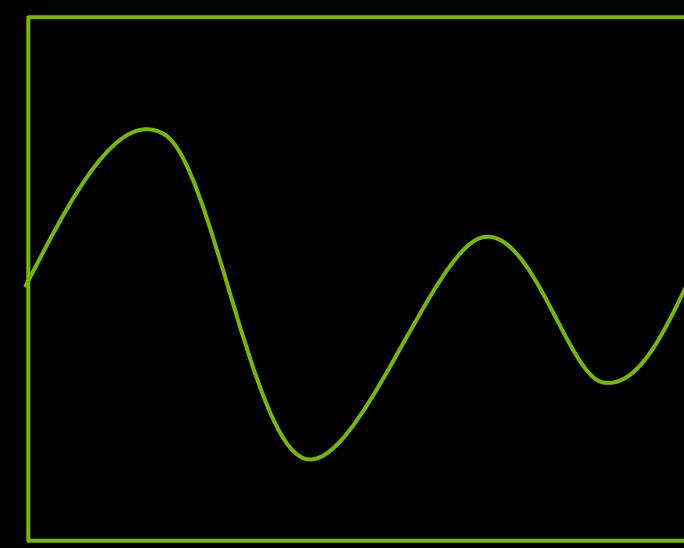
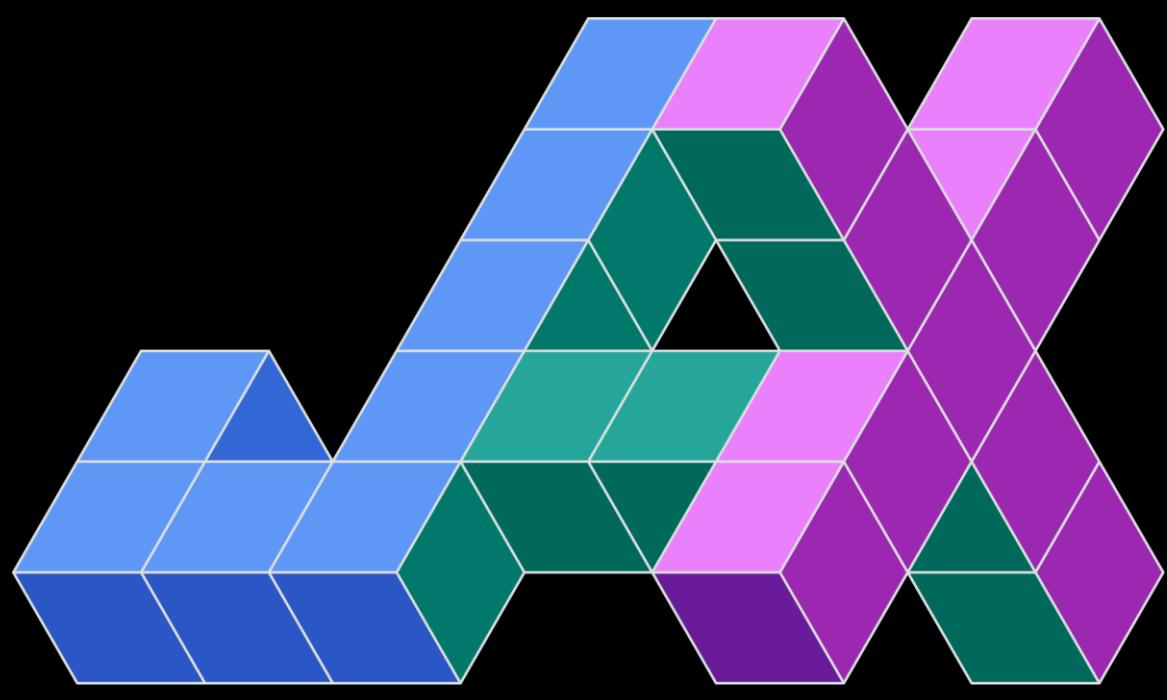
1.35x vs. ELPA1
on 1024 GPUs



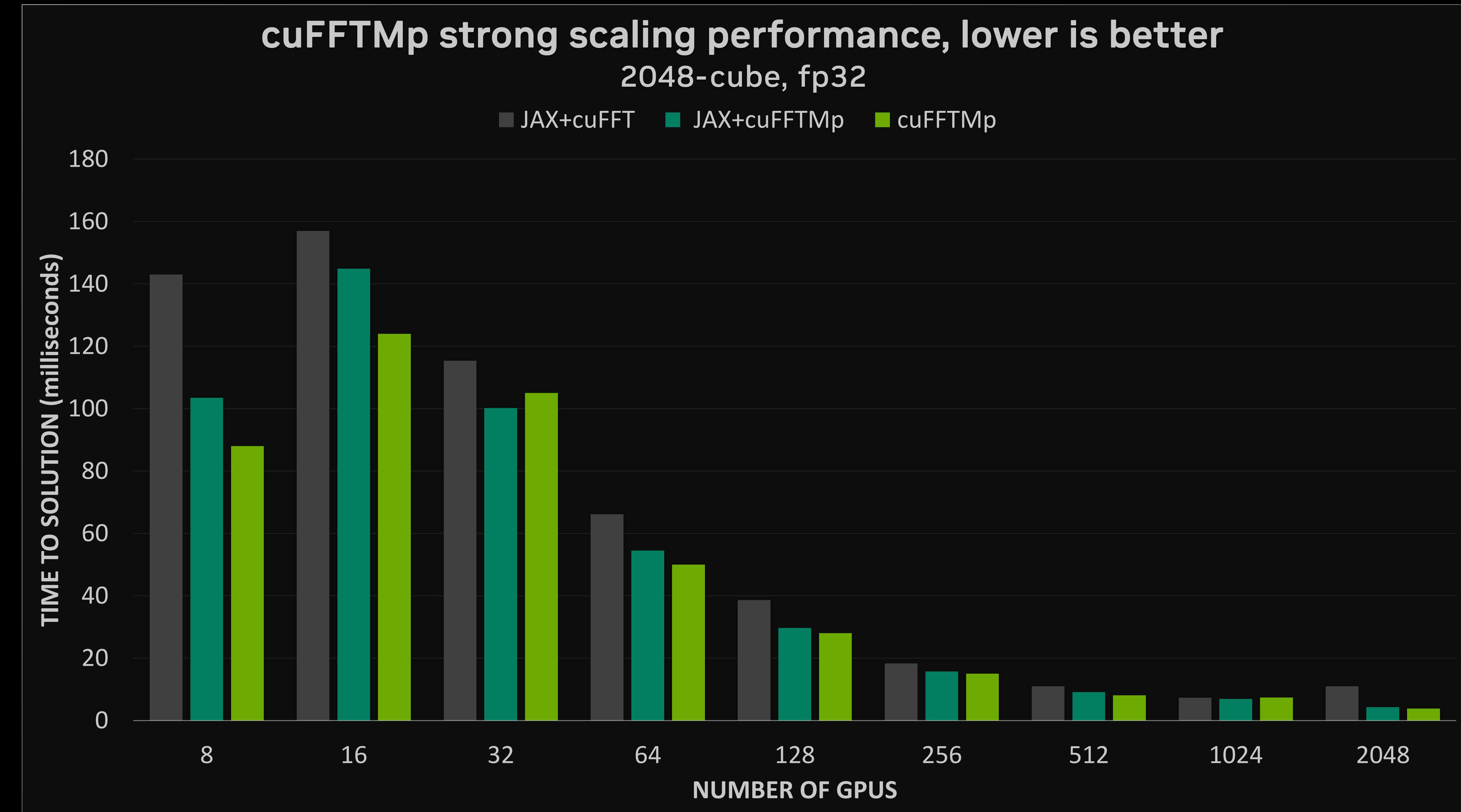
cuFFTMp + JAX



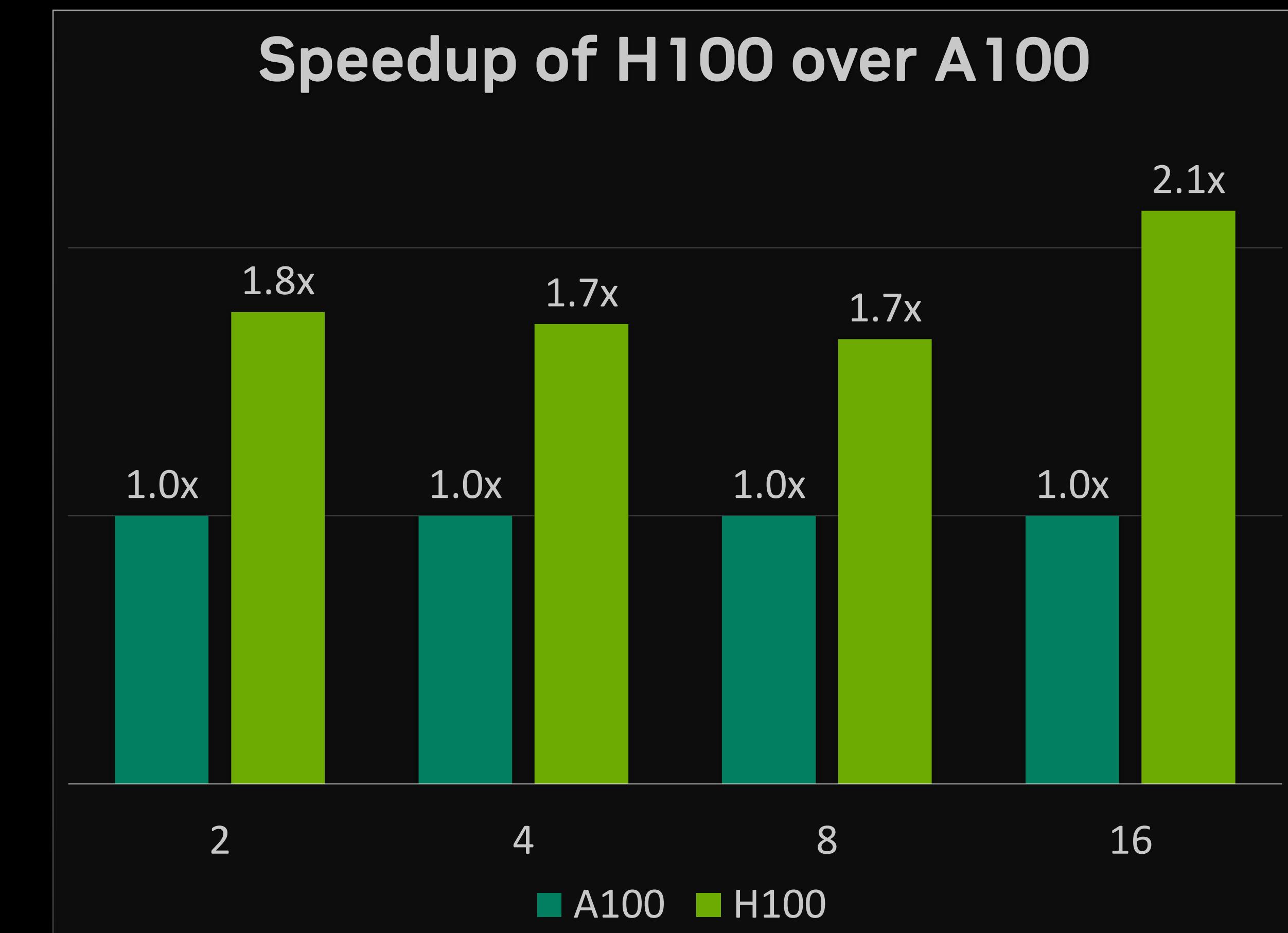
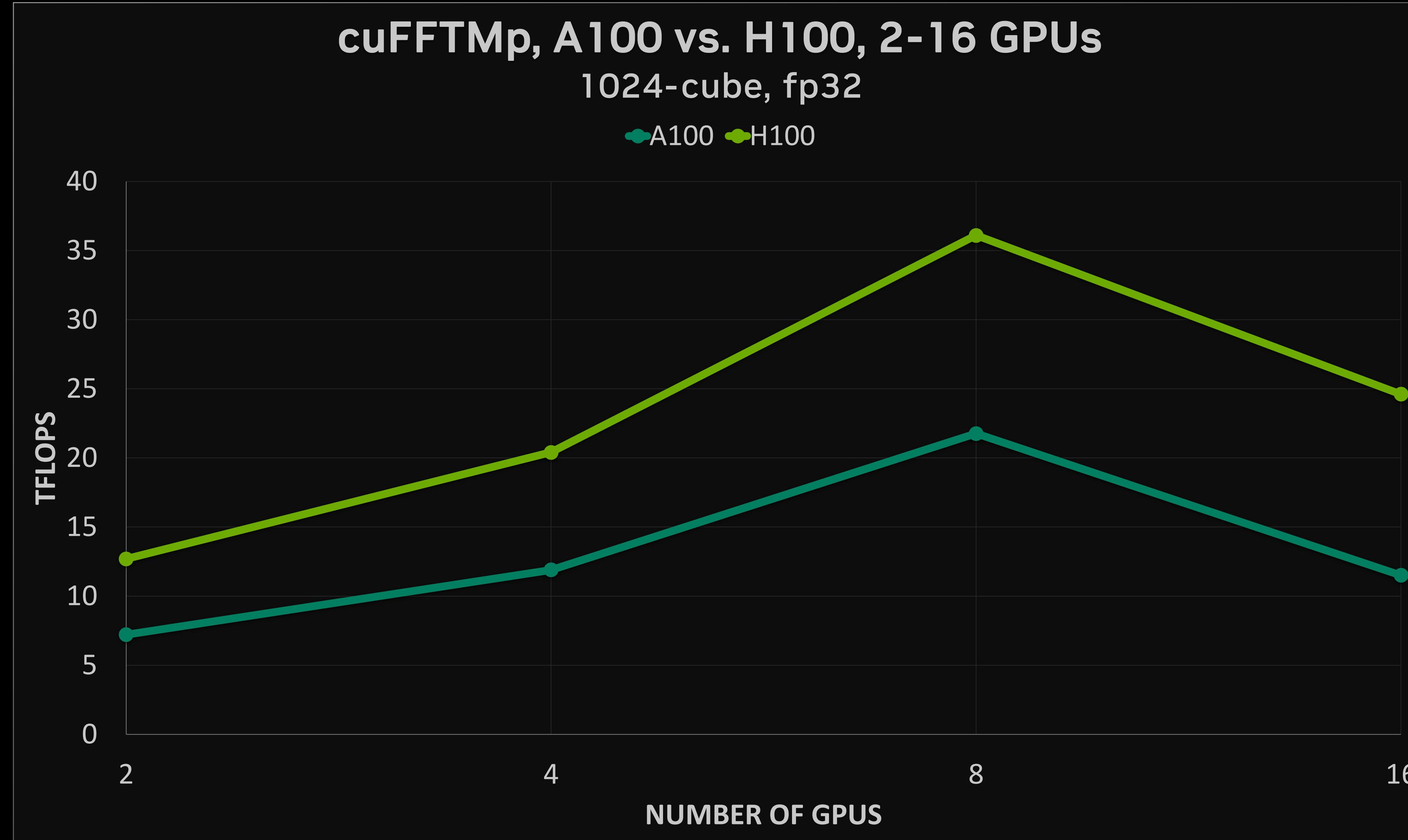
cuFFTMp + JAX



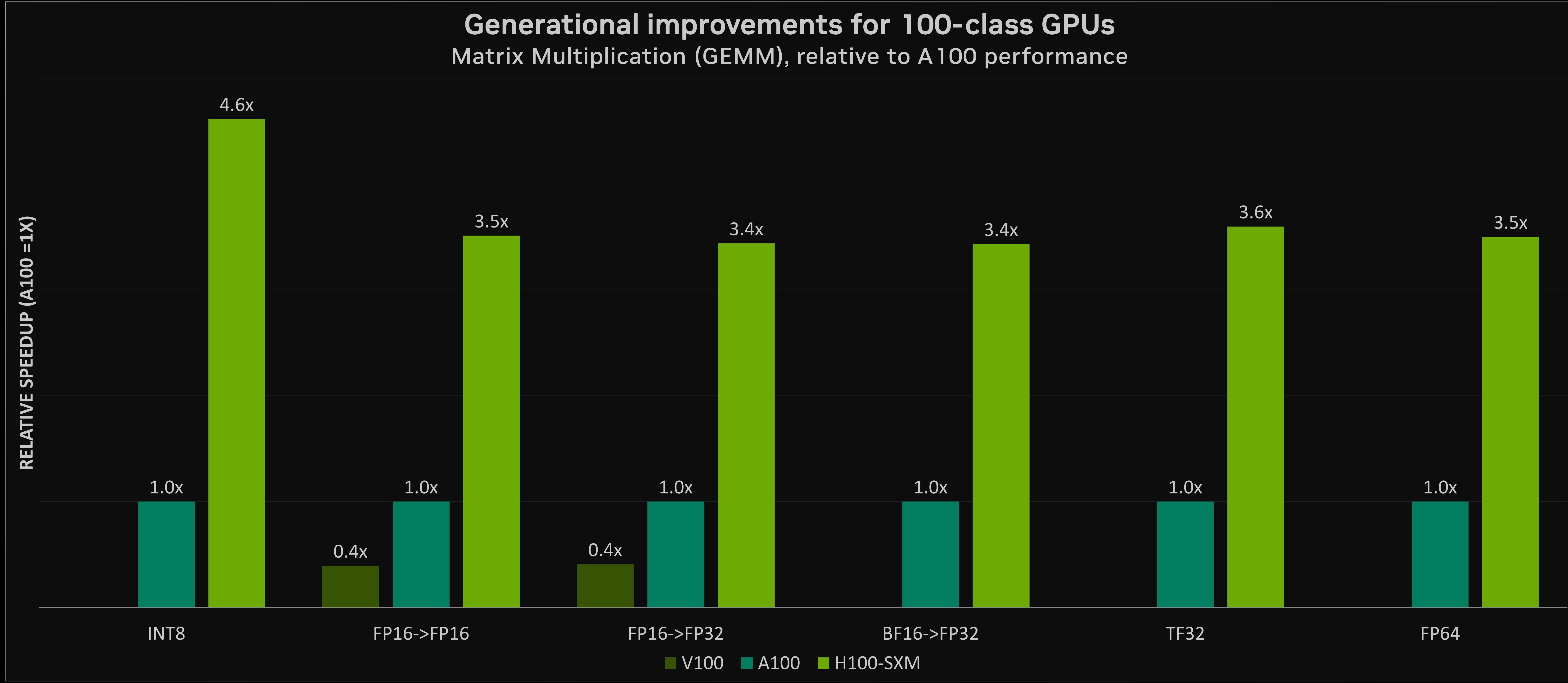
cuFFTMp



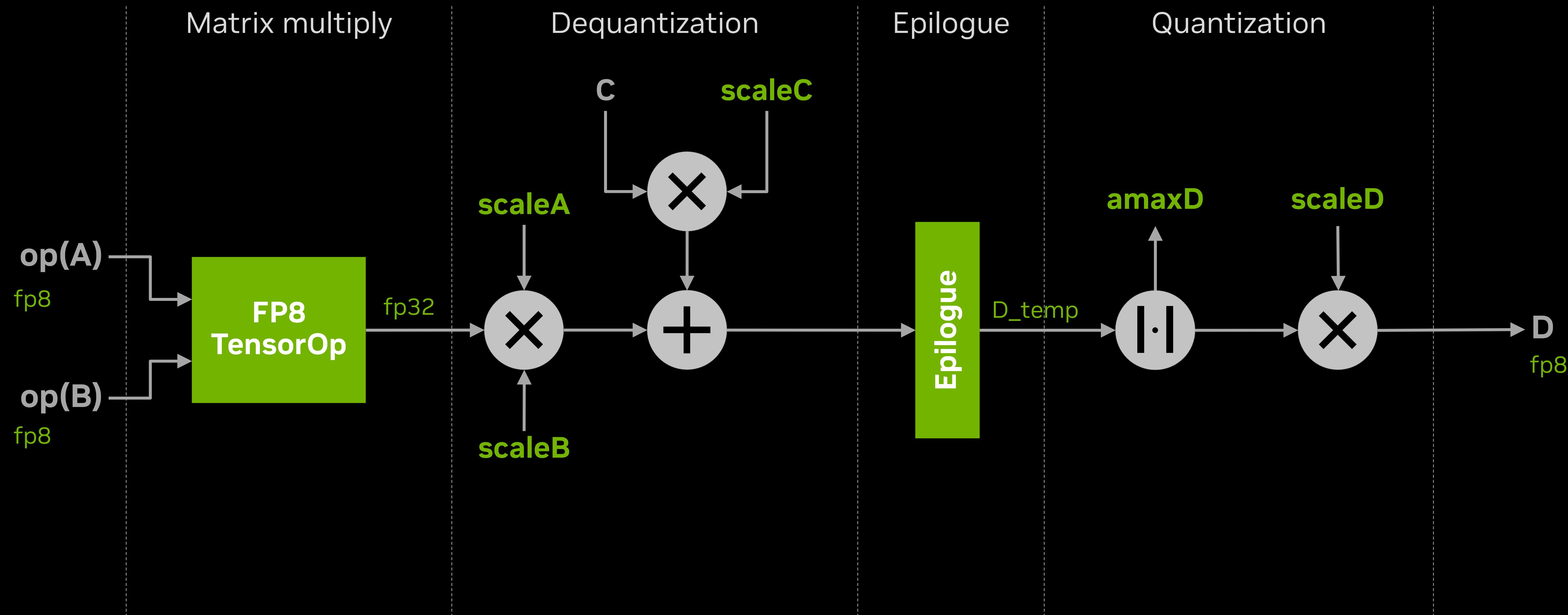
cuFFTMp Performance on H100 GPU vs. A100



Tensor Core Performance Across Generations

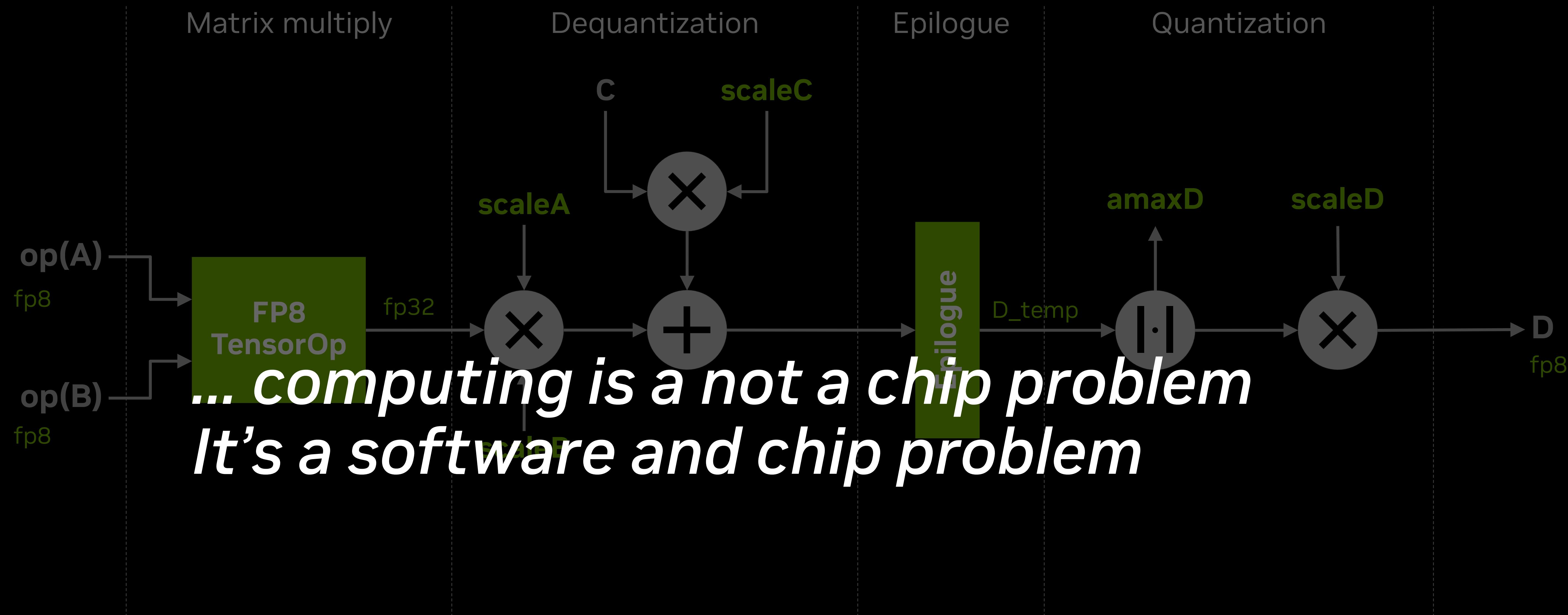


FP8 Matrix Multiplication in cuBLASLt



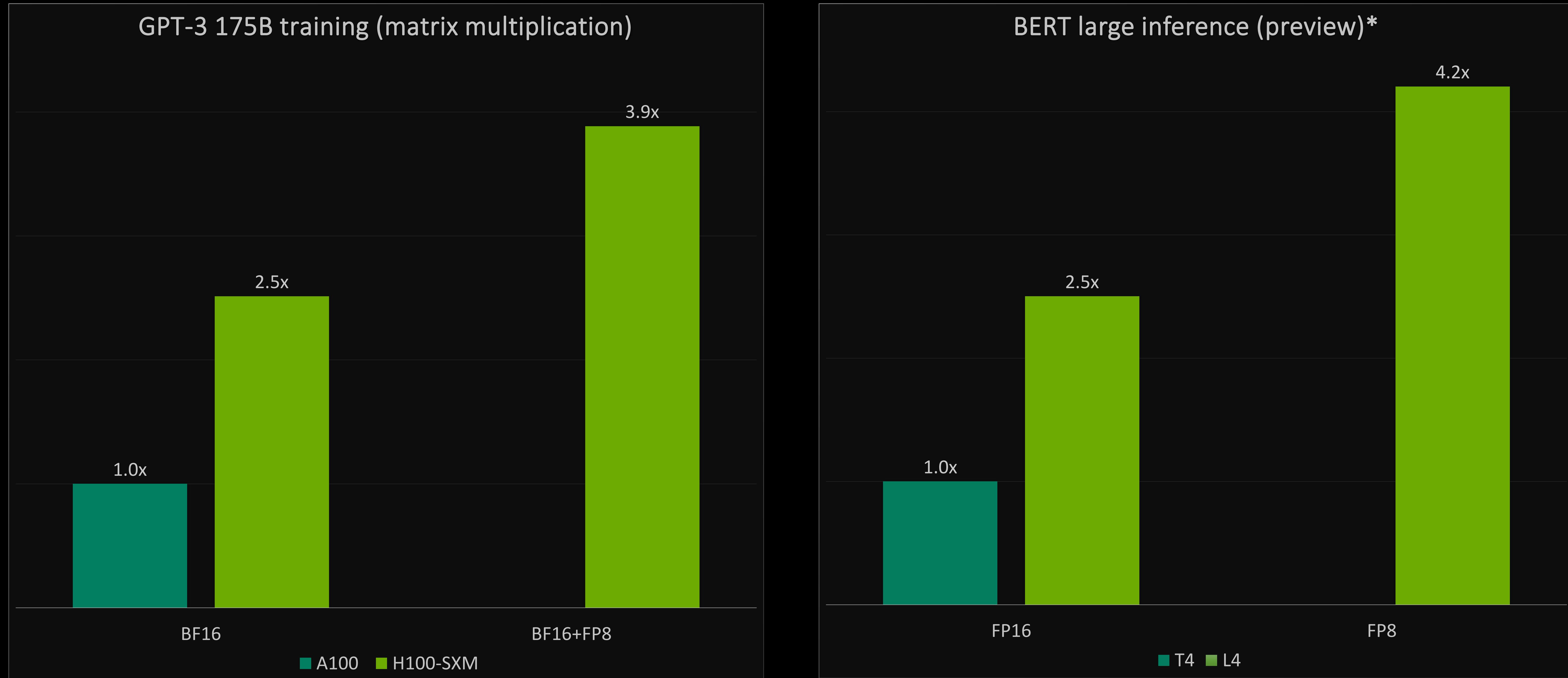
Block diagram of FP8 matrix multiplication & scaling
(can also directly output fp16, bf16 & fp32)

FP8 Matrix Multiplication in cuBLASLt



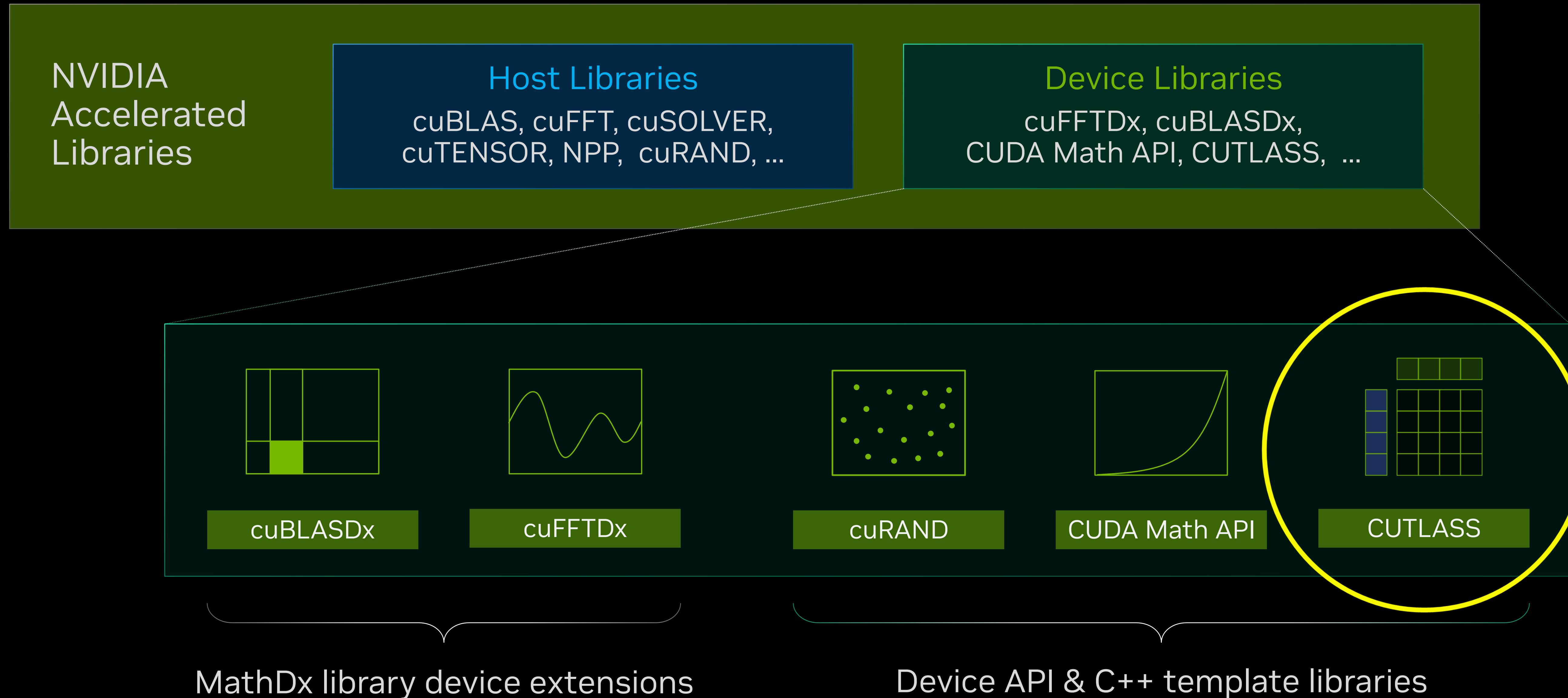
Block diagram of FP8 matrix multiplication & scaling
(can also directly output fp16, bf16 & fp32)

cuBLAS FP8 on NVIDIA H100 and L4 GPUs



* Average over a range of batch sizes and average sequence lengths

Device Libraries



[Download MathDx Today](#)



CUTLASS 3.0

Tensor Core programming model for linear algebra kernels in CUDA

CUTLASS 3.0

New major version of CUTLASS re-written with **CuTe**

Tiled MMA and Copy micro-kernels represent all GPU arch inner loops

Freely compose micro-kernels with any outer loop schedule

Ideal abstraction for custom kernel development targeting Hopper

Hopper support

Optimal computations for 4th generation Tensor Cores in H100

GEMM kernels targeting Hopper WGMMA + TMA + Threadblock Clusters

Warp-specialized persistent grid scheduling for peak utilization

	0	1	2	3	4	5	6	7
0	T0 V0	T4 V0	T8 V0	T12 V0	T16 V0	T20 V0	T24 V0	T28 V0
1	T1 V0	T5 V0	T9 V0	T13 V0	T17 V0	T21 V0	T25 V0	T29 V0
2	T2 V0	T6 V0	T10 V0	T14 V0	T18 V0	T22 V0	T26 V0	T30 V0
3	T3 V0	T7 V0	T11 V0	T15 V0	T19 V0	T23 V0	T27 V0	T31 V0
4	T0 V1	T4 V1	T8 V1	T12 V1	T16 V1	T20 V1	T24 V1	T28 V1
5	T1 V1	T5 V1	T9 V1	T13 V1	T17 V1	T21 V1	T25 V1	T29 V1
6	T2 V1	T6 V1	T10 V1	T14 V1	T18 V1	T22 V1	T26 V1	T30 V1
7	T3 V1	T7 V1	T11 V1	T15 V1	T19 V1	T23 V1	T27 V1	T31 V1
	0	1	2	3	4	5	6	7
0	T0 V0	T1 V0	T2 V0	T3 V0	T0 V2	T1 V2	T2 V2	T3 V2
1	T4 V0	T5 V0	T6 V0	T7 V0	T4 V2	T5 V2	T6 V2	T7 V2
2	T8 V0	T9 V0	T10 V0	T11 V0	T8 V2	T9 V2	T10 V2	T11 V2
3	T12 V0	T13 V0	T14 V0	T15 V0	T12 V2	T13 V2	T14 V2	T15 V2
4	T16 V0	T17 V0	T18 V0	T19 V0	T16 V2	T17 V2	T18 V2	T19 V2
5	T20 V0	T21 V0	T22 V0	T23 V0	T20 V2	T21 V2	T22 V2	T23 V2
6	T24 V0	T25 V0	T26 V0	T27 V0	T24 V2	T25 V2	T26 V2	T27 V2
7	T28 V0	T29 V0	T30 V0	T31 V0	T28 V2	T29 V2	T30 V2	T31 V2
8	T0 V1	T1 V1	T2 V1	T3 V1	T0 V3	T1 V3	T2 V3	T3 V3
9	T4 V1	T5 V1	T6 V1	T7 V1	T4 V3	T5 V3	T6 V3	T7 V3
10	T8 V1	T9 V1	T10 V1	T11 V1	T8 V3	T9 V3	T10 V3	T11 V3
11	T12 V1	T13 V1	T14 V1	T15 V1	T12 V3	T13 V3	T14 V3	T15 V3
12	T16 V1	T17 V1	T18 V1	T19 V1	T16 V3	T17 V3	T18 V3	T19 V3
13	T20 V1	T21 V1	T22 V1	T23 V1	T20 V3	T21 V3	T22 V3	T23 V3
14	T24 V1	T25 V1	T26 V1	T27 V1	T24 V3	T25 V3	T26 V3	T27 V3
15	T28 V1	T29 V1	T30 V1	T31 V1	T28 V3	T29 V3	T30 V3	T31 V3
	0	1	2	3	4	5	6	7
0	T0 V0	T0 V1	T1 V0	T1 V1	T2 V0	T2 V1	T3 V0	T3 V1
1	T4 V0	T4 V1	T5 V0	T5 V1	T6 V0	T6 V1	T7 V0	T7 V1
2	T8 V0	T8 V1	T9 V0	T9 V1	T10 V0	T10 V1	T11 V0	T11 V1
3	T12 V0	T12 V1	T13 V0	T13 V1	T14 V0	T14 V1	T15 V0	T15 V1
4	T16 V0	T16 V1	T17 V0	T17 V1	T18 V0	T18 V1	T19 V0	T19 V1
5	T20 V0	T20 V1	T21 V0	T21 V1	T22 V0	T22 V1	T23 V0	T23 V1
6	T24 V0	T24 V1	T25 V0	T25 V1	T26 V0	T26 V1	T27 V0	T27 V1
7	T28 V0	T28 V1	T29 V0	T29 V1	T30 V0	T30 V1	T31 V0	T31 V1
	0	1	2	3	4	5	6	7
0	T0 V2	T0 V3	T1 V2	T1 V3	T2 V2	T2 V3	T3 V2	T3 V3
1	T4 V2	T4 V3	T5 V2	T5 V3	T6 V2	T6 V3	T7 V2	T7 V3
2	T8 V2	T8 V3	T9 V2	T9 V3	T10 V2	T10 V3	T11 V2	T11 V3
3	T12 V2	T12 V3	T13 V2	T13 V3	T14 V2	T14 V3	T15 V2	T15 V3
4	T16 V2	T16 V3	T17 V2	T17 V3	T18 V2	T18 V3	T19 V2	T19 V3
5	T20 V2	T20 V3	T21 V2	T21 V3	T22 V2	T22 V3	T23 V2	T23 V3
6	T24 V2	T24 V3	T25 V2	T25 V3	T26 V2	T26 V3	T27 V2	T27 V3
7	T28 V2	T28 V3	T29 V2	T29 V3	T30 V2	T30 V3	T31 V2	T31 V3

Layout of threads & data for a single Hopper double precision (FP64) tensor core operation

[Download CUTLASS 3.0 from the GitHub repo](#)



CuTe: Robust Representation for Thread and Data Layouts

New programming model and API

RowMajor
ColumnMajor
RowMajorInterleaved
ColumnMajorInterleaved
PitchLinear
TensorNCHW

VoltaTensorOpMultiplicandCongruous
ColumnMajorVoltaTensorOpMultiplicandCongruous
RowMajorVoltaTensorOpMultiplicandCongruous
VoltaTensorOpMultiplicandBCongruous
ColumnMajorVoltaTensorOpMultiplicandBCongruous
RowMajorVoltaTensorOpMultiplicandBCongruous
VoltaTensorOpMultiplicandCrosswise
ColumnMajorVoltaTensorOpMultiplicandCrosswise
RowMajorVoltaTensorOpMultiplicandCrosswise
Many, many more ...

CUTLASS 2.x Layout Definitions

CuTe Layout

Hierarchical and multi-dimensional layout representation
Formal algebra for manipulating layouts
Layouts for both threads and data
Used in CUTLASS 3.0 to build a unified GPU micro-kernel abstraction



Layout<Shape, Stride>

CUTLASS 3.0 **CuTe** Layout Definition

CUTLASS 3: Towards Better Developer and User Experience

Simplified but powerful kernel instantiation

Hopper-optimized GEMM mainloop

```
using CollectiveOp =
typename gemm::collective::CollectiveBuilder<
    arch::Sm90, arch::OpClassTensorOp,
    half_t, LayoutA, 8,
    half_t, LayoutB, 8,
    float,
    Shape<_128,_128,_64>, Shape<_1, 2,_1>,
    gemm::collective::StageCountAuto,
    gemm::collective::KernelScheduleAuto
>::CollectiveOp;
```

Hopper-optimized GEMM mainloop
with persistent scheduling and 5 stages

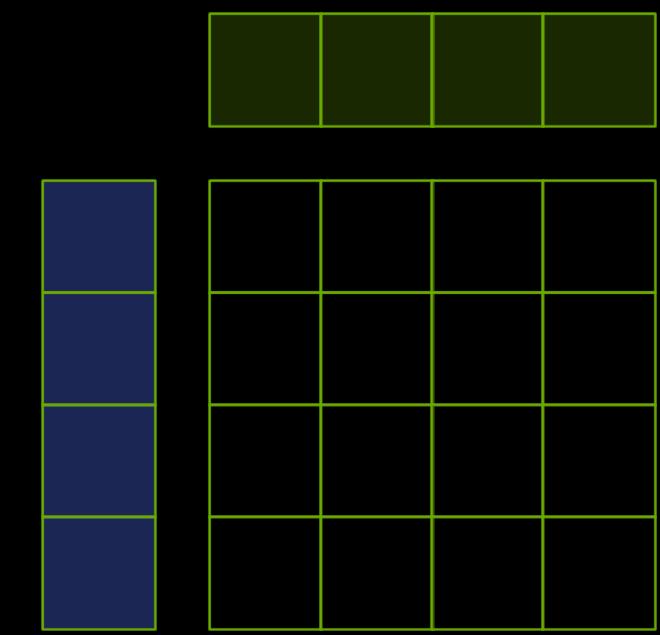
```
using CollectiveOp =
typename gemm::collective::CollectiveBuilder<
    arch::Sm90, arch::OpClassTensorOp,
    half_t, LayoutA, 8,
    half_t, LayoutB, 8,
    float,
    Shape<_128,_128,_64>, Shape<_1, 2,_1>,
    gemm::collective::StageCount<5>,
    gemm::KernelTmaWarpSpecializedPersistent
>::CollectiveOp;
```

CUTLASS 3: Towards Better Developer and User Experience

Python interface coming in CUTLASS 3.1 (experimental)

Simplified kernel declaration

```
plan = cutlass.GroupedGemm(  
    element=torch.float16,  
    layout=cutlass.LayoutType.RowMajor)
```



CUTLASS

Simplified porting to DL frameworks

```
cutlass.emit.pytorch(  
    plan.construct(), name='grouped_gemm', cc=80)
```

PY

C++

CU

```
> python setup.py install
```

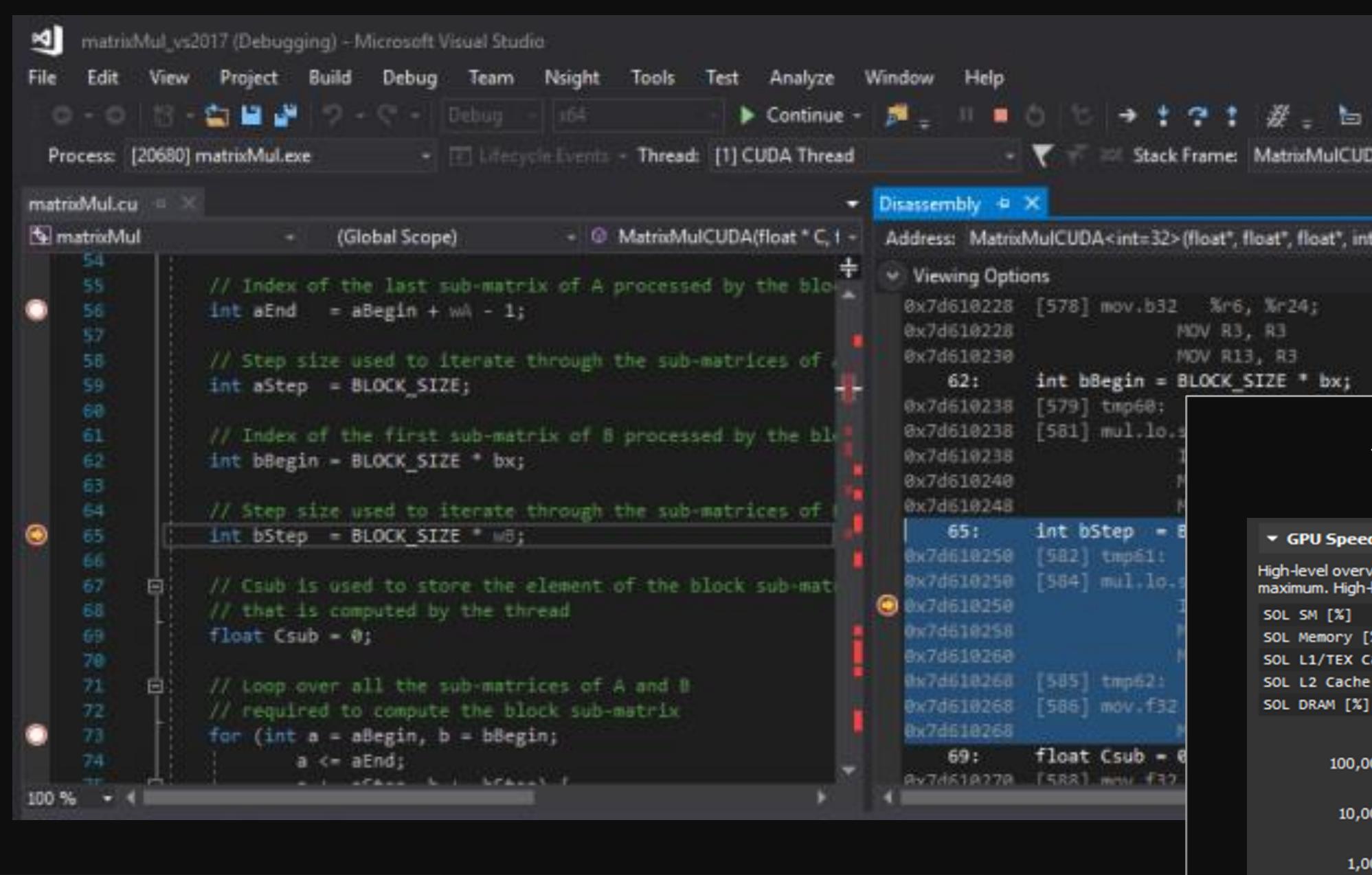
```
# Import and run generated PyTorch extension  
import grouped_gemm  
Ds = grouped_gemm.run([A0, A1],  
                      [B0, B1])
```

NVIDIA GPU Computing Tools

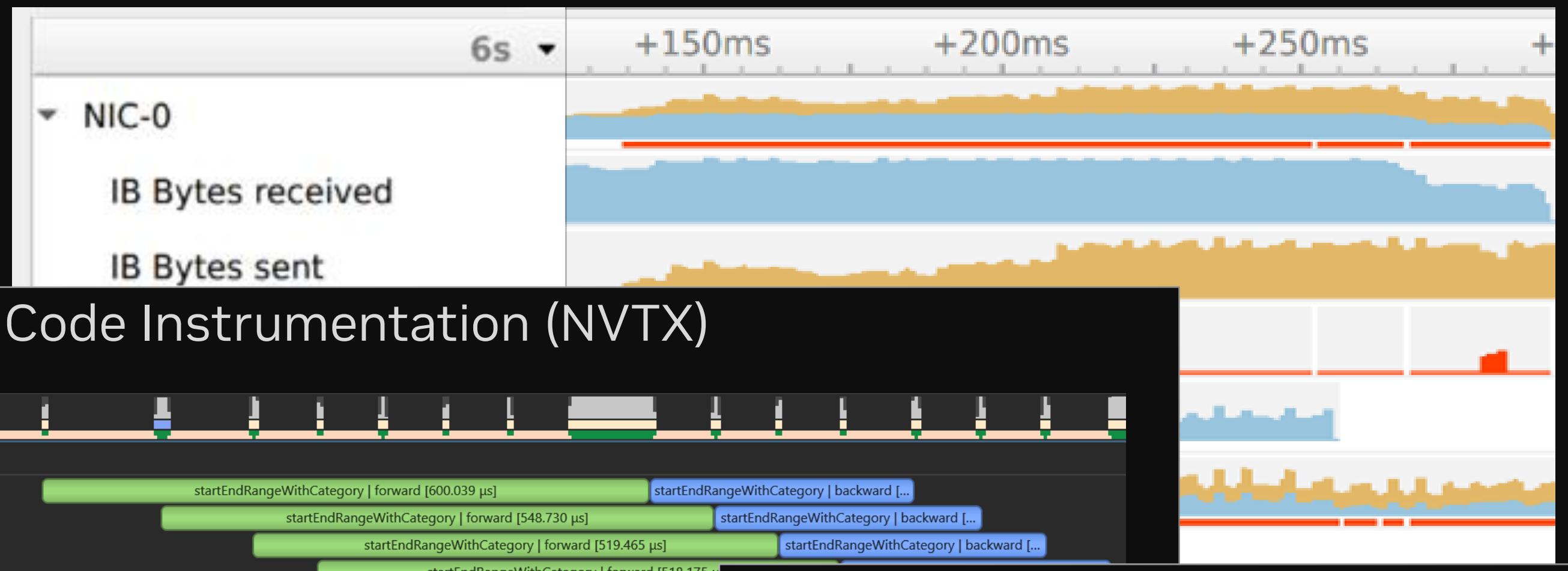
Correctness Tools

```
$ compute-sanitizer --leak-check full memcheck_demo  
===== COMPUTE-SANITIZER  
Mallocing memory  
Running unaligned_kernel  
Ran unaligned_kernel: no error  
Sync: no error  
Running out_of_bounds_kernel  
Ran out_of_bounds_kernel: no error  
Sync: no error
```

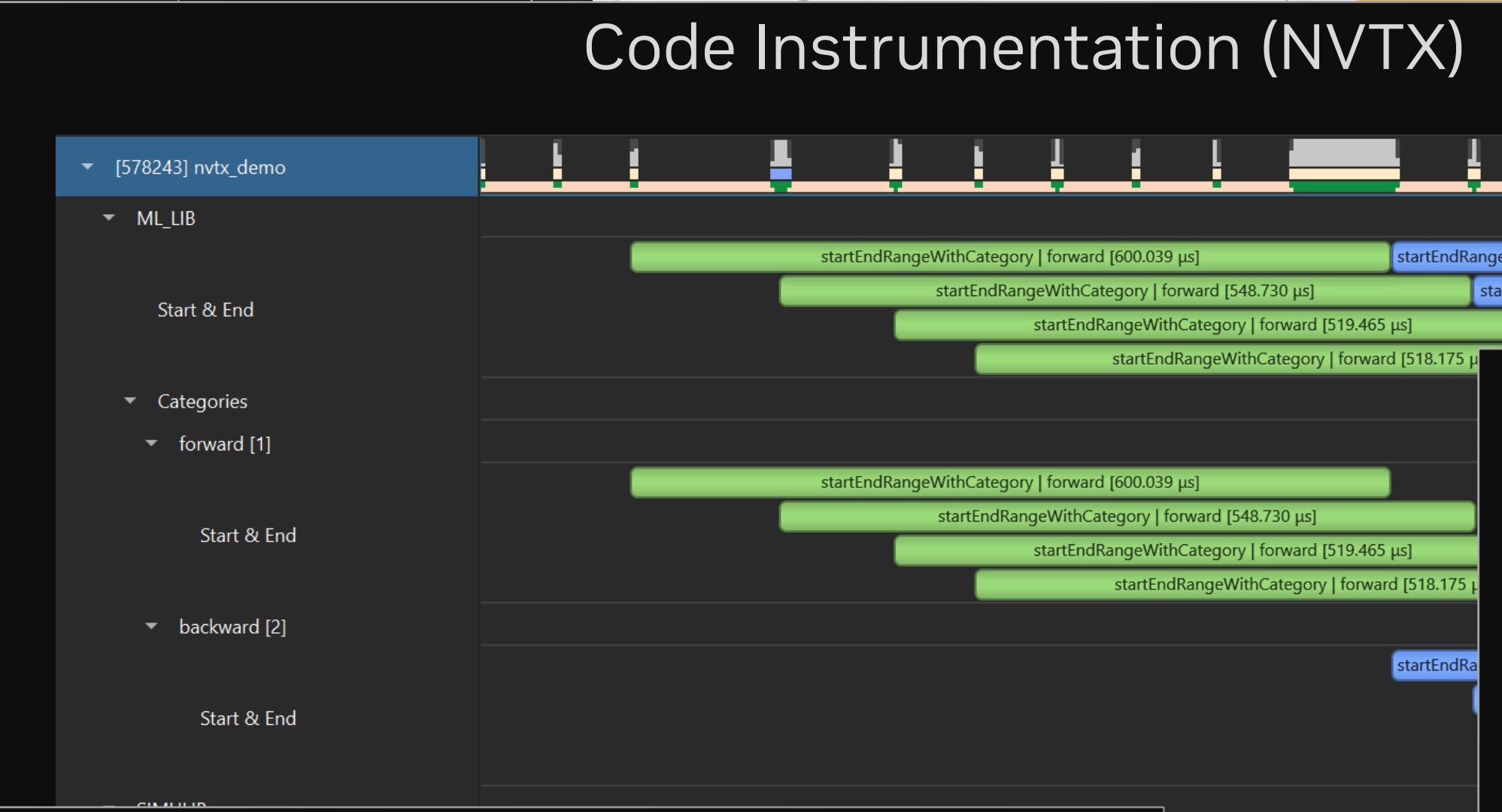
Debuggers, Profilers, IDE Integration



Networking & Communications Sampling



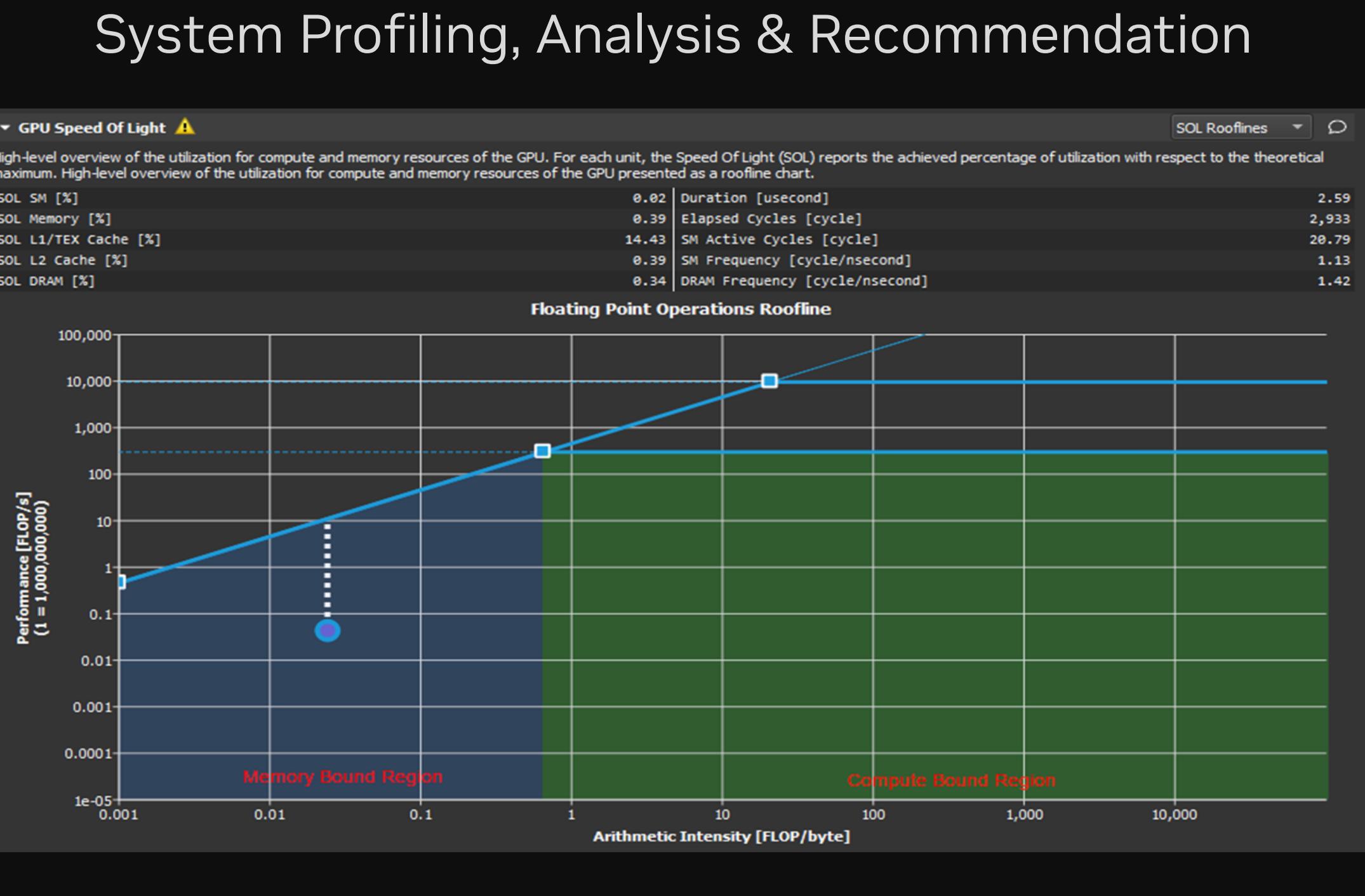
Code Instrumentation (NVTX)



Online Monitoring (DCGM)



System Profiling, Analysis & Recommendation



GPU Speed of Light
High-level overview of the utilization for compute and memory resources of the GPU. For each unit, the Speed Of Light (SOL) reports the achieved percentage of utilization with respect to the theoretical maximum. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

SOL Unit	Duration [usecond]	Elapsed Cycles [cycle]	SM Active Cycles [cycle]	SM Frequency [cycle/second]	DRAM Frequency [cycle/second]
SOL SM [%]	0.02	2,993	28.79	1.13	1.42
SOL Memory [%]	0.39				
SOL L1/TEX Cache [%]	14.43				
SOL L2 Cache [%]	0.39				
SOL DRAM [%]	0.34				

Floating Point Operations Roofline

Performance [FLOP/s] (1 = 1,000,000,000)

Arithmetic Intensity [FLOP/byte]

Memory Bound Region

Compute Bound Region

From the Macro to the Micro: CUDA Developer Tools Find and Fix Problems at Any Scale [S51205]



Nsight Systems Multi-Node Analysis

Literally tools for your tools

Recipes to analyze existing reports

Collect reports from multiple sources/nodes

Reduces “getting started” time

Analyze thousands of reports

Scale-out processing

Cluster-wide conclusions

Recipes generate results various formats

Jupyter Notebook

Parquet, CSV, ...

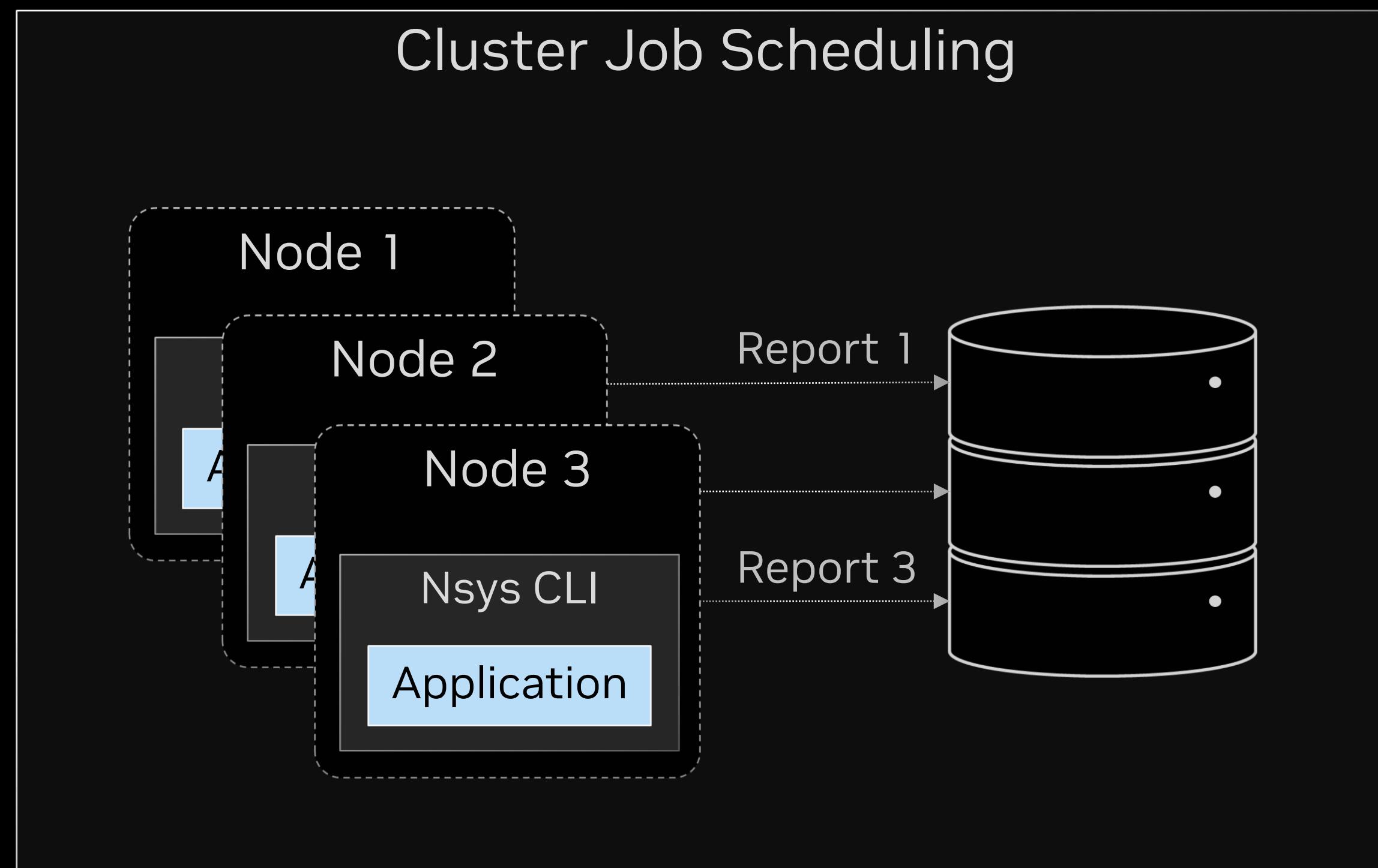


Multi-Report Analysis Workflow

Preview in Nsight Systems 2023.2



Select or create custom recipes
Recipes built on extensible Python library
Set recipe parameters & arguments

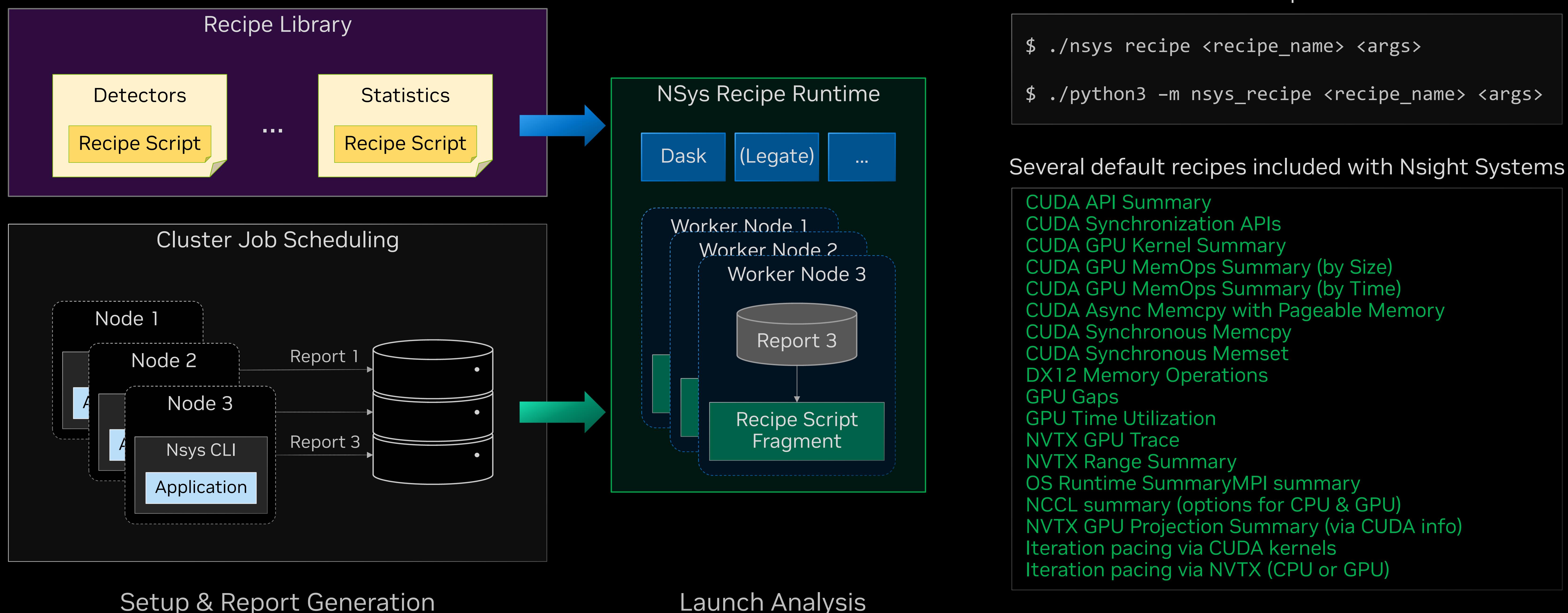


Generate reports per-node / per-rank as usual
Data is output to storage during run
Can schedule for analysis immediately or later

Setup & Report Generation

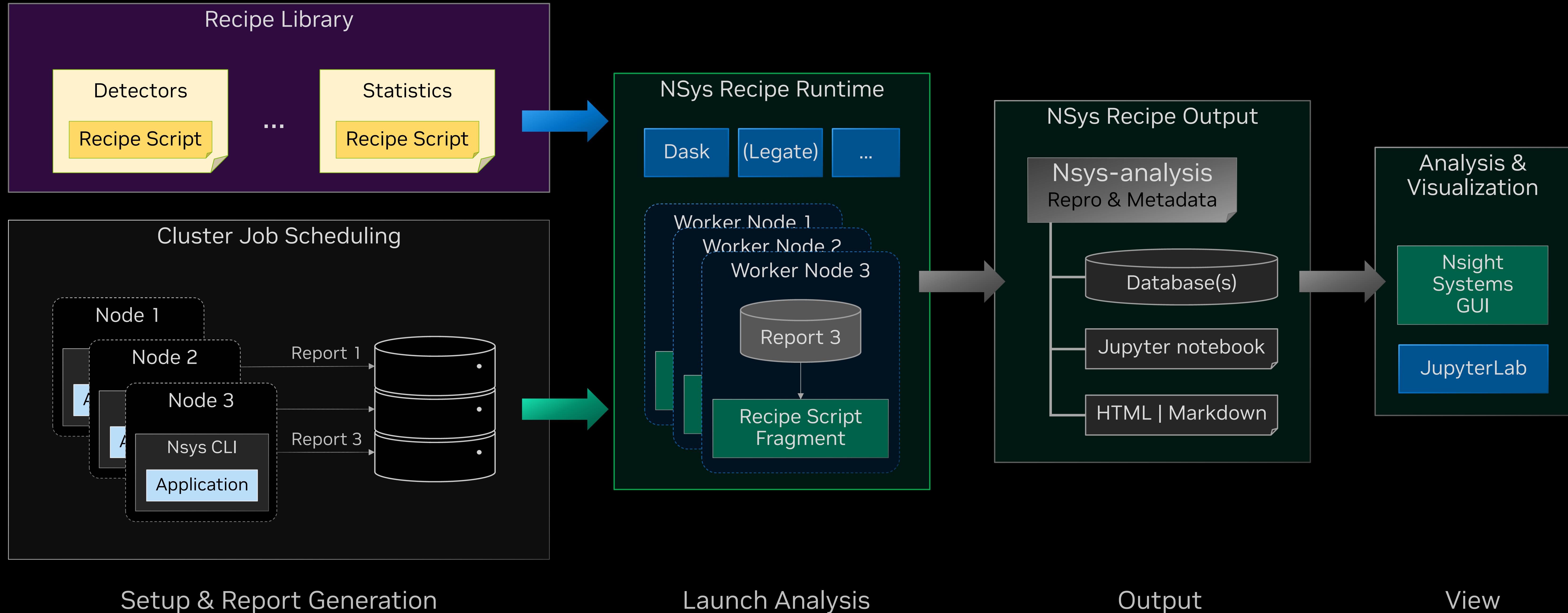
Multi-Report Analysis Workflow

Preview in Nsight Systems 2023.2

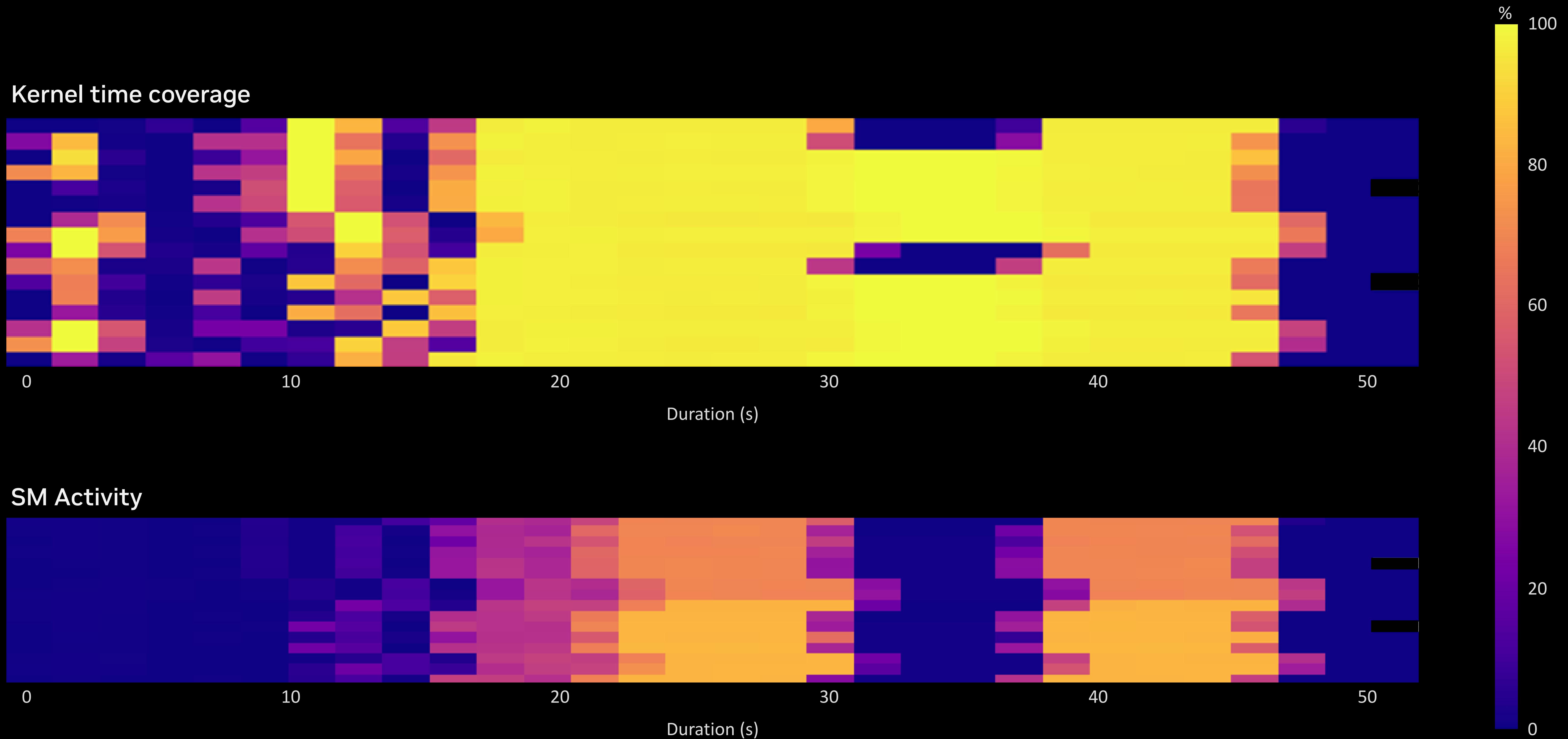


Multi-Report Analysis Workflow

Preview in Nsight Systems 2023.2



Example: Heatmap Analysis



Datacenter Scale Computing



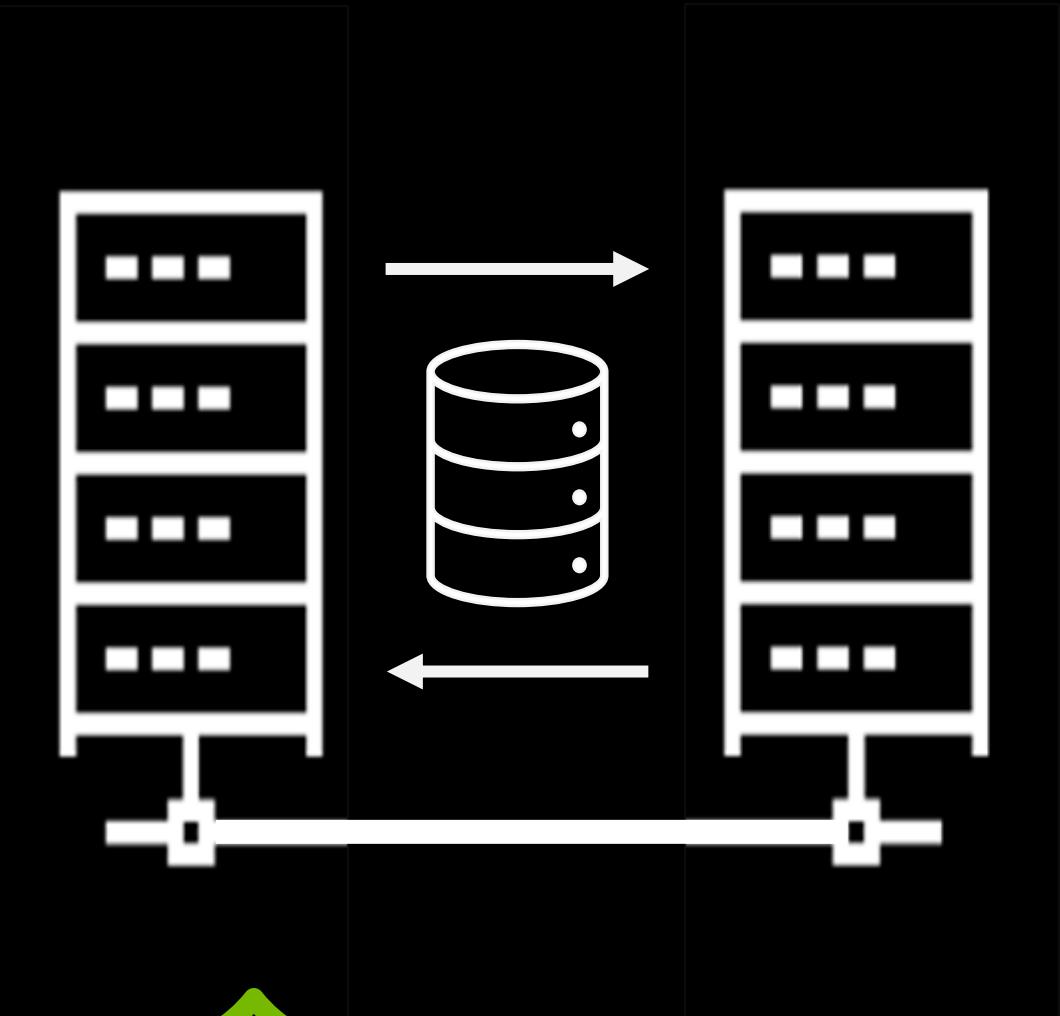
Datacenter Data Security

Data at Rest
(In Storage)



 STRONG PROTECTION

Data in Transit
(Across a Network)



 STRONG PROTECTION

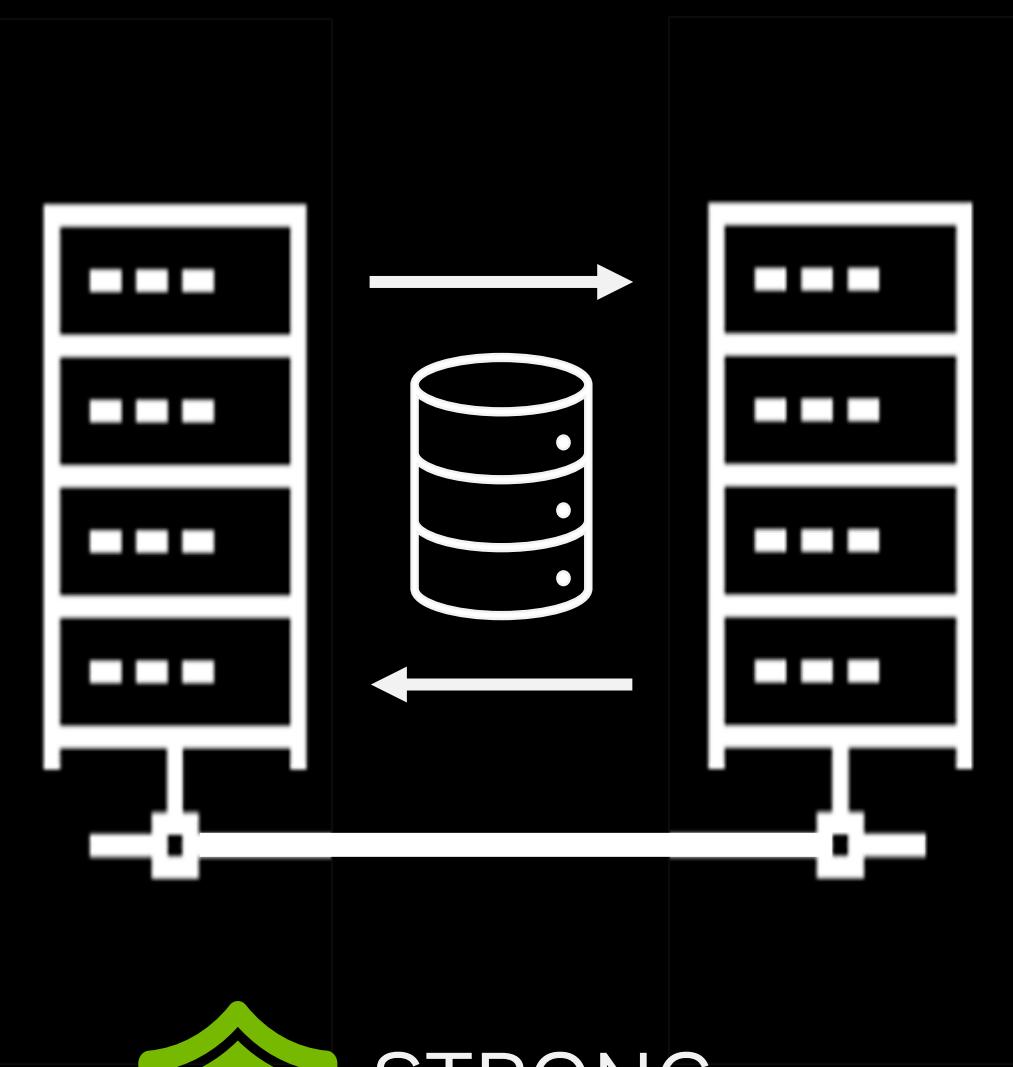
Datacenter Data Security

Data at Rest
(In Storage)



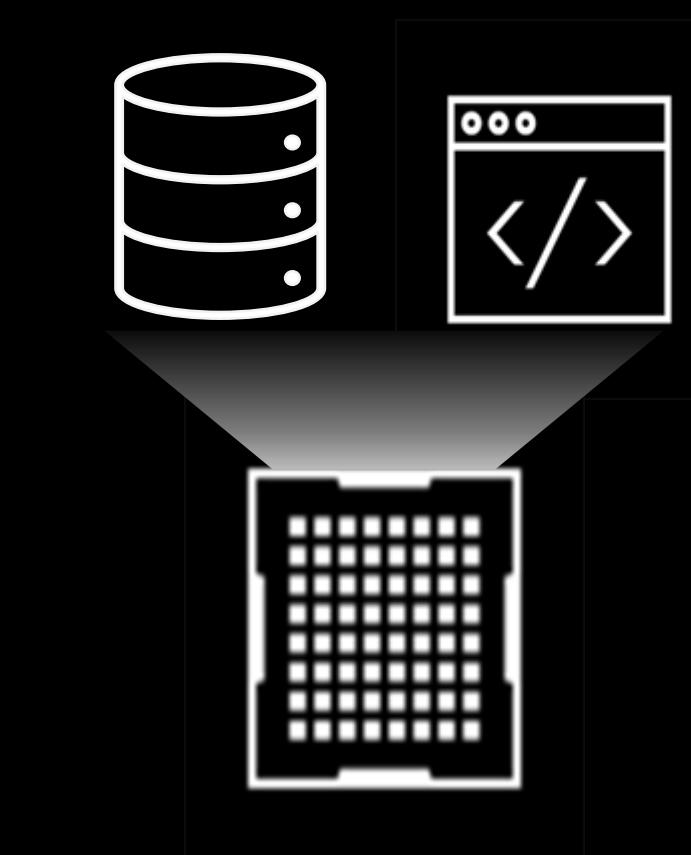
STRONG PROTECTION

Data in Transit
(Across a Network)

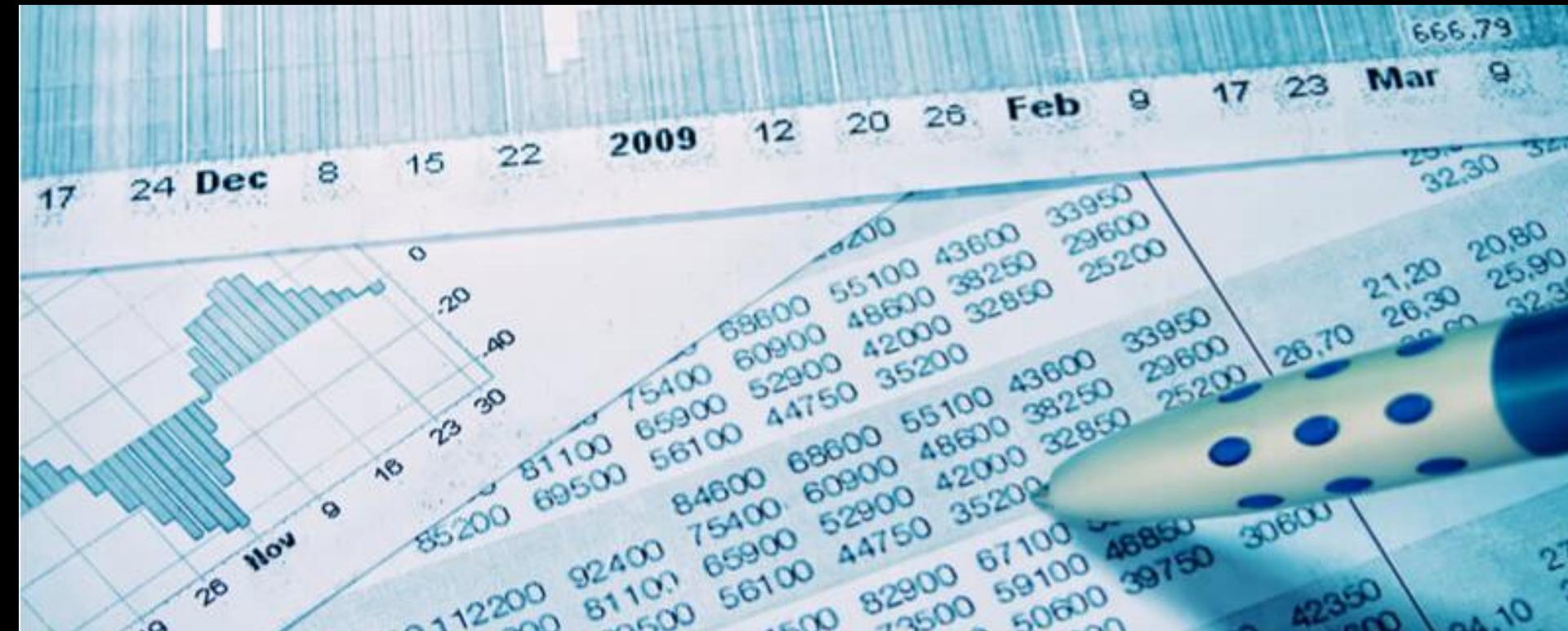


STRONG PROTECTION

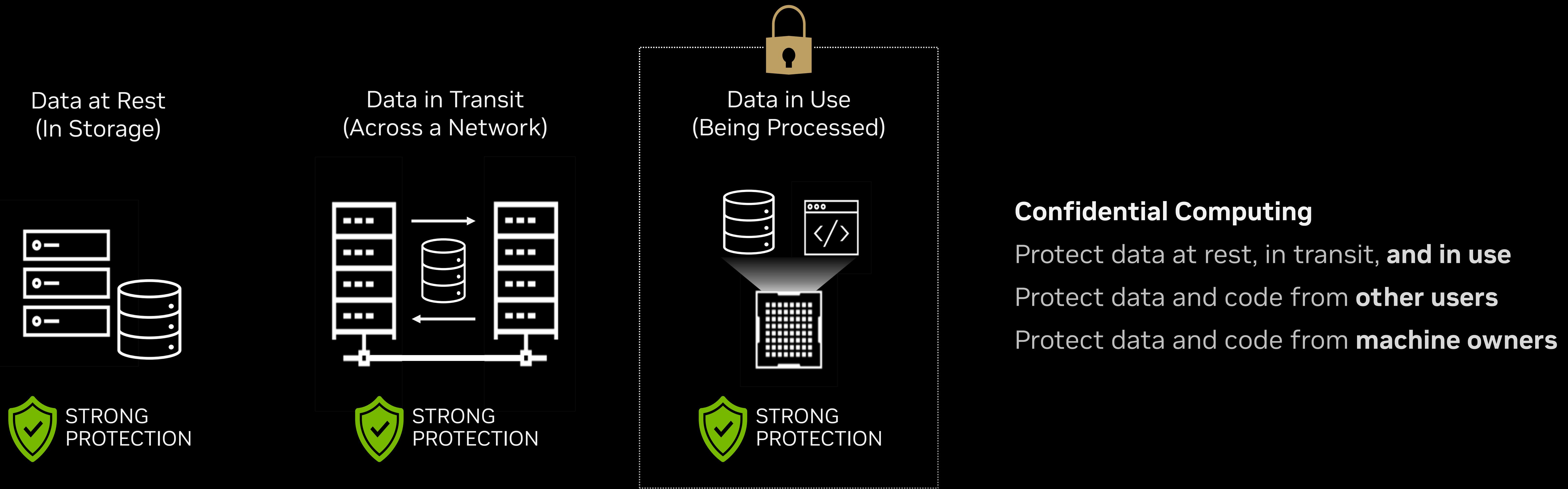
Data in Use
(Being Processed)

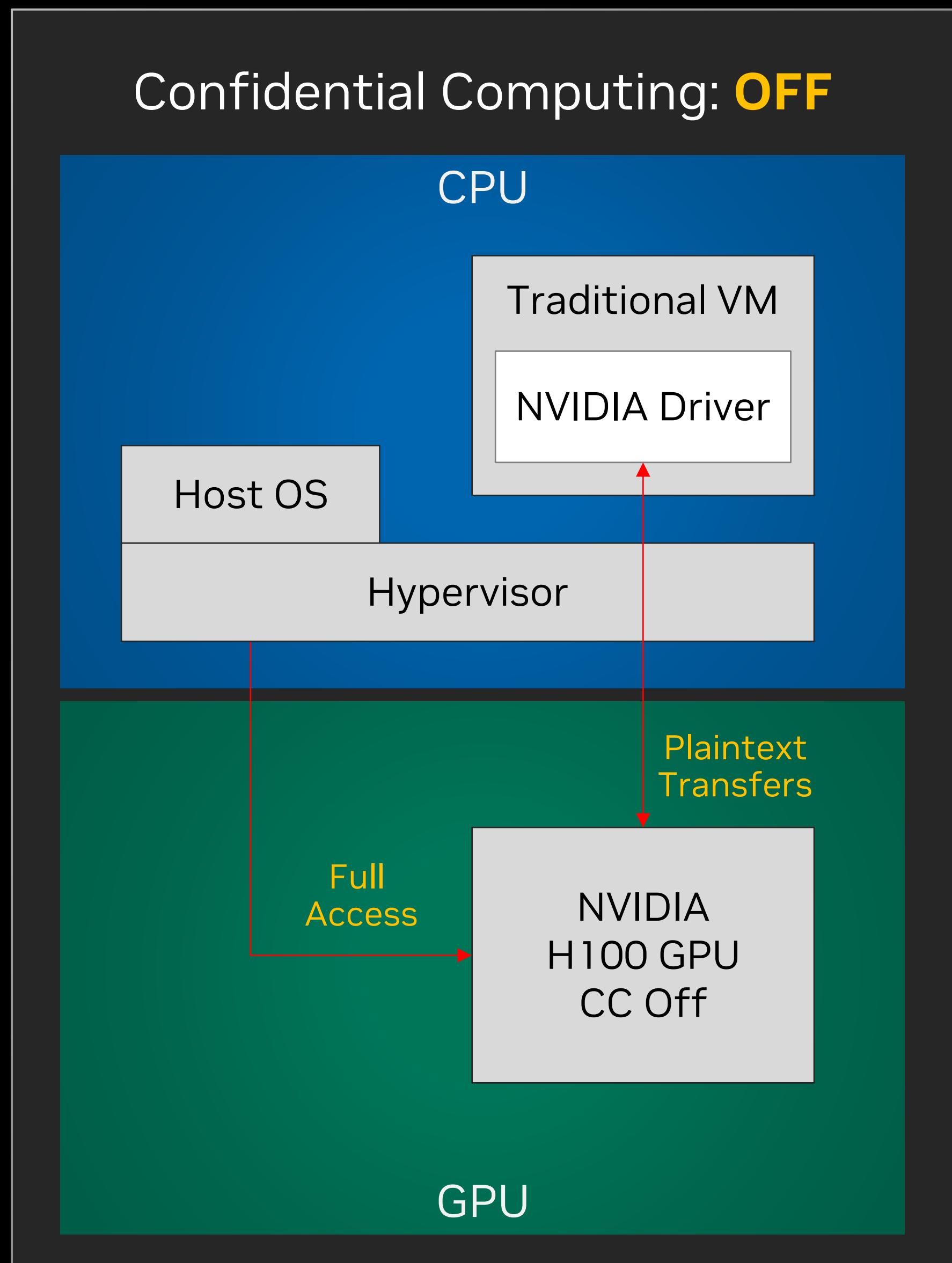


WEAK PROTECTION



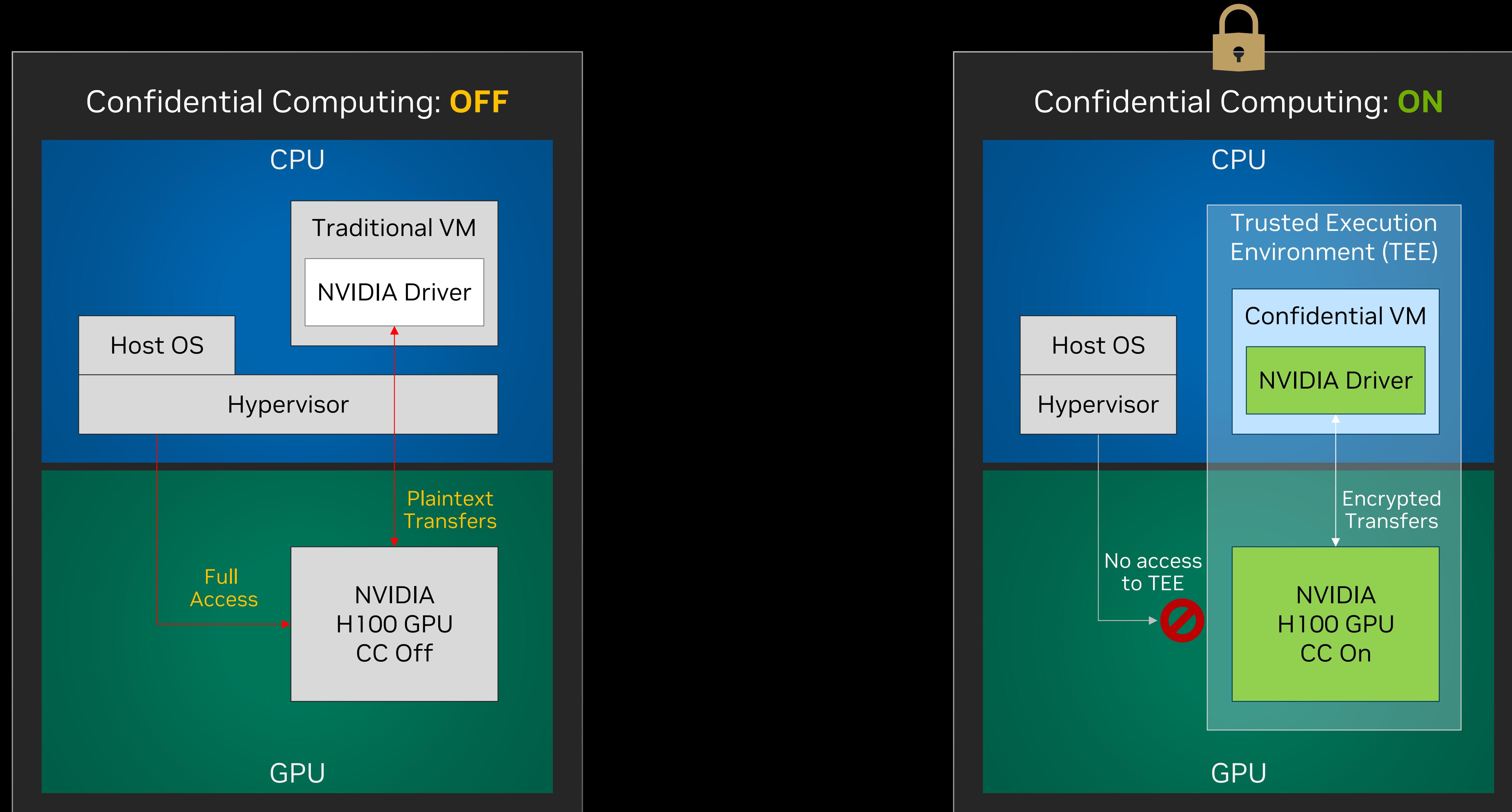
What is Confidential Computing?





Hypervisor has full access to all of system memory and all of GPU memory

GPU With VM-Based Confidential Computing



Hypervisor has full access to all of system memory and all of GPU memory

Hypervisor is blocked from accessing the Confidential VM in system memory and blocked from reading GPU memory

GPU With VM-Based Confidential Computing

Applications run unchanged in confidential mode

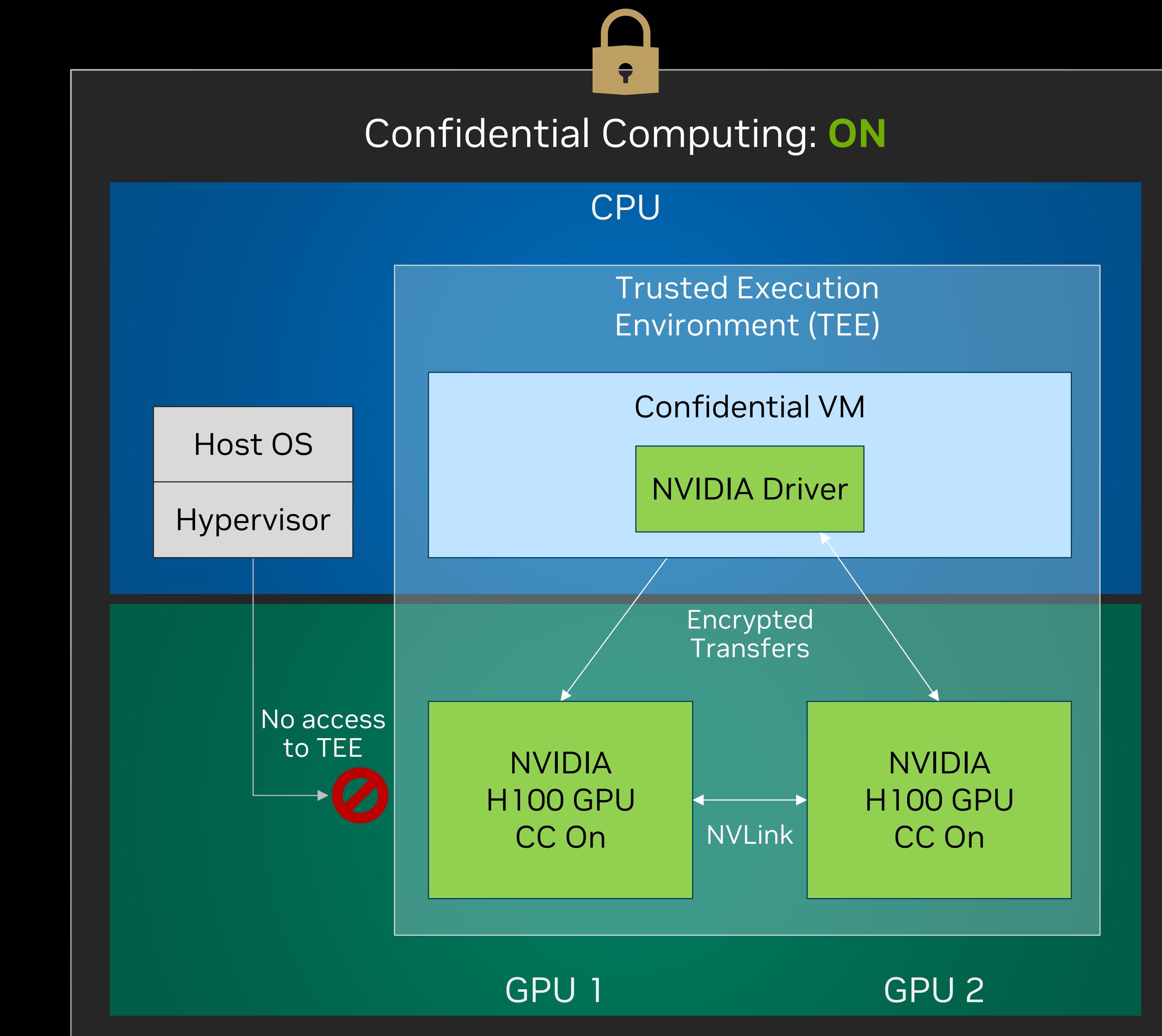
NVIDIA drivers, libraries, APIs run in the VM

VM-based TEE protects the entire workload

Each GPU is passed through to just one VM

VMs protected from each other

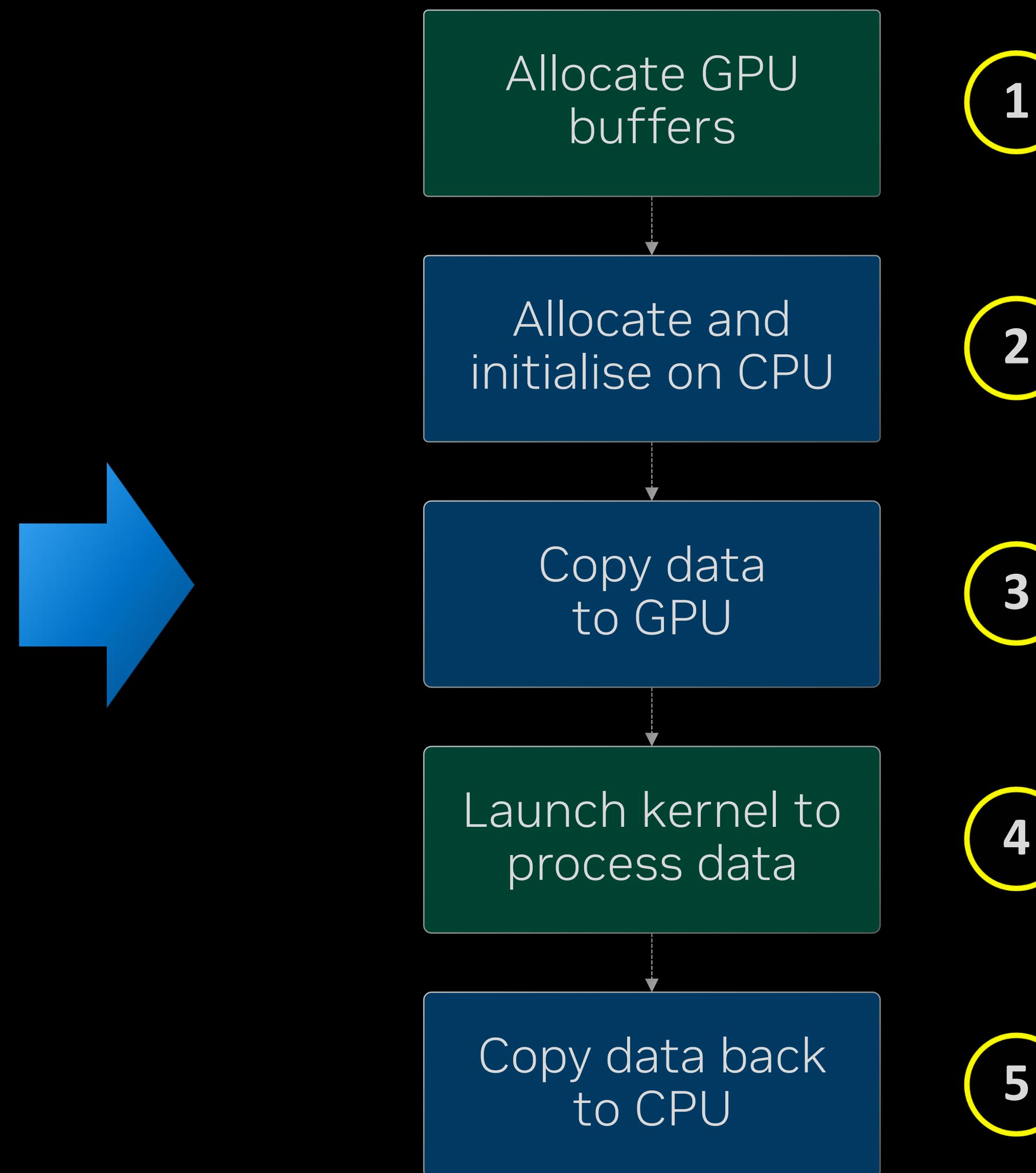
One VM may have multiple GPUs



Two NVLink-connected GPUs within the same TEE

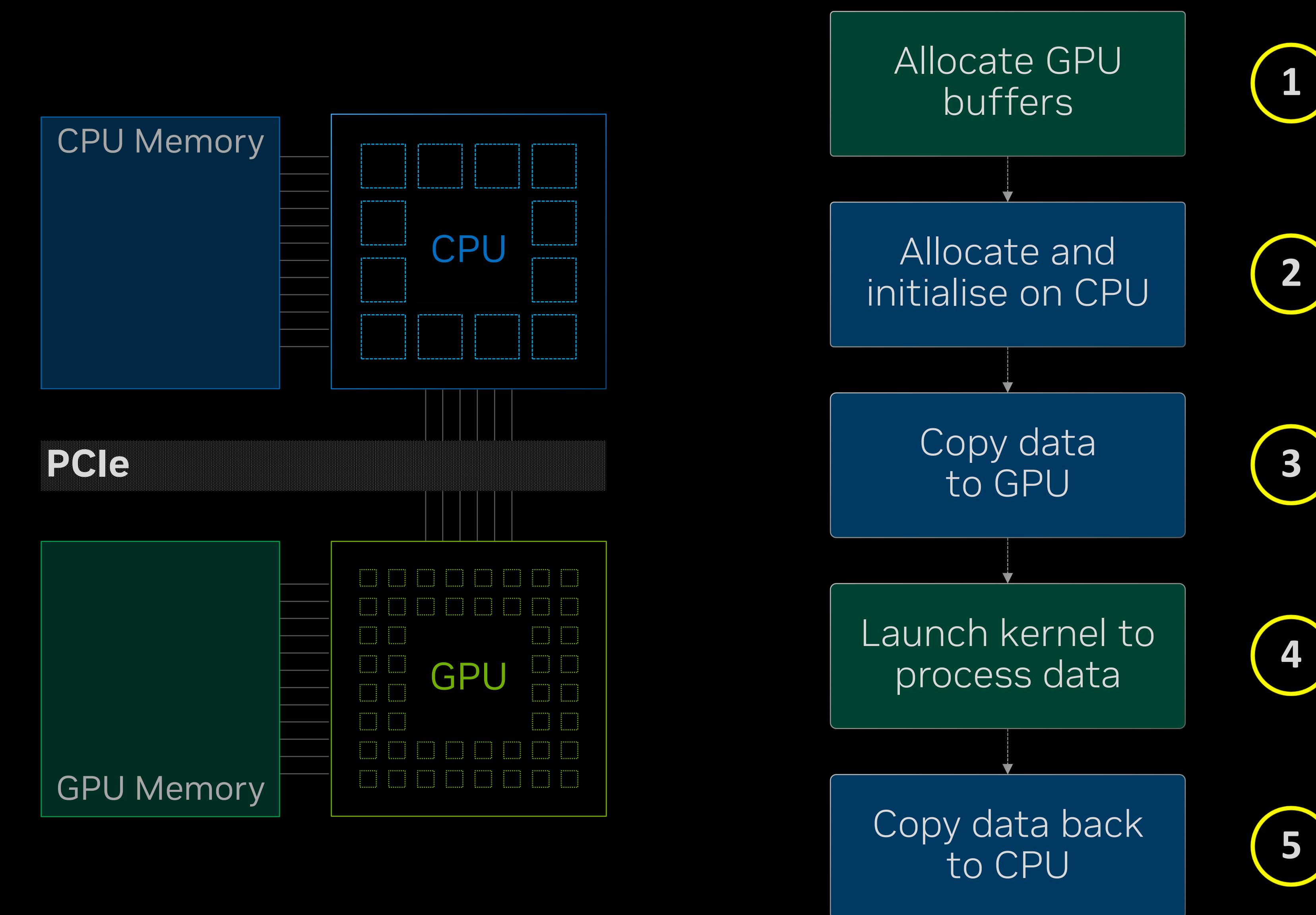
CUDA Applications Run Unmodified Under Confidential Computing

```
int *gpu_compute(size_t size) {  
    int *d_data, *d_result, *h_data;  
  
    // Allocate GPU buffers  
    cudaMalloc(&d_data, size);  
    cudaMalloc(&d_result, size);  
  
    // Allocate and initialise on CPU  
    cudaMallocHost(h_data, size);  
    init(h_data, size);  
  
    // Copy to GPU and process it, then copy back  
    cudaMemcpy(d_data, h_data, size);  
    kernel<<< ... >>>(d_data, d_result);  
    cudaMemcpy(h_data, d_data, size);  
  
    // Return the result  
    return h_data;  
}
```



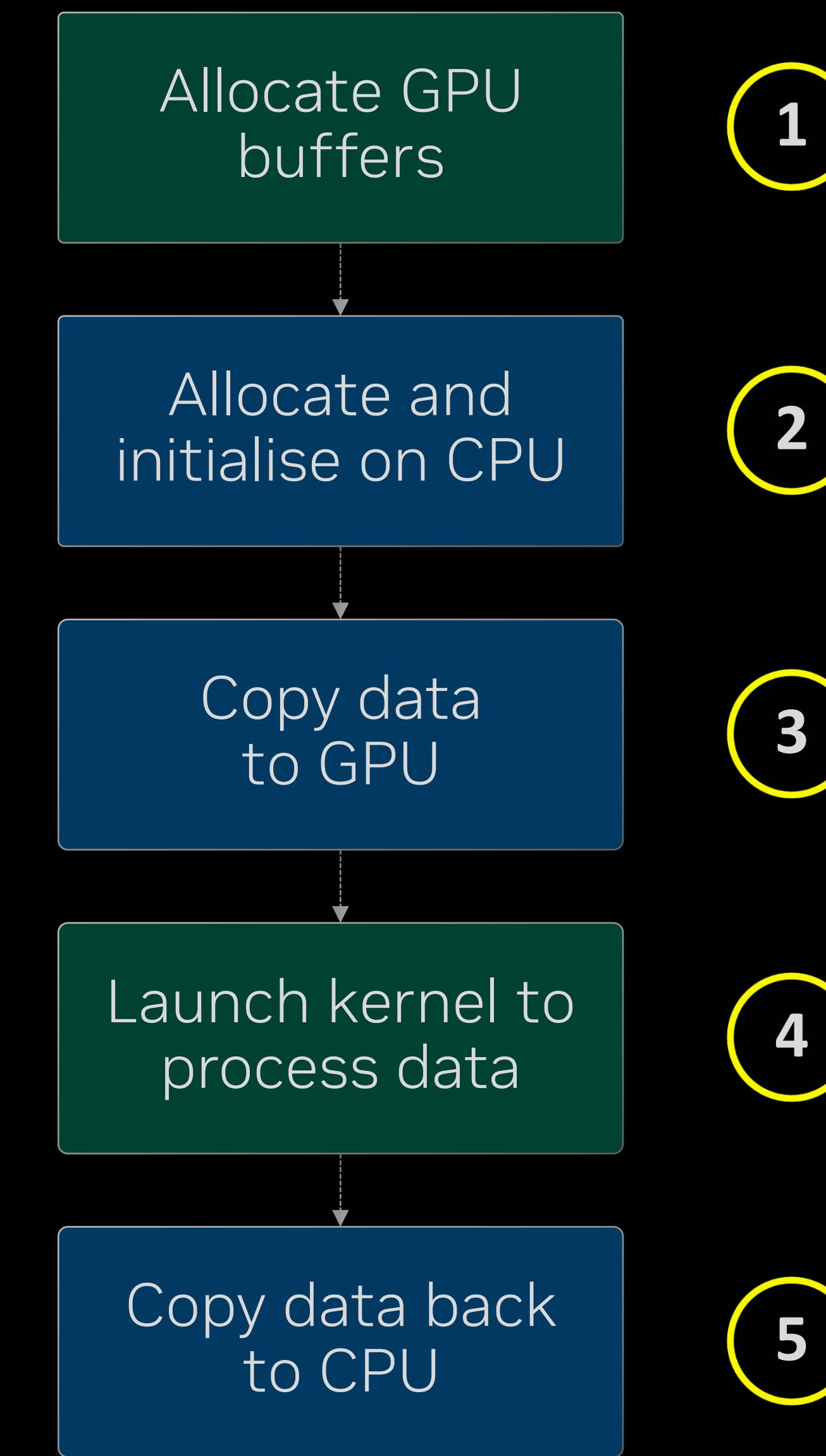
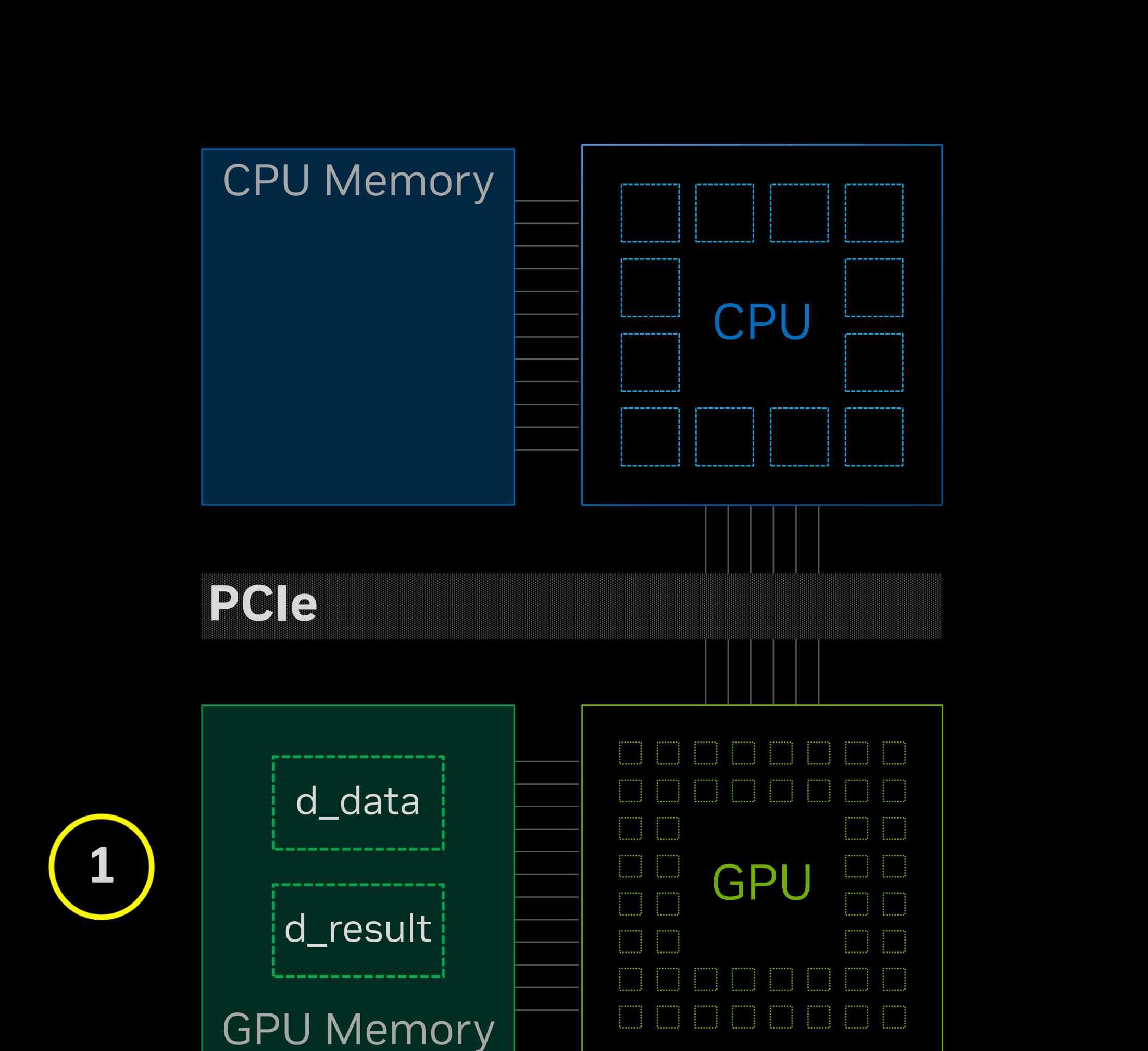
CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Allocate GPU buffers	Allocate GPU buffers
2	Allocate and initialise on CPU	Allocate and initialise on CPU
3	Copy data to GPU	Copy data to GPU
4	Launch kernel to process data	Launch kernel to process data
5	Copy data back to CPU	Copy data back to CPU



CUDA Applications Run Unmodified Under Confidential Computing

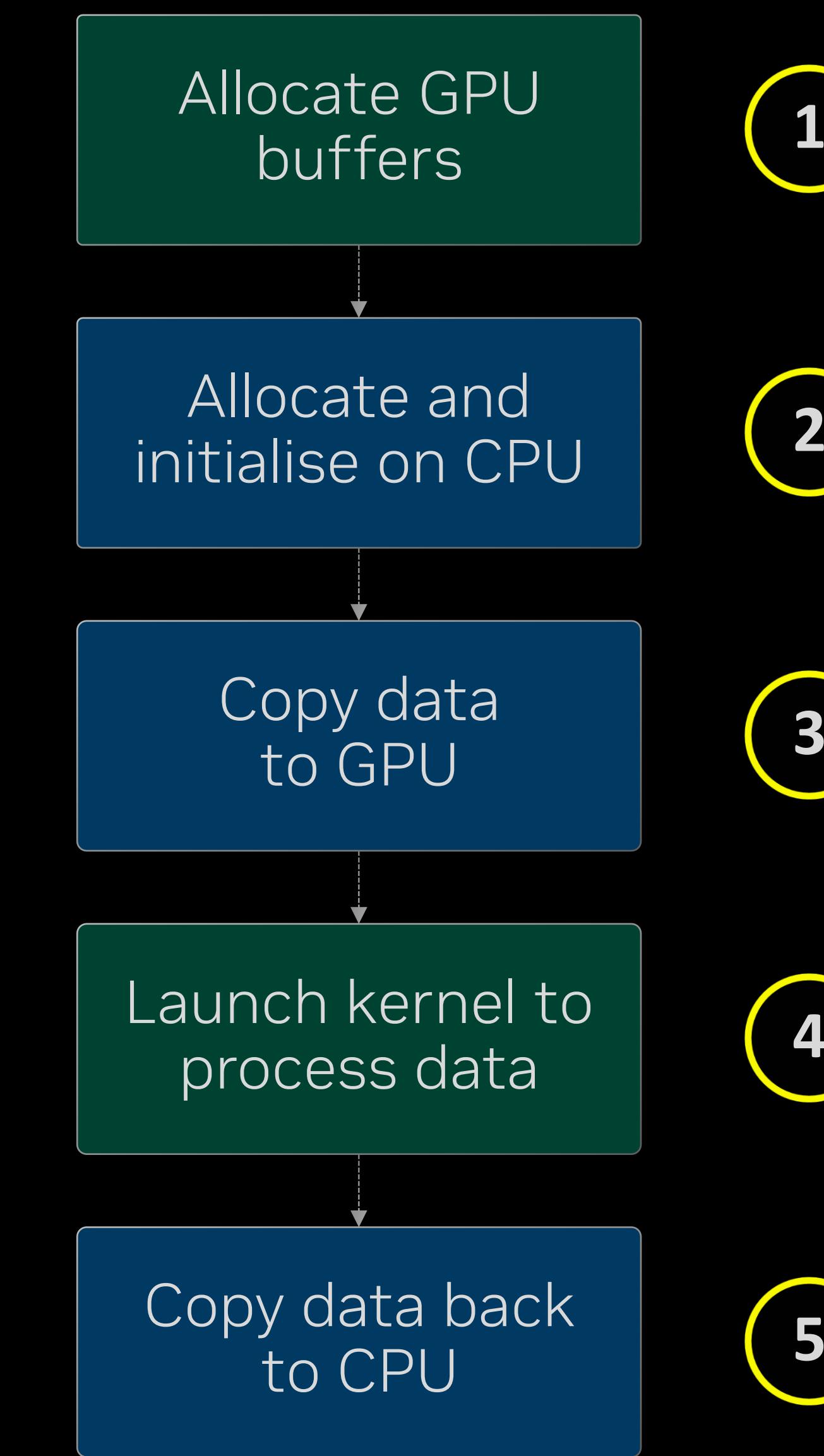
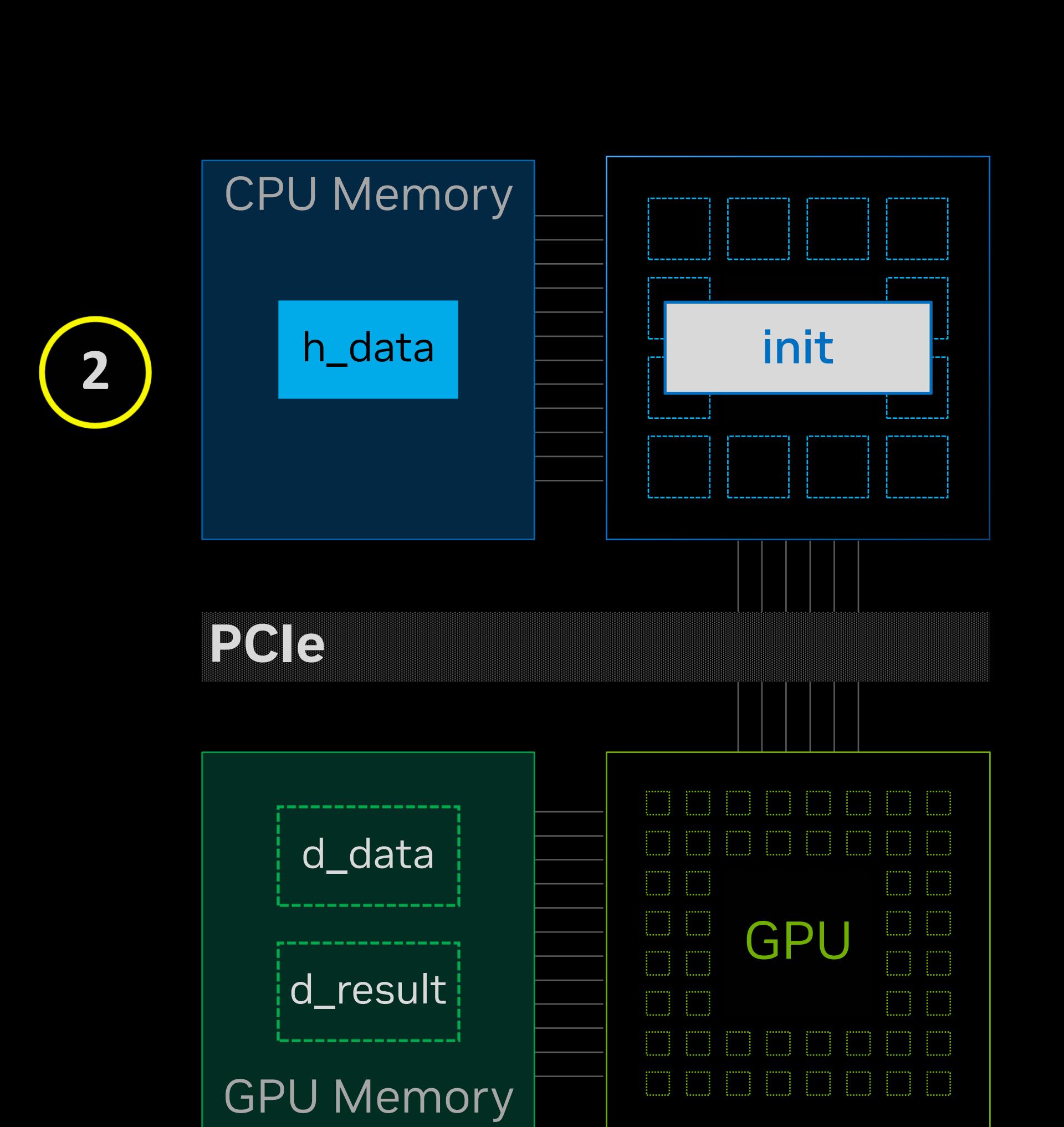
Step	CC Off	CC On
1	Normal GPU allocation	



```
cudaMalloc(&d_data, size);  
cudaMalloc(&d_result, size);
```

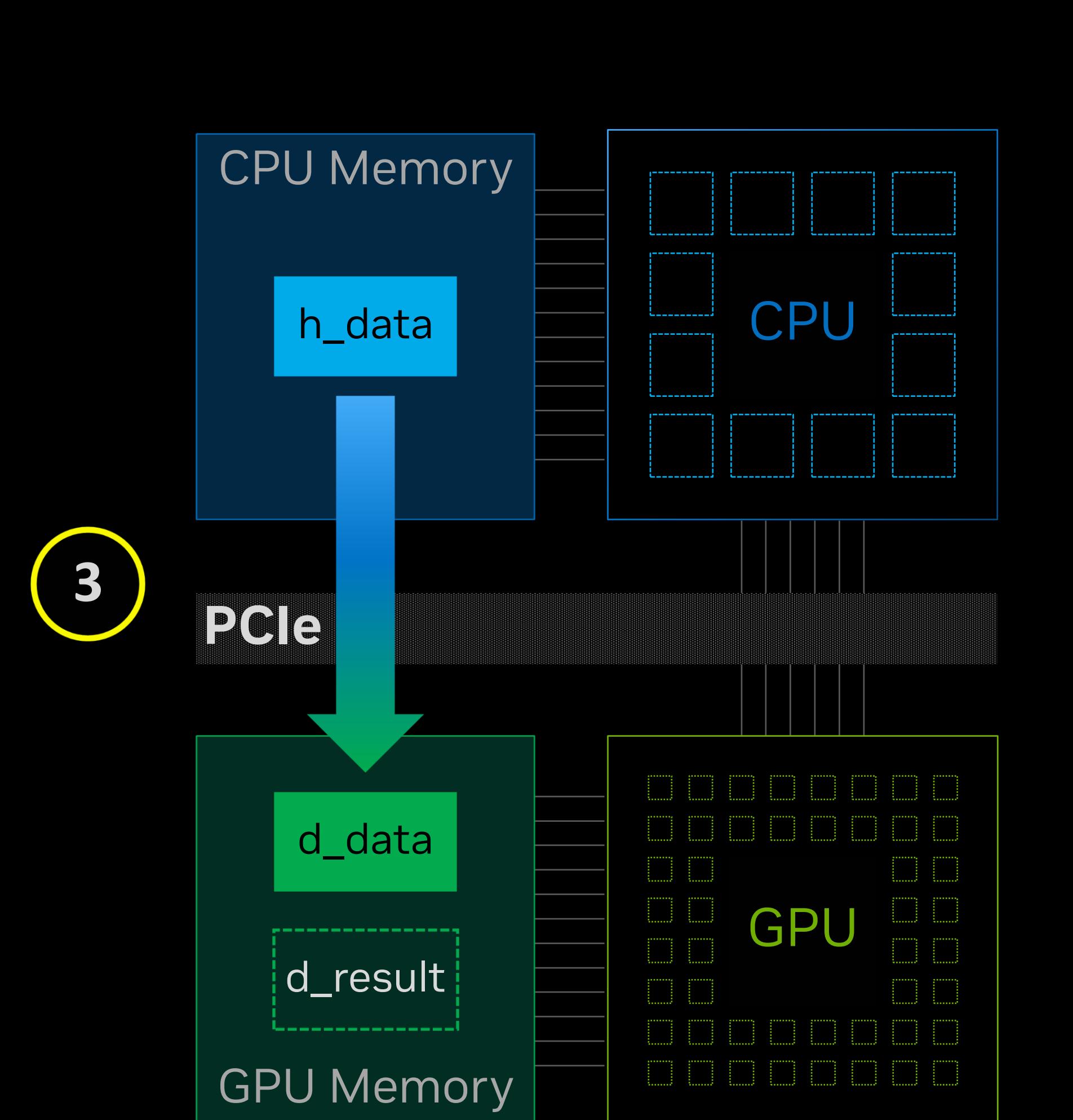
CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	
2	Pinned system allocation & init	

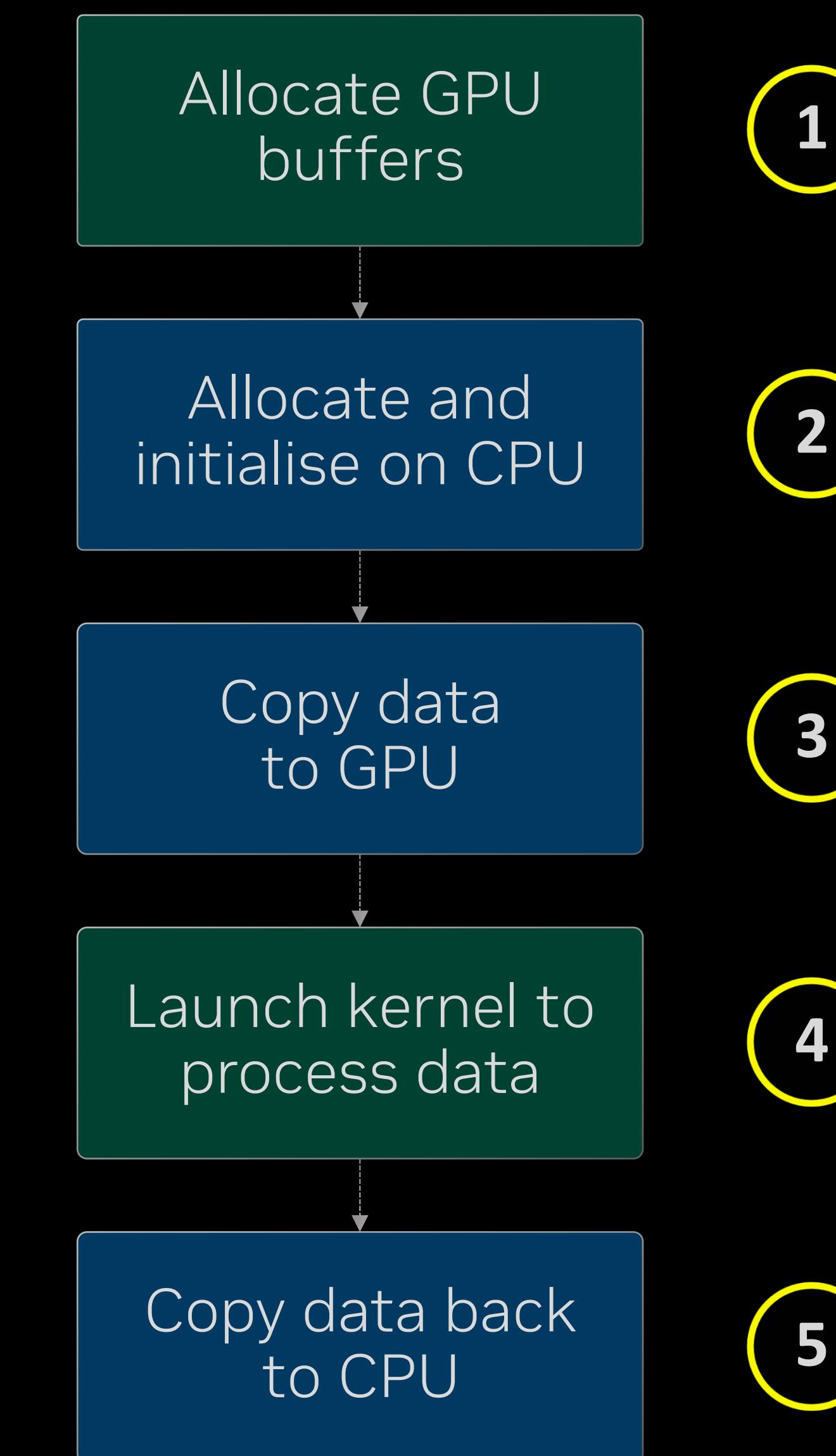


CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	
2	Pinned system allocation & init	
3	DMA transfer	

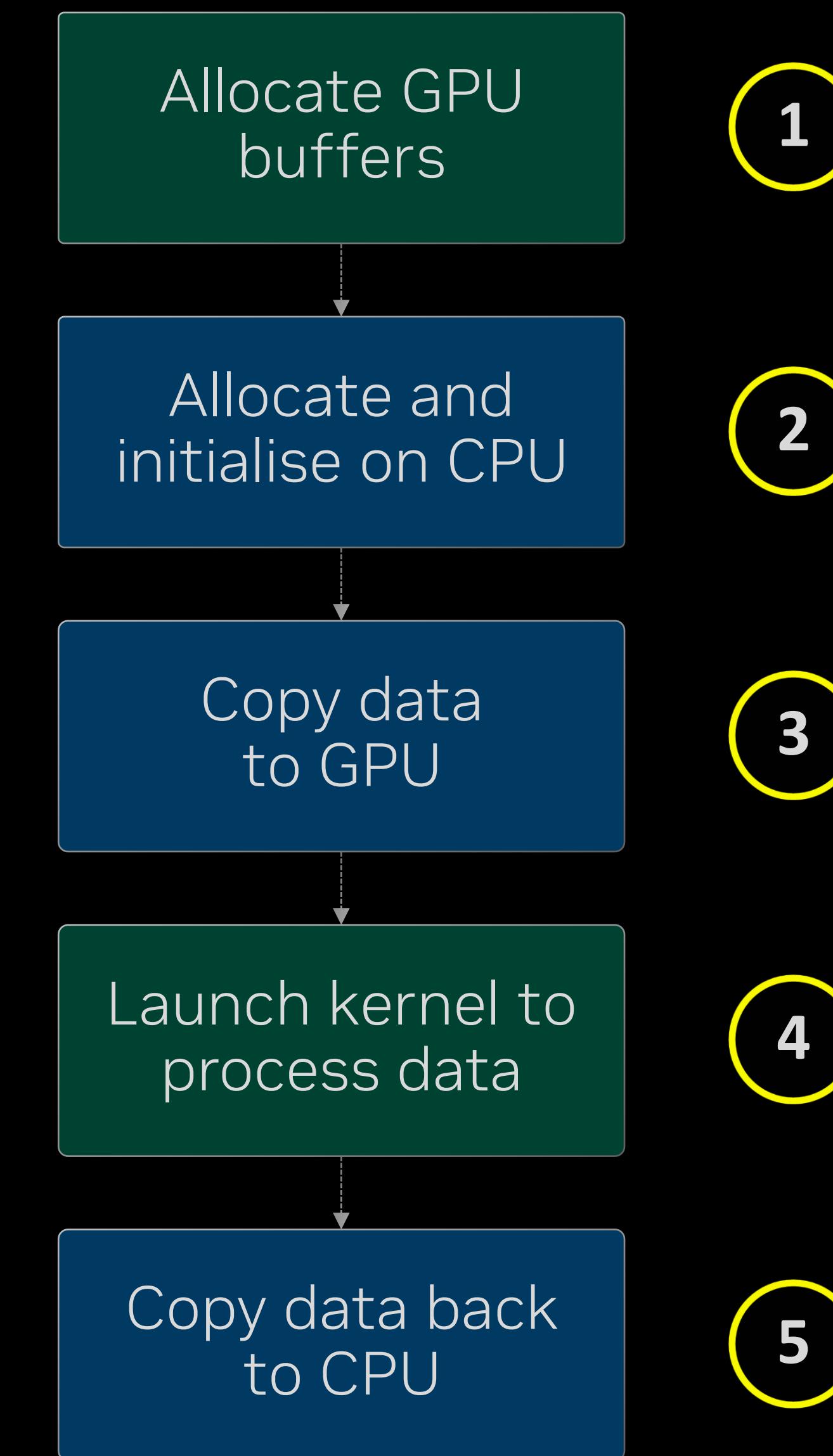
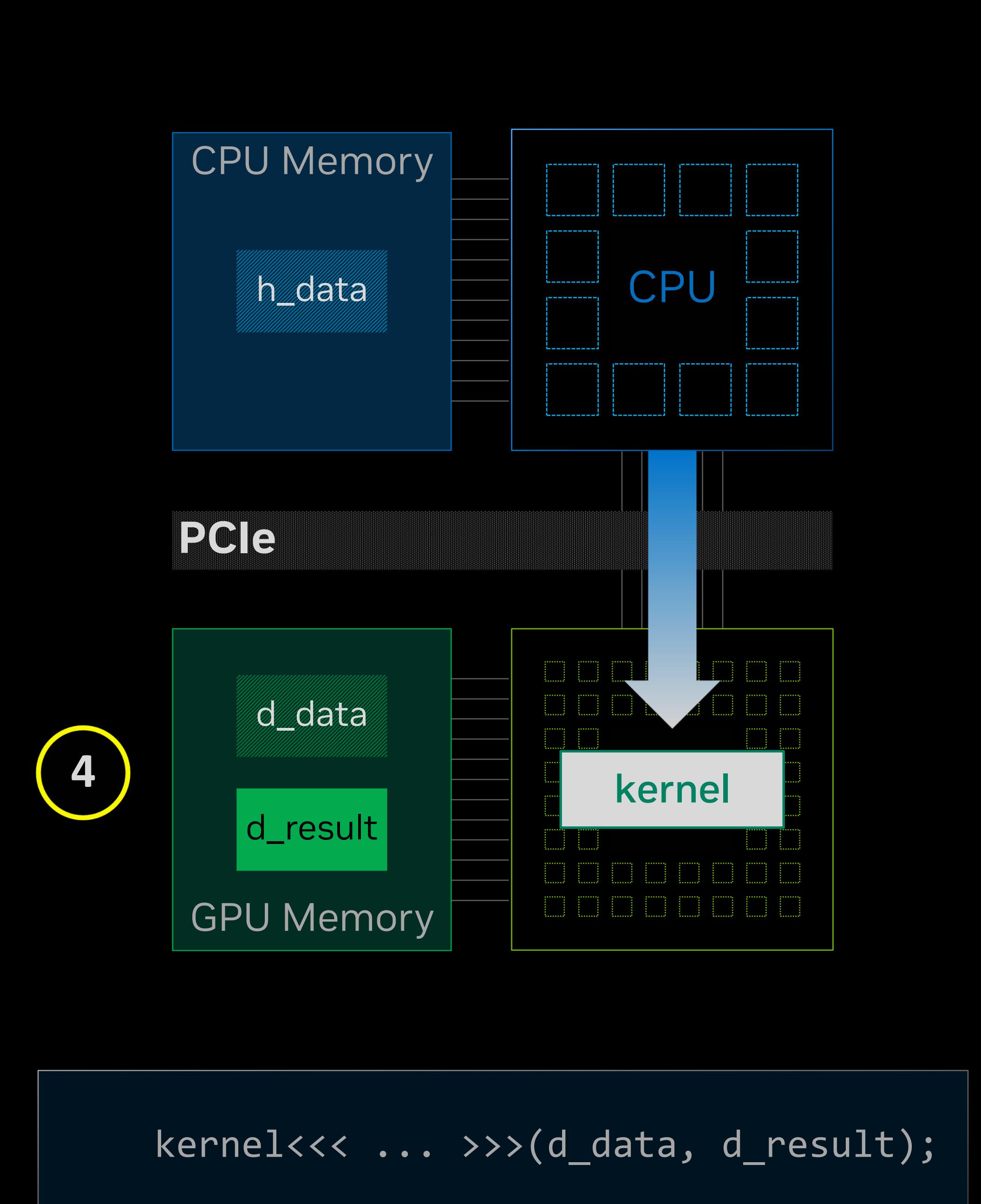


```
cudaMemcpy(d_data, h_data, size);
```



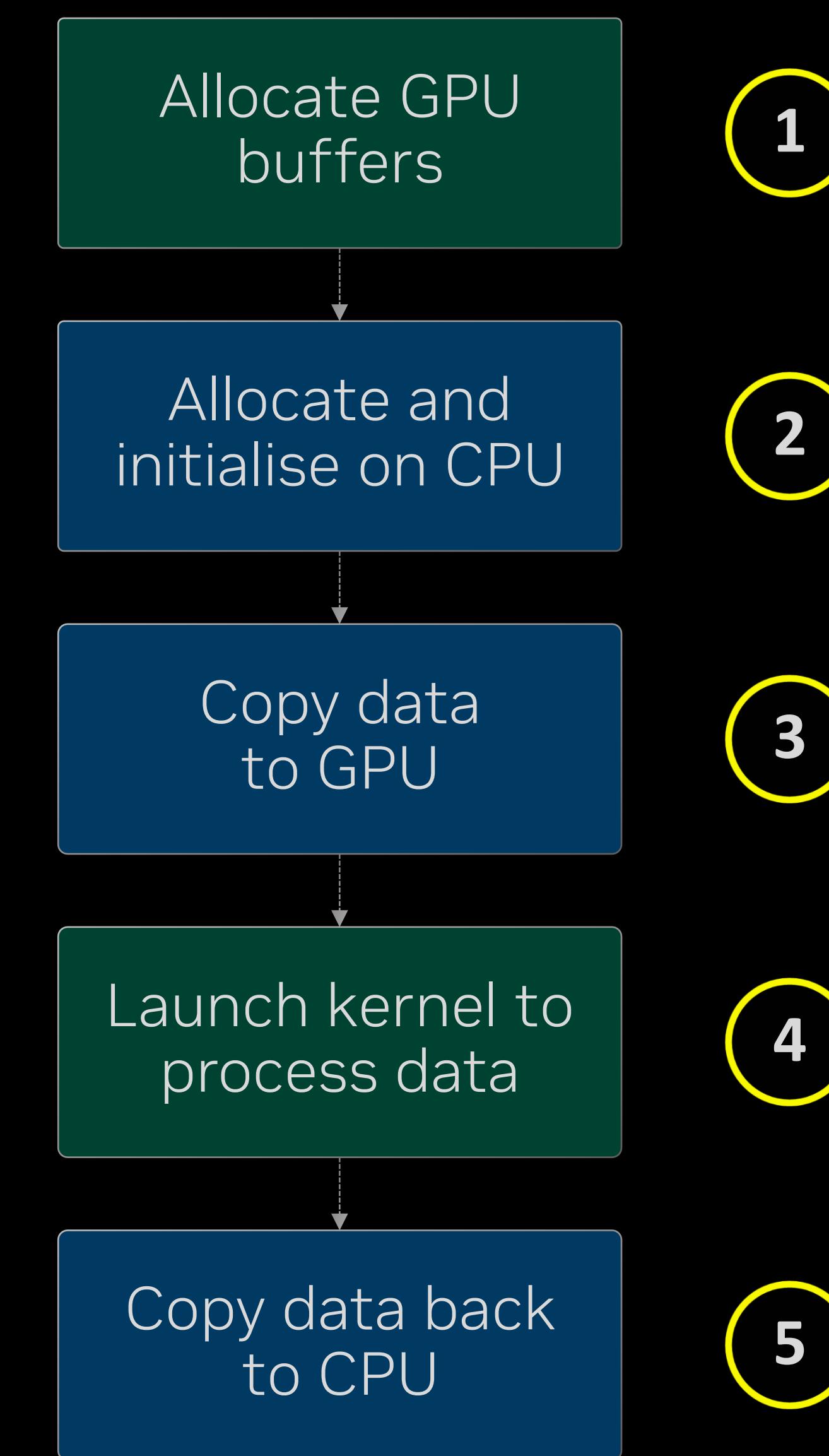
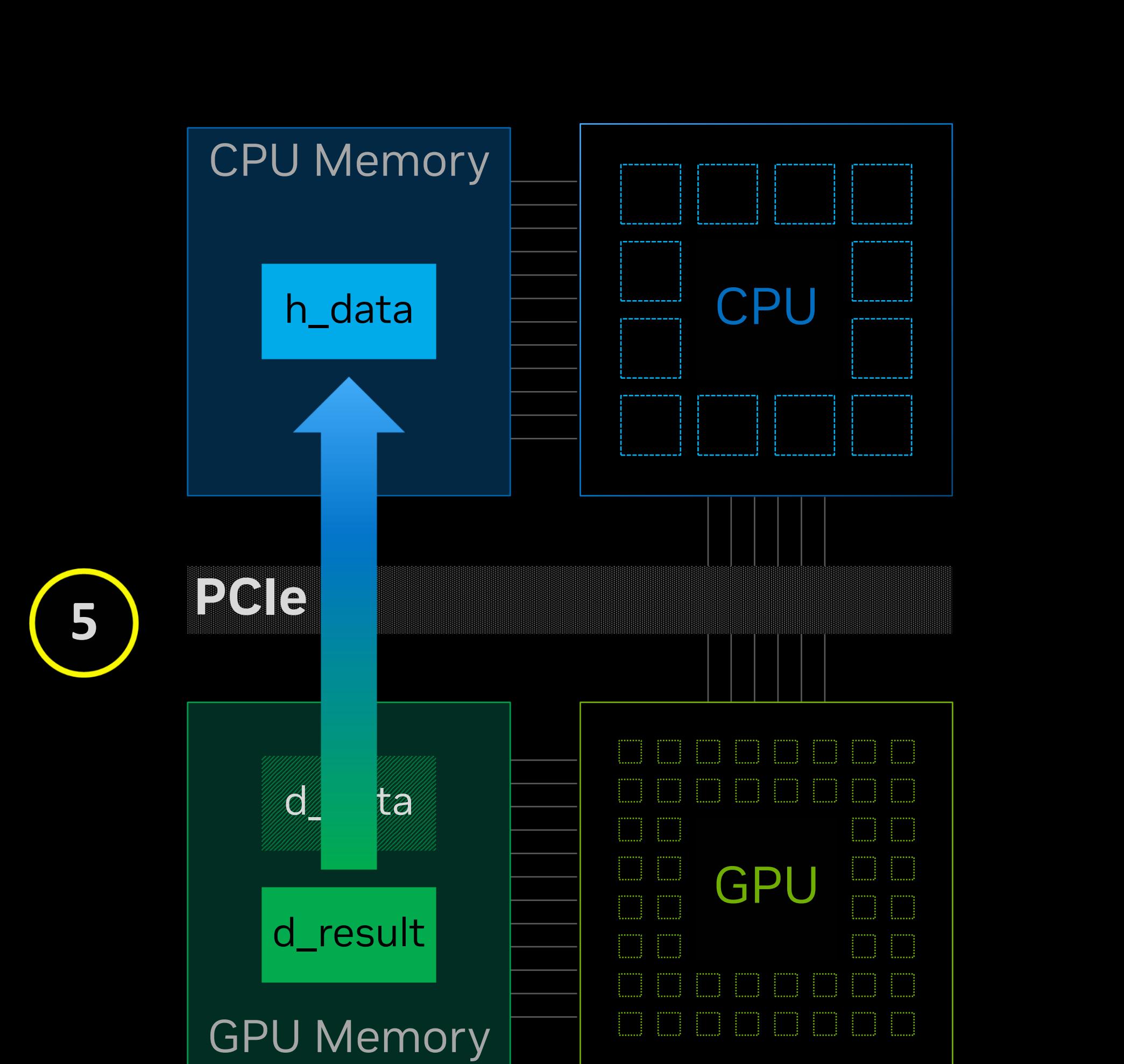
CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	
2	Pinned system allocation & init	
3	DMA transfer	
4	Kernel launch	



CUDA Applications Run Unmodified Under Confidential Computing

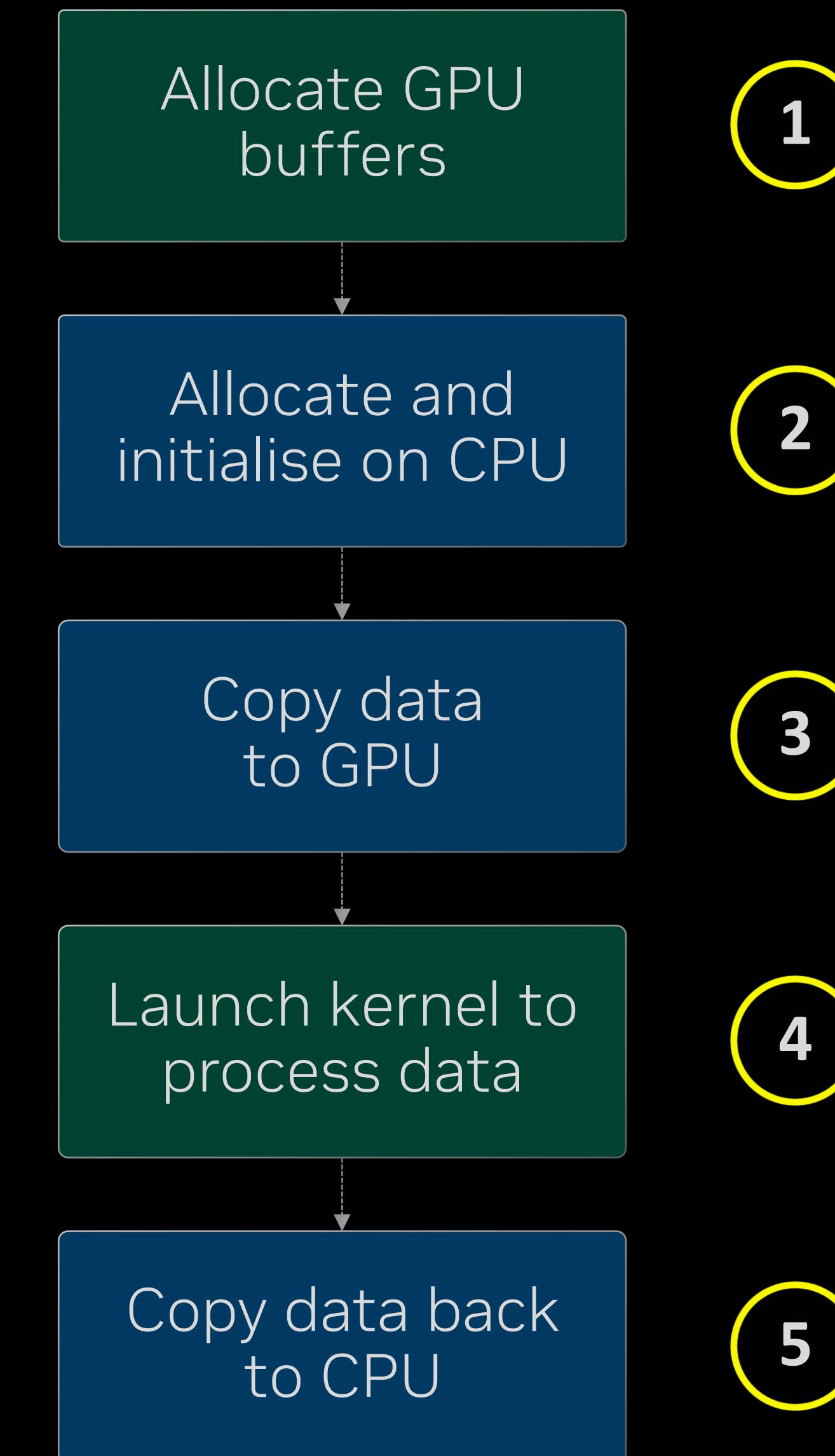
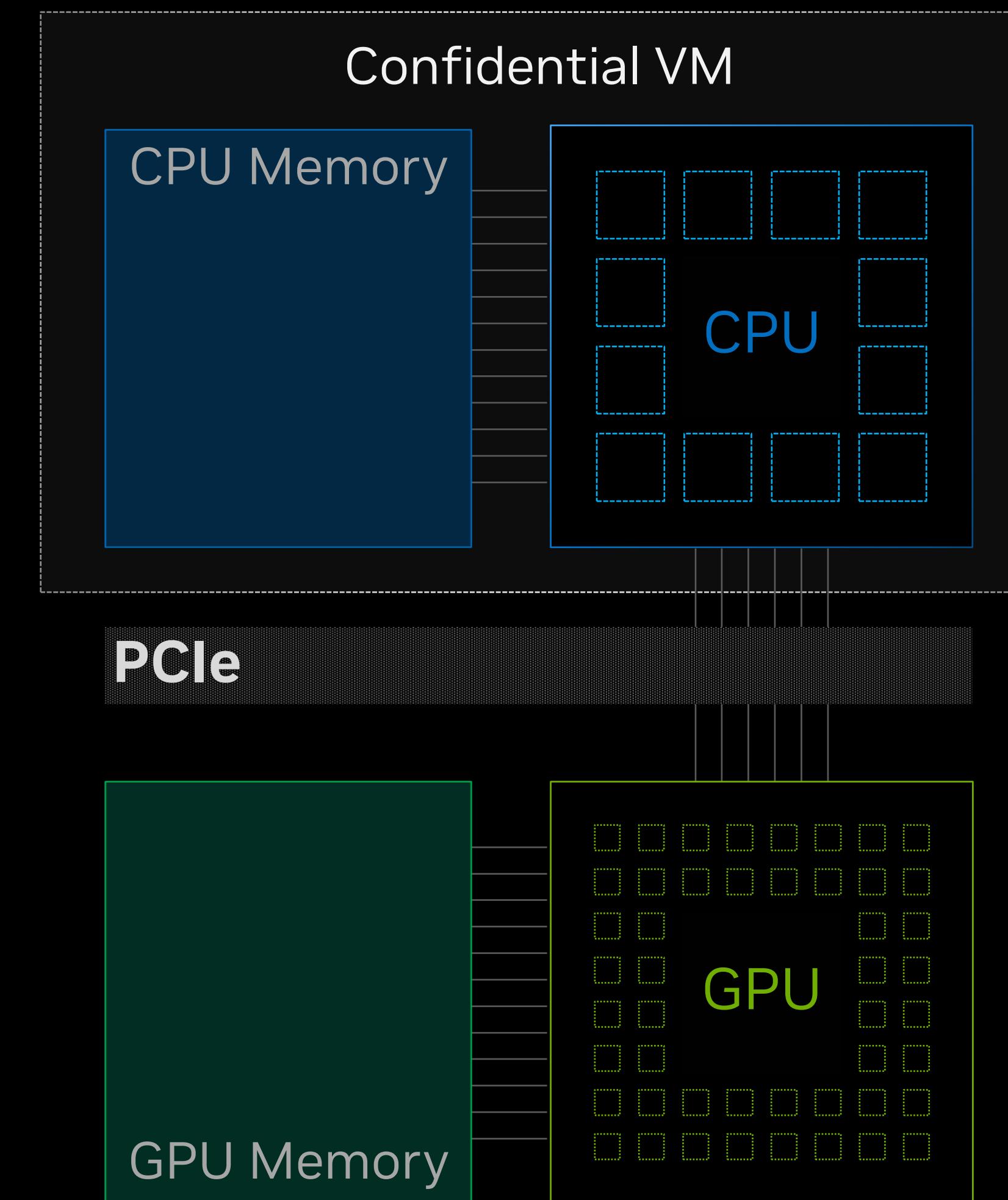
Step	CC Off	CC On
1	Normal GPU allocation	
2	Pinned system allocation & init	
3	DMA transfer	
4	Kernel launch	
5	DMA transfer	



```
cudaMemcpy(h_data, d_data, size);
```

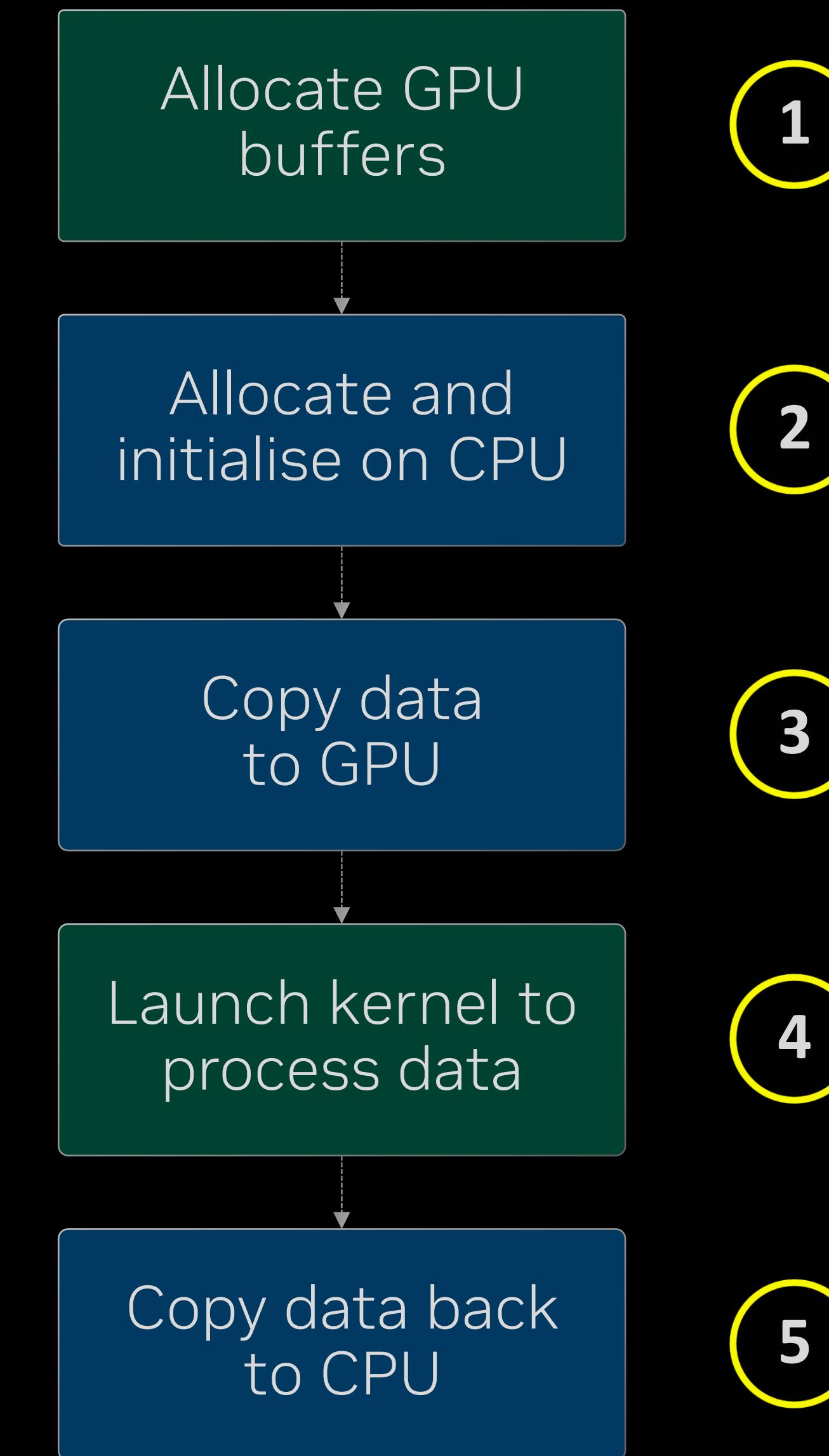
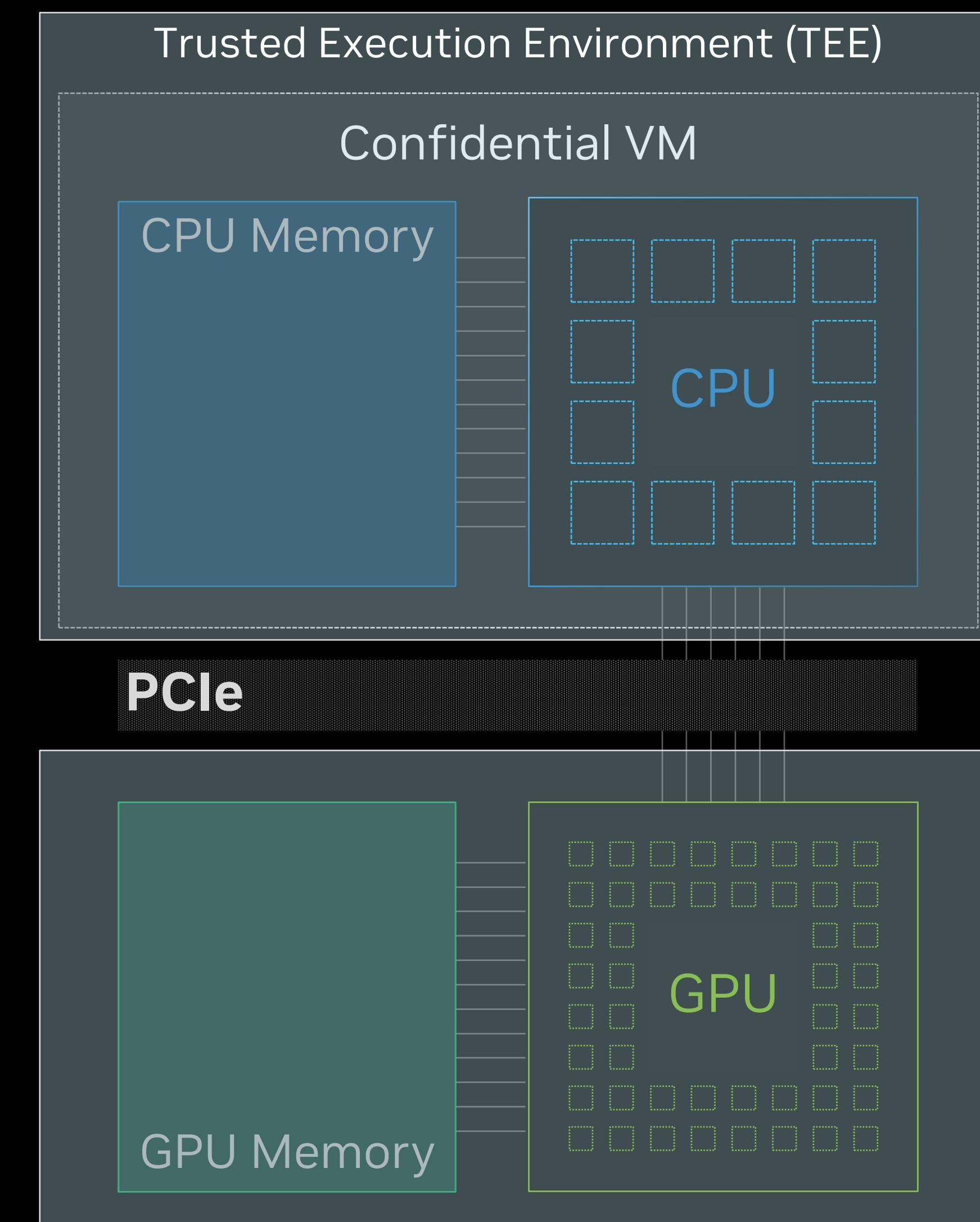
CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
	Normal GPU allocation	
	Pinned system allocation & init	
	DMA transfer	
	Kernel launch	
	DMA transfer	



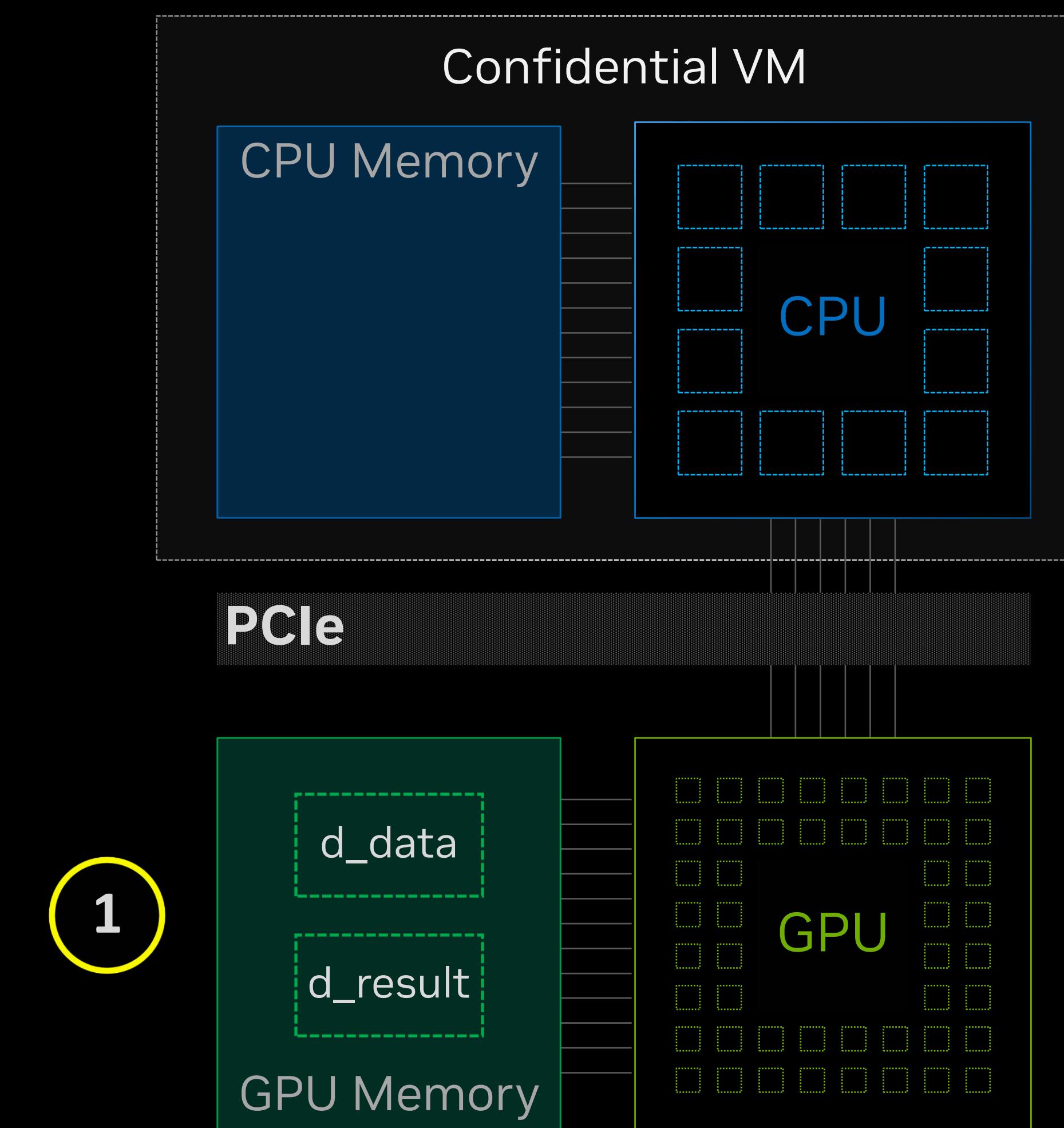
CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
	Normal GPU allocation	
	Pinned system allocation & init	
	DMA transfer	
	Kernel launch	
	DMA transfer	

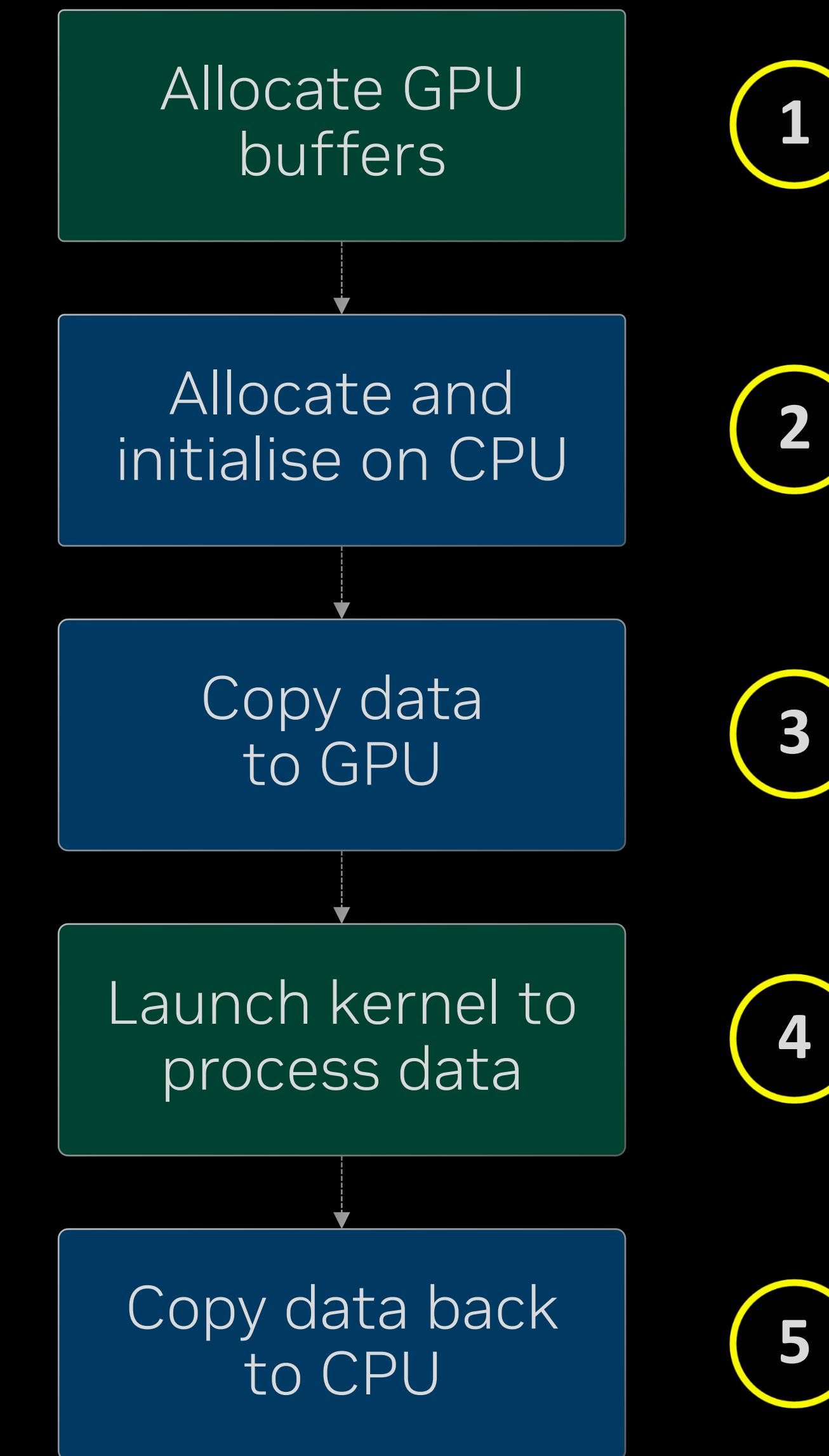


CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	Normal GPU allocation
	Pinned system allocation & init	
	DMA transfer	
	Kernel launch	
	DMA transfer	

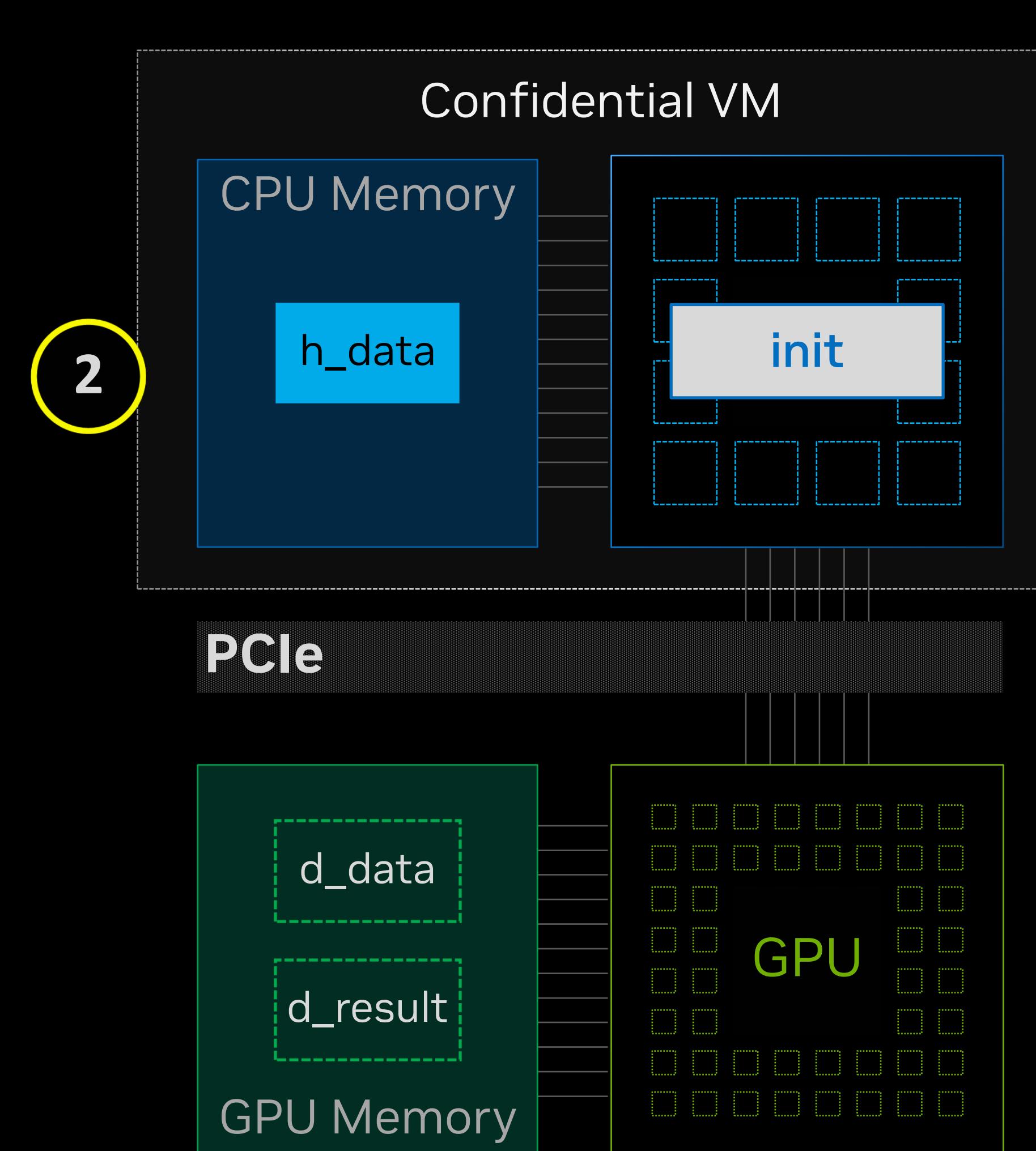


```
cudaMalloc(&d_data, size);  
cudaMalloc(&d_result, size);
```

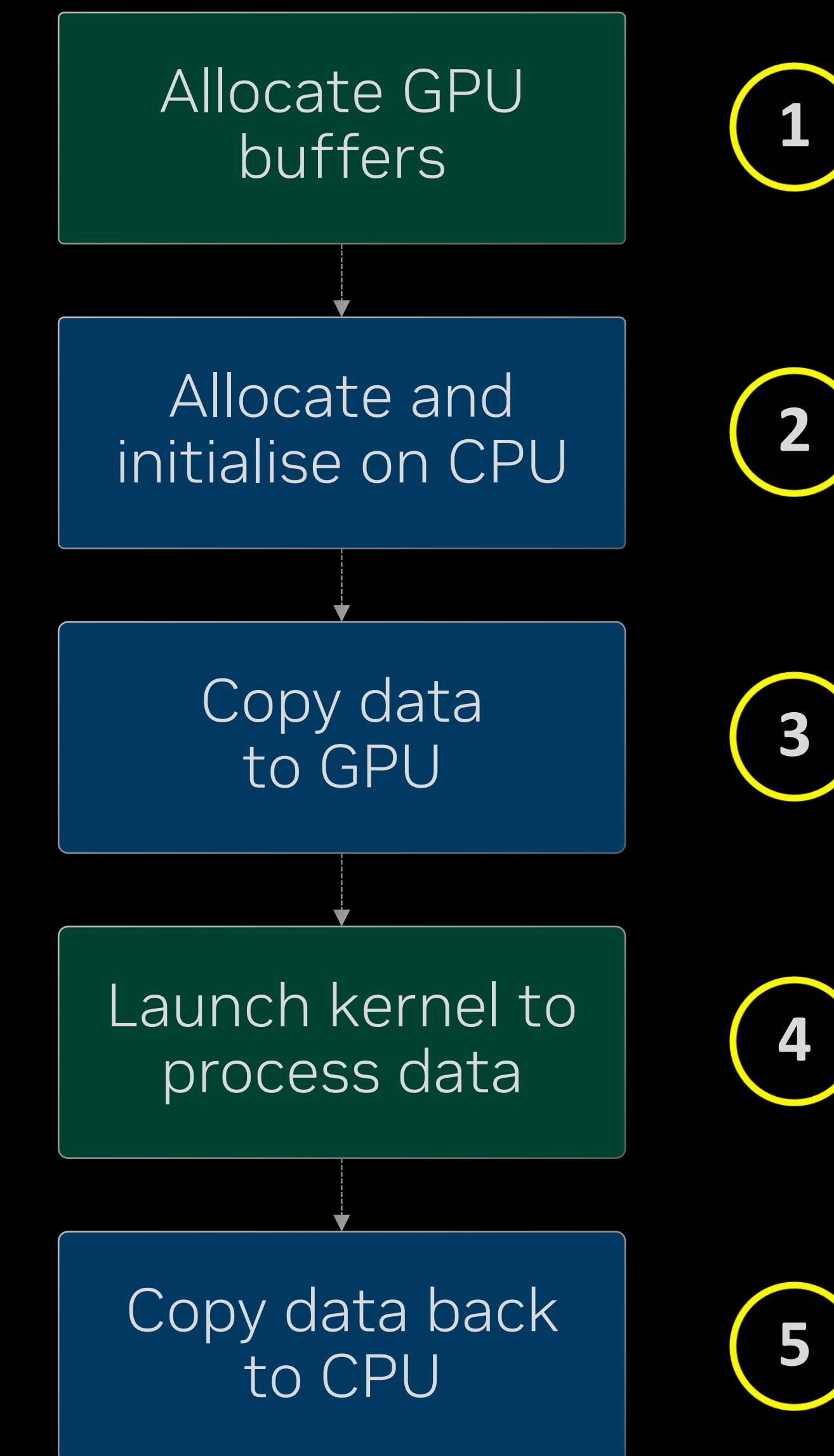


CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	Normal GPU allocation
2	Pinned system allocation & init	Managed allocation & init
	DMA transfer	
	Kernel launch	
	DMA transfer	

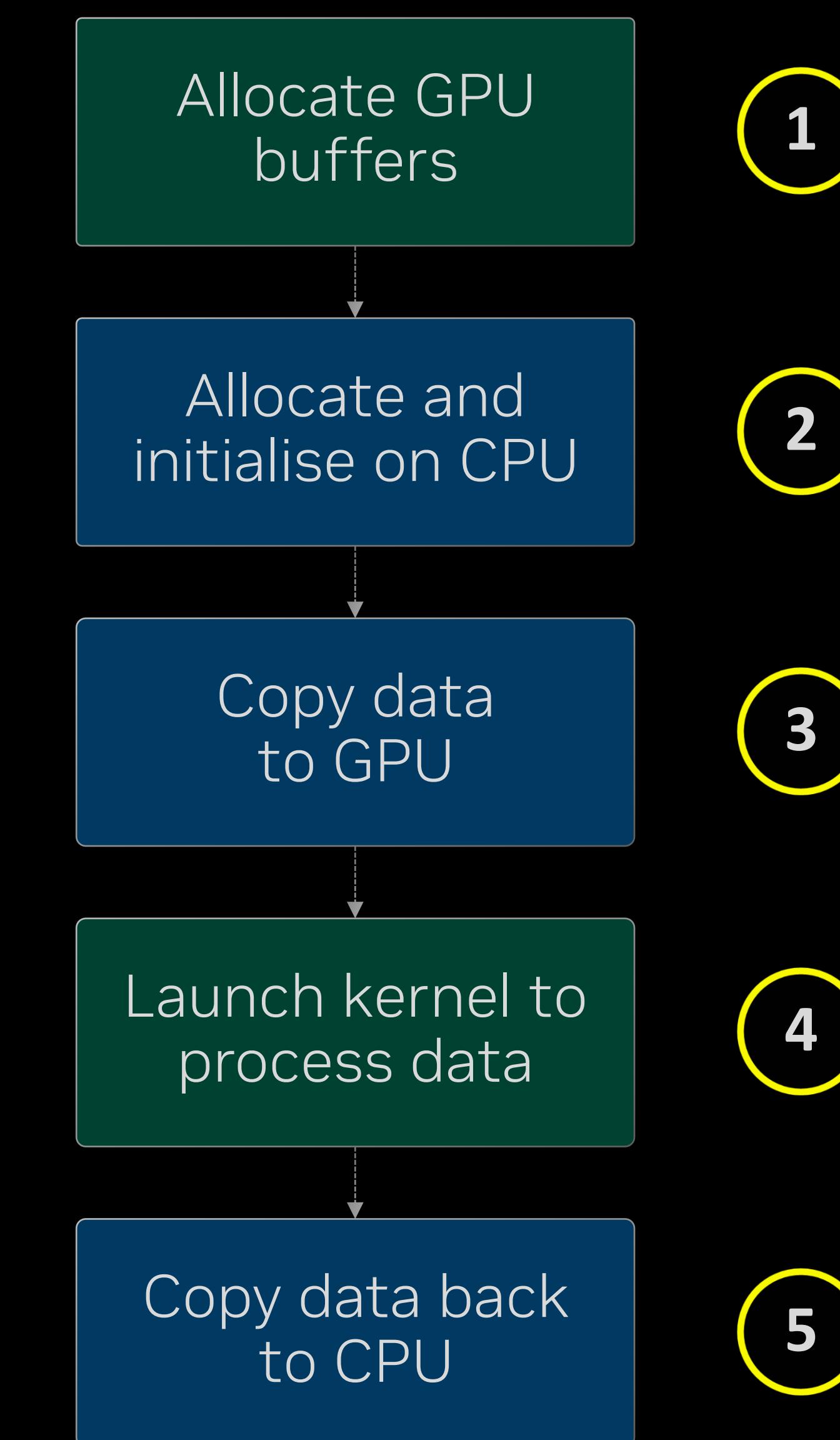
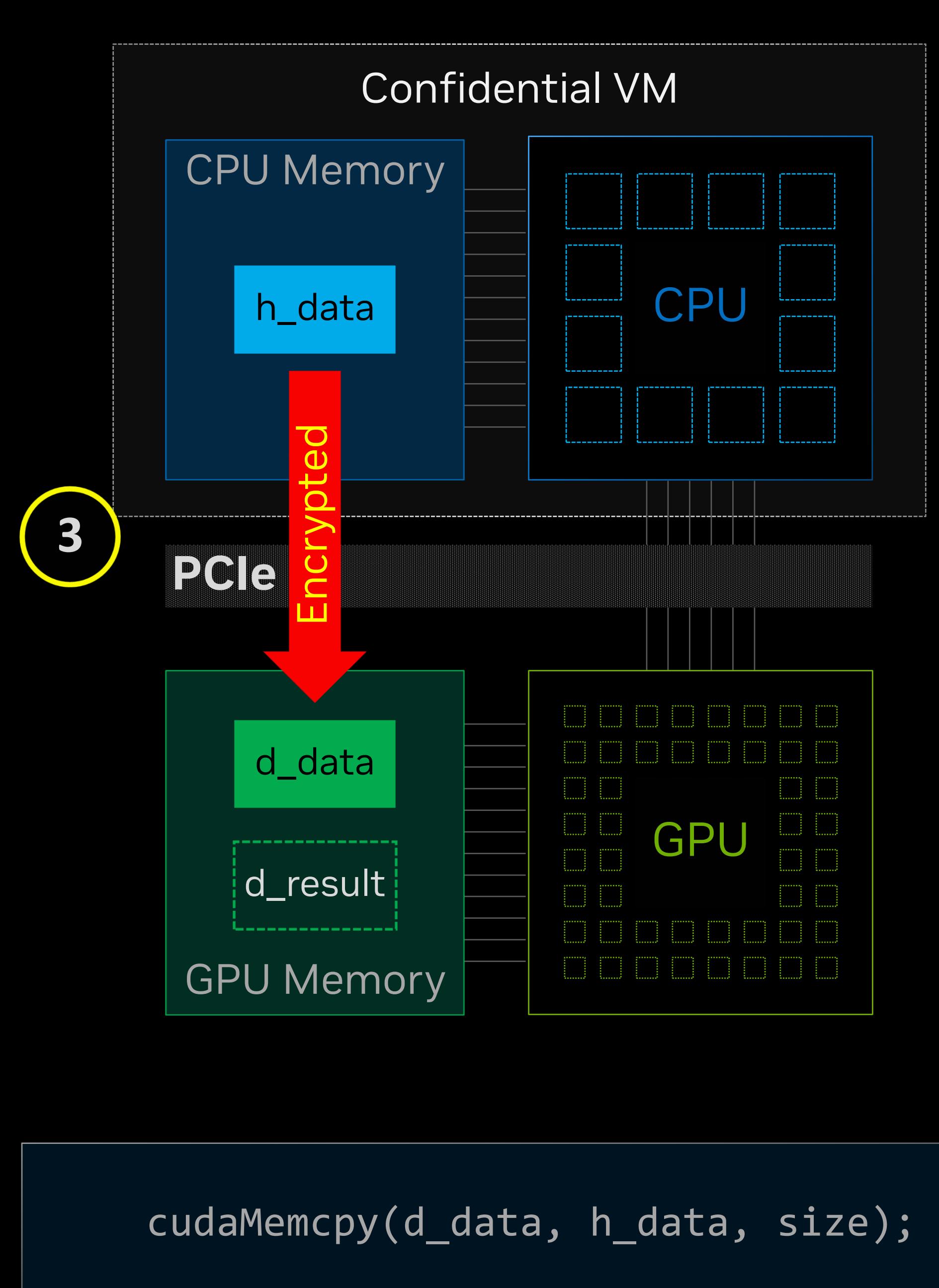


```
cudaMallocHost(h_data, size);  
init(h_data, size);
```



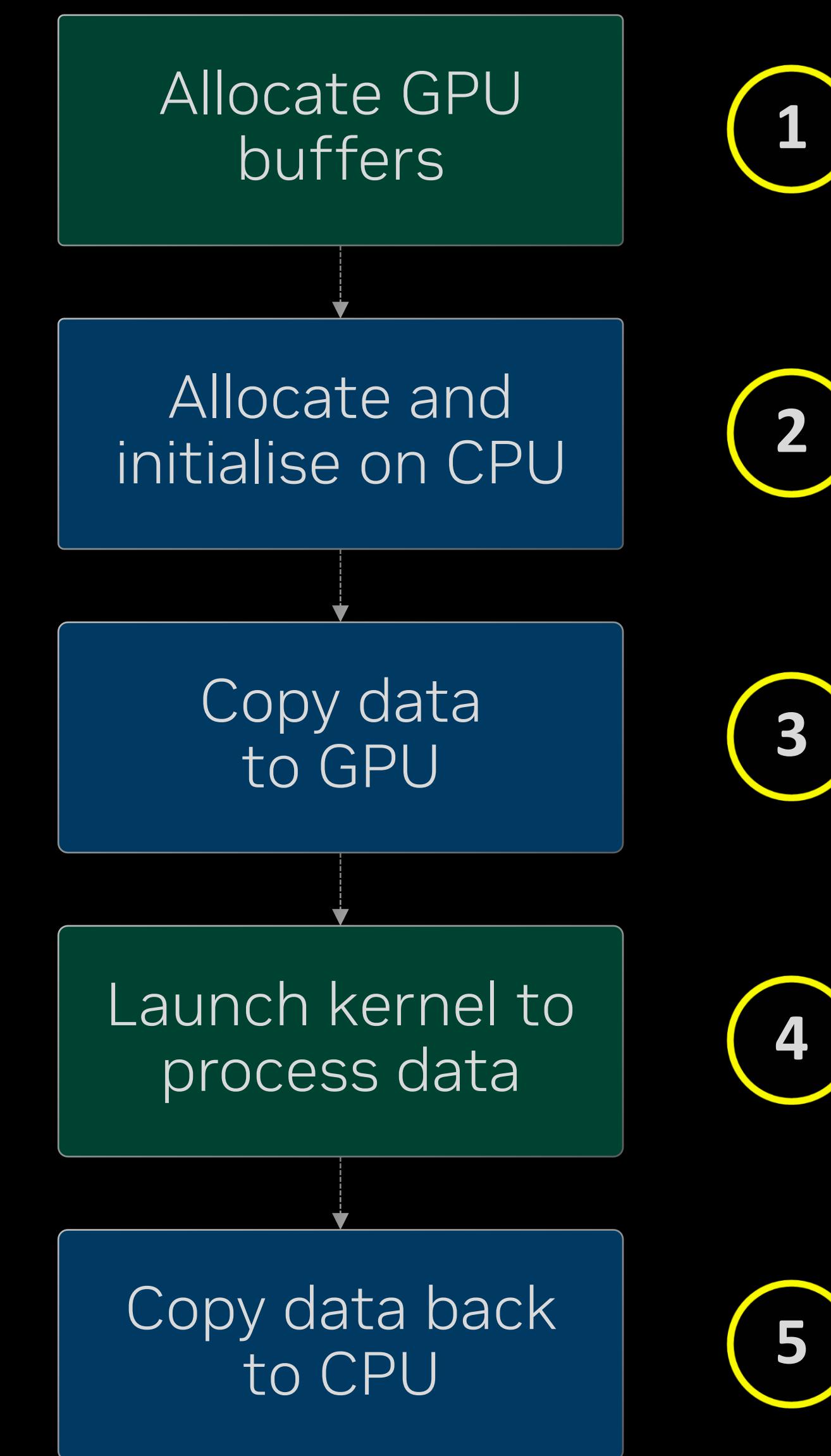
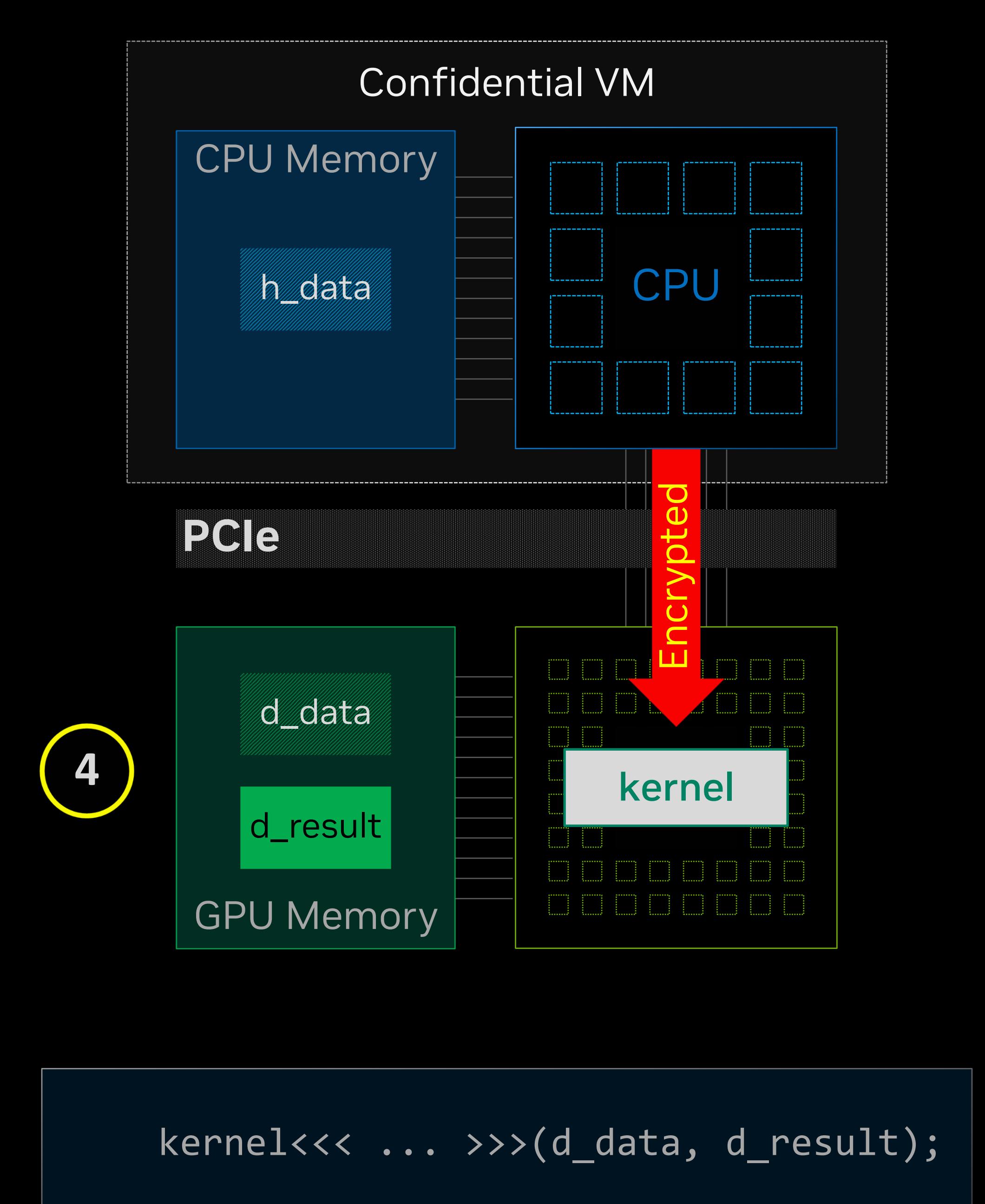
CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	Normal GPU allocation
2	Pinned system allocation & init	Managed allocation & init
3	DMA transfer	Encrypted DMA transfer
	Kernel launch	
	DMA transfer	



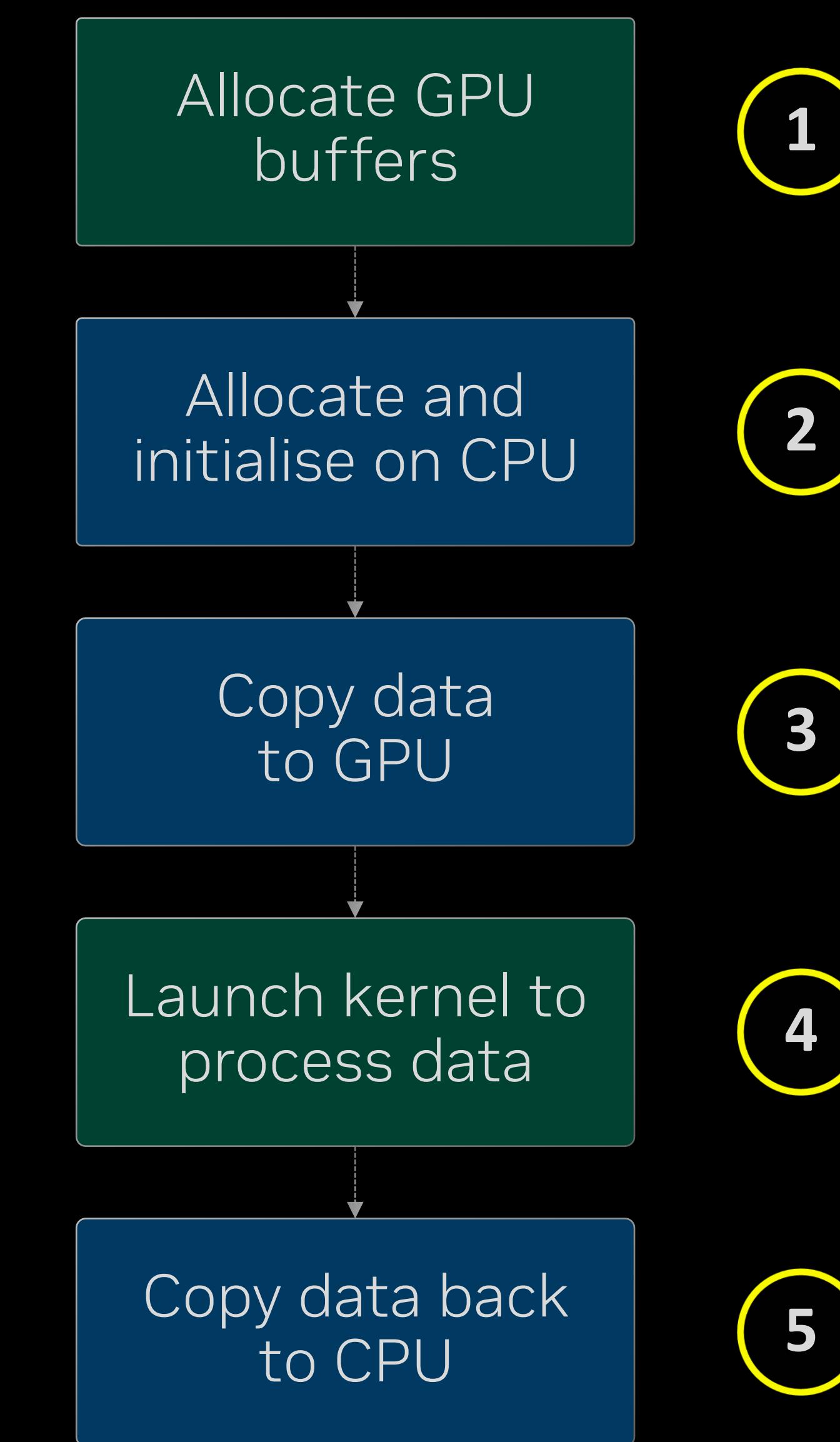
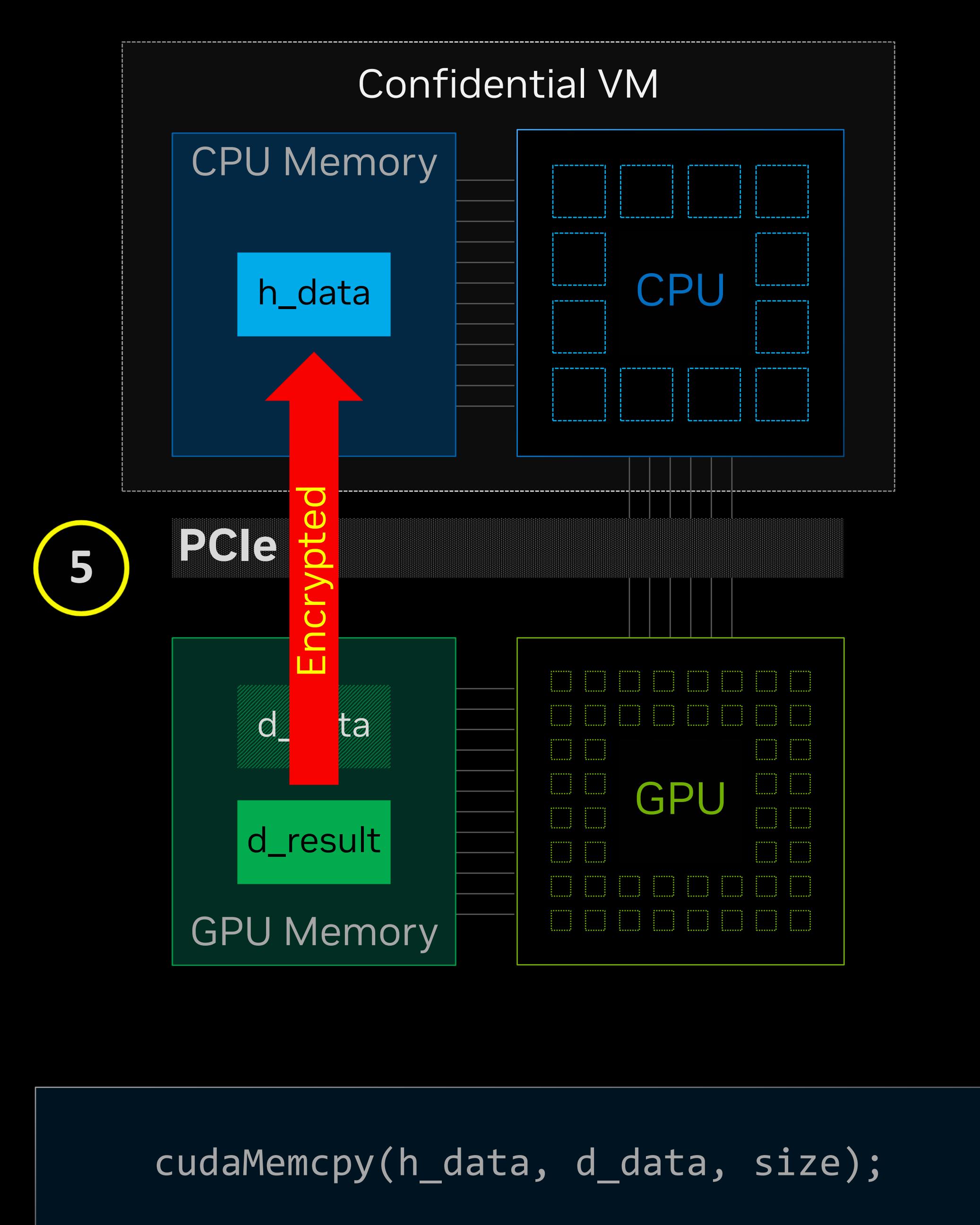
CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	Normal GPU allocation
2	Pinned system allocation & init	Managed allocation & init
3	DMA transfer	Encrypted DMA transfer
4	Kernel launch	Encrypted kernel launch
	DMA transfer	



CUDA Applications Run Unmodified Under Confidential Computing

Step	CC Off	CC On
1	Normal GPU allocation	Normal GPU allocation
2	Pinned system allocation & init	Managed allocation & init
3	DMA transfer	Encrypted DMA transfer
4	Kernel launch	Encrypted kernel launch
5	DMA transfer	Encrypted DMA transfer



References

- S51120 [Programming Model and Applications for the Grace Hopper Superchip](#)
- S51054 [Accelerating HPC applications with ISO C++ on Grace Hopper](#)
- S51755 [C++ Standard Parallelism](#)
- S51210 [How to Write a CUDA Program](#)
- S51705 [How to Streamline Shared Memory Space With the NVSHMEM Communication Library](#)
- S51789 [cuNumeric and Legate: How to Create a Distributed GPU Accelerated Library](#)
- S51176 [Recent Developments in NVIDIA Math Libraries](#)
- S51413 [Developing Optimal CUDA Kernels on Hopper Tensor Cores](#)
- S51205 [From the Macro to the Micro: CUDA Developer Tools Find and Fix Problems at Any Scale](#)
- S51421 [Optimizing at Scale: Investigating and Resolving Hidden Bottlenecks for Multi-Node Workloads](#)
- S51709 [Hopper Confidential Computing: How it Works under the Hood](#)
- S51684 [Confidential Computing: The Developer's View to Secure an Application and Data on NVIDIA H100](#)
- S51074 [A Deep Dive into the Latest HPC Software](#)
- S51880 [A Demonstration of AI and HPC Applications for NVIDIA Grace CPU](#)
- CWES52036 [Connect with the Experts: What's in Your CUDA Toolbox? CUDA Profiling, Optimization, and Debugging Tools for the Latest Architectures](#)
- Whitepaper [NVIDIA Grace Hopper Superchip Architecture](#)
- Blog [CUDA 12.0 Compiler Support for Runtime LTO Using nvJitLink Library](#)
- Package [Download MathDx Today](#)
- Repo [Download CUTLASS 3.0 from the GitHub repo](#)