# Work Assignment - Phase 3

Parallel Computing

José Miguel Ferreira Barbosa
PG52689
Braga, Portugal

Tiago Adriano Gomes Moreira
PG52704
Braga, Portugal

*Abstract*—**This document is the written report of the third and final phase of work assignments of the curricular unit, Parallel Computing. The focus of this phase was how to design and implement an efficient parallel version of the case study.**

## I. INTRODUCTION

The final phase of this work assignment aims to understand how to design and implement an efficient parallel version of the case study, eventually incorporating accelerators, with the main goal still being of reducing the execution time. In this report, we will explore the adopted strategies, the challenges encountered, and the proposed solutions to optimize the system's performance, emphasizing the crucial importance of parallel approaches in the pursuit of computational efficiency. It's worth noting that because of the issues we had in the two previous phases regarding results, we choose to pick up the original code and work from there again.

## II. THE APPLICATION HOT-SPOTS

As we figured out in the previous phases of this work assignment, we have two code blocks that call our hot-spots, those being:

- **Potential**
- **computeAccelarations**

It is imperative to acknowledge that both these functions were identified as hot-spots not only in the original code, that we will pick up again when working on this final phase, but also in the sequential version we had made and worked on during phase 2. So our main efforts will once again be directed towards optimizing and enhancing the performance of these components.

## III. THE CHOOSEN PATH

In this final work assignment, we had the freedom to choose one of three paths to improve the execution. After a discussing upon ourselves, we decided to look into the approach where we would be challenged with designing and implementing a new version for accelerators (GPU with CUDA). We choose this one because we find more suitable than the other two, as it also was the focus of the last classes of the semester of this curricular unit. We will then, work with CUDA, on the functions **Potential** and **computeAccelarations** specifically.

## IV. CHANGES AND SEQUENTIAL CODE

As we mentioned in the introduction, we started from scratch again. Because of that, the code is different from the one we delivered in the versions prior to this one. First of all, to enhance our code's efficiency, we've shifted to Struct of Arrays (SoA). We made this change because it offers advantages like improved memory access patterns and better alignment with GPU parallelism.

We kept all the functions that are not the two hot spots, pretty much the same as they were in the original code (of course we adapted them to fit SoA), since they are not going to be worked on using CUDA. Also, we changed them in the first phase, but since we encountered problems with the simulation results we also felt like making more changes than required gave us a bigger chance of getting wrong results again, so we focused solely on **Potential** and **computeAccelarations** and worked carefully to make sure the same thing didn't happen again. We added a function called **power** (fig 11) with the sole purpose of replacing the **pow**, making the code a bit more efficient in the process. We only did this in natural powers, because non-natural powers are more complex to calculate, and it could lead to incorrect results. The function **Potential**, follows the same logic we applied in phase 1, when we worked on it, where we made slight changes in the operations made when calculating the potential, to improve the execution time, removing some time-consuming operations in the process, for example divisions and the call of the function **srqt**. When it comes to **computeAccelarations** we also followed the same logic we did in phase 1, where we removed of for cycles and also made slight changes to the operations made, for example we eliminated some redundant operations, and we introduced some intermediate variables to optimize the calculations and reduce repetition.

## V. PHASE II CORRECTIONS

As we talked about, we started from scratch and made changes to the sequential code we had worked on previously. Since we changed the code we are working on and are working with CUDA for this final phase, we didn't pick up OpenMP again. However, we were informed that we had a data race in our last assignment, that went by unnoticed, and we identified it and decided to correct it and that correction can be seen in fig 1. To resolve this issue, we added a barrier like shown in the correction that will also guarantee that all the acceleration

Fig. 1. Correction of data race issue in phase 2


Fig. 2. Side by Side comparison of the CUDA and Sequential version of Potential

arrays are initialized correctly before the calculation starts, and we believe this small addiction solves our problem.

## VI. APPROACH WITH CUDA

After finishing and making sure that the simulation results from the sequential version were correct, we went ahead and started our CUDA implementation.

In regards to the CUDA adaptation of the function **Potential** (fig 2) and to really explain what we did, let's just quickly get into what it does. The potential function calculates the potential energy of each particle in the system and sums them all, thus obtaining the total potential energy of the particle system in question. In the sequential version, the function performs these two steps almost simultaneously, where the accumulated value of the system's potential energy is stored in the variable Pot. So we came to the conclusion, that this function has two distinct parts: a mapping, corresponding to the stage of calculating the potential energy of all particles, and a reduction, corresponding to the sum of all these potential energies to obtain the potential energy of the system. It's important to note that the calculation of potential energies for different particles is embarrassingly parallel, as these calculations have no dependencies among them. The same can be said for the reduction phase, which aims to obtain the total potential energy of the system in a parallel fashion. With this understanding of the function in mind, we decided to divide this function into two different kernels, one for each component of this function: **Potential_Map**, which will fill an array with the potential energy values of each particle in the system, and **Potential_Reduce**, which will sum all the values calculated by **Potential_Map**, cumulatively placing the total sum value in the first position of the array. Additionally, we also launched a device called **power_cuda** (fig 12), which will solely perform the exponentiation of various values.

The **Potential_Map** Kernel will receive arrays with the positions of each particle (drx, dry, and drz), as well as their size (N) and sigma. It is worth mentioning that we are creating one thread for each particle in the system, that is, each thread will fill in the Pot array with the potential energy of its corresponding particle at the position corresponding to the number of that particle. For example, the potential energy of particle 0 should be located at position 0 in the Pot array. Each thread needs a specific identifier, 'i'. This way, it places the calculated potential energy in the correct location in the Pot array. This identifier should be calculated as follows: the identifier of each thread (ranging from 0 to the number of threads defined for each block), added to the identifier of each block (ranging from 0 to the number of blocks we decided on) times its dimension (which is equal to the total number of threads in the block). This formula allows us to correctly identify each thread independently of its block because we need to add the overhead given by the product between the block it is in and its dimension to its identifier within the block. In other words, we are positioning ourselves in the specific location.

example:

with 256 threads per block, thread 7 in block 8, 'i' = 2055. It is worth noting that the calculation itself has not changed since our changes in phase 1, that we mentioned briefly in the previous section, so we will write the value of the potential energy of the respective particle into the Pot array at the correct position, as mentioned earlier. The other auxiliary variables (r3, r2, r6, pot, sig6), being local to the kernel, will be private to each thread.

**Potential_Reduce** will receive the array of previously calculated potential energies, as well as the number of elements that should undergo an initial reduction. The first reduction involves reducing the number of elements from 5000 to a number where we can have one thread per particle in a single block, in this case, 256 elements. In other words, this initial reduction involves going from 5000 elements to 256 elements. This necessity arises from how the reduction operation functions. The **Potential_Reduce** will perform accumulated sums of different elements. However, to do this, we

must ensure synchronization at the thread level; otherwise, we would encounter data races. For this purpose, we are using the synchronization primitive **__syncthreads()**, which ensures that within the same block, threads cannot overtake each other. In other words, thread 1 cannot finish before thread 2. However, given the maximum block size of 1024 threads, to have one thread per element, we proceed with the initial reduction described above. Thus, with 256 threads, we can handle 5000 elements. It is worth noting that this initial **__syncthreads()** ensures that all threads write to the shared array before moving on to the next phase, it functions as a barrier in other words. Next, we will enter the loop of the second reduction, which is essentially a sequence of cascade sums, where in each iteration we will successively need half of the threads. For each iteration, a thread will sum two elements. For instance, thread 0 adds values 0 and 1, thread 1 adds values 2 and 3. In other words, with each iteration, we will halve the number of necessary threads until we only have to sum two values. At that point, thread 0 will add them and store the final value in the first position of the result array, Pot_r.

A small additional point to mention is that in this function, we leverage the concept of shared memory. The use of shared memory (__shared__ double partial_sum[256];) is implemented to minimize the number of accesses to global memory, which is slower compared to shared memory. Each thread block maintains a local copy of the partial_sum variable, which is utilized to accumulate the partial sum of the Pot vector elements for each thread.

These two kernels are launched by the **PotentialPrepare** function, which will be responsible for sending the values of the three position arrays (rx, ry, rz) to the GPU. It is worth noting that the array where the potential energies of each particle will be stored has a size of 5120, and not the expected 5000. This increase in size is done so that in the reduction phase, we have the correct calculations regarding the number of values to sum in the initial reduction phase which, in this case is, 20 values at a time. These additional 120 positions will have no impact since this array is initialized with 0, meaning their presence will not affect the calculation of the system's potential energy in any way.
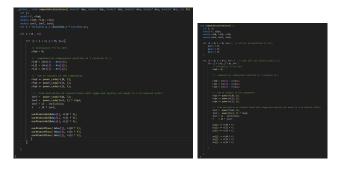


Fig. 3. Side by Side comparison of the CUDA and Sequential version of computeAccelarations

Now let's get in to the other hot spot function **computeAccelerations** (fig 3). In our new version, the **computeAccel-** **erationsPrepare** function is responsible for launching the kernel, initializing all accelerations to 0, and copying these values to the GPU. The **computeAccelerations** function in the sequential version calculates the accelerations of each particle from scratch every time it is called. For the CUDA version, the procedure is similar. To perform this task, we create one thread per particle and launch a single kernel that will perform this calculation for each particle. Just like in the potential function, all threads need a universal identifier calculated. The calculation itself is not altered from the sequential version to this CUDA version; what changes is the complexity of the function, decreasing from $O(N^2)$ to $O(N)$. In this function, we use the **ourAtomicAdd** and **ourAtomicMinus** (fig 7) functions from the device, which will ensure the absence of data races. This is the solution we came up with, and it works making sure that everything goes well and the one that is in the final code sent. However, we have another possible approach that we thought of, that we didn't implement, but we think is worth mentioning.

In this alternative approach (fig 8), instead of using the **ourAtomicAdd** and **ourAtomicMinus** devices, we employ the synchronization primitive **__syncthreads()**. As mentioned earlier, it serves as a barrier, preventing threads from the same block from overlapping. Each thread will have a different 'i', so we wouldn't have problems with data races. The value of 'j' depends on 'i', where 'j = i+1', and within the same block, thanks to **__syncthreads()**, no thread would surpass another. Therefore, for a single block, we wouldn't encounter any issues. However, as we are using more than one block, to achieve this, we would need to use an inter-block synchronization mechanism, ensuring that all threads from all blocks cannot surpass each other (it would be a kind of general barrier). We attempted to use **cooperative_groups**, but it didn't work.

Nevertheless, we believe that this approach would be very efficient, but since we couldn't guarantee the absence of "data races" we went with the strategy mentioned above and that gives us the assurances we desired.

## VII. RESULTS, ANALYSIS AND DISCUSSION

In this section we will talk about the results we obtained in the phase of the work assignment and also talk about some tests we conducted and some metrics we think are valuable analyzing.

In this work assignment, we design and ran two different tests to see how our CUDA solution worked under some different circumstances. The tests were:

1) Changing the value of N in both the sequential and CUDA versions
2) Changing the number of threads per block in the CUDA version
3) Comparing the time and percentage of time for each hotspot function in both versions

## A. Test 1

For the first test we conducted had a very simple principle. We wanted to see how the execution time changed in both the CUDA and the C++ versions when we changed the value of N, or in other words, we changed the number of particles. We tested 6 different numbers: 1000, 2000, 5000, 7500, 10000 and 15000 for both the sequential version(fig 5) and the CUDA version (fig 6).

We can see the results of the execution time of both versions in the graph (fig 4) and TABLE I

### TABLE I
#### EXECUTION TIME FOR BOTH VERSIONS

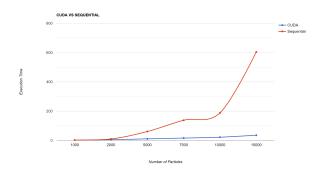| N = | 1000 | 2000 | 5000 | 7500 | 10000 | 15000 |
|---|---|---|---|---|---|---|
| CUDA | 2.0309s | 4.0700s | 10.2138s | 15.3932s | 21.2099s | 35.7484s |
| Sequential | 2.43169s | 9.71843s | 59.8467s | 137.473s | 187.933s | 604.508s |



Fig. 4. Graph of execution time of CUDA and Sequential versions with different number of particles

By at looking and analyzing both of these we can see very easily, that as the value of the number of particles increases so does the difference in execution time between both versions. But we can get the full picture by calculating the speedup. Remembering the concept of speedup, it is a metric that indicates the performance gain achieved by utilizing parallelism compared to a sequential implementation. So will we use to really grasp, the difference between both versions.

In TABLE II, we can take a look at results of the calculations made to get this metric for each different value of N. Speedup is obtained by dividing the sequential time for parallel time, we can see the formula be applied in exemple below.

$$N = 5000 : 59.8467/10.2138 = 5.8594$$

### TABLE II
#### SPEEDUP FOR EACH VALUE OF N

| N = | 1000 | 2000 | 5000 | 7500 | 10000 | 15000 |
|---|---|---|---|---|---|---|
| SpeedUP = | 1.1973 | 2.2052 | 5.8594 | 8.9308 | 8.8606 | 16.9101 |

We can notice once again, by analyzing this metric, that the difference rises with the number of particles (with the exception of 7500 and 10000). If talk about complexity, we can assess that the sequential version was a complexity of about $O(N^2)$for both the **Potential** and the **computeAccelarations**

functions, while in the CUDA version the complexity of doing the same operations lowers to about $O(N)$, and so, with the input set at 'ar', this means each function is called about 200 times. So that this means that while in the CUDA version, the rise is pretty much linear, in the sequential version it's quadratic.

Let's talk about scalability analysis, where we examine how the performance of a system changes as the problem size or workload increases and what we can conclude with this test of changing the value of N. In the CUDA, execution time increases gradually as N increases, indicating a relatively good scalability. While in the sequential version, the execution time increases significantly as N increases, demonstrating poor scalability. The rate of increase is much steeper compared to the CUDA version. This shows us that the CUDA version is able to handle the workload more efficiently, whereas the sequential version struggles and becomes increasingly inefficient. Furthermore, upon observing the values presented for the number of particles, specifically 7500 and 10000, it is evident that they are quite close, with a slight decrease in the metric. Upon closer analysis, it becomes apparent that this decrease is part of a variation considered normal. This phenomenon typically precedes the next significant leap in Speedup, which is expected to occur between 10000 particles and 150000 particles.

The analysis consistently reveals that the CUDA approach demonstrates superior efficiency, particularly when dealing with considerably large values, aligning well with our intended goals. This finding underlines the scalability of our solution, highlighting that Speedup rises with the rise of the number of particles (N).

## B. Test 2

The second test was changing the number of threads used in our CUDA version, with the number of particles set at 5000. We did this by changing the value of the variable 'TB_SIZE' in both **computeAccelerationsPrepare** and **PotentialPrepare**, which is the variable that defines the number of threads per block in our solution. The values we tested were: 64, 128, 256, 512 and 1024. 256, was the number that was delivered in the final code. It's important to mention that the number of threads is directly influences the number of blocks.

```
Example
const int TB_SIZE = 256;
int GRID_SIZE =
ceil((N * 1.0) / (TB_SIZE * 1.0));
// GRID_SIZE = 20
computeAccelerations
<<<GRID_SIZE, TB_SIZE>>>
```

Another quick thing worth mentioning is that in CUDA, we should always run the numbers of thereds as a multiple of 32 since GPU hardware is designed to perform operations in parallel in warps of 32 threads. If we have a total number of threads that is not a multiple of 32, the GPU will still execute the threads, but there will be a number of "ghost threads" that

will be filled in to complete a warp. These additional threads may not be efficiently utilized and can result in some resource wastage. It's something that we have to be careful about when managing threads.

| Nº of Blocks | 78 | 39 | 20 | 10 | 5 |
|---|---|---|---|---|---|
| Threads per Block | 64 | 128 | 256 | 512 | 1024 |
| Execution Time | 10.07s | 10.15s | 10.21s | 10.45s | 12.21s |

Analyzing Table III we can see that the execution time is lower in the cases with more blocks and fewer threads, but they all remain around the 10 seconds mark, with the notable exception being the last example, where we run 1024 threads per block, with 5 blocks where it rises to 12 seconds. 1024 threads is the maximum amount we can run in a single block. We believe that the reason this case is such an outlier in this test is that the huge amount of threads in this block causes a slight overhead when trying to synchronize them all. In regards to the other values, as we just noted as well, they aren't that far apart from each other. We can attribute the fact the execution time is lower with more blocks to a couple of potential reasons. With a lower number of threads per block, it is possible that the GPU is using available resources more efficiently. This may occur because the GPU can allocate and manage blocks more effectively, avoiding underutilization of its resources. Another reason that may cause this, is that a lower number of threads per block can lead to better load balancing, ensuring that each thread has a significant amount of work to perform which can be important to avoid situations where some threads finish quickly while others are still processing, which could result in delays, and we believe these two reasons are why we obtained the results we did in regards to this specific test.

*C. Test 3*

Our goal with this third test is to do a more detailed comparison of the sequential version of the code with our CUDA solution in order to have a more in depth comprehension of how the implementations we made really affects performance. In this test, we compare the time taken by each hotspot function, taking into consideration as well the percentage of time taken regarding the total execution time. They are the only two relevant functions to analyze since, they are the only ones that were changed and also are the ones that have a relevant impact in the execution time of the program, in both versions. To do this test we used the command *'gprof'* for the sequential version (fig 17) and for the CUDA version we analyzed the profiling results, while in TABLE IV we have the values for time taken by each hotspot function, while in TABLE V we have the percentage of time each function in regards to the total time ((please note that in the CUDA version the value for 'Potential' includes both **Potential_Map** and **Potencial_Reduce**).

As we have discussed before, the total time between both versions is very different, with the parallel version being

| | Potential | computeAccelarations | Total Time |
|---|---|---|---|
| Sequential Version | 30.61s | 29.24s | 59.86s |
| CUDA Version | 2.4387s | 7.75105s | 10.2138s |

| | Potential | computeAccelarations | Total Time |
|---|---|---|---|
| Sequential Version | 51.16% | 48.87% | 59.86s |
| CUDA Version | 23.89% | 75.89% | 10.2138s |

almost 6 times faster than the sequential one. However, this is not the point of this test and taking a look at both tables there is one thing that stands out immediately, that being the fact that while in the sequential version both functions take almost similar times (30.61s and 29.24s) and because of that the percentage of both is almost 50/50, in the CUDA version, there is a very big difference. The combination of both Potential functions that run on GPU (**Potential_Map** and **Potencial_Reduce**) only takes 23.89% of the total execution time, while **computeAccelarations** takes 75.89% of the time.

This can be explained by how both functions were implemented. Our implementation of the **computeAccelarations** function, works well and does exactly what it's required, and guarantees the absence of data races, but like we discussed in the previous section, isn't the most efficient solution possible, since it uses **ourAtomicAdd** and **ourAtomicMinus** devices that, while they are a good and safe way to guarantee functionality, it comes at a performance cost. By contrast, our **Potential** solution, implemented via **Potential_Map** and **Potencial_Reduce** (and of course **PotencialPrepare**), is very efficient. This is a lot due to the fact, that in both functions that run on the GPU, we don't really have any "critical zones", because it's extremely parallel, and adding the fact that we use shared memory while in the reduce function, it's going to come at a very good performance gain, explaining why it gives us such a low execution time.

## VIII. CONCLUSION

To conclude this phase 3 and final report in regards to this work assignment, we feel like the work we did on this project as a whole, helped us understand a lot better how parallel programming works and how we can improve it and that we should be aware of when working on it. We faced a lot of difficulties during the whole work assignment, but we believe that these helped us see first hand what to avoid doing and what mistakes can happen. We feel that this final phase has helped us to better consolidate our knowledge in the realm of CUDA and to gain a deeper understanding of how this aspect of programming operates.

ATTACHMENTS



Fig. 5. Results of C++ version with different Number of Particles



Fig. 6. Results of CUDA version with different Number of Particles

```
__device__ double ourAtomicAdd(double* address, double val)
{
    unsigned long long int* address_as_ull =
                                (unsigned long long int*)address;
    unsigned long long int old = *address_as_ull, assumed;

    do {
        assumed = old;
        old = atomicCAS(address_as_ull, assumed,
                        __double_as_longlong(val +
                        __longlong_as_double(assumed)));

    // Note: uses integer comparison to avoid hang in case of NaN (since NaN != NaN)
    } while (assumed != old);

    return __longlong_as_double(old);
}

__device__ double ourAtomicMinus(double* address, double val){
    unsigned long long int* address_as_ull =
                                (unsigned long long int*)address;
    unsigned long long int old = *address_as_ull, assumed;

    do {
        assumed = old;
        old = atomicCAS(address_as_ull, assumed,
                        __double_as_longlong(__longlong_as_double(assumed) - val));

    // Note: uses integer comparison to avoid hang in case of NaN (since NaN != NaN)
    } while (assumed != old);

    return __longlong_as_double(old);
}
```

Fig. 7. ourAtomicAdd e ourAtomicMinus

```
__global__ void computeAccelerations( double* dax, double* day, double* daz, double* drx, double* dry, double* drz, int N){
    int j;
    double f, rSqd;
    double rij0, rij1, rij2;
    double inv3, inv7, inv2;
    int i = threadIdx.x + (blockIdx.x * blockDim.x);

    if( i < N - 1){

        for (j = i + 1; j < N; j++){

            // initialize r^2 to zero
            rSqd = 0;

            // component-by-componenet position of i relative to j
            rij0 = (drx[i] - drx[j]);
            rij1 = (dry[i] - dry[j]);
            rij2 = (drz[i] - drz[j]);

            // sum of squares of the components
            rSqd += power_cuda(rij0, 2);
            rSqd += power_cuda(rij1, 2);
            rSqd += power_cuda(rij2, 2);

            // From derivative of lennard-jones with sigma and epsilon set equal to 1 in natural units!
            inv3 = power_cuda(rSqd, 3);
            inv2 = power_cuda(inv3, 2) * rSqd;
            inv7 = (2 - inv3)/inv2;
            f    = 24 * inv7;

            dax[i] += rij0 * f;
            day[i] += rij1 * f;
            daz[i] += rij2 * f;

            __syncthreads();

            dax[j] -= rij0 * f;
            day[j] -= rij1 * f;
            daz[j] -= rij2 * f;
        }
    }
}
```

Fig. 8. computeAccelerations CUDA Alternative

```
void computeAccelerationsPrepare(int N){
    double temp[N] = {0};
    size_t bytes = N * sizeof(double);

    double * drx;
    double * dry;
    double * drz;

    double * dax;
    double * day;
    double * daz;

    cudaMalloc(&dax, bytes);
    cudaMalloc(&day, bytes);
    cudaMalloc(&daz, bytes);
    cudaMemcpy(dax, temp, bytes, cudaMemcpyHostToDevice);
    cudaMemcpy(day, temp, bytes, cudaMemcpyHostToDevice);
    cudaMemcpy(daz, temp, bytes, cudaMemcpyHostToDevice);

    cudaMalloc(&drx, bytes);
    cudaMalloc(&dry, bytes);
    cudaMalloc(&drz, bytes);
    cudaMemcpy(drx, rx, bytes, cudaMemcpyHostToDevice);
    cudaMemcpy(dry, ry, bytes, cudaMemcpyHostToDevice);
    cudaMemcpy(drz, rz, bytes, cudaMemcpyHostToDevice);

    const int TB_SIZE = 256;
    int GRID_SIZE = ceil((N * 1.0) / (TB_SIZE * 1.0));

    computeAccelerations<<<GRID_SIZE, TB_SIZE>>>(dax,day,daz, drx,dry,drz, N);

    cudaMemcpy(ax, dax, bytes, cudaMemcpyDeviceToHost);
    cudaMemcpy(ay, day, bytes, cudaMemcpyDeviceToHost);
    cudaMemcpy(az, daz, bytes, cudaMemcpyDeviceToHost);

    cudaFree(drx);
    cudaFree(dry);
    cudaFree(drz);
    cudaFree(dax);
    cudaFree(day);
    cudaFree(daz);
}
```

Fig. 9. computeAccelerationsPrepare

Fig. 10. Potencial Prepare

```
double power(double base, int expoent){
    double res = 1.;
    int i = 0;
    while (i < expoent){
        res *= base;
        i++;
    }

    return res;
}
```

Fig. 11. Power

```
__device__ double power_cuda(double base, int expoent) {
    double res = 1.0;
    int i = 0;
    while (i < expoent) {
        res *= base;
        i++;
    }
    return res;
}
```

Fig. 12. Power Cuda

Fig. 13. Results of changing to 64 threads

```
AVERAGE TEMPERATURE (K):              131.4550641503

AVERAGE PRESSURE  (Pa):               131001228.4507598877

PV/nT (J * mol^-1 K^-1):               28.4727885201

PERCENT ERROR of pV/nT AND GAS CONSTANT: 242.4490490666

THE COMPRESSIBILITY (unitless):         3.4244904654

TOTAL VOLUME (m^3):                    2.3721986292e-25

NUMBER OF PARTICLES (unitless):      5000
==20511== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   75.67%  7.68723s       202  38.056ms  37.995ms  38.117ms  computeAccelerations(double*, double*, doub
 e*, double*, double*, double*, int)
                   24.05%  2.44374s       201  12.158ms  11.973ms  12.259ms  Potential_Map(double*, double*, double*, do
 ble*, int, double)
                    0.18%  17.908ms      2217  8.0770us  7.9680us  17.472us  [CUDA memcpy HtoD]
                    0.05%  5.3040ms       807  6.5720us  6.0480us  15.712us  [CUDA memcpy DtoH]
                    0.05%  5.2129ms       201  25.934us  25.280us  26.721us  Potential_Reduce(double*, double*, int)
      API calls:   96.91%  10.2171s      3024  3.3787ms  17.023us  38.274ms  cudaMemcpy
                    2.57%  270.62ms      2217  122.06us  2.5510us  215.86ms  cudaMalloc
                    0.44%  46.702ms      2217  21.065us  2.5220us  151.15us  cudaFree
                    0.07%  7.1928ms       604  11.908us  6.7760us  91.740us  cudaLaunchKernel
                    0.01%  586.39us       101  5.8050us    282ns  373.35us  cuDeviceGetAttribute
                    0.00%  357.16us         1  357.16us  357.16us  357.16us  cuDeviceTotalMem
                    0.00%  31.901us         1  31.901us  31.901us  31.901us  cuDeviceGetName
                    0.00%  14.727us         1  14.727us  14.727us  14.727us  cuDeviceGetPCIBusId
                    0.00%  2.9350us         3    978ns    580ns  1.6430us  cuDeviceGetCount
                    0.00%  1.7720us         2    886ns    386ns  1.3860us  cuDeviceGet
                    0.00%    860ns         1    860ns    860ns    860ns  cuDeviceGetUuid
```

Fig. 14.  Results of changing to 128 threads

```
pg52689@search7edu2:~/testescuda
AVERAGE TEMPERATURE (K):              131.4550641503

AVERAGE PRESSURE  (Pa):               131001228.4507598877

PV/nT (J * mol^-1 K^-1):               28.4727885201

PERCENT ERROR of pV/nT AND GAS CONSTANT: 242.4490490666

THE COMPRESSIBILITY (unitless):         3.4244904654

TOTAL VOLUME (m^3):                    2.3721986292e-25

NUMBER OF PARTICLES (unitless):      5000
==20584== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   76.34%  7.98307s       202  39.520ms  39.435ms  39.603ms  computeAccelerations(double*, double*, doub
 e*, double*, double*, double*, int)
                   23.40%  2.44692s       201  12.174ms  11.981ms  12.299ms  Potential_Map(double*, double*, double*, do
 ble*, int, double)
                    0.17%  17.927ms      2217  8.0860us  7.9680us  17.376us  [CUDA memcpy HtoD]
                    0.05%  5.2960ms       807  6.5620us  6.0480us  7.5520us  [CUDA memcpy DtoH]
                    0.04%  4.1138ms       201  20.466us  19.584us  21.408us  Potential_Reduce(double*, double*, int)
      API calls:   96.99%  10.5148s      3024  3.4771ms  16.730us  39.753ms  cudaMemcpy
                    2.50%  271.41ms      2217  122.42us  2.5920us  216.59ms  cudaMalloc
                    0.43%  46.814ms      2217  21.115us  2.5400us  154.29us  cudaFree
                    0.07%  7.3053ms       604  12.094us  6.8750us  87.485us  cudaLaunchKernel
                    0.01%  570.57us         1  570.57us  570.57us  570.57us  cuDeviceTotalMem
                    0.00%  366.41us       101  3.6270us    390ns  142.18us  cuDeviceGetAttribute
                    0.00%  219.97us         1  219.97us  219.97us  219.97us  cuDeviceGetName
                    0.00%  14.886us         1  14.886us  14.886us  14.886us  cuDeviceGetPCIBusId
                    0.00%  3.0940us         3  1.0310us    410ns  1.7920us  cuDeviceGetCount
                    0.00%  1.7010us         2    850ns    434ns  1.2670us  cuDeviceGet
                    0.00%    983ns         1    983ns    983ns    983ns  cuDeviceGetUuid
```

Fig. 15.  Results of changing to 512threads

```
'cp_average.txt':

AVERAGE TEMPERATURE (K):              131.4550641503

AVERAGE PRESSURE  (Pa):               131001228.4507599026

PV/nT (J * mol^-1 K^-1):               28.4727885201

PERCENT ERROR of pV/nT AND GAS CONSTANT: 242.4490490666

THE COMPRESSIBILITY (unitless):         3.4244904654

TOTAL VOLUME (m^3):                    2.3721986292e-25

NUMBER OF PARTICLES (unitless):      5000
Profiling application: ./bin/stencil
==20660== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   79.70%  9.73244s       202  48.180ms  48.084ms  48.266ms  computeAccelerations(double*, double*, doub
 le*, double*, double*, double*, int)
                   20.07%  2.45087s       201  12.193ms  12.075ms  12.254ms  Potential_Map(double*, double*, double*, do
 uble*, int, double)
                    0.15%  17.904ms      2217  8.0750us  7.9680us  17.760us  [CUDA memcpy HtoD]
                    0.04%  5.4096ms       807  6.7030us  6.0490us  7.8730us  [CUDA memcpy DtoH]
                    0.04%  4.4058ms       201  21.919us  20.608us  22.720us  Potential_Reduce(double*, double*, int)
      API calls:   97.28%  12.3006s      3024  4.0677ms  26.120us  48.296ms  cudaMemcpy
                    2.29%  289.76ms      2217  130.70us  2.8090us  236.24ms  cudaMalloc
                    0.36%  45.648ms      2217  20.589us  2.5410us  154.32us  cudaFree
                    0.06%  7.3821ms       604  12.221us  7.4290us  78.824us  cudaLaunchKernel
                    0.00%  411.20us         1  411.20us  411.20us  411.20us  cuDeviceTotalMem
                    0.00%  312.03us       101  3.0890us    297ns  124.86us  cuDeviceGetAttribute
                    0.00%  47.939us         1  47.939us  47.939us  47.939us  cuDeviceGetName
                    0.00%  25.379us         1  25.379us  25.379us  25.379us  cuDeviceGetPCIBusId
                    0.00%  4.0920us         3  1.3640us    588ns  2.6020us  cuDeviceGetCount
                    0.00%  2.0780us         2  1.0390us    457ns  1.6210us  cuDeviceGet
                    0.00%    798ns         1    798ns    798ns    798ns  cuDeviceGetUuid
```

Fig. 16.  Results of changing to 1024 threads

```
Each sample counts as 0.01 seconds.
  %   cumulative   self              self     total
 time   seconds   seconds    calls  ms/call  ms/call  name
 51.16    30.61    30.61                               Potential()
 48.87    59.84    29.24      201   145.46   145.46   computeAccelerations()
  0.03    59.86     0.02                               VelocityVerlet(double, int, _IO_FILE*)
```

Fig. 17.  gprof for sequential version