

GPU parallel simulation algorithm of Brownian particles with excluded volume using Delaunay triangulations

Francisco Carter^{a,*}, Nancy Hitschfeld^a, Cristóbal Navarro^b, Rodrigo Soto^c

^a*Department of Computer Science, FCFM, Universidad de Chile, Santiago, Chile*

^b*Institute of Informatics, Universidad Austral de Chile, Valdivia, Chile*

^c*Physics Department, FCFM, Universidad de Chile, Santiago, Chile*

Abstract

A novel parallel simulation algorithm on the GPU, implemented in CUDA and C++, is presented for the simulation of Brownian particles that display excluded volume repulsion and interact with long and short range forces. When an explicit Euler-Maruyama integration step is performed to take into account the pairwise forces and Brownian motion, particle overlaps can appear. The excluded volume property brings up the need for correcting these overlaps as they happen, since predicting them is not feasible due to the random displacement of Brownian particles. The proposed solution handles, at each time step, a Delaunay triangulation of the particle positions because it allows us to efficiently solve overlaps between particles by checking just their neighborhood. The algorithm starts by generating a Delaunay triangulation of the particle initial positions on CPU, but after that the triangulation is always kept on GPU memory. We used a parallel edge-flip implementation to keep the triangulation updated during each time step, checking previously that the triangulation was not rendered invalid due to the particle displacements. We designed and implemented an exact long range force simulation with an all-pairs N -body simulation, tiling the particle interaction computations based on the warp size of the target device architecture. For the short range force simulation, we devel-

*Corresponding author

Email address: `francisco.carter@ug.uchile.cl` (Francisco Carter)

oped a parallel algorithm that builds and uses Verlet lists in order to handle the particle neighborhood in parallel. The algorithm is validated with two models of active colloidal particles. Upon testing the parallel implementation of a long range forces simulation, the results show a performance improvement of up to two orders of magnitude when compared to the previously existing sequential solution. The algorithm for the short range force presents a similar performance improvement regarding the parallel long range implementation.

Keywords: Parallel computing, Particle dynamics, Brownian dynamics, Overlap correction, Delaunay Triangulations, CUDA, GPGPU, N-body simulation

1. Introduction

A colloidal suspension is a mixture of microscopical insoluble particles dispersed throughout a continuous fluid, where particle sizes range from 1 nm to 10 μ m. Colloidal suspensions appear in several natural and artificial substances as the milk, mud, inks, cosmetics or latex paint, for example. Also, they are used in many intermediate industrial processes. The interactions between colloidal particles of various kinds [1] have effects on the physical and chemical properties of the mixture such as its viscosity or light dispersion. To study these and other properties it is necessary to simulate particle systems of growing numbers ($N \geq 10^4$). Also, colloids are being used as models for active systems, to describe the motion of self-propelled microorganisms [2, 3, 4].

Colloids can be modelled as hard bodies subject to Brownian diffusive motion. Colloidal particles can typically interact through the fluid in what is called hydrodynamic interactions, via electrostatic forces for charged colloids, which can be screened in an electrolyte, or with van der Waals forces [1]. In out of equilibrium conditions, phoretic forces also appear [5]. Except for the hydrodynamic forces, these interactions can be modelled with good approximation as pairwise additive forces, which in out of equilibrium conditions can eventually break the action-reaction symmetry.

The simulation of colloidal dispersions, can be divided on two main problems executed in sequence: updating the positions of the particles due to the interparticle interactions, according to some integration rule and ensuring that the bodies do not overlap because of their movement, in order to respect the excluded volume interaction. These problems are specific instances of the n-body simulation and collision detection respectively [6]. In some contexts, the simulation of colloidal particles is referred as Brownian dynamics.

There are two main methods for solving overlaps between particles: correcting all of them at once after they happen or use an event-driven approach, integrating the system until the collision instant, process the involved particles and repeat until the system reaches the target time step. The last method, is particularly useful when inertia is important and collisions result in rebounces as in granular materials [7, 8]. It requires knowing the positions of the involved bodies at the time of collision, which becomes difficult when random Brownian motion is present. For the simulation of colloidal particles, which lack of inertia and excluded volume acts like a boundary condition rather than producing collisions, the first method is more suited.

This work focuses on designing and implementing a novel parallel simulation algorithm for 2D colloidal particle interacting with short and long range pairwise forces, with periodical boundaries, excluded volume and Brownian motion. The algorithm implementation takes advantage of the data-parallel computing capabilities of the GPU architecture, which have proven to be effective at accelerating the simulation process of several computational physics problems [9, 10, 11]. The interactions forces are allowed to be non-reciprocal as in the case of active particles [12, 13]. The main contribution of this work consists of a new and efficient method of resolving particle overlaps by using Delaunay triangulations, which are maintained periodically and fully on the graphics card. The starting positions and triangulation are initialized on the host while all the simulation code is executed on the device. The random values are also generated on the graphics card, both on the initialization and simulation phases. The algorithm uses a GPU edge-flip implementation to keep the triangulation

fulfilling the Delaunay condition during each time step and to correct inverted triangles in case they are generated due to the particle displacements. For the short range force simulation, we developed a parallel algorithm that builds and uses Verlet lists in order to handle the particle neighborhood in parallel. The algorithm is validated with two models of active colloidal particles. Upon testing the parallel implementation of a long range forces simulation, the results show a performance improvement of up to two orders of magnitude when compared to the previously existing sequential solution. The algorithm for the short range force presents a similar performance improvement regarding the parallel long range implementation.

The paper is organized as follows: Section 2 describes the specific conditions and properties that the simulated systems must operate under. Section 3 lists previous related work used to solve similar problems. Section 4 details the designed solution with its subcomponents, data structures, and optimizations. The implementation of the algorithm is described in Section 5. Sections 6 and 7 cover the tests, benchmarks, validation and used methodology, presenting the running time and performance results when compared to the other implemented solutions. Finally, section 8 rounds up the obtained results.

2. Description of the model

This section contains the description of the involved concepts and properties of this problem that may differentiate it from other body simulation problems, such as the excluded volume and stochastic component of particle movement.

2.1. Preliminaries

Let $P = \{p_1, p_2, \dots, p_N\}$ be a set of N bodies on a d -dimensional space. The n -body simulation is the computation of the interactions over each body in P , where F_i corresponds to the interaction over p_i by effect of $P_i = P \setminus \{p_i\}$, the set of all other bodies in the system. The interaction forces F typically depend on the distance r_{ij} between two bodies as $F \sim r_{ij}^{-q}$. If $q > d$ the force is said to be

short ranged, while if $q \leq d$, it is considered a long range force. When the forces are long ranged, the set P_i cannot be reduced in n-body simulations and an exact evaluation of the forces has a cost $O(N^2)$. Approximate solutions for long range interactions as the Barnes-Hut algorithm reduce the cost to $O(N \log N)$ [14]. But for short range forces, the interactions can be truncated and, therefore, P_i can be reduced to the neighborhood of particles close to p_i . In this case, the evaluation of the forces costs $O(N * N_{NL})$ on average, where N_{NL} is the average number of neighbors of a body [15]. Since the construction of the list is $O(N^2)$ for evaluating all pairwise distances between bodies, it is possible to partition the simulation domain in cells so that close bodies get binned together in the same cells. Assignment of bodies to their respective cells takes $O(N)$ time [16, 17].

The simulation domain is a two-dimensional $L \times L$, periodical box across the X and Y axes, meaning that the particles wrap around the box as they move across its boundaries. For the distance calculations between particles, including force calculations, we follow the minimum-image convention, in which a particle interacts with another via its real position or its image depending on which is the shortest.

2.2. Particle interaction without excluded volume

Microscopic particles move in an overdamped regime, with no inertia. When subject to a force \vec{F} , the equation of motion is simplified to $d\vec{r}/dt = \gamma\vec{F}$, where γ is the mobility. Absorbing the mobility coefficient into the force, which will then have velocity units, in a time step Δt , the integration rule for updating the position of a particle over time is performed using the Euler-Maruyama method:

$$\vec{r}_i(t + \Delta t) = \vec{r}_i(t) + \vec{F}_i(t)\Delta t + \vec{\xi}\sqrt{D\Delta t}, \quad (1)$$

where \vec{F}_i is the deterministic velocity obtained from the interactions between the particle i and P_i , D is the diffusion coefficient, and $\vec{\xi}$ is a random vector, where the components follow a normal distribution of zero mean and unit variance,

and corresponds to a noise added that takes into account the diffusive Brownian motion [18].

The force model we use for the simulations describes the interaction of self-diffusiophoretic active particles [12]. In this model, particles can be of different type, characterized by two charges, α and μ ; the former is responsible of creating the concentration field, while the second describes the response of a particle to the field, leading to the following interaction law:

$$\vec{F}_i = \sum_{k \neq i} \mu_i \alpha_k \vec{f}(\vec{r}_i - \vec{r}_k), \quad (2)$$

where $\vec{f}(\vec{r}) = \vec{r}/r^3$ for the studied long range force, while $\vec{f}(\vec{r}) = \vec{r}/r^7$ for the short range interaction. Note that if $\alpha_i \neq \mu_i$, the action-reaction symmetry is broken and self-motion is possible. Charged colloidal particles are included in this model if $\alpha_i = \mu_i = q_i$, equal to the electric charge of the particles.

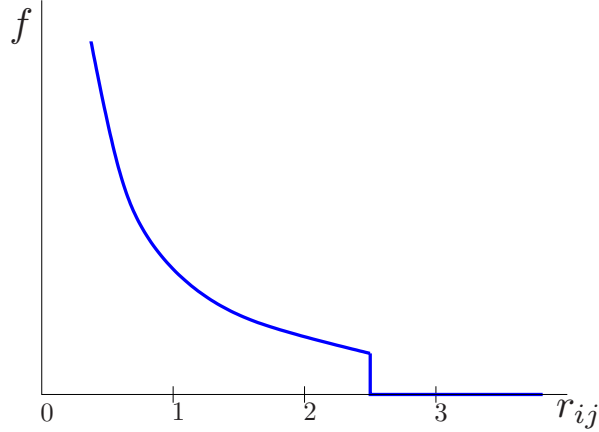


Figure 1: Cutoff of the short range forces. For distances larger than r_{cutoff} , the force is small and therefore is set to zero to speed up calculations. The jump at r_{cutoff} has been exaggerated for illustration purposes.

Since the short range force decays much faster with distance compared to the long range force, its calculation considers a cutoff radius from which the value of the force is considered zero, as shown on Figure 1. The short range

force is then computed as:

$$\vec{F}_{ij}(\vec{r}_{ij}) = \begin{cases} \mu_i \alpha_k \vec{f}(\vec{r}_{ij}), & \text{if } r_{ij} \leq r_{\text{cutoff}} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

We used $r_{\text{cutoff}} = 2.5\sigma$ for the simulated short range force in our experiments, where σ is the particle diameter.

2.3. Excluded volume

The simulated particles are represented as hard disks with a uniform diameter σ . Although here we consider only monodisperse colloids, it is direct to extend the method to polydisperse systems where radii do not differ too much. Since the integration rule (1) ignores the excluded volume condition, it can happen that the updated positions produce overlaps between two or more particles, resulting in a physical impossibility. In order to ensure this property, at the end of each time step, the members of all overlapping pairs (p_i, p_j) are moved apart from each other in a way that corrects the overlaps:

$$\vec{r}_1' = \vec{r}_1 - \delta^* \frac{(\vec{r}_2 - \vec{r}_1)}{|\vec{r}_{12}|} \quad \vec{r}_2' = \vec{r}_2 - \delta^* \frac{(\vec{r}_1 - \vec{r}_2)}{|\vec{r}_{12}|}, \quad (4)$$

where \vec{r}_1 and \vec{r}_2 are the original positions and \vec{r}_1', \vec{r}_2' the updated positions. If $\delta^* = (\sigma - |\vec{r}_{12}|)/2$, the particles would move in opposite directions from each other along $\hat{n} = (\vec{r}_2 - \vec{r}_1)/|\vec{r}_{12}|$, leaving the particles in tangential contact. If $\delta^* = \sigma - |\vec{r}_{12}|$, the movement is proportional to the magnitude of the previously existing overlap, simulating a bounce effect resulting from the collision at some instant $t^* \leq t + \Delta t$. This last value is the one used for processing the overlaps in the simulation and guarantees that no accumulation is produced at the contact distance.

2.4. Stochastic displacements

The particle displacements on (1) have a random noise component $\vec{\xi}$, modeled as a random variable with standard normal (or Gaussian) distribution with zero mean and standard deviation $\hat{\sigma} = \sqrt{D\Delta t}$. Reducing Δt , the deterministic and stochastic displacements in each time step are reduced. However, for

a Gaussian distribution, it is always possible that large values are generated (at the tail of the distribution), leading to excessively large displacements (see Figure 2). To avoid these problems, the simulation ignores values larger than 3 standard deviations. We considered two methods in order to achieve this, as shown on Figure 3:

- (a) Reroll the values outside the range $[-3\hat{\sigma}, 3\hat{\sigma}]$.
- (b) Truncate the values to the range $[-3\hat{\sigma}, 3\hat{\sigma}]$.

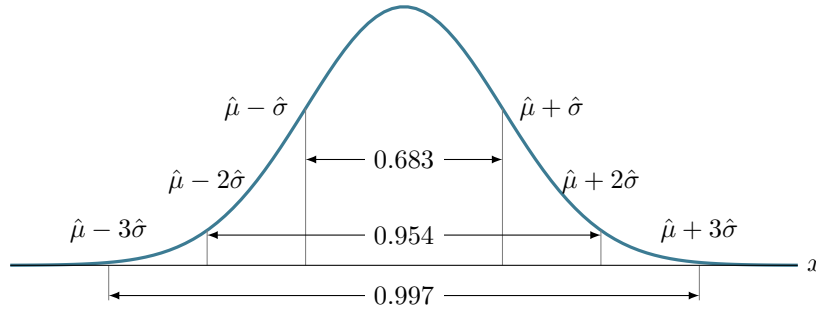


Figure 2: Normal distribution with mean $\hat{\mu}$ and standard deviation $\hat{\sigma}$. The probability to get numbers in the ranges $[\hat{\mu} - \hat{\sigma}, \hat{\mu} + \hat{\sigma}]$, $[\hat{\mu} - 2\hat{\sigma}, \hat{\mu} + 2\hat{\sigma}]$, and $[\hat{\mu} - 3\hat{\sigma}, \hat{\mu} + 3\hat{\sigma}]$ are 0.683, 0.954, and 0.997, respectively.

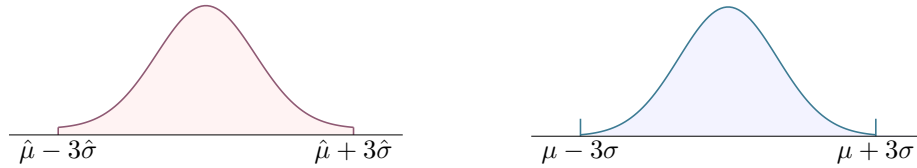


Figure 3: Distributions that result after discarding values larger than $3\hat{\sigma}$ from the original Gaussian distribution. Two methods are used. Left: Rerolling values out of range. Right: Truncating the generated values. In this case, Dirac-delta contributions of small amplitude appear at $\hat{\mu} \pm 3\hat{\sigma}$.

Both methods produce probability distributions different from each other and from the original; while the first alternative raises the probability of all values in range, the second one raises the probability at the edges. These modifications do not generate a noticeable statistic distortion, since the considered

range includes 99.7% of the possible values. In our simulations, we opted for the second method, which turns out to be faster and better suited for parallel execution, since it needs to generate a single random number instead of a variable quantity of random values in the first method.

3. Related work

For short-range forces calculation, the standard technique is the use of Verlet lists [19, 16, 17]. The authors in [20, 21] parallelize the list construction by having a $O(N^2)$ list of all possible pairs of bodies. A predicate checking closeness between the pair members is evaluated over all elements of the list, which can then be used for a key-value sort to group all the neighboring pairs consecutively in the array. A parallel scan operation allows to get the number of elements that must be copied to the neighbor list. The authors then combine this algorithm with fixed cell partitioning in order to replace distance calculations with less-expensive cell neighborhood checks.

For the parallel n-body simulation, with full calculation of the $O(N^2)$ forces, Nyland et al. [22] developed a grid-style tiling algorithm, reading the particles from the global space and storing them on GPU shared memory, increasing performance as multiple threads read from that space at a lower latency. Partitioning the load/store process on groups of p particles allows fitting an arbitrary input size on the hardware-limited shared memory size. Burtcher and Pingali [23] parallelized the Barnes-Hut simulation [14], which computes an approximation for the force, representing the cell hierarchy kd -tree as multiple arrays for each node field. It uses atomic lock operations to build the tree in parallel, throttling the threads that failed to get the lock so they do not waste bandwidth with unsuccessful lock requests. The tree is then filled with the center of mass data, starting from lower nodes in the tree according to the order of allocation for the scan. Bedorf et al. [11] uses a Z-order curve to sort the particles spatially. Each thread is assigned to a particle, applying a mask value to it to determine the octree cell the particle should be assigned. The linking

of the tree is made by assigning a thread to each cell node and then doing a binary search over the corresponding Z-order key to find both the parent and child nodes, if appropriate.

To detect and process collisions, Hawick and Playne [24] developed a multi-GPU algorithm with a tiling scheme similar to the one used by Nyland et al. [22]. If a pair of particles overlap, the associated threads store the index of its colliding neighbor and the time at which the collision occurred. The collisions then are resolved iteratively starting from the earliest, redoing the previous process in order to find possible new collisions.

Finally, for overlap correction, Strating [25] describes a brute-force sequential algorithm that checks all pairs of bodies for possible overlaps and corrects them following equation (4). The algorithm may need to iterate an unbounded number of times at each time step because some corrections may generate new overlaps with neighboring particles.

4. Algorithm

This section describes the parallel algorithm in detail. It includes the generation of the initial data, data structures, overlapping detection and correction, and Delaunay condition updates, among others, putting emphasis in what threads are doing at each time step.

4.1. Overview

The simulation consists of two phases: (i) sequential initialization of the simulation data, followed by a host to device transfer and (ii) a parallel simulation phase. The initial positions are initialized over a triangular mesh with $N^* \geq N$ vertex, where each vertex represents a particle and their types are assigned randomly according to the specified concentrations. A sample of N particles is selected from the mesh by Reservoir Sampling [26], resulting in the input particle set, which is homogeneous in space.

The n-body algorithm for the long range force is based on a grid-style tiling, which uses the shared memory of the multiprocessor assigned to each thread

block to store the particles in groups. In this algorithm each thread is mapped to exactly one particle in the system, and since it is possible to lack action-reaction symmetry on the force, no redundant computation is done unlike the cases where $\vec{f}_{ij} = -\vec{f}_{ji}$.

Once the forces are calculated, the particles are advanced one time step using eqn (1). As a result, particle overlaps can appear. When Δt is small enough, for hard disks of similar or equal radii, only neighbor particles can overlap. Then, to detect and correct overlaps, instead of a brute-force algorithm that would check all $O(N^2)$ pairs, only neighbors are checked. The Delaunay triangulation [27] is particularly well suited to detect neighbors for monodisperse or slightly polydisperse disks. In dense systems, the overlap corrections can be highly non-local, as the correction of one pair can generate a sequence of other overlaps that need correction. It is therefore not clear a priori the computational cost of this stage, which is the reason why we consider both short and long range interaction forces.

The Delaunay triangulation can be built constructively or from an existing triangulation. Lawson’s algorithm [28] accepts a triangulation as input and transforms it into a Delaunay triangulation via a finite sequence of edge-flip operations [27]. Based on the Lawson algorithm, Navarro et al. [29] developed a parallel implementation for generating quasi-Delaunay triangulations, so it is possible to keep the triangulation updated without need of host-device memory transfers. These are quasi-Delaunay because exact predicates are too expensive on the GPU; nevertheless, this approximate construction is sufficient for our problem. However, since the input for the Navarro et al. algorithm must be a valid triangulation, we must first correct potential triangle inversions with invalid edge intersections, which can result from the particle displacements. With this strategy, the Delaunay triangulation of the particle positions is built only once from scratch on the host [30, 31], which is transferred to the device. Thereafter, the triangulation is maintained updated after each time step on the device.

Algorithm 1 Particle system simulation

Require: $P = \{p_1, \dots, p_N\}$ list of particle positions

Ensure: $P = \{p_1, \dots, p_N\}$ list of positions updated to current time

```
1: procedure RUNSIMULATION( $P$ )
2:   Generate starting position of  $N$  particles
3:   Build the Delaunay triangulation
4:   for  $t \leftarrow 0$  to  $T_f$  do
5:     Integrate the  $N$  particles on  $t + \Delta t$ 
6:     Correct inverted triangles
7:     Update Delaunay triangulation
8:     Correct overlaps between particles
9:   end for
10: end procedure
```

4.2. Data structures

We store the simulation data as a Structure of Arrays (SoA) on global device memory, using total $O(N)$ space. The particle data consists of their position (x_i, y_i) and their charges (α_i, μ_i) , stored as floating point vector types¹ in order to increase bandwidth utilization [32]. We use an additional buffer array for positions so that writes are not done at the same addresses for reads, avoiding a synchronization step. We store the simulation parameters that remain unchanged during a same instance on a constant device memory structure [32], such as $N, D, \Delta t, \sigma$ and derived constants σ^2 and $\sqrt{D\Delta t}$. Additionally, we store the triangulation data using the same scheme as [29].

4.3. Inverted triangle detection

There are two possible reasons for a triangle inversion: an edge gets inverted because the distance vector between its terminal vertex changes orientation or because one of its vertex crosses the edge opposite to it. The first problem means

¹float2 or double2.

that the particles went through each other, which is a physical impossibility and must not be allowed, and can only take place for large Δt . The criterion $\vec{r}_{ij}^0 \cdot \vec{r}_{ij}^1 < 0$ is used to determine if the above situation happened on the current time step, where $\vec{r}_{ij} = \vec{r}_j - \vec{r}_i$ is the distance vector between the compared particles, and \vec{r}_{ij}^0 and \vec{r}_{ij}^1 are evaluated with current and previous positions, respectively. The buffer array allows to compare the distances before and after integration, and the edges of the triangulation show the pairs that need checking. Once an invalid movement is detected, the last positions are discarded and integration is repeated with a lower Δt value than the currently used.

Finally, the inverted triangle detection becomes equal to checking if a vertex crossed (towards a neighboring triangle) any of the edges that enclose it, which is equivalent to point-in-triangle detection. Using the barycentric coordinates d , s , and t on triangles (see Figure 4),

$$\vec{e}_0 = \vec{v}_2 - \vec{v}_1 \quad d_2 = \vec{e}_2 \times \vec{e}_0 \quad (5)$$

$$\vec{e}_1 = \vec{v}_3 - \vec{v}_1 \quad s_2 = \vec{e}_1 \times \vec{e}_0 \quad (6)$$

$$\vec{e}_2 = \vec{v}_4 - \vec{v}_1 \quad t_2 = \vec{e}_2 \times \vec{e}_1 \quad (7)$$

the criterion used to detect the edges that must be flipped is

$$\text{flip}(\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4) = \begin{cases} (s_2 \leq 0) \wedge (t_2 \leq 0) \wedge (s_2 + t_2 \geq d_2), & \text{if } d_2 < 0 \\ (s_2 \geq 0) \wedge (t_2 \geq 0) \wedge (s_2 + t_2 \leq d_2), & \text{otherwise} \end{cases} \quad (8)$$

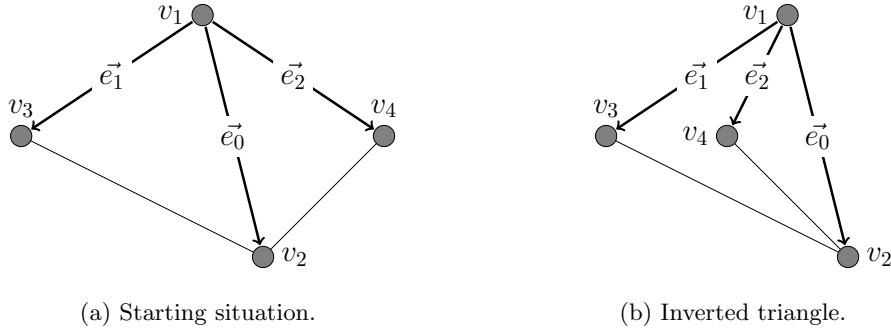
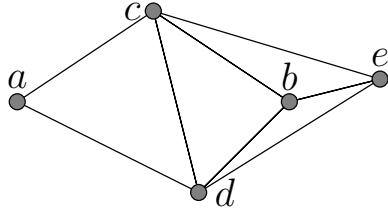
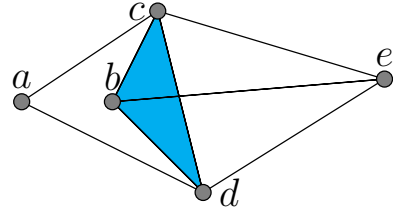


Figure 4: Inverted triangle detection using barycentric coordinates. Particle sizes are scaled down compared to distances.

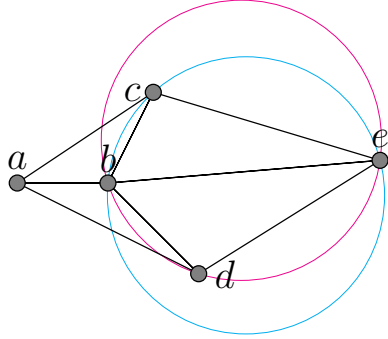
Figure 4 shows the vectors used on the predicate that checks if the original edge must be flipped. The predicate becomes true when applied to edge (v_1, v_2) in (b), because point v_4 lies inside the triangle (v_1, v_2, v_3) . This means that (v_1, v_2) must be replaced with (v_3, v_4) , as in a common edge-flip operation. In this case, v_3 and v_4 are stored as opposite vertices to edge (v_1, v_2) in the triangulation data structure in GPU device memory.



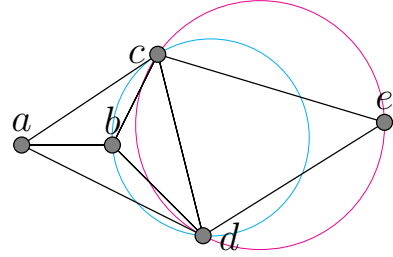
(a) Starting situation.



(b) Particle b moves over the edge bc .



(c) Edge flip between ab and cd .



(d) Edge flip between cd and be .

Figure 5: Inverted triangle correction. The cyan shaded triangle was inverted by the movement of particle b . Particle sizes are scaled down compared to distances.

It is worth noting on Figure 5 that the movement of point b across edge (c, d) creates an intersection between it and edge (b, e) . The edge flip between (a, b) and (c, d) removes the inverted triangle, restoring the local triangulation. The triangulation may still not satisfy the Delaunay property, so additional edge

flips may be needed on further steps. For this stage, we use the Navarro et al. algorithm [29].

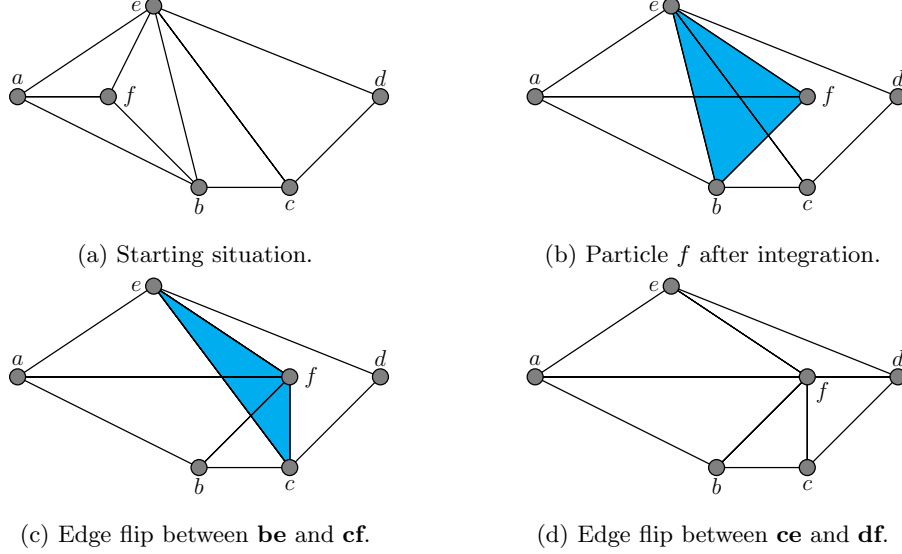


Figure 6: Inverted triangle correction with two edge flips. Particle sizes are scaled down compared to distances.

Figure 6 shows a situation where particle f moves across edges (b, e) and (c, e) , needing two consecutive flips in order to restore the local triangulation. While function 8 cannot properly evaluate this case or similar movements across further distances, this kind of inversion can be detected allowing to revert the step. Anyway, small time steps guarantee that this situation is extremely unlikely to happen.

4.4. Overlap correction

The overlap correction uses the topological information contained in the edges of the Delaunay triangulation, which allows for each particle fast access to the neighborhood of particles that may be overlapping with it. The algorithm maps threads to edges in such a way that each thread handles one edge of the triangulation. A thread gets the positions of the particles that form the edge,

checking if there exists an overlap between them ($r_{ij} < \sigma$). If the check is positive, the algorithm computes the displacements of the involved particles according to (3). Since the same particle can be part of many edges at once, the algorithm sums atomically the displacements in a global array, in order to avoid concurrency hazards.

Algorithm 2 Overlap correction

Require: P_0 starting positions, E triangulation edges

Ensure: P_1 displacements over each particle

```

1: procedure CORRECTOVERLAPS( $P_0, P_1, E$ )
2:   for  $i \leftarrow 0$  to  $|E|$  do
3:      $e_i \leftarrow E[i]$ 
4:      $b_i \leftarrow P_0[e_i.\text{first}]$ 
5:      $b_j \leftarrow P_0[e_i.\text{second}]$ 
6:      $r_{ij} \leftarrow \text{dist}(b_i, b_j)$ 
7:     if  $r_{ij} < \sigma$  then
8:        $\delta \leftarrow \sigma - r_{ij}$ 
9:       atomicAdd( $P_1[e_i.x], -\delta \vec{r}_{ij}$ )
10:      atomicAdd( $P_1[e_i.y], \delta \vec{r}_{ij}$ )
11:    end if
12:  end for
13: end procedure

```

Once the algorithm computes the total displacements, it maps each thread with a particle in the same way as described previously for edges. Each thread then updates the position of its particle, applying the periodic boundary conditions when necessary. It is possible that the updated positions may still have some of the previous overlaps or even have some newly generated ones. In this case, the algorithm repeats the previous process until no overlaps are present.

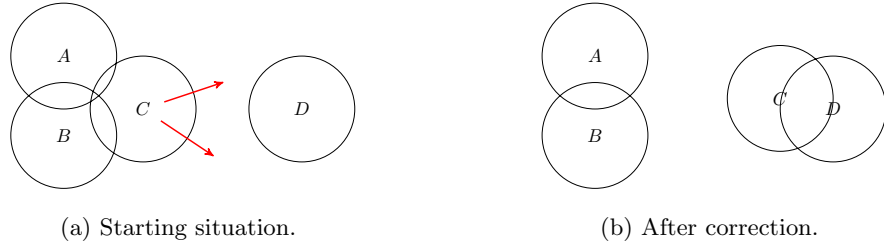


Figure 7: Possible instability with the parallel overlap correction. The displacement over C caused by the overlaps with A and B can be greater than needed, which can generate a larger overlap with a neighboring particle D . This problem is prevented by truncating the displacements to a maximum amount.

When adding the partial displacements on a particle, it may happen, if these point in the same direction, that the resulting total displacement is excessively large (see Figure 7). These displacements, larger to what is needed to solve the overlap, can generate new overlaps. Eventually, the correction of the new overlaps can result in an instability, where the displacements increase with alternating sign and the iterative procedure does not converge. A solution is truncating the displacement with the heuristic value $\sigma/4$ (half the particle radius), preventing the emergence of the instability.

Finally, it is worth mentioning that the parallel correction algorithm presented here does not correspond to a parallelization of the sequential algorithm of Strating [25], which displaces particles sequentially, while in our case the displacements are added and performed in parallel. Hence, due to the chaotic dynamics of the system, the small differences in these algorithms will produce different outputs for finite Δt .

4.5. Long range forces

An improvement to the long range force calculation consists on using the intrinsic warp shuffle instruction, which allows a thread access to the registers of other threads belonging to the same warp. Each thread is assigned a lane number that identifies it from the other warp members, allowing them to read different particles from global memory. Then, each warp member takes turns in

propagating the data of its corresponding particle to the other threads, who can read it via the `_shfl()` instruction by passing as argument the lane number of the thread currently in turn. Once the whole warp has shared the data among its members, each member reads a particle from global memory and repeats the same process until all particles have been visited. The main advantage of this optimization is a greater efficiency of memory accesses, since most of the time the threads are sharing data at registry level instead of more expensive load requests on global memory. Also, the concurrent execution of the warp members makes unnecessary the explicit synchronization of the inner warp shuffle loop.

Algorithm 3 Particle system integration

Require: P_0 particle array $(x_i, y_i, \alpha_i, \mu_i)$

Ensure: P_1 particles with updated positions

```

1: procedure INTEGRATE( $P_0, P_1$ )
2:   for  $i \leftarrow 0$  to  $|P|$  do
3:      $l_i \leftarrow \text{threadIdx.x} \ \& \ (\text{warpSize} - 1)$ 
4:      $\vec{b}_i \leftarrow P_0[i]$ 
5:      $\vec{v}_i \leftarrow 0$ 
6:     for  $j \leftarrow 0$  to  $|P|$ ;  $j \leftarrow j + \text{warpSize}$  do
7:        $\vec{b}_j \leftarrow P_i[j + l_i]$ 
8:       for  $k \leftarrow 0$  to  $\text{warpSize}$  do
9:          $\vec{b}_k \leftarrow \text{\_shfl}(\vec{b}_j, k)$ 
10:         $\vec{r}_{ik} \leftarrow \text{dist}(\vec{b}_i, \vec{b}_k)$ 
11:         $\vec{v}_i \leftarrow \vec{v}_i - \vec{r}_{ik} \cdot (\alpha_i \mu_k) / r_{ik}^3$ 
12:      end for
13:       $\text{\_syncthreads}()$ 
14:    end for
15:  end for
16:   $P_1[i].x \leftarrow \vec{b}_i.x + \vec{v}_i.x \Delta t + \vec{\xi}_i.x \sqrt{D \Delta t}$ 
17:   $P_1[i].y \leftarrow \vec{b}_i.y + \vec{v}_i.y \Delta t + \vec{\xi}_i.y \sqrt{D \Delta t}$ 
18: end procedure
```

The `.x` and `.y` operators reference the data of the vectorized CUDA structures for each respective variable. For example, the noise $\vec{\xi}_i$ has \hat{x} and \hat{y} components, so it is grouped as a single vector for increased memory performance [32].

5. Implementation

The parallel algorithms described in the previous section were implemented on CUDA 7.5 and C++ 11, using function templates to choose between `float` and `double` precision formats at compile time. We use the CGAL library [30] to create the 2D periodic Delaunay triangulation, which is then sent to device memory alongside the particle data before starting the simulation. The random numbers used on the parallel implementations when initializing the starting positions and generating noise during integration are created with the XOR-WOW pseudorandom number generation algorithm of the `cuRAND` library [33], using the host and device APIs respectively. Each configuration has two particle types, although the program can support a variable number of particle types for simulating. For comparison purposes, we also implemented a fully sequential long range forces algorithm with the overlap correction discussed on [25], and a parallel short range forces algorithm using Verlet lists and a discrete grid over the simulation box. The neighbor list computation during the Verlet lists construction in parallel is similar to the one described on [34], grouping together all the particles that belong in the same cells.

When the simulation finishes, the final positions are brought back to host memory and written to an output file. The visualizations on Figures 8, 13, 14 and 15 were generated reading the respective output files.

6. Performance results

Parameters:. We generated inputs for 11 different values of N and 5 parameter configurations, as described on Table 1. The starting positions are generated semi-randomly as described on section 4, keeping the same seed value for the random number generator across all simulation instances.

config	ϕ_1	α_1	μ_1	ϕ_2	α_2	μ_2	ρ
0	0.7	1	1	0.3	1	-1	0.79
1	0.5	1	1	0.5	1	-1	0.52
2	0.5	1	1	0.5	-1	-2	0.52
3	0.5	-1	1	0.5	1	-2	0.79
4	0.5	1	1	0.5	-1	-4	$8.73 \cdot 10^{-2}$

Table 1: Parameters used for the tests, where each configuration is identified by a digit and all of them contain two types of particles. ϕ_i is the concentration of particles of type i , α_i, μ_i are the charges used in the force calculation, and ρ is the packing fraction of particles on the simulation box.

Each configuration has two types of particles with charges α_i, μ_i . The fraction of particles of each type is given by $\phi_i = N_i/N$, where N_i is the number of particles of each type and $N = N_1 + N_2$ is the total number of particles. The area fraction $\rho = N(\sigma/2)^2/L^2$ is a measure of the particle density. To study scaling times, we change the length of the simulation box to keep density constant when increasing N :

$$L(N) = \sqrt{\frac{N\pi(\sigma/2)^2}{\rho}} \quad (9)$$

The values for N start at 2^{10} , raising the exponent by 1 until $N = 2^{20}$. Finally, we kept constant the values for $\sigma = 1, \Delta t = 0.01, D = 0.01, \delta = 1.0$ for all configurations and input sizes.

Figure 8 displays the particle positions after 10^4 time steps for each configuration. The election of the parameters help to test the algorithms under different conditions of fluidity, density and homogeneity. For c0, there is an asymmetric attraction between particles of type 1 and 2, resulting in an homogeneous mixture, with fluid-like motion. In c1, there is a larger concentration of type 2 particles, which self-attract forming a dense cluster, segregated from type 1 particles, which self-repel forming in a gas-like state. In configuration c2, equal particles repel, while dissimilar particles attract, favouring the formation

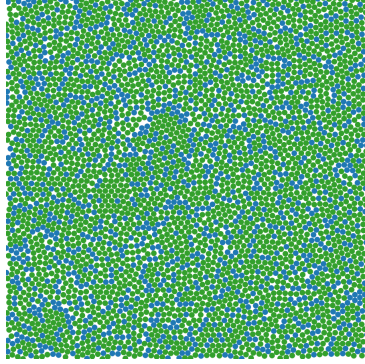
of chain-like structures, where 1 and 2 particles alternate. In c3, the situation is the opposite, where equal particles attract, while dissimilar particles repel, leading to the formation of dense segregated clusters. Finally, the interactions in c4 are analogous to those of c2, in a dilute regime, resulting in the formation of small clusters.

System:. We ran the tests on a machine with a Tesla K40c GPU and a Intel(R) Xeon(R) CPU E5-2640 v3 @ 2.60GHz. The tests for both the sequential and parallel implementations were made on the same machine.

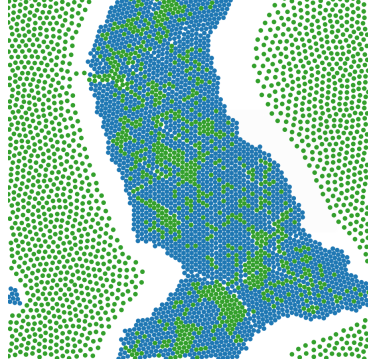
Compilation:. We compiled the program using `nvcc V7.5.17` with compiler options `--std=c++11 -gencode arch=compute_30,code=sm_30`. For the sequential code, we used `g++ 4.9.2` with options `--std=c++11 -O3`.

Metric:. We ran simulations for 100 iterations, long enough to ensure that particle collisions happen frequently, except for the first iterations where the bodies are separated from each other. There we compute the average execution time and iteration averages for overlap correction cycles and edge flips.

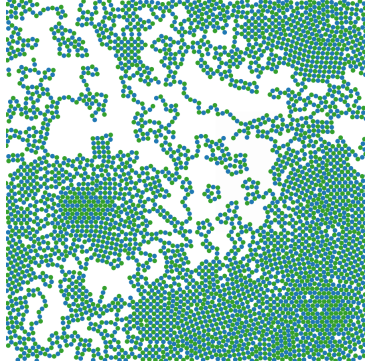
The average execution time per time step, presented in Figure 9, shows two interesting features. First, for the case with long range forces the execution time is $O(N^2)$, while for short range forces it is $O(N)$. Since both include the overlap correction algorithm, this implies that the execution time for the later is $O(N)$. Second, except for small systems, in the case of long range forces the execution time does not depend on the configuration, while for the short range forces, there is a clear dependence, with increasing complexity for c_4, c_1, c_2, c_0, c_3 (the same order of complexity is observed for long range forces at small N). This result is consistent with a cost $O(N)$ for the overlap correction, with a prefactor that may depend on the density and extension of the clusters.



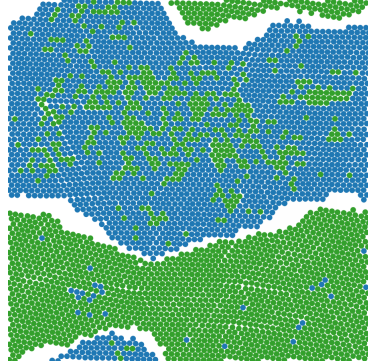
(a) Configuration 0



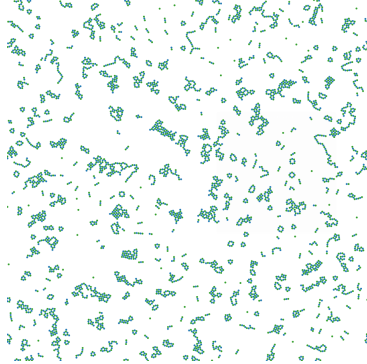
(b) Configuration 1



(c) Configuration 2



(d) Configuration 3



(e) Configuration 4

Figure 8: Particle positions after 10^4 iterations, with $N = 4096$ and $\Delta t = 0.01$. Type 1 particles are in green and type 2 in blue. The configurations and particle types for each system are described on Table 1.

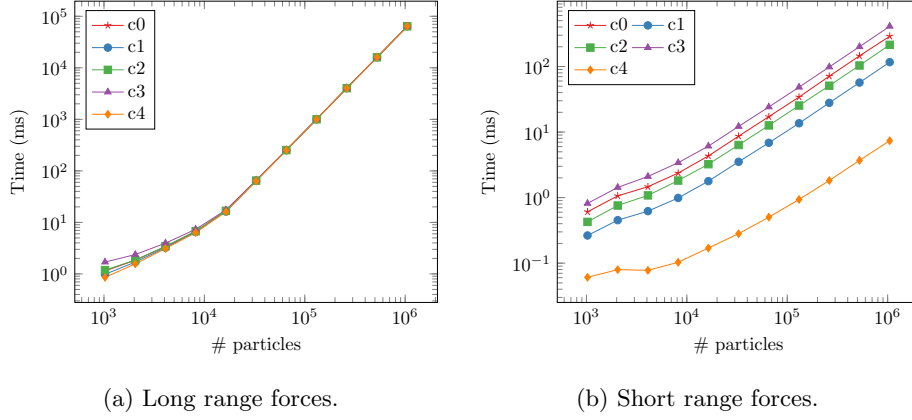


Figure 9: Average execution time per time step for simulations using long and short range forces, using the parameters shown on Table 1.

To study the dependence of different configurations on the complexity of the overlap correction, in Figure 10 we plot the average per time step of iterations needed to correct all overlaps. The increasing complexity for c_4, c_1, c_2, c_0, c_3 is consistent with the previous results on Figure 9, because in c_3 most of the particles participate in corrections, while in c_2 almost half of the particles are excluded due to repulsion between same-type particles. However, this does not explain why c_4 has the least complexity factor, even though half of the particles overlap. This happens because all the overlaps on c_4 are corrected on the first iteration, which is probably due to the small size of the clusters. The number of iterations follow the same order in complexity as the execution time. Except for c_4 where clusters are disconnected, the number of iterations grow with N . This effect is due to the percolation of the large clusters, which cover the entire box and, therefore, the corrections become non-local and system size dependent. This growth is nevertheless weak, following an approximate logarithmic law. It is also noteworthy that the curves for long range forces are constant on c_4, c_1 and c_0 , slightly grow on c_2 and is relatively greater on c_3 .

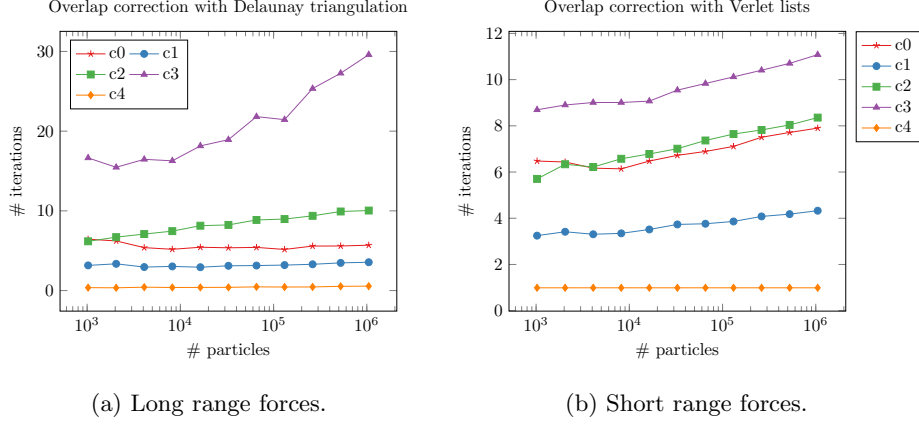


Figure 10: Average overlap correction iterations per time step for short and long range forces simulation, using the parameters shown on Table 1.

We also measured the performance of the Delaunay triangulation update algorithm, reporting the average of edge flip iterations made for both inverted triangle corrections and Lawson's algorithm. The curves obtained on Figure 11 are less regular than the previous results, but keep the same general tendency. Unlike the curves for overlap correction, on where c_4 shows much smaller values than the other configurations, here the curve is comparable to c_0 and c_1 . This happens because the underlying triangulation for c_4 has a great number of slivers, formed by the small particle density that forms relatively long edges. Then, according to the inverted triangle condition on section 4, it is more likely for c_4 to produce inverted triangle than the other configuration, whose triangles are more equilateral. The average for edge-flip iterations for long range forces has linear growth for all configurations, noting that c_2 and c_3 are the hardest cases to solve, as is the case on Figure 10. Though the number of iterations grows with N , it still remains negligible regarding the total time of a time step, so it is not a priority target for optimization.

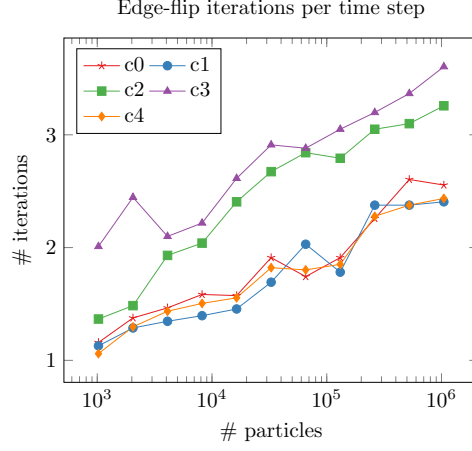


Figure 11: Average edge-flip iterations per time step for long range forces integration, using the simulation parameters shown on Table 1. The short range forces algorithm is not analyzed, since it does not use Delaunay triangulations.

Finally, we compared the n-body algorithms for long range forces, used on the different implementations without considering overlap corrections. `shuffle` is the presented optimization using warp-shuffle, while `sharedMem` is the GPU device memory algorithm described on section 4, observing a performance improvement of up to 2.4 times from optimizing the parallel implementation for all tested values of N . It is also noteworthy that the optimized n-body algorithm allows simulation of $N = 10^6$ particles at the same time that the sequential implementation solves the problem for 10^5 bodies. For input sizes relevant to this study ($N \geq 10^4$), the time used by the sequential implementation is two orders of magnitude higher than the parallel solution, which allows the simulation of bigger particle systems for a longer physical time.

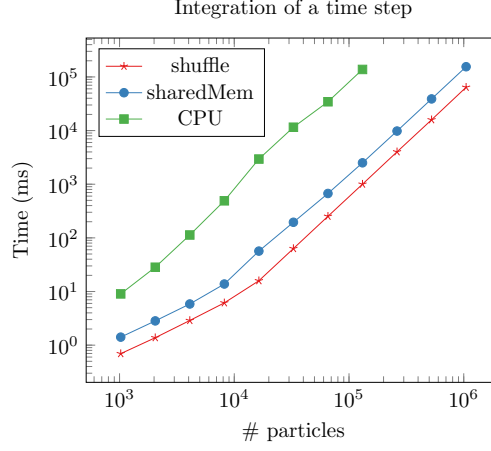


Figure 12: Comparison between execution times in milliseconds for both implementations of the quadratic n-body algorithm, using configuration c_0 described on Table 1.

7. Validation

In order to verify that the developed overlap correction algorithm is efficient enough, we made two validation experiments. First, we test the locality of the correction, that is, how far it propagates through the system. For configurations c_0 and c_3 we print the result after 10^5 time step iterations, painting with red the particles that took part in overlaps during the last simulated time step. Particles that did not take part in overlaps were painted green, so that every particle has a color. We repeat the process for decreasing values of Δt , expecting that the number of overlaps will decrease as the time step produces smaller movements. The results on Figure 13 allows us to verify the complexity factor associated to each configuration that shows up on the previous performance curves. Configuration c_0 involves much less particles on overlap corrections than c_3 upon lowering the time step. This decrease in execution time by reducing Δt does not compensate, however, for the larger number of steps that are needed to achieve an specified physical time.

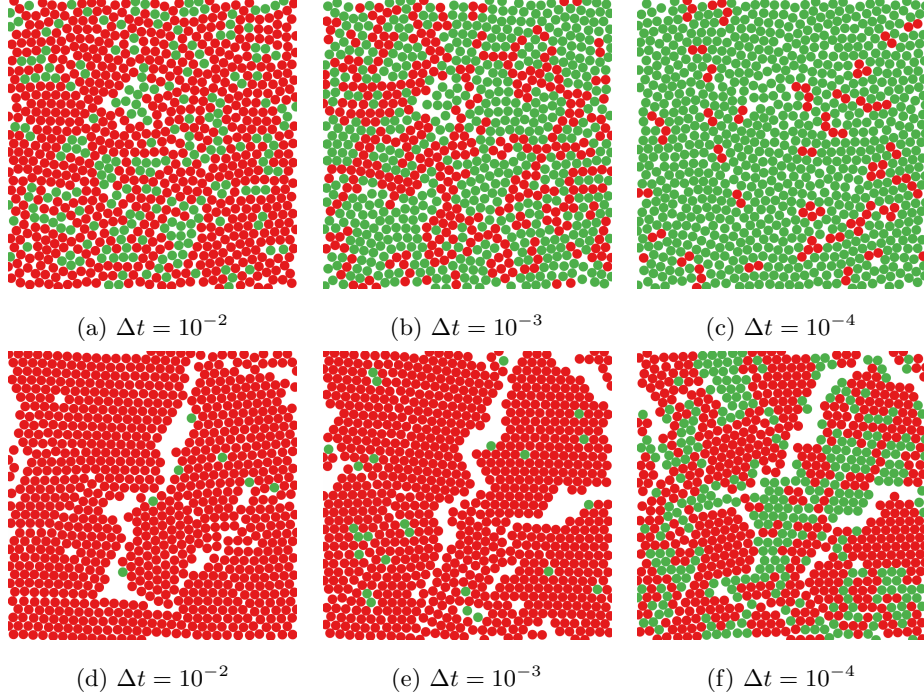


Figure 13: Overlap correction locality visualizations of configurations *c0* and *c3*. The red particles participated in at least one overlap correction on the same time step, while green particles did not.

The second validation consists in testing the overlap correction on another colloidal model. We consider the Active Brownian Particle (ABP) model [4], where particles move in 2D with velocities of fixed magnitude V_0 , with a direction that is specified by the director angle θ_i . The integration rule for the positions after an interval Δt is:

$$\vec{r}_i(t + \Delta t) = \vec{r}_i(t) + V_0(\cos \theta_i \hat{x} + \sin \theta_i \hat{y})\Delta t. \quad (10)$$

In the same time interval, the angles θ_i are subjected to diffusive rotational Brownian motion, of amplitude D , and therefore evolve as:

$$\theta_i(t + \Delta t) = \theta_i(t) + \sqrt{2D\Delta t} n_i, \quad (11)$$

where n_i is a random Gaussian variable of zero mean and unit variance.

We simulate the system with the same parameters used in Ref. [35], for two different packing fractions, obtaining the same phenomenology. At large packing fractions, the system evolves to the formation of a dense percolating dynamic cluster (see Fig. 14). Reducing the packing fraction, small clusters form, which merge in a slow coarsening process in the course of time as shown in Fig. 15.

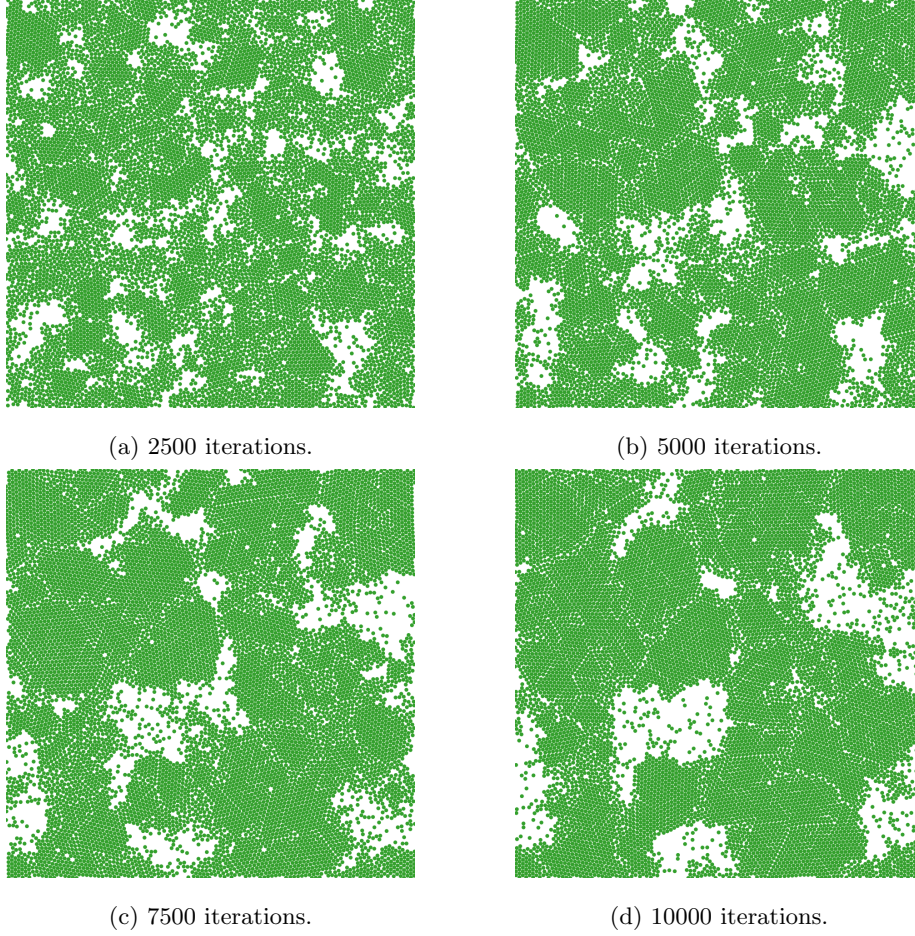


Figure 14: Snapshots of an ABP system $N = 10^4$, $L = 105.9$, $D = 0.01$ and packing fraction $\rho = 0.7$.

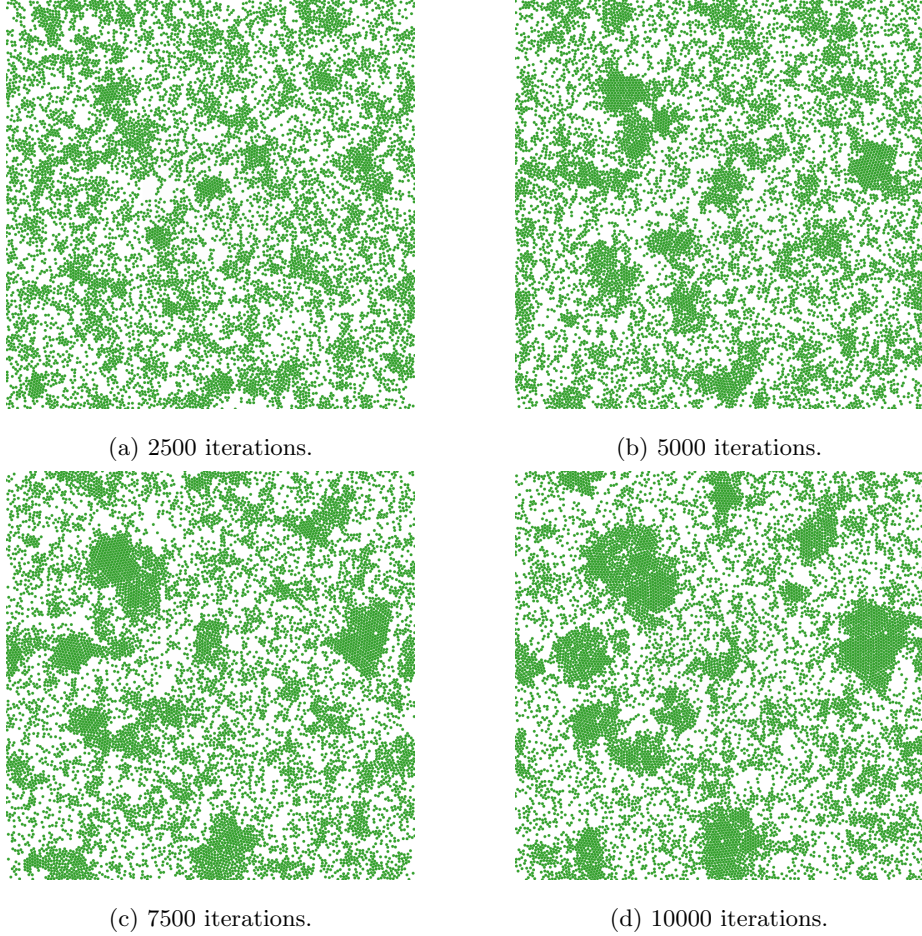


Figure 15: Snapshots of an ABP system rescaled to fit the packing fraction $\rho = 0.4$. The other parameters are the same as presented on Figure 14.

8. Discussion

We presented algorithms for simulating colloidal particles subject to Brownian motion, interacting with short or long range force interactions, and presenting excluded volume. The overlap correction algorithm using Delaunay triangulations is a novel method. The algorithms implemented in CUDA for simulation are fully parallel, transferring data back to host only for measurements or outputs. The overlap correction algorithm can be used independently from

the forces calculation, allowing to simulate different colloidal models including changed particles or self-propelled active systems. The Delaunay triangulation and the parallel edge-flip algorithm proved to be useful for solving overlaps efficiently. This opens the possibility for using the Delaunay triangulation for solving related problems in the simulation, such as short range force calculation or approximated n-body simulations. The parallel n-body implementation was also successfully adapted and optimized to the particular conditions of colloidal particles, which opens up simulations of up to two orders of magnitude the number of particles used on the previous sequential implementation.

Acknowledgements

The authors would like to thank the NVIDIA GPU Research Center of the Department of Computer Science of the Universidad de Chile for supplying the equipment used for the tests presented here. This work was partially supported by the FONDECYT projects No. 1140778 and No. 3160182, and by project No. ENL009/15, VID, Universidad de Chile.

References

References

- [1] W. B. Russel, D. A. Saville, W. R. Schowalter, Colloidal dispersions, Cambridge university press, 1989.
- [2] T. Vicsek, A. Zafeiris, Collective motion, Physics Reports 517 (3) (2012) 71–140.
- [3] M. C. Marchetti, J. Joanny, S. Ramaswamy, T. Liverpool, J. Prost, M. Rao, R. A. Simha, Hydrodynamics of soft active matter, Reviews of Modern Physics 85 (3) (2013) 1143.
- [4] P. Romanczuk, M. Bär, W. Ebeling, B. Lindner, L. Schimansky-Geier, Active brownian particles, The European Physical Journal Special Topics 202 (1) (2012) 1–162.

- [5] J. L. Anderson, Colloid transport by interfacial forces, *Annual review of fluid mechanics* 21 (1) (1989) 61–99.
- [6] C. A. Navarro, N. Hitschfeld-Kahler, L. Mateu, A survey on parallel computing and its applications in data-parallel problems using GPU architectures, *Communications in Computational Physics* 15 (2) (2014) 285–329.
- [7] B. Andreotti, Y. Forterre, O. Pouliquen, *Granular Media: Between Fluid and Solid*, Cambridge University Press, 2013.
- [8] T. Pöschel, T. Schwager, *Computational Granular Dynamics: Models and Algorithms*, Springer, 2005.
- [9] M. Weigel, Performance potential for simulating spin models on {GPU}, *Journal of Computational Physics* 231 (8) (2012) 3064–3082.
- [10] M. Weigel, Simulating spin models on {GPU}, *Computer Physics Communications* 182 (9) (2011) 1833–1836.
- [11] J. Bédorf, E. Gaburov, S. Portegies Zwart, A sparse octree gravitational N-body code that runs entirely on the GPU processor, *Journal of Computational Physics* 231 (7) (2012) 2825–2839.
- [12] R. Soto, R. Golestanian, Self-assembly of catalytically active colloidal molecules: Tailoring activity through surface chemistry, *Physical Review Letters* 112 (6).
- [13] R. Soto, R. Golestanian, Self-assembly of active colloidal molecules with dynamic function, *Physical Review E* 91 (5) (2015) 052304.
- [14] J. Barnes, P. Hut, A hierarchical $O(N \log N)$ force-calculation algorithm, *Nature* 324 (6096) (1986) 446–449.
- [15] Z. Yao, J. S. Wang, G. R. Liu, M. Cheng, Improved neighbor list algorithm in molecular simulations using cell decomposition and data

- sorting method, *Computer Physics Communications* 161 (1-2) (2004) 27–35.
- [16] M. P. Allen, D. J. Tildesley, *Computer simulation of liquids*, Oxford university press, 1989.
 - [17] D. Frenkel, B. Smit, *Understanding molecular simulation: from algorithms to applications*, Vol. 1, Academic press, 2001.
 - [18] M. S. Miguel, R. Toral, *Stochastic Effects in Physical Systems, Instabilities and Nonequilibrium Structures VI* 5 (2000) 35–120.
`arXiv:9707147`.
 - [19] L. Verlet, Computer "experiments" on classical fluids. I. Thermodynamical properties of Lennard-Jones molecules, *Physical Review* 159 (1) (1967) 98–103.
 - [20] T. J. Lipscomb, S. S. Cho, Parallel Verlet Neighbor List Algorithm for GPU-Optimized MD Simulations Categories and Subject Descriptors, *ACM Conference on Bioinformatics, Computational Biology and Biomedicine* (2012) 321–328.
 - [21] A. J. Proctor, T. J. Lipscomb, A. Zou, J. A. Anderson, S. S. Cho, Performance analyses of a parallel verlet neighbor list algorithm for gpu-optimized MD simulations, in: *Proceedings of the 2012 ASE International Conference on BioMedical Computing, BioMedCom 2012*, 2013, pp. 14–19.
 - [22] L. Nyland, M. Harris, J. Prins, Fast N-Body Simulation with CUDA, *Simulation* 3 (1) (2007) 677–696.
 - [23] M. Burtcher, K. Pingali, An efficient CUDA implementation of the tree-based barnes hut n-body algorithm, in: *GPU Computing Gems Emerald Edition*, 2011.

- [24] K. A. Hawick, D. P. Playne, Hard-sphere collision simulations with multiple GPUs, PCIe extension buses and GPU-GPU communications, *Conferences in Research and Practice in Information Technology Series* 127 (2012) 13–21.
- [25] P. Strating, Brownian dynamics simulation of a hard-sphere suspension, *Physical Review E* 59 (2) (1999) 2175–2187.
- [26] J. S. Vitter, Random sampling with a reservoir, *ACM Transactions on Mathematical Software* 11 (1) (1985) 37–57.
- [27] M. De Berg, O. Cheong, M. Van Kreveld, M. Overmars, *Computational Geometry: Algorithms and Applications*, Vol. 17, 2008.
- [28] C. L. Lawson, Transforming triangulations, *Discrete Mathematics* 3 (4) (1972) 365–372.
- [29] C. Navarro, N. Hitschfeld, E. Scheihing, Quasi-delaunay triangulations using gpu-based edge-flips, *Communications in Computer and Information Science* 458 (2014) 36–49.
- [30] The CGAL Project, *CGAL User and Reference Manual*, 4.8.1 Edition, CGAL Editorial Board, 2016.
- [31] N. Kruithof, 2D periodic triangulations, in: *CGAL User and Reference Manual*, 4.8.1 Edition, CGAL Editorial Board, 2016.
- [32] NVIDIA, *Cuda C Programming Guide*, Programming Guides (September) (2015) 1–261.
- [33] NVIDIA, *CURAND Library: Programming Guide*, Version 7.0 (2015).
- [34] S. Green, Particle simulation using cuda, *NVIDIA whitepaper* 6 (2010) 121–128.
- [35] Y. Fily, M. C. Marchetti, Athermal phase separation of self-propelled particles with no alignment, *Physical Review Letters* 108 (23) (2012) 00319007.