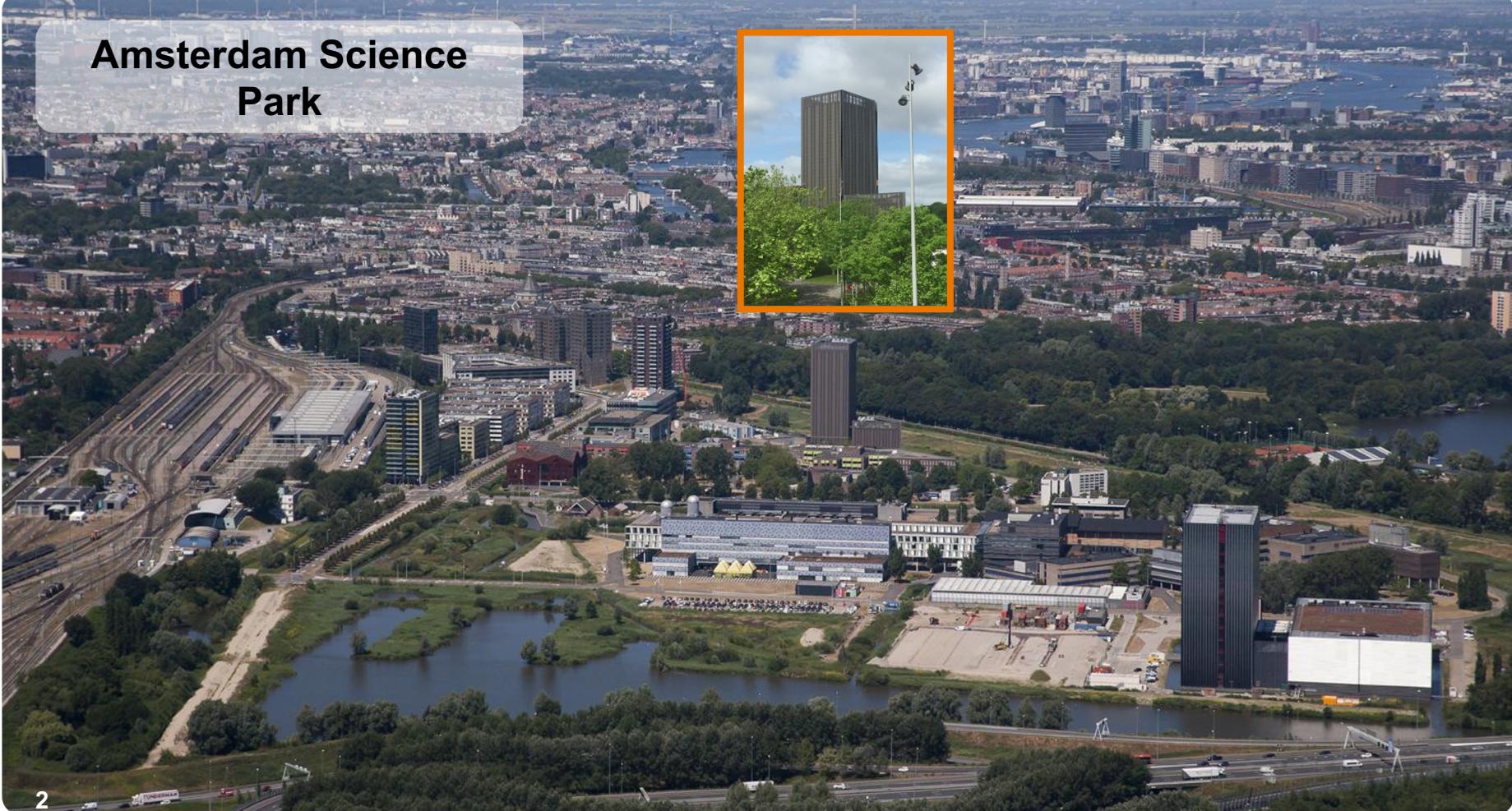

Introduction to accelerated computing ecosystem

PRACE Parallel programming course with Python

Sagar Dolas
Program Lead, Innovation Lab
SURF

Amsterdam Science Park



Agenda

- Overview of accelerated computing
- GPU Architecture
- Software ecosystem & applications
- Alternative computing paradigms

Go to menti.com and enter the code : **9671 7210**

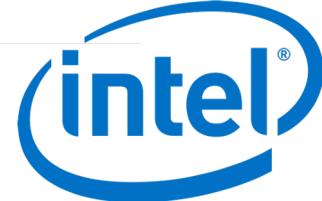


Main message



- Difference between CPU vs GPU
- State of GPUs internationally
- Different GPU architecture
- Programming GPUs
- New computing paradigms

Snellius : Dutch National Supercomputer



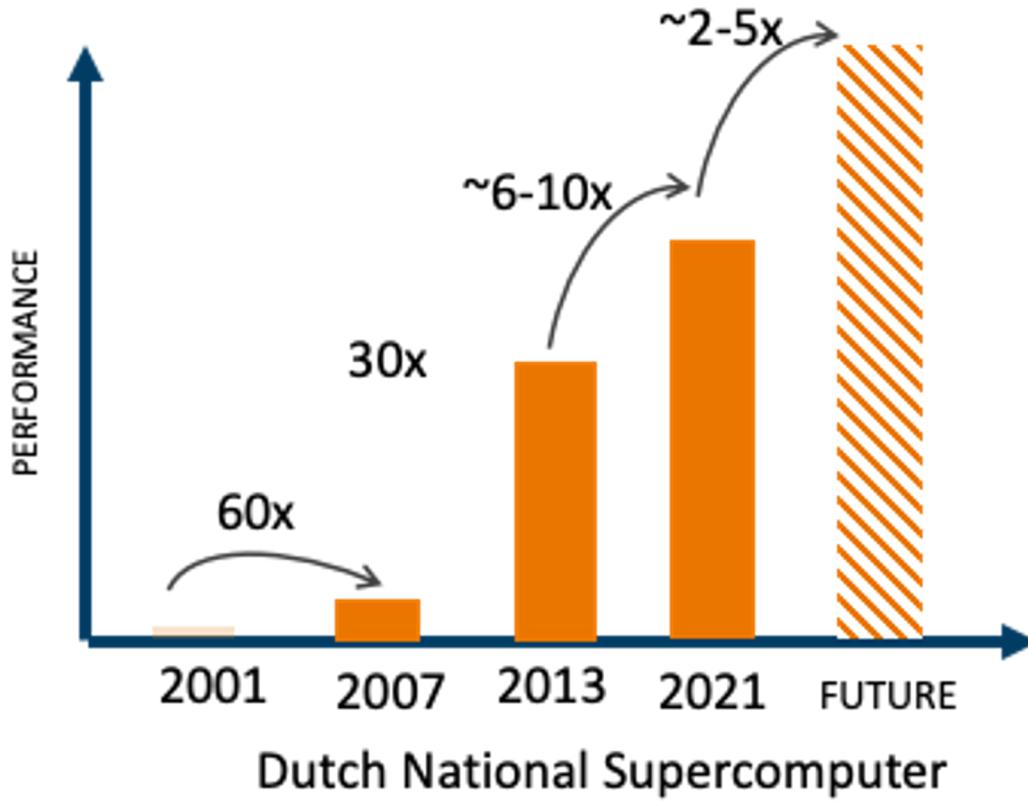
Snellius

Phase 1:

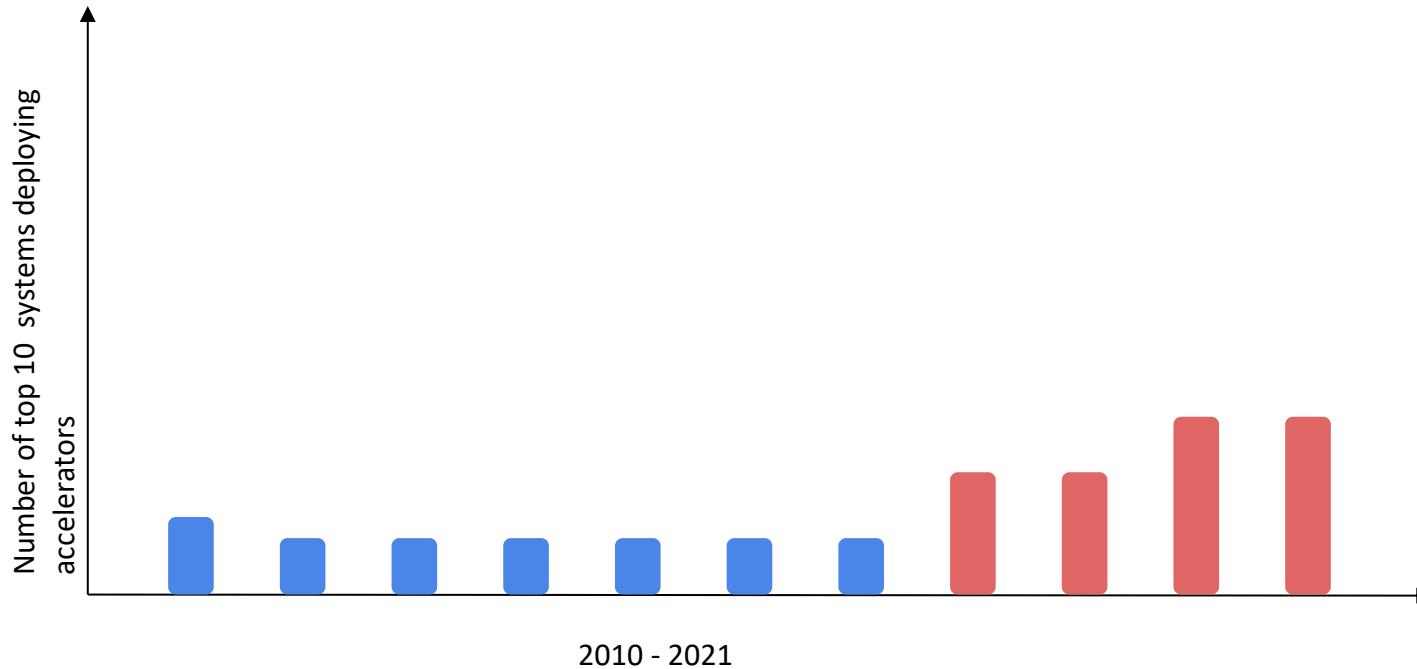
- 76,832 cores (1.6 ×)
- 144 GPUs (3 Pflop/s, 14 ×)
- Total peak: 6.1 Pflop/s (3.4 ×)

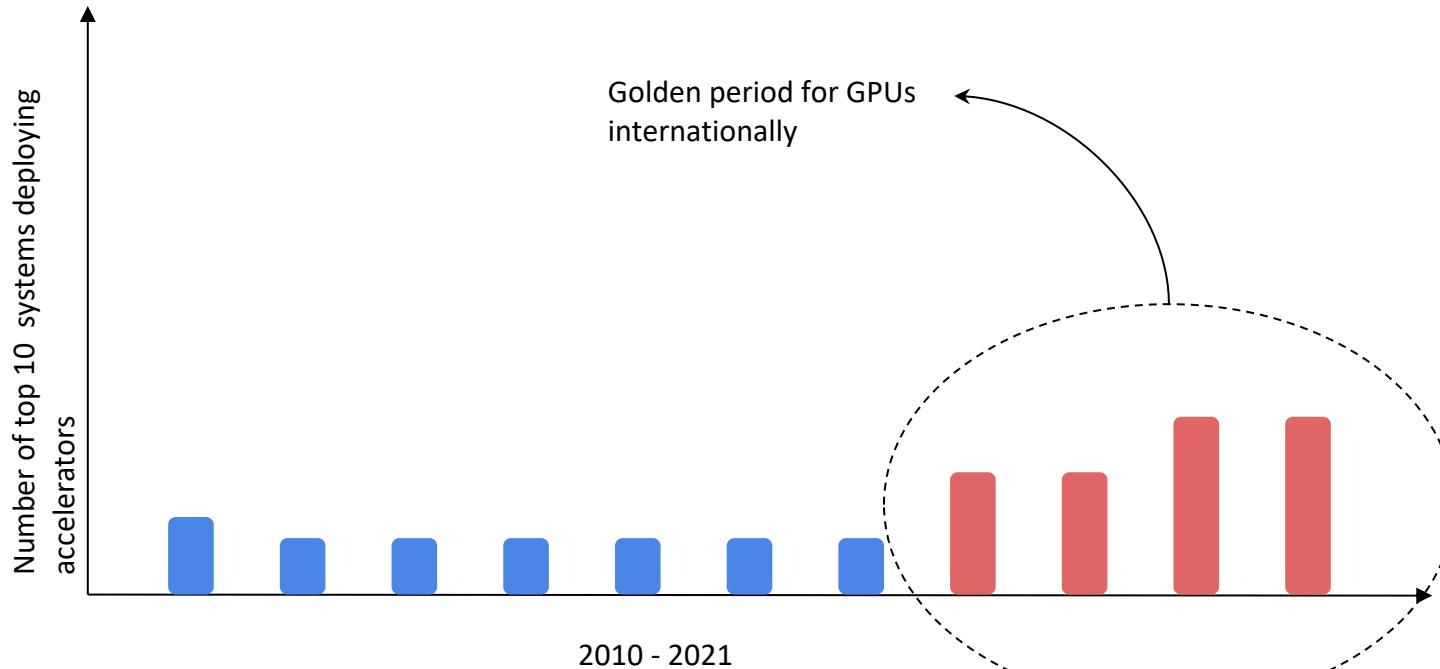
Phase 1 and 2

- Total peak: 11.2 Pflop/s (6.2 ×)
- Full system, based on choice for Phase 3:
 - > 200,000 cores (> 4 ×)
 - Total peak: 13.6 – 21.5 Pflop/s (7.6 – 11.9 ×)



Let's look at the supercomputing trends
with respect to GPUs





Let's look at the pre-exascale systems & exascale system architecture trends



1.5 ExFLOPs/s : AMD CPUs
and GPUs



2 ExFLOPs/s : AMD CPUs and GPUs

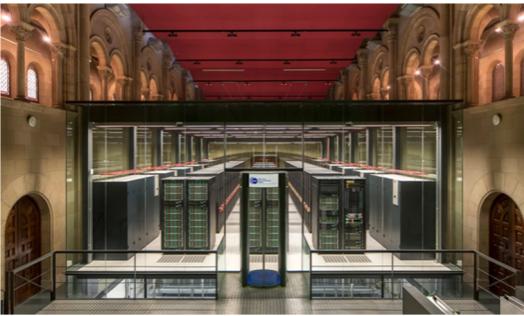


~2 ExFLOPs/s : Intel CPUs and
GPUs



1.024 x 512

~350 PetaFLOPs: AMD CPUs
and GPUs



N/A

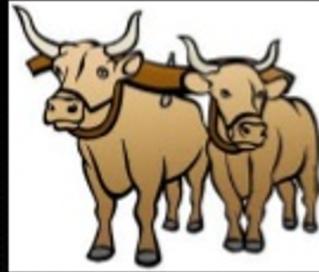
~250 PetaFLOPs : Intel CPUs
and NVIDIA GPUs



What's the difference between CPU & GPU

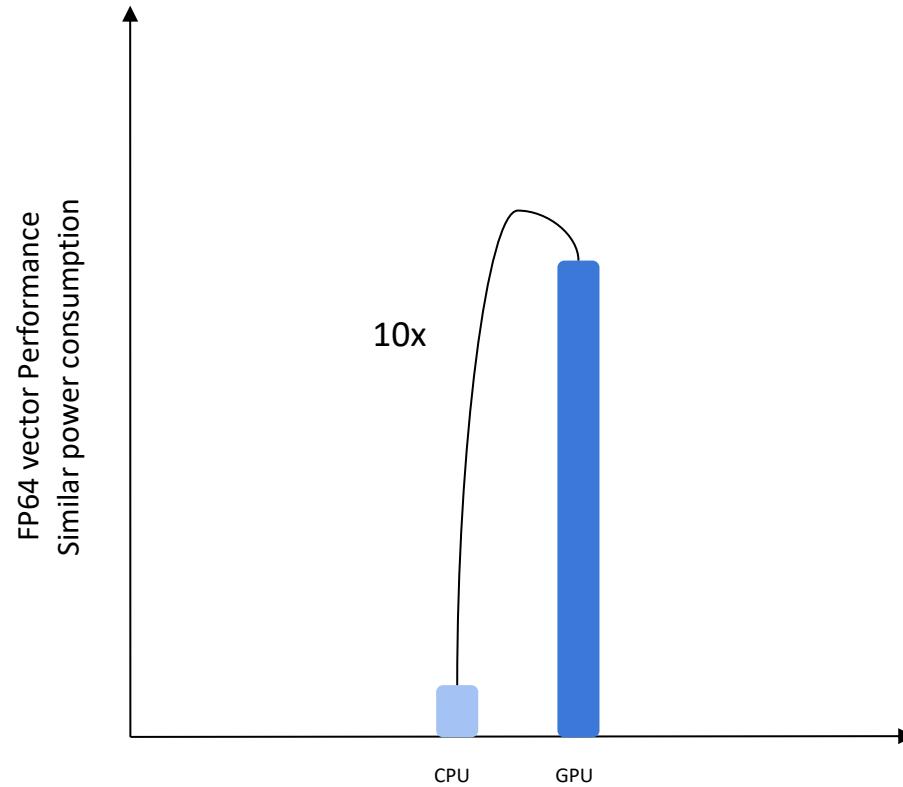
“If you were plowing a field, which would you rather use? Two strong oxen or 1024 chickens?”

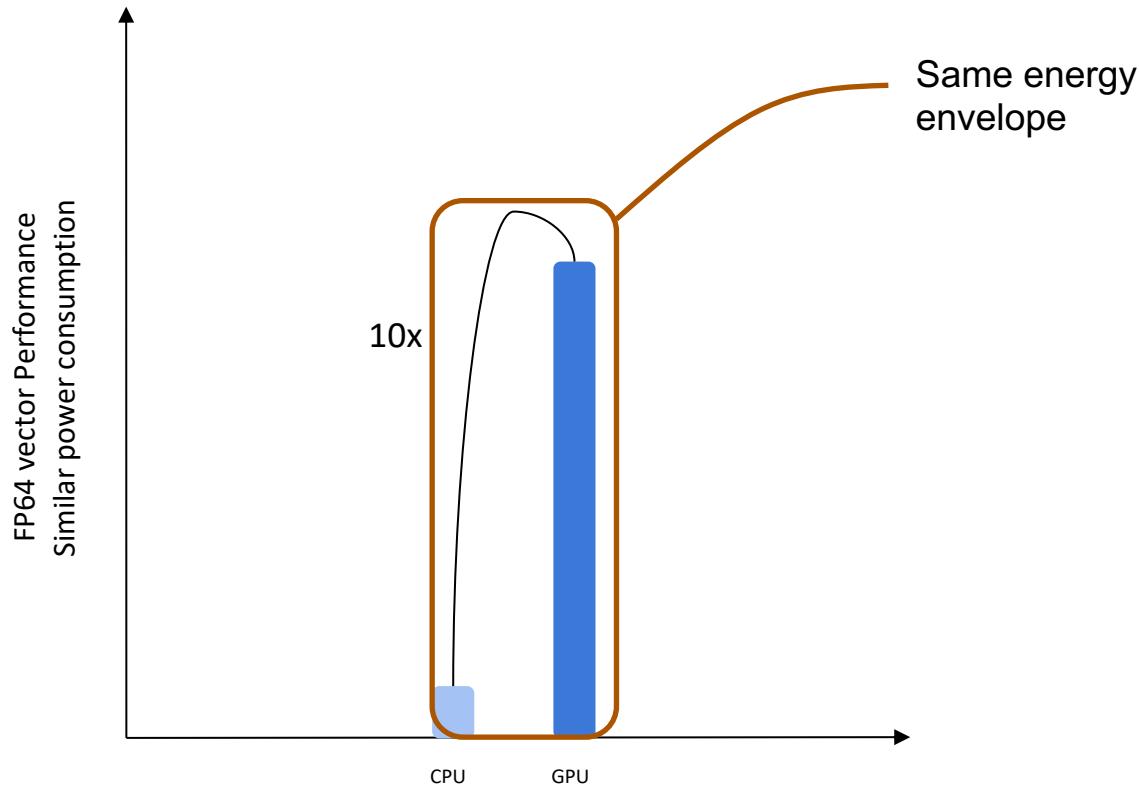
—Seymour Cray



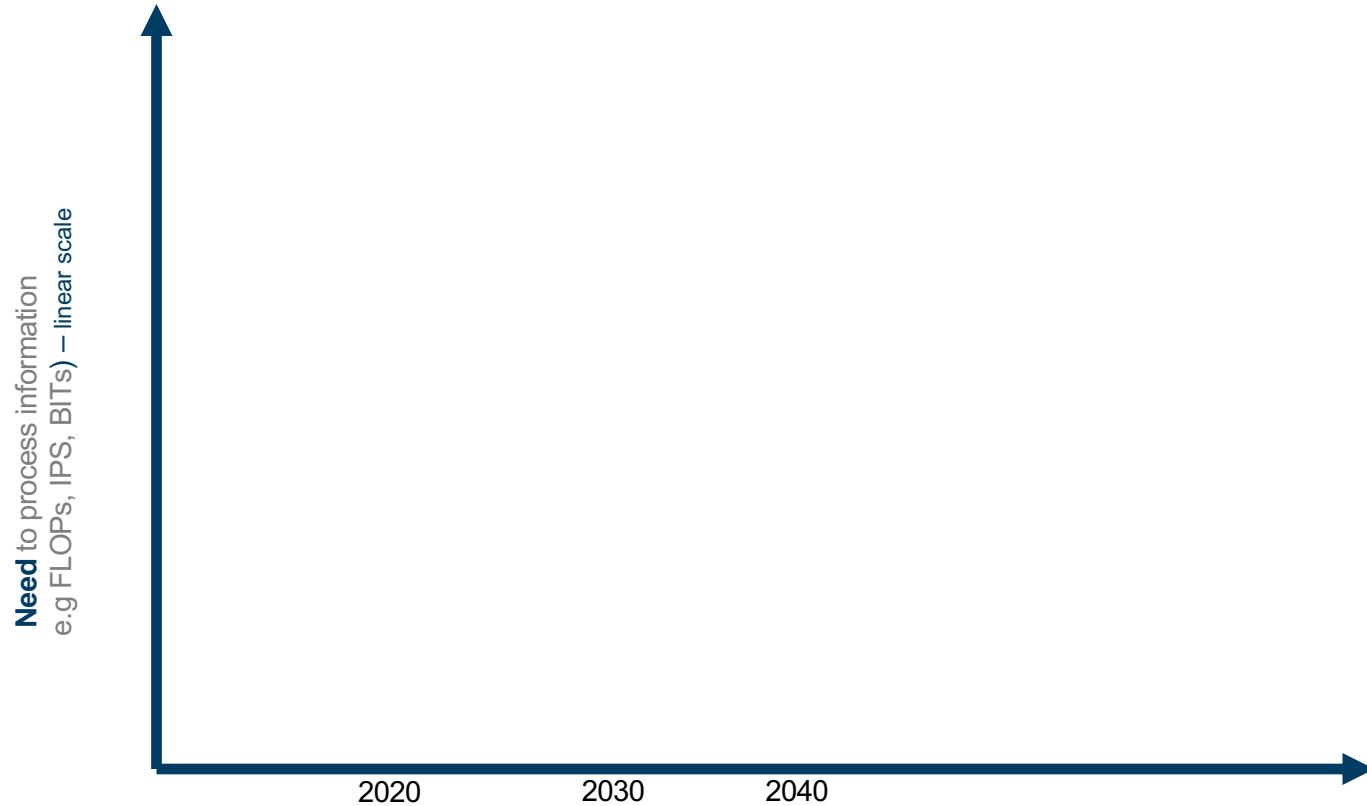
CPU are much better at managing lot of tasks

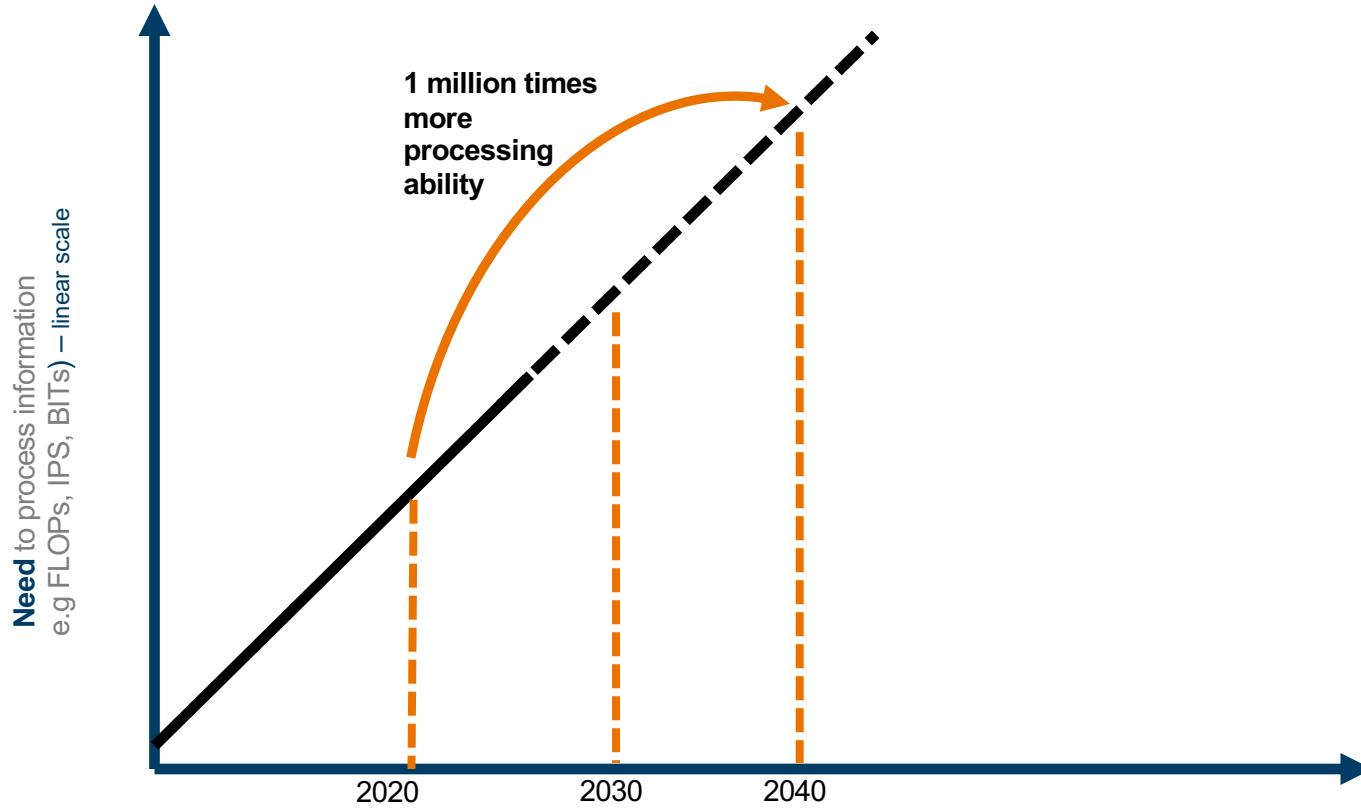
GPUs are specialised for particular workloads involving lot of threads & parallelism

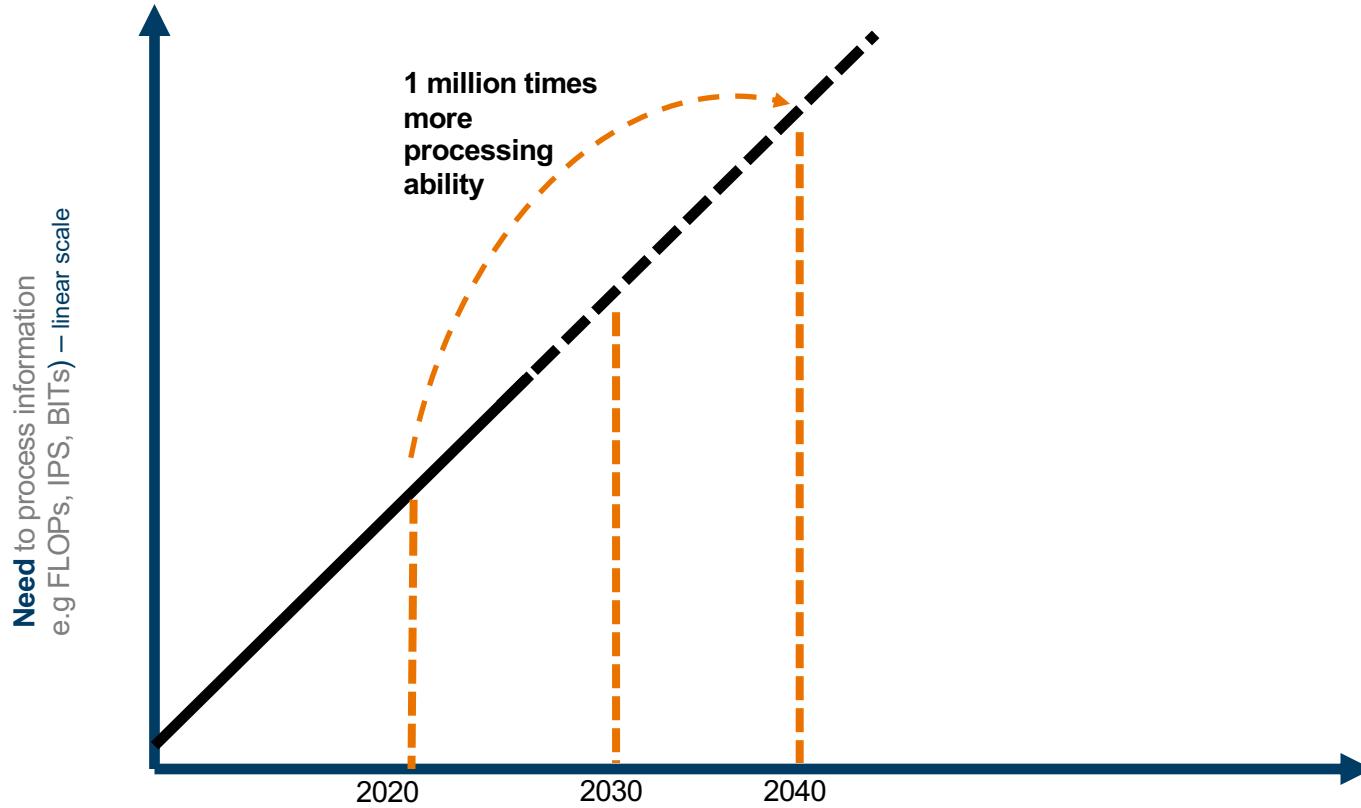




Let's look at the need for computing

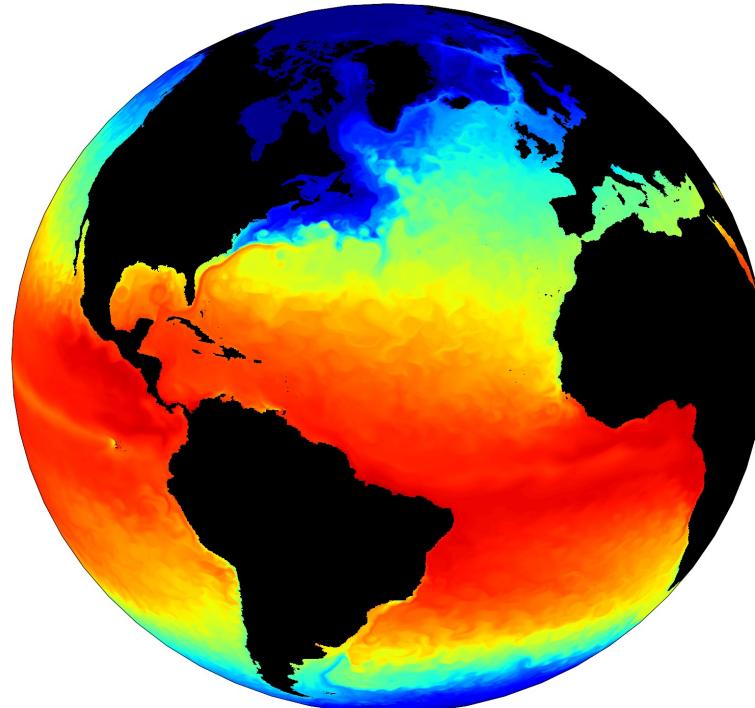






Think about **global**
climate simulations

1km x 1km x 1km
resolution



Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity

William Fedus*

LIAMFEDUS@GOOGLE.COM

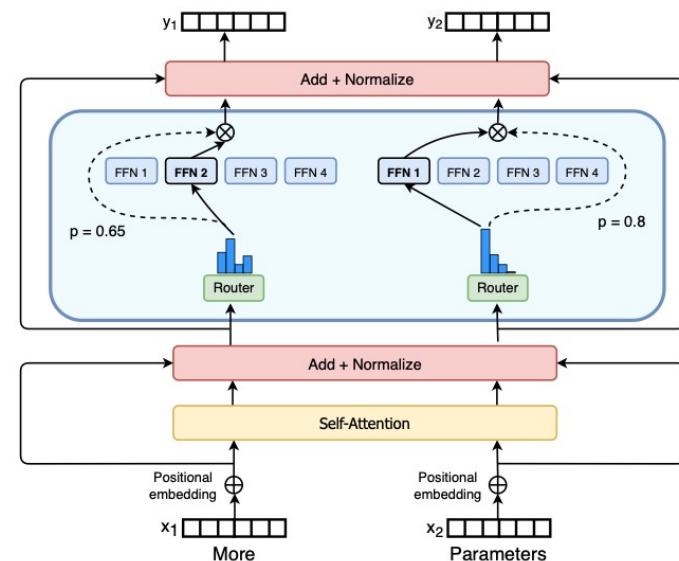
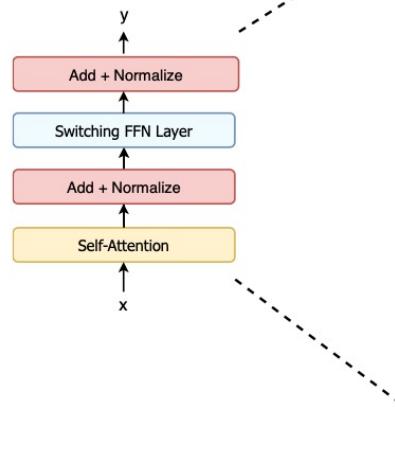
Barret Zoph*

BARRETZOPH@GOOGLE.COM

Noam Shazeer

NOAM@GOOGLE.COM

Google, Mountain View, CA 94043, USA



Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity

William Fedus*

LIAMFEDUS@GOOGLE.COM

Barret Zoph*

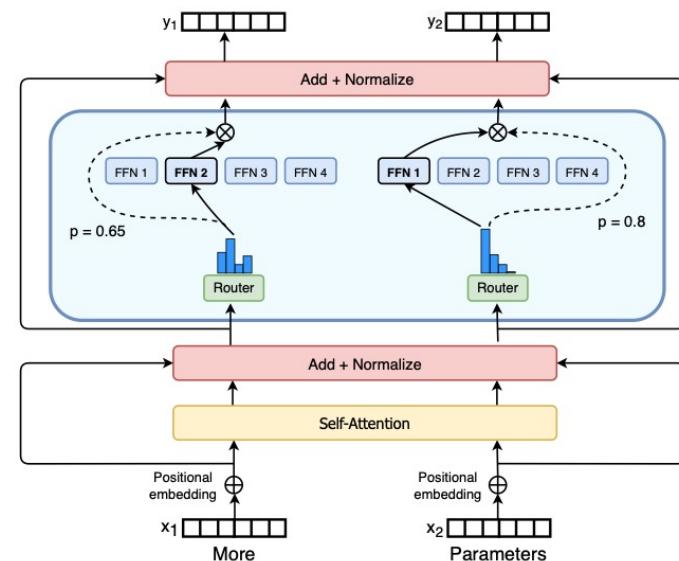
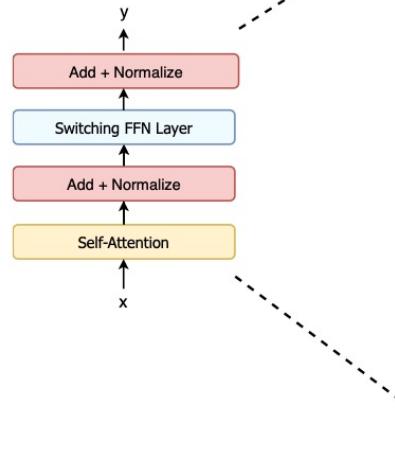
BARRETZOPH@GOOGLE.COM

Noam Shazeer

NOAM@GOOGLE.COM

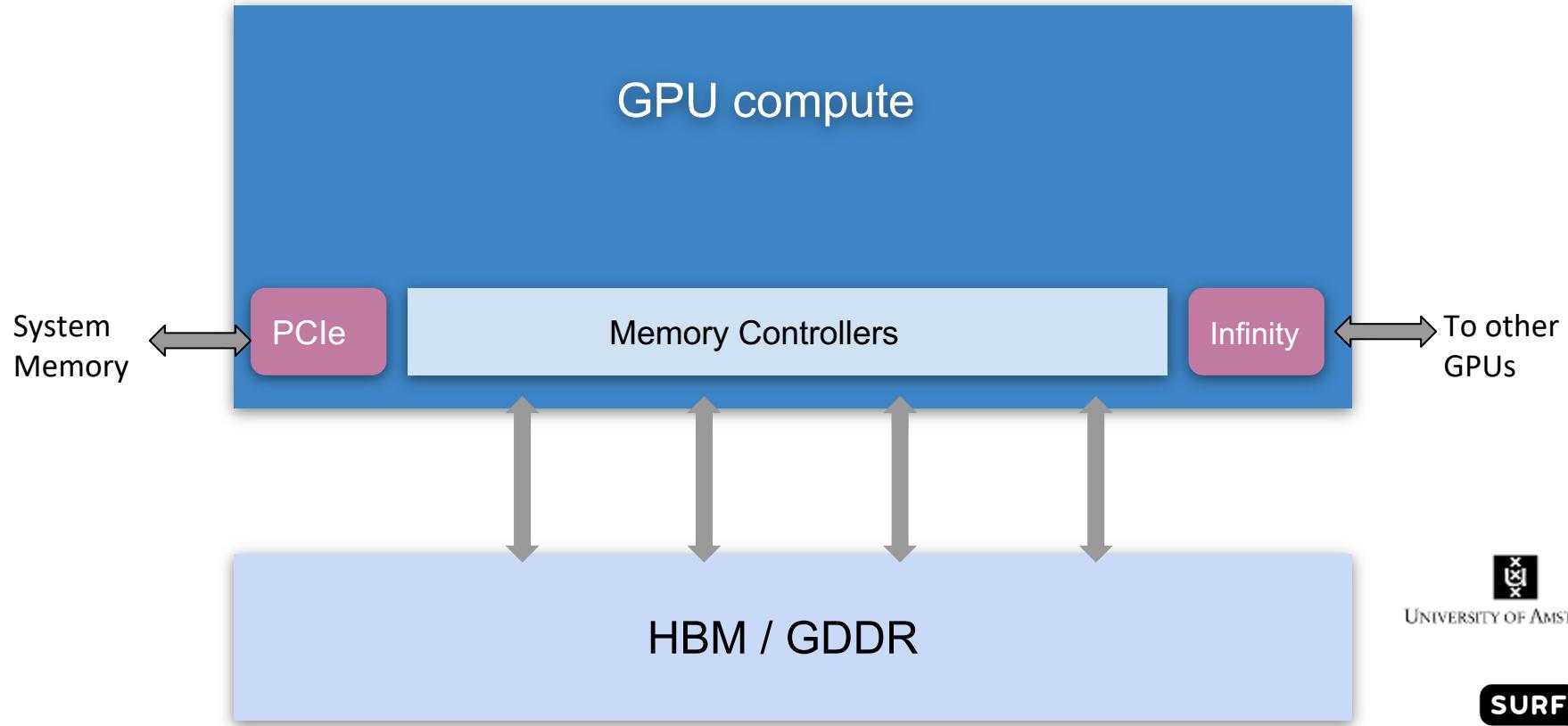
Google, Mountain View, CA 94043, USA

Model parameters
doubling every 4
months or so !!!!

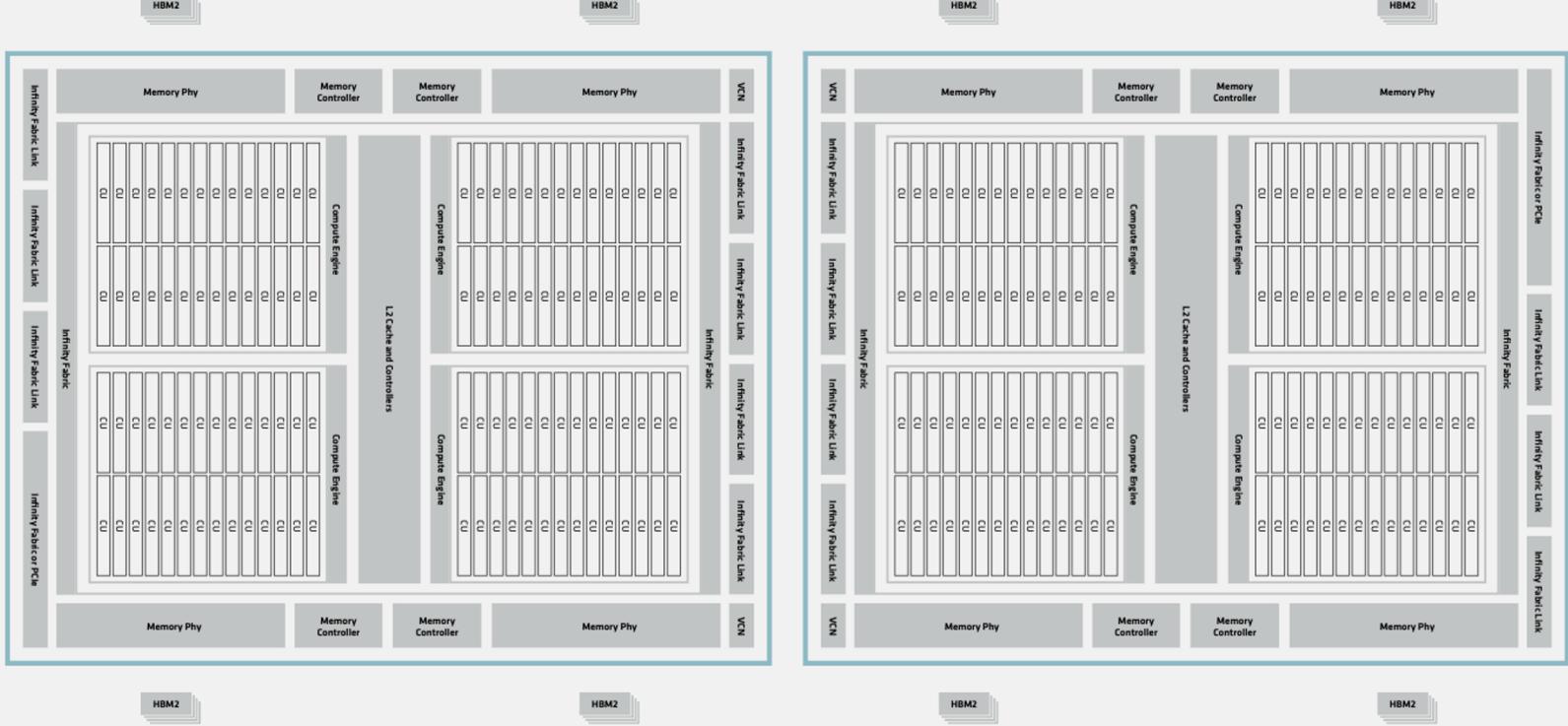


Let's look at the GPU architecture

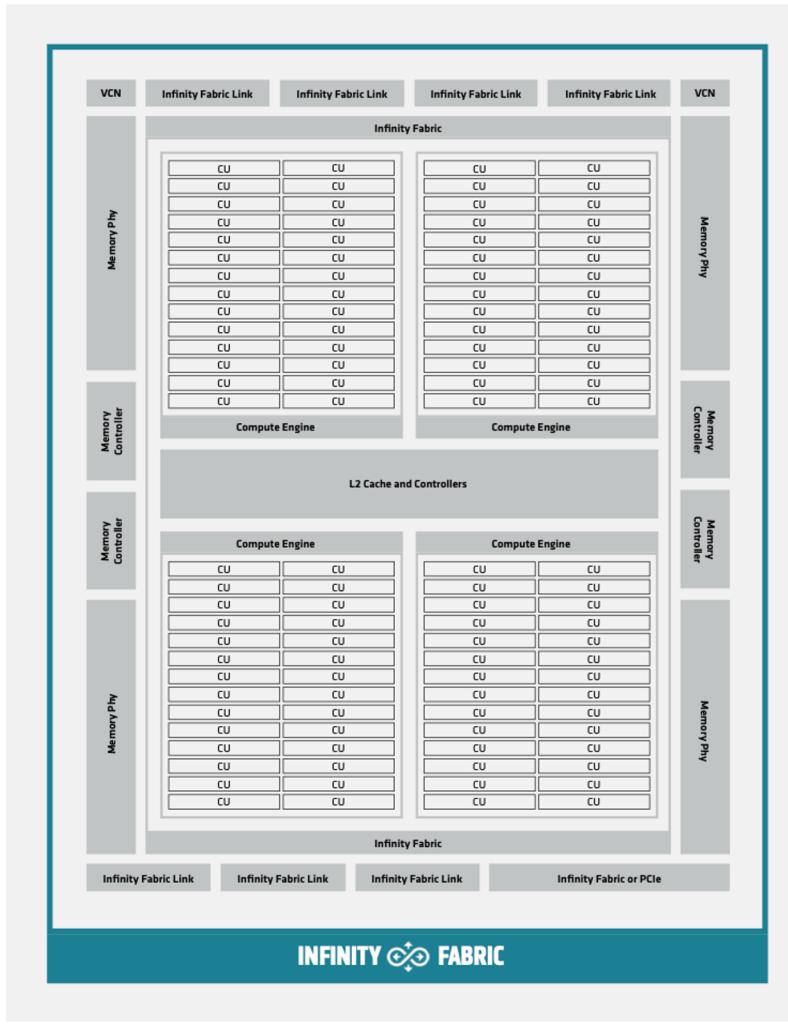
GPU Memory subsystem



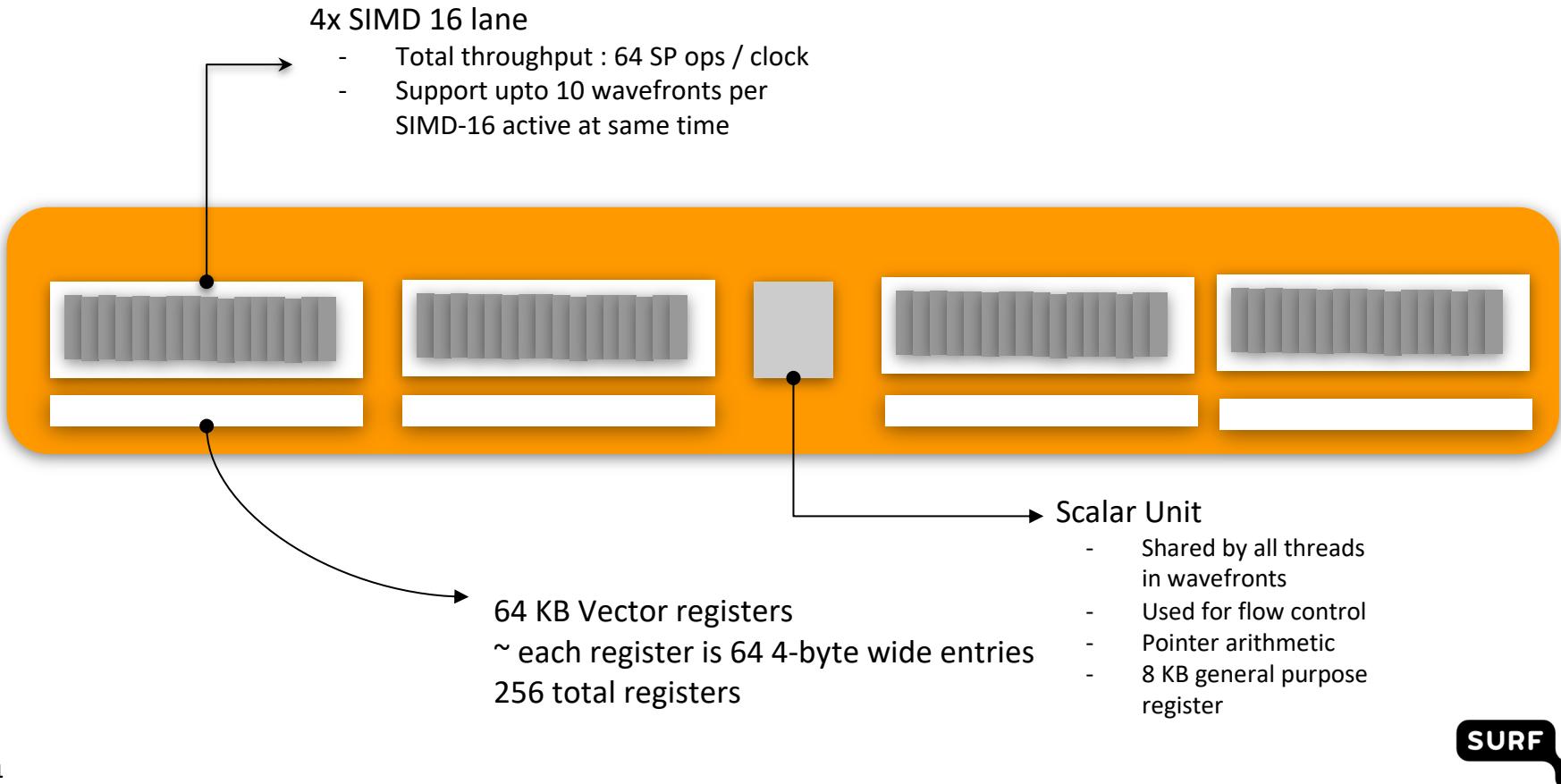
AMD GPU architecture



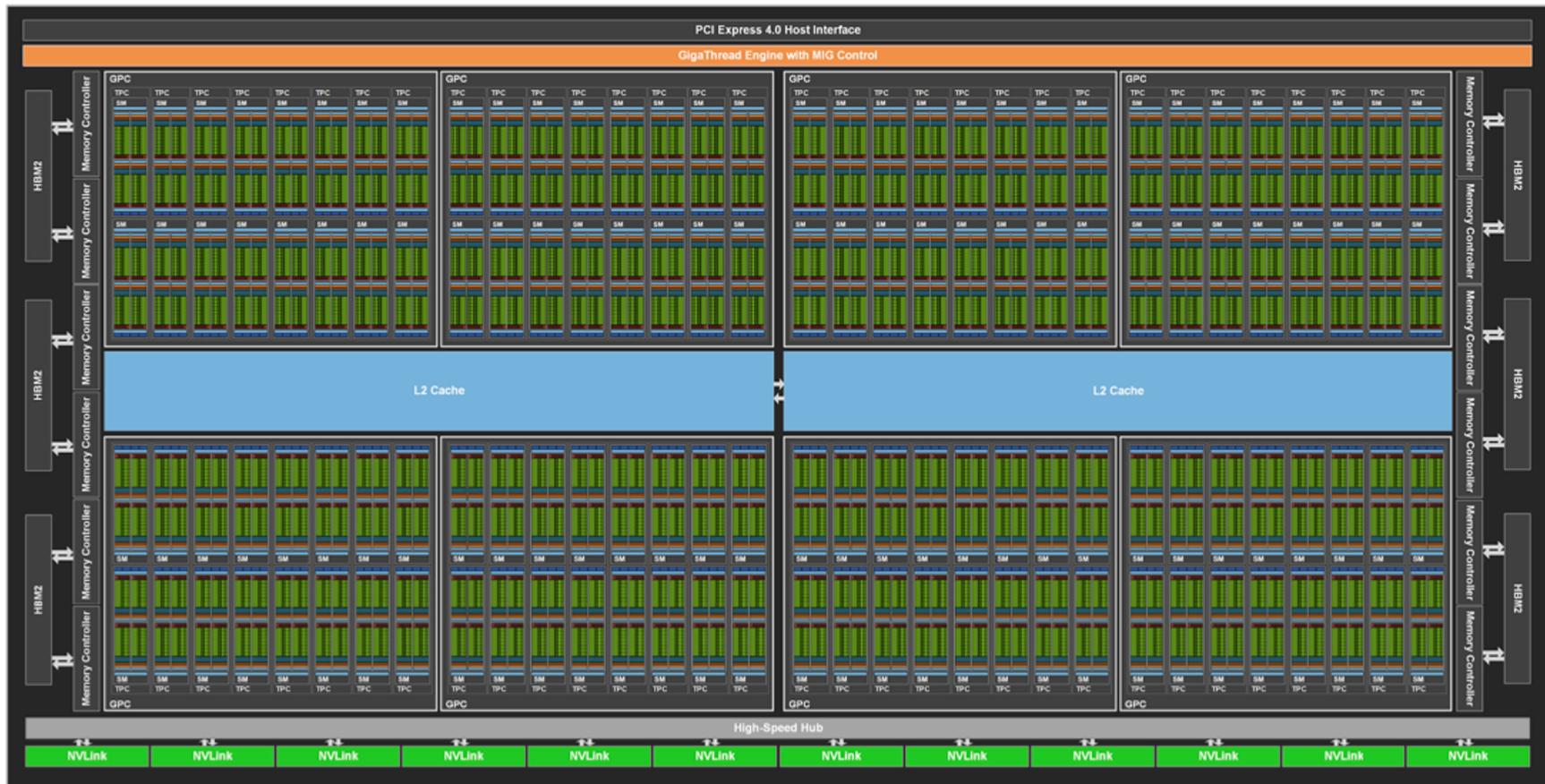
INFINITY FABRIC



Compute Unit



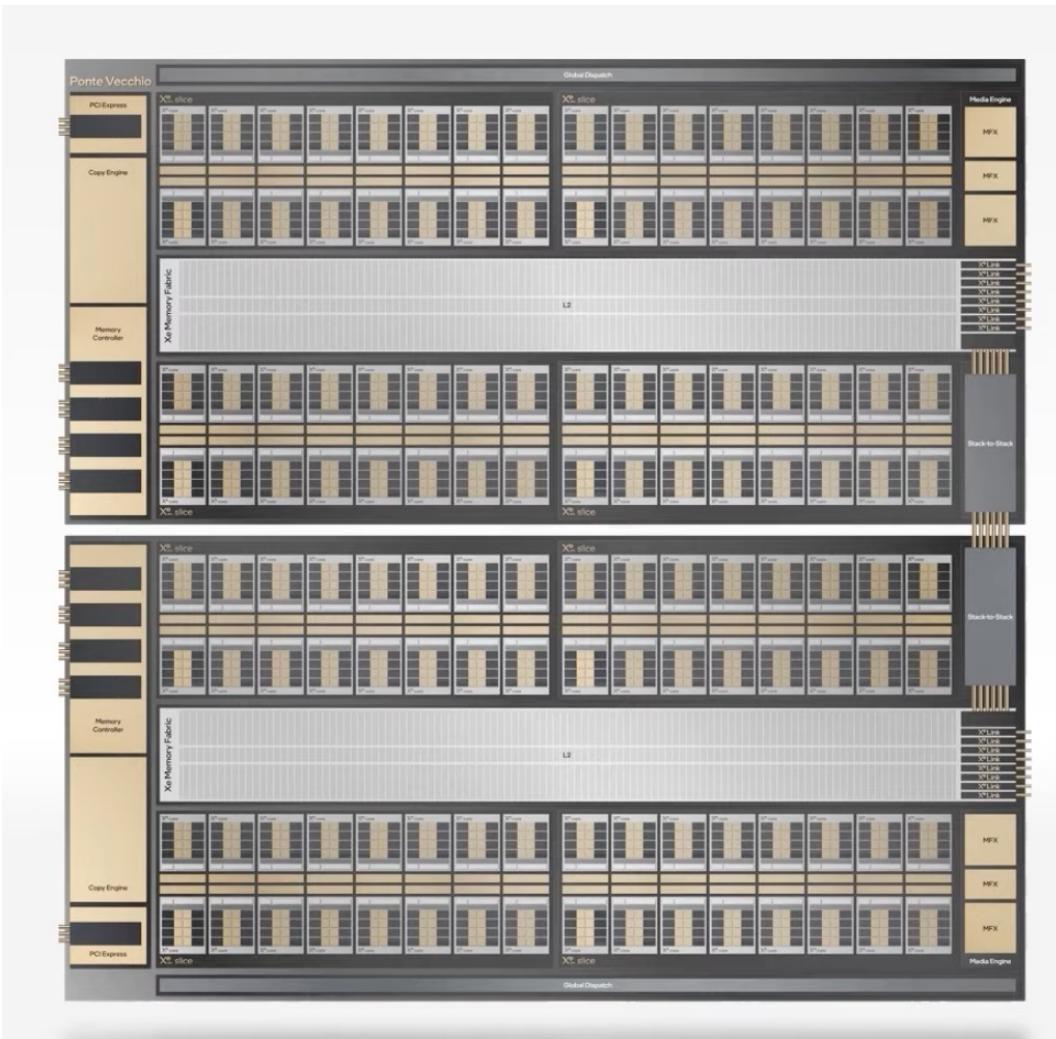
NVIDIA GPU architecture

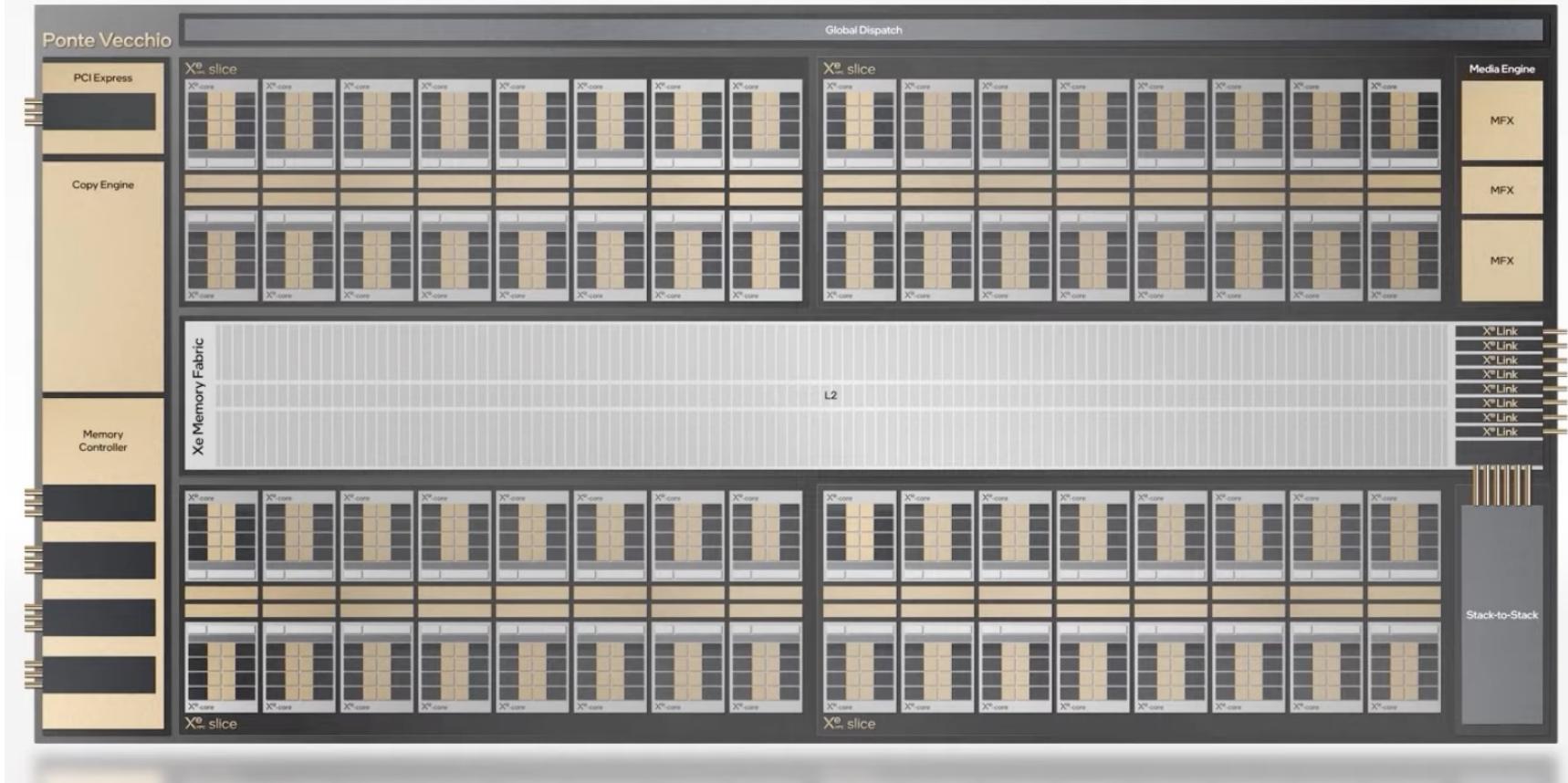


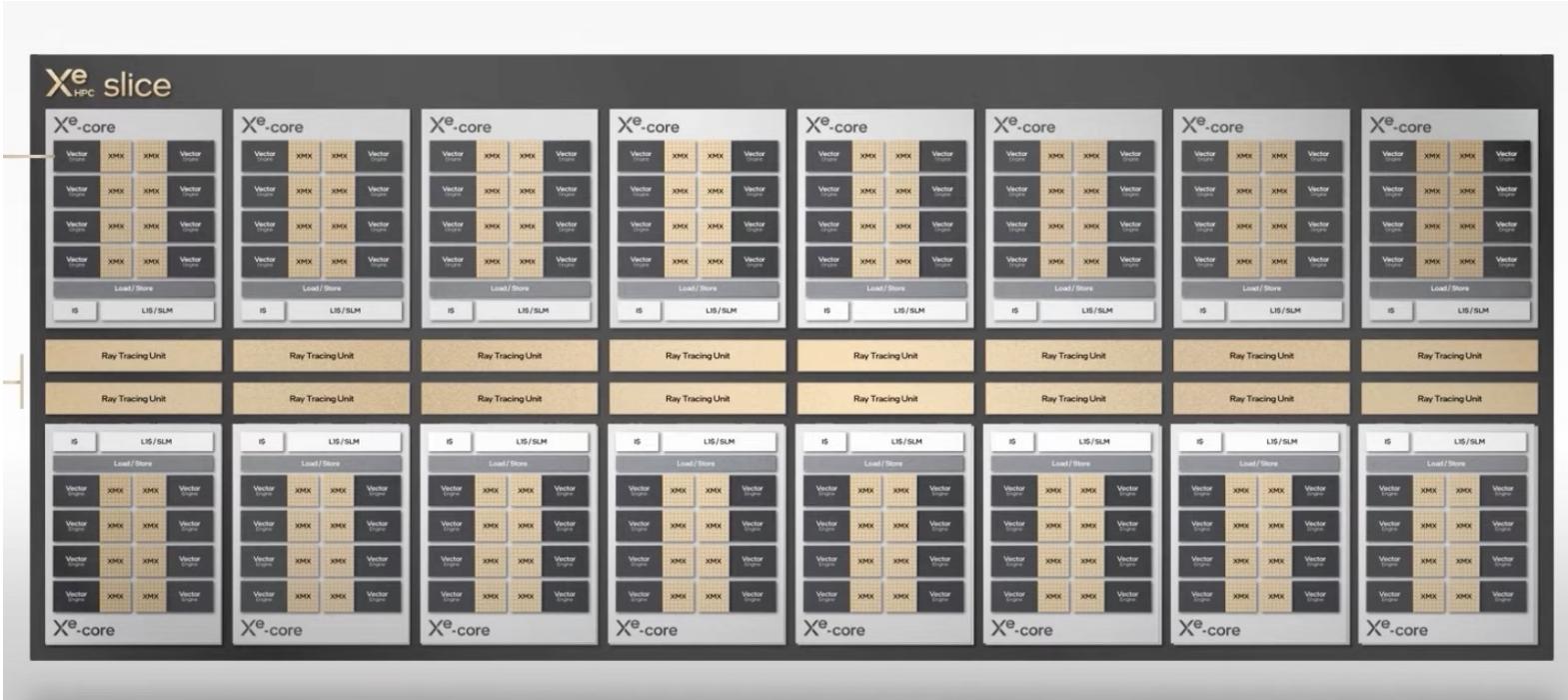
SM

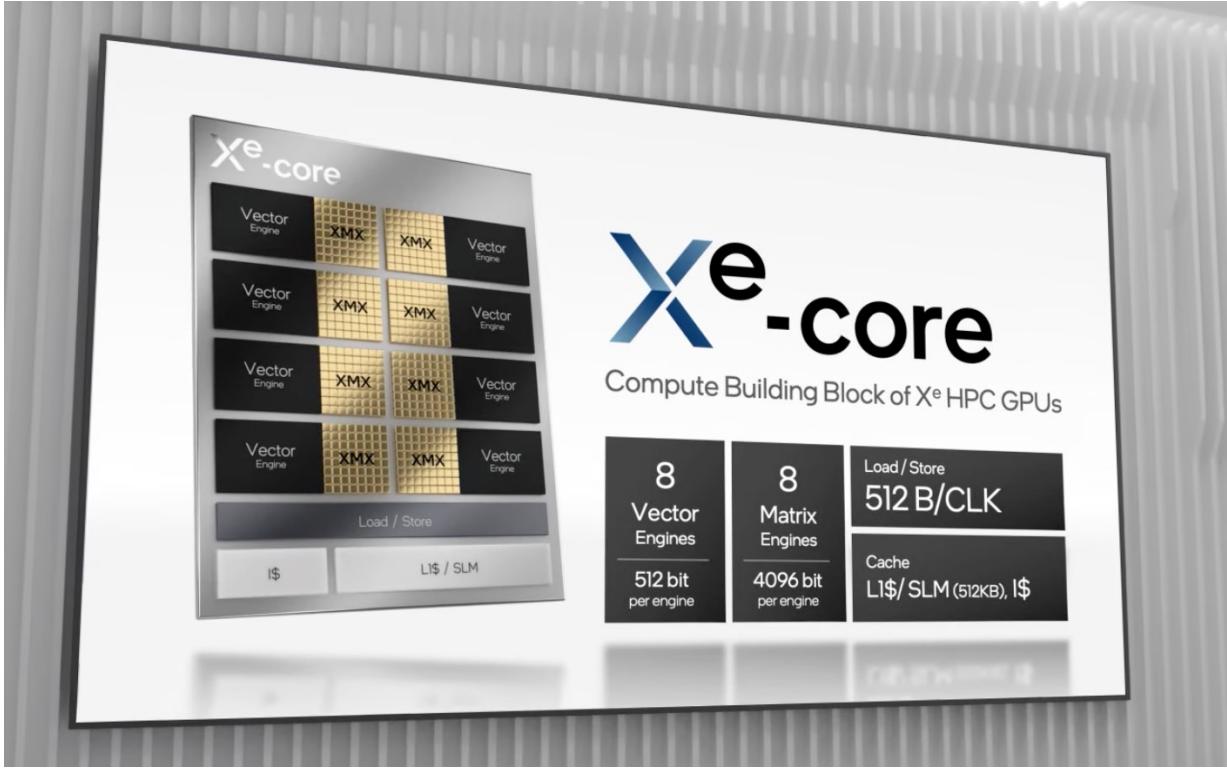


Intel GPU architecture

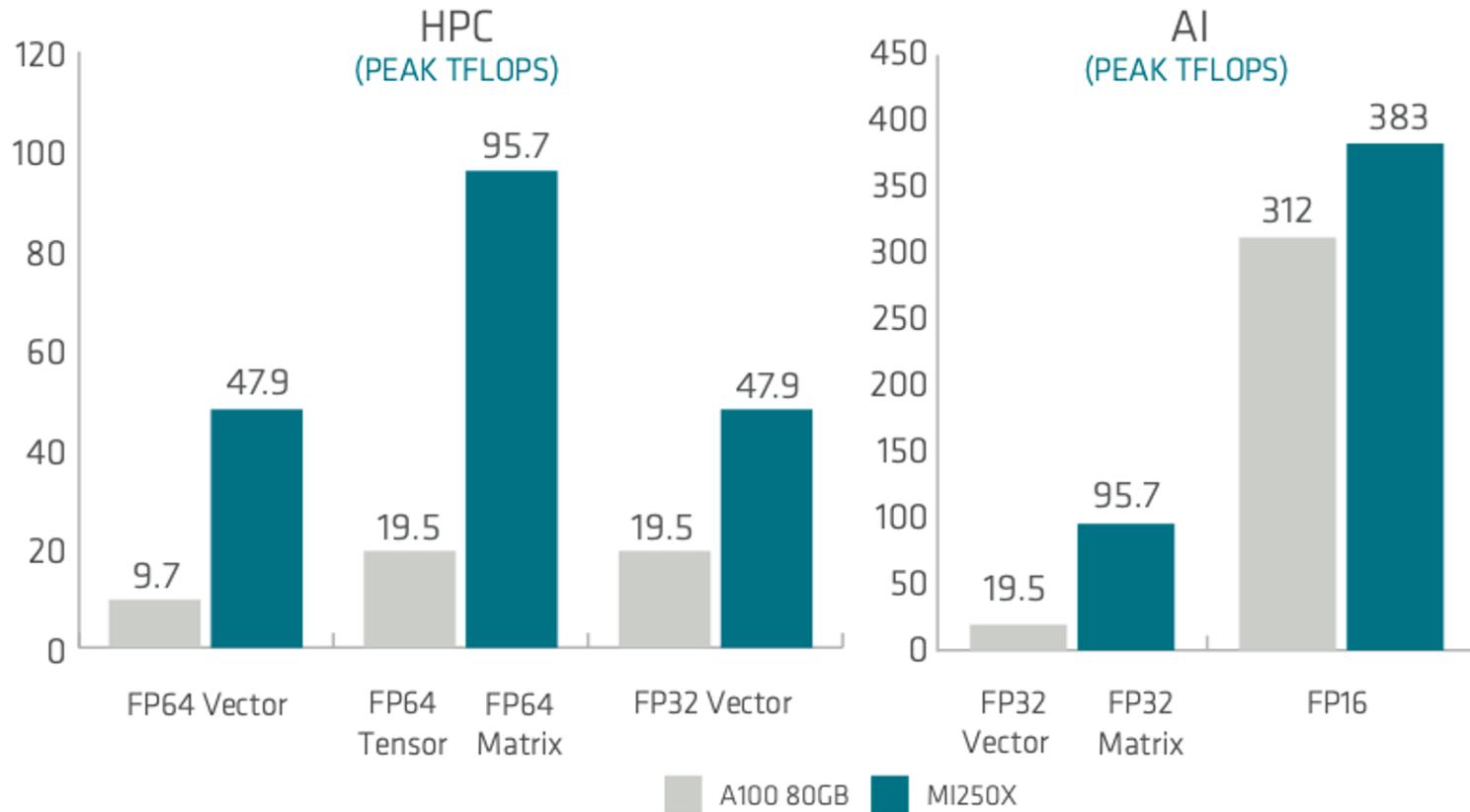








Some Performance Numbers



Ponte Vecchio

Execution Progress



A0 Silicon Current Status

> 45 TFLOPS

FP32 Throughput

> 5 TBps

Memory Fabric Bandwidth

> 2 TBps

Connectivity Bandwidth

For workloads and configurations visit www.intel.com/ArchDay21Claims. Results may vary.

Architecture Day 2021

intel. 10

Programming GPUs

3 Ways to Accelerate Applications

Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily Accelerate
Applications

Programming
Languages

Maximum
Flexibility

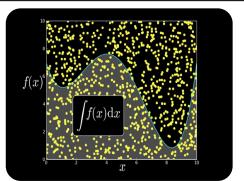
Libraries: Easy, High-Quality Acceleration

- **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 (replacing MKL/IPP/FFTW/...)
- **Quality:** Libraries offer high-quality implementations of functions encountered in a broad range of applications
- **Performance:** NVIDIA libraries are tuned by experts

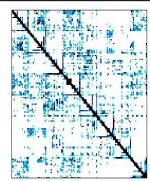
Some GPU-accelerated Libraries



NVIDIA cuBLAS



NVIDIA cuRAND



NVIDIA cuSPARSE



NVIDIA NPP



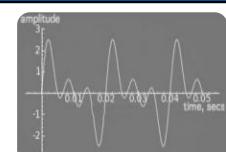
Vector Signal
Image Processing



GPU
Accelerated
Linear Algebra



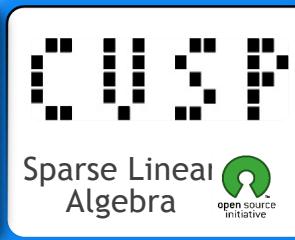
Matrix Algebra
on GPU and
Multicore
open source
initiative



NVIDIA cuFFT



ArrayFire
Matrix
Computations



Sparse Linear
Algebra
open source
initiative



C++ STL
Features for
CUDA
open source
initiative

3 Steps to CUDA-accelerated application

- **Step 1:** Substitute library calls with equivalent CUDA library calls

saxpy (...) ▶ cublasSaxpy (...)

- **Step 2:** Manage data locality

- with CUDA: `cudaMalloc()`, `cudaMemcpy()`, etc.
 - with CUBLAS: `cublasAlloc()`, `cublasSetVector()`, etc.

- **Step 3:** Rebuild and link your CUDA Library-accelerated application

`nvcc myobj.o -l cublas`

3 Ways to Accelerate Applications

Applications

Libraries

“Drop-in”
Acceleration

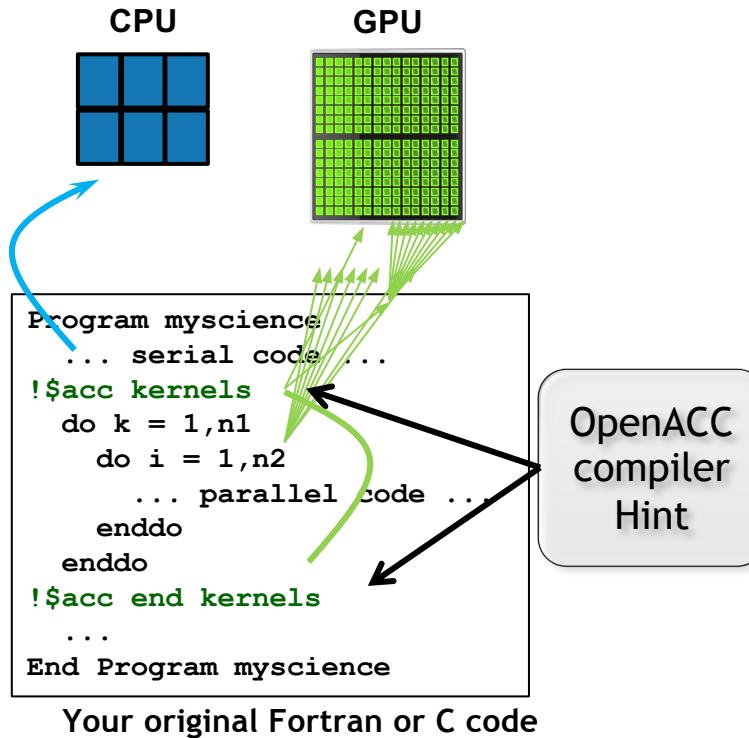
OpenACC Directives

Easily Accelerate
Applications

Programming Languages

Maximum
Flexibility

OpenACC Directives



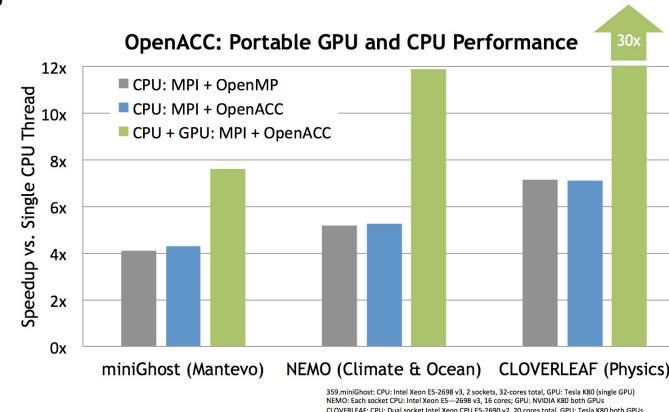
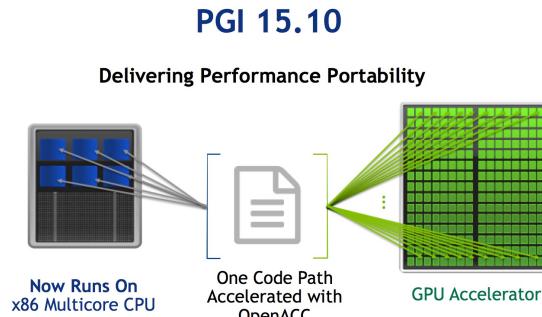
Simple Compiler hints

Compiler Parallelizes code

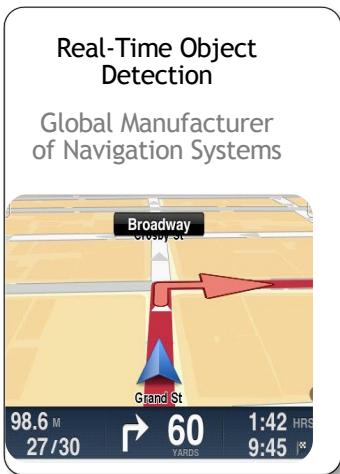
Works on many-core GPUs & multicore CPUs

OpenACC: The Standard for GPU Directives

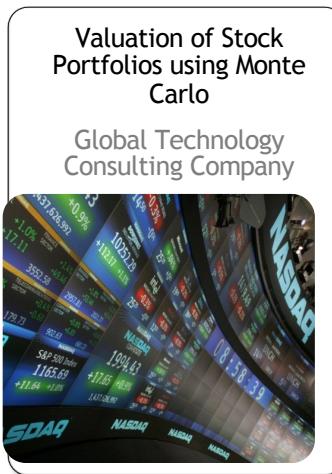
- **Easy:** Directives are the easy path to accelerate compute intensive applications
- **Open:** OpenACC is an open GPU directives standard, making GPU programming straightforward and portable across parallel and multi-core processors
- **Powerful:** GPU Directives allow complete access to the massive parallel power of a GPU



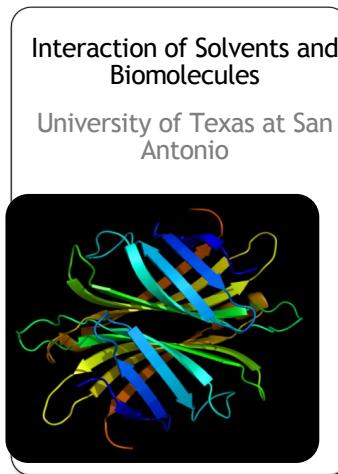
Directives: Easy & Powerful



5x in 40 Hours



2x in 4 Hours



5x in 8 Hours

“Optimizing code with directives is quite easy, especially compared to CPU threads or writing CUDA kernels.
The most important thing is avoiding restructuring of existing code for production applications.”

-- Developer at the Global Manufacturer of Navigation Systems

SURF

3 Ways to Accelerate Applications

Applications

Libraries

“Drop-in”
Acceleration

OpenACC Directives

Easily Accelerate
Applications

Programming Languages

Maximum
Flexibility

GPU Programming Languages

Numerical analytics ➤

MATLAB, Mathematica, LabVIEW

Fortran ➤

CUDA Fortran, OpenACC,
OpenMP4.5

C ➤

CUDA C, OpenCL, OpenACC,
OpenMP4.5

C++ ➤

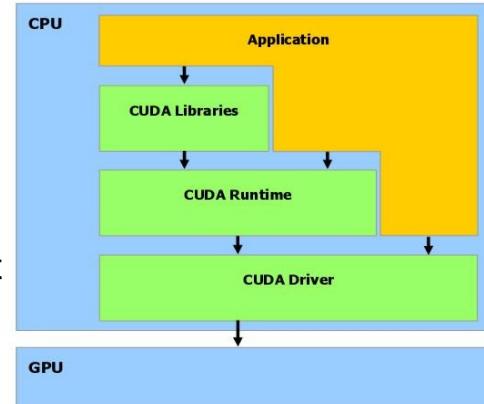
CUDA C++, Thrust, OpenCL,
OpenACC/OpenMP4.5

Python ➤

PyCUDA/PyOpenCL, Numba, ...

CUDA

- “Compute Unified Device Architecture”
- General purpose programming model
 - User kicks off batches of threads on the GPU
 - GPU = dedicated super-threaded, massively data parallel co-processor
- Targeted software stack
- Driver for loading computation programs into GPU
 - Standalone Driver - Optimized for computation
 - Explicit GPU memory management

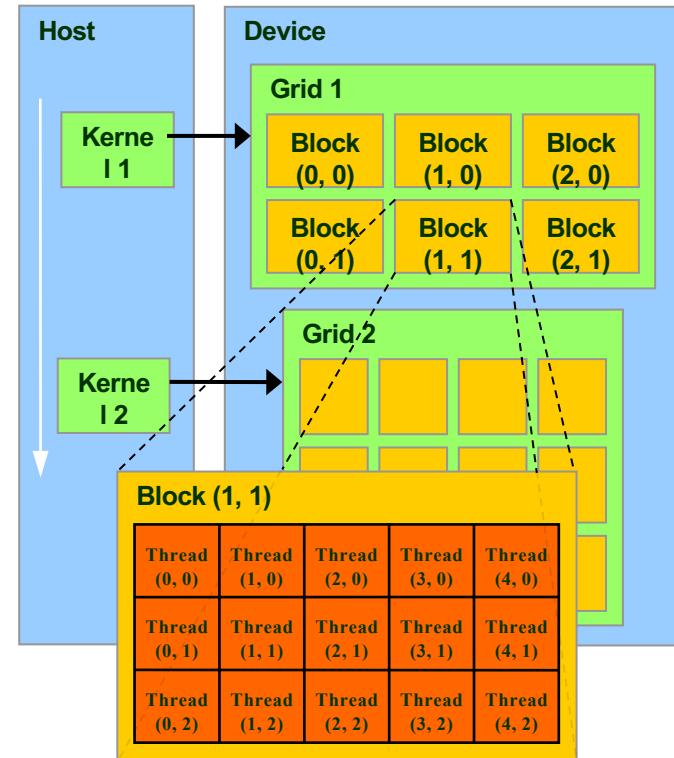


CUDA Programming Model

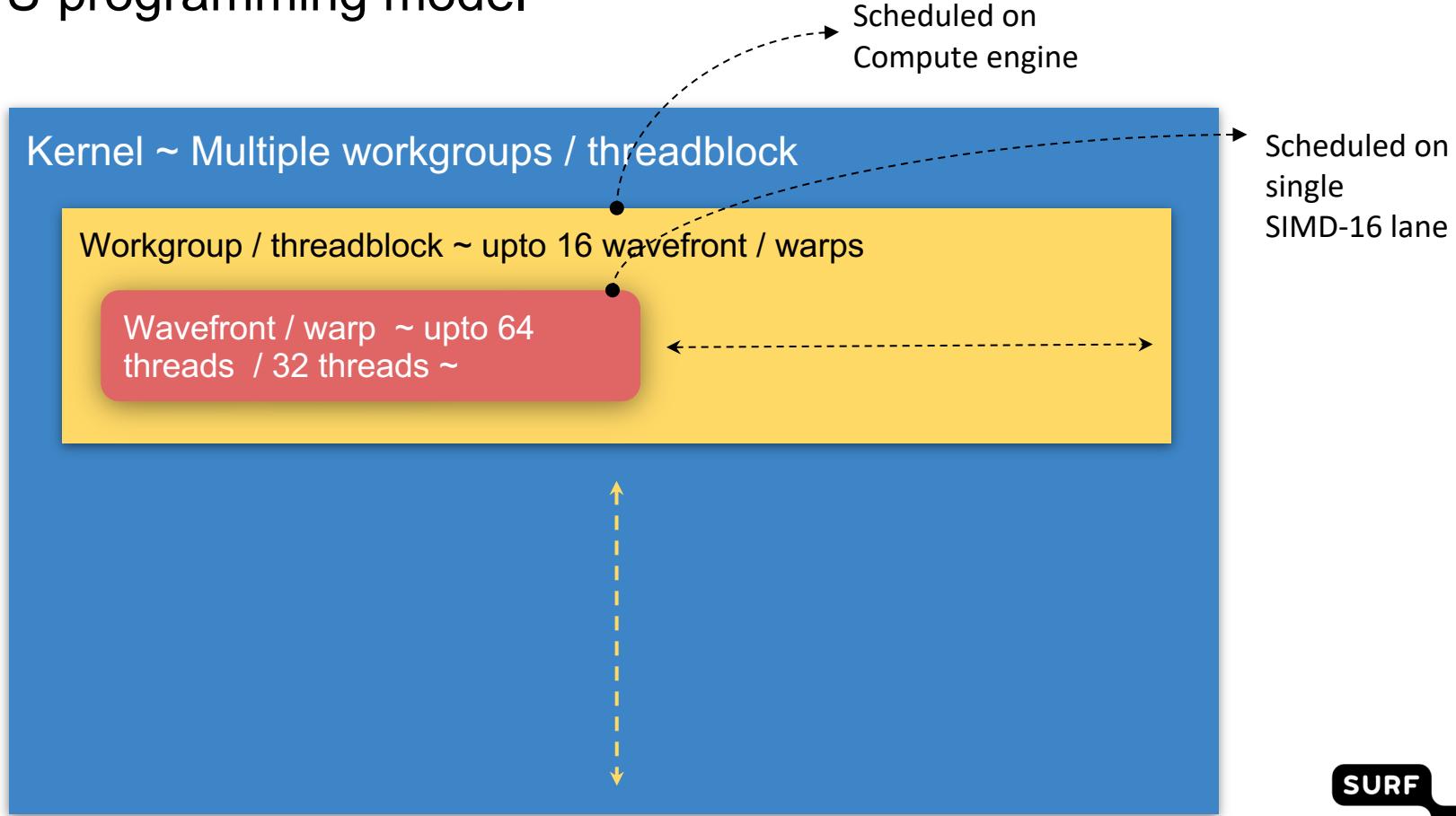
- The GPU is viewed as a compute **device** that:
 - Is a coprocessor to the CPU or **host**
 - Has its own DRAM (**device memory**)
 - Runs many **threads in parallel**
 - - Hardware switching between threads (in 1 cycle) on long-latency memory reference
 - - **Overprovision** (1000s of threads) → hide latencies
 - Data-parallel portions of an application are executed on the device as kernels which run in parallel on many threads
 - Differences between GPU and CPU threads
 - - GPU threads are extremely lightweight
 - - Very little creation overhead
 - - GPU needs 1000s of threads for full efficiency
 - - Multi-core CPU needs only a few

Thread Batching: Grids and Blocks

- Kernel executed as a **grid of thread blocks**
 - - All threads share data memory space
- **Thread block** is a batch of threads, can **cooperate** with each other by:
 - - Synchronizing their execution:
 - For hazard-free shared memory accesses
 - - Efficiently sharing data through the low latency **shared memory**
- Two threads from two different blocks cannot cooperate (**until CUDA8/9 and Volta**)
 - - Unless thru slow global memory
- Threads and blocks have IDs



GPU programming model



Tensorflow, PyTorch

Gromacs
Molecular Dynamics

Libraries

Compiler toolchain

Programming models (CUDA, HIP)

Driver (runtime system)

Novel / Unconventional accelerated computing paradigm

- Quantum Technology (Qubits)
- Brain-inspired computing (Neuromorphic paradigm)
- Mixed precision algorithm (floating point + Integer precision + bfloat + etc)
- Analog computing
- Specialized computing (Google TPUs, SambaNova, Cerebras, Groq, Graphcore, Habana AI, Intel Loihi, IBM Analog systems,)
- Application design (AI + HPC)
- Edge + cloud + HPC : Computing continuum

- **All images available in the presentation are present in the public domain**

SURF

Thank you