# INTRODUCTION TO
# GPU PROGRAMMING

# EXPECTATIONS

- A very broad introduction to the subject.

# EXPECTATIONS

- A very broad introduction to the subject.
- Not a specifically CUDA workshop.

# EXPECTATIONS

- A very broad introduction to the subject.
- <u>Not</u> a specifically CUDA workshop.
- <u>Not</u> a machine learning workshop.

# EXPECTATIONS

- A very broad introduction to the subject.
- <u>Not</u> a specifically CUDA workshop.
- <u>Not</u> a machine learning workshop.
- <u>Not</u> an ad for GPUs.

# EXPECTATIONS

- A very broad introduction to the subject.
- <u>Not</u> a specifically CUDA workshop.
- <u>Not</u> a machine learning workshop.
- <u>Not</u> an ad for GPUs.
- Focusing on GPUs in high-performance computing (HPC).

# INTRODUCTION

# WHAT IS A GPU AND WHAT IS IT GOOD FOR?

- An electronic circuit.

# WHAT IS A GPU AND WHAT IS IT GOOD FOR?

- An electronic circuit.
- Originally a *graphics* processing unit.
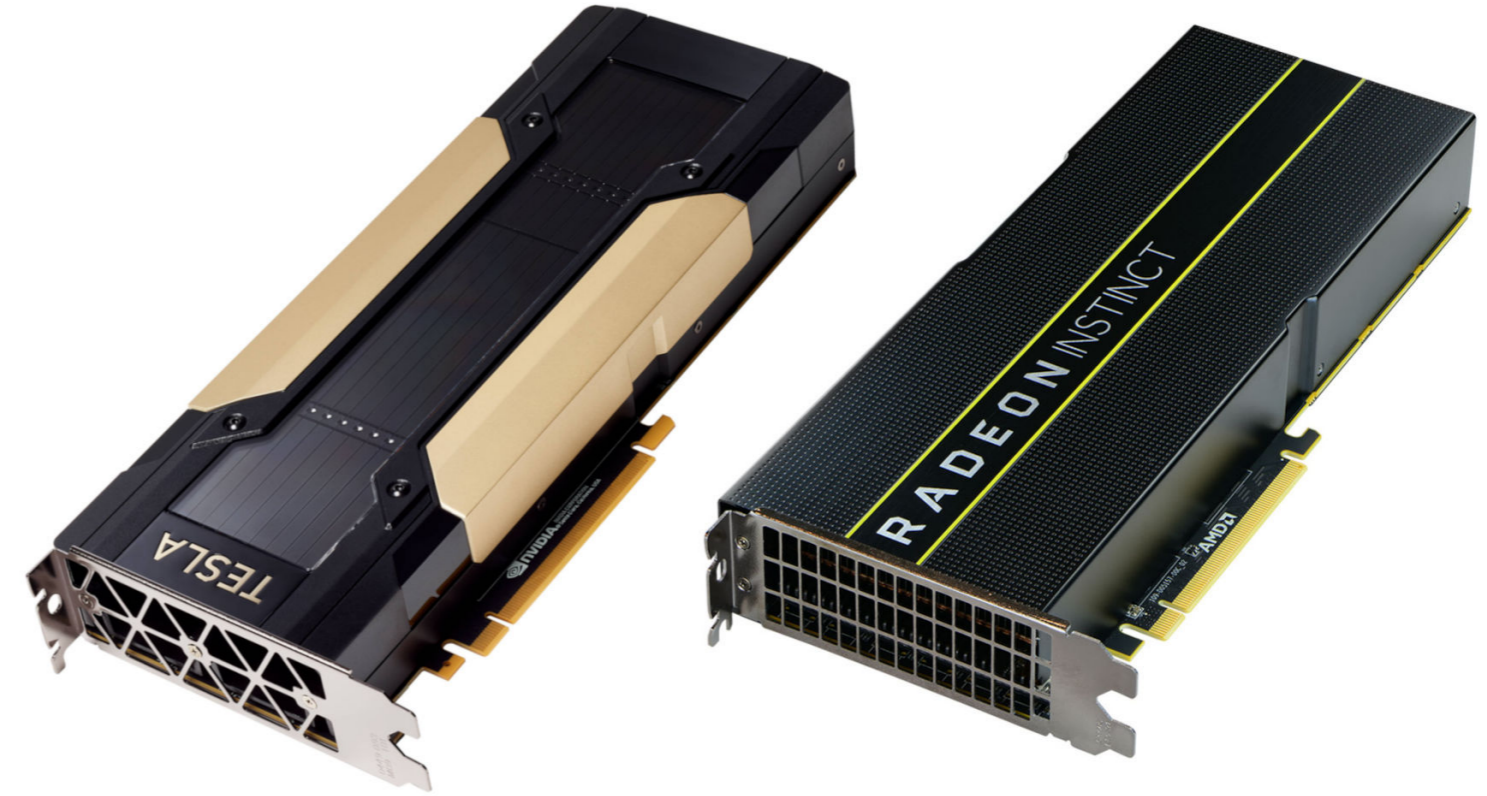
# WHAT IS A GPU AND WHAT IS IT GOOD FOR?

- An electronic circuit.
- Originally a *graphics* processing unit.
- Accelerates computation through *parallelization*.

# WHAT IS A GPU AND WHAT IS IT GOOD FOR?

- An electronic circuit.
- Originally a *graphics* processing unit.
- Accelerates computation through *parallelization*.
- Supports *general purpose* computations.

# WHAT IS A GPU AND WHAT IS IT GOOD FOR?

- An electronic circuit.
- Originally a *graphics* processing unit.
- Accelerates computation through *parallelization*.
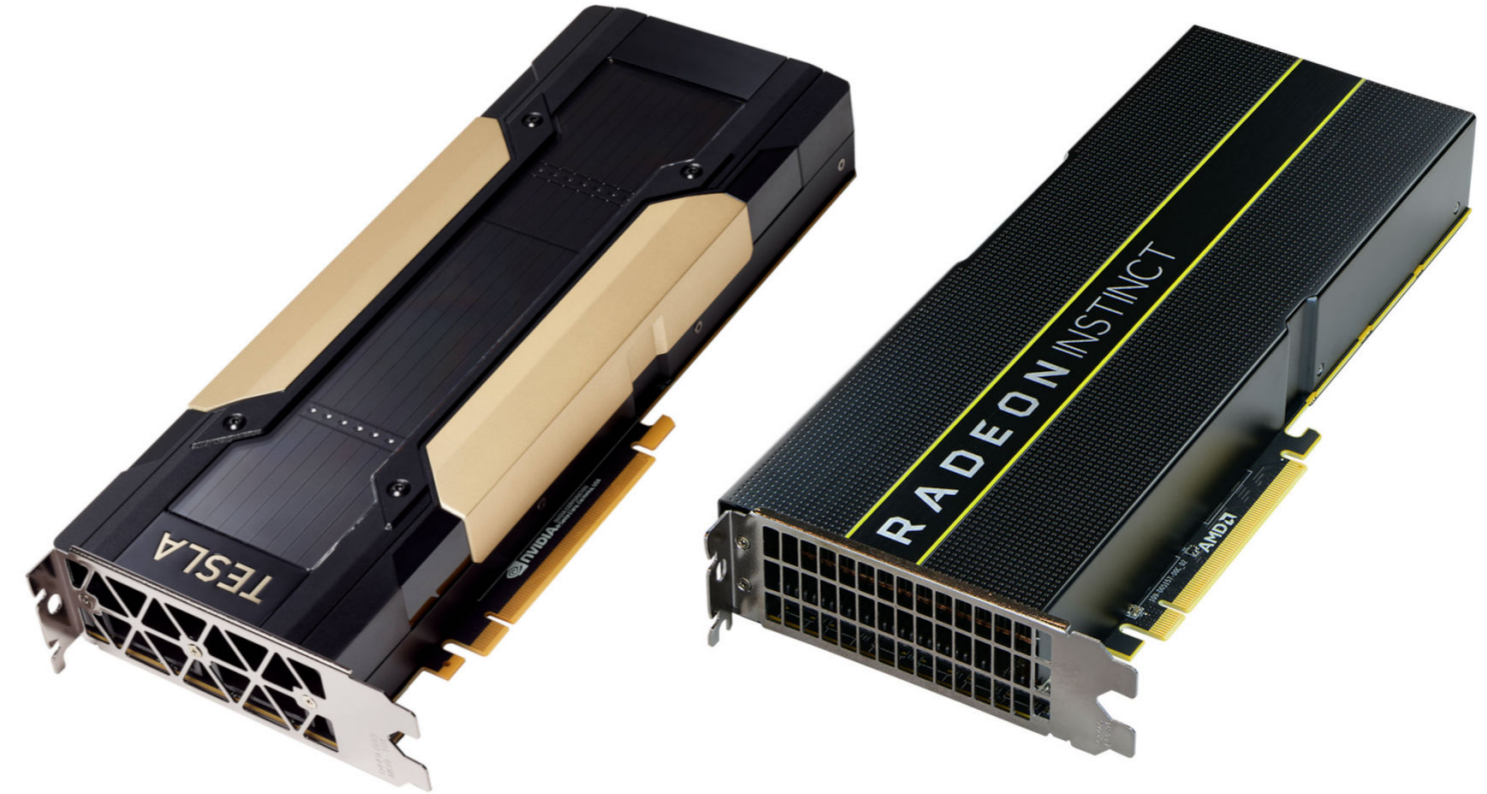- Supports *general purpose* computations.

# WHAT IS A GPU AND WHAT IS IT GOOD FOR?

- An electronic circuit.
- Originally a *graphics* processing unit.
- Accelerates computation through *parallelization*.
- Supports *general purpose* computations.

Wait, where are the fans? 🤨

…and the graphics output???

# VECTORIZATION

## C CODE

```c
void add(float* restrict a, float* restrict b)
{
  for (int i=0; i<16; i++)
    b[i] = a[i] + b[i];
}
```

How many *add* instructions?

# VECTORIZATION

## C CODE

```c
void add(float* restrict a, float* restrict b)
{
  for (int i=0; i<16; i++)
    b[i] = a[i] + b[i];
}
```

How many *add* instructions?

Compile with: `gcc -O3 -mavx512f`

# VECTORIZATION

## C CODE

```c
void add(float* restrict a, float* restrict b)
{
  for (int i=0; i<16; i++)
    b[i] = a[i] + b[i];
}
```

How many *add* instructions?

Compile with: `gcc -O3 -mavx512f`

## GENERATED CPU INSTRUCTIONS

```
vmovups zmm0, ZMMWORD PTR [rdi]
vaddps  zmm0, zmm0, ZMMWORD PTR [rsi]
vmovups ZMMWORD PTR [rsi], zmm0
; ...
```

On **Niagara**: *one* add instruction!

# VECTORIZATION

## C CODE

```c
void add(float* restrict a, float* restrict b)
{
  for (int i=0; i<16; i++)
    b[i] = a[i] + b[i];
}
```

How many *add* instructions?

Compile with: `gcc -O3 -mavx512f`

## GENERATED CPU INSTRUCTIONS

```asm
vmovups zmm0, ZMMWORD PTR [rdi]
vaddps  zmm0, zmm0, ZMMWORD PTR [rsi]
vmovups ZMMWORD PTR [rsi], zmm0
; ...
```

On **Niagara**: *one* add instruction!

single instruction, multiple data (*SIMD*)

# VECTORIZATION

## C CODE

```c
void add(float* restrict a, float* restrict b)
{
  for (int i=0; i<16; i++)
    b[i] = a[i] + b[i];
}
```

How many *add* instructions?

Compile with: `gcc -O3 -mavx512f`

## GENERATED CPU INSTRUCTIONS

```
vmovups zmm0, ZMMWORD PTR [rdi]
vaddps  zmm0, zmm0, ZMMWORD PTR [rsi]
vmovups ZMMWORD PTR [rsi], zmm0
; ...
```

On **Niagara**: *one* add instruction!

single instruction, multiple data (*SIMD*)

- CPU programming: have to *convince the compiler* to vectorize.
- GPU programming: vectorization *by default*.

# CPU VS. GPU

| Intel Xeon Gold 6148 *CPU* | Nvidia Tesla V100 *GPU* |
|---|---|
| 20 cores | 84 SMs |
| 2 × 512 bit SIMD units / core | 4 × 512 bit SIMD units / SM |
| 2.4 GHz base clock (3.7 GHz turbo) | 0.9 GHz base clock (1.3 GHz boost) |
| 128 GB/s max. bandwidth | 900 GB/s max. bandwidth |

# PROGRAMMING WITH A CO-PROCESSOR

- Examples of hardware that can be used as a co-processors:
  - GPU, MIC, FPGA, ASIC, VE...

# PROGRAMMING WITH A CO-PROCESSOR

- Examples of hardware that can be used as a co-processors:
  - GPU, MIC, FPGA, ASIC, VE...
- Architecture agnostic code $\longrightarrow$ architecture specific instructions.
  - A program has *host code* and *device code*.

# PROGRAMMING WITH A CO-PROCESSOR

- Examples of hardware that can be used as a co-processors:
  - GPU, MIC, FPGA, ASIC, VE...
- Architecture agnostic code ⟶ architecture specific instructions.
  - A program has *host code* and *device code*.
- Host and device memories are *separately addressed*.
  - *Copying* data between memory spaces is (usually) required.

# THE THREAD MODEL

~~~ thread ~~~

A *thread* is a serial stream of instructions.

Threads should <u>oversubscribe</u> the "CUDA cores" (stream processors).

# THE THREAD MODEL

~~~ thread ~~~

Multiple threads are arranged in a *block*.

Threads in one block run on the same SM (compute unit).

Within a block they can <u>synchronize</u> and <u>share</u> cache memory.

# THE THREAD MODEL

$\sim\!\!\sim\!\!\sim$ thread $\sim\!\!\sim\!\!\sim$

Multiple blocks are arranged in a *grid*.

Threads in different blocks <u>cannot</u> synchronize.

# THE THREAD MODEL

$$\sim\!\sim\!\sim \text{thread} \sim\!\sim\!\sim$$

A *host* may have different grids running on separate GPUs.

The CPU cores can work simultaneously.

# THE THREAD MODEL

$\sim\sim$ thread $\sim\sim$

The *cluster* has multiple hosts.

Use MPI or NCCL/RCCL for collective communication.

# THE THREAD MODEL (SUMMARY)

- Threads running in lockstep in a warp.
- One or more warps in block.
- Multiple blocks in a grid.
- Multiple grids on a host.
- The host has CPU threads as well.
- Multiple hosts in the cluster.

# SETUP

# COMPUTING RESOURCES FOR THE EXERCISE

# COMPUTING RESOURCES FOR THE EXERCISE

- ComputeCanada users: use one of the national systems.

# COMPUTING RESOURCES FOR THE EXERCISE

- ComputeCanada users: use one of the national systems.
- Non-ComputeCanada users: use your guest account on Graham.

# COMPUTING RESOURCES FOR THE EXERCISE

- ComputeCanada users: use one of the national systems.
- Non-ComputeCanada users: use your guest account on Graham.
- You can also use your own workstation,

# COMPUTING RESOURCES FOR THE EXERCISE

- ComputeCanada users: use one of the national systems.
- Non-ComputeCanada users: use your guest account on Graham.
- You can also use your own workstation,
- or free GPU cloud resources like Colab, Kaggle, or SageMaker Studio Lab

# COMPUTING RESOURCES FOR THE EXERCISE

- ComputeCanada users: use one of the national systems.
- Non-ComputeCanada users: use your guest account on Graham.
- You can also use your own workstation,
- or free GPU cloud resources like Colab, Kaggle, or SageMaker Studio Lab

Exact type of GPU generally doesn't matter, but **not all platforms support all frameworks**.

# COMPUTING RESOURCES FOR THE EXERCISE

- ComputeCanada users: use one of the national systems.
- Non-ComputeCanada users: use your guest account on Graham.
- You can also use your own workstation,
- or free GPU cloud resources like Colab, Kaggle, or SageMaker Studio Lab

Exact type of GPU generally doesn't matter, but **not all platforms support all frameworks**.

Useful links:

https://docs.computecanada.ca/wiki/Using_GPUs_with_Slurm
https://docs.computecanada.ca/wiki/Python#Creating_and_using_a_virtual_environment
https://docs.scinet.utoronto.ca/index.php/Mist

# SETTING UP THE ENVIRONMENT (PYTHON)

```
mkdir $SCRATCH/scinet-hpc133
cd $SCRATCH/scinet-hpc133
```

**Graham**
```
module load cuda python scipy-stack
virtualenv --no-download virtualenv
source virtualenv/bin/activate
pip install --no-index --upgrade pip numba cupy
```

**Mist**
```
module load anaconda3
conda create -p condaenv -c conda-forge \
    python=3.8.14 numba cupy cudatoolkit=11.0.3 -y
source activate ./condaenv
conda clean --all -y
rm -rf ~/.conda/pkgs/*
```

```
wget scinet.courses/mod/2253 -O src.tar.bz2
tar xf src.tar.bz2
```

**Graham**

*Running from login node*

```
srun --time=00:01:00 --gres=gpu:p100:1 \
    python src/vector_add/08_numba.py
```

*Interactive job*

```
salloc --time=00:15:00 --gres=gpu:p100:1
```

\* you may have to specify --account

**Mist**

*Running from login node*

```
python src/vector_add/08_numba.py
```

*Interactive job*

```
debugjob -g 1
```

# GPU PROGRAMMING FRAMEWORKS

# THE FRAMEWORKS



## + programs & libraries!

# THE FRAMEWORKS

NVIDIA CUDA · ROCm · OpenCL™ · **OpenACC** More Science, Less Programming · OpenMP® · SYCL™ · Numba

<u>+ programs & libraries!</u>

- Theoretically, it doesn't matter which one you choose.
  - In practice, performance can vary.

# THE FRAMEWORKS



## + programs & libraries!

- Theoretically, it doesn't matter which one you choose.
  - In practice, performance can vary.
- Threads execute small programs called *kernels*.

# THE FRAMEWORKS



## + programs & libraries!

- Theoretically, it doesn't matter which one you choose.
  - In practice, performance can vary.
- Threads execute small programs called *kernels*.
- *Memory transfer* may be explicit or implicit.

```cpp
#include <algorithm>
#include <iostream>

__global__ void add(float *a, float *b)
{
  int i = blockDim.x * blockIdx.x + threadIdx.x;
  b[i] = a[i] + b[i];
}

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  float *a_dev, *b_dev;
  cudaMalloc((void**)&a_dev, n*sizeof(float));
  cudaMalloc((void**)&b_dev, n*sizeof(float));
  cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);

  add<<<n/64,64>>>(a_dev, b_dev);

  cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  cudaFree(a_dev);
  cudaFree(b_dev);
}
```

CUDA

- Oldest and most mature.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  __global__ void add(float *a, float *b)
5  {
6    int i = blockDim.x * blockIdx.x + threadIdx.x;
7    b[i] = a[i] + b[i];
8  }
9
10 int main()
11 {
12   constexpr size_t n{1'000'000};
13   float *a{new float[n]}, *b{new float[n]};
14   std::fill(a, a+n, 2.);
15   std::fill(b, b+n, 3.);
16
17   float *a_dev, *b_dev;
18   cudaMalloc((void**)&a_dev, n*sizeof(float));
19   cudaMalloc((void**)&b_dev, n*sizeof(float));
20   cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
21   cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);
22
23   add<<<n/64,64>>>(a_dev, b_dev);
24
25   cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);
26
27   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
28     std::cout << "Success!\n";
29   else std::cout << "Failure :(\n";
30
31   delete[] a;
32   delete[] b;
33   cudaFree(a_dev);
34   cudaFree(b_dev);
35 }
```

- Oldest and most mature.
- Proprietary & Nvidia-specific.

```cpp
#include <algorithm>
#include <iostream>

__global__ void add(float *a, float *b)
{
  int i = blockDim.x * blockIdx.x + threadIdx.x;
  b[i] = a[i] + b[i];
}

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  float *a_dev, *b_dev;
  cudaMalloc((void**)&a_dev, n*sizeof(float));
  cudaMalloc((void**)&b_dev, n*sizeof(float));
  cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);

  add<<<n/64,64>>>(a_dev, b_dev);

  cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  cudaFree(a_dev);
  cudaFree(b_dev);
}
```

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
 1  #include <algorithm>
 2  #include <iostream>
 3
 4  __global__ void add(float *a, float *b)
 5  {
 6    int i = blockDim.x * blockIdx.x + threadIdx.x;
 7    b[i] = a[i] + b[i];
 8  }
 9
10  int main()
11  {
12    constexpr size_t n{1'000'000};
13    float *a{new float[n]}, *b{new float[n]};
14    std::fill(a, a+n, 2.);
15    std::fill(b, b+n, 3.);
16
17    float *a_dev, *b_dev;
18    cudaMalloc((void**)&a_dev, n*sizeof(float));
19    cudaMalloc((void**)&b_dev, n*sizeof(float));
20    cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
21    cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);
22
23    add<<<n/64,64>>>(a_dev, b_dev);
24
25    cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);
26
27    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
28      std::cout << "Success!\n";
29    else std::cout << "Failure :(\n";
30
31    delete[] a;
32    delete[] b;
33    cudaFree(a_dev);
34    cudaFree(b_dev);
35  }
```

NVIDIA
CUDA

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
#include <algorithm>
#include <iostream>

__global__ void add(float *a, float *b)
{
  int i = blockDim.x * blockIdx.x + threadIdx.x;
  b[i] = a[i] + b[i];
}

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  float *a_dev, *b_dev;
  cudaMalloc((void**)&a_dev, n*sizeof(float));
  cudaMalloc((void**)&b_dev, n*sizeof(float));
  cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);

  add<<<n/64,64>>>(a_dev, b_dev);

  cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  cudaFree(a_dev);
  cudaFree(b_dev);
}
```

**NVIDIA CUDA**

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  __global__ void add(float *a, float *b)
5  {
6    int i = blockDim.x * blockIdx.x + threadIdx.x;
7    b[i] = a[i] + b[i];
8  }
9
10 int main()
11 {
12   constexpr size_t n{1'000'000};
13   float *a{new float[n]}, *b{new float[n]};
14   std::fill(a, a+n, 2.);
15   std::fill(b, b+n, 3.);
16
17   float *a_dev, *b_dev;
18   cudaMalloc((void**)&a_dev, n*sizeof(float));
19   cudaMalloc((void**)&b_dev, n*sizeof(float));
20   cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
21   cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);
22
23   add<<<n/64,64>>>(a_dev, b_dev);
24
25   cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);
26
27   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
28     std::cout << "Success!\n";
29   else std::cout << "Failure :(\n";
30
31   delete[] a;
32   delete[] b;
33   cudaFree(a_dev);
34   cudaFree(b_dev);
35 }
```

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
#include <algorithm>
#include <iostream>

__global__ void add(float *a, float *b)
{
  int i = blockDim.x * blockIdx.x + threadIdx.x;
  b[i] = a[i] + b[i];
}

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  float *a_dev, *b_dev;
  cudaMalloc((void**)&a_dev, n*sizeof(float));
  cudaMalloc((void**)&b_dev, n*sizeof(float));
  cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);

  add<<<n/64,64>>>(a_dev, b_dev);

  cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  cudaFree(a_dev);
  cudaFree(b_dev);
}
```

**NVIDIA CUDA**

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  __global__ void add(float *a, float *b)
5  {
6    int i = blockDim.x * blockIdx.x + threadIdx.x;
7    b[i] = a[i] + b[i];
8  }
9
10 int main()
11 {
12   constexpr size_t n{1'000'000};
13   float *a{new float[n]}, *b{new float[n]};
14   std::fill(a, a+n, 2.);
15   std::fill(b, b+n, 3.);
16
17   float *a_dev, *b_dev;
18   cudaMalloc((void**)&a_dev, n*sizeof(float));
19   cudaMalloc((void**)&b_dev, n*sizeof(float));
20   cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
21   cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);
22
23   add<<<n/64,64>>>(a_dev, b_dev);
24
25   cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);
26
27   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
28     std::cout << "Success!\n";
29   else std::cout << "Failure :(\n";
30
31   delete[] a;
32   delete[] b;
33   cudaFree(a_dev);
34   cudaFree(b_dev);
35 }
```

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
#include <algorithm>
#include <iostream>

__global__ void add(float *a, float *b)
{
  int i = blockDim.x * blockIdx.x + threadIdx.x;
  b[i] = a[i] + b[i];
}

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  float *a_dev, *b_dev;
  cudaMalloc((void**)&a_dev, n*sizeof(float));
  cudaMalloc((void**)&b_dev, n*sizeof(float));
  cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);

  add<<<n/64,64>>>(a_dev, b_dev);

  cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  cudaFree(a_dev);
  cudaFree(b_dev);
}
```

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
#include <algorithm>
#include <iostream>

__global__ void add(float *a, float *b)
{
  int i = blockDim.x * blockIdx.x + threadIdx.x;
  b[i] = a[i] + b[i];
}

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  float *a_dev, *b_dev;
  cudaMalloc((void**)&a_dev, n*sizeof(float));
  cudaMalloc((void**)&b_dev, n*sizeof(float));
  cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);

  add<<<n/64,64>>>(a_dev, b_dev);

  cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  cudaFree(a_dev);
  cudaFree(b_dev);
}
```

NVIDIA CUDA

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  __global__ void add(float *a, float *b)
5  {
6    int i = blockDim.x * blockIdx.x + threadIdx.x;
7    b[i] = a[i] + b[i];
8  }
9
10 int main()
11 {
12   constexpr size_t n{1'000'000};
13   float *a{new float[n]}, *b{new float[n]};
14   std::fill(a, a+n, 2.);
15   std::fill(b, b+n, 3.);
16
17   float *a_dev, *b_dev;
18   cudaMalloc((void**)&a_dev, n*sizeof(float));
19   cudaMalloc((void**)&b_dev, n*sizeof(float));
20   cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
21   cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);
22
23   add<<<n/64,64>>>(a_dev, b_dev);
24
25   cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);
26
27   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
28     std::cout << "Success!\n";
29   else std::cout << "Failure :(\n";
30
31   delete[] a;
32   delete[] b;
33   cudaFree(a_dev);
34   cudaFree(b_dev);
35 }
```

- Oldest and most mature.
- Proprietary & Nvidia-specific.
- Extension of C++.

```cpp
#include <algorithm>
#include <iostream>

__global__ void add(float *a, float *b)
{
  int i = blockDim.x * blockIdx.x + threadIdx.x;
  b[i] = a[i] + b[i];
}

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  float *a_dev, *b_dev;
  cudaMalloc((void**)&a_dev, n*sizeof(float));
  cudaMalloc((void**)&b_dev, n*sizeof(float));
  cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);

  add<<<n/64,64>>>(a_dev, b_dev);

  cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  cudaFree(a_dev);
  cudaFree(b_dev);
}
```

```cpp
1  #include <algorithm>
2  #include <iostream>
3  #include <thrust/host_vector.h>
4  #include <thrust/device_vector.h>
5  #include <thrust/functional.h>
6  #include <thrust/transform.h>
7
8  int main()
9  {
10   constexpr size_t n{1'000'000};
11   thrust::host_vector<float> a(n, 2.), b(n, 3.);
12
13   thrust::device_vector<float> a_dev{a}, b_dev{b};
14
15   thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
16                     b_dev.begin(), thrust::plus<float>());
17
18   b = b_dev;
19
20   if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
21     std::cout << "Success!\n";
22   else std::cout << "Failure :(\n";
23 }
```

- Thrust is included in CUDA.

```cpp
#include <algorithm>
#include <iostream>
#include <thrust/host_vector.h>
#include <thrust/device_vector.h>
#include <thrust/functional.h>
#include <thrust/transform.h>

int main()
{
  constexpr size_t n{1'000'000};
  thrust::host_vector<float> a(n, 2.), b(n, 3.);

  thrust::device_vector<float> a_dev{a}, b_dev{b};

  thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
                    b_dev.begin(), thrust::plus<float>());

  b = b_dev;

  if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";
}
```

- Thrust is included in CUDA.
- Provides STL-like abstractions.

```cpp
1  #include <algorithm>
2  #include <iostream>
3  #include <thrust/host_vector.h>
4  #include <thrust/device_vector.h>
5  #include <thrust/functional.h>
6  #include <thrust/transform.h>
7
8  int main()
9  {
10    constexpr size_t n{1'000'000};
11    thrust::host_vector<float> a(n, 2.), b(n, 3.);
12
13    thrust::device_vector<float> a_dev{a}, b_dev{b};
14
15    thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
16                      b_dev.begin(), thrust::plus<float>());
17
18    b = b_dev;
19
20    if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
21      std::cout << "Success!\n";
22    else std::cout << "Failure :(\n";
23  }
```

- Thrust is included in CUDA.
- Provides STL-like abstractions.
- Useful library functions for reduction, sorting, etc.

```cpp
 1  #include <algorithm>
 2  #include <iostream>
 3  #include <thrust/host_vector.h>
 4  #include <thrust/device_vector.h>
 5  #include <thrust/functional.h>
 6  #include <thrust/transform.h>
 7
 8  int main()
 9  {
10    constexpr size_t n{1'000'000};
11    thrust::host_vector<float> a(n, 2.), b(n, 3.);
12
13    thrust::device_vector<float> a_dev{a}, b_dev{b};
14
15    thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
16                      b_dev.begin(), thrust::plus<float>());
17
18    b = b_dev;
19
20    if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
21      std::cout << "Success!\n";
22    else std::cout << "Failure :(\n";
23  }
```

- Thrust is included in CUDA.
- Provides STL-like abstractions.
- Useful library functions for reduction, sorting, etc.

```cpp
#include <algorithm>
#include <iostream>
#include <thrust/host_vector.h>
#include <thrust/device_vector.h>
#include <thrust/functional.h>
#include <thrust/transform.h>

int main()
{
  constexpr size_t n{1'000'000};
  thrust::host_vector<float> a(n, 2.), b(n, 3.);

  thrust::device_vector<float> a_dev{a}, b_dev{b};

  thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
                    b_dev.begin(), thrust::plus<float>());

  b = b_dev;

  if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";
}
```

- Thrust is included in CUDA.
- Provides STL-like abstractions.
- Useful library functions for reduction, sorting, etc.

```cpp
1  #include <algorithm>
2  #include <iostream>
3  #include <thrust/host_vector.h>
4  #include <thrust/device_vector.h>
5  #include <thrust/functional.h>
6  #include <thrust/transform.h>
7
8  int main()
9  {
10   constexpr size_t n{1'000'000};
11   thrust::host_vector<float> a(n, 2.), b(n, 3.);
12
13   thrust::device_vector<float> a_dev{a}, b_dev{b};
14
15   thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
16                     b_dev.begin(), thrust::plus<float>());
17
18   b = b_dev;
19
20   if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
21     std::cout << "Success!\n";
22   else std::cout << "Failure :(\n";
23 }
```

- Thrust is included in CUDA.
- Provides STL-like abstractions.
- Useful library functions for reduction, sorting, etc.

```cpp
1  #include <algorithm>
2  #include <iostream>
3  #include <thrust/host_vector.h>
4  #include <thrust/device_vector.h>
5  #include <thrust/functional.h>
6  #include <thrust/transform.h>
7
8  int main()
9  {
10   constexpr size_t n{1'000'000};
11   thrust::host_vector<float> a(n, 2.), b(n, 3.);
12
13   thrust::device_vector<float> a_dev{a}, b_dev{b};
14
15   thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
16                     b_dev.begin(), thrust::plus<float>());
17
18   b = b_dev;
19
20   if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
21     std::cout << "Success!\n";
22   else std::cout << "Failure :(\n";
23  }
```

- Thrust is included in CUDA.
- Provides STL-like abstractions.
- Useful library functions for reduction, sorting, etc.

```cpp
#include <algorithm>
#include <iostream>
#include <thrust/host_vector.h>
#include <thrust/device_vector.h>
#include <thrust/functional.h>
#include <thrust/transform.h>

int main()
{
  constexpr size_t n{1'000'000};
  thrust::host_vector<float> a(n, 2.), b(n, 3.);

  thrust::device_vector<float> a_dev{a}, b_dev{b};

  thrust::transform(a_dev.begin(), a_dev.end(), b_dev.begin(),
                    b_dev.begin(), thrust::plus<float>());

  b = b_dev;

  if (std::all_of(begin(b), end(b), [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";
}
```
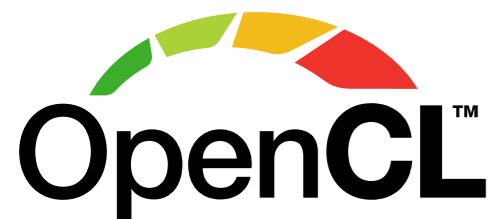
# HIP

- AMD's *clone* of CUDA (mutually intelligible).
- But it's *open source*!
- Automatic conversion tools provided.

```cpp
1   #include <algorithm>
2   #include <iostream>
3
4   __global__ void add(float *a, float *b)
5   {
6     int i = blockDim.x * blockIdx.x + threadIdx.x;
7     b[i] = a[i] + b[i];
8   }
9
10  int main()
11  {
12    constexpr size_t n{1'000'000};
13    float *a{new float[n]}, *b{new float[n]};
14    std::fill(a, a+n, 2.);
15    std::fill(b, b+n, 3.);
16
17    float *a_dev, *b_dev;
18    cudaMalloc((void**)&a_dev, n*sizeof(float));
19    cudaMalloc((void**)&b_dev, n*sizeof(float));
20    cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);
21    cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);
22
23    add<<<n/64,64>>>(a_dev, b_dev);
24
25    cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);
26
27    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
28      std::cout << "Success!\n";
29    else std::cout << "Failure :(\n";
30
31    delete[] a;
32    delete[] b;
33    cudaFree(a_dev);
34    cudaFree(b_dev);
35  }
```

```cpp
1   #include "hip/hip_runtime.h"
2   #include <algorithm>
3   #include <iostream>
4
5   __global__ void add(float *a, float *b)
6   {
7     int i = blockDim.x * blockIdx.x + threadIdx.x;
8     b[i] = a[i] + b[i];
9   }
10
11  int main()
12  {
13    constexpr size_t n{1'000'000};
14    float *a{new float[n]}, *b{new float[n]};
15    std::fill(a, a+n, 2.);
16    std::fill(b, b+n, 3.);
17
18    float *a_dev, *b_dev;
19    hipMalloc((void**)&a_dev, n*sizeof(float));
20    hipMalloc((void**)&b_dev, n*sizeof(float));
21    hipMemcpy(a_dev, a, n*sizeof(float), hipMemcpyHostToDevice);
22    hipMemcpy(b_dev, b, n*sizeof(float), hipMemcpyHostToDevice);
23
24    hipLaunchKernelGGL(add, dim3(n/64), dim3(64), 0, 0, a_dev, b_dev);
25
26    hipMemcpy(b, b_dev, n*sizeof(float), hipMemcpyDeviceToHost);
27
28    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
29      std::cout << "Success!\n";
30    else std::cout << "Failure :(\n";
31
32    delete[] a;
33    delete[] b;
34    hipFree(a_dev);
35    hipFree(b_dev);
36  }
```

```
 1  #include <algorithm>                                        1+  #include "hip/hip_runtime.h"
 2  #include <iostream>                                         2   #include <algorithm>
 3                                                              3   #include <iostream>
 4  __global__ void add(float *a, float *b)                     4
 5  {                                                           5   __global__ void add(float *a, float *b)
 6    int i = blockDim.x * blockIdx.x + threadIdx.x;            6   {
 7    b[i] = a[i] + b[i];                                       7     int i = blockDim.x * blockIdx.x + threadIdx.x;
 8  }                                                           8     b[i] = a[i] + b[i];
 9                                                              9   }
10  int main()                                                 10
11  {                                                          11  int main()
12    constexpr size_t n{1'000'000};                           12  {
13    float *a{new float[n]}, *b{new float[n]};                13    constexpr size_t n{1'000'000};
14    std::fill(a, a+n, 2.);                                   14    float *a{new float[n]}, *b{new float[n]};
15    std::fill(b, b+n, 3.);                                   15    std::fill(a, a+n, 2.);
16                                                             16    std::fill(b, b+n, 3.);
17    float *a_dev, *b_dev;                                    17
18    cudaMalloc((void**)&a_dev, n*sizeof(float));             18    float *a_dev, *b_dev;
19    cudaMalloc((void**)&b_dev, n*sizeof(float));        19+  hipMalloc((void**)&a_dev, n*sizeof(float));
20    cudaMemcpy(a_dev, a, n*sizeof(float), cudaMemcpyHostToDevice);  20+  hipMalloc((void**)&b_dev, n*sizeof(float));
21    cudaMemcpy(b_dev, b, n*sizeof(float), cudaMemcpyHostToDevice);  21+  hipMemcpy(a_dev, a, n*sizeof(float), hipMemcpyHostToDevice);
22                                                         22+  hipMemcpy(b_dev, b, n*sizeof(float), hipMemcpyHostToDevice);
23    add<<<n/64,64>>>(a_dev, b_dev);                          23
24                                                         24+  hipLaunchKernelGGL(add, dim3(n/64), dim3(64), 0, 0, a_dev, b_dev);
25    cudaMemcpy(b, b_dev, n*sizeof(float), cudaMemcpyDeviceToHost);  25
26                                                         26+  hipMemcpy(b, b_dev, n*sizeof(float), hipMemcpyDeviceToHost);
27    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))   27
28      std::cout << "Success!\n";                          28    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
29    else std::cout << "Failure :(\n";                     29      std::cout << "Success!\n";
30                                                         30    else std::cout << "Failure :(\n";
31    delete[] a;                                          31
32    delete[] b;                                          32    delete[] a;
33    cudaFree(a_dev);                                     33    delete[] b;
34    cudaFree(b_dev);                                34+  hipFree(a_dev);
35  }                                                   35+  hipFree(b_dev);
                                                        36  }
```

# OpenCL

- A *standard* rather than a piece of software.
- Depends on *vendor implementations* (and that's a mess).
- Targets *multiple* "platforms", not just GPUs.

```cpp
1  #define CL_TARGET_OPENCL_VERSION 120
2  #include <algorithm>
3  #include <CL/opencl.h>
4  #include <iostream>
5
6  const char *kernel_source =
7  "__kernel void add(__global float *a, __global float *b) \n"
8  "{                                                        \n"
9  "  int i = get_global_id(0);                              \n"
10 "  b[i] = a[i] + b[i];                                    \n"
11 "}                                                        \n";
12
13 int main()
14 {
15   constexpr size_t n{1'000'000};
16   float *a{new float[n]}, *b{new float[n]};
17   std::fill(a, a+n, 2.);
18   std::fill(b, b+n, 3.);
19
20   cl_platform_id platform;
21   clGetPlatformIDs(1, &platform, NULL);
22   cl_device_id device_id;
23   clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NUL
24   cl_int err;
25   cl_context context = clCreateContext(0, 1, &device_id, NULL, NU
26   cl_command_queue queue = clCreateCommandQueue(context, device_i
27
28   cl_mem a_dev;
29   cl_mem b_dev;
30   a_dev = clCreateBuffer(context, CL_MEM_READ_ONLY,  n*sizeof(flo
31   b_dev = clCreateBuffer(context, CL_MEM_READ_WRITE, n*sizeof(flo
32   clEnqueueWriteBuffer(queue, a_dev, CL_TRUE, 0, n*sizeof(float),
33   clEnqueueWriteBuffer(queue, b_dev, CL_TRUE, 0, n*sizeof(float),
34
35   cl_program program = clCreateProgramWithSource(context, 1, &ker
36   clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
37   cl_kernel kernel = clCreateKernel(program, "add", &err);
38   clSetKernelArg(kernel, 0, sizeof(cl_mem), &a_dev);
39   clSetKernelArg(kernel, 1, sizeof(cl_mem), &b_dev);
40   const size_t global_size{n}, local_size{64};
```

OpenCL™

- A *standard* rather than a piece of software.
- Depends on *vendor implementations* (and that's a mess).
- Targets *multiple* "platforms", not just GPUs.
- Has tonnes of *boilerplate* and a steep learning curve.

```cpp
1  #define CL_TARGET_OPENCL_VERSION 120
2  #include <algorithm>
3  #include <CL/opencl.h>
4  #include <iostream>
5
6  const char *kernel_source =
7  "__kernel void add(__global float *a, __global float *b) \n"
8  "{                                                        \n"
9  "  int i = get_global_id(0);                              \n"
10 "  b[i] = a[i] + b[i];                                    \n"
11 "}                                                        \n";
12
13 int main()
14 {
15   constexpr size_t n{1'000'000};
16   float *a{new float[n]}, *b{new float[n]};
17   std::fill(a, a+n, 2.);
18   std::fill(b, b+n, 3.);
19
20   cl_platform_id platform;
21   clGetPlatformIDs(1, &platform, NULL);
22   cl_device_id device_id;
23   clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NUL
24   cl_int err;
25   cl_context context = clCreateContext(0, 1, &device_id, NULL, NU
26   cl_command_queue queue = clCreateCommandQueue(context, device_i
27
28   cl_mem a_dev;
29   cl_mem b_dev;
30   a_dev = clCreateBuffer(context, CL_MEM_READ_ONLY,  n*sizeof(flo
31   b_dev = clCreateBuffer(context, CL_MEM_READ_WRITE, n*sizeof(flo
32   clEnqueueWriteBuffer(queue, a_dev, CL_TRUE, 0, n*sizeof(float),
33   clEnqueueWriteBuffer(queue, b_dev, CL_TRUE, 0, n*sizeof(float),
34
35   cl_program program = clCreateProgramWithSource(context, 1, &ker
36   clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
37   cl_kernel kernel = clCreateKernel(program, "add", &err);
38   clSetKernelArg(kernel, 0, sizeof(cl_mem), &a_dev);
39   clSetKernelArg(kernel, 1, sizeof(cl_mem), &b_dev);
40   const size_t global_size{n}, local_size{64};
```

- A *standard* rather than a piece of software.
- Depends on *vendor implementations* (and that's a mess).
- Targets *multiple* "platforms", not just GPUs.
- Has tonnes of *boilerplate* and a steep learning curve.
- Strongly inspired by *CUDA*.

```cpp
#define CL_TARGET_OPENCL_VERSION 120
#include <algorithm>
#include <CL/opencl.h>
#include <iostream>

const char *kernel_source =
"__kernel void add(__global float *a, __global float *b) \n"
"{                                                        \n"
"  int i = get_global_id(0);                              \n"
"  b[i] = a[i] + b[i];                                    \n"
"}                                                        \n";

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  cl_platform_id platform;
  clGetPlatformIDs(1, &platform, NULL);
  cl_device_id device_id;
  clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NUL
  cl_int err;
  cl_context context = clCreateContext(0, 1, &device_id, NULL, NU
  cl_command_queue queue = clCreateCommandQueue(context, device_i

  cl_mem a_dev;
  cl_mem b_dev;
  a_dev = clCreateBuffer(context, CL_MEM_READ_ONLY,  n*sizeof(flo
  b_dev = clCreateBuffer(context, CL_MEM_READ_WRITE, n*sizeof(flo
  clEnqueueWriteBuffer(queue, a_dev, CL_TRUE, 0, n*sizeof(float),
  clEnqueueWriteBuffer(queue, b_dev, CL_TRUE, 0, n*sizeof(float),

  cl_program program = clCreateProgramWithSource(context, 1, &ker
  clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
  cl_kernel kernel = clCreateKernel(program, "add", &err);
  clSetKernelArg(kernel, 0, sizeof(cl_mem), &a_dev);
  clSetKernelArg(kernel, 1, sizeof(cl_mem), &b_dev);
  const size_t global_size{n}, local_size{64};
```

# OpenCL

- A *standard* rather than a piece of software.
- Depends on *vendor implementations* (and that's a mess).
- Targets *multiple* "platforms", not just GPUs.
- Has tonnes of *boilerplate* and a steep learning curve.
- Strongly inspired by *CUDA*.
- Kernels compiled *just in time* from strings.

```cpp
 1  #define CL_TARGET_OPENCL_VERSION 120
 2  #include <algorithm>
 3  #include <CL/opencl.h>
 4  #include <iostream>
 5
 6  const char *kernel_source =
 7  "__kernel void add(__global float *a, __global float *b) \n"
 8  "{                                                        \n"
 9  "  int i = get_global_id(0);                              \n"
10  "  b[i] = a[i] + b[i];                                    \n"
11  "}                                                        \n";
12
13  int main()
14  {
15    constexpr size_t n{1'000'000};
16    float *a{new float[n]}, *b{new float[n]};
17    std::fill(a, a+n, 2.);
18    std::fill(b, b+n, 3.);
19
20    cl_platform_id platform;
21    clGetPlatformIDs(1, &platform, NULL);
22    cl_device_id device_id;
23    clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NUL
24    cl_int err;
25    cl_context context = clCreateContext(0, 1, &device_id, NULL, NU
26    cl_command_queue queue = clCreateCommandQueue(context, device_i
27
28    cl_mem a_dev;
29    cl_mem b_dev;
30    a_dev = clCreateBuffer(context, CL_MEM_READ_ONLY,  n*sizeof(flo
31    b_dev = clCreateBuffer(context, CL_MEM_READ_WRITE, n*sizeof(flo
32    clEnqueueWriteBuffer(queue, a_dev, CL_TRUE, 0, n*sizeof(float),
33    clEnqueueWriteBuffer(queue, b_dev, CL_TRUE, 0, n*sizeof(float),
34
35    cl_program program = clCreateProgramWithSource(context, 1, &ker
36    clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
37    cl_kernel kernel = clCreateKernel(program, "add", &err);
38    clSetKernelArg(kernel, 0, sizeof(cl_mem), &a_dev);
39    clSetKernelArg(kernel, 1, sizeof(cl_mem), &b_dev);
40    const size_t global_size{n}, local_size{64};
```

**OpenCL**

- A *standard* rather than a piece of software.
- Depends on *vendor implementations* (and that's a mess).
- Targets *multiple* "platforms", not just GPUs.
- Has tonnes of *boilerplate* and a steep learning curve.
- Strongly inspired by *CUDA*.
- Kernels compiled *just in time* from strings.
- Before HIP, was the only way to program *AMD* GPUs.

```cpp
 1  #define CL_TARGET_OPENCL_VERSION 120
 2  #include <algorithm>
 3  #include <CL/opencl.h>
 4  #include <iostream>
 5
 6  const char *kernel_source =
 7  "__kernel void add(__global float *a, __global float *b) \n"
 8  "{                                                        \n"
 9  "  int i = get_global_id(0);                              \n"
10  "  b[i] = a[i] + b[i];                                    \n"
11  "}                                                        \n";
12
13  int main()
14  {
15    constexpr size_t n{1'000'000};
16    float *a{new float[n]}, *b{new float[n]};
17    std::fill(a, a+n, 2.);
18    std::fill(b, b+n, 3.);
19
20    cl_platform_id platform;
21    clGetPlatformIDs(1, &platform, NULL);
22    cl_device_id device_id;
23    clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NUL
24    cl_int err;
25    cl_context context = clCreateContext(0, 1, &device_id, NULL, NU
26    cl_command_queue queue = clCreateCommandQueue(context, device_i
27
28    cl_mem a_dev;
29    cl_mem b_dev;
30    a_dev = clCreateBuffer(context, CL_MEM_READ_ONLY,  n*sizeof(flo
31    b_dev = clCreateBuffer(context, CL_MEM_READ_WRITE, n*sizeof(flo
32    clEnqueueWriteBuffer(queue, a_dev, CL_TRUE, 0, n*sizeof(float),
33    clEnqueueWriteBuffer(queue, b_dev, CL_TRUE, 0, n*sizeof(float),
34
35    cl_program program = clCreateProgramWithSource(context, 1, &ker
36    clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
37    cl_kernel kernel = clCreateKernel(program, "add", &err);
38    clSetKernelArg(kernel, 0, sizeof(cl_mem), &a_dev);
39    clSetKernelArg(kernel, 1, sizeof(cl_mem), &b_dev);
40    const size_t global_size{n}, local_size{64};
```

# OpenCL

- A *standard* rather than a piece of software.
- Depends on *vendor implementations* (and that's a mess).
- Targets *multiple* "platforms", not just GPUs.
- Has tonnes of *boilerplate* and a steep learning curve.
- Strongly inspired by *CUDA*.
- Kernels compiled *just in time* from strings.
- Before HIP, was the only way to program *AMD* GPUs.
- Nvidia's OpenCL implementation:
  - is historically slow compared to the CUDA runtime.
  - is historically behind the standard.
  - doesn't even work on PPC64LE (Mist).

```cpp
 1  #define CL_TARGET_OPENCL_VERSION 120
 2  #include <algorithm>
 3  #include <CL/opencl.h>
 4  #include <iostream>
 5
 6  const char *kernel_source =
 7  "__kernel void add(__global float *a, __global float *b) \n"
 8  "{                                                        \n"
 9  "  int i = get_global_id(0);                              \n"
10  "  b[i] = a[i] + b[i];                                    \n"
11  "}                                                        \n";
12
13  int main()
14  {
15    constexpr size_t n{1'000'000};
16    float *a{new float[n]}, *b{new float[n]};
17    std::fill(a, a+n, 2.);
18    std::fill(b, b+n, 3.);
19
20    cl_platform_id platform;
21    clGetPlatformIDs(1, &platform, NULL);
22    cl_device_id device_id;
23    clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NUL
24    cl_int err;
25    cl_context context = clCreateContext(0, 1, &device_id, NULL, NU
26    cl_command_queue queue = clCreateCommandQueue(context, device_i
27
28    cl_mem a_dev;
29    cl_mem b_dev;
30    a_dev = clCreateBuffer(context, CL_MEM_READ_ONLY,  n*sizeof(flo
31    b_dev = clCreateBuffer(context, CL_MEM_READ_WRITE, n*sizeof(flo
32    clEnqueueWriteBuffer(queue, a_dev, CL_TRUE, 0, n*sizeof(float),
33    clEnqueueWriteBuffer(queue, b_dev, CL_TRUE, 0, n*sizeof(float),
34
35    cl_program program = clCreateProgramWithSource(context, 1, &ker
36    clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
37    cl_kernel kernel = clCreateKernel(program, "add", &err);
38    clSetKernelArg(kernel, 0, sizeof(cl_mem), &a_dev);
39    clSetKernelArg(kernel, 1, sizeof(cl_mem), &b_dev);
40    const size_t global_size{n}, local_size{64};
```

# OpenCL

- A *standard* rather than a piece of software.
- Depends on *vendor implementations* (and that's a mess).
- Targets *multiple* "platforms", not just GPUs.
- Has tonnes of *boilerplate* and a steep learning curve.
- Strongly inspired by *CUDA*.
- Kernels compiled *just in time* from strings.
- Before HIP, was the only way to program *AMD* GPUs.
- Nvidia's OpenCL implementation:
  - is historically slow compared to the CUDA runtime.
  - is historically behind the standard.
  - doesn't even work on PPC64LE (Mist).
- <u>Not recommended</u> for beginners (or anyone) in 2023.

```cpp
 1  #define CL_TARGET_OPENCL_VERSION 120
 2  #include <algorithm>
 3  #include <CL/opencl.h>
 4  #include <iostream>
 5
 6  const char *kernel_source =
 7  "__kernel void add(__global float *a, __global float *b) \n"
 8  "{                                                        \n"
 9  "  int i = get_global_id(0);                              \n"
10  "  b[i] = a[i] + b[i];                                    \n"
11  "}                                                        \n";
12
13  int main()
14  {
15      constexpr size_t n{1'000'000};
16      float *a{new float[n]}, *b{new float[n]};
17      std::fill(a, a+n, 2.);
18      std::fill(b, b+n, 3.);
19
20      cl_platform_id platform;
21      clGetPlatformIDs(1, &platform, NULL);
22      cl_device_id device_id;
23      clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NUL
24      cl_int err;
25      cl_context context = clCreateContext(0, 1, &device_id, NULL, NU
26      cl_command_queue queue = clCreateCommandQueue(context, device_i
27
28      cl_mem a_dev;
29      cl_mem b_dev;
30      a_dev = clCreateBuffer(context, CL_MEM_READ_ONLY,  n*sizeof(flo
31      b_dev = clCreateBuffer(context, CL_MEM_READ_WRITE, n*sizeof(flo
32      clEnqueueWriteBuffer(queue, a_dev, CL_TRUE, 0, n*sizeof(float),
33      clEnqueueWriteBuffer(queue, b_dev, CL_TRUE, 0, n*sizeof(float),
34
35      cl_program program = clCreateProgramWithSource(context, 1, &ker
36      clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
37      cl_kernel kernel = clCreateKernel(program, "add", &err);
38      clSetKernelArg(kernel, 0, sizeof(cl_mem), &a_dev);
39      clSetKernelArg(kernel, 1, sizeof(cl_mem), &b_dev);
40      const size_t global_size{n}, local_size{64};
```

```cpp
#include <algorithm>
#include <iostream>

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  #pragma acc kernels
  #pragma acc loop independent
  for (int i = 0; i < n; i++)
    b[i] = a[i] + b[i];

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
}
```

Example of OpenACC.

# OpenMP

# OpenACC
More Science, Less Programming

- *Directive*-based approaches.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  int main()
5  {
6    constexpr size_t n{1'000'000};
7    float *a{new float[n]}, *b{new float[n]};
8    std::fill(a, a+n, 2.);
9    std::fill(b, b+n, 3.);
10
11   #pragma acc kernels
12   #pragma acc loop independent
13   for (int i = 0; i < n; i++)
14     b[i] = a[i] + b[i];
15
16   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
17     std::cout << "Success!\n";
18   else std::cout << "Failure :(\n";
19
20   delete[] a;
21   delete[] b;
22 }
```

Example of OpenACC.

# OpenMP   OpenACC

More Science, Less Programming

- *Directive*-based approaches.
- "Easy" to accelerate existing CPU code.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  int main()
5  {
6    constexpr size_t n{1'000'000};
7    float *a{new float[n]}, *b{new float[n]};
8    std::fill(a, a+n, 2.);
9    std::fill(b, b+n, 3.);
10
11   #pragma acc kernels
12   #pragma acc loop independent
13   for (int i = 0; i < n; i++)
14     b[i] = a[i] + b[i];
15
16   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
17     std::cout << "Success!\n";
18   else std::cout << "Failure :(\n";
19
20   delete[] a;
21   delete[] b;
22 }
```

Example of OpenACC.

# OpenMP    # OpenACC

**More Science, Less Programming**

- *Directive*-based approaches.
- "Easy" to accelerate existing CPU code.
- Single code base for CPU and GPU.

```cpp
 1  #include <algorithm>
 2  #include <iostream>
 3
 4  int main()
 5  {
 6    constexpr size_t n{1'000'000};
 7    float *a{new float[n]}, *b{new float[n]};
 8    std::fill(a, a+n, 2.);
 9    std::fill(b, b+n, 3.);
10
11    #pragma acc kernels
12    #pragma acc loop independent
13    for (int i = 0; i < n; i++)
14      b[i] = a[i] + b[i];
15
16    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
17      std::cout << "Success!\n";
18    else std::cout << "Failure :(\n";
19
20    delete[] a;
21    delete[] b;
22  }
```

Example of OpenACC.

# OpenMP  OpenACC

**More Science, Less Programming**

- *Directive*-based approaches.
- "Easy" to accelerate existing CPU code.
- Single code base for CPU and GPU.
- Shared origin and mutually intelligible.

```cpp
 1  #include <algorithm>
 2  #include <iostream>
 3
 4  int main()
 5  {
 6    constexpr size_t n{1'000'000};
 7    float *a{new float[n]}, *b{new float[n]};
 8    std::fill(a, a+n, 2.);
 9    std::fill(b, b+n, 3.);
10
11    #pragma acc kernels
12    #pragma acc loop independent
13    for (int i = 0; i < n; i++)
14      b[i] = a[i] + b[i];
15
16    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
17      std::cout << "Success!\n";
18    else std::cout << "Failure :(\n";
19
20    delete[] a;
21    delete[] b;
22  }
```

Example of OpenACC.

# OpenMP  OpenACC

More Science, Less Programming

- *Directive*-based approaches.
- "Easy" to accelerate existing CPU code.
- Single code base for CPU and GPU.
- Shared origin and mutually intelligible.
- Kernels are generated *from loops*.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  int main()
5  {
6    constexpr size_t n{1'000'000};
7    float *a{new float[n]}, *b{new float[n]};
8    std::fill(a, a+n, 2.);
9    std::fill(b, b+n, 3.);
10
11   #pragma acc kernels
12   #pragma acc loop independent
13   for (int i = 0; i < n; i++)
14     b[i] = a[i] + b[i];
15
16   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
17     std::cout << "Success!\n";
18   else std::cout << "Failure :(\n";
19
20   delete[] a;
21   delete[] b;
22 }
```

Example of OpenACC.

# OpenMP  OpenACC

More Science, Less Programming

- *Directive*-based approaches.
- "Easy" to accelerate existing CPU code.
- Single code base for CPU and GPU.
- Shared origin and mutually intelligible.
- Kernels are generated *from loops*.
- Memory copies [kinda] automatic.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  int main()
5  {
6    constexpr size_t n{1'000'000};
7    float *a{new float[n]}, *b{new float[n]};
8    std::fill(a, a+n, 2.);
9    std::fill(b, b+n, 3.);
10
11   #pragma acc kernels
12   #pragma acc loop independent
13   for (int i = 0; i < n; i++)
14     b[i] = a[i] + b[i];
15
16   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
17     std::cout << "Success!\n";
18   else std::cout << "Failure :(\n";
19
20   delete[] a;
21   delete[] b;
22 }
```

Example of OpenACC.

# OpenMP / OpenACC
*More Science, Less Programming*

- *Directive*-based approaches.
- "Easy" to accelerate existing CPU code.
- Single code base for CPU and GPU.
- Shared origin and mutually intelligible.
- Kernels are generated *from loops*.
- Memory copies [kinda] automatic.

```cpp
1  #include <algorithm>
2  #include <iostream>
3
4  int main()
5  {
6    constexpr size_t n{1'000'000};
7    float *a{new float[n]}, *b{new float[n]};
8    std::fill(a, a+n, 2.);
9    std::fill(b, b+n, 3.);
10
11   #pragma omp target data map(to: a[:n]) map(tofrom: b[:n])
12   #pragma omp target parallel for
13   for (int i = 0; i < n; i++)
14     b[i] = a[i] + b[i];
15
16   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
17     std::cout << "Success!\n";
18   else std::cout << "Failure :(\n";
19
20   delete[] a;
21   delete[] b;
22 }
```

Example of OpenMP offloading.

```cpp
#include <algorithm>
#include <CL/sycl.hpp>
#include <iostream>
using namespace cl;

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  sycl::device device{sycl::default_selector()};
  sycl::context context{device};
  sycl::queue queue{context, device};

  float* a_dev{sycl::malloc_device<float>(n, device, context)};
  float* b_dev{sycl::malloc_device<float>(n, device, context)};

  queue.copy(a, a_dev, n);
  queue.copy(b, b_dev, n);
  queue.wait();

  queue.parallel_for(sycl::range<1>{n},
    [=](const auto i) { b_dev[i] = a_dev[i] + b_dev[i]; }
  );
  queue.wait();

  queue.copy(b_dev, b, n);
  queue.wait();

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  sycl::free(a_dev, context);
  sycl::free(b_dev, context);
}
```

# SYCL™

- Another standard developed by Khronos Group.

```cpp
1  #include <algorithm>
2  #include <CL/sycl.hpp>
3  #include <iostream>
4  using namespace cl;
5
6  int main()
7  {
8    constexpr size_t n{1'000'000};
9    float *a{new float[n]}, *b{new float[n]};
10   std::fill(a, a+n, 2.);
11   std::fill(b, b+n, 3.);
12
13   sycl::device device{sycl::default_selector()};
14   sycl::context context{device};
15   sycl::queue queue{context, device};
16
17   float* a_dev{sycl::malloc_device<float>(n, device, context)};
18   float* b_dev{sycl::malloc_device<float>(n, device, context)};
19
20   queue.copy(a, a_dev, n);
21   queue.copy(b, b_dev, n);
22   queue.wait();
23
24   queue.parallel_for(sycl::range<1>{n},
25     [=](const auto i) { b_dev[i] = a_dev[i] + b_dev[i]; }
26   );
27   queue.wait();
28
29   queue.copy(b_dev, b, n);
30   queue.wait();
31
32   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
33     std::cout << "Success!\n";
34   else std::cout << "Failure :(\n";
35
36   delete[] a;
37   delete[] b;
38   sycl::free(a_dev, context);
39   sycl::free(b_dev, context);
40 }
```

- Another standard developed by Khronos Group.
- Favoured by *Intel* for their future GPUs.

```cpp
#include <algorithm>
#include <CL/sycl.hpp>
#include <iostream>
using namespace cl;

int main()
{
  constexpr size_t n{1'000'000};
  float *a{new float[n]}, *b{new float[n]};
  std::fill(a, a+n, 2.);
  std::fill(b, b+n, 3.);

  sycl::device device{sycl::default_selector()};
  sycl::context context{device};
  sycl::queue queue{context, device};

  float* a_dev{sycl::malloc_device<float>(n, device, context)};
  float* b_dev{sycl::malloc_device<float>(n, device, context)};

  queue.copy(a, a_dev, n);
  queue.copy(b, b_dev, n);
  queue.wait();

  queue.parallel_for(sycl::range<1>{n},
    [=](const auto i) { b_dev[i] = a_dev[i] + b_dev[i]; }
  );
  queue.wait();

  queue.copy(b_dev, b, n);
  queue.wait();

  if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
    std::cout << "Success!\n";
  else std::cout << "Failure :(\n";

  delete[] a;
  delete[] b;
  sycl::free(a_dev, context);
  sycl::free(b_dev, context);
}
```

**SYCL**™

- Another standard developed by Khronos Group.
- oneAPI Favoured by *Intel* for their future GPUs.
- Based on *standard* C++17.

```cpp
1  #include <algorithm>
2  #include <CL/sycl.hpp>
3  #include <iostream>
4  using namespace cl;
5
6  int main()
7  {
8    constexpr size_t n{1'000'000};
9    float *a{new float[n]}, *b{new float[n]};
10   std::fill(a, a+n, 2.);
11   std::fill(b, b+n, 3.);
12
13   sycl::device device{sycl::default_selector()};
14   sycl::context context{device};
15   sycl::queue queue{context, device};
16
17   float* a_dev{sycl::malloc_device<float>(n, device, context)};
18   float* b_dev{sycl::malloc_device<float>(n, device, context)};
19
20   queue.copy(a, a_dev, n);
21   queue.copy(b, b_dev, n);
22   queue.wait();
23
24   queue.parallel_for(sycl::range<1>{n},
25     [=](const auto i) { b_dev[i] = a_dev[i] + b_dev[i]; }
26   );
27   queue.wait();
28
29   queue.copy(b_dev, b, n);
30   queue.wait();
31
32   if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
33     std::cout << "Success!\n";
34   else std::cout << "Failure :(\n";
35
36   delete[] a;
37   delete[] b;
38   sycl::free(a_dev, context);
39   sycl::free(b_dev, context);
40 }
```

# SYCL™

- Another standard developed by Khronos Group.
- 🔬 **oneAPI** Favoured by *Intel* for their future GPUs.
- Based on *standard* C++17.
- *DPC++* and *OpenSYCL* support multiple backends.

```cpp
 1  #include <algorithm>
 2  #include <CL/sycl.hpp>
 3  #include <iostream>
 4  using namespace cl;
 5
 6  int main()
 7  {
 8    constexpr size_t n{1'000'000};
 9    float *a{new float[n]}, *b{new float[n]};
10    std::fill(a, a+n, 2.);
11    std::fill(b, b+n, 3.);
12
13    sycl::device device{sycl::default_selector()};
14    sycl::context context{device};
15    sycl::queue queue{context, device};
16
17    float* a_dev{sycl::malloc_device<float>(n, device, context)};
18    float* b_dev{sycl::malloc_device<float>(n, device, context)};
19
20    queue.copy(a, a_dev, n);
21    queue.copy(b, b_dev, n);
22    queue.wait();
23
24    queue.parallel_for(sycl::range<1>{n},
25      [=](const auto i) { b_dev[i] = a_dev[i] + b_dev[i]; }
26    );
27    queue.wait();
28
29    queue.copy(b_dev, b, n);
30    queue.wait();
31
32    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
33      std::cout << "Success!\n";
34    else std::cout << "Failure :(\n";
35
36    delete[] a;
37    delete[] b;
38    sycl::free(a_dev, context);
39    sycl::free(b_dev, context);
40  }
```

# SYCL™

- Another standard developed by Khronos Group.
- **oneAPI** Favoured by *Intel* for their future GPUs.
- Based on *standard* C++17.
- *DPC++* and *OpenSYCL* support multiple backends.
- *Still niche* as of 2023.

```cpp
 1  #include <algorithm>
 2  #include <CL/sycl.hpp>
 3  #include <iostream>
 4  using namespace cl;
 5
 6  int main()
 7  {
 8    constexpr size_t n{1'000'000};
 9    float *a{new float[n]}, *b{new float[n]};
10    std::fill(a, a+n, 2.);
11    std::fill(b, b+n, 3.);
12
13    sycl::device device{sycl::default_selector()};
14    sycl::context context{device};
15    sycl::queue queue{context, device};
16
17    float* a_dev{sycl::malloc_device<float>(n, device, context)};
18    float* b_dev{sycl::malloc_device<float>(n, device, context)};
19
20    queue.copy(a, a_dev, n);
21    queue.copy(b, b_dev, n);
22    queue.wait();
23
24    queue.parallel_for(sycl::range<1>{n},
25      [=](const auto i) { b_dev[i] = a_dev[i] + b_dev[i]; }
26    );
27    queue.wait();
28
29    queue.copy(b_dev, b, n);
30    queue.wait();
31
32    if (std::all_of(b, b+n, [](const auto x){ return x == 5.; }))
33      std::cout << "Success!\n";
34    else std::cout << "Failure :(\n";
35
36    delete[] a;
37    delete[] b;
38    sycl::free(a_dev, context);
39    sycl::free(b_dev, context);
40  }
```

```python
import numpy as np
from numba import vectorize

@vectorize('float32(float32, float32)', target='cuda')
def add(a, b):
    return a + b

if __name__ == '__main__':
    n = 1_000_000
    a = np.ones(n, dtype=np.float32)*2
    b = np.ones(n, dtype=np.float32)*3

    b = add(a, b)

    if all([x == 5. for x in b]):
        print('Success!')
    else: print('Failure :(')
```

# Numba

- *Just-in-time* (JIT) compiler for Python.

```python
import numpy as np
from numba import vectorize

@vectorize('float32(float32, float32)', target='cuda')
def add(a, b):
  return a + b

if __name__ == '__main__':
  n = 1_000_000
  a = np.ones(n, dtype=np.float32)*2
  b = np.ones(n, dtype=np.float32)*3

  b = add(a, b)

  if all([x == 5. for x in b]):
    print('Success!')
  else: print('Failure :(')
```

# Numba

- *Just-in-time* (JIT) compiler for Python.
- Supports Nvidia ~~and AMD~~.

```python
import numpy as np
from numba import vectorize

@vectorize('float32(float32, float32)', target='cuda')
def add(a, b):
    return a + b

if __name__ == '__main__':
    n = 1_000_000
    a = np.ones(n, dtype=np.float32)*2
    b = np.ones(n, dtype=np.float32)*3

    b = add(a, b)

    if all([x == 5. for x in b]):
        print('Success!')
    else: print('Failure :(')
```

# Numba

- *Just-in-time* (JIT) compiler for Python.
- Supports Nvidia ~~and AMD~~.
- Kernels, ufuncs, and reductions.

```python
import numpy as np
from numba import vectorize

@vectorize('float32(float32, float32)', target='cuda')
def add(a, b):
  return a + b

if __name__ == '__main__':
  n = 1_000_000
  a = np.ones(n, dtype=np.float32)*2
  b = np.ones(n, dtype=np.float32)*3

  b = add(a, b)

  if all([x == 5. for x in b]):
    print('Success!')
  else: print('Failure :(')
```

# Numba

- *Just-in-time* (JIT) compiler for Python.
- Supports Nvidia ~~and AMD~~.
- Kernels, ufuncs, and reductions.
- Device code is a *restricted subset* of Python.

```python
import numpy as np
from numba import vectorize

@vectorize('float32(float32, float32)', target='cuda')
def add(a, b):
    return a + b

if __name__ == '__main__':
    n = 1_000_000
    a = np.ones(n, dtype=np.float32)*2
    b = np.ones(n, dtype=np.float32)*3

    b = add(a, b)

    if all([x == 5. for x in b]):
        print('Success!')
    else: print('Failure :(')
```

# Numba

- *Just-in-time* (JIT) compiler for Python.
- Supports Nvidia ~~and AMD~~.
- Kernels, ufuncs, and reductions.
- Device code is a *restricted subset* of Python.
- Strongly associated with *NumPy*.

```python
import numpy as np
from numba import vectorize

@vectorize('float32(float32, float32)', target='cuda')
def add(a, b):
  return a + b

if __name__ == '__main__':
  n = 1_000_000
  a = np.ones(n, dtype=np.float32)*2
  b = np.ones(n, dtype=np.float32)*3

  b = add(a, b)

  if all([x == 5. for x in b]):
    print('Success!')
  else: print('Failure :(')
```

# EXERCISE

# HOMEWORK EXERCISE

## 2D DIFFUSION (HEAT) EQUATION

$$\frac{\partial T}{\partial t} = D \left( \frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} \right)$$

Implement an explicit *finite difference*, *time marching* solution using a GPU programming framework.

# HOMEWORK EXERCISE

## 2D DIFFUSION (HEAT) EQUATION

$$\frac{\partial T}{\partial t} = D\left(\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2}\right)$$

Implement an explicit *finite difference*, *time marching* solution using a GPU programming framework.

Assumptions:

- Domain is a <u>square</u>.
- Arbitrary <u>initial</u> conditions: $T_0(x, y)$.
- Simple <u>boundary</u> conditions: $T = 0$ on edges.

▶ 0:00 / 0:10 🔊 ⛶ ⋮

# FINITE DIFFERENCE METHOD

Discretize the domain: $(x_j, y_i) = (j\Delta x, i\Delta y)$ for integer $j$ and $i$.

# FINITE DIFFERENCE METHOD

Discretize the domain: $(x_j, y_i) = (j\Delta x, i\Delta y)$ for integer $j$ and $i$.

Discrete first and second derivatives of some function $f(x, \ldots)$:

$$\left.\frac{\partial f}{\partial x}\right|_{x_j} \approx \frac{f(x_{j+1}) - f(x_j)}{\Delta x}$$

$$\left.\frac{\partial^2 f}{\partial x^2}\right|_{x_j} \approx \frac{f(x_{j+1}) - 2f(x_j) + f(x_{j-1})}{\Delta x^2}$$

# FINITE DIFFERENCE METHOD

Discretize the domain: $(x_j, y_i) = (j\Delta x, i\Delta y)$ for integer $j$ and $i$.

Discrete first and second derivatives of some function $f(x, \ldots)$:

$$\left.\frac{\partial f}{\partial x}\right|_{x_j} \approx \frac{f(x_{j+1}) - f(x_j)}{\Delta x}$$

$$\left.\frac{\partial^2 f}{\partial x^2}\right|_{x_j} \approx \frac{f(x_{j+1}) - 2f(x_j) + f(x_{j-1})}{\Delta x^2}$$

A single *step* of the diffusion equation ($t_k \rightarrow t_{k+1} \equiv t_k + \Delta t$, also assuming $\Delta y = \Delta x$):

$$T(x_j, y_i, t_{k+1}) = T(x_j, y_i, t_k) + \frac{D\Delta t}{\Delta x^2}[\quad T(x_{j+1}, y_i, t_k) + T(x_{j-1}, y_i, t_k) +$$

$$T(x_j, y_{i+1}, t_k) + T(x_j, y_{i-1}, t_k) -$$

$$4T(x_j, y_i, t_k) \quad ]$$

# LAPLACIAN AS A 5-POINT STENCIL

$$\nabla^2 f(x, y) \equiv \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

$$\approx f(x_{j+1}, y_i) + f(x_{j-1}, y_i)$$
$$+ f(x_j, y_{i+1}) + f(x_j, y_{i-1})$$
$$- 4f(x_j, y_i)$$

```
laplacian[i][j] =   f[ i ][j+1] + f[ i ][j-1]
              +  f[i+1][ j ] + f[i-1][ j ]
              -4*f[ i ][ j ]
```

# HOME EXERCISE INSTRUCTIONS

Submit your attempt by 2023 March 11 00:00.

- Serial CPU-based solutions are provided in Python and C++.
  You can start by modifying the one in your language, but don't have to.
- You need to identify the *bottleneck* and accelerate it using the GPU.
- There is more than one right answer.
- Bonus (1): the smaller $\Delta x$, the more *accurate* and *computationally heavy* the solution. Plot the timing for your solution and of the serial CPU-based solution (and possibly improved CPU-based solutions) as a function of $\Delta x$.

Bonus 2 & 3 are beyond the scope of this workshop:

- Bonus (2): decompose the domain and solve the problem with multiple GPUs on the same node.
- Bonus (3): use a distributed memory library to deploy your solution on multiple nodes.

**Hint:** for a single node you could use `multiprocessing` in Python and `thread` or *OpenMP* in C++.
For multiple nodes you could use `mpi4py` (Python) or *MPI* (C++).

# CLASS EXERCISE

1. Problem overview.

2. A naïve solution in Python.

3. Successively improving the solution.

4. A GPU-accelerated solution with Numba.
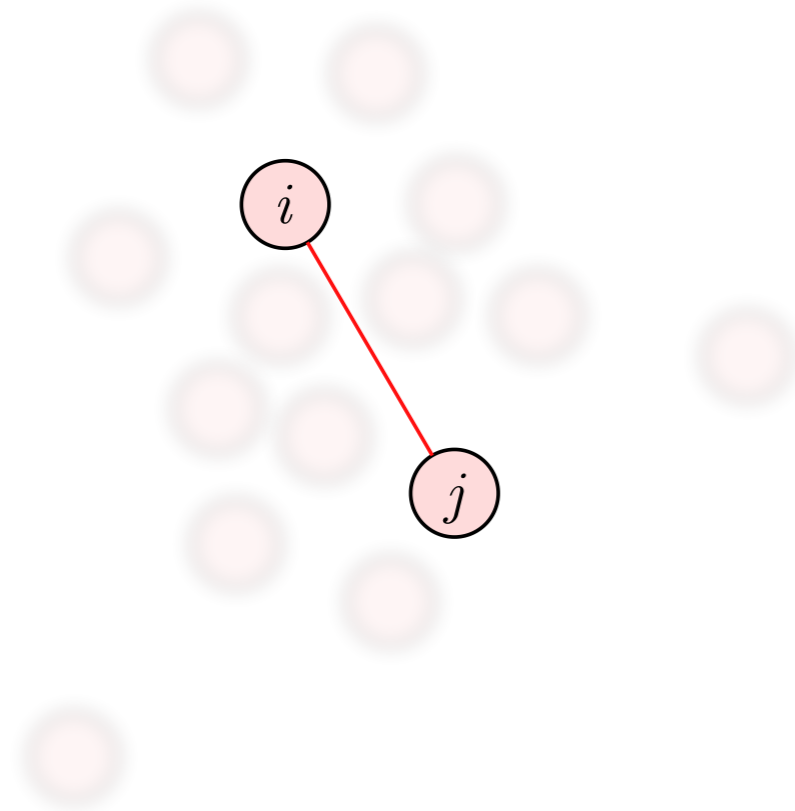
5. Comparing with a professional *N*-body library.

# PROBLEM OVERVIEW

Given a system of *N* particles, calculate the *gravitational potential* on each one.

# PROBLEM OVERVIEW

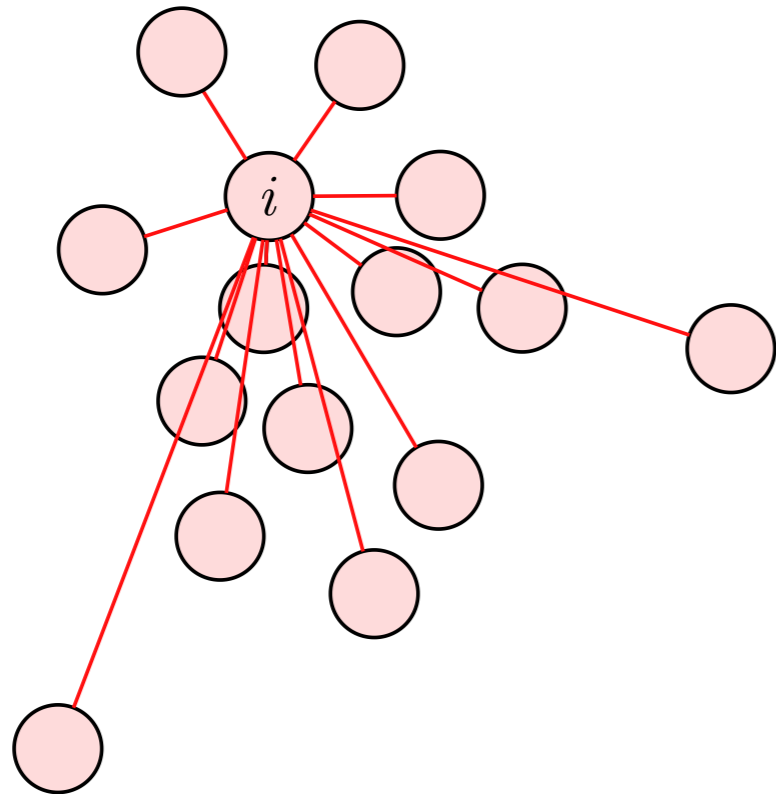Given a system of *N* particles, calculate the *gravitational potential* on each one.

Two-body potential

$$\Phi_i = -\frac{Gm_j}{r_{ij}}$$

# PROBLEM OVERVIEW

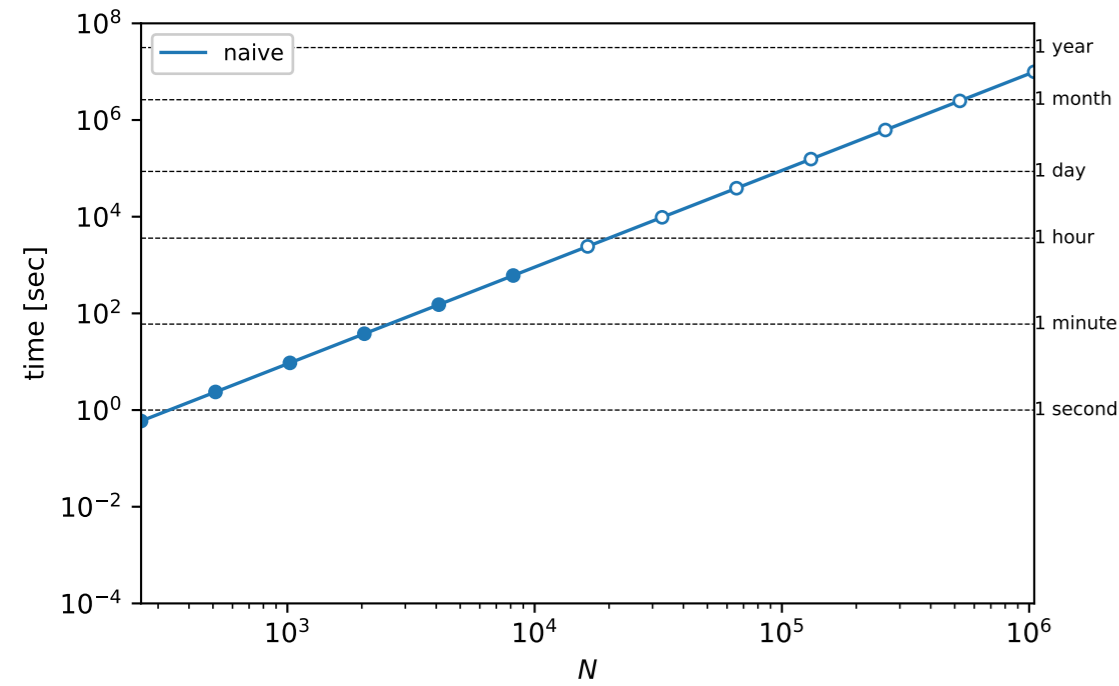Given a system of $N$ particles, calculate the *gravitational potential* on each one.



Many-body potential

$$\Phi_i = -\sum_{i \neq j} \frac{Gm_j}{r_{ij}}$$

# PROBLEM OVERVIEW

Given a system of *N* particles, calculate the *gravitational potential* on each one.

Many-body potential



$$\Phi_i = -\sum_{i \neq j} \frac{Gm_j}{r_{ij}}$$

- Calculate $\Phi_i$ *for every* $i$.
- The number of pairs is $N(N-1)/2$.
- The complexity is $\mathcal{O}(N^2)$.
- Note on alternative algorithms.
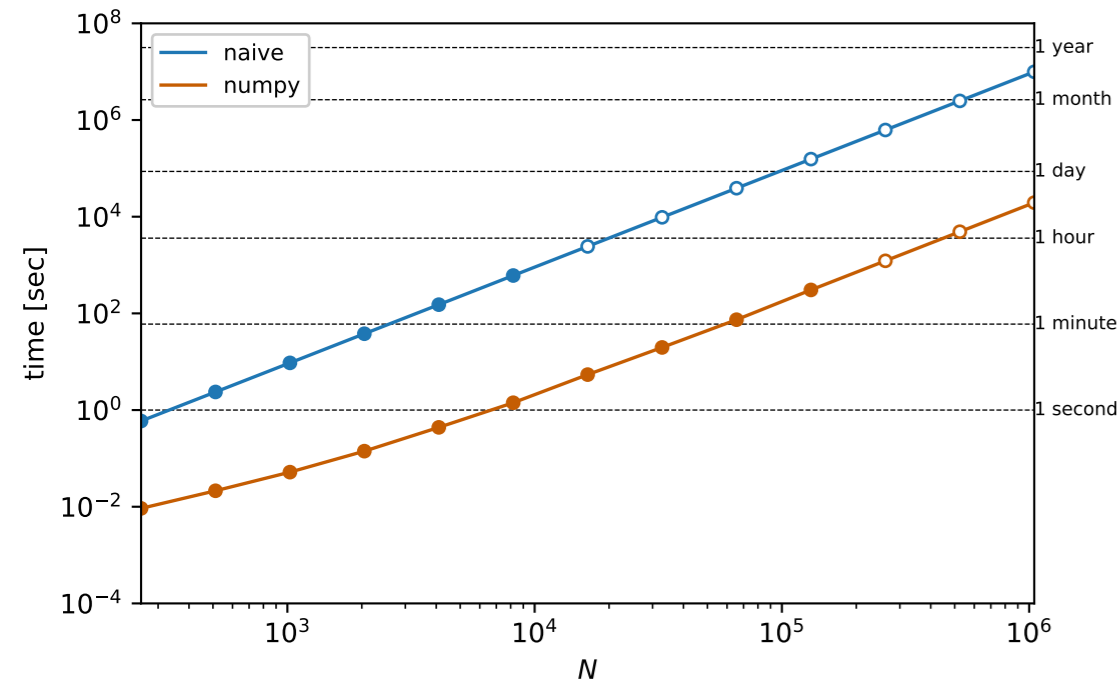- We'll assume $G = 1$ and $m_j = 1/N$.

# NAÏVE SOLUTION



| Naïve | 10 000 000 | sec | (~4 months) |
|-------|------------|-----|-------------|

```python
import numpy as np

def calculate_potential(position : np.ndarray) -> np.ndarray:
    N = len(position)
    mass = 1 / N
    potential = np.empty(N)
    for i in range(N):
        potential[i] = 0
        for j in range(N):
            if j == i: continue
            dx, dy, dz = position[i,:] - position[j,:]
            r = np.sqrt(dx**2 + dy**2 + dz**2)
            potential[i] += - mass / r
    return potential
```
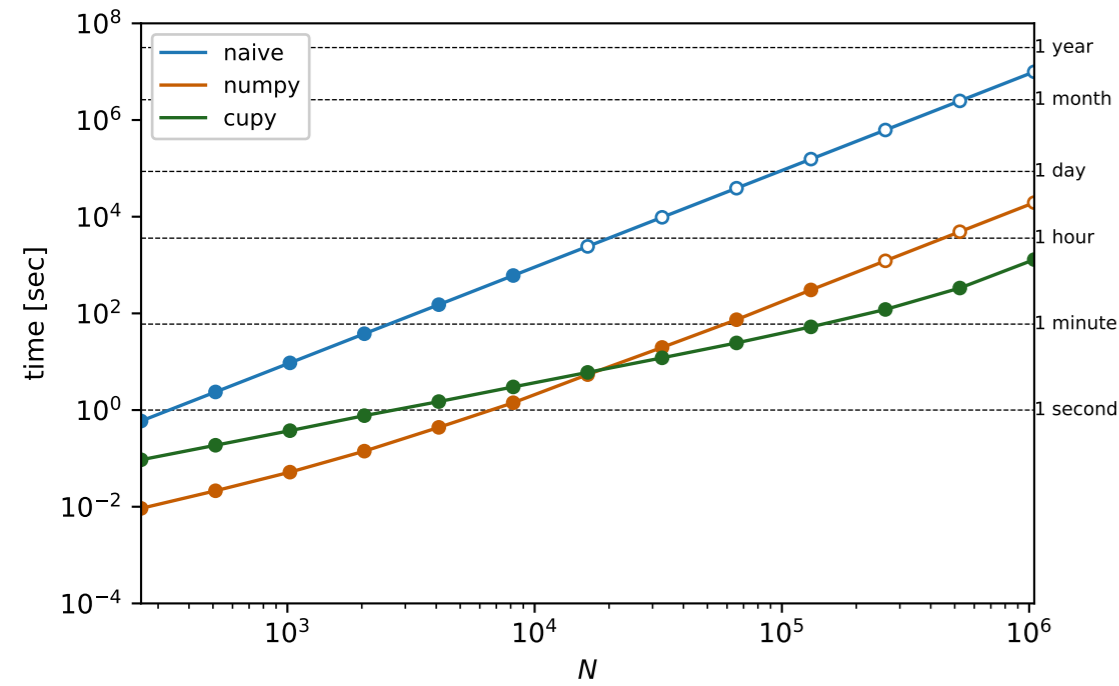
# USING NUMPY



```python
import numpy as np

def calculate_potential(position : np.ndarray) -> np.ndarray:
    N = len(position)
    mass = 1 / N
    potential = np.empty(N)
    for i in range(N):
        dx, dy, dz = (position[i,:] - position).T
        r = np.sqrt(dx**2 + dy**2 + dz**2)
        r[i] = np.inf
        potential[i] = np.sum(- mass / r)
    return potential
```
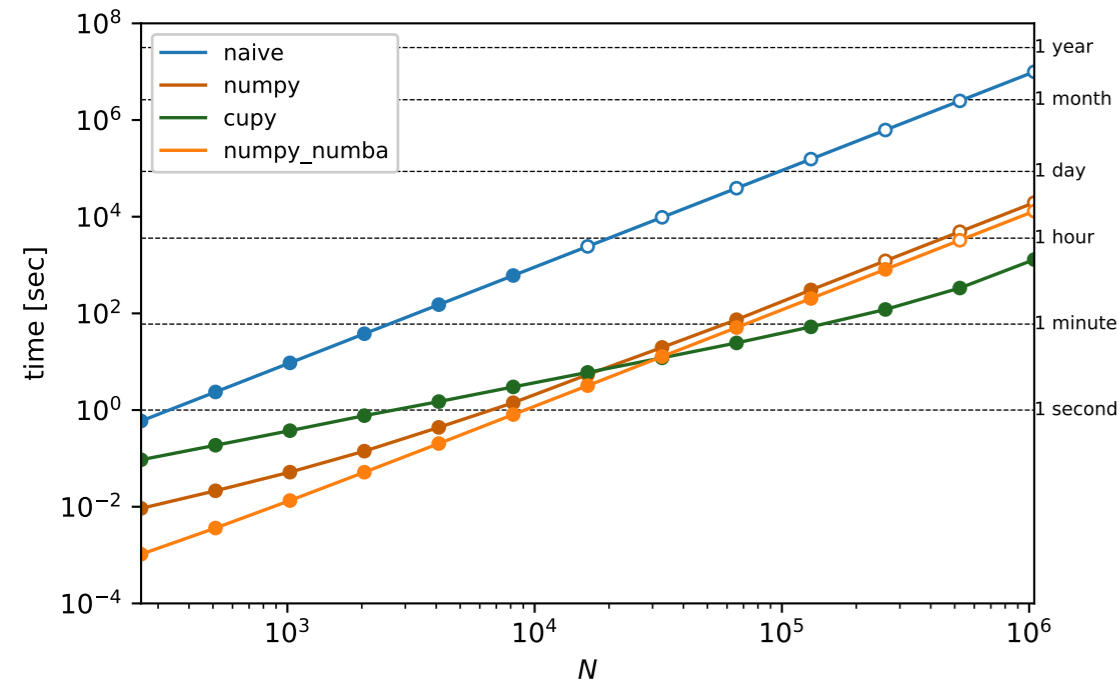
| | | | |
|---|---|---|---|
| Naïve | 10 000 000 | sec | (~4 months) |
| NumPy | 20 000 | sec | (~5 hours) |

# REPLACING NUMPY WITH CUPY



| Naïve | 10 000 000 | sec | (~4 months) |
|---|---|---|---|
| NumPy | 20 000 | sec | (~5 hours) |
| CuPy | 1 300 | sec | (~20 minutes) |

```python
1  import numpy as np, cupy as cp
2
3  def calculate_potential(position : np.ndarray) -> np.ndarray:
4      position = cp.array(position)
5      N = len(position)
6      mass = 1 / N
7      potential = np.empty(N)
8      for i in range(N):
9          dx, dy, dz = (position[i,:] - position).T
10         r = cp.sqrt(dx**2 + dy**2 + dz**2)
11         r[i] = cp.inf
12         potential[i]  = cp.sum(- mass / r)
13     return potential
```

# NUMPY + NUMBA



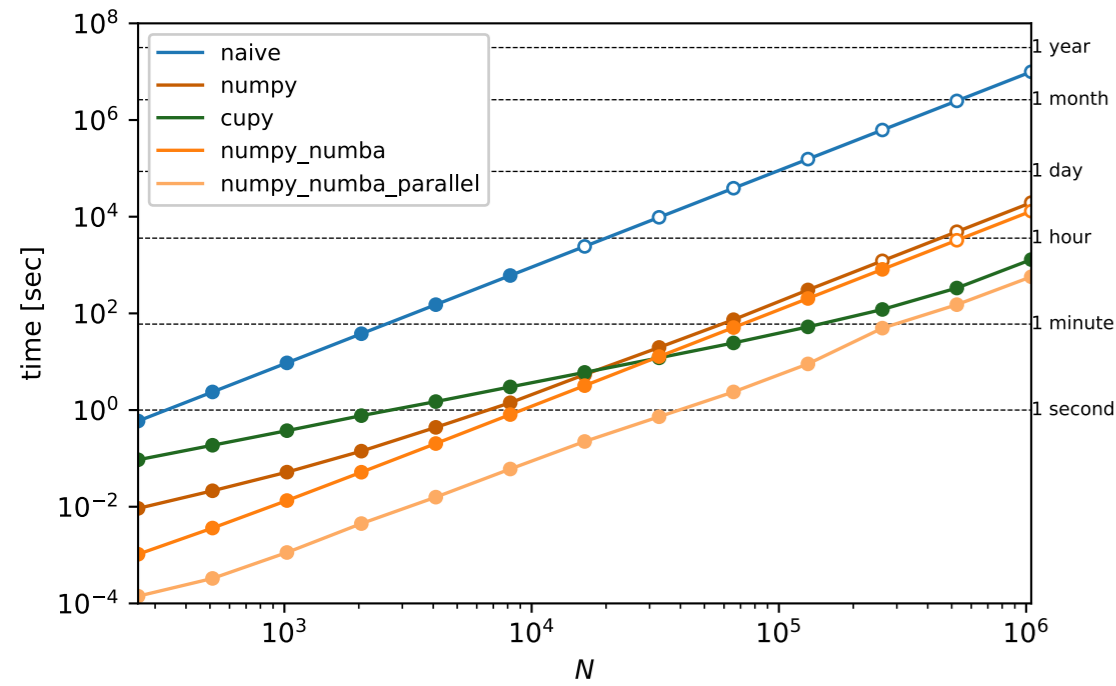| Naïve | 10 000 000 | sec | (~4 months) |
|---|---|---|---|
| NumPy | 20 000 | sec | (~5 hours) |
| CuPy | 1 300 | sec | (~20 minutes) |
| NumPy + Numba | 13 000 | sec | (~4 hours) |

```python
1  import numpy as np
2  import numba
3
4  @numba.njit(numba.float64[:](numba.float64[:,:]))
5  def calculate_potential(position : np.ndarray) -> np.ndarray:
6      N = len(position)
7      mass = 1 / N
8      potential = np.empty(N)
9      for i in range(N):
10         dx, dy, dz = (position[i,:] - position).T
11         r = np.sqrt(dx**2 + dy**2 + dz**2)
12         r[i] = np.inf
13         potential[i] = np.sum(- mass / r)
14     return potential
```
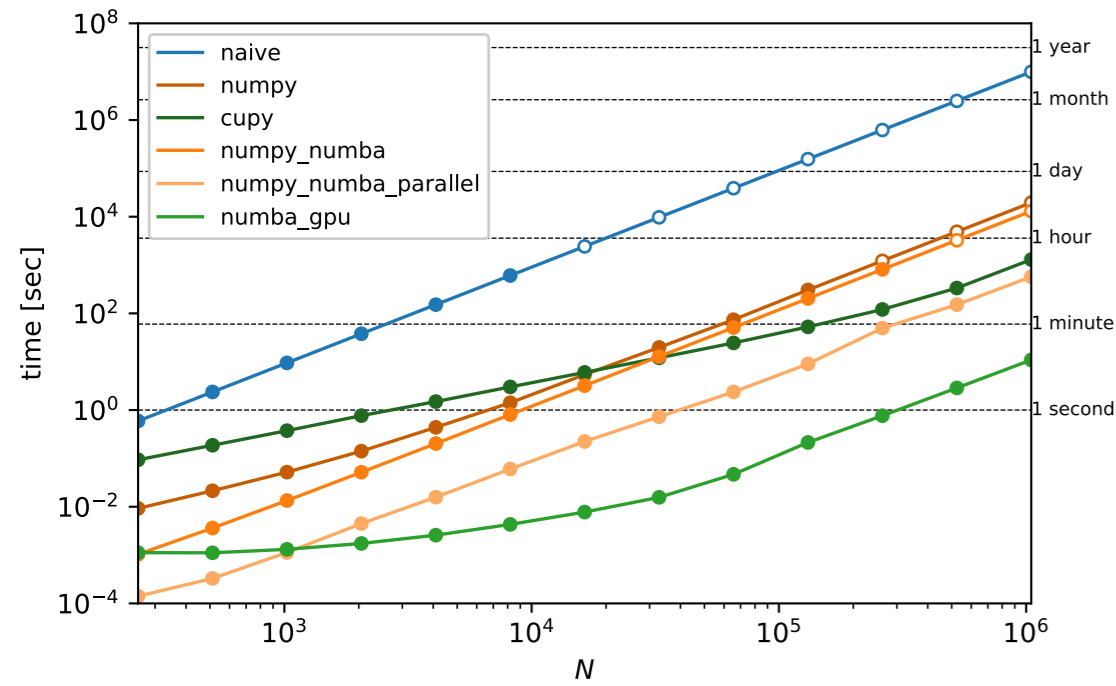
# NUMPY + NUMBA (PARALLEL)



| Naïve | 10 000 000 | sec | (~4 months) |
|---|---|---|---|
| NumPy | 20 000 | sec | (~5 hours) |
| CuPy | 1 300 | sec | (~20 minutes) |
| NumPy + Numba | 13 000 | sec | (~4 hours) |
| NumPy + Numba (parallel) | 560 | sec | (~9 minutes) |

```python
1  import numpy as np
2  import numba
3
4  @numba.njit(numba.float64[:](numba.float64[:,:]), parallel=True)
5  def calculate_potential(position : np.ndarray) -> np.ndarray:
6      N = len(position)
7      mass = 1 / N
8      potential = np.empty(N)
9      for i in numba.prange(N):
10         dx, dy, dz = (position[i,:] - position).T
11         r = np.sqrt(dx**2 + dy**2 + dz**2)
12         r[i] = np.inf
13         potential[i] = np.sum(- mass / r)
14     return potential
```
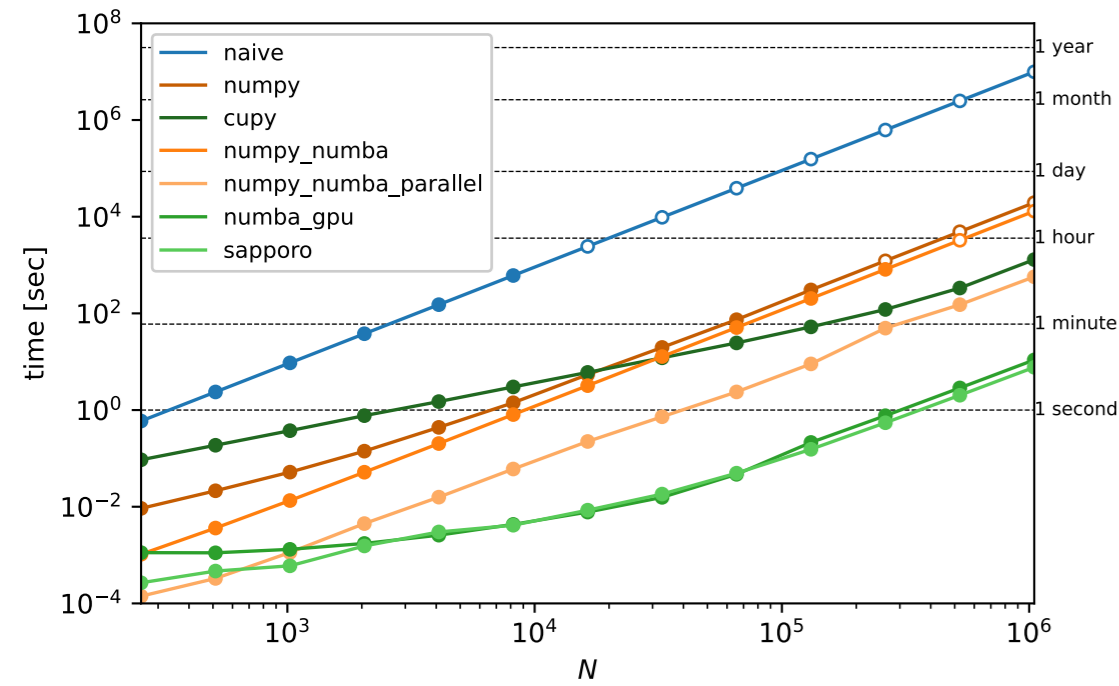
# NUMBA ON THE GPU



| | | | |
|---|---|---|---|
| Naïve | 10 000 000 | sec | (~4 months) |
| NumPy | 20 000 | sec | (~5 hours) |
| CuPy | 1 300 | sec | (~20 minutes) |
| NumPy + Numba | 13 000 | sec | (~4 hours) |
| NumPy + Numba (parallel) | 560 | sec | (~9 minutes) |
| Numba (GPU) | 11 | sec | |

```python
import numpy as np, cupy as cp
from numba import cuda
import math

@cuda.jit
def kernel(position, potential):
    i = cuda.threadIdx.x + cuda.blockIdx.x * cuda.blockDim.x
    N = len(position)
    if i >= N: return
    mass = 1 / N
    potential_i = 0
    for j in range(N):
        if i != j:
            dx = position[i,0] - position[j,0]
            dy = position[i,1] - position[j,1]
            dz = position[i,2] - position[j,2]
            r = math.sqrt(dx**2 + dy**2 + dz**2)
            potential_i += - mass / r
    potential[i] = potential_i

def calculate_potential(position : np.ndarray) -> np.ndarray:
    threads_per_block = 32
    N = len(position)
    blocks_per_grid = int(np.ceil(N/threads_per_block))
    position  = cp.array(position)
    potential = cp.empty(N)
    kernel[blocks_per_grid, threads_per_block](position, potentia
    return cp.asnumpy(potential)
```

# LIBRARY FUNCTION

## 🍎 COMPARISON WITH SAPPORO 🍊

| Our solution | Sapporo |
|---|---|
| full double precision | "double-single" |
| `sqrt` and division | `rsqrt` and multiplication |
| potential only | potential, acceleration, & jerk |
| optimized for large $n_i$ | also for small $n_i$ |



| | | | |
|---|---:|---|---|
| Naïve | 10 000 000 | sec | (~4 months) |
| NumPy | 20 000 | sec | (~5 hours) |
| CuPy | 1 300 | sec | (~20 minutes) |
| NumPy + Numba | 13 000 | sec | (~4 hours) |
| NumPy + Numba (parallel) | 560 | sec | (~9 minutes) |
| Numba (GPU) | 11 | sec | |
| Sapporo | 7.7 | sec | |

# TIPS FOR HOMEWORK EXERCISE

- Adding types to decorator can help:

```
@cuda.jit('void(float64[:,:], float64[:])')
```

- Instead of a CuPy array, we could do

```
position = cuda.to_device(position)
potential = cuda.device_array(N, dtype=np.float64)
...
return potential.copy_to_host()
```

- Remember that threads and blocks can be indexed in 2D.