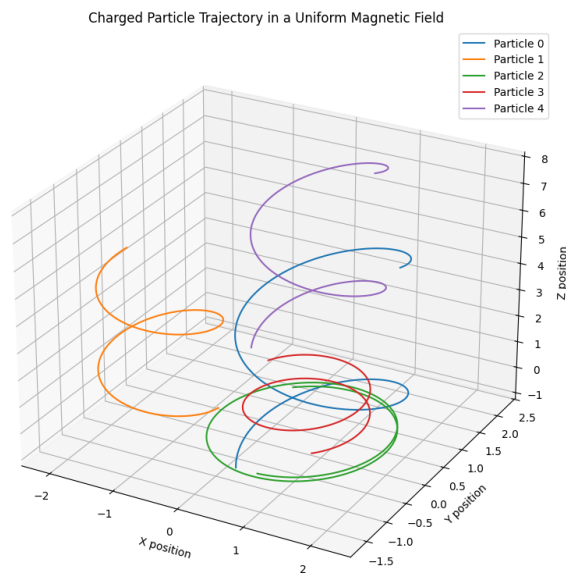


INTRODUCTION TO

High-performance test particle tracing code



Author:

El Ghaib, Adam

Contents

	Page
1 Introduction	1
1.1 Objective	1
1.2 Considerations	1
1.3 Environment	1
2 Simulations	2
2.1 Pure CUDA	2
2.2 Thrust library	6
2.3 Python	8
3 Performance	9
4 Conclusion	10
Bibliography	i
A Appendix	ii
A.1 Python Code	ii

List of Figures

1	Trajectory of 10 particles simulated in CUDA.	6
2	Trajectory of 10 particles simulated in CUDA using Thrust library.	7
3	Execution time in <i>ms</i> vs the <i>number of particles</i> for each program.	9

1 Introduction

1.1 Objective

This project aims to introduce the author into parallel ion tracing programming. To do so, three simple programs have been written. Two of them were written in CUDA language, where one of these uses Thrust library, to enhance development productivity. And the third one has been written in Python.

Moreover, the execution time of the three programs is compared to determine the optimal approach to build a high-performance test particle tracing simulator.

1.2 Considerations

To initially maintain the simulator simple:

- The electric $\mathbf{E}(\mathbf{r},t)$ and magnetic $\mathbf{B}(\mathbf{r},t)$ fields are considered constant through space and time.
- Elastic or inelastic collisions between particles are not taken into account.
- Only Lorentz force will act upon the particles and the 3D Newton's equation of motion is solved numerically using a fourth order Runge-Kutta algorithm.

1.3 Environment

The simulations have been executed in a Windows 11 operating system using Microsoft Visual Studio 2022. The GPU used was an Nvidia GeForce RTX 4070, the CPU was an Intel(R) Core(TM) i9-14900HX @ 2.20 GHz, and the RAM is of 32 GB.

2 Simulations

As mentioned before, to keep the programs as simple as possible, only the Lorentz force eq.1 will be considered.

$$m \frac{d\vec{v}}{dt} = q(\vec{E} + \vec{v} \times \vec{B}) \quad (1)$$

Furthermore, the mass is set to $\mathbf{m} = 1 \text{ kg}$, particles charge to $\mathbf{q} = 1 \text{ C}$, electric field to $\vec{\mathbf{E}} = \mathbf{0} \text{ N/C}$ and magnetic field to $\vec{\mathbf{B}} = 1 \text{ T}$. Finally, the initial velocity components of each particle will be a random integer between -5 and 5 m/s .

2.1 Pure CUDA

The first program works as follows. First, basic CUDA and C/C++ libraries are included and some operator functions like addition (+), multiplication (*), and others as (+ =), (− =) are defined, since in CUDA do not come by default at the device. Moreover, a Particle structure is created, only the position and velocity are fields of the structure as the other variables (q, m) will be the same for every particle.

```

1 #include "cuda_runtime.h"
2 #include "device_launch_parameters.h"
3 #include <stdio.h>
4 #include <iostream>
5 #include <fstream>
6 #include <vector>
7
8 struct Particle {
9     float3 position;
10    float3 velocity;
11 };
12
13 __device__ float3 operator+(const float3 &a, const float3 &b) {
14     return make_float3(a.x+b.x, a.y+b.y, a.z+b.z);
15 }
16
17 __device__ float3 operator*(const float3& a, const float3& b) {
18     return make_float3(a.x * b.x, a.y * b.y, a.z * b.z);
19 }

```

Then a device function is written to calculate the Lorentz force:

```

1 __device__ float3 lorentz_force(float3 velocity, float3 position) {
2     float3 E = make_float3(0.0f, 0.0f, 0.0f); // Electric field
3     float3 B = make_float3(0.0f, 0.0f, 1.0f); // Magnetic field
4
5     float3 v_cross_B = make_float3(
6         velocity.y * B.z - velocity.z * B.y,
7         velocity.z * B.x - velocity.x * B.z,
8         velocity.x * B.y - velocity.y * B.x
9     );
10
11     float q = 1.0f; // Charge of the particle
12     float m = 1.0f; // Mass of the particle
13     return (q/m) * (E + v_cross_B);
14 }

```

The kernel function that will be executed at each thread is a fourth-order Runge-Kutta algorithm. Which gives a decent approximation to the solution of the differential equation of motion eq.1.

```

1 // CUDA kernel to update particle positions
2 __global__ void updateParticles(Particle* particles, int numParticles, float dt
3 ) {
4     int idx = blockIdx.x * blockDim.x + threadIdx.x;
5     if (idx < numParticles) {
6         Particle p = particles[idx];
7
8         float3 k1_v = lorentz_force(p.velocity, p.position);
9         float3 k1_p = p.velocity;
10
11         float3 temp_velocity = p.velocity + 0.5f * k1_v * dt;
12         float3 temp_position = p.position + 0.5f * k1_p * dt;
13         float3 k2_v = lorentz_force(p.velocity + 0.5f * k1_v * dt, p.position +
14             0.5f * k1_p * dt);
15         float3 k2_p = temp_velocity;
16
17         temp_velocity = p.velocity + 0.5f * k2_v * dt;
18         temp_position = p.position + 0.5f * k2_p * dt;
19         float3 k3_v = lorentz_force(p.velocity + 0.5f * k2_v * dt, p.position +
20             0.5f * k2_p * dt);
21         float3 k3_p = temp_velocity;
22
23         temp_velocity = p.velocity + k3_v * dt;
24         temp_position = p.position + k3_p * dt;

```

```

22     float3 k4_v = lorentz_force(p.velocity + k3_v * dt, p.position + k3_p *
    dt);
23     float3 k4_p = temp_velocity;
24
25     p.velocity += (k1_v + 2.0f * k2_v + 2.0f * k3_v + k4_v) * (dt / 6.0f);
26     p.position += (k1_p + 2.0f * k2_p + 2.0f * k3_p + k4_p) * (dt / 6.0f);
27
28     particles[idx] = p;
29 }
30 }

```

In the main() function, the number of ions, simulation time, and time steps are set. Then two vectors with the structure particle are created, one is initialized with initial positions and velocities at the host (CPU) and the other copies this data and allocates it to the device (GPU) memory with cudaMalloc(). In addition, a thrird vector is created to store the xyz positions of the particles to later be exported as a binary file and plot it in Python using matplotlib.

```

1     const int numParticles = 10;
2     const int numTimesteps = 10000;
3     const float dt = 0.001f;
4
5     // Vector to store positions for plotting
6     std::vector<float> positions;
7     positions.reserve(numParticles * numTimesteps * 3);
8
9     // Create host and device vectors
10    Particle* h_particles = new Particle[numParticles];
11    Particle* d_particles;
12    cudaMalloc(&d_particles, numParticles * sizeof(Particle));
13
14    // Initialize particles
15    for (int i = 0; i < numParticles; ++i) {
16        h_particles[i].position = make_float3(1.0f, 0.0f + i, 0.0f);
17        h_particles[i].velocity = make_float3(i, 1.0f, 0.0f); // Constant
            velocity
18    }
19
20    // Copy particles to device
21    cudaMemcpy(d_particles, h_particles, numParticles * sizeof(Particle),
        cudaMemcpyHostToDevice);

```

Next, the number of blocks and threads that are going to be executed parallelly in the device is defined according to the number of simulated particles. The kernel function will be executed "*numTimesteps* = 10000" times. For each time step, the result of the `__global__` function is copied back to the host vector and for each particle their positions are saved in the *positions* vector. Finally, the memory is cleared and the positions file is exported.

```

1  int blockSize = 256;
2  int numBlocks = (numParticles + blockSize - 1) / blockSize;
3
4  // Run simulation for numTimesteps
5  for (int t = 0; t < numTimesteps; ++t) {
6      updateParticles << < numBlocks, blockSize >> > (d_particles,
7          numParticles, dt);
8      cudaMemcpy(h_particles, d_particles, numParticles * sizeof(Particle),
9          cudaMemcpyDeviceToHost);
10
11     // Store positions for plotting
12     for (int i = 0; i < numParticles; ++i) {
13         positions.push_back(h_particles[i].position.x);
14         positions.push_back(h_particles[i].position.y);
15         positions.push_back(h_particles[i].position.z);
16     }
17 }
18
19 // Free memory
20 cudaFree(d_particles);
21 delete[] h_particles;
22
23 // Save positions to a file
24 const char* path = "C:/Users/adame/OneDrive/Escritorio/CUDA/
25     Particle_Tracing/.venv/particle_positions.bin";
26 std::ofstream outFile(path, std::ios::binary);
27 outFile.write(reinterpret_cast<char*>(positions.data()), positions.size() *
28     sizeof(float));
29 outFile.close();
30
31 return 0;

```


The result of the simulation is shown in the next figure fig.1.

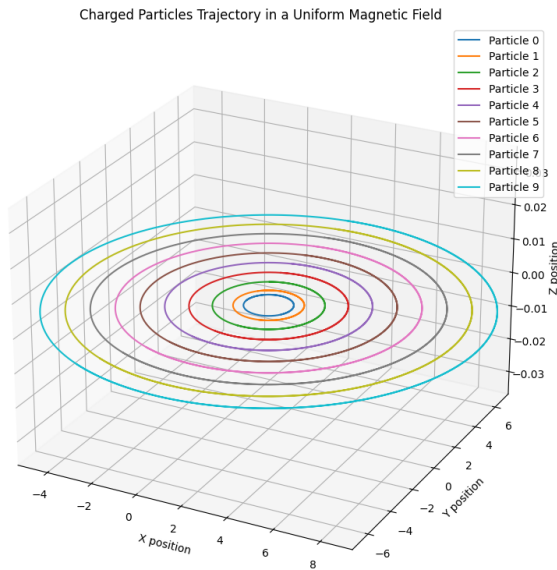


Figure 1: Trajectory of 10 particles simulated in CUDA.

2.2 Thrust library

The Thrust library is a high-level parallel algorithms library in CUDA, designed to resemble the C++ Standard Template Library (STL). Thrust enables developers to leverage GPU acceleration without needing to manage low-level details of CUDA programming explicitly [2].

Nonetheless, using Thrust can sometimes lead to longer compilation times due to the heavy use of templates and metaprogramming.

The main differences between this program and the previous one lie in `main()`. Here, the data translation between the host and device vectors can be done more intuitively. Furthermore, a built-in random number generator has been used to give random initial velocities to each particle.

```

1  thrust::host_vector<Particle> h_particles(numParticles);
2  // Initialize particles with random positions and velocities
3  thrust::default_random_engine rng;
4  thrust::uniform_real_distribution<float> dist(-5.0f, 5.0f);
5  for (int i = 0; i < numParticles; ++i) {
6      h_particles[i].position = make_float3(0.0f, 0.0f, 0.0f);
7      h_particles[i].velocity = make_float3(dist(rng), dist(rng), dist(rng));
8  }
9  thrust::device_vector<Particle> d_particles = h_particles;
10 for (int t = 0; t < numTimesteps; ++t) {
11     thrust::for_each(thrust::device, d_particles.begin(), d_particles.end()
12                     , UpdateParticles(dt));
13     thrust::copy(d_particles.begin(), d_particles.end(), h_particles.begin
14                 ());
15     for (int i = 0; i < numParticles; ++i) {
16         positions.push_back(h_particles[i].position.x);
17         positions.push_back(h_particles[i].position.y);
18         positions.push_back(h_particles[i].position.z);
19     }
20 }

```

The result of this simulation is shown in the next figure fig.2.

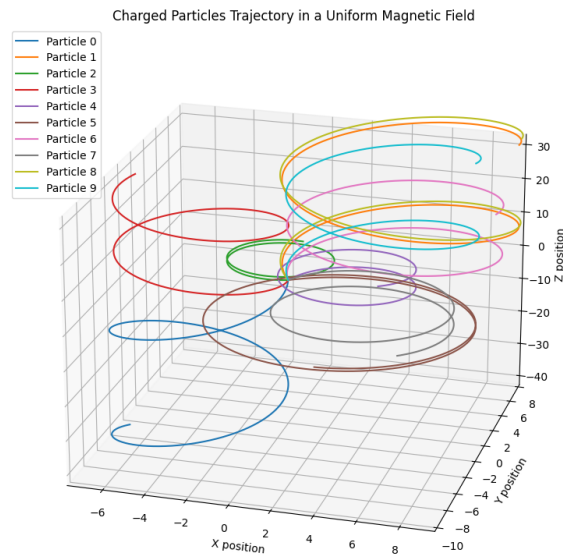


Figure 2: Trajectory of 10 particles simulated in CUDA using Thrust library.

2.3 Python

Finally, a Python simulation has also been written to compare its execution time, in section 3 with the last two simulations. The code will not be explained in this report, but it can be found in appendix A.1.

3 Performance

In this section, the execution times of the three proposed programs are compared. It is worth mentioning that the programs have been slightly modified to only measure the elapsed time that took for them to run the threads and simulate the particles state, that is to say, for the first program (2.1), the time between lines 4 and 8. Thus, the time used to initialize the vectors or copy the data and export it into a file has not been taken into account.

The results are shown in the next figure fig.3.

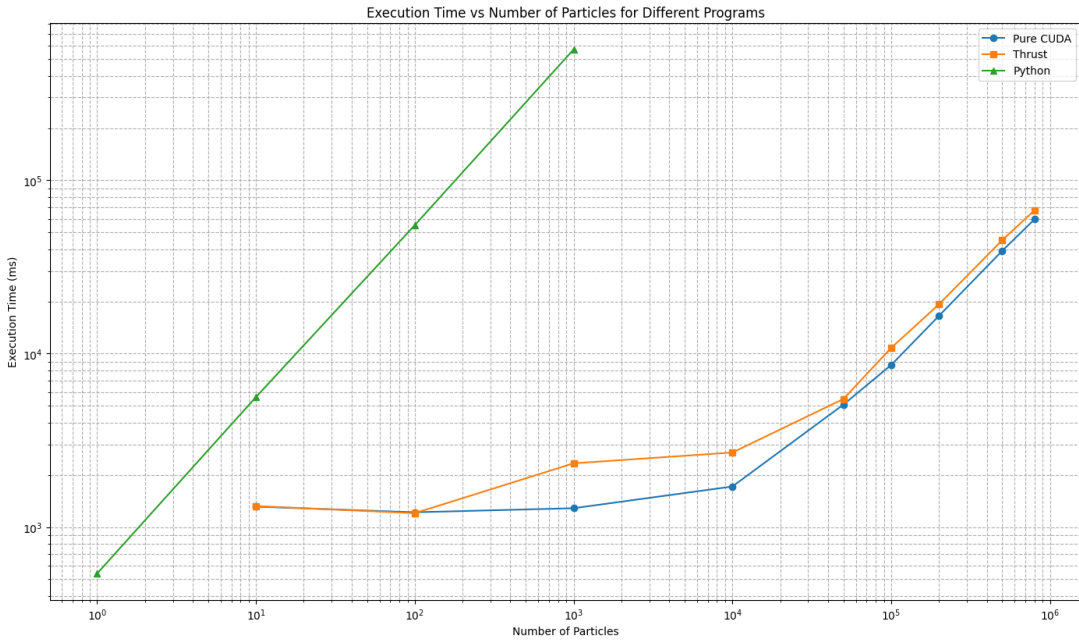


Figure 3: Execution time in *ms* vs the *number of particles* for each program.

It can be noticed that when the number of particles n is low, e.g $n \lesssim 20$, the python program is the fastest as it uses the high-frequency cores of the CPU. However, the clear superiority of the GPU when the number of particles is increased does not take long to appear. And at $n = 1000$ the CUDA code is already more than **400x faster**.

On the other hand, the execution time difference between the pure CUDA code and the one using the Thrust library is not very remarkable. It indeed takes more time for the Thrust program to compile (5s over 2s of the pure CUDA) and to execute. Nevertheless, the simulations have omitted the effect of collisions and have been executed under constant fields, which greatly simplified the programs and shared memory processes between cores. For a more complex simulation, the slightly higher execution time cost over easier memory management for the developer that Thrust offers, could be worth it.

4 Conclusion

With all the data previously presented, it can be concluded that the programs written in CUDA perform remarkably better than the Python one. Moreover, for a higher complexity simulation, the Thrust library could facilitate the developer, memory management between cores and, make the program more understandable, at the expense of a slightly higher execution time.

Finally, as future steps for this didactic project. A variable magnetic field could be added, to do so and maintain the magnetic field divergence null, Chebyshev polynomials would be used. Furthermore, the Focker-Planck theory could also be used to implement the effect of elastic collisions. An elaborated CUDA particle simulator that would be followed as an example can be found in [1].

References

- [1] C.F. Clauser, R. Farengo, H.E. Ferrari. Focus: A full-orbit cuda solver for particle simulations in magnetized plasmas. DOI: <https://doi.org/10.1016/j.cpc.2018.07.018>, year: 2018.
- [2] Nvidia. Thrust library. <https://developer.nvidia.com/thrust>.

A Appendix

A.1 Python Code

Particle simulator written in Python which uses threads.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import threading
4 import time
5
6 # Constants
7 q = 1.0 # Charge of the particle
8 m = 1.0 # Mass of the particle
9 Bz = 1.0 # Magnetic field in the z-direction
10
11 # Time parameters
12 t0 = 0.0 # Initial time
13 tf = 10.0 # Final time
14 dt = 0.001 # Time step
15 num_steps = int((tf - t0) / dt) # Number of time steps
16
17 def lorentz_force(v, B):
18     return q * np.cross(v, B)
19
20 def rk4_step(r, v, dt):
21     B = np.array([0, 0, Bz])
22     k1_r = v
23     k1_v = lorentz_force(v, B) / m
24
25     k2_r = v + 0.5 * dt * k1_v
26     k2_v = lorentz_force(v + 0.5 * dt * k1_v, B) / m
27
28     k3_r = v + 0.5 * dt * k2_v
29     k3_v = lorentz_force(v + 0.5 * dt * k2_v, B) / m
30
31     k4_r = v + dt * k3_v
32     k4_v = lorentz_force(v + dt * k3_v, B) / m
33
34     r_next = r + (dt / 6) * (k1_r + 2*k2_r + 2*k3_r + k4_r)
35     v_next = v + (dt / 6) * (k1_v + 2*k2_v + 2*k3_v + k4_v)
36
37     return r_next, v_next
38
39 def simulate_particle(index, r0, v0, positions_all, velocities_all):
40     r = r0
41     v = v0

```

```

42 positions = np.zeros((num_steps, 3))
43 velocities = np.zeros((num_steps, 3))
44
45 for i in range(num_steps):
46     positions[i] = r
47     velocities[i] = v
48
49     r, v = rk4_step(r, v, dt)
50
51 positions_all[index] = positions
52 velocities_all[index] = velocities
53
54 # Function to measure execution time for a given number of particles
55 def measure_execution_time(num_particles):
56     positions_all = np.zeros((num_particles, num_steps, 3))
57     velocities_all = np.zeros((num_particles, num_steps, 3))
58
59     # Generate unique initial conditions for each particle
60     initial_conditions = [(np.array([1.0, 0.0, 0.0]) + 0.1*np.random.randn(3),
61                             np.array([0.0, 1.0, 0.0]) + 0.1*np.random.randn(3))
62                             for _ in range(num_particles)]
63
64     # Measure start time
65     start_time = time.time()
66
67     # Create threads for each particle
68     threads = []
69     for i in range(num_particles):
70         r0, v0 = initial_conditions[i]
71         thread = threading.Thread(target=simulate_particle, args=(i, r0, v0,
72                             positions_all, velocities_all))
73         threads.append(thread)
74         thread.start()
75
76     # Wait for all threads to finish
77     for thread in threads:
78         thread.join()
79
80     # Measure end time
81     end_time = time.time()
82
83     # Calculate and return the execution time
84     execution_time = end_time - start_time
85     return execution_time
86
87 # Number of particles to test
88 num_particles_list = [1, 10, 100, 1000, 10000]

```



```
87 execution_times = []
88
89 # Measure execution time for each number of particles
90 for num_particles in num_particles_list:
91     execution_time = measure_execution_time(num_particles)
92     execution_times.append(execution_time)
93     print(f"Execution Time for {num_particles} particles: {execution_time:.2f}
94           seconds")
95
96 # Plotting the results
97 plt.figure(figsize=(10, 5))
98 plt.plot(num_particles_list, execution_times, marker='o')
99 plt.xlabel('Number of Particles')
100 plt.ylabel('Execution Time (seconds)')
101 plt.title('Execution Time vs Number of Particles')
102 plt.grid(True)
103 plt.show()
```