

Homework 2.3

CS 267

3/17/2024

Giovanni Battista Alteri, Jeffy Jeffy, Elizabeth Gilson

Group 36

Introduction:

The assignment is an extension of the particle simulation performed in Assignment 2.1. The previous assignment was a naive implementation of a looping algorithm which calculates the force on particles in a closed system. Based on the force and acceleration of each particle, the particle moves within the box, and interacts with each other. The number of particles is constant, and so is the volume of the system (defined by the variable “size”, given by the square root of the total number of particles). We have chosen a density sufficiently low so that we will have $O(n)$ interactions with n particles. This means that we could achieve the same level of accuracy of the naive code even if we consider only the interactions among particles that are sufficiently close to each other.

Parallel Implementation using CUDA

For this homework we aim to parallelize the code by targeting Perlmutter's A100 GPUs. First, we attempted to implement the OpenMP solution into CUDA, however the arrays that we used in our OpenMP solution were too large for the GPUs. After facing this limitation, we opted to refactor the code to leverage CUDA's parallel computing capabilities efficiently. We also tried utilizing linkedList as a data structure since they produce results. However, due to the ambiguity presented by using a linked list, it became much harder to implement and debug. Additionally, due to its dynamic nature, it is also harder to perform shared memory implementations of the program.

Parallelization with CUDA Kernels: CUDA allows executing kernels (functions that run on the GPU) in parallel on multiple threads. In this code, both `compute_forces_gpu` and `move_gpu` kernels are invoked with multiple threads, where each thread handles a different particle or a group of particles. CUDA enables concurrent execution of kernels across multiple threads, which is essential for handling large numbers of particles efficiently. This parallel execution significantly speeds up the computation compared to sequential execution on a CPU.

Memory Allocation Optimization: Memory allocation for `bin_ids`, `particle_ids`, and `bin_counts` is at first done on the GPU using `cudaMalloc`. This ensures that memory is directly accessible by GPU kernels, eliminating the need for data transfers between the CPU and GPU during simulation execution. However, utilizing this method was causing illegal access and resulting in a segmentation fault. Thus, we tried implementing thrust vectors to see if we could avoid running into this issue and prevent any segmentation fault caused by dynamic memory access.

```
// Initialize thrust vectors
thrust::device_vector<int> bin_counts(tot_bins);
thrust::device_vector<int> bin_ids(tot_bins);
thrust::device_vector<int> particle_ids(num_parts);
thrust::device_vector<int> particle_counter(tot_bins);
```

Thread and Block Management: The simulation distributes the workload across multiple CUDA threads and blocks. The number of threads per block is defined as NUM_THREADS, and the number of blocks is calculated based on the total number of particles (num_parts). This optimization ensures that the GPU's resources are effectively utilized, and the workload is evenly distributed across available CUDA cores.

```
#define NUM_THREADS 256
```

Atomic Operations: Atomic operations like atomicAdd are used to safely update shared memory locations accessed by multiple threads concurrently. In the assign_particles_to_bins kernel, atomic addition is employed to update bin_counts and assign particle indices to appropriate bins without race conditions.

```
// Function to update the number of particles per bin using atomicAdd
__global__ void update_bin_counts_gpu(particle_t* parts, int num_parts, int* bin_counts, int num_bins, double bin_size) {
    int tid = threadIdx.x + blockIdx.x * blockDim.x;
    if (tid >= num_parts)
        return;

    int part_x = parts[tid].x / bin_size;
    int part_y = parts[tid].y / bin_size;
    int index = part_x + part_y * num_bins;
    atomicAdd(&bin_counts[index], 1);
}
```



Thrust Library: The Thrust library is utilized for performing exclusive scan operations (`thrust::exclusive_scan`). Exclusive scan is applied to `bin_counts` after counting the number of particles per bin. This operation is essential for calculating the starting index of particles within each bin, enabling efficient access during force computation.

```
void simulate_one_step(particle_t* parts, int num_parts, double size) {
    // parts live in GPU memory
    // Rewrite this function

    // Initialize thrust vectors
    thrust::device_vector<int> bin_counts(tot_bins);
    thrust::device_vector<int> bin_ids(tot_bins);
    thrust::device_vector<int> particle_ids(num_parts);
    thrust::device_vector<int> particle_counter(tot_bins);

    // Update bin_counts and particle_ids
    update_bin_counts_gpu<<<blks, NUM_THREADS>>>>(parts, num_parts, bin_counts.data().get(), num_bins, bin_size);

    // Perform exclusive scan on bin_counts
    thrust::exclusive_scan(thrust::device, bin_counts.begin(), bin_counts.end(), bin_ids.begin());

    // Copy the result back to bin_counts
    // cudaMemcpy(bin_counts, bin_ids, num_bins * num_bins * sizeof(int), cudaMemcpyDeviceToDevice);

    // Update particle_ids
    update_particle_ids_gpu<<<blks, NUM_THREADS>>>>(parts, num_parts, particle_ids.data().get(), bin_ids.data().get(), bin_counts.data().get(), particle_counter.data().get(), num_bins, bin_size);

    // Compute forces
    compute_forces_gpu<<<blks, NUM_THREADS>>>>(parts, num_parts, particle_ids.data().get(), bin_ids.data().get(), bin_counts.data().get(), num_bins, bin_size);

    // Move particles
    move_gpu<<<blks, NUM_THREADS>>>>(parts, num_parts, size);
}
```

Kernel Fusion: Kernel fusion is employed by combining multiple computational tasks into a single kernel invocation. In the `simulate_one_step` function, multiple CUDA kernels (`count_particles_per_bin`, `thrust::exclusive_scan`, and `assign_particles_to_bins`) are executed sequentially within the same CUDA kernel invocation. This reduces kernel launch overhead and improves computational efficiency.

Parallelization of Particle Interactions: In the `compute_forces_gpu` kernel, each thread is responsible for computing the forces acting on a single particle due to interactions with its neighboring particles. By assigning a unique thread to handle each particle, the computation of forces between particles is parallelized. This means that multiple particles can be processed simultaneously by different threads on the GPU. For example, if there are N particles in the simulation, N threads will be launched, each responsible for computing the forces acting on a specific particle. This parallel approach enables the GPU to exploit its massive parallelism efficiently, as thousands of threads can execute concurrently on the GPU cores. As a result, the computation of forces between particles can be completed much faster compared to a sequential implementation, especially when dealing with a large number of particles. We also ensured that our `compute_forces_gpu` function is modified so that it can identify the index of each particle. Then loop over the 8 surrounding bins, and identify the indices of neighboring particles and their bins. Using this information, we iteratively go through each bin and compute force. Using CUDA to parallelize this process significantly speeds up the computation time.

```

__global__ void compute_forces_gpu(particle_t* particles, int num_parts, int* particle_ids, int* bin_ids, int* bin_counts, int num_bins, double bin_size) {
    // Get thread (particle) ID
    int tid = threadIdx.x + blockIdx.x * blockDim.x;
    if (tid >= num_parts)
        return;

    // set the acceleration to zero
    particles[tid].ax = particles[tid].ay = 0;
    // find the bin index for the particles
    int part_x = particles[tid].x / bin_size;
    int part_y = particles[tid].y / bin_size;

    // loop over the neighboring bins
    for (int i = -1; i <= 1; i++) {
        for (int j = -1; j <= 1; j++) {

            // get the neighboring bin id
            int neighbor_x = part_x + i;
            int neighbor_y = part_y + j;

            // iterate through all the particles in the neighboring bin
            if (neighbor_y >= 0 && neighbor_y < num_bins && neighbor_x >= 0 && neighbor_x < num_bins) {
                int index = neighbor_x + neighbor_y * num_bins;

                // print the index
                // printf("Index: %d\n", index);

                int neighbor_bin_start = bin_ids[index];
                int neighbor_bin_end = neighbor_bin_start + bin_counts[index];

                for (int k = neighbor_bin_start; k < neighbor_bin_end; k++) {
                    int neighbor_id = particle_ids[k];
                    if (neighbor_id != tid){
                        apply_force_gpu(particles[tid], particles[neighbor_id]);
                    }
                }
            }
        }
    }
}

```

Parallelization of Particle Movement: Similarly, in the `move_gpu` kernel, each thread is assigned to update the position of a single particle based on its velocity and acceleration. This process is also parallelized, allowing for efficient computation of particle movement across multiple particles simultaneously. Just like with particle interactions, if there are N particles in the simulation, N threads will be launched in the `move_gpu` kernel. Each thread will independently update the position of its assigned particle based on the particle's velocity and acceleration, without relying on the computation of other particles. This parallel approach enables the movement computation of multiple particles to be performed concurrently, leveraging the GPU's parallel architecture to achieve significant speedup. By parallelizing both particle interactions and movement, the simulation can fully utilize the computational power of the GPU, resulting in faster and more efficient simulations, especially when dealing with large-scale particle systems. We did not change the `move_gpu` function.

```

__global__ void move_gpu(particle_t* particles, int num_parts, double size) {

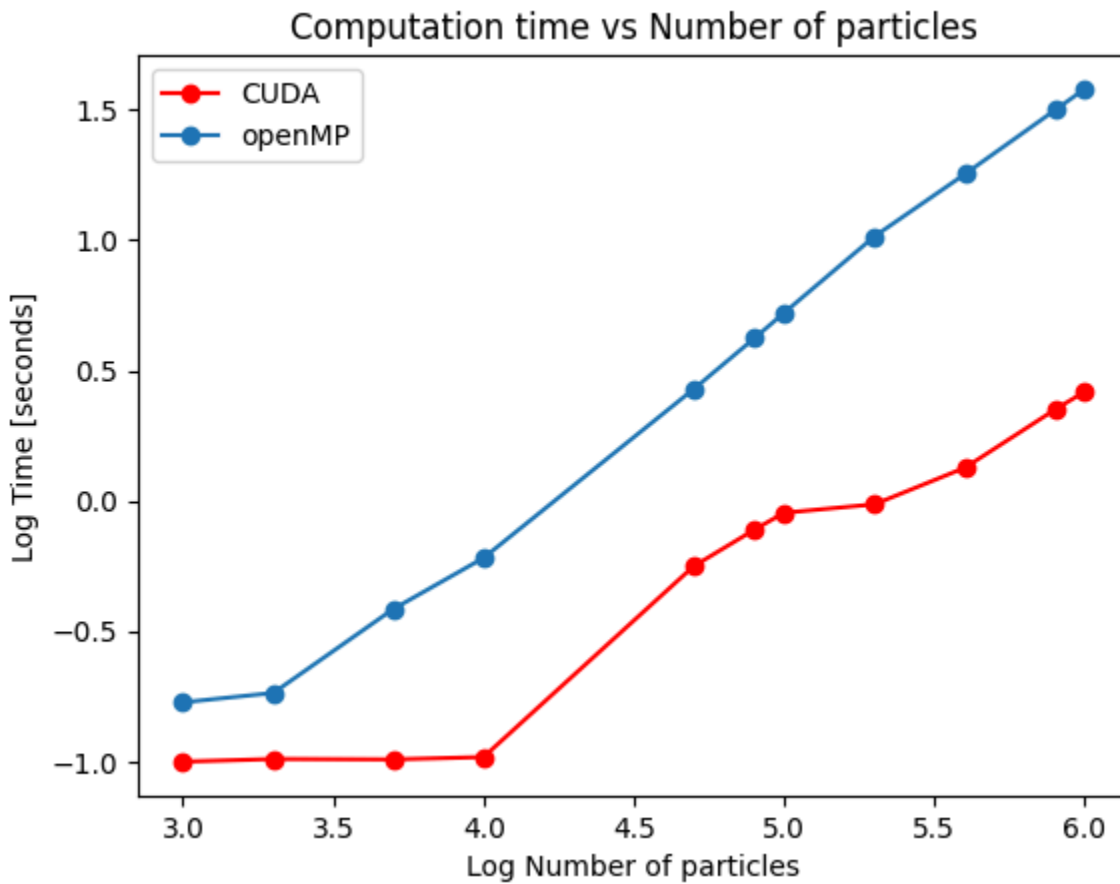
    // Get thread (particle) ID
    int tid = threadIdx.x + blockIdx.x * blockDim.x;
    if (tid >= num_parts)
        return;

    particle_t* p = &particles[tid];
    //
    //  slightly simplified Velocity Verlet integration
    //  conserves energy better than explicit Euler method
    //
    p->vx += p->ax * dt;
    p->vy += p->ay * dt;
    p->x += p->vx * dt;
    p->y += p->vy * dt;

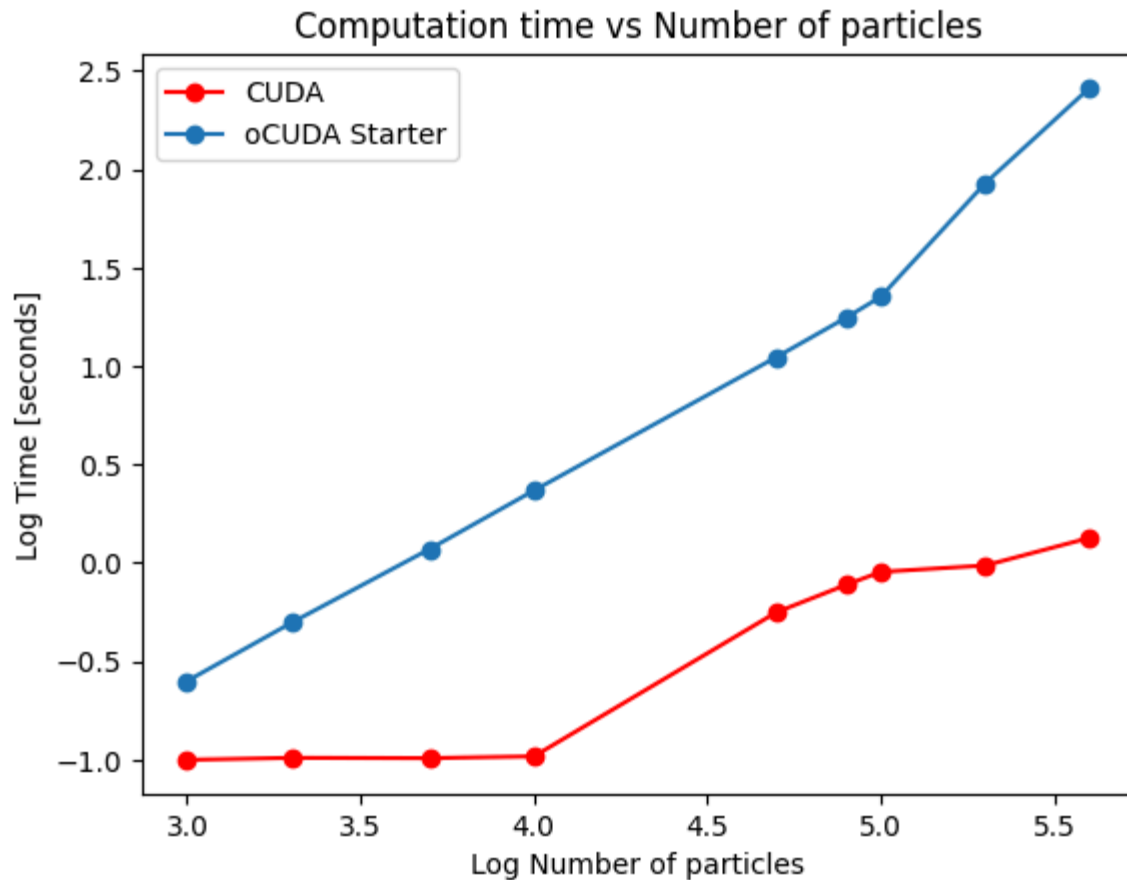
    //
    //  bounce from walls
    //
    while (p->x < 0 || p->x > size) {
        p->x = p->x < 0 ? -(p->x) : 2 * size - p->x;
        p->vx = -(p->vx);
    }
    while (p->y < 0 || p->y > size) {
        p->y = p->y < 0 ? -(p->y) : 2 * size - p->y;
        p->vy = -(p->vy);
    }
}

```

Results



Above is a graph showing the log of the time in seconds vs the log of the number of particles for our CUDA and OpenMP implementations of the particle simulation. As noted in the graph, the CUDA gpu code is faster than the openMP code for all instances presented. It is interesting that the speed is comparable for a small number of particles but there is a significant difference in computation time as the number of particles increases. This could also be justified by overhead costs that accumulate with a smaller number of particles which is probably not “worth” the number of kernels and threads that are released to perform the operations; however, with a large number of particles, this overhead is balanced by the amount of “work” that needs to be performed. It is interesting to note the fluctuations in the cuda computation time as observed around $\log(5.0) - \log(5.5)$ particles, which is present in the graph below as well.



In this graph we have the log of time in seconds vs the log in the number of particles for both the original naive CUDA implementation and our implementation. As you can see, for any number of particles, our code is faster than the original starter code provided. This is especially after we accounted for the neighboring regions while computing forces.

Time Components

To study the different components of the total computational time, we used a modified version of the job-gpu file that employs NVIDIA Nsight (from the CUDA Toolkits). This is the code used:

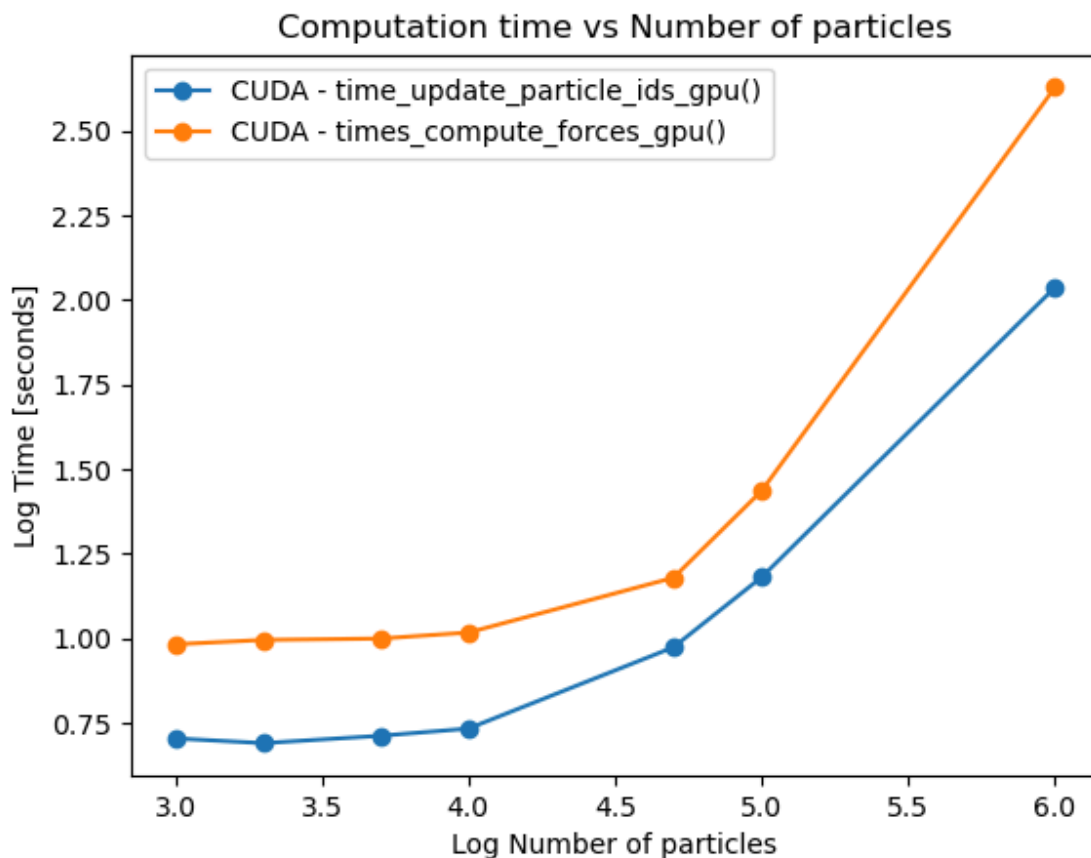
```
#!/bin/bash
#SBATCH -N 1
#SBATCH -C gpu
#SBATCH -A mp309
#SBATCH -t 00:30:00
#SBATCH --gres=gpu:1

srun --ntasks-per-node 1 dcgmi profile --pause
```

```
size=(2000)
output_file="profile${size}.ncu-rep"
srun ncu --kernel-id :::1 -o "${output_file}" --call-stack ./gpu -s 1 -n
"${size}"
srun --ntasks-per-node 1 dcgmi profile --resume
```

The toolkit generates an output file that can be opened using a Desktop version of NVIDIA NSight Compute; this is a typical result:

We can visualize the total computational time spent in the CUDA section of the code, and we can highlight the time required by each different function speeded-up using CUDA. For example, in this result we can see that the `compute_forces_gpu` method took the longest amount of time and memory. This trend is repeated for all the experiments we tried, using different numbers of particles (1000,2000, 5000, 10000, 50000, 100000, 1000000).



As you can see the compute forces function took the most amount of time followed by the update particle ids function. The difference between them appears to increase as the number of particles gets above 100,000, however they both follow the same general pattern of a logarithmic curve.

Can we improve this code?

There is a lot of room for improvement within this code. Some of the considerations that we wanted to include but were not able to successfully implement are as follows:

Implementation of linked list: We tried creating a ListNode struct where each node holds the index of a particle in the 'particles_t*' particles array. The general function of the linked list is very similar to the function of the array in the code provided above. In this implementation of the code, once the linked list are initialized, the kernel inserts particles indices into the appropriate linked list based on their bin ids

Addition of padding: In a grid-based approach for spatial partitioning, padding is appropriate to handle particles that are near the boundary of a bin. This is important because particles near the edge of a bin may actually overlap with neighboring bins. Without padding, these particles would be assigned to incorrect neighboring bins.

We believe implementing these improvements can strongly improve our code and the speed of the particle simulation.

Contributions

Jeffy- helped work on the code, optimize the code, helped with data collection and helped with the write-up

Giovanni- helped work with the timings, on editing and testing the code, creating the graphs and helped with the write-up

Elizabeth- helped with the write-up, on editing the code, and with the data for timings and nsight timings