



# CUDA Introduction PATC GPU Programming Course 2021

26 April 2021 | Andreas Herten | Forschungszentrum Jülich

# Outline

Introduction

  GPU History

  Architecture Comparison

  Jülich Systems

  App Showcase

Platform

  Overview

  3 Core Features

    Memory

    Asynchronicity

    SIMT

High Throughput

Summary

Programming GPUs

  Libraries

  GPU programming models

  Directives

  Thrust

  CUDA C/C++

    Kernels

    Grid, Blocks

    Memory Management

    Unified Memory

# History of GPUs

A short but unparalleled story

- 1999 Graphics computation pipeline implemented in dedicated *graphics hardware*
  - Computations using OpenGL graphics library [2]
  - »GPU« coined by NVIDIA [3]

# History of GPUs

A short but unparalleled story

- 1999 Graphics computation pipeline implemented in dedicated *graphics hardware*
  - Computations using OpenGL graphics library [2]
  - »GPU« coined by NVIDIA [3]
- 2001 NVIDIA GeForce 3 with *programmable* shaders (instead of fixed pipeline) and floating-point support; 2003: DirectX 9 at ATI

# History of GPUs

A short but unparalleled story

- 1999 Graphics computation pipeline implemented in dedicated *graphics hardware*
  - Computations using OpenGL graphics library [2]
  - »GPU« coined by NVIDIA [3]
- 2001 NVIDIA GeForce 3 with *programmable* shaders (instead of fixed pipeline) and floating-point support; 2003: DirectX 9 at ATI
- 2007 CUDA

# History of GPUs

A short but unparalleled story

- 1999 Graphics computation pipeline implemented in dedicated *graphics hardware*
  - Computations using OpenGL graphics library [2]
  - »GPU« coined by NVIDIA [3]
- 2001 NVIDIA GeForce 3 with *programmable* shaders (instead of fixed pipeline) and floating-point support; 2003: DirectX 9 at ATI
- 2007 CUDA
- 2009 OpenCL

# History of GPUs

A short but unparalleled story

- 1999 Graphics computation pipeline implemented in dedicated *graphics hardware*
  - Computations using OpenGL graphics library [2]
  - »GPU« coined by NVIDIA [3]
- 2001 NVIDIA GeForce 3 with *programmable* shaders (instead of fixed pipeline) and floating-point support; 2003: DirectX 9 at ATI
- 2007 CUDA
- 2009 OpenCL
- 2020 Top 500: 25 % with NVIDIA GPUs (#2, #3) [4], Green 500: 6 of top 10 with GPUs [5]

# History of GPUs

A short but unparalleled story

- 1999 Graphics computation pipeline implemented in dedicated *graphics hardware*
  - Computations using OpenGL graphics library [2]
  - »GPU« coined by NVIDIA [3]
- 2001 NVIDIA GeForce 3 with *programmable* shaders (instead of fixed pipeline) and floating-point support; 2003: DirectX 9 at ATI
- 2007 CUDA
- 2009 OpenCL
- 2020 Top 500: 25 % with NVIDIA GPUs (#2, #3) [4], Green 500: 6 of top 10 with GPUs [5]
- 2021 : Leonardo (250 PFLOP/s\*, Italy), NVIDIA GPUs;  LUMI (552 PFLOP/s, Finland), AMD GPUs  
: Frontier (> 1.5 EFLOP/s, ORNL), AMD GPUs

\*: Effective FLOP/s, not theoretical peak

# History of GPUs

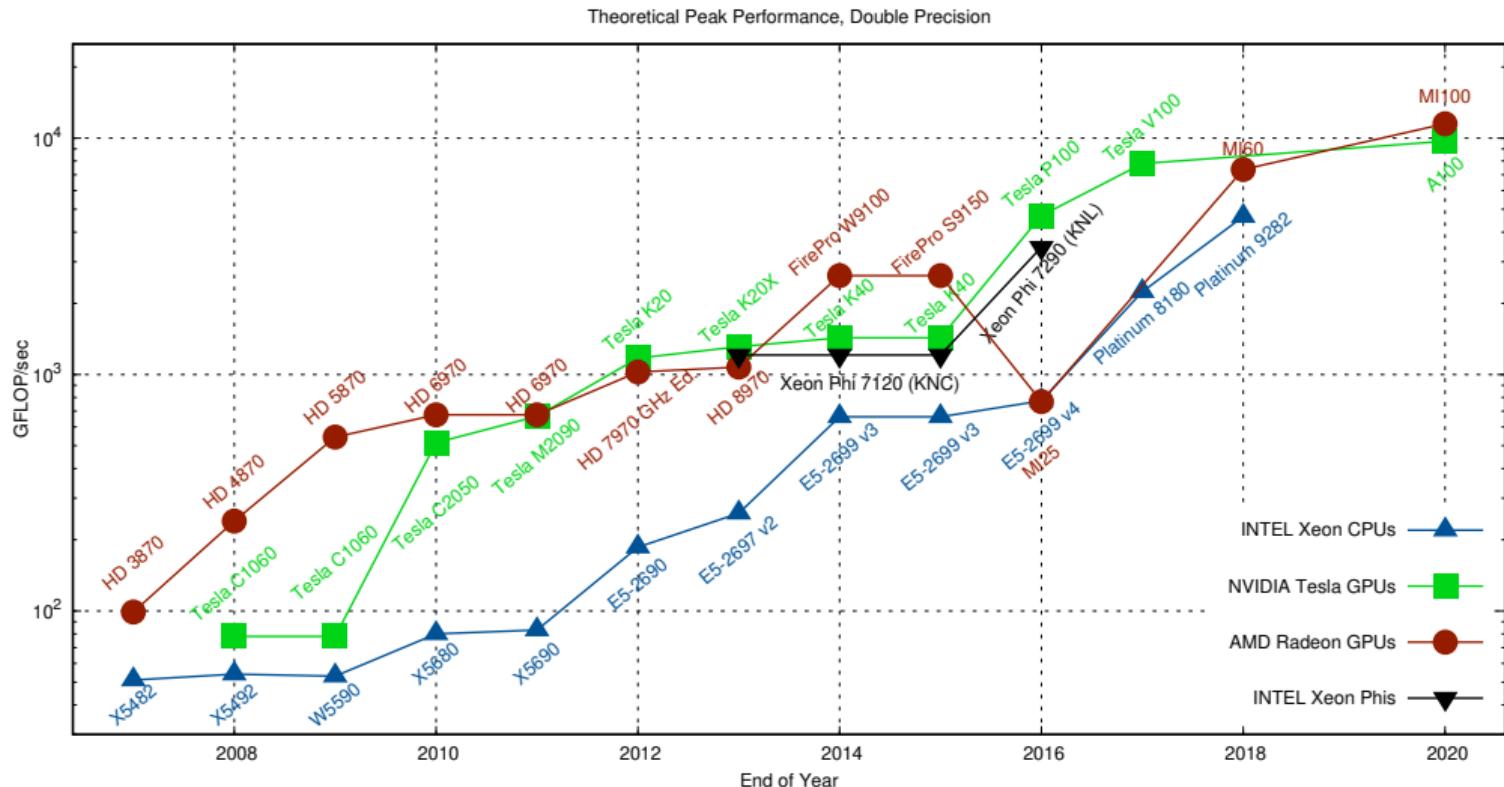
A short but unparalleled story

- 1999 Graphics computation pipeline implemented in dedicated *graphics hardware*
  - Computations using OpenGL graphics library [2]
  - »GPU« coined by NVIDIA [3]
- 2001 NVIDIA GeForce 3 with *programmable* shaders (instead of fixed pipeline) and floating-point support; 2003: DirectX 9 at ATI
- 2007 CUDA
- 2009 OpenCL
- 2020 Top 500: 25 % with NVIDIA GPUs (#2, #3) [4], Green 500: 6 of top 10 with GPUs [5]
- 2021 : Leonardo (250 PFLOP/s\*, Italy), NVIDIA GPUs; LUMI (552 PFLOP/s, Finland), AMD GPUs  
: Frontier (> 1.5 EFLOP/s, ORNL), AMD GPUs
- Future : ???
  - : Aurora ( $\approx$  1 EFLOP/s, Argonne), Intel GPUs; El Capitan ( $\approx$  2 EFLOP/s, LLNL), AMD GPUs

\*: Effective FLOP/s, not theoretical peak

# Status Quo Across Architectures

## Performance

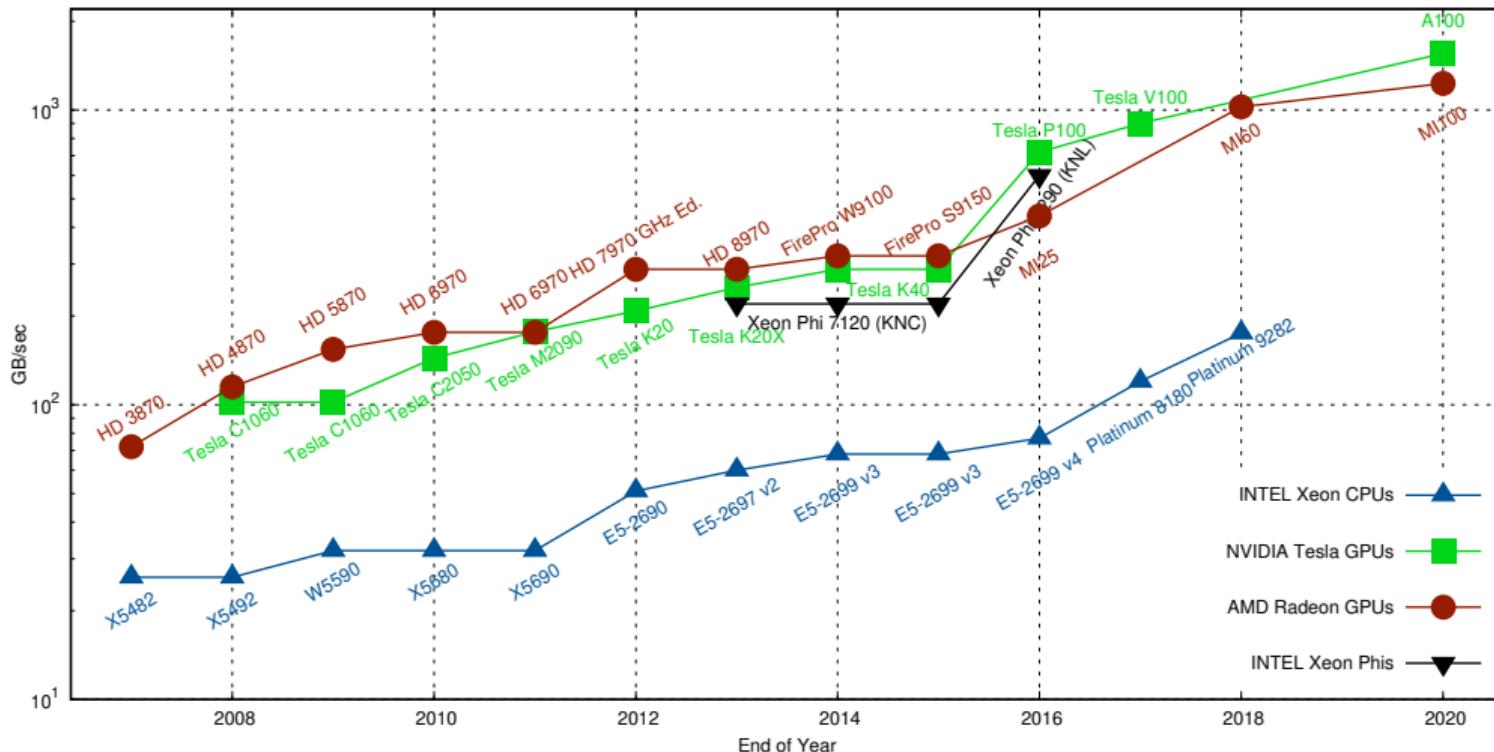


Graphic: Rupp [6]

# Status Quo Across Architectures

## Memory Bandwidth

Theoretical Peak Memory Bandwidth Comparison





## JUWELS Cluster – Jülich's Scalable System

- 2500 nodes with Intel Xeon CPUs ( $2 \times 24$  cores)
- 46 + 10 nodes with 4 NVIDIA Tesla V100 cards (16 GB memory)
- 10.4 (CPU) + 1.6 (GPU) PFLOP/s peak performance (Top500: #44)



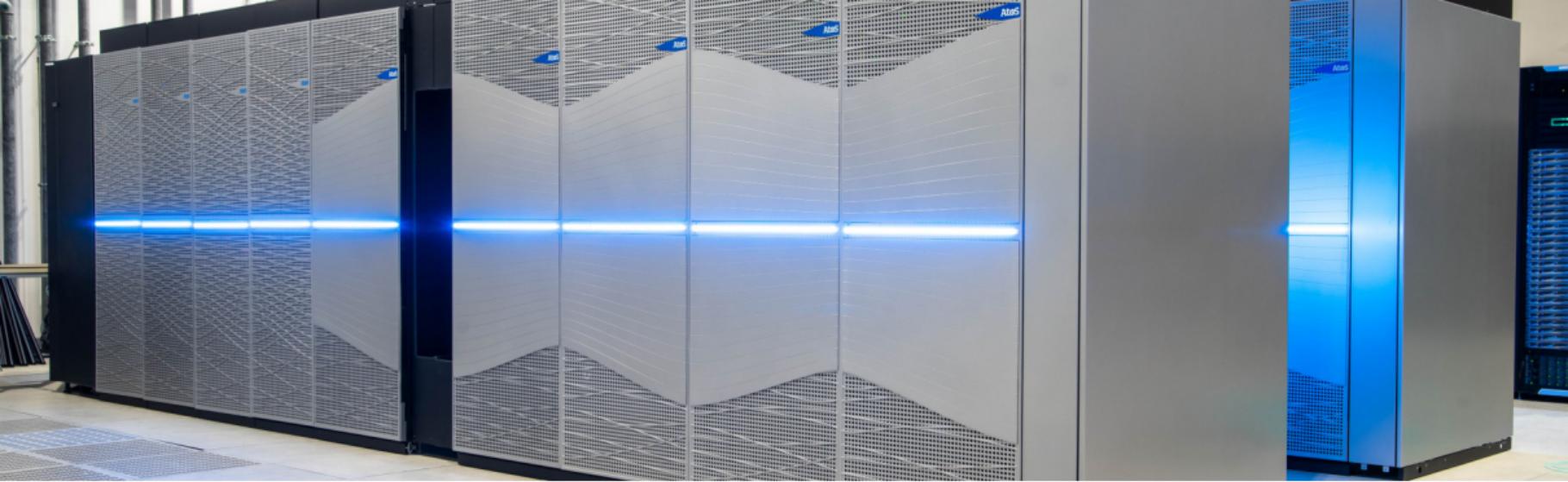
## JUWELS Booster – Scaling Higher!

- 936 nodes with AMD EPYC Rome CPUs ( $2 \times 24$  cores)
- Each with 4 NVIDIA A100 Ampere GPUs (each:  $\text{FP64TC: } 19.5$  TFLOP/s,  $40$  GB memory)  
 $\text{FP64: } 9.7$
- InfiniBand DragonFly+ HDR-200 network;  $4 \times 200$  Gbit/s per node



## JUWELS Booster – Scaling Higher!

- 936 nodes with AMD EPYC Rome CPUs ( $2 \times 24$  cores)
- Each with 4 NVIDIA A100 Ampere GPUs (each:  $\text{FP64TC: } 19.5$  TFLOP/s, 40 GB memory)  
 $\text{FP64: } 9.7$
- InfiniBand DragonFly+ HDR-200 network;  $4 \times 200$  Gbit/s per node



## JURECA DC – Multi-Purpose

- 768 nodes with AMD EPYC Rome CPUs ( $2 \times 64$  cores)
- 192 nodes with 4 NVIDIA A100 Ampere GPUs
- InfiniBand DragonFly+ HDR-100 network
- Also: JURECA Booster: 1640 nodes with Intel Xeon Phi *Knights Landing*

# Getting GPU-Acquainted

## Some Applications

TASK

Location of Code:

`~/CUDA-Course/1-Basics/exercises/tasks/01-Getting-Started`

See `Instructions.ipynb` for hints.

*Make sure to have sourced the course environment!*

# Getting GPU-Acquainted

TASK

## Some Applications

GEMM

N-Body

Location of Code:

`~/CUDA-Course/1-Basics/exercises/tasks/01-Getting-Started`

See `Instructions.ipynb` for hints.

*Make sure to have sourced the course environment!*

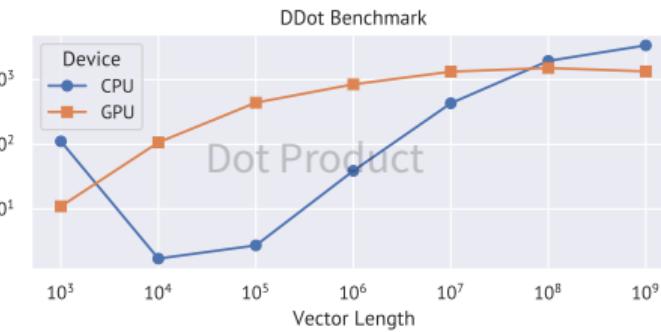
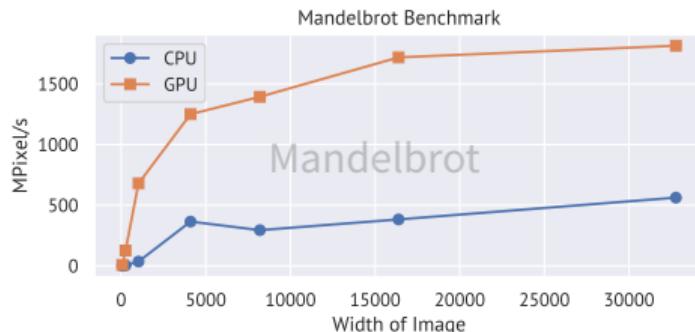
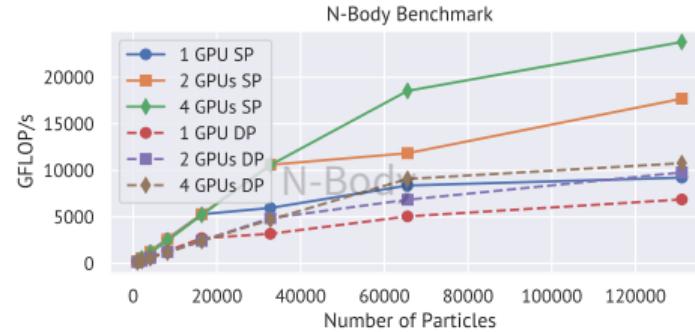
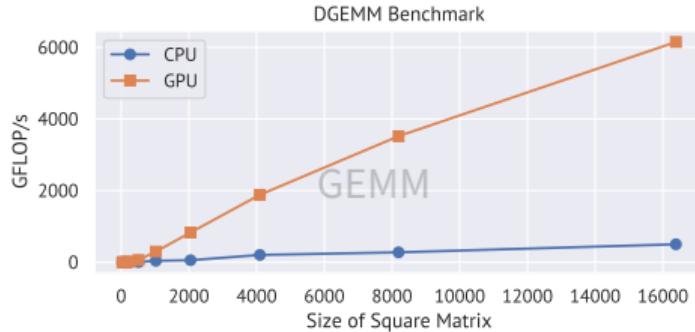
Mandelbrot

Dot Product

# Getting GPU-Acquainted

TASK

## Some Applications



# Platform

# CPU vs. GPU

A matter of specialties



# CPU vs. GPU

A matter of specialties



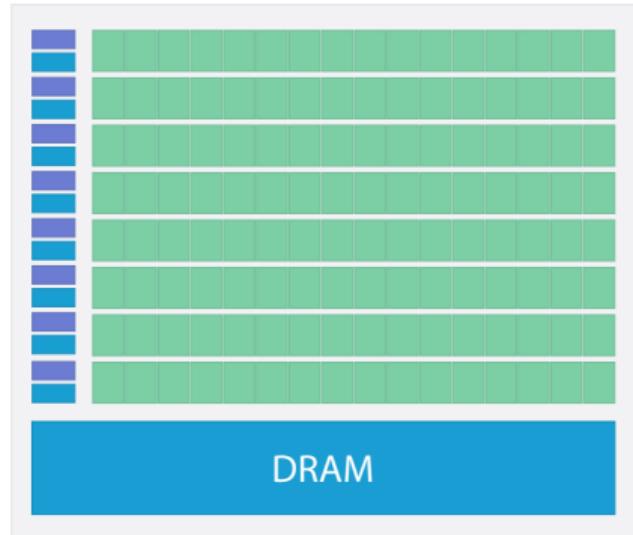
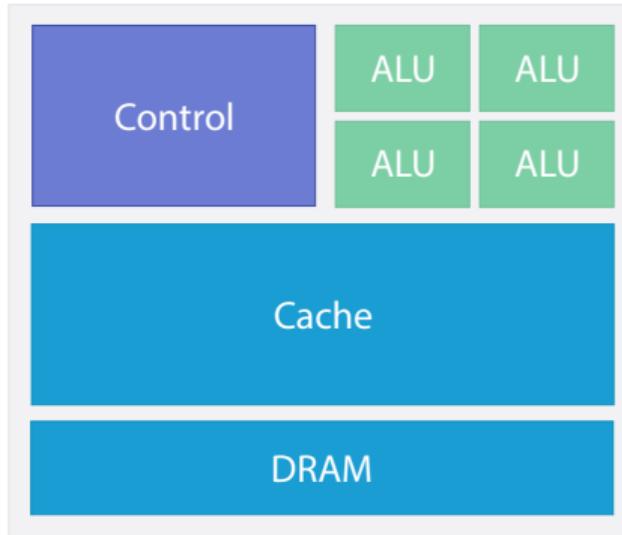
Transporting one



Transporting many

# CPU vs. GPU

## Chip



# GPU Architecture

## Overview

Aim: Hide Latency  
*Everything else follows*

# GPU Architecture

## Overview

Aim: Hide Latency  
*Everything else follows*

SIMT

Asynchronicity

Memory

# GPU Architecture

## Overview

Aim: Hide Latency  
*Everything else follows*

SIMT

Asynchronicity

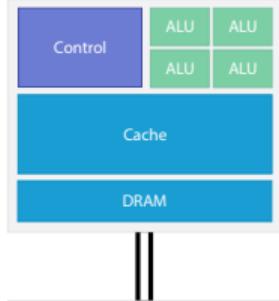
Memory

# Memory

GPU memory ain't no CPU memory

- GPU: accelerator / extension card
- Separate device from CPU

Host



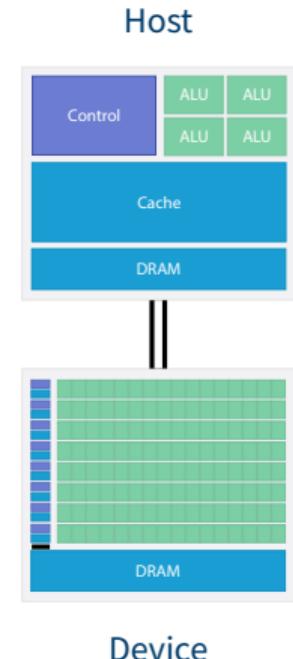
Device

# Memory

GPU memory ain't no CPU memory

*Unified Virtual Addressing*

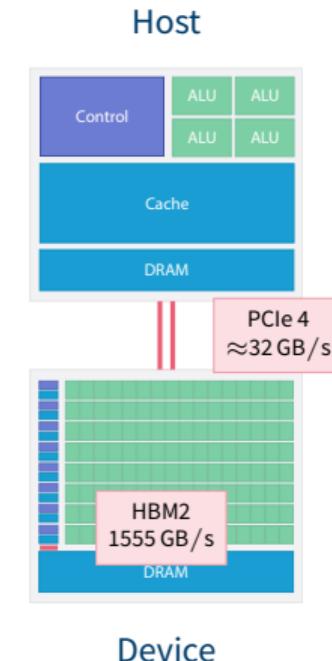
- GPU: accelerator / extension card
  - Separate device from CPU
- Separate memory, but UVA**



# Memory

GPU memory ain't no CPU memory

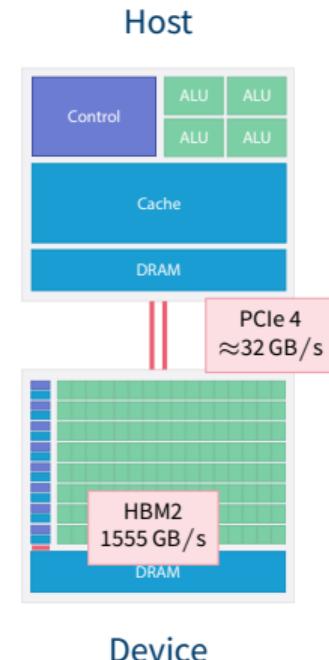
- GPU: accelerator / extension card
- Separate device from CPU
- Separate memory, but UVA**



# Memory

GPU memory ain't no CPU memory

- GPU: accelerator / extension card
- Separate device from CPU
- Separate memory, but UVA**
- Memory transfers need special consideration!  
*Do as little as possible!*

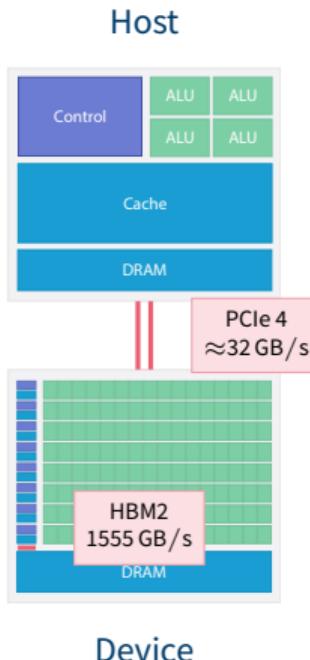


# Memory

GPU memory ain't no CPU memory

Unified Memory

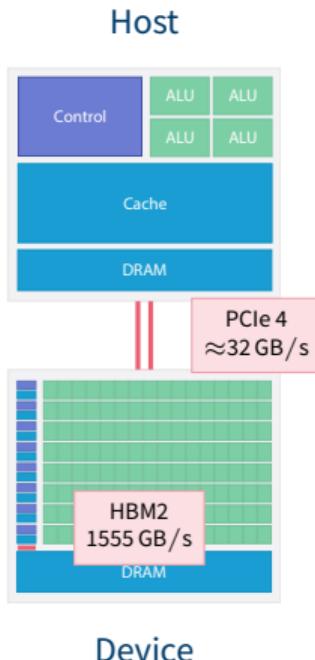
- GPU: accelerator / extension card
- Separate device from CPU
- Separate memory, but UVA and UM**
- Memory transfers need special consideration!  
*Do as little as possible!*
- Choice: automatic transfers (convenience) or manual transfers (control)



# Memory

GPU memory ain't no CPU memory

- GPU: accelerator / extension card
- Separate device from CPU
- Separate memory, but UVA and UM**
- Memory transfers need special consideration!  
*Do as little as possible!*
- Choice: automatic transfers (convenience) or manual transfers (control)



# Memory

GPU memory ain't no CPU memory

- GPU: accelerator / extension card
- Separate device from CPU
- Separate memory, but UVA and UM**
- Memory transfers need special consideration!  
*Do as little as possible!*
- Choice: automatic transfers (convenience) or manual transfers (control)

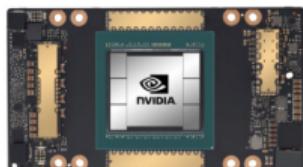
V100

32 GB RAM, 900 GB/s

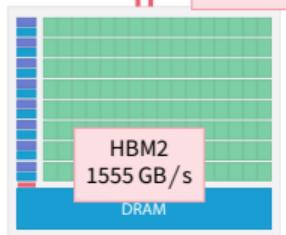
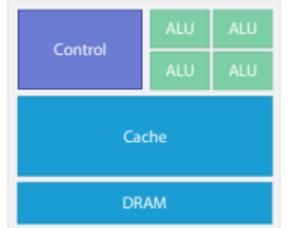


A100

40 GB RAM, 1555 GB/s



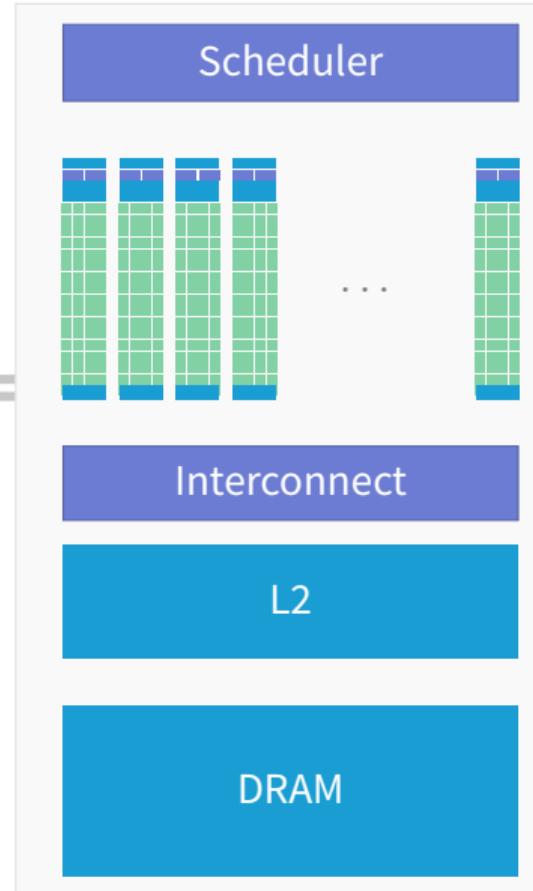
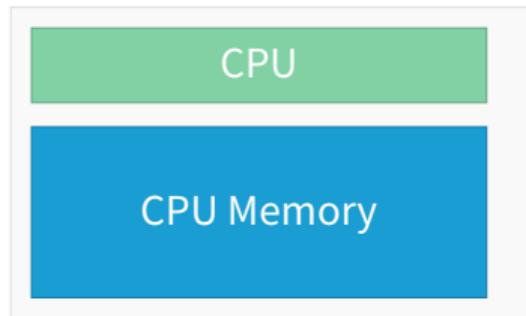
Host



Device

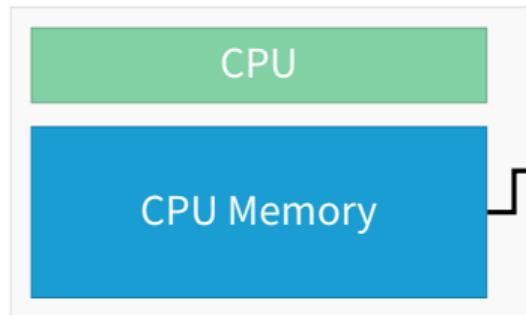
# Processing Flow

CPU → GPU → CPU

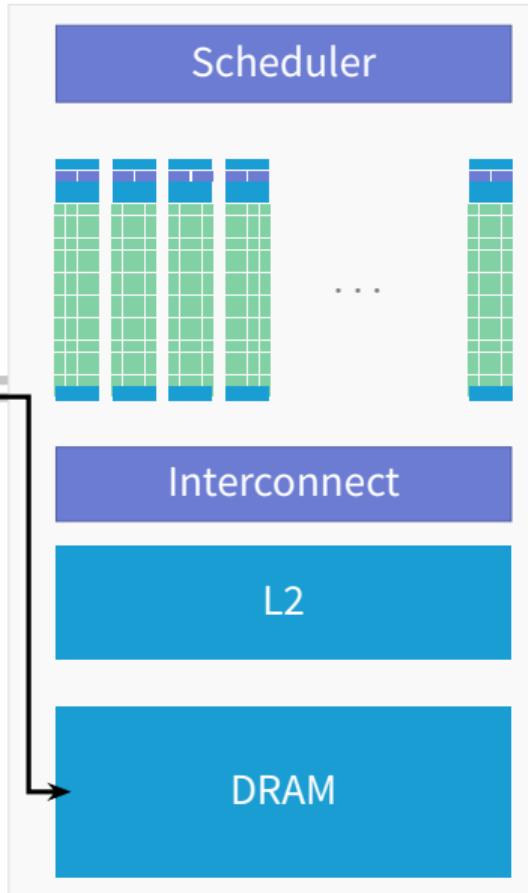


# Processing Flow

CPU → GPU → CPU

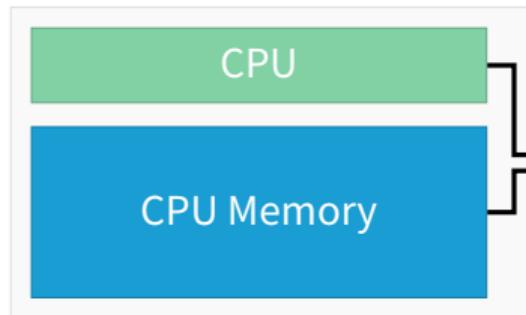


- 1 Transfer data from CPU memory to GPU memory



# Processing Flow

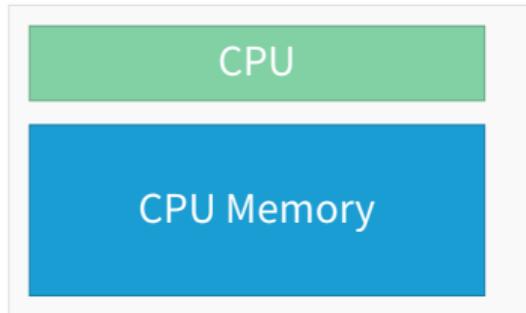
CPU → GPU → CPU



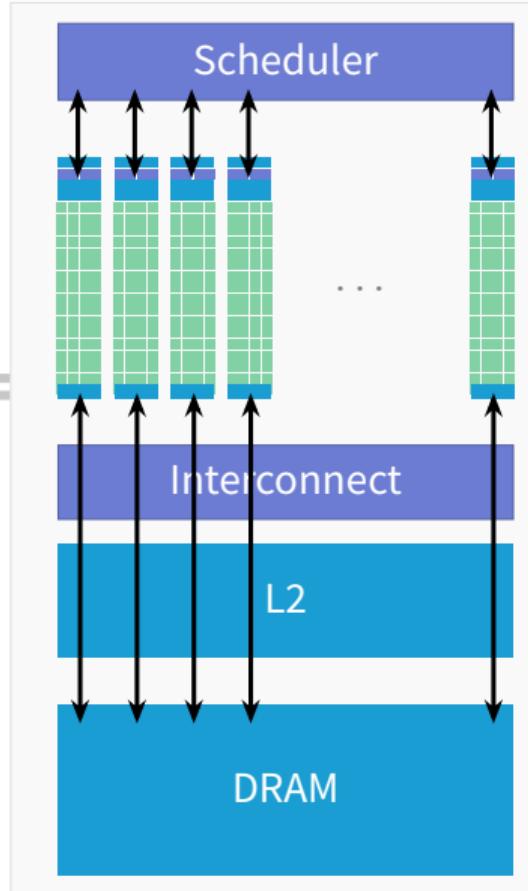
- 1 Transfer data from CPU memory to GPU memory, transfer program

# Processing Flow

CPU → GPU → CPU

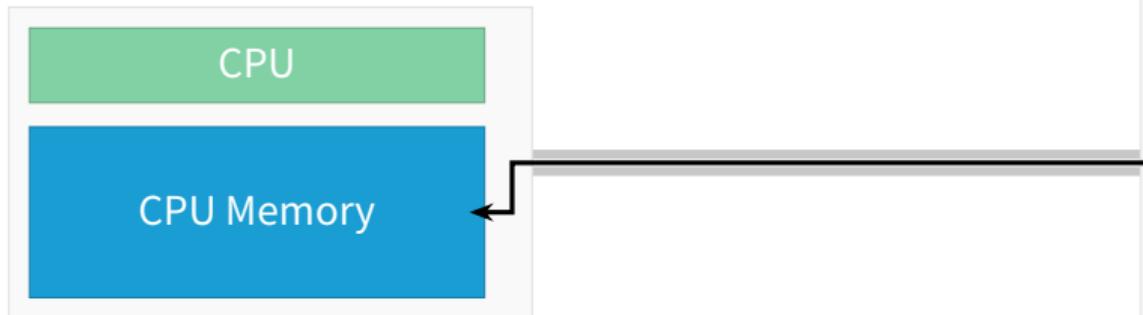


- 1 Transfer data from CPU memory to GPU memory, transfer program
- 2 Load GPU program, execute on SMs, get (cached) data from memory; write back

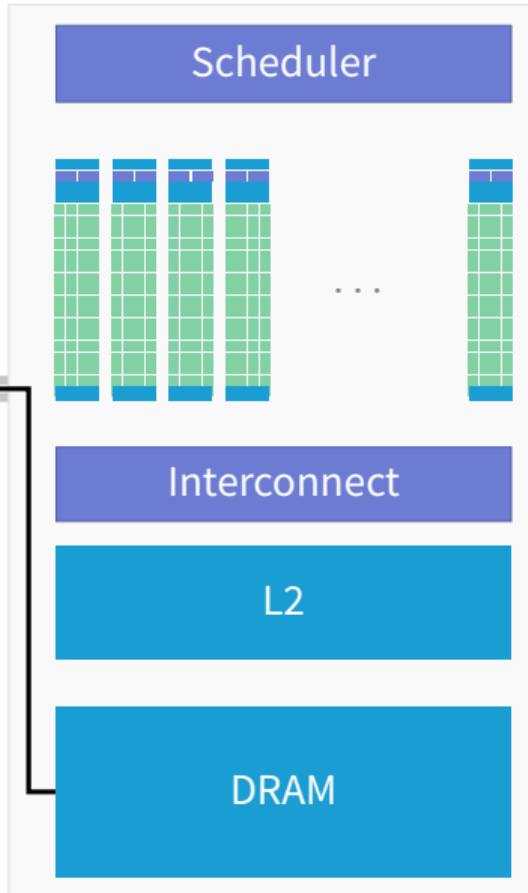


# Processing Flow

CPU → GPU → CPU



- 1 Transfer data from CPU memory to GPU memory, transfer program
- 2 Load GPU program, execute on SMs, get (cached) data from memory; write back
- 3 Transfer results back to host memory



# GPU Architecture

## Overview

Aim: Hide Latency  
*Everything else follows*

SIMT

Asynchronicity

Memory

# GPU Architecture

## Overview

Aim: Hide Latency  
*Everything else follows*

SIMT

Asynchronicity

Memory

# Async

## Following different streams

- Problem: Memory transfer is comparably slow  
Solution: Do something else in meantime (**computation**)!
- Overlap tasks
- Copy and compute engines run separately (*streams*)



- GPU needs to be fed: Schedule many computations
- CPU can do other work while GPU computes; synchronization

# GPU Architecture

## Overview

Aim: Hide Latency  
*Everything else follows*

SIMT

Asynchronicity

Memory

# GPU Architecture

## Overview

Aim: Hide Latency  
*Everything else follows*

SIMT

Asynchronicity

Memory

# SIMT

**SIMT = SIMD  $\oplus$  SMT**

- CPU:
  - Single Instruction, Multiple Data (SIMD)

*Scalar*

$$\begin{array}{l} A_0 + B_0 = C_0 \\ A_1 + B_1 = C_1 \\ A_2 + B_2 = C_2 \\ A_3 + B_3 = C_3 \end{array}$$

# SIMT

**SIMT = SIMD  $\oplus$  SMT**

- CPU:
  - Single Instruction, Multiple Data (SIMD)

*Vector*

$$\begin{array}{c} A_0 \\ A_1 \\ A_2 \\ A_3 \end{array} + \begin{array}{c} B_0 \\ B_1 \\ B_2 \\ B_3 \end{array} = \begin{array}{c} C_0 \\ C_1 \\ C_2 \\ C_3 \end{array}$$

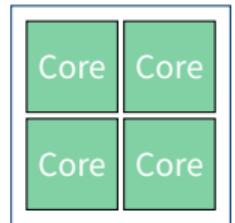
# SIMT

**SIMT = SIMD  $\oplus$  SMT**

- CPU:
  - Single Instruction, Multiple Data (SIMD)
  - Simultaneous Multithreading (SMT)

*Vector*

$$\begin{array}{c} A_0 \\ A_1 \\ A_2 \\ A_3 \end{array} + \begin{array}{c} B_0 \\ B_1 \\ B_2 \\ B_3 \end{array} = \begin{array}{c} C_0 \\ C_1 \\ C_2 \\ C_3 \end{array}$$



# SIMT

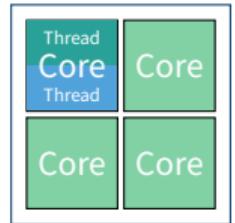
**SIMT = SIMD  $\oplus$  SMT**

- CPU:
  - Single Instruction, Multiple Data (SIMD)
  - Simultaneous Multithreading (SMT)

*Vector*

$$\begin{array}{c|c|c} A_0 & B_0 & C_0 \\ \hline A_1 & B_1 & C_1 \\ \hline A_2 & B_2 & C_2 \\ \hline A_3 & B_3 & C_3 \end{array} + \begin{array}{c|c|c} & & \\ & & \\ & & \\ & & \end{array} = \begin{array}{c|c|c} & & \\ & & \\ & & \\ & & \end{array}$$

*SMT*



# SIMT

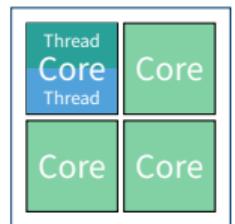
**SIMT = SIMD  $\oplus$  SMT**

- CPU:
  - Single Instruction, Multiple Data (SIMD)
  - Simultaneous Multithreading (SMT)
- GPU: Single Instruction, Multiple Threads (SIMT)

*Vector*

$$\begin{array}{c} A_0 \\ A_1 \\ A_2 \\ A_3 \end{array} + \begin{array}{c} B_0 \\ B_1 \\ B_2 \\ B_3 \end{array} = \begin{array}{c} C_0 \\ C_1 \\ C_2 \\ C_3 \end{array}$$

*SMT*



# SIMT

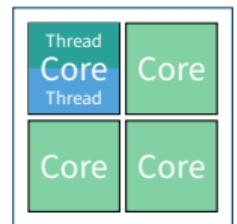
**SIMT = SIMD  $\oplus$  SMT**

- CPU:
  - Single Instruction, Multiple Data (SIMD)
  - Simultaneous Multithreading (SMT)
- GPU: Single Instruction, Multiple Threads (SIMT)

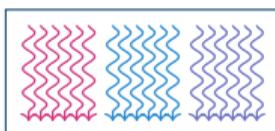
*Vector*

$$\begin{array}{ccc} A_0 & B_0 & C_0 \\ A_1 & B_1 & C_1 \\ A_2 & B_2 & C_2 \\ A_3 & B_3 & C_3 \end{array} + \begin{array}{c} \\ \\ \\ \end{array} = \begin{array}{ccc} & & \\ & & \\ & & \\ & & \end{array}$$

*SMT*



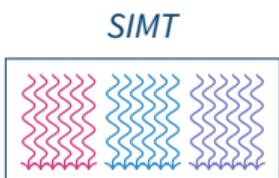
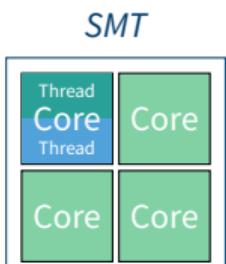
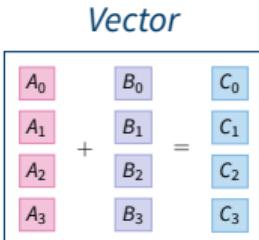
*SIMT*



SIMT

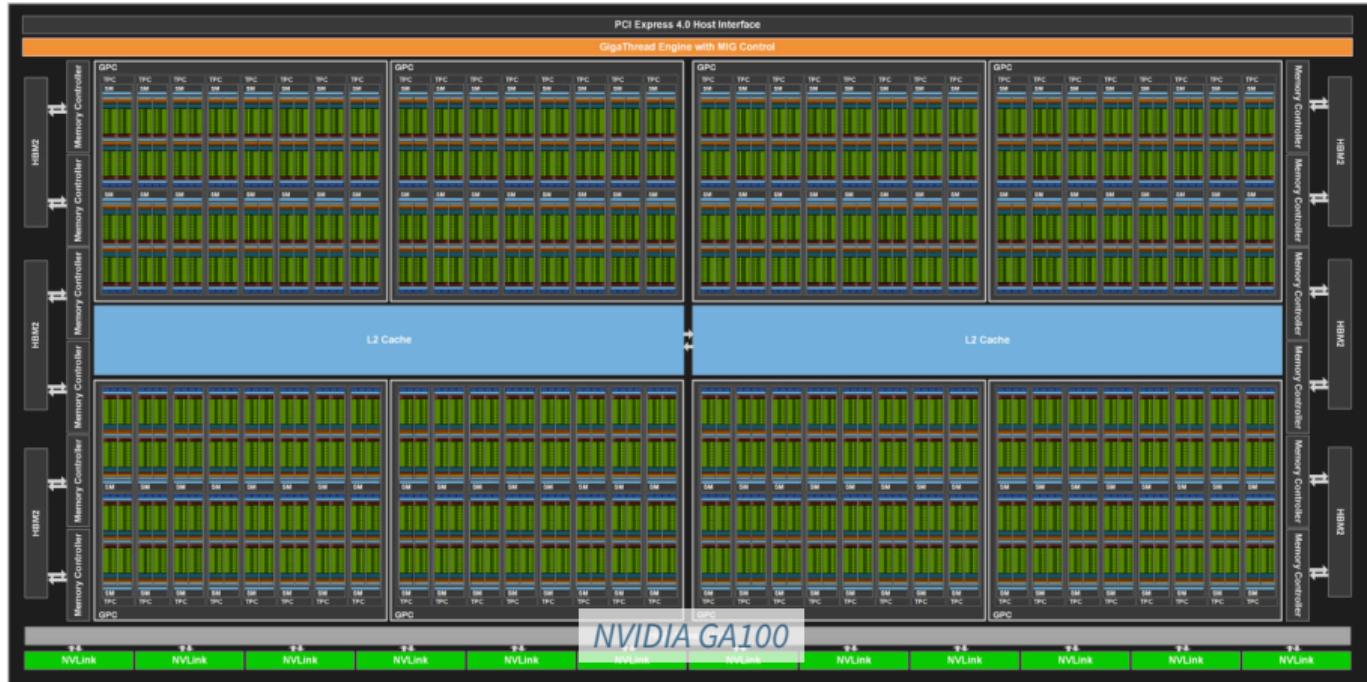
$$\text{SIMT} = \text{SIMD} \oplus \text{SMT}$$

- CPU:
    - Single Instruction, Multiple Data (**SIMD**)
    - Simultaneous Multithreading (**SMT**)
  - GPU: Single Instruction, Multiple Threads (**SIMT**)
    - CPU core  $\approx$  GPU multiprocessor (**SM**)
    - Working unit: set of threads (32, a *warp*)
    - Fast switching of threads (large register file)
    - Branching 



# SIMT

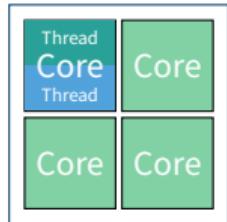
**SIMT = SIMD  $\oplus$  SMT**



Vector

$$\begin{matrix} A_0 & & B_0 & & C_0 \\ A_1 & + & B_1 & = & C_1 \\ A_2 & & B_2 & & C_2 \\ A_3 & & B_3 & & C_3 \end{matrix}$$

SMT



Graphics: img:ampere/pictures

SIMT



# SIMT

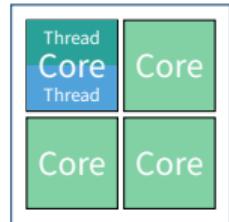
**SIMT = SIMD  $\oplus$  SMT**



Vector

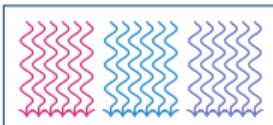
$$\begin{matrix} A_0 & & B_0 & & C_0 \\ A_1 & + & B_1 & = & C_1 \\ A_2 & & B_2 & & C_2 \\ A_3 & & B_3 & & C_3 \end{matrix}$$

SMT



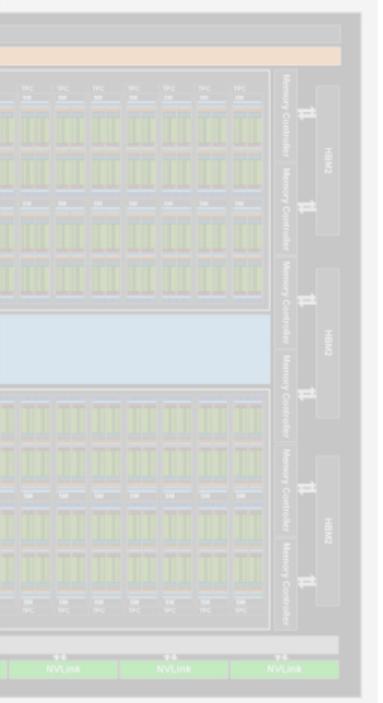
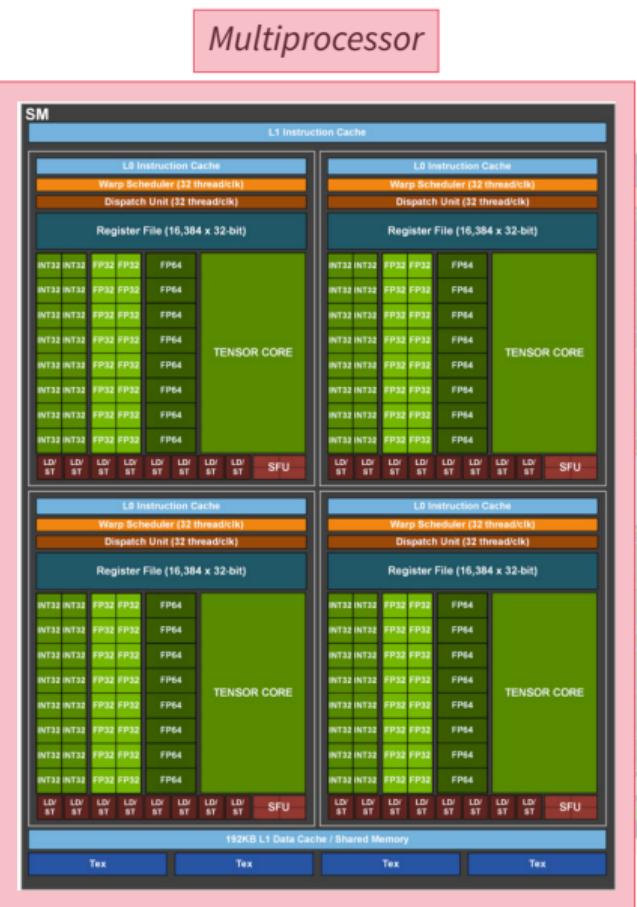
Graphics: img:ampere/pictures

SIMT



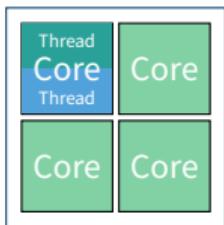
# SIMT

**SIMT = SIMD  $\oplus$  SMT**



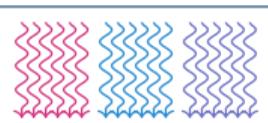
$$\begin{matrix} A_0 & & B_0 & & C_0 \\ A_1 & + & B_1 & = & C_1 \\ A_2 & & B_2 & & C_2 \\ A_3 & & B_3 & & C_3 \end{matrix}$$

**SMT**



Graphics: img:ampere/pictures

**SIMT**



# Low Latency vs. High Throughput

Maybe GPU's ultimate feature

CPU Minimizes latency within each thread

GPU Hides latency with computations from other thread warps

# Low Latency vs. High Throughput

Maybe GPU's ultimate feature

CPU Minimizes latency within each thread

GPU Hides latency with computations from other thread warps

CPU Core: Low Latency



- Thread/Warp
- Processing
- Context Switch
- Ready
- Waiting

# Low Latency vs. High Throughput

Maybe GPU's ultimate feature

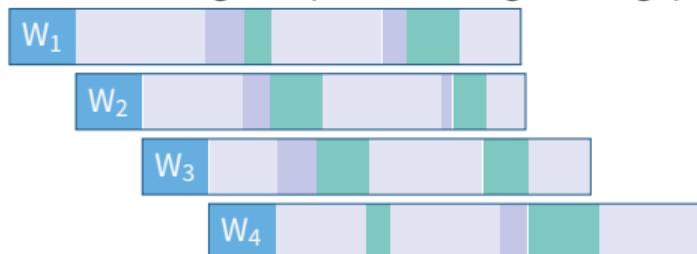
CPU Minimizes latency within each thread

GPU Hides latency with computations from other thread warps

CPU Core: Low Latency



GPU Streaming Multiprocessor: High Throughput



- █ Thread/Warp
- █ Processing
- █ Context Switch
- █ Ready
- █ Waiting

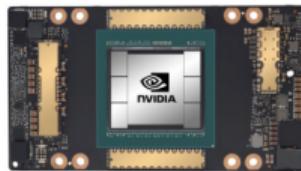
# CPU vs. GPU

Let's summarize this!



Optimized for **low latency**

- + Large main memory
- + Fast clock rate
- + Large caches
- + Branch prediction
- + Powerful ALU
- Relatively low memory bandwidth
- Cache misses costly
- Low performance per watt



Optimized for **high throughput**

- + High bandwidth main memory
- + Latency tolerant (parallelism)
- + More compute resources
- + High performance per watt
- Limited memory capacity
- Low per-thread performance
- Extension card

# Programming GPUs

# Preface: CPU

A simple CPU program!

SAXPY:  $\vec{y} = a\vec{x} + \vec{y}$ , with single precision

Part of LAPACK BLAS Level 1

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

int a = 42;
int n = 10;
float x[n], y[n];
// fill x, y

saxpy(n, a, x, y);
```

# Summary of Acceleration Possibilities



*Drop-in*  
Acceleration

*Easy*  
Acceleration

*Flexible*  
Acceleration

# Summary of Acceleration Possibilities



*Drop-in*  
Acceleration

*Easy*  
Acceleration

*Flexible*  
Acceleration

# Libraries

Programming GPUs is easy: Just don't!

# Libraries

Programming GPUs is easy: Just don't!

*Use applications & libraries*

# Libraries

Programming GPUs is easy: Just don't!

*Use applications & libraries*



# Libraries

Programming GPUs is easy: Just don't!

*Use applications & libraries*



cuBLAS



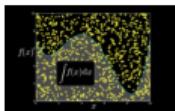
cuSPARSE



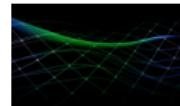
cuDNN



cuFFT



cuRAND



CUDA Math



{ } ARRAYFIRE

Numba



theano

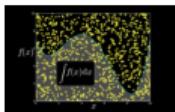
# Libraries

Programming GPUs is easy: Just don't!

*Use applications & libraries*



cuFFT



cuRAND



{} ARRAYFIRE



CUDA Math



theano



JÜLICH  
SUPERCOMPUTING  
CENTRE

# cuBLAS

## Parallel algebra



- GPU-parallel BLAS (all 152 routines)
  - Single, double, complex data types
  - Constant competition with Intel's MKL
  - Multi-GPU support
- <https://developer.nvidia.com/cublas>  
<http://docs.nvidia.com/cuda/cublas>

# cuBLAS

## Code example

```
int a = 42;  int n = 10;
float x[n], y[n];
// fill x, y

cublasHandle_t handle;
cublasCreate(&handle);

float * d_x, * d_y;
cudaMallocManaged(&d_x, n * sizeof(x[0]));
cudaMallocManaged(&d_y, n * sizeof(y[0]));
cublasSetVector(n, sizeof(x[0]), x, 1, d_x, 1);
cublasSetVector(n, sizeof(y[0]), y, 1, d_y, 1);

cublasSaxpy(n, a, d_x, 1, d_y, 1);

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

cudaFree(d_x); cudaFree(d_y);
cublasDestroy(handle);
```

# cuBLAS

## Code example

```
int a = 42;  int n = 10;  
float x[n], y[n];  
// fill x, y
```

```
cublasHandle_t handle;  
cublasCreate(&handle);
```

```
float * d_x, * d_y;  
cudaMallocManaged(&d_x, n * sizeof(x[0]));  
cudaMallocManaged(&d_y, n * sizeof(y[0]));  
cublasSetVector(n, sizeof(x[0]), x, 1, d_x, 1);  
cublasSetVector(n, sizeof(y[0]), y, 1, d_y, 1);  
  
cublasSaxpy(n, a, d_x, 1, d_y, 1);  
  
cublasGetVector(n, sizeof(y[0]), d_y, 1, y, 1);  
  
cudaFree(d_x); cudaFree(d_y);  
cublasDestroy(handle);
```

Initialize

# cuBLAS

## Code example

```
int a = 42;  int n = 10;  
float x[n], y[n];  
// fill x, y
```

```
cublasHandle_t handle;  
cublasCreate(&handle);
```

Initialize

```
float * d_x, * d_y;  
cudaMallocManaged(&d_x, n * sizeof(x[0]));  
cudaMallocManaged(&d_y, n * sizeof(y[0]));  
cublasSetVector(n, sizeof(x[0]), x, 1, d_x, 1);  
cublasSetVector(n, sizeof(y[0]), y, 1, d_y, 1);
```

```
cublasSaxpy(n, a, d_x, 1, d_y, 1);
```

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

```
cudaFree(d_x); cudaFree(d_y);  
cublasDestroy(handle);
```

Allocate GPU memory

# cuBLAS

## Code example

```
int a = 42;  int n = 10;  
float x[n], y[n];  
// fill x, y
```

```
cublasHandle_t handle;  
cublasCreate(&handle);
```

Initialize

```
float * d_x, * d_y;
```

```
cudaMallocManaged(&d_x, n * sizeof(x[0]));  
cudaMallocManaged(&d_y, n * sizeof(y[0]));
```

```
cublasSetVector(n, sizeof(x[0]), x, 1, d_x, 1);  
cublasSetVector(n, sizeof(y[0]), y, 1, d_y, 1);
```

Allocate GPU memory

```
cublasSaxpy(n, a, d_x, 1, d_y, 1);
```

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

```
cudaFree(d_x); cudaFree(d_y);  
cublasDestroy(handle);
```

Copy data to GPU

# cuBLAS

## Code example

```
int a = 42;  int n = 10;  
float x[n], y[n];  
// fill x, y  
  
cublasHandle_t handle;  
cublasCreate(&handle);
```

Initialize

```
float * d_x, * d_y;  
cudaMallocManaged(&d_x, n * sizeof(x[0]));  
cudaMallocManaged(&d_y, n * sizeof(y[0]));  
cublasSetVector(n, sizeof(x[0]), x, 1, d_x, 1);  
cublasSetVector(n, sizeof(y[0]), y, 1, d_y, 1);
```

Allocate GPU memory

```
cublasSaxpy(n, a, d_x, 1, d_y, 1);
```

Copy data to GPU

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

Call BLAS routine

```
cudaFree(d_x); cudaFree(d_y);  
cublasDestroy(handle);
```

# cuBLAS

## Code example

```
int a = 42;  int n = 10;  
float x[n], y[n];  
// fill x, y
```

```
cublasHandle_t handle;  
cublasCreate(&handle);
```

```
float * d_x, * d_y;  
cudaMallocManaged(&d_x, n * sizeof(x[0]));  
cudaMallocManaged(&d_y, n * sizeof(y[0]));  
cublasSetVector(n, sizeof(x[0]), x, 1, d_x, 1);  
cublasSetVector(n, sizeof(y[0]), y, 1, d_y, 1);
```

```
cublasSaxpy(n, a, d_x, 1, d_y, 1);
```

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

```
cudaFree(d_x); cudaFree(d_y);  
cublasDestroy(handle);
```

Initialize

Allocate GPU memory

Copy data to GPU

Call BLAS routine

Copy result to host

# cuBLAS

## Code example

```
int a = 42;  int n = 10;  
float x[n], y[n];  
// fill x, y
```

```
cublasHandle_t handle;  
cublasCreate(&handle);
```

```
float * d_x, * d_y;  
cudaMallocManaged(&d_x, n * sizeof(x[0]));  
cudaMallocManaged(&d_y, n * sizeof(y[0]));  
cublasSetVector(n, sizeof(x[0]), x, 1, d_x, 1);  
cublasSetVector(n, sizeof(y[0]), y, 1, d_y, 1);
```

```
cublasSaxpy(n, a, d_x, 1, d_y, 1);
```

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

```
cudaFree(d_x); cudaFree(d_y);  
cublasDestroy(handle);
```

Initialize

Allocate GPU memory

Copy data to GPU

Call BLAS routine

Copy result to host

Finalize

# cuBLAS Task

## Implement a matrix-matrix multiplication

- Location of code: 1-Basics/exercises/tasks/02-cuBLAS
- Look at Instructions.ipynb Notebook for instructions
  - 1 Implement call to double-precision GEMM of cuBLAS
  - 2 Build with make (load modules of this task via `source setup.sh!`)
  - 3 Run with `make run`
- Check [cuBLAS documentation](#) for details on `cublasDgemm()`

# Summary of Acceleration Possibilities

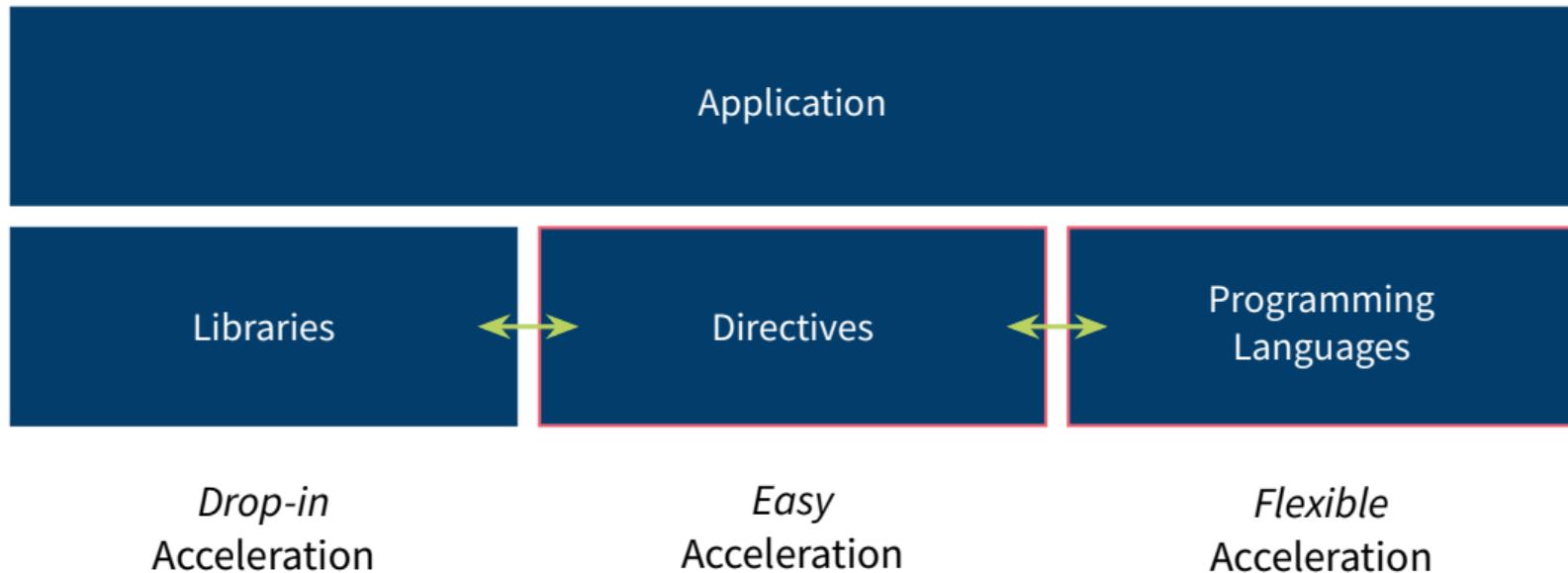


*Drop-in*  
Acceleration

*Easy*  
Acceleration

*Flexible*  
Acceleration

# Summary of Acceleration Possibilities



# ⚠ Parallelism

Libraries are not enough?

You think you want to write your own GPU code?

# Primer on Parallel Scaling

## Amdahl's Law

Possible maximum speedup for  
 $N$  parallel processors

**Total Time**  $t = t_{\text{serial}} + t_{\text{parallel}}$

# Primer on Parallel Scaling

## Amdahl's Law

Possible maximum speedup for  
 $N$  parallel processors

**Total Time**  $t = t_{\text{serial}} + t_{\text{parallel}}$

**$N$  Processors**  $t(N) = t_s + t_p/N$

# Primer on Parallel Scaling

## Amdahl's Law

Possible maximum speedup for  
 $N$  parallel processors

Total Time  $t = t_{\text{serial}} + t_{\text{parallel}}$

$N$  Processors  $t(N) = t_s + t_p/N$

Speedup  $s(N) = t/t(N) = \frac{t_s+t_p}{t_s+t_p/N}$

# Primer on Parallel Scaling

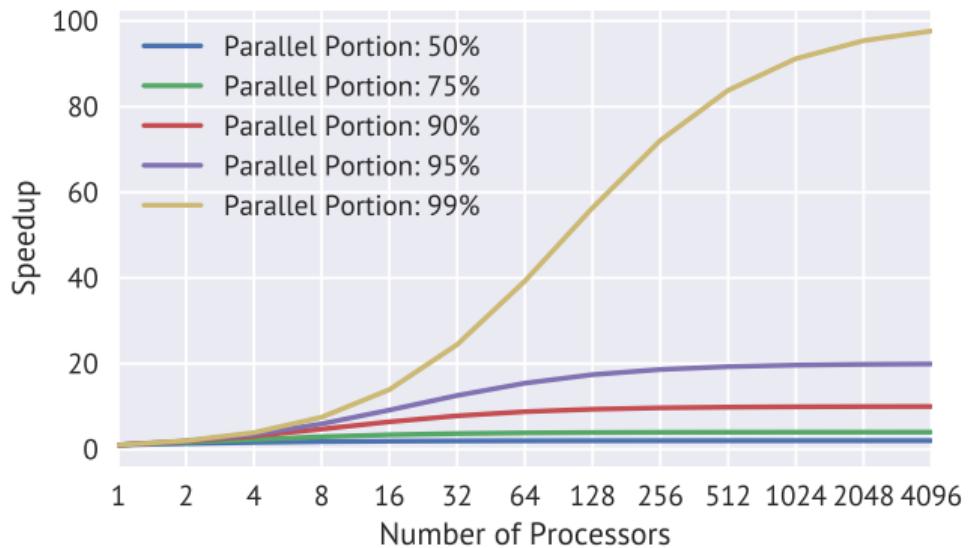
## Amdahl's Law

Possible maximum speedup for  $N$  parallel processors

$$\text{Total Time } t = t_{\text{serial}} + t_{\text{parallel}}$$

$$N \text{ Processors } t(N) = t_s + t_p/N$$

$$\text{Speedup } s(N) = t/t(N) = \frac{t_s+t_p}{t_s+t_p/N}$$



# ⚠ Parallelism

Parallel programming is not easy!

Things to consider:

- Is my application **computationally intensive** enough?
- What are the levels of **parallelism**?
- How much **data** needs to be **transferred**?
- Is the **gain** worth the **pain**?

# Alternatives

## The twilight

There are alternatives to CUDA C, which **can** ease the *pain*...

- OpenACC, OpenMP
- Thrust
- Kokkos, RAJA, ALPAKA, SYCL, DPC++, pSTL
- PyCUDA, Cupy, Numba

Other alternatives

- CUDA Fortran
- HIP
- OpenCL

# Alternatives

## The twilight

There are alternatives to CUDA C, which **can** ease the *pain*...

- OpenACC, OpenMP
- Thrust
- Kokkos, RAJA, ALPAKA, SYCL, DPC++, pSTL
- PyCUDA, Cupy, Numba

Other alternatives

- CUDA Fortran
- HIP
- OpenCL

# Programming GPUs

## Directives

# GPU Programming with Directives

Keepin' you portable

- Annotate serial source code by directives

```
#pragma acc loop  
for (int i = 0; i < 1; i++) {};
```

# GPU Programming with Directives

Keepin' you portable

- Annotate serial source code by directives

```
#pragma acc loop  
for (int i = 0; i < 1; i++) {};
```

- OpenACC:** Especially for GPUs; **OpenMP:** Has GPU support
- Compiler interprets directives, creates according instructions

# GPU Programming with Directives

Keepin' you portable

- Annotate serial source code by directives

```
#pragma acc loop  
for (int i = 0; i < 1; i++) {};
```

- OpenACC:** Especially for GPUs; **OpenMP:** Has GPU support
- Compiler interprets directives, creates according instructions

## Pro

- Portability
  - Other compiler? No problem! To it, it's a serial program
  - Different target architectures from same code
- Easy to program

## Con

- Only few compilers
- Not all the raw power available
- A little harder to debug

# GPU Programming with Directives

The power of... two.

OpenMP Standard for multithread programming on CPU, GPU since 4.0, better since 4.5

```
#pragma omp target map(tofrom:y), map(to:x)
#pragma omp teams num_teams(10) num_threads(10)
#pragma omp distribute
for ( ) {
    #pragma omp parallel for
    for ( ) {
        // ...
    }
}
```

OpenACC Similar to OpenMP, but more specifically for GPUs

For C/C++ and Fortran

# OpenACC

## Code example

```
void saxpy_acc(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];  
}  
  
int a = 42;  
int n = 10;  
float x[n], y[n];  
// fill x, y  
  
saxpy_acc(n, a, x, y);
```

# OpenACC

## Code example

```
void saxpy_acc(int n, float a, float * x, float * y) {  
    #pragma acc parallel loop copy(y) copyin(x)  
    for (int i = 0; i < n; i++)  
        y[i] = a * x[i] + y[i];  
}  
  
int a = 42;  
int n = 10;  
float x[n], y[n];  
// fill x, y  
  
saxpy_acc(n, a, x, y);
```

# Programming GPUs

## Thrust

# Thrust

Iterators! Iterators everywhere!



- $\frac{\text{Thrust}}{\text{CUDA}} = \frac{\text{STL}}{\text{C++}}$
- Template library
- Based on iterators
- Data-parallel primitives (`scan()`, `sort()`, `reduce()`, ...)
- Fully compatible with plain CUDA C (comes with **CUDA Toolkit**)
- Great with `[](){} lambdas!`

→ <http://thrust.github.io/>  
<http://docs.nvidia.com/cuda/thrust/>

# Thrust

## Code example

```
int a = 42;
int n = 10;
thrust::host_vector<float> x(n), y(n);
// fill x, y

thrust::device_vector d_x = x, d_y = y;

using namespace thrust::placeholders;
thrust::transform(d_x.begin(), d_x.end(), d_y.begin(), d_y.begin(), a * _1 + _2);

x = d_x;
```

# Thrust Task

Let's sort some randomness

- Location of code: 1-Basics/exercises/tasks/03-Thrust
- Look at Instructions.ipynb for instructions
  - 1 Sort random numbers with Thrust on CPU and GPU
  - 2 Build with make  
Reset environment to original; call `source setup.sh` or re-login!
  - 3 Run with `make run`
- Check Thrust documentation for details on `thrust::sort()`

# Summary of Acceleration Possibilities



*Drop-in*  
Acceleration

*Easy*  
Acceleration

*Flexible*  
Acceleration

# Programming GPUs

## CUDA C/C++

# CUDA SAXPY

With runtime-managed data transfers

```
__global__ void saxpy_cuda(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];
}

int a = 42;
int n = 10;
float x[n], y[n];
// fill x, y
cudaMallocManaged(&x, n * sizeof(float));
cudaMallocManaged(&y, n * sizeof(float));

saxpy_cuda<<<2, 5>>>(n, a, x, y);

cudaDeviceSynchronize();
```

# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:

# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:

- Thread



# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:
  - Threads



# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:

-  Threads → Block



# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:

- Threads → Block
- Block

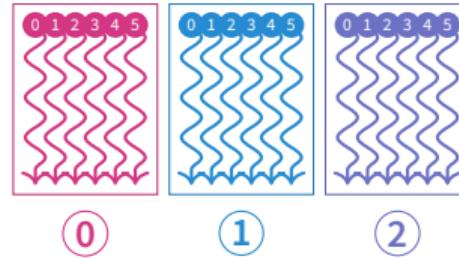


# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:

- Threads → Block
- Blocks

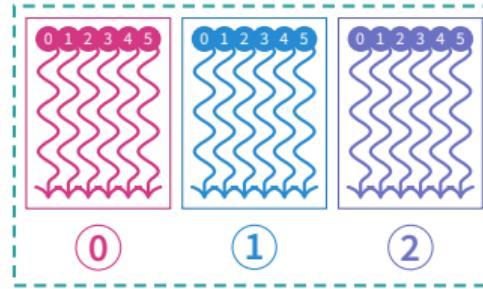


# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:

- Threads → Block
- Block → Grid

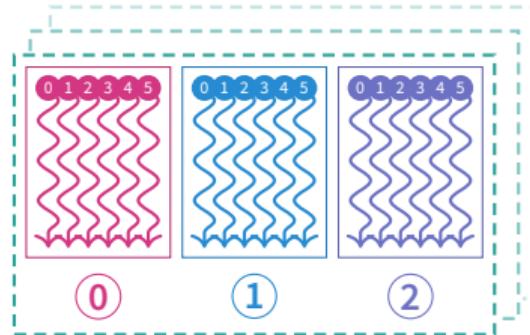


# CUDA's Parallel Model

In software: Threads, Blocks

- Methods to exploit parallelism:

- Threads → Block
- Blocks → Grid
- Threads & blocks in 3D 



# CUDA's Parallel Model

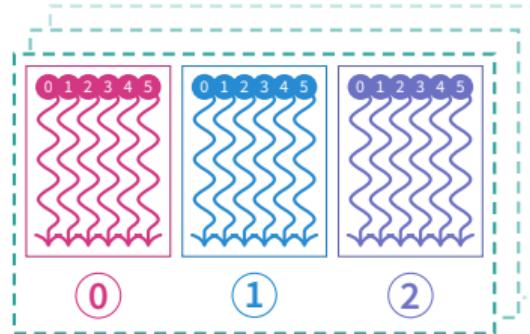
In software: Threads, Blocks

- Methods to exploit parallelism:

- Threads → Block

- Blocks → Grid

- Threads & blocks in 



- Parallel function: **kernel**

- **\_\_global\_\_ kernel(int a, float \* b) { }**

- Access own ID by global variables **threadIdx.x**, **blockIdx.y**, ...

- Execution entity: **threads**

- Lightweight → fast switching!

- 1000s threads execute simultaneously → order non-deterministic!

# Kernel Functions

- Kernel: Parallel GPU function
  - Executed by each thread
  - In parallel
  - Called from host or device

# Kernel Functions

- Kernel: Parallel GPU function
  - Executed by each thread
  - In parallel
  - Called from host or device
- All threads execute same code; but can take different paths in program flow (some penalty)

# Kernel Functions

- Kernel: Parallel GPU function
  - Executed by each thread
  - In parallel
  - Called from host or device
- All threads execute same code; but can take different paths in program flow (some penalty)
- Info about thread: local, global IDs

```
int currentThreadId = threadIdx.x;  
float x = input[currentThreadId];  
output[currentThreadId] = x*x;
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

## Identify Loops

```
void scale(float scale, float * in, float * out, int N) {  
    for (int i = 0; i < N; i++)  
        out[i] = scale * in[i];  
}
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

## Identify Loops

```
void scale(float scale, float * in, float * out, int N) {  
    for (  
        int i = 0;  
        i < N;  
        i++)  
        out[i] = scale * in[i];  
}
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

Identify Loops | Extract Index

```
void scale(float scale, float * in, float * out, int N) {  
    int i = 0  
    for ( ;  
        i < N;  
        i++)  
        out[i] = scale * in[i];  
}
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

Identify Loops | Extract Index | Extract Termination Condition

```
void scale(float scale, float * in, float * out, int N) {  
    int i = 0  
    for ( ;  
          ;  
          i++)  
        if (i < N)  
            out[i] = scale * in[i];  
}
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

Identify Loops | Extract Index | Extract Termination Condition | Remove for

```
void scale(float scale, float * in, float * out, int N) {  
    int i = 0  
  
    if (i < N)  
        out[i] = scale * in[i];  
}
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

Identify Loops | Extract Index | Extract Termination Condition | Remove for | Add global

```
__global__ void scale(float scale, float * in, float * out, int N) {
    int i = 0

    if (i < N)
        out[i] = scale * in[i];
}
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

Identify Loops | Extract Index | Extract Termination Condition | Remove for | Add global

Replace i by threadIdx.x

```
__global__ void scale(float scale, float * in, float * out, int N) {
    int i = threadIdx.x;

    if (i < N)
        out[i] = scale * in[i];
}
```

# Kernel Conversion

Recipe for C Function → CUDA Kernel

Identify Loops | Extract Index | Extract Termination Condition | Remove for | Add global

Replace i by threadIdx.x | ... including block configuration

```
__global__ void scale(float scale, float * in, float * out, int N) {
    int i = threadIdx.x + blockIdx.x * blockDim.x;

    if (i < N)
        out[i] = scale * in[i];
}
```

# Kernel Conversion

## Summary

- C function with explicit loop

```
void scale(float scale, float * in, float * out, int N) {  
    for (int i = 0; i < N; i++)  
        out[i] = scale * in[i];  
}
```

- CUDA kernel with implicit loop

```
__global__ void scale(float scale, float * in, float * out, int N) {  
    int i = threadIdx.x + blockIdx.x * blockDim.x;  
    if (i < N)  
        out[i] = scale * in[i];  
}
```

# Kernel Launch

```
kernel<<<int blockDim, int blockDim>>>(...)
```

- Parallel threads of kernel launched with *triple-chevron syntax*
- Total number of threads, divided into
  - Number of blocks on the grid (gridDim)
  - Number of threads per block (blockDim)

# Kernel Launch

```
kernel<<<int blockDim, int blockDim>>>(...)
```

- Parallel threads of kernel launched with *triple-chevron syntax*
- Total number of threads, divided into
  - Number of blocks on the grid (gridDim)
  - Number of threads per block (blockDim)

# Kernel Launch

```
kernel<<<int blockDim, int blockDim>>>(...)
```

- Parallel threads of kernel launched with *triple-chevron syntax*
- Total number of threads, divided into
  - Number of blocks on the grid (gridDim)
  - Number of threads per block (blockDim)
- Call returns immediately; kernel launch is **asynchronous!**

# Kernel Launch

```
kernel<<<int blockDim, int blockDim>>>(...)
```

- Parallel threads of kernel launched with *triple-chevron syntax*
- Total number of threads, divided into
  - Number of blocks on the grid (gridDim)
  - Number of threads per block (blockDim)
- Call returns immediately; kernel launch is **asynchronous!**
- Example:

```
int nThreads = 32;  
scale<<<N/nThreads, nThreads>>>(23, in, out, N)
```

# Kernel Launch

```
kernel<<<int blockDim, int blockDim>>>(...)
```

- Parallel threads of kernel launched with *triple-chevron syntax*
- Total number of threads, divided into
  - Number of blocks on the grid (gridDim)
  - Number of threads per block (blockDim)
- Call returns immediately; kernel launch is **asynchronous!**
- Example:

```
int nThreads = 32;
scale<<<N/nThreads, nThreads>>>(23, in, out, N)
```
- Possibility for too many threads; include termination condition into kernel!

# Full Kernel Launch

## For Reference

```
kernel<<<dim3 gD, dim bD, size_t shared, cudaStream_t stream>>>(...)
```

- 2 additional, optional parameters

# Full Kernel Launch

## For Reference

```
kernel<<<dim3 gD, dim bD, size_t shared, cudaStream_t stream>>>(...)
```

- 2 additional, optional parameters

### shared Dynamic **shared memory**

- Small GPU memory space; share data in block (high bandwidth)
- Shared memory: allocate statically (compile time) or dynamically (run time)
- **size\_t shared**: bytes of shared memory allocated per block (in addition to static shared memory)

# Full Kernel Launch

## For Reference

```
kernel<<<dim3 gD, dim bD, size_t shared, cudaStream_t stream>>>(...)
```

- 2 additional, optional parameters

### shared Dynamic **shared memory**

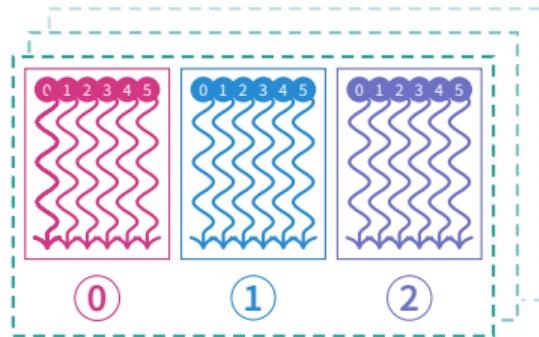
- Small GPU memory space; share data in block (high bandwidth)
- Shared memory: allocate statically (compile time) or dynamically (run time)
- **size\_t shared**: bytes of shared memory allocated per block (in addition to static shared memory)

### stream Associated **CUDA stream**

- CUDA streams enable different channels of communication with GPU
- Can overlap in some cases (communication, computation)
- **cudaStream\_t stream**: ID of stream to use for this kernel launch

# Grid Dimensions

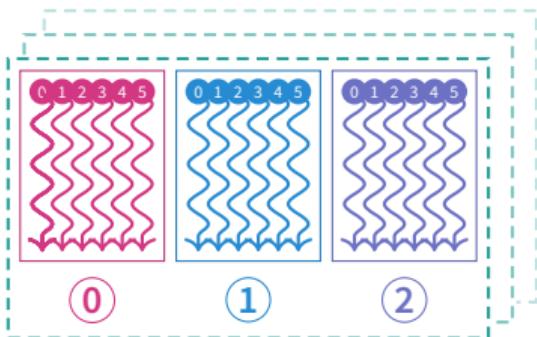
- Threads & blocks in 3D



# Grid Dimensions

- Threads & blocks in 3D
- Create 3D configurations with `struct dim3`

```
dim3 blockOrGridDim(size_t dimX, size_t dimY, size_t dimZ)
```



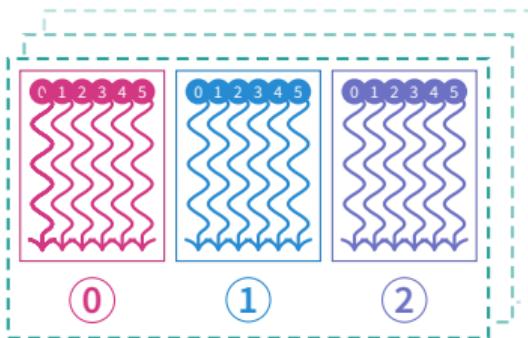
# Grid Dimensions

- Threads & blocks in 3D
- Create 3D configurations with `struct dim3`

```
dim3 blockDim(size_t dimX, size_t dimY, size_t dimZ)
```

- Example:

```
dim3 blockDim(32, 32);
dim3 gridDim = {1000, 100};
```



# Grid Dimensions

- Threads & blocks in 3D
- Create 3D configurations with `struct dim3`

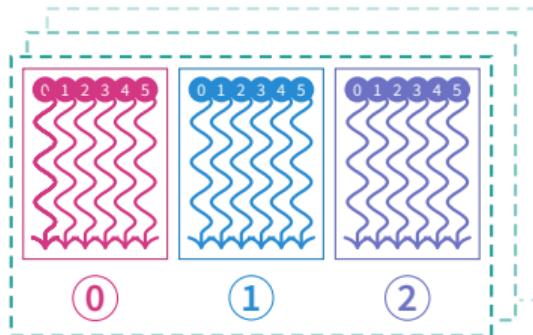
```
dim3 blockDim(size_t dimX, size_t dimY, size_t dimZ)
```

- Example:

```
dim3 blockDim(32, 32);  
dim3 gridDim = {1000, 100};
```

- Kernel call with `dim3`

```
kernel<<<dim3 gridDim, dim3 blockDim>>>(...)
```



# Grid Sizes

- Block and grid sizes are hardware-dependent

# Grid Sizes

- Block and grid sizes are hardware-dependent
- For JSC GPUs: Tesla V100, A100

Block      ■  $\vec{N}_{\text{Thread}} \leq (1024_x, 1024_y, 64_z)$   
             ■  $|\vec{N}_{\text{Thread}}| = N_{\text{Thread}} \leq 1024$

# Grid Sizes

- Block and grid sizes are hardware-dependent
- For JSC GPUs: Tesla V100, A100

Block      ■  $\vec{N}_{\text{Thread}} \leq (1024_x, 1024_y, 64_z)$

■  $|\vec{N}_{\text{Thread}}| = N_{\text{Thread}} \leq 1024$

Grid      ■  $\vec{N}_{\text{Blocks}} \leq (2147483647_x, 65535_y, 65535_z) = (2^{31}, 2^{16}, 2^{16}) - \vec{1}$

# Grid Sizes

- Block and grid sizes are hardware-dependent
- For JSC GPUs: Tesla V100, A100
  - **Block**      ■  $\vec{N}_{\text{Thread}} \leq (1024_x, 1024_y, 64_z)$
  - $|\vec{N}_{\text{Thread}}| = N_{\text{Thread}} \leq 1024$
  - **Grid**      ■  $\vec{N}_{\text{Blocks}} \leq (2147483647_x, 65535_y, 65535_z) = (2^{31}, 2^{16}, 2^{16}) - \vec{1}$
- Find out yourself: deviceQuery example from CUDA Samples

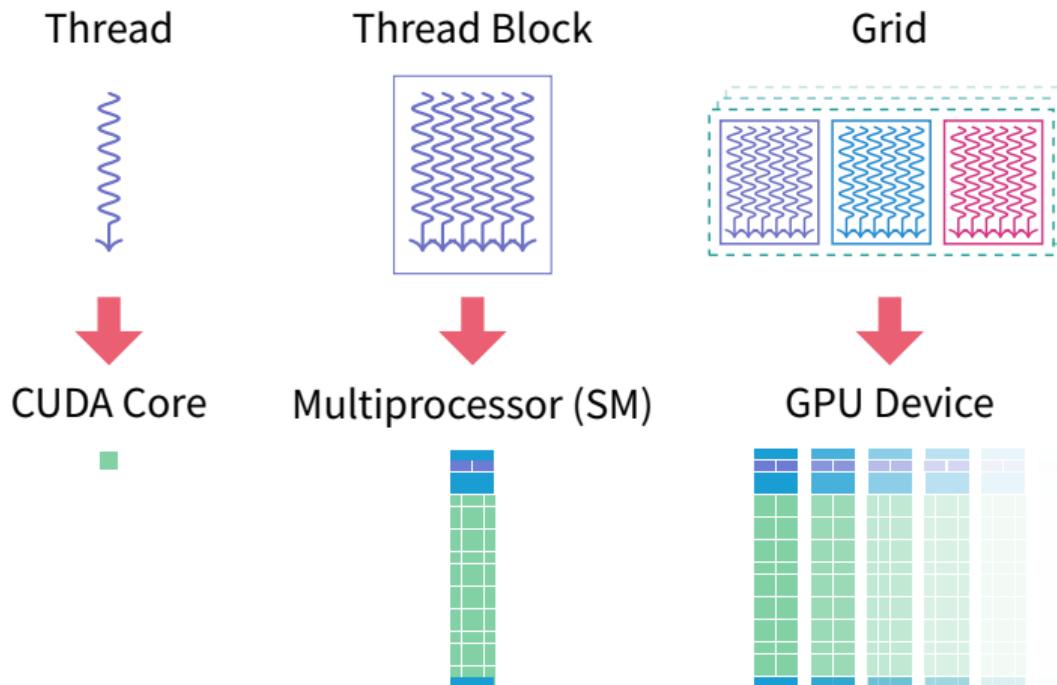
# Grid Sizes

- Block and grid sizes are hardware-dependent
- For JSC GPUs: Tesla V100, A100
  - Block      ■  $\vec{N}_{\text{Thread}} \leq (1024_x, 1024_y, 64_z)$
  - $|\vec{N}_{\text{Thread}}| = N_{\text{Thread}} \leq 1024$
- Grid      ■  $\vec{N}_{\text{Blocks}} \leq (2147483647_x, 65535_y, 65535_z) = (2^{31}, 2^{16}, 2^{16}) - 1$
- Find out yourself: deviceQuery example from CUDA Samples
- Workflow: Choose 128 or 256 as block dim; calculate grid dim from problem size

```
int Nx = 1000, Ny = 1000;
dim3 blockDim(16, 16);
int gx = (Nx % blockDim.x == 0) ? Nx / blockDim.x : Nx / blockDim.x + 1;
int gy = (Ny % blockDim.y == 0) ? Ny / blockDim.y : Ny / blockDim.y + 1;
dim3 gridDim(gx, gy);
kernel<<<gridDim, blockDim>>>();
```

# Hardware Threads

Mapping Software Threads to Hardware



# Memory Management

## With Automated Transfers

- Allocate memory to be used on GPU or CPU

```
cudaMallocManaged(T** ptr, size_t nBytes)
```

- Data is copied to GPU or to CPU automatically (*managed*)

# Memory Management

## With Automated Transfers

- Allocate memory to be used on GPU or CPU

```
cudaMallocManaged(T** ptr, size_t nBytes)
```

- Data is copied to GPU or to CPU automatically (*managed*)
- Example:

```
float * a;  
int N = 2048;  
cudaMallocManaged(&a, N * sizeof(float));
```

# Memory Management

## With Automated Transfers

- Allocate memory to be used on GPU or CPU

```
cudaMallocManaged(T** ptr, size_t nBytes)
```

- Data is copied to GPU or to CPU automatically (*managed*)

- Example:

```
float * a;  
int N = 2048;  
cudaMallocManaged(&a, N * sizeof(float));
```

- Free device memory

```
cudaFree(void* ptr)
```

# Memory Management

## With Manual Transfers

- Allocate memory to be used on GPU

```
cudaMalloc(T** ptr, size_t nBytes)
```

# Memory Management

## With Manual Transfers

- Allocate memory to be used on GPU

```
cudaMalloc(T** ptr, size_t nBytes)
```

- Copy data between host ↔ device

```
cudaMemcpy(void* dst, void* src, size_t nByte, enum cudaMemcpyKind dir)
```

# Memory Management

## With Manual Transfers

- Allocate memory to be used on GPU

```
cudaMalloc(T** ptr, size_t nBytes)
```

- Copy data between host ↔ device

```
cudaMemcpy(void* dst, void* src, size_t nByte, enum cudaMemcpyKind dir)
```

- Example:

```
float * a, * a_d;  
int N = 2048;  
// fill a  
cudaMalloc(&a_d, N * sizeof(float));  
cudaMemcpy(a_d, a, N * sizeof(float), cudaMemcpyHostToDevice);  
kernel<<<1,1>>>(a_d, N);  
cudaMemcpy(a , a_d, N * sizeof(float), cudaMemcpyDeviceToHost);
```

# Task: Scale Vector

Work on an Array of Data

- Location of code: 1-Basics/exercises/tasks/04-Scale-Vector
- Look at Instructions.ipynb for instructions
  - 1 Implement the whole CUDA flow (allocation, kernel configuration, kernel launch)
  - 2 Build with make
  - 3 Run with make run
- Additional task: Look at the version with explicit transfers (\_et)

# Task: Jacobi

## Implement Manual Memory Handling

- Location of code: 1-Basics/exercises/tasks/05-Jacobi-Explicit-Transfers
- Look at Instructions.ipynb for instructions
  - 1 Port the application from Unified Memory to manual memory handling
  - 2 Build with make
  - 3 Run with make run

# Unified Memory

## Overview

- Everything started with manual data management
- First Unified Memory since CUDA 6.0
- Better Unified Memory better since CUDA 8.0

# Manual Memory vs. Unified Memory

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    char *data_d;  
  
    data = (char *)malloc(N);  
    cudaMalloc(&data_d, N);  
  
    fread(data, 1, N, fp);  
  
    cudaMemcpy(data_d, data, N, cudaMemcpyHostToDevice);  
    kernel<<<...>>>(data, N);  
  
    cudaMemcpy(data, data_d, N, cudaMemcpyDeviceToHost);  
    host_func(data)  
    cudaFree(data_d); free(data);  
}
```

```
void sortfile(FILE *fp, int N) {  
    char *data;  
  
    cudaMallocManaged(&data, N);  
  
    fread(data, 1, N, fp);  
  
    kernel<<<...>>>(data, N);  
    cudaDeviceSynchronize();  
  
    host_func(data);  
    cudaFree(data);  
}
```

# Implementation Details

## Under the hood

```
cudaMallocManaged(&ptr, ...);
```

```
*ptr = 1;
```

```
kernel<<<...>>>(ptr);
```

# Implementation Details

## Under the hood

```
cudaMallocManaged(&ptr, ...); ← • Empty! No pages anywhere yet (like malloc())
```

```
*ptr = 1;
```

```
kernel<<<...>>>(ptr);
```

# Implementation Details

## Under the hood

```
cudaMallocManaged(&ptr, ...);
```

Empty! No pages anywhere yet (like malloc())

```
*ptr = 1;
```

CPU page fault: data allocates on CPU

```
kernel<<<...>>>(ptr);
```

# Implementation Details

## Under the hood

`cudaMallocManaged(&ptr, ...);`  Empty! No pages anywhere yet (like `malloc()`)

`*ptr = 1;`  CPU page fault: data allocates on CPU

`kernel<<<...>>>(ptr);`  GPU page fault: data migrates to GPU

# Implementation Details

## Under the hood

`cudaMallocManaged(&ptr, ...);`  Empty! No pages anywhere yet (like `malloc()`)

`*ptr = 1;`  CPU page fault: data allocates on CPU

`kernel<<<...>>>(ptr);`  GPU page fault: data migrates to GPU

- Pages populate on **first touch**
- Pages migrate on-demand
- GPU memory over-subscription possible
- Concurrent access from CPU and GPU to memory (page-level)

# Performance Analysis

Comparing `scale_vector_um` (Unified Memory) and `scale_vector` (manual copy) for 20 480 float elements.

UM

Time(%)	Total Time (ns)	Name
---------	-----------------	------

100.0	463,286	<code>scale(float, float*, float*, int)</code>
-------	---------	--

Manual

Time(%)	Total Time (ns)	Name
---------	-----------------	------

100.0	4,792	<code>scale(float, float*, float*, int)</code>
-------	-------	--

# Performance Analysis

Comparing `scale_vector_um` (Unified Memory) and `scale_vector` (manual copy) for 20 480 float elements.

UM

Time(%)	Total Time (ns)	Name
---------	-----------------	------

100.0	463,286	<code>scale(float, float*, float*, int)</code>
-------	---------	--

*100× slower?!*

What's going wrong here?

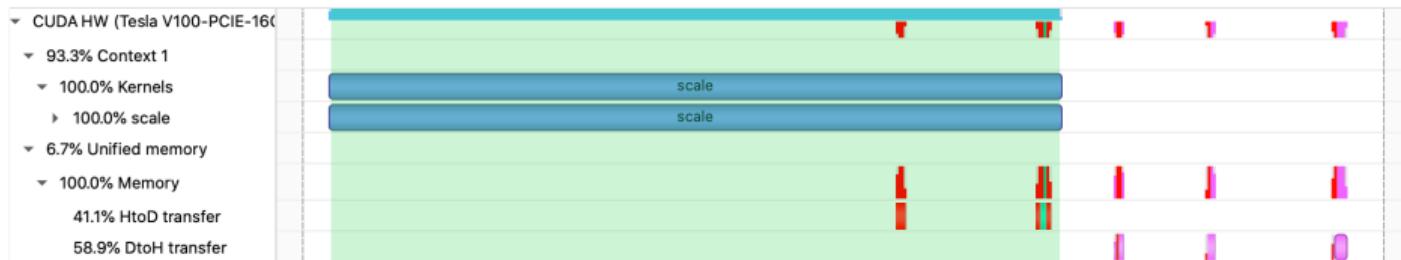
Manual

Time(%)	Total Time (ns)	Name
---------	-----------------	------

100.0	4,792	<code>scale(float, float*, float*, int)</code>
-------	-------	--

# Performance Analysis

Comparing `scale_vector_um` (Unified Memory) and `scale_vector` (manual copy) for 20 480 float elements.



UM

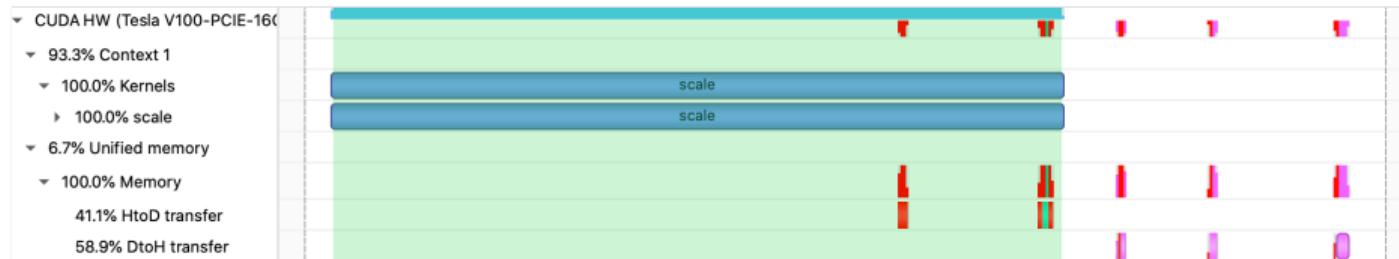
Manual

Time(%)	Total Time (ns)	Name
100.0	4,792	scale(float, float*, float*, int)

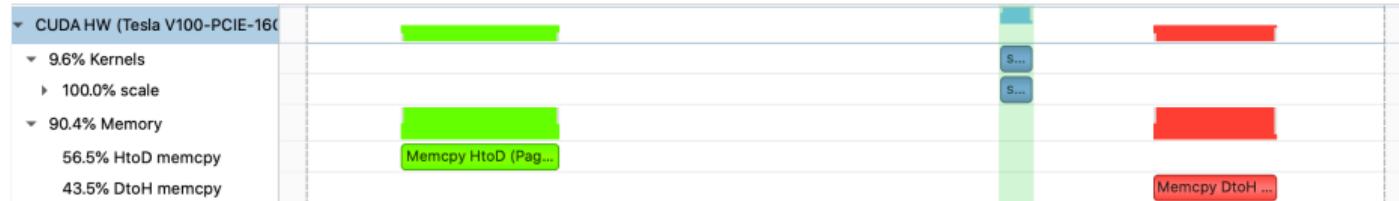
# Performance Analysis

Comparing `scale_vector_um` (Unified Memory) and `scale_vector` (manual copy) for 20 480 float elements.

UM



Manual



# Comparing UM and Explicit Transfers

**UM** Kernel is launched, data is needed by kernel, data migrates host→device  
⇒ Run time of kernel **incorporates** time for data transfers

**Explicit** Data will be needed by kernel – data migrates host→device **before** kernel launch  
⇒ Run time of **kernel** without any transfers

# Comparing UM and Explicit Transfers

**UM** Kernel is launched, data is needed by kernel, data migrates host→device  
⇒ Run time of kernel **incorporates** time for data transfers

**Explicit** Data will be needed by kernel – data migrates host→device **before** kernel launch  
⇒ Run time of **kernel** without any transfers

- UM more convenient
- Total run time of whole program does not principally change  
*Except: Fault handling costs  $\mathcal{O}(10 \mu\text{s})$ , stalls execution*
- But data transfers sometimes sorted to kernel launch

# Comparing UM and Explicit Transfers

**UM** Kernel is launched, data is needed by kernel, data migrates host→device  
⇒ Run time of kernel **incorporates** time for data transfers

**Explicit** Data will be needed by kernel – data migrates host→device **before** kernel launch  
⇒ Run time of **kernel** without any transfers

- UM more convenient
  - Total run time of whole program does not principally change  
*Except: Fault handling costs  $\mathcal{O}(10 \mu\text{s})$ , stalls execution*
  - But data transfers sometimes sorted to kernel launch
- ⇒ Improve UM behavior with performance hints!

# Performance Hints for UM

## New API routines

API calls to augment data location knowledge of runtime

- `cudaMemPrefetchAsync(data, length, device, stream)`  
Prefetches data to device (on stream) asynchronously

# Performance Hints for UM

## New API routines

API calls to augment data location knowledge of runtime

- `cudaMemPrefetchAsync(data, length, device, stream)`  
Prefetches data to device (on stream) asynchronously
- `cudaMemAdvise(data, length, advice, device)`  
Advise about usage of given data, advice:

# Performance Hints for UM

## New API routines

API calls to augment data location knowledge of runtime

- `cudaMemPrefetchAsync(data, length, device, stream)`  
Prefetches data to device (on stream) asynchronously
- `cudaMemAdvise(data, length, advice, device)`  
Advise about usage of given data, advice:
  - `cudaMemAdviseSetReadMostly`: Read-only copy is kept

# Performance Hints for UM

## New API routines

API calls to augment data location knowledge of runtime

- `cudaMemPrefetchAsync(data, length, device, stream)`  
Prefetches data to device (on stream) asynchronously
- `cudaMemAdvise(data, length, advice, device)`  
Advise about usage of given data, advice:
  - `cudaMemAdviseSetReadMostly`: Read-only copy is kept
  - `cudaMemAdviseSetPreferredLocation`: Set preferred location to avoid migrations; first access will establish mapping

# Performance Hints for UM

## New API routines

API calls to augment data location knowledge of runtime

- `cudaMemPrefetchAsync(data, length, device, stream)`  
Prefetches data to device (on stream) asynchronously
- `cudaMemAdvise(data, length, advice, device)`  
Advise about usage of given data, advice:
  - `cudaMemAdviseSetReadMostly`: Read-only copy is kept
  - `cudaMemAdviseSetPreferredLocation`: Set preferred location to avoid migrations; first access will establish mapping
  - `cudaMemAdviseSetAccessedBy`: Data is accessed by *this* device; will pre-map data to avoid page fault

# Performance Hints for UM

## New API routines

API calls to augment data location knowledge of runtime

- `cudaMemPrefetchAsync(data, length, device, stream)`  
Prefetches data to device (on stream) asynchronously
- `cudaMemAdvise(data, length, advice, device)`  
Advise about usage of given data, advice:
  - `cudaMemAdviseSetReadMostly`: Read-only copy is kept
  - `cudaMemAdviseSetPreferredLocation`: Set preferred location to avoid migrations; first access will establish mapping
  - `cudaMemAdviseSetAccessedBy`: Data is accessed by *this* device; will pre-map data to avoid page fault
- Use `cudaCpuDeviceId` for device CPU, or use `cudaGetDevice()` as usual to retrieve current GPU device id (default: 0)

# Hints in Code

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    // ...  
    cudaMallocManaged(&data, N);  
  
    fread(data, 1, N, fp);  
  
    cudaMemPrefetchAsync(data, N, device);  
    kernel<<<...>>>(data, N);  
    cudaDeviceSynchronize();  
  
    host_func(data);  
    cudaFree(data); }
```

# Hints in Code

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    // ...  
    cudaMallocManaged(&data, N);  
  
    fread(data, 1, N, fp);  
  
    cudaMemPrefetchAsync(data, N, device);  
    kernel<<<...>>>(data, N);  
    cudaDeviceSynchronize();  
  
    host_func(data);  
    cudaFree(data); }
```

Prefetch data to avoid expensive GPU page faults

# Hints in Code

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    // ...  
    cudaMallocManaged(&data, N);  
  
    fread(data, 1, N, fp);  
  
    cudaMemAdvise(data, N, cudaMemAdviseSetReadMostly, device);  
    cudaMemPrefetchAsync(data, N, device);  
    kernel<<<...>>>(data, N);  
    cudaDeviceSynchronize();  
  
    host_func(data);  
    cudaFree(data); }
```

Read-only copy of data  
is created on GPU during  
prefetch  
→ CPU and GPU reads will  
not fault

Prefetch data to avoid ex-  
pensive GPU page faults

# Tuning scale\_vector\_um

Express data movement

- Location of code: [1-Basics/exercises/tasks/06-Scale-Vector-Hints/](#)
- Look at `Instructions.ipynb` for instructions
  - 1 Task: Advise CUDA runtime that data should be migrated to GPU before kernel call
  - 2 Build with `make`
  - 3 Run with `make run`
  - 4 Glimpse at profile with `make profile`
- See also [CUDA C programming guide \(L.3.\)](#) for details on data performance tuning

# Conclusions

- GPUs achieve performance by specialized hardware
- Acceleration can be done by different means
- Libraries are the easiest
- Thrust, OpenACC can give first entry point
- Full power with CUDA
- Threads, Blocks to expose parallelism for a kernel
- Several API routines exist
- Unified Memory productive, possibly with hints

# Conclusions

- GPUs achieve performance by specialized hardware
- Acceleration can be done by different means
- Libraries are the easiest
- Thrust, OpenACC can give first entry point
- Full power with CUDA
- Threads, Blocks to expose parallelism for a kernel
- Several API routines exist
- Unified Memory productive, possibly with hints

Thank you  
for your attention!  
[a.herten@fz-juelich.de](mailto:a.herten@fz-juelich.de)

# Appendix

# Appendix

## Glossary

## References

# Glossary I

- AMD Manufacturer of **CPUs** and **GPUs**. [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#)
- Ampere **GPU** architecture from **NVIDIA** (announced 2019). [13](#), [14](#), [15](#)
- API A programmatic interface to software by well-defined functions. Short for application programming interface. [176](#)
- ATI Canada-based **GPUs** manufacturing company; bought by AMD in 2006. [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#)
- CUDA Computing platform for **GPUs** from NVIDIA. Provides, among others, CUDA C/C++. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [83](#), [84](#), [93](#), [97](#), [98](#), [99](#), [100](#), [101](#), [102](#), [103](#), [104](#), [105](#), [106](#), [107](#), [125](#), [126](#), [127](#), [132](#), [133](#), [134](#), [135](#), [136](#), [144](#), [169](#), [170](#), [171](#), [175](#)
- JSC Jülich Supercomputing Centre, the supercomputing institute of Forschungszentrum Jülich, Germany. [175](#)

# Glossary II

JURECA A multi-purpose supercomputer at JSC. [15](#)

JUWELS Jülich's new supercomputer, the successor of JUQUEEN. [12](#), [13](#), [14](#)

NVIDIA US technology company creating GPUs. [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [12](#), [13](#), [14](#), [15](#), [50](#), [51](#), [52](#),  
[174](#), [175](#), [176](#), [177](#)

NVLink NVIDIA's communication protocol connecting CPU ↔ GPU and GPU ↔ GPU with high bandwidth. [177](#)

OpenACC Directive-based programming, primarily for many-core machines. [83](#), [84](#), [86](#), [87](#),  
[88](#), [89](#), [90](#), [91](#), [170](#), [171](#)

OpenCL The *Open Computing Language*. Framework for writing code for heterogeneous architectures (CPU, GPU, DSP, FPGA). The alternative to CUDA. [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#),  
[83](#), [84](#)

# Glossary III

- OpenGL** The *Open Graphics Library*, an API for rendering graphics across different hardware architectures. [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#)
- OpenMP** Directive-based programming, primarily for multi-threaded machines. [83](#), [84](#), [86](#), [87](#), [88](#), [89](#)
- SAXPY** Single-precision  $A \times X + Y$ . A simple code example of scaling a vector and adding an offset. [58](#), [98](#)
- Tesla** The GPU product line for general purpose computing computing of NVIDIA. [12](#), [132](#), [133](#), [134](#), [135](#), [136](#)
- Thrust** A parallel algorithms library for (among others) GPUs. See <https://thrust.github.io/>. [83](#), [84](#), [93](#), [95](#), [170](#), [171](#)

# Glossary IV

**V100** A large **GPU** with the **Volta** architecture from **NVIDIA**. It employs **NVLink 2** as its interconnect and has fast **HBM2** memory. Additionally, it features **Tensorcores** for Deep Learning and Independent Thread Scheduling. [132](#), [133](#), [134](#), [135](#), [136](#)

**Volta** **GPU** architecture from **NVIDIA** (announced 2017). [177](#)

**CPU** Central Processing Unit. [12](#), [15](#), [20](#), [21](#), [22](#), [26](#), [27](#), [28](#), [29](#), [30](#), [31](#), [32](#), [33](#), [34](#), [35](#), [36](#), [37](#), [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#), [58](#), [89](#), [95](#), [138](#), [139](#), [140](#), [148](#), [149](#), [150](#), [151](#), [152](#), [160](#), [161](#), [162](#), [163](#), [164](#), [165](#), [168](#), [174](#), [175](#)

# Glossary V

- GPU** Graphics Processing Unit. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#), [12](#), [13](#), [14](#), [15](#), [16](#), [17](#), [18](#), [20](#), [21](#), [22](#), [23](#), [24](#), [25](#), [26](#), [27](#), [28](#), [29](#), [30](#), [31](#), [32](#), [33](#), [34](#), [35](#), [36](#), [37](#), [38](#), [39](#), [41](#), [42](#), [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#), [53](#), [54](#), [55](#), [57](#), [61](#), [62](#), [63](#), [64](#), [65](#), [66](#), [77](#), [85](#), [86](#), [87](#), [88](#), [89](#), [92](#), [95](#), [97](#), [108](#), [109](#), [110](#), [125](#), [126](#), [127](#), [132](#), [133](#), [134](#), [135](#), [136](#), [138](#), [139](#), [140](#), [141](#), [142](#), [143](#), [148](#), [149](#), [150](#), [151](#), [152](#), [160](#), [161](#), [162](#), [163](#), [164](#), [165](#), [167](#), [168](#), [169](#), [170](#), [171](#), [174](#), [175](#), [176](#), [177](#)
- SIMD** Single Instruction, Multiple Data. [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#)
- SIMT** Single Instruction, Multiple Threads. [23](#), [24](#), [25](#), [38](#), [39](#), [41](#), [42](#), [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#)
- SM** Streaming Multiprocessor. [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#)
- SMT** Simultaneous Multithreading. [43](#), [44](#), [45](#), [46](#), [47](#), [48](#), [49](#), [50](#), [51](#), [52](#)

# References I

- [2] Kenneth E. Hoff III et al. “Fast Computation of Generalized Voronoi Diagrams Using Graphics Hardware.” In: *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*. SIGGRAPH ’99. New York, NY, USA: ACM Press/Addison-Wesley Publishing Co., 1999, pp. 277–286. ISBN: 0-201-48560-5. DOI: [10.1145/311535.311567](https://doi.org/10.1145/311535.311567). URL: <http://dx.doi.org/10.1145/311535.311567> (pages 3–9).
- [3] Chris McClanahan. “History and Evolution of GPU Architecture.” In: *A Survey Paper* (2010). URL: <http://mcclanahoochie.com/blog/wp-content/uploads/2011/03/gpu-hist-paper.pdf> (pages 3–9).
- [4] Jack Dongarra et al. *TOP500*. Nov. 2016. URL: <https://www.top500.org/lists/2016/11/> (pages 3–9).
- [5] Jack Dongarra et al. *Green500*. Nov. 2016. URL: <https://www.top500.org/green500/lists/2016/11/> (pages 3–9).

## References II

- [6] Karl Rupp. *Pictures: CPU/GPU Performance Comparison*. URL:  
<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/> (pages 10, 11).
- [9] Wes Brezell. *Picture: Wizard*. URL:  
<https://thenounproject.com/wes13/collection/its-a-wizards-world/> (pages 61–65).

# References: Images, Graphics I

- [1] Héctor J. Rivas. *Color Reels*. Freely available at Unsplash. URL:  
<https://unsplash.com/photos/87hFrPk3V-s>.
- [7] Mark Lee. *Picture: kawasaki ninja*. URL:  
<https://www.flickr.com/photos/pochacco20/39030210/> (pages 20, 21).
- [8] Shearings Holidays. *Picture: Shearings coach 636*. URL:  
<https://www.flickr.com/photos/shearings/13583388025/> (pages 20, 21).