# Introduction to GPU Computing

N-ways bootcamp



# Introduction to gpu computing

## What to expect?

- Overview of GPU Architecture
- Approaches to getting started with GPUs in your work
- Intro to programming GPUs



2



### Universe of GPU Computing

## **Rise of GPU Computing**



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

# ACCELERATED COMPUTING

- -



## **GPU Accelerator**

Optimized for Parallel Tasks





# SILICON BUDGET

The three components of any processor





# CPU IS A LATENCY REDUCING ARCHITECTURE

### **CPU** Optimized for Serial Tasks



## GPU CPU Strengths tor

- Very large main memory
- Very fast clock speeds
- Latency optimized via large caches
- Small number of threads can run very quickly

#### CPU Weaknesses

- Relatively low memory bandwidth
- Cache misses very costly
- Low performance/watt



# GPU IS ALL ABOUT HIDING LATENCY

#### **GPU** Strengths

- High bandwidth main memory
- Significantly more compute resources
- Latency tolerant via parallelism
- High throughput
- High performance/watt

GPU Weaknesses

- Relatively low memory capacity
- Low per-thread performance

## **GPU Accelerator**

Optimized for Parallel Tasks





# LOW LATENCY VS HIGH THROUGHPUT

- CPU architecture must minimize latency within each thread
- GPU architecture hides latency with computation (data-parallelism, to 30k threads!)





# SPEED V. THROUGHPUT

## Speed

## Throughput





Which is better depends on your needs...





#### \*Images from Wikimedia Commons via Creative Commons

# HOW GPU ACCELERATION WORKS

**Application Code** 



## Heterogeneous Programming



## Execution FLOW - H2D



## Execution FLOW – kernel launch



## Execution FLOW – D2H



## **GH100 GPU architecture**

#### https://resources.nvidia.com/en-us-tensor-core/gtc22-whitepaper-hopper





## GH100 GPU architecture

#### https://resources.nvidia.com/en-us-tensor-core/gtc22-whitepaper-hopper





17









### Accelerating at all scales







20



### What is CUDA?





### It's Not Just C++

OpenACC	OpenMP	Standard C++	Python	Julia	MATLAB	
CUDA FORTRAN	CUDA C++	OpenCL	Ada	Haskell	R	



## Target the Abstraction Layer That Works Best For Your Application

Developer & Application Ecosystem	Frameworks PyTorch, TensorFlow, Jax, Modulus, Triton,	SDKs Medical Devices, Energy, Autonomous Vehicles,	
NVIDIA Accelerated Libraries		adlukiti aufit Math API	
Parallel Languages	OpenACC OpenAIP Standard C++ Pyth CUDA FDRTBAN CUDA C++ OpenCL Ad	on Julia MATLAB a Haskell &	
Compilation Stack	NVVM / I PTX Asse		



## Ways to ACCELERATION





24

## GPU accelerated APPS and frameworks

- All major DL frameworks PyTorch, TensorFlow etc
- Top 15 most used HPC apps globally
- Apps in a huge range of fields
  - Over 3000 apps in total catalogue: <u>link</u>
- Domain-specific frameworks robotics, vis, healthcare, smart cities etc
- <u>GROMACS</u>, <u>VASP</u>, LAMMPS, RELION, QE, NAMD, SPECFEM3D ...

https://developer.nvidia.com/hpc-application-performance



25

# **GPU ACCELERATED MATH LIBRARIES**



KATHONS

# Programming the NVIDIA Platform

CPU, GPU, and Network

ACCELERATED STANDARD ISO C++, ISO Fortran	INCREMENTAL PORTABLE OpenACC, OpenMP	PLATFORM SPECIALIZATION CUDA						
<pre>std::transform(par, x, x+n, y, y, [=](float x, float y){ return y + a*x; } ); do concurrent (i = 1:n) y(i) = y(i) + a*x(i) enddo import cunumeric as np  def saxpy(a, x, y): y[:] += a*x</pre>	<pre>#pragma acc data copy(x,y) { std::transform(par, x, x+n, y, y,    [=](float x, float y){     return y + a*x; }); } #pragma omp target data map(x,y) { std::transform(par, x, x+n, y, y,    [=](float x, float y){     return y + a*x; }); }</pre>	<pre>global void saxpy(int n, float a, float *x, float *y) { int i = blockIdx.x*blockDim.x +</pre>						
ACCELERATION LIBRARIES								
Core Math	Communication Data Analytics	Al Quantum						

## **NVIDIA HPC SDK**

Download at developer.nvidia.com/hpc-sdk

## **NVIDIA HPC SDK**



Develop for the NVIDIA HPC Platform: GPU, CPU and Interconnect

HPC Libraries | GPU Accelerated C++ and Fortran | Directives | CUDA

**OpenACC** 

KATHONS

# Programming Model Interoperability

All programming models are interoperable:

- > Can be used in the same source file.
- > ABI compatible with each other and NVCC.

nvc++ -stdpar -cuda -acc -mp



**OpenACC** 



CUDA

```
__global__ void
times_two(float* first, uint64_t n) {
  auto t = blockIdx.x*blockDim.x+threadIdx.x;
  if (t < n) first[t] *= 2;
}
```

void compute(std::vector<float>& x) {
 auto const b = ((x.size() - 1)/ 64) + 1;
 times\_two<<<b, 64>>>(x.data(), x.size());
 cudaDeviceSynchronize();

```
#pragma omp target teams loop
for (auto& e : x) e *= e;
```

```
#pragma acc parallel loop
for (auto& e : x) e += e;
```