# Fundamentals of CUDA C/C++

## Introduction

<u>CUDA</u> is a parallel computing platform developed by NVIDIA that allows general-purpose computing on GPUs (<u>GPGPU</u>). It is widely used in the fields related to high-performance computing, such as machine learning and computer simulation.

Typically, the parallelism and performance on GPUs are much higher than on CPUs. Let's look at the <u>FLOPs</u> of some latest CPU and GPU chips. With less than double prize, NVIDIA A100 offers nearly 26 times better performance than today's most powerful CPU.

Chips	Performance (TFLOPs)	Price in 2024
AMD Ryzen™ Threadripper™ PRO 7995WX	12.16	US\$ 10,000
NVIDIA Tesla A100	312	US\$ 18,000

In this article, we will discuss the basic methodology on CUDA programming.

## **Before You Start**

NVIDIA offers a command line tool nvidia-smi (system management interface) to show the information of the GPU installed on the computer. Don't worry if you do not have a NVIDIA GPU, simply turn to <u>Google Colab</u> for a free GPU environment. The rest of the article will be done on Google Colab.

After we connect to a GPU runtime on Colab, we can now run nvidia-smi in the interactive environment to see the information of the GPU installed. Note that we are in a Python environment and ! before system commands tells the environment to run it in a Linux shell.

Fri Jul	. 503	3:05:50 202	4					
NVIDI	A-SMI	535.104.05		Driver	Version: 5	535.104.05	CUDA Versio	on: 12.2
GPU   Fan	Name Temp	Perf	Persiste Pwr:Usag	ence-M ge/Cap	Bus-Id N	Disp.A Memory-Usage	Volatile   GPU-Util 	Uncorr. ECC Compute M. MIG M.
   0   N/A 	Tesla 66C	T4 P8	10W /	Off 70W	00000000 0Mie	:00:04.0 Off 3 / 15360MiB	   0% 	e Default N/A
+   Proce   GPU	esses:	ст.	PTD Type	Proces	s name			GPU Memory
	ID	ID	i ib iype	11000	55 Hanc			Usage

Great! We now have a Tesla T4 GPU! Make sure you have seen similar outputs and we will now move on to our next topic.

## **CUDA Programming Basics**

The prefix of a file in the CUDA language is .cu, and a CUDA file can be compiled along with C/C++/Fortran easily with the help of CMake. However, it is not today's topic and we will simply stop here.

### **Our First CUDA Program**

Before we talk about CUDA, let's first look at a C program that we are already familiar with.

```
#include <stdio.h>
1
2
   void printHelloworld() {
3
     printf("Hello, world!\n");
4
5
   }
6
7
  int main() {
     printHelloWorld();
8
9
   }
```

Straightforward, right? After rename the file as first.cu, we can compile it with CUDA compiler, just like what we usually do with a C compiler.

1 | !nvcc first.cu -o first -run

**nvcc** is the CUDA compiler, **-o** specified the name of the output executable file, **-run** executes the binary file right after the compilation process finishes. Now we can see the output from the program.





The output is expected. However, it is nothing different from what we usually do in C, since what we have written is exactly a C program. Now we will refactor the program so that it runs on GPU. In the context of CUDA programming, CPU is usually called host and GPU device. A function that runs on device is called **kernel function**. The return value of a kernel function **must** be void. To enable some function to run on device, we must add some special qualifier before the function. Here we list common qualifiers and explained their meanings.

Qualifiers	Meanings
global	The function is executed on device, called by host.
device	The function is executed on device, called by device.
host	The function is executed on host. This is the default choice when we omit the qualifier.

In this article, tasks are assigned by the host to the device, so we will only use the \_\_\_\_global\_\_\_\_ qualifier and default qualifier (no qualifier).

To launch a kernel function (don't forget what is a kernel), we use triple angle brackets to specify the configuration of the kernel function. Let's look at an example. The kernel function is defined here,

```
1 __global__ void printHelloworld() {
2 printf("Hello, world!\n");
3 }
```

and called here.

```
1 int main() {
2 printHelloWorld<<<1, 2>>>();
3 }
```

The configuration of the kernel will be explained later. Note that a kernel function is asynchronous, which is to mean that the host won't wait for the kernel function to finish. Rather, the host will continue to execute the codes below. Function <a href="mailto:cudaDeviceSynchronize">cudaDeviceSynchronize()</a> let host wait for **all** the kernels to finish. Let's put it together and look at the results. The complete CUDA program should look at this:

```
#include <stdio.h>
1
2
   __global___ void printHelloworld() {
3
4
      printf("Hello, world!\n");
   }
5
6
7
   int main() {
8
      printHelloworld<<<1, 2>>>();
9
      cudaDeviceSynchronize();
10
   }
```

Now we compile and run the program

1 !nvcc first.cu -o first -run

And we get the results

[5] !nvcc ./first.cu -o first -run → Hello, world! Hello, world!

which shows our program are truly running on GPU!

### **CUDA Thread Hierarchy**

#### **Kernel Configuration**

In the past section, we haven't really talked about what does the parameters in the triple angle brackets actually denote. Let's look at the CUDA thread hierarchy.



CUDA follows a grid-block-thread hierarchy. The overall structure is called **grid**, which is in black and contains blocks. The first parameter denotes **number of blocks** and the second parameter denotes **threads per block**. **Blocks** are painted blue and **threads** white in the slide. So, kernel function performwork<<<2, 4>>> actually says, the host assigns work to 2 blocks, with 4 threads each. Now you can explain why we have seen two echoes in our first CUDA program.

A thread can get its **block index** and **thread index within the block** by variables **blockIdx.x** and **threadIdx.x**, which are available directly in the definition of a kernel function.

Let's look at an example,



which gives the output:

### [6] !nvcc ./first.cu -o first -run

₹₹	Thread	0	from	block	1
	Thread	1	from	block	1
	Thread	2	from	block	1
	Thread	3	from	block	1
	Thread	0	from	block	0
	Thread	1	from	block	0
	Thread	2	from	block	0
	Thread	3	from	block	0

There are two more variables gridDim.x and blockDim.x, meaning the number of blocks in a grid and the number of threads in a block respectively and **corresponding** to the two parameters we passed when launched the kernel function.

You might wonder why there is a .x after each variable. The truth is that both the blocks and threads can be place in 3 dimensions. With the help of class dim3, we can pass a multidimensional configuration into a kernel. This trick is only for convenience in some particular applications, for example matrix operations, and would do nothing to the performance.

performWork<<<dim3(2,2,1), dim3(2,2,1)>>>();

The hierarchy of the configuration above will look like this:



Now we can access these values with gridDim.x, gridDim.y, gridDim.z, blockDim.x, blockDim.y, blockDim.z(), blockIdx.x, blockIdx.y, blockIdx.z, threadIdx.x, threadIdx.y, threadIdx.z.

These values are important for many calculations, as we will discuss later with memory.

#### **Streaming Multiprocessors**

In the end this section, we explain **streaming multiprocessors** (SMs) and offer a simple rule for picking proper numbers for numberofBlocks and threadsPerBlock.



SMs are basic units that execute tasks. As shown in the figure above, blocks are scheduled to run on SMs. When the number of blocks is multiple of the number of SMs, the execution will be efficient. Therefore, we cannot hardcode numberofBlocks and may use an API to get a value for a particular machine instead. Usually, taking numberOfBlocks as numberOfSMs times 16, 32, or 64 would be a good choice.

As for threadsPerBlock, it can be any integer between 1 and 1024.512 usually becomes a good choice.

Here is an example that we call a CUDA API to get the number of SMs. The official reference for struct cudaDeviceProp is <u>here</u>.

```
1 int deviceId;
2 cudaGetDevice(&deviceId);
3 
4 cudaDeviceProp prop;
5 cudaGetDeviceProperties(&prop, deviceId);
6 
7 int numberOfSMs = prop.multiProcessorCount;
```

Then, we can create a kernel configuration like this

```
1 int numberOfBlocks = numberOfSMs * 32;
2 int threadsPerBlock = 512;
```

3 some\_kernel<<<numberOfBlocks, threadsPerBlock>>>();

The default stream is special: it blocks all kernels in all other streams



DEEP LEARNING INSTITUTE

GPU tasks are scheduled in **streams** and the kernels in one stream is **serial**. Before, we didn't specify the stream and the kernels are executed in the **default stream**. Default streaming is **blocking**. That is to mean, when there is a task scheduled in the default stream, all the tasks on the **non-default streams** will be blocked, while all the other non-default streams are **non-blocking**, allowing concurrent kernel execution. CUDA offers a stream class cudaStream\_t and they are created with cudaStreamCreate() and destroyed with cudaStreamDestroy(). As we didn't mention before, a kernel configuration actually accepts 4 parameters, numberOfBlocks, threadsPerBlock, sharedMemoryBytes, and stream. We just left the latter 2 parameter in defaults before. Here is an example with concurrent streams.

```
1
    cudaStream_t stream1;
 2
    cudaStream_t stream2;
 3
    cudaStreamCreate(&stream1);
4
    cudaStreamCreate(&stream2);
 5
    some_kernel<<<<numberOfBlocks, threadsPerblock, 0, stream1>>>();
6
7
    some_kernel<<<numberOfBlocks, threadsPerblock, 0, stream2>>>();
8
9
    cudaStreamDestroy(stream1);
10
    cudaStreamDestroy(stream2);
```

In this part of program, the kernels in stream1 and stream2 will be concurrent. Since we don't need shared memory here, we just left them in 0. Note that although cudaStreamDestroy() returns immediately after calling, the streams are not actually destroyed until the kernels in the stream finish, so we don't need to worry about the side effect of destroying a stream.

## **CUDA Memory Management**

So far, we have discussed how to launch a kernel running on device. However, such kernels could only access memory automatically allocated. That is, local variables in kernel function definition. We wish to allocate memory that can be accessed by both the host and the device. CUDA offers a unified memory model so that we don't need to worry about memory access.

### **Unified Memory**

#### **Allocation and Freeing**

To allocate unified memory that can be accessed both on the device and the host, we call cudaMallocManaged() function.



Thus, the integer array can be accessed both on the device and the host. To free allocated unified memory, we call cudaFree() utility.

```
1 cudaFree(array);
```

#### **Reduce Page Faults**



DEEP LEARNING

When UM was initially allocated, it may not be resident on CPU or GPU. If the memory was first initialized by CPU then GPU, a <u>page fault</u> occurs. Memory will be transferred from host to device and slow down the tasks. Similarly, if UM was initialized on GPU and then accessed by CPU, a page fault occurs and memory transfer begins. The place memory is resident on depends on the last access. When a page fault is present, memory is transferred is small batch size. If we can predict a page fault, we can transfer the corresponding memory in advance with bigger batch size to increase the efficiency. CUDA offers an API called cudaPrefetch to perform such behaviors. Here is an example.

```
1 int deviceId;
2 cudaGetDevice(&deviceId);
3 
4 // Allocate unified memory
5 int numberOfSMs = prop.multiProcessorCount;
6 int N = 2<<20;
7 int* array;
```

```
8 int size = N * sizeof(int);
 9 cudaMallocManaged(&array, size);
 10
 11 // Initialize the UM on CPU, then access it on GPU
 12
     init_on_cpu(array, size);
    cudaPrefetchAsync(array, size, deviceId); // Trasfer the memory from host to
 13
     device
     access_memory<<<numberOfBlocks, threadsPerBlock>>>();
 14
 15
     cudaDeviceSynchronize();
 16
 17
     // Free memory
 18 cudaFree(array);
```

The third parameter of cudaPrefetchAsync specifies the direction of memory transfer. When filled with deviceId, the memory is transferred from host to device (HtoD), while cudaCpuDeviceId specifies a transfer from device to host (DtoH). Note that the variable cudaCpuDeviceId can directly be accessed globally and we need no APIs to get this variable.

#### **Index Calculation**

With unified memory, we can finally schedule some tasks on GPU. Let's look at an example. Suppose we have two vectors containing integers, each with length 2^22, and we want to accelerate vector addition with CUDA. Our number of threads might no be bigger enough to establish a bijection between thread indices and vector indices, so a thread must execute more than one addition. The convention is to define a variable stride that equals to the total number of threads in a (also the only) grid, and increase the loop index by stride each time. We assume the kernel configuration is 1-d and here is how we calculate index. Note that we are checking index boundary each time to avoid unexpected memory access.

```
1 __global__ void vectorAdd(int* a, int* b, int* res, int N) {
2 int idx = blockDim.x * blockIdx.x + threadIdx.x;
3 int stride = gridDim.x * blockDim.x;
4 for (int i = idx; i < N; i += stride) {
5 res[i] = a[i] + b[i];
6 }
7 }</pre>
```

### **Manual Memory Management and Streaming**

#### **Basic Commands**

Although UM is powerful enough, we may still want to manage the memory ourselves to further optimize the efficiency. Here are some commands for manually managing the memory.

- cudaMalloc will allocate the memory on GPU. The memory is **not** accessible on CPU.
- cudaFree frees the memory allocated on device.
- cudaMallocHost will allocate the memory on CPU just like what a normal malloc does. The memory will be page locked on host. Too many page locked memory would reduce CPU performance.
- cudaFreeHost frees the memory allocated on host.
- cudaMemcpy copies the memory DtoH or HtoD, instead of transferring.

• cudaMemcpyAsync allows asynchronously memory copying (explained later).

Let's look at an example.

```
1
   int N = 2 << 20;
 2
    int size = N * sizeof(int);
 3
   int* array_device;
   int* array_host;
 4
 5
    cudaMalloc(&array_device, size);
    cudaMallocHost(&array_host, size);
 6
 7
8
    init_on_cpu(array_host, N);
9
    cudaMemcpy(array_device, array_host, size, cudaMemcpyHostToDevice);
    some_kernel<<<blocks, threads, 0, stream>>>();
10
11
    cudaMemcpy(array_host, array_device, size, cudaMemcpyDeviceToHost);
12
    cudaFree(array_device);
13
    cudaFreeHost(array_host);
14
```

Note that cudaMemcpy first accepts the destination pointer then source pointer, unlike Linux cp first-source-then-destination convention. cudaMemcpyDeviceToHost and cudaMemcpyHostToDevice are two variables that can be directly and globally accessed that specify the direction of memory copying.

#### Asynchronous Memory Copying

In the last example, memory copy starts after kernel finishes. To optimize the process, we can start asynchronous (or concurrent) memory copying right after a part of kernel finishes.



An simple approach is to divide the kernel in several segments. A segment of memory copying starts right after a segment of kernel finishes. We will refactor the former vector addition program and take it as an example.



```
int segmentSize = size / numberOfSegments;
4
5
      for (int i = 0; i < numberOfSegments; i++) {</pre>
6
        cudaStream_t stream;
7
        cudaStreamCreate(&stream);
8
        int offset = i * segmentN;
9
        vectorAdd<<<blocks, threads, 0, stream1>>>(a_device + offset, b_device +
    offset.
10
                                                     c_device + offset, segmentN);
        cudaMemcpyAsync(c_host + offset, c_device + offset, segmentSize,
11
    cudaMemcpyDeviceToHost, stream);
12
        cudaStreamDestroy(stream);
13
      }
14
      cudaDeviceSynchronize();
```

The piece of program above divides the kernel vectorAdd in 4 segments. After one segment of kernel finishes, asynchronous memory copying starts. Every stream first executes the kernel then copying the corresponding memory to the host. Note that we must **carefully** handle all the indices here to avoid illegal memory access. After all of theses are done, the stream is destroyed (recall stream destruction behavior). The whole accelerated program is pasted below for your reference. Don't forget to destroy unused streams and free unused memories.

```
1
    #include <stdio.h>
 2
    __global__ void init_on_gpu(int* a, int init_val, int N) {
3
 4
      int idx = blockDim.x * blockIdx.x + threadIdx.x;
5
      int stride = gridDim.x * blockDim.x;
 6
      for (int i = idx; i < N; i += stride) {
 7
        a[i] = init_val;
8
      }
9
    }
10
    __global__ void vectorAdd(int* a, int* b, int* res, int N) {
11
12
      int idx = blockDim.x * blockIdx.x + threadIdx.x;
13
      int stride = gridDim.x * blockDim.x;
      for (int i = idx; i < N; i += stride) {
14
15
        res[i] = a[i] + b[i];
16
      }
    }
17
18
    int verify_on_cpu(int* a, int res, int N) {
19
      for (int i = 0; i < N; i++) {
20
21
        if (a[i] != res) {
22
          return 0;
23
        }
24
      }
25
      return 1;
26
    }
27
28
    int main() {
29
      // Variable declarations
30
      int N = 2 << 20;
      int size = N * sizeof(int);
31
32
      int* a_device;
33
      int* b_device;
```

```
34
      int* c_device;
35
      int* c_host;
      cudaMalloc(&a_device, size);
36
37
      cudaMalloc(&b_device, size);
38
      cudaMalloc(&c_device, size);
39
      cudaMallocHost(&c_host, size);
40
41
      // Kernel configuration
42
      int deviceId;
43
      cudaGetDevice(&deviceId);
44
      cudaDeviceProp prop;
45
      cudaGetDeviceProperties(&prop, deviceId);
      int SMs = prop.multiProcessorCount;
46
47
      int blocks = SMs * 32;
      int threads = 512;
48
49
50
      // Stream creation
51
      cudaStream_t stream1;
52
      cudaStream_t stream2;
53
      cudaStream_t stream3;
54
      cudaStreamCreate(&stream1);
55
      cudaStreamCreate(&stream2);
56
      cudaStreamCreate(&stream3);
57
58
      // Initialize arrays on device
      init_on_gpu<<<blocks, threads, 0, stream1>>>(a_device, 2, N);
59
      init_on_gpu<<<blocks, threads, 0, stream2>>>(b_device, 3, N);
60
      init_on_gpu<<<<blocks, threads, 0, stream3>>>(c_device, 0, N);
61
62
      cudaDeviceSynchronize();
63
      // Perform vector addition in segments
64
      int numberOfSegments = 4;
65
66
      int segmentN = N / numberOfSegments;
      int segmentSize = size / numberOfSegments;
67
68
      for (int i = 0; i < numberOfSegments; i++) {</pre>
69
        cudaStream_t stream;
70
        cudaStreamCreate(&stream);
71
        int offset = i * segmentN;
        vectorAdd<<<<blocks, threads, 0, stream1>>>(a_device + offset, b_device +
72
    offset, c_device + offset, segmentN);
73
        cudaMemcpyAsync(c_host + offset, c_device + offset, segmentSize,
    cudaMemcpyDeviceToHost, stream);
74
        cudaStreamDestroy(stream);
75
      }
76
      cudaDeviceSynchronize();
77
      // Verify results
78
79
      if (verify_on_cpu(c_host, 5, N)) {
80
        printf("Results are correct!\n");
81
      } else {
82
        printf("Results are incorrect!\n");
83
      }
84
85
      // Free memories
86
      cudaFree(a_device);
87
      cudaFree(b_device);
```

```
88 cudaFree(c_device);
89 cudaFreeHost(c_host);
90 }
```

## **Error Handling**

CUDA usually does not report runtime errors, so we must detect and record them manually. CUDA offers a class cudaError\_t to handle errors. The return values of CUDA APIs are cudaError\_t, allowing us to catch error directly.

```
#include <stdio.h>
1
2
   int main() {
3
     int* array;
4
     int size = -1;
5
     cudaError_t err;
      err = cudaMallocManaged(&array, size);
6
7
     if (err != cudaSuccess) {
      fprintf(stderr, "Error: %s\n", cudaGetErrorString(err));
8
9
      }
10 }
```

This program will throw an error because -1 is not a valid size.

1 Error: out of memory

However, the return values of custom kernels are void, meaning we cannot catch error in the same way when launching those kernels. CUDA offers cudaGetLastError to catch the last error thrown. Also, we can create a utility function to encapsulate such processes.

```
inline void checkCudaError() {
1
2
    cudaError_t err;
3
     err = cudaGetLastError();
    if (err != cudaSuccess) {
4
       fprintf(stderr, "Error: %s\n", cudaGetErrorString(err));
5
6
       assert(err == cudaSuccess);
7
     }
   }
8
```

Therefore, we can check the errors when launching kernels.

```
1 some_kernel<<<-1, -1>>>();
2 checkCudaError();
```

which yields

1 Error: invalid configuration argument

Such error handling may help debug the CUDA program.

## **Performance Profiling**

NVIDIA Nsight Systems command line tool (nsys) is a command line profiler that will gather following information:

- Profile configuration details
- Report file(s) generation details
- CUDA API Statistics
- CUDA Kernel Statistics
- CUDA Memory Operation Statistics (time and size)
- OS Runtime API Statistics

To install nsys in Colab, run the following commands:

```
1 !wget
https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64/nsi
ght-systems-2023.2.3_2023.2.3.1001-1_amd64.deb
2 !apt update
3 !apt install ./nsight-systems-2023.2.3_2023.2.3.1001-1_amd64.deb
4 !apt --fix-broken install
```

Then, we can profile the running information of our CUDA program. Here we take vector addition as an example. profile means we want to profile the executable and --stats=true tells nsys to print all the information.

#### 1 Insys profile --stats=true vector-add

A part of printed information is pasted below:

```
[5/8] Executing 'cuda_api_sum' stats report
  Time (%) Total Time (ns) Num Calls Avg (ns)
                                                                                        Med (ns) Min (ns) Max (ns)
                                                                                                                                                    StdDev (ns)
                                                                                                                                                                                                   Name

        Image (ns)
        Mode Calls
        Avg (ns)
        Med (ns)
        Med (ns)
        Med (ns)

        192,376,747
        3
        64,125,582.3
        102,474.0
        94,203

        4,932,889
        1
        4,932,889.0
        4,932,889.0
        4,932,889.0
        4,932,889.0

        1,727,437
        1
        1,727,437.0
        1,727,437.0
        1,727,437.0
        1,727,437.0

        946,023
        3
        315,341.0
        389,557.0
        163,210

        718,243
        2
        359,121.5
        359,121.5
        97,247

        280,743
        7
        40,106.1
        6,209.0
        5,624

        134,190
        1
        134,190.0
        134,190.0
        134,190.0

        73,363
        7
        10,480.4
        4,131.0
        2,853

        27,767
        4
        6,941.8
        5,908.5
        4,612

        15,667
        3,916.8
        3,328.5
        3,175

                                                                                 -- --
                                                                                                       - -
        96.0
                                                                                                                94,203 192,180,070 110,898,439.5 cudaMalloc
                                                                                                                                    4,032,889
          2.0
         0.9
                                                                                                                                 1,727,437
          0.5
          0.4
          0.1
          0.1
                                                                                                               134,190
                                                                                                                                                                   0.0 cudaGetDeviceProperties v2 v12000
          0.0
                                                                                                                  2,853
          0.0
                                                                                                                   4,612
          0.0
                                15,667
                                                                        3,916.8
                                                                                             3,328.5
                                                                                                                   3,175
                                  2,002
          0.0
                                                                       2,002.0
                                                                                            2,002.0
                                                                                                                   2,002
                                                          1
[6/8] Executing 'cuda_gpu_kern_sum' stats report
  Time (%) Total Time (ns) Instances Avg (ns) Med (ns) Min (ns) Max (ns) StdDev (ns)
                                                                                                                                 -----
                                                                                                                -----
                                                                                                                                  241.0 vectorAdd(int *, int *, int *, int)
3,701.8 init_on_gpu(int *, int, int)
                              141,309435,327.335,311.5116,989338,996.340,479.0
                                                                                                                35,615
41,727
        54.7
45.3
                                                                                                   35,071
                              116,989
                                                                                                  34,783
[7/8] Executing 'cuda_gpu_mem_time_sum' stats report
 Time (%) Total Time (ns) Count Avg (ns) Med (ns) Min (ns) Max (ns) StdDev (ns)
                                                                                                                                                         Operation
                              660,782
                                                4 165,195.5 163,995.5 160,699 172,092
      100.0
                                                                                                                                5,186.2 [CUDA memcpy DtoH]
[8/8] Executing 'cuda_gpu_mem_size_sum' stats report
  Total (MB) Count Avg (MB) Med (MB) Min (MB) Max (MB) StdDev (MB)
                                                                                                                             Operation
                                                                                                 -----
                                    2.097
         8.389
                       4
                                                    2.097
                                                                     2.097
                                                                                   2.097
                                                                                                          0.000 [CUDA memcpy DtoH]
```

Here we can trace the running time of kernels and the memory behavior. You can compare the results of different numberofBlocks and see what kind of memcpy page faults will bring about.

## **Final Exercise**

An <u>n-body</u> simulator predicts the individual motions of a group of objects interacting with each other gravitationally. Below is a simple, though working, n-body simulator for bodies moving through 3 dimensional space.

In its current CPU-only form, this application takes about 5 seconds to run on 4096 particles, and **20 minutes** to run on 65536 particles. Your task is to GPU accelerate the program, retaining the correctness of the simulation.

### **Considerations to Guide Your Work**

Here are some things to consider before beginning your work:

- Especially for your first refactors, the logic of the application, the bodyForce function in particular, can and should remain largely unchanged: focus on accelerating it as easily as possible.
- The code base contains a for-loop inside main for integrating the interbody forces calculated by bodyForce into the positions of the bodies in the system. This integration both needs to occur after bodyForce runs, and, needs to complete before the next call to bodyForce. Keep this in mind when choosing how and where to parallelize.
- Use a **profile driven** and iterative approach.
- You are not required to add error handling to your code, but you might find it helpful, as you are responsible for your code working correctly.

```
1 #include <math.h>
 2 #include <stdio.h>
   #include <stdlib.h>
 3
   #include "timer.h"
 4
   #include "files.h"
 5
 6
 7
    #define SOFTENING 1e-9f
 8
9
     * Each body contains x, y, and z coordinate positions,
10
11
    * as well as velocities in the x, y, and z directions.
12
     */
13
   typedef struct { float x, y, z, vx, vy, vz; } Body;
14
15
16
   /*
17
    * Calculate the gravitational impact of all bodies in the system
18
    * on all others.
19
     */
20
21
   void bodyForce(Body *p, float dt, int n) {
22
     for (int i = 0; i < n; ++i) {</pre>
23
        float Fx = 0.0f; float Fy = 0.0f; float Fz = 0.0f;
24
        for (int j = 0; j < n; j++) {
25
26
         float dx = p[j].x - p[i].x;
27
          float dy = p[j].y - p[i].y;
28
          float dz = p[j].z - p[i].z;
```

```
float distSqr = dx*dx + dy*dy + dz*dz + SOFTENING;
29
30
          float invDist = rsqrtf(distSqr);
31
          float invDist3 = invDist * invDist * invDist;
32
          Fx += dx * invDist3; Fy += dy * invDist3; Fz += dz * invDist3;
33
34
        }
35
36
        p[i].vx += dt*Fx; p[i].vy += dt*Fy; p[i].vz += dt*Fz;
37
      }
    }
38
39
40
41
    int main(const int argc, const char** argv) {
42
      // The assessment will test against both 2<11 and 2<15.</pre>
43
      // Feel free to pass the command line argument 15 when you generate
44
    ./nbody report files
45
      int nBodies = 2 < <11;
      if (argc > 1) nBodies = 2<<atoi(argv[1]);</pre>
46
47
48
      // The assessment will pass hidden initialized values to check for
    correctness.
      // You should not make changes to these files, or else the assessment
49
    will not work.
50
      const char * initialized_values;
      const char * solution_values;
51
52
53
     if (nBodies == 2<<11) {
54
        initialized_values = "files/initialized_4096";
        solution_values = "files/solution_4096";
55
      } else { // nBodies == 2<<15</pre>
56
57
        initialized_values = "files/initialized_65536";
58
        solution_values = "files/solution_65536";
59
      }
60
61
      if (argc > 2) initialized_values = argv[2];
62
      if (argc > 3) solution_values = argv[3];
63
      const float dt = 0.01f; // Time step
64
      const int nIters = 10; // Simulation iterations
65
66
      int bytes = nBodies * sizeof(Body);
67
      float *buf;
68
69
70
      buf = (float *)malloc(bytes);
71
      Body *p = (Body*)buf;
72
73
74
      read_values_from_file(initialized_values, buf, bytes);
75
76
      double totalTime = 0.0;
77
78
      /*
       * This simulation will run for 10 cycles of time, calculating
79
    gravitational
80
       * interaction amongst bodies, and adjusting their positions to reflect.
```

```
81 */
 82
 83
       for (int iter = 0; iter < nIters; iter++) {</pre>
         StartTimer();
 84
 85
 86
       /*
 87
        * You will likely wish to refactor the work being done in `bodyForce`,
        * and potentially the work to integrate the positions.
 88
        */
 89
 90
 91
         bodyForce(p, dt, nBodies); // compute interbody forces
 92
 93
       /*
 94
        * This position integration cannot occur until this round of `bodyForce`
     has completed.
        * Also, the next round of `bodyForce` cannot begin until the integration
 95
     is complete.
 96
        */
 97
         for (int i = 0 ; i < nBodies; i++) { // integrate position</pre>
 98
 99
           p[i].x += p[i].vx*dt;
           p[i].y += p[i].vy*dt;
100
           p[i].z += p[i].vz*dt;
101
         }
102
103
104
         const double tElapsed = GetTimer() / 1000.0;
         totalTime += tElapsed;
105
106
       }
107
       double avgTime = totalTime / (double)(nIters);
108
       float billionsOfOpsPerSecond = 1e-9 * nBodies * nBodies / avgTime;
109
110
       write_values_to_file(solution_values, buf, bytes);
111
       // You will likely enjoy watching this value grow as you accelerate the
112
     application,
113
       // but beware that a failure to correctly synchronize the device might
     result in
114
       // unrealistically high values.
       printf("%0.3f Billion Interactions / second", billionsOfOpsPerSecond);
115
116
117
       free(buf);
118 }
```