

An efficient solution for hazardous geophysical flows simulation using GPUs

A. Lacasta¹, C. Juez¹, J. Murillo¹, P. García-Navarro¹,

Abstract

The movement of poorly sorted material over steep areas constitutes a hazardous environmental problem. Computational tools help in the understanding and predictions of such landslides. The main drawback is the high computational effort required for obtaining accurate numerical solutions due to the high number of cells involved in the calculus. In order to overcome this problem, this work proposes the use of GPUs for decreasing significantly the CPU simulation time. The numerical scheme implemented in GPU is based on a finite volume scheme and it was validated in previous work with exact solutions and experimental data. The computational cost time obtained with the Graphical Hardware technology, GPU, are compared against single-Core (sequential) and multi-Core (parallel) CPU implementations. The GPU implementation allows to reduce the computational cost time in two orders of magnitude.

Keywords: CUDA; GPU; Landslides; Numerical modeling; Shallow Flow; Coulomb forces

1. Introduction

2 Landslides play an important role on the evolution of landscape and con-
3 stitute an important environmental risk. They can be responsible for dra-
4 matic civilian damages and that is the reason why the building of defenses
5 and barriers is required. The computational tools are a suitable partner for

Email addresses: alacasta@unizar.es (A. Lacasta), carmelo@unizar.es (C. Juez),
javier.murillo@unizar.es (J. Murillo), pigar@unizar.es (P. García-Navarro)

¹LIFTEC-EINA, CSIC-Universidad de Zaragoza, Spain

6 developing a careful design of such elements. Over recent years, reliable pre-
7 dictions of the spreading of granular material have been obtained (Pirulli
8 et al., 2007; Pirulli and Mangeney, 2008; Moretti et al., 2012; Bouchut et al.,
9 2008) and numerical results have been validated against series of experiments
10 based on granular dry flows (Savage and Hutter, 1989; Iverson and Denlinger,
11 2001; Pouliquen and Forterre, 2002; Lajeunesse et al., 2004; Mangeney et al.,
12 2010). In particular, Murillo and García-Navarro (2012); Juez et al. (2013)
13 have recently presented a robust finite volume upwind scheme which includes
14 the presence of steep slopes leading to obtain promising results.

15 Once the forecasting capability of the computational tool has been achieved
16 another important concern is the improvement of the efficiency in terms of
17 the computational cost. This type of geophysical flow involves the study of
18 huge domains, such as catchments, mountains or gorges, where an accurate
19 Digital Terrain Model is required in order to mimic the complex topography
20 of the terrain. For this reason, implementations based on the most recent
21 GPU hardware (Graphics Processing Unit) emerges as a promising strategy
22 for handling this environmental and up to date problem. This hardware con-
23 sists of thousands of simple arithmetic processing units originally designed
24 to pixel-based computations. In the recent years, this technology has been
25 successfully extended and applied to more general problems in order to ac-
26 celerate them. This concept is also known as General Purpose Computing
27 on Graphics Processing Units (GPGPU) (Owens et al., 2007).

28 In terms of scientific computation, the last four decades have followed
29 Moore’s law (Moore, 2003), i.e. the number of transistors on a chip increases
30 exponentially. This integration allowed to obtain faster and faster applica-
31 tions, by recompiling the code for these new processors. Unfortunately, power
32 has become the primary design constraint for chip designers, where both en-
33 ergy and power dissipation create a technological barrier for the integration
34 capacity (Dreslinski et al., 2010). Nevertheless, Multi-Core micro architec-
35 ture together with an adequate programming model (OpenMP is one of the
36 most extended) brings a chance to exploit the parallelism of some parts of
37 the code (Sharma and Gupta, 2013). Moreover, the Multi-Core paradigm has
38 a large power consumption rate when performing small tasks and, for these
39 purposes, Many-Core systems appear to be a very interesting option (Borkar,
40 2007). Many-Core architectures are those composed of smaller and not so
41 complex cores that usually have special purposes. Industrial implementation
42 of this solution has been obtained in the field of Graphical Processing, where
43 several efforts have been devoted to make more powerful devices. Indeed,

44 this technology has been historically oriented to the very particular task of
 45 performing shading operations when rendering graphics. Following Lacasta
 46 et al. (2014), the purpose of the present work is to apply this hardware to
 47 the simulation of hazardous and very high time consuming geophysical flows.

48 The outline of this work is as follows: Section 2 is devoted to explain the
 49 mathematical model and numerical scheme used for modeling and solving the
 50 landslides behavior. In Section 3, the implementation on GPU is described.
 51 Section 4 gathers several experimental and realistic cases considered for test-
 52 ing the GPU performance. The differences on the computational cost time
 53 between the sequential, parallel and GPU strategies is discussed in Section
 54 5. Finally in Section 6 the conclusions are summarized.

55 2. Mathematical model and numerical scheme

56 2.1. Mathematical model

57 The mathematical model considered for reproducing the landslides phe-
 58 nomenon is based on the shallow flow equations, where the general three-
 59 dimensional conservation laws are depth averaged. The pressure distribution
 60 is considered hydrostatic and as frictional terms, only Coulomb type friction
 61 forces are assumed, Juez et al. (2013). Bearing in mind these hypothesis, the
 62 2D equations are written in global coordinates as follows:

$$\frac{\partial \mathbf{U}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{U})}{\partial x} + \frac{\partial \mathbf{G}(\mathbf{U})}{\partial y} = \mathbf{S}_\tau + \mathbf{S}_b \quad (1)$$

63 where

$$\mathbf{U} = (h, hu, hv)^T \quad (2)$$

64 are the conserved variables with h representing granular material depth par-
 65 allel to the z coordinate and (u, v) the depth averaged components of the
 66 velocity vector. The fluxes are given by

$$\begin{aligned} \mathbf{F} &= \left(hu, hu^2 + \frac{1}{2}g_\psi h^2, huv \right)^T \\ \mathbf{G} &= \left(hv, huv, hv^2 + \frac{1}{2}g_\psi h^2 \right)^T \end{aligned} \quad (3)$$

67 with $g_\psi = g \cos^2 \psi$ and ψ the direction cosine of the bed normal with respect
 68 to the vertical axis. The physical basis of this gravity projection is explained

69 in detail in Juez et al. (2013) and it is of utmost importance for keeping
70 accurate numerical predictions when the simulation involves the presence of
71 steep slopes and consequently, the direction cosines become relevant.

72 The term \mathbf{S}_τ notes the frictional effects in the bed, and is defined as

$$\mathbf{S}_\tau = \left(0, -\frac{\tau_{b,x}}{\rho}, -\frac{\tau_{b,y}}{\rho} \right)^T \quad (4)$$

73 with $\tau_{b,x}, \tau_{b,y}$ the bed shear stress in the x and y directions respectively and
74 ρ the density of the granular mass. In this work, only dense granular flows
75 are considered, therefore, the main rheological properties are governed by the
76 frictional forces. These interactions between the sand grains are computed
77 by means of the Coulomb law. This formula is based on the internal friction
78 angle of the material, θ_b .

79 On the other hand, the term \mathbf{S}_b is defined for gathering the information
80 relative to the pressure force exerted over the bottom.

81 Thanks to the hyperbolic character of (1) it is possible to obtain a Jaco-
82 bian matrix, \mathbf{J}_n , which is built by means of the flux normal to a direction
83 given by the unit vector \mathbf{n} , $\mathbf{E}_n = \mathbf{F}n_x + \mathbf{G}n_y$,

$$\mathbf{J}_n = \frac{\partial \mathbf{E}_n}{\partial \mathbf{U}} = \frac{\partial \mathbf{F}}{\partial \mathbf{U}} n_x + \frac{\partial \mathbf{G}}{\partial \mathbf{U}} n_y \quad (5)$$

84 whose components are

$$\mathbf{J}_n = \begin{pmatrix} 0 & n_x & n_y \\ (g_\psi h - u^2)n_x - uvn_y & vn_y + 2un_x & un_y \\ (g_\psi h - v^2)n_y - uvn_x & vn_x & un_x + 2vn_y \end{pmatrix} \quad (6)$$

85 The eigenvalues of this Jacobian matrix constitute the basis of the upwind
86 technique which is detailed in the next subsection.

87 2.2. Numerical scheme

88 System in (1) is integrated in a grid cell Ω_i and the Gauss theorem is
89 applied, being the normal vector \mathbf{n} outward to the cell Ω_i , as displayed in
90 Figure 1

$$\frac{\partial}{\partial t} \int_{\Omega_i} \mathbf{U} d\Omega + \oint_{\partial\Omega_i} \mathbf{E}_n dl = \int_{\Omega_i} (\mathbf{S}_\tau + \mathbf{S}_b) d\Omega \quad (7)$$

91 The second integral in (7) can be explicitly expressed as a sum over the
92 cell edges,

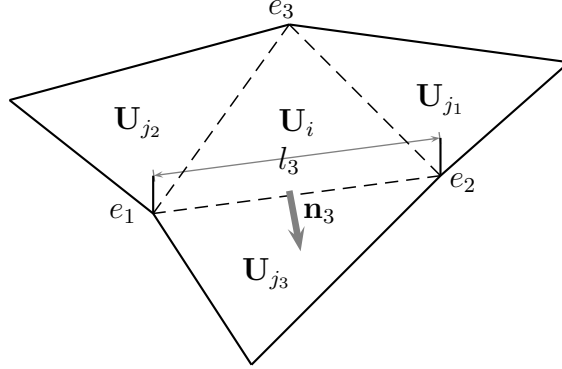


Figure 1: Cell parameters

$$\frac{\partial}{\partial t} \int_{\Omega_i} \mathbf{U} d\Omega + \sum_{k=1}^{NE} (\mathbf{E}_{\mathbf{n}})_k l_k = \sum_{k=1}^{NE} \mathbf{S}_{\mathbf{n}\tau} l_k + \sum_{k=1}^{NE} \mathbf{S}_{\mathbf{n}b} l_k \quad (8)$$

93 where l_k is the corresponding edge length and $\mathbf{S}_{\mathbf{n}b}$ and $\mathbf{S}_{\mathbf{n}\tau}$ are suitable inte-
 94 grals of the bed slope and friction source terms (Juez et al., 2013),

$$\mathbf{S}_{\mathbf{n}b} = \left(0, -g_\psi h \frac{\partial z}{\partial x} n_x, -g_\psi h \frac{\partial z}{\partial y} n_y \right)^T \quad (9)$$

$$\mathbf{S}_{\mathbf{n}\tau} = \left(0, \rho g_\psi h \tan \theta_b n_x, \rho g_\psi h \tan \theta_b n_y \right)^T \quad (10)$$

95 where n_x and n_y are the components of the unit vector \mathbf{n} of each edge of each
 96 computational cell.

97 The numerical scheme is constructed by defining a local linearization in
 98 terms of an approximate Jacobian matrix $\tilde{\mathbf{J}}$ at each k edge between neigh-
 99 boring cells defined through the normal flux $\mathbf{E}_{\mathbf{n}}$

$$(\delta \mathbf{E}_{\mathbf{n}})_k = \tilde{\mathbf{J}}_{\mathbf{n},k} \delta \mathbf{U}_k \quad (11)$$

100 with $\delta(\mathbf{E}_{\mathbf{n}})_k = (\mathbf{E}_j - \mathbf{E}_i)_{\mathbf{n}_k}$, $\delta \mathbf{U}_k = \mathbf{U}_j - \mathbf{U}_i$, and \mathbf{U}_i and \mathbf{U}_j the initial values
 101 at cells i and j sharing edge k .

102 From this approximate Jacobian matrix a set of three real eigenvalues $\tilde{\lambda}_k^m$
 103 and eigenvectors $\tilde{\mathbf{e}}_k^m$ are obtained. The vector of conserved variables, $\delta \mathbf{U}$, is
 104 projected onto the matrix eigenvectors basis, $\tilde{\mathbf{P}}$, as

$$\delta \mathbf{U}_k = \tilde{\mathbf{P}}_k \tilde{\mathbf{A}}_k = \sum_{m=1}^3 (\tilde{\alpha} \tilde{\mathbf{e}})_k^m \quad (12)$$

105 with

$$\tilde{\mathbf{P}}_k = (\mathbf{e}^1, \mathbf{e}^2, \mathbf{e}^3)_k \quad \tilde{\mathbf{A}}_k = (\alpha^1 \quad \alpha^2 \quad \alpha^3)_k^T \quad (13)$$

106 The source terms are also projected onto the matrix eigenvectors basis,
 107 $\tilde{\mathbf{P}}$, to guarantee the exact equilibrium between fluxes and source terms, i.e.
 108 well-balanced property, Murillo and García-Navarro (2010).

$$(\mathbf{S}_{\mathbf{n},b}, \mathbf{S}_{\mathbf{n},\tau})_k = \tilde{\mathbf{P}}_k \tilde{\mathbf{B}}_k = \sum_{m=1}^3 (\tilde{\beta} \tilde{\mathbf{e}})_k^m \quad (14)$$

109 with

$$\tilde{\mathbf{B}}_k = (\beta^1 \quad \beta^2 \quad \beta^3)_k^T \quad (15)$$

110 Additionally, it is worth remarking that two numerical fixes are also
 111 needed for computing physical solutions (Murillo and García-Navarro, 2012):
 112 they are necessary for avoiding unphysical solutions by means of decreasing
 113 the wave source strengths, \mathbf{B}_k .

114 Gathering all the previous information the volume integral in the cell at
 115 time t^{n+1} is expressed as

$$\mathbf{U}_i^{n+1} = \mathbf{U}_i^n - \sum_{k=1}^{NE} \sum_{m=1}^3 (\tilde{\lambda}^- \tilde{\alpha} - \tilde{\beta}^-)_k^m \mathbf{e}_k^m l_k \frac{\Delta t}{A_i} \quad (16)$$

116 The minus sign in (16) implies that only the incoming waves are con-
 117 sidered for updating the values of each cell. The numerical scheme is only
 118 computed over the cells where there is a sand depth.

119 In order to guarantee stability the time step Δt has to be taken small
 120 enough so that there are no interactions of waves between neighboring cells.

$$\Delta t \leq CFL \Delta t^{\tilde{\lambda}} \quad \Delta t^{\tilde{\lambda}} = \frac{\min(\chi_i, \chi_j)}{\max|\tilde{\lambda}^m|} \quad (17)$$

121 with $CFL=1/2$ in the case of triangular unstructured grids. The term χ_i is
 122 the relevant distance, which in a 2D framework must consider the volume of
 123 the cell i and the length of the shared k edges (Murillo and García-Navarro,
 124 2010),

$$\chi_i = \frac{A_i}{\max_{k=1,NE} l_k} \quad (18)$$

125 Due to the fact that the numerical scheme is controlled by a explicit global
 126 time step, the only way of reducing the computational cost is by means of
 127 using a high performance computing technique: the implementation of the
 128 numerical scheme on a GPU, as it is explained in the next section.

129 3. Implementation of the numerical scheme

130 The implementation of the model follows a time-stepping process where,
 131 until the simulation time is achieved, the numerical scheme (16) is applied
 132 over the whole domain. The way this process is performed is illustrated in
 133 Figure 2. The main operations of the numerical scheme are displayed in green
 134 and pink. The calculation of the fluxes is performed looping edgewise and the
 135 updating process as well as the boundary calculation are performed looping
 136 over the cells. Both processes go from 1 to n_{edges} and n_{cells} respectively.

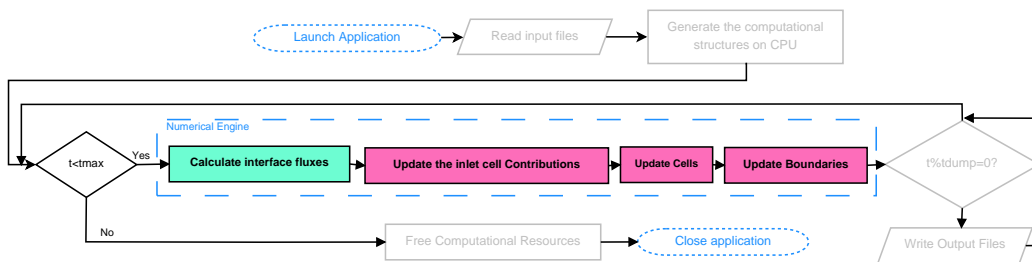


Figure 2: Flowchart of the simulation process. The green function is performed looping over cell edges, while the pink is a loop over the cells.

137 In order to make the process faster, some of the parts can be paralellized
 138 or even adapted to perform the operations in the GPU. For succeeding in this
 139 task, we know that most of the processors are designed with multi-threading
 140 capabilities. Despite these capabilities, it is not direct to perform operations
 141 using every single core, since the calculus operations have to be shared among
 142 all the cores. Figure 3 displays a typical iterative process using a single core
 143 of a Multi-Core processor. It represents how one core is activated and it has
 144 all the workload while the rest are not active and they do not participate in
 145 the calculations.

146 When possible, the iterative process may be split into different sub-
 147 processes allowing the processing of different elements at the same time by
 148 different processing units. This is the concept of parallelization. It can be di-
 149 rectly applied when the loop contains straight-line code (a single basic block

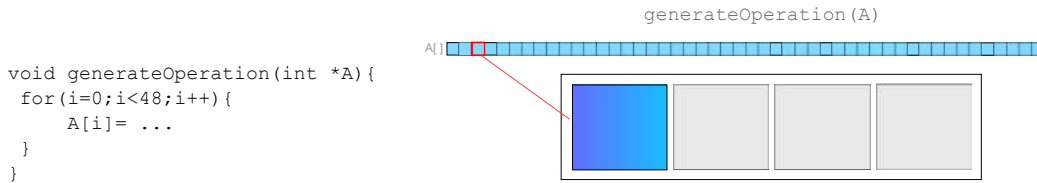


Figure 3: Sequential iteration over vector $A[]$

150 with no jumps) and when there is no data dependency between iterations
 151 (element $i - 1$ is not required for the element i). If data flow does not satisfy
 152 the previous condition, the manner of processing the elements to make them
 153 parallelizable must be re-structured.

154 OpenMP is a very useful standard to perform shared-memory paralleliza-
 155 tion (Chapman et al., 2007). This standard implements parallel primitives
 156 in many programming languages such as C or Fortran making it very simple
 157 to implement a parallel solution in a reasonable time. The main disadvan-
 158 tage of this kind of parallelization is that the Multi-Core processors have a
 159 physical boundary, since the number of transistors which can be used within
 160 the space of a chip are limited. Contrary, with the Many-Core architectures
 161 devices, such as the GPUs, a lot of specialized processing units allow to make
 162 many more operations than in the multi-core unit. An example of a parallel
 163 loop using OpenMP is displayed in Figure 4.

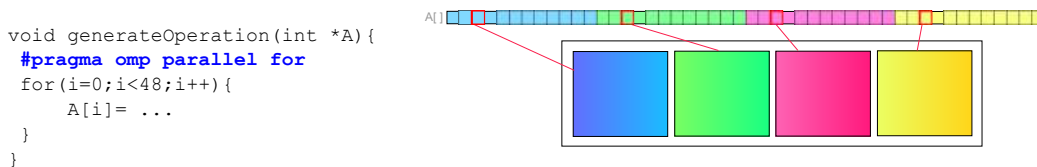


Figure 4: Parallel iteration over vector $A[]$ using 4 cores

164 Many-Core processors provide a hardware architecture to perform a large
 165 number of independent operations in a parallel manner. The main difference
 166 between Multi-Core threads and Many-Core threads is the *lightweight* of the
 167 second solution as well as the hardware support of a larger number of the
 168 latter ($> 1000s$). Unlike the Multi-Core thread, where each thread processes
 169 a group of elements, in the Many-Core paradigm, each thread will process
 170 just one element. Moreover, the creation and management of the Many-Core
 171 threads as well as their capacity to perform operations are simpler than in
 172 the case of the common Multi-Core thread.

173 CUDA (Compute Unified Device Architecture) was released in order to
 174 perform this kind of parallelization on NVIDIA GPUs (NVIDIA, 2011). It
 175 provides both a software model and a set of compilation tools that support
 176 the NVIDIA Many-Core GPUs. An example of a CUDA instance for the
 177 previous loop is displayed in Figure 5.

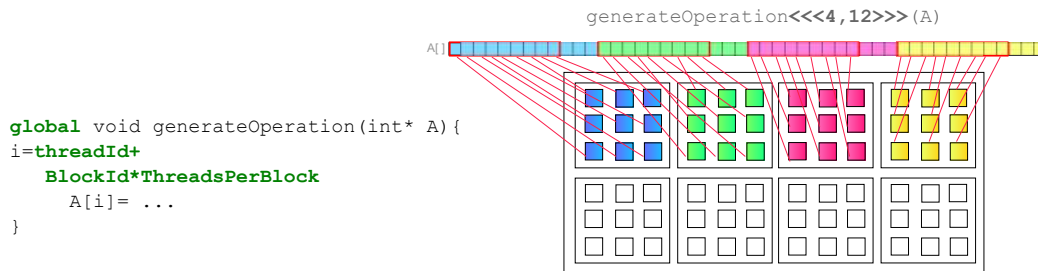


Figure 5: Massively parallel iteration over vector $A[]$ using a Many-Core architecture.

178 The improvement on the implementation of the numerical solver will be
 179 achieved by means of the translation of the processes sketched in Figure
 180 2 into both solutions: OpenMP for the Shared-Memory level parallelism in
 181 Multi-Core hardware and the CUDA based solution for the Many-Core GPUs
 182 architecture.

183 3.1. Implementation on GPU

184 The GPU contains a large number of processors working all together ap-
 185 plying the same operation over different elements. In order to obtain high
 186 performance in the GPU implementation it is necessary to understand the
 187 way the model works. NVIDIA GPUs are formed by Streaming Multipro-
 188 cessors (SMs). Each Multiprocessor is composed by Streaming Processors
 189 (SPs) that represent the minimum processing unit. In the case of the Tesla
 190 Series GPUs, there are between 14 and 16 SMs and each one is composed
 191 by 32 SPs (Glaskowsky, 2009). From the software point of view, there are
 192 groups of blocks that contain threads. More specifically, each block contains
 193 $blockDim$ threads and the number of blocks $gridDim$ must be larger than
 194 the number of elements (n_{elem}) to be processed ($n_{elem} \leq blockDim * gridDim$).
 195 Each block will be processed by a SM in groups of 32 elements which form
 196 a *warp* taking into account that each Streaming Processor will process each
 197 element of the warp. More details about CUDA and NVIDIA GPUs can be
 198 found in (NVIDIA, 2011).

199 Implementation for GPU has been developed using the approach of La-
 200 casta et al. (2014) for unstructured meshes and adapting the required ele-
 201 ments into this new model. A general overview of the CUDA implementation
 202 of the simulation process described in Figure 2 is highlighted in Listing 1.

Listing 1: Overview of the CUDA implementation.

```

203 1 ...
204 2 // Configuration of the parameters
205 3 Threads=256;
206 4 edgeBlocks=nEdge/Threads;
207 5 cellBlocks=nCell/Threads;
208 6 while(t<tmax){
209 7   // Calculate the fluxes
210 8   calculateEdgeFluxes<<<edgeBlocks,Threads,0,executionStream>>>(...);
211 9   // Establish the minimum dt obtaining the ID of the
212 0   // minimum dt
213 1   // (*) Explained at Listing 3
214 2   cublasIdamin(...,nEdge,vDt,1,id);
215 3   // And assign it
216 4   newDt<<<1,1,0,executionStream>>>(dt,vDt,id);
217 5   // Update the elapsed time (in GPU)
218 6   updateT<<<1,1,0,executionStream>>>(cuda_t,dt);
219 7   // Retrieves the value of t to CPU
220 8   cudaMemcpy(t,cuda_t,sizeof(double),cudaMemcpyDeviceToHost);
221 9   // Update the cell values
222 0   assignFluxes<<<cellBlocks,Threads,0,executionStream>>>(...);
223 1   // Verify whether it is necessary to dump data and
224 2   // if so, process it.
225 3   // (*) Detailed in Listing 4
226 4   if(t<t_dump){
227 5     // Transfer data from GPU to CPU and write it
228 6     // into an output file
229 7   }
230 8 }

```

231 As it was justified in Juez et al. (2013) the topology of the mesh plays
 232 an important role on the quality of the numerical results. Only unstructured
 233 meshes avoid the presence of preferential directions. Because of this necessity,
 234 it is not straightforward to obtain an efficient solution for the GPU process-
 235 ing. GPUs are designed to work efficiently with ordered information. This
 236 means that, those structures without ordered pattern may reduce the overall
 237 performance of the implementation (Lacasta et al., 2014). Moreover, special
 238 attention on those parts that are not intrinsic parallel is required (such as
 239 the selection of the global time step).

240 To obtain an efficient solution three special considerations are highlighted.

- 241 • Using Structure of Arrays (SoA) instead of Arrays of Structures (AoS)
242 to store data. Because of the manner of processing the elements, SoA
243 provides improvements when accessing consecutive elements and this
244 will be maximized if conserved variables (h, hu, hv) are stored consec-
245 utively as $h_0 \dots h_{ncells} hu_0 \dots hu_{ncells} hv_0 \dots hv_{ncells}$ instead of $h_0 hu_0 hv_0 \dots h_n hu_n hv_n$.
246 When a thread within a warp reads the variable $h(i)$ it is likely that
247 the following thread requires $h(i + 1)$. In other words, this approach
248 improves the data locality and allows to the GPU to obtain coalesced
249 memory access.

- 250 • To maximize coalescence when accessing over cells it is required to re-
251 order the mesh (Lacasta et al., 2014). It is applied by means of modify-
252 ing the way the cells are connected among themselves. One strategy is
253 to fit the connectivity matrix into a banded matrix. Figure 6 shows the
254 effect of applying a reordering method in the connectivity matrix when
255 using unstructured meshes. In order to obtain the previous ordering,
256 bandwidth reduction methods are very useful. RCM algorithm is one
257 of them. These methods have been traditionally used to obtain better
258 performance in iterative methods by means of reducing the bandwidth
259 of the matrix. Similarly, the idea can be applied to the connectivity
260 matrix and this provides a reordered cells indexing that improves the
261 spatial locality and provides coalescence when accessing by cells.

- 262 • The ordering when processing over edges. In the previous point, RCM
263 achieves coalesced access when processing over cells, but it is not enough
264 to process over edges. Actually, it is likely that when calculating by
265 edges (when accessing the two neighboring cells of the edge) this re-
266 ordered mesh may introduce penalty (Lacasta et al., 2014). To avoid
267 this, reordering the edges (reordering the pair of cells) will strongly
268 increase the performance of the GPU memory access when accessing
269 over cells to the primitives values. The edge ordering is graphically
270 detailed in Figure 7. This reordering will allow to process the pair of
271 cells $(0, 1)$ for the thread which process edge 0 and the pair $(0, 2)$ for
272 the thread which process edge 1 instead of the pair $(19, 23)$ that would
273 be processed in the edge 1 without the edge ordering.

274 When applying all these previous improvements it is possible to achieve

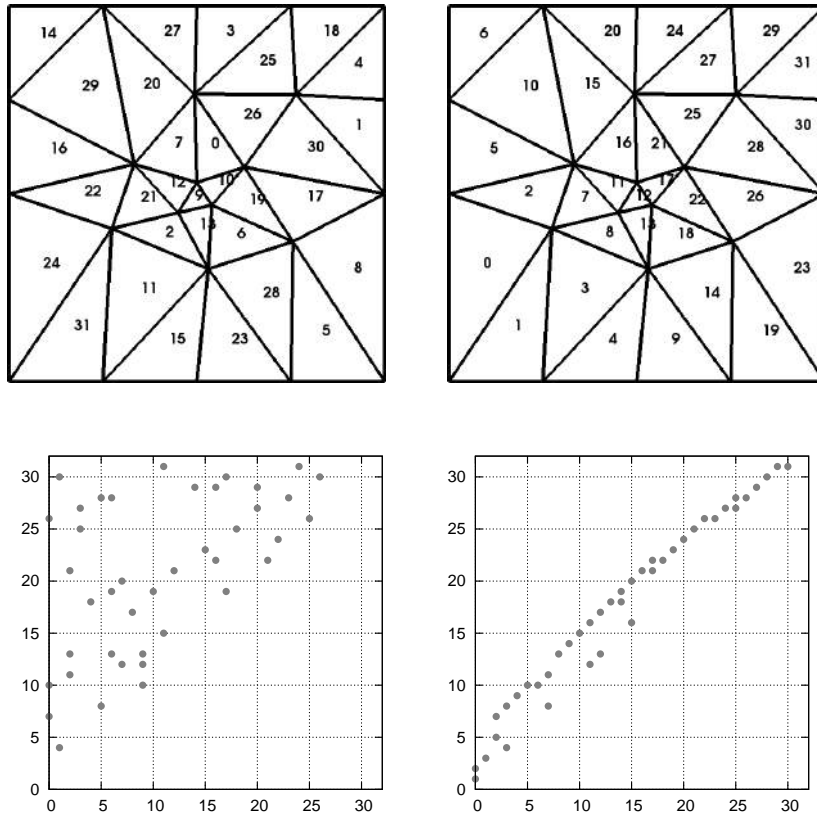


Figure 6: Upper: Unstructured mesh (left) and unstructured mesh postprocessed with RCM algorithm (right). Lower: Connectivity matrix for the original unstructured mesh (left) and for the reordered unstructured mesh (right). (Note that there exists a point if $m(i, j) = 1$).

275 great speed-ups comparing with a sequential or a multi-core implementation
 276 of the numerical scheme. This is discussed in Section 5.

277 4. Test cases

278 This section is devoted to describe the test cases considered as bench-
 279 mark for evaluating the performance of GPU implementation with respect
 280 to Single-Core and Multi-Core CPU implementations. The first three tests
 281 are based on experimental works previously used for validating the numer-
 282 ical scheme (Juez et al., 2013). The last test case considered constitutes a

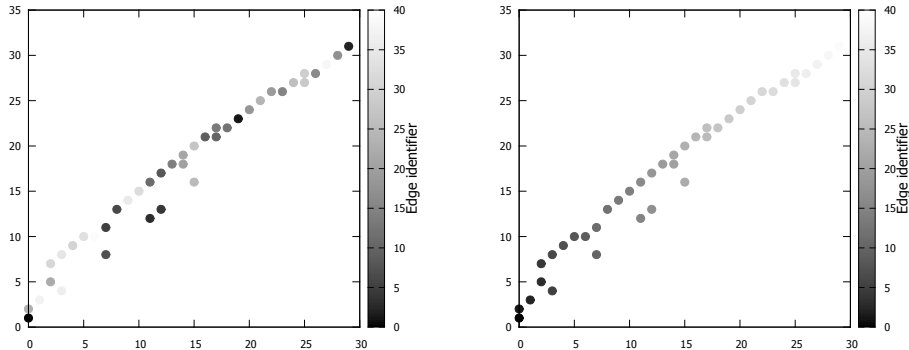


Figure 7: Unstructured mesh RCM processed edges ordering. Original ordering (left), `qsort` postprocessed ordering (right)

283 practical application, based on a real topography which is located in Spain
 284 (García-Ruiz et al., 2005).

285 All the numerical results have been computed using unstructured meshes
 286 to avoid numerical effects associated to the presence of preferential flow prop-
 287 agation directions as it was justified in Juez et al. (2013). Additionally, the
 288 CFL imposed in the numerical scheme has been equal to 0.5 in all the cases.

289 4.1. Spreading of cylindrical granular mass

290 The numerical simulation described in this test case is related with the
 291 experimental work developed by Lajeunesse et al. (2004). It has been chosen
 292 due to the simplicity of the initial configuration. During this experiment
 293 the spreading of a granular mass, originally retained within a cylinder, was
 294 recorded over a flat plane. The material was characterized with an internal
 295 friction angle $\theta_b=32^\circ$. The test case considered is Test A, where the initial
 296 height of the sand was 39.48 mm and the radius was 70.5 mm. The number
 297 of cells involved in the simulation has been 320000.

298 Figure 8 displays a 3D view of the temporal evolution of the sand depth
 299 and velocity, from the beginning, when the mass is enclosed within the cylin-
 300 der up to the final stage when all the granular mass has achieved an equilib-
 301 rium condition and it was stopped.

302 On the other hand, if we try to achieve a CPU simulation with a compu-
 303 tational cost time similar to the one provided by the GPU, the mesh has to

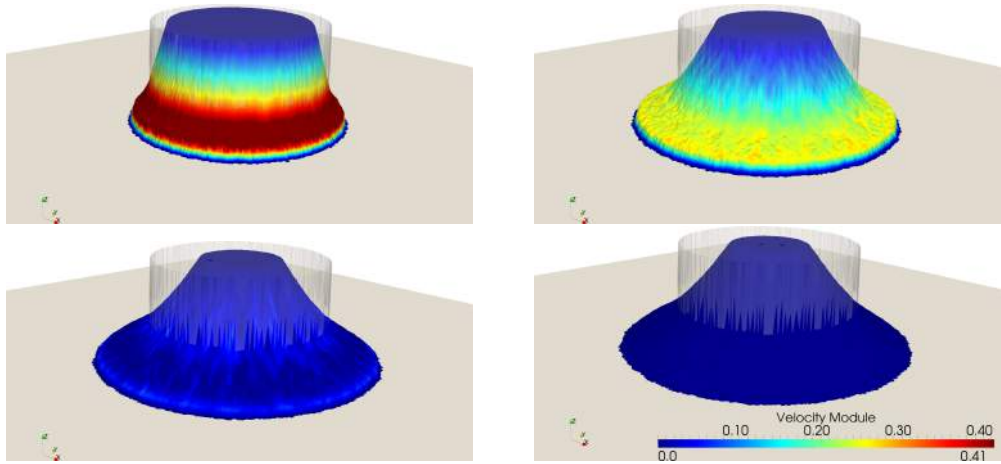


Figure 8: Evolution of the Sand Depth with the velocity (m/s) at time $t = 0.05s$. (top-left), $t = 0.15s$. (top-right), $t = 0.25s$. (bottom-left) and $t = 0.4s$. (bottom-right). The transparent element corresponds to the initial condition.

304 be consequently coarser and the results suffer from a lack of accuracy. This
 305 idea is showed in Figure 9, where the coarse mesh is composed by 16000 cells.

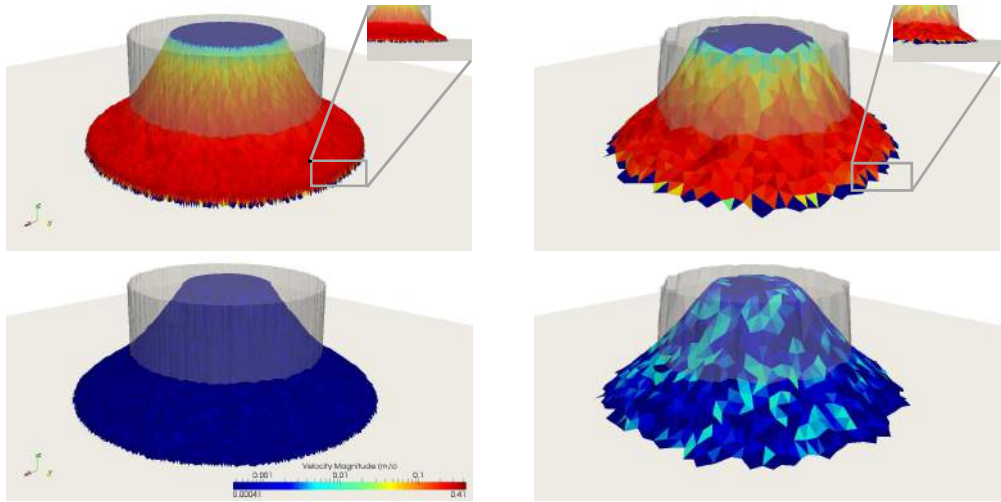


Figure 9: Comparison of fine mesh (left) and coarse mesh (right) of velocity magnitude at $t = 0.1s$. (upper) and $t = 0.4s$. (lower)

306 *4.2. Spreading of a granular mass including the presence of an obstacle*

307 The development of this experiment was motivated by the presence of
 308 obstacles in real events, i.e. avalanches. The presence of these obstacles
 309 causes the birth and propagation of quick moving shocks. The experimental
 310 work was carried out by Juez et al. (2014). The experimental setup consists
 311 of a rough inclined plane with a changing slope and with an obstacle in the
 312 middle of the flow path. The initial condition is given by a semi spherical
 313 cap full of sand. When this cap is pulled out, the granular material is free to
 314 evolve downwards from its original position until surrounding the obstacle.
 315 The internal friction angle of the material was characterized as $\theta_b = 26^\circ$ and
 316 the number of cells included in the simulation was 460000.

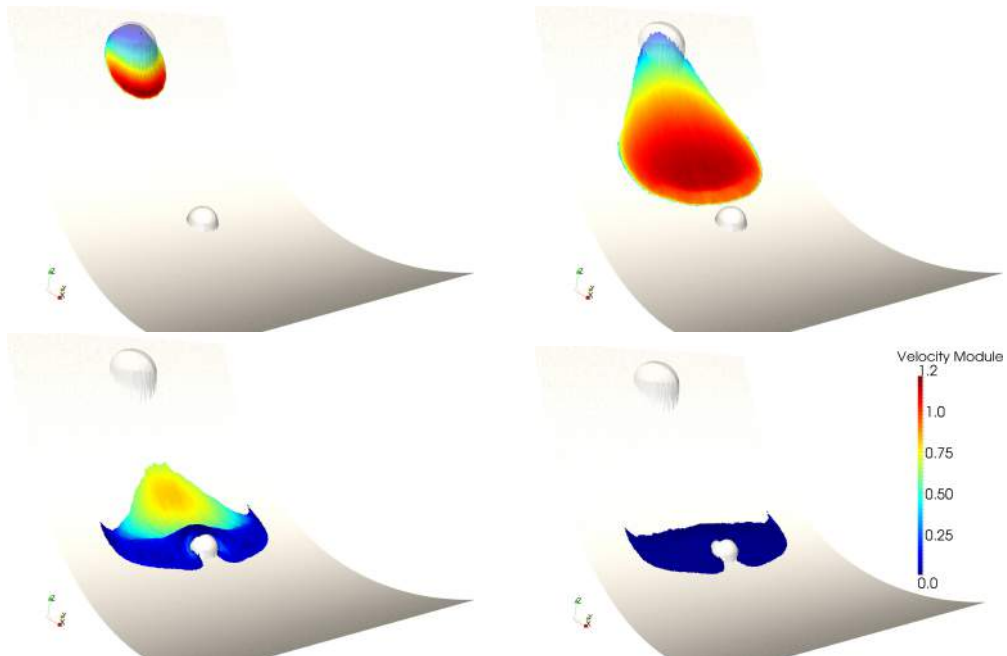


Figure 10: Evolution of the Sand Depth with the velocity (m/s) at time $t = 0.05s$. (top-left), $t = 0.4s$. (top-right) $t = 0.8s$. (bottom-left) and $t = 1.2s$. (bottom-right). The transparent element corresponds to the initial condition.

317 The temporal evolution of the granular flow is displayed in Figure 10.
 318 The front of the avalanches moves quickly until reaching the obstacle. Then,
 319 the flow is divided in two directions up to achieve a static equilibrium stage.

320 *4.3. Spreading of a granular mass over a parabolic chute*

321 This test case is another step in complexity, since the bed slope changes
322 in the longitudinal and transversal direction. The original experiment was
323 carried out by (Gray et al., 1999) and it was based on a semi spherical cap full
324 of sand which was suddenly removed over a parabolic chute. The material
325 was characterized with a internal friction angle $\theta_b = 30^\circ$. The number of cells
326 considered in the calculation has been 670000.

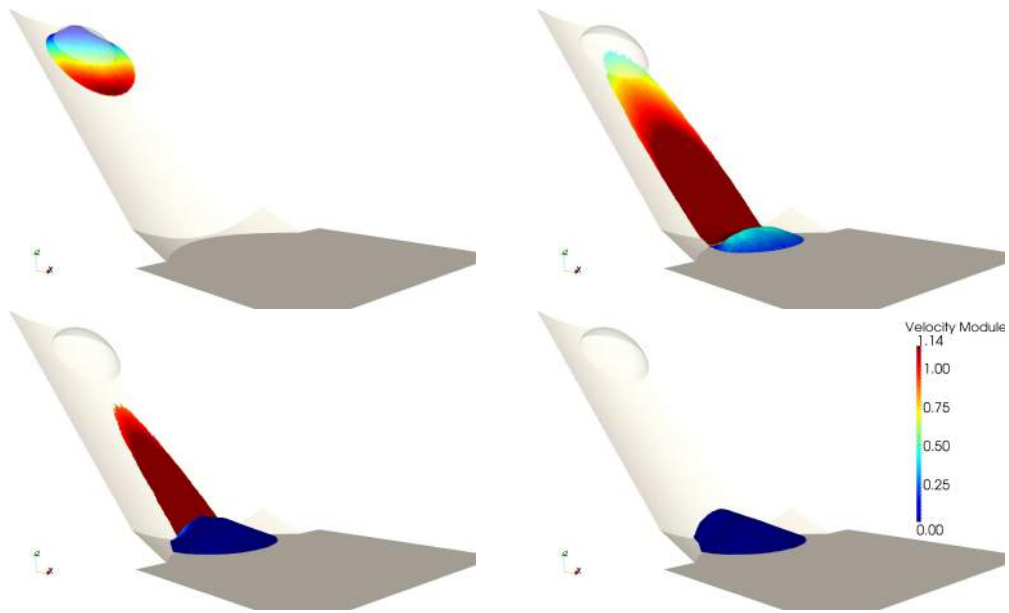


Figure 11: Evolution of the depth with the velocity (m/s) at time $t = 0.2s$. (top-left), $t = 1.4s$. (top-right), $t = 2s$. (bottom-left) and $t = 2.6s$. (bottom-right). The transparent element corresponds to the initial condition.

327 Numerical results are plotted in Figure 11. The material flows quickly
328 downstream until it reaches the horizontal area. At that position, the velocity
329 of the front is reduced by the lower bed slope. Hence, the granular mass starts
330 to spread transversally and finally, an uneven surface level is obtained at the
331 equilibrium stage.

332 *4.4. Landslide in a practical application*

333 The final test is based on a real topography of a catchment (García-
334 Ruiz et al., 2005). The Arnas catchment is located in the northern Spain

335 Pyrenees, in the Borau valley, and has a surface of 2.84 km^2 , ranging in alti-
 336 tude from around 900 to 1340 meters. Geologically, the catchment lies over
 337 Eocene flysch formations and has suffered land use and coverage changes
 338 in recent decades, generating a mixed vegetation cover which ranges from
 339 forest patches, dense and open shrubs, grassland cover and bare land. The
 340 assumption made by the authors is that the soil is composed by poorly sorted
 341 material and the idea is to verify the maximum run out and potential con-
 342 sequences of a massive mobilization of that material. Figure 12 displays the
 343 fixed bed rock and the soil mass which can be potentially mobilized. Due
 344 to the large dimensions of the catchment, the number of cells involved in
 345 the calculation is over 869000, making this type of phenomena a suitable
 346 candidate for GPU strategy.

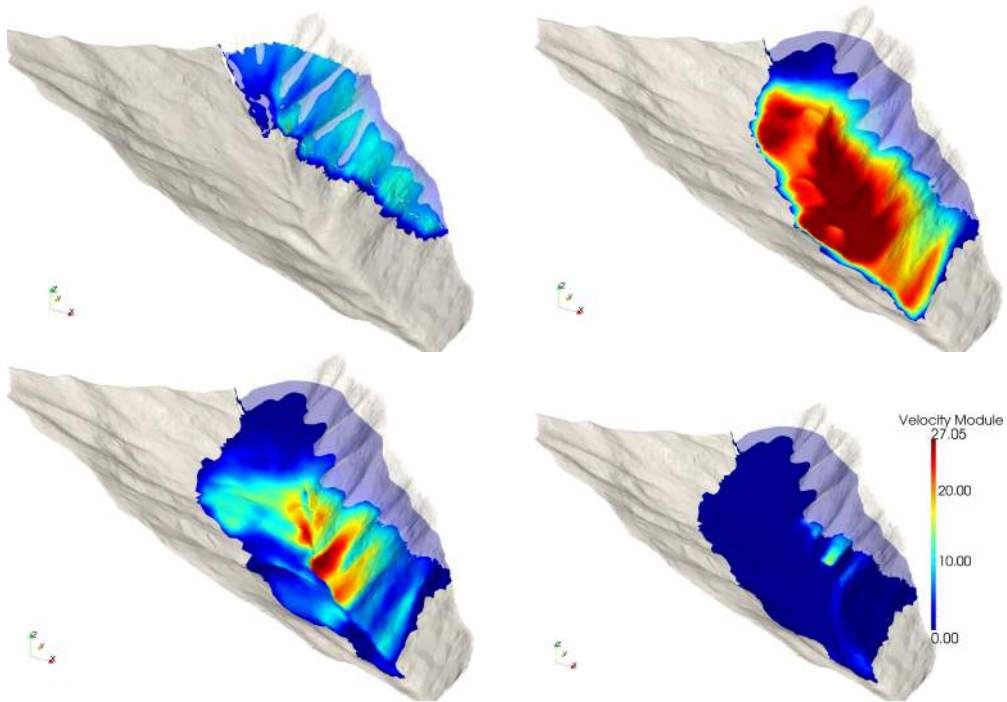


Figure 12: Evolution of the depth with the velocity (m/s) at time $t = 1s$. (top-left), $t = 2s$. (top-right), $t = 4s$. (bottom-left) and $t = 54s$. (bottom-right). The transparent element corresponds to the initial condition.

347 5. Results and computational load

348 As it has been mentioned before, Tests 1, 2 and 3 were previously used
349 by the authors for CPU computations. In such situations, the computational
350 cost was huge in comparison with the simulation time, ≈ 100 times the phys-
351 ical time. The cause of this problem has a twofold nature: first, the number
352 of cells involved in the computation is large, since the scale of interest in
353 these tests was in the order of millimeters. Additionally, due to the explicit
354 numerical scheme considered, there is a restriction in the time step used for
355 updating the solution. In these test cases, the time step is highly penalized
356 due to the small area of the cells and the fast flows, as it has been explained
357 in 17. On the other hand, the last test case does not require such high
358 refinement because of the uncertainty of the data acquisition. Indeed, the
359 mesh has been designed using the most accurate available LIDAR data for
360 the topography, with the resolution of 5mx5m.

361 In order to analyze the computational load of the numerical engine as well
362 as the gain obtained using both Single-Core and Multi-Core approaches, a
363 sequential version using an *Intel Core i7 3770k@3.5 GHz* is compared against
364 an OpenMP (4 Threads) parallel version and a GPU version without taking
365 into account the improvements proposed on the paper (i.e. without using
366 Structure of Arrays, the RCM technique for cells and edges) and against
367 an optimized version running on GPU. Both GPU versions, non-optimized
368 and optimized, have been executed with a *NVIDIA Tesla c2075* GPU. All
369 the implementations have been tested in the previous four cases and the
370 computational cost time as well as the speed-ups of the parallel versions are
371 highlighted in Table 1 and in Figure 13.

372 Taking into account both factors, small time-step size and high number of
373 cells, the computational cost time for the sequential version of the numerical
374 engine for tests cases 1, 2, and 3 is three orders of magnitude larger than
375 the physical time. On the other hand, the parallel Multi-Core version, with
376 4 cores of the same CPU, can accelerate the computation between 2.25 and
377 2.65 times, depending on the test chosen. The GPU improves the simulation
378 cost between 34.88 and 49.40 times compared with the sequential version.
379 Moreover, with the optimized GPU version, the computational cost time is
380 smaller and a speed-up between 49.96 and 59.85 is obtained.

381 The main reason for differences on the speed-up of cases 1, 2 and 3 against
382 test case 4 is that the number of calculations in the latter is higher compared
383 with the number of calculations for the other cases, i.e. there are more cells

384 which have a soil volume able of being mobilized. Consequently, the more
 385 cells are involved in the calculus, the capabilities of the GPU are exploited in
 386 a more more efficient way, leading to larger speed-ups (Lacasta et al., 2014).

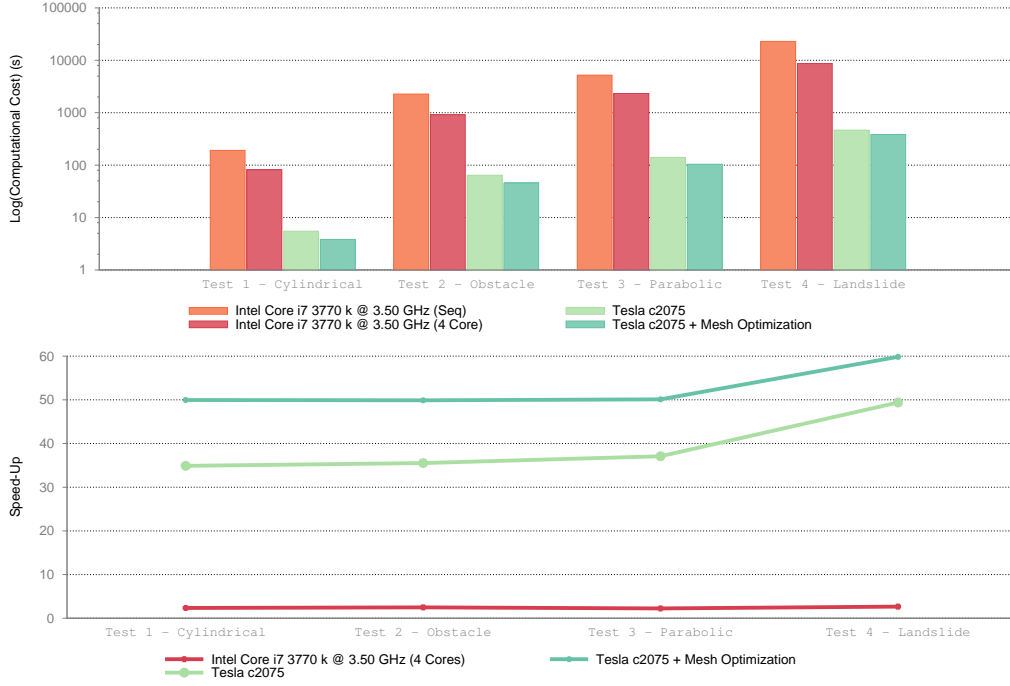


Figure 13: Computational time (upper) using logarithmic scale and speed-up (lower) for the compared implementations

Case	n_{cells}	Seq.	4 Cores		GPU Std.		GPU Opt	
		$t(s)$	$t(s)$	sup	$t(s)$	sup	$t(s)$	sup
Test 1	319354	191.15	81.63	2.34	5.48	34.88	3.83	49.96
Test 2	458684	2274.13	912.94	2.49	64.03	35.52	45.59	49.89
Test 3	670940	5217.70	2319.48	2.25	140.68	37.09	104.12	50.11
Test 4	869149	22929.15	8657.59	2.65	464.16	49.40	383.13	59.85

Table 1: Detail of computational cost and speed-up for the compared implementations using the four cases.

387 6. Conclusions

388 A finite volume scheme for modeling landslides has been implemented on
389 GPU. Unstructured meshes have been considered, since this grid topology
390 is the only one which avoid misleading preferential flow directions. The
391 use of unstructured meshes generates difficulties when moving to the GPU
392 architecture since an optimization process is required. In this work, the RCM
393 technique is described for overcoming this difficulty.

394 The computational cost times have been compared with those obtained
395 when considering Single-Core and Multi-Core processors. The GPU imple-
396 mentation and specially, the application of cell and wall ordering algorithms,
397 have driven to obtain noticeable improvements in the speed-up of the test
398 cases. The GPU use allows to obtain speed-ups ranging from 49-59, depend-
399 ing on the test case.

400 Future work may include the research of the implementation of these
401 methods using distributed computing. These ideas can also be extended to
402 develop a multi-GPU solution where the limitation on the problem size can
403 be relaxed.

404 References

- 405 Borkar, S., 2007. Thousand core chips: A technology perspective, in: Pro-
406 ceedings of the 44th Annual Design Automation Conference, ACM, New
407 York, NY, USA. pp. 746–749.
- 408 Bouchut, F., Fernández-Nieto, E.D., Mangeney, A., Lagrée, P.Y., 2008. On
409 new erosion models of Savage-Hutter type for avalanches. *Acta Mechanica*
410 199, 181–208.
- 411 Chapman, B., Jost, G., Pas, R.v.d., 2007. *Using OpenMP: Portable Shared*
412 *Memory Parallel Programming (Scientific and Engineering Computation)*.
413 The MIT Press.
- 414 Dreslinski, R., Wieckowski, M., Blaauw, D., Sylvester, D., Mudge, T., 2010.
415 Near-threshold computing: Reclaiming moore’s law through energy effi-
416 cient integrated circuits. *Proceedings of the IEEE* 98, 253–266.
- 417 García-Ruiz, J., Arnaez, J., Beguería, S., Seeger, M., Marti-Bono, C., Regues,
418 D., Lana-Renault, N., White, S., 2005. Runoff generation in an inten-

- 419 sively disturbed, abandoned farmland catchment, Central Spanish Pyrenees. *Catena* 59, 79–92.
- 420
- 421 Glaskowsky, P.N., 2009. Nvidia’s fermi: The first complete GPU computing architecture. A white paper prepared under contract with NVIDIA Corporation , 126.
- 422
- 423
- 424 Gray, J., Wieland, K., Hutter, K., 1999. Gravity-driven free surface flow of granular avalanches over complex basal topography. *Proc. Royal Soc. London Ser. A*, 455, 1841.
- 425
- 426
- 427 Iverson, R., Denlinger, R., 2001. Flow of variably fluidized granular masses across three-dimensional terrain. A Coulomb mixture theory. *Journal of Geophysical Research* 106, 537–552.
- 428
- 429
- 430 Juez, C., Caviedes, D., Murillo, J., García-Navarro, P., 2014. Experimental study of 2D granular flow over a complex topography with obstacles using a consumer-grade RGB-D sensor. Under review in *Journal of Fluid Mechanics* 73, 177–197.
- 431
- 432
- 433
- 434 Juez, C., Murillo, J., García-Navarro, P., 2013. 2D simulation of granular flow over irregular steep slopes using global and local coordinates. *Journal of Computational Physics* 255, 166–204.
- 435
- 436
- 437 Lacasta, A., Morales-Hernández, M., Murillo, J., García-Navarro, P., 2014. An optimized GPU implementation of a 2d free surface simulation model on unstructured meshes. *Advances in Engineering Software* 78, 1 – 15.
- 438
- 439
- 440 Lajeunesse, E., Mangeney-Castelnau, A., Villote, J.P., 2004. Spreading of a granular mass on a horizontal plane. *Physics of Fluids* 16, 2371–2381.
- 441
- 442 Mangeney, A., Roche, O., Hungr, O., Mangold, N., Faccanoni, G., Lucas, A., 2010. Erosion and mobility in granular collapse over sloping beds. *Journal of Geophysical Research* 115, F03040.
- 443
- 444
- 445 Moore, G., 2003. No exponential is forever: but ”forever” can be delayed! [semiconductor industry], in: *Solid-State Circuits Conference, 2003. Digest of Technical Papers. ISSCC. 2003 IEEE International*, pp. 20–23 vol.1.
- 446
- 447
- 448 Moretti, L., Mangeney, A., Capdeville, Y., Stutzmann, E., Huggel, C., Schneider, D., Bouchut, F., 2012. Numerical modeling of the Mount Steller
- 449

- 450 landslide flow history and of the generated long period seismic waves. Geo-
451 physical Research Letters 39, L16402.
- 452 Murillo, J., García-Navarro, P., 2010. Weak solutions for partial differential
453 equations with source terms: Application to the shallow water equations.
454 Journal of Computational Physics 229, 4327–4368.
- 455 Murillo, J., García-Navarro, P., 2012. Wave Riemann description of friction
456 terms in unsteady shallow flows: Application to water and mud/debris
457 floods. Journal of Computational Physics 231, 1963–2001.
- 458 NVIDIA, 2011. NVIDIA CUDA C Programming Guide.
- 459 Owens, J.D., Luebke, D., Govindaraju, N., Harris, M., Krger, J., Lefohn, A.,
460 Purcell, T.J., 2007. A survey of general-purpose computation on graphics
461 hardware. Computer Graphics Forum 26, 80–113.
- 462 Pirulli, M., Bristeau, M., Mangeney-Castelnau, A., Scavia, C., 2007. The
463 effect of the earth pressure coefficients on the runout of granular material.
464 Environmental Modelling and Software 22, 1437–1454.
- 465 Pirulli, M., Mangeney, A., 2008. Results of Back-Analysis of the Propaga-
466 tion of Rock Avalanches as a Function of the Assumed Rheology. Rock
467 Mechanics and Rock Engineering 41, 59–84.
- 468 Pouliquen, O., Forterre, Y., 2002. Friction law for dense granular flows:
469 application to the motion of a mass down a rough inclined plane. Journal
470 of Fluid Mechanics 453, 133–151.
- 471 Savage, S., Hutter, K., 1989. The motion of a finite mass of granular material
472 down a rough incline. Journal of Fluid Mechanics 199, 177–215.
- 473 Sharma, S., Gupta, K., 2013. Parallel performance of numerical algorithms
474 on multi-core system using openmp, in: Meghanathan, N., Nagamalai, D.,
475 Chaki, N. (Eds.), Advances in Computing and Information Technology.
476 Springer Berlin Heidelberg. volume 177 of *Advances in Intelligent Systems
477 and Computing*, pp. 279–288.