



Article

Developing Efficient Discrete Simulations on Multi-Core and GPU Architectures

Daniel Cagigas-Muñiz¹, Fernando Diaz-del-Rio¹, M.R. López-Torres¹, F. Jiménez-Morales²
and J.L. Guisado^{1,*}

¹ Department of Computer Architecture and Technology, University of Seville, Spain; e-mail@e-mail.com

² Department of Condensed Matter Physics, University of Seville, Spain; e-mail@e-mail.com

* Correspondence: jlguisado@us.es

Abstract: In this paper we show how to efficiently implement parallel discrete simulations on Multi-Core and GPU architectures through a real example of application: a cellular automata model of laser dynamics. We describe the techniques employed to build and optimize the implementations using OpenMP and CUDA frameworks. We have evaluated the performance on two different hardware platforms that represent different target market segments: high-end platforms for scientific computing, using an Intel Xeon Platinum 8259CL server with 48 cores and also an NVIDIA Tesla V100 GPU, both running on Amazon Web Server (AWS) Cloud, and on a consumer-oriented platform, using an Intel Core i9 9900k CPU and an NVIDIA GeForce GTX 1050 TI GPU. Performance results are compared and analysed in detail. We show that excellent performance and scalability can be obtained in both platforms, and we extract some important issues that imply a performance degradation for them. We also found that current Multi-Core CPUs with large core numbers can bring a performance very near to that of GPUs, even similar in some cases.

Keywords: laser dynamics; parallel computing; cellular automatas; GPUs and Multi-Core processors performance

1. Introduction

Discrete simulation methods encompass a family of modeling techniques which employ entities that inhabit discrete states and evolve in discrete time steps. Examples include models with an intrinsic discrete nature, such as cellular automata (CA) and related lattice automata, like lattice gas automata (LGA) or lattice Boltzmann method (LBM), and also discretization of continuous models like many stencil-based partial differential equation (PDE) solvers and particle methods based on fixed neighbor lists. They are a powerful tool that has been widely used to simulate complex systems of very different kinds (in which a global behaviour results from the collective action of many simple components that interact locally) and to solve systems of differential equations.

To accurately simulate real systems, the quality of the computed results very often depends on the number of data points used for the computations and the complexity of the model. As a result, realistic simulations often involve too large runtime and memory requirements for a sequential computer. Therefore, efficient parallel implementation of this kind of discrete simulations is extremely important. But this type of discrete algorithms has a strong parallel nature, because they are composed of many individual components or cells that are simultaneously updated. They also have a local nature, since the evolution of cells is determined by strictly local rules, i.e. each cell only interacts with a number of neighboring cells. Thanks to this, they are very suitable to be implemented efficiently on parallel computers [1,2].

34 In this paper, we study the efficient parallel implementation of a real application of this type, a CA
35 model of laser dynamics, on Multi-Core and GPU architectures, employing the most commonly used
36 software frameworks for these platforms today: OpenMP and CUDA, respectively. In both cases, we
37 describe code optimizations that can speed-up the computation and reduce memory usage. In either
38 case we have evaluated the performance on two different hardware platforms that represent different
39 target market segments: on a high-end chip intended for scientific computing or for servers and on a
40 consumer-oriented one. In the case of the Multi-Core architecture, the performance has been evaluated
41 on a dual socket server with 2 high-end Intel Xeon Platinum 8259CL processors (completing 48 cores
42 between them) running on Amazon Web Server (AWS) Cloud, and also on a PC market Intel Core i9
43 9900k processor. For the GPU architecture, we present performance evaluation results on a high-end
44 GPGPU NVIDIA Tesla V100 GPU running on AWS Cloud and on a consumer-oriented NVIDIA
45 GeForce GTX 1050 TI. In all cases, we reported speedups compared to a sequential implementation.
46 The aim of this work is to extract lessons that may be helpful for practitioners trying to implement
47 discrete simulations of real systems in parallel.

48 The considered application uses cellular automata, a class of discrete, spatially-distributed
49 dynamical systems with the following characteristics: spatial and temporal discrete character, local
50 interaction and synchronous parallel dynamical evolution [3,4]. They can be described as a set
51 of identical finite state machines (cells) arranged along a regular spatial grid, whose states are
52 simultaneously updated by a uniformly applied state-transition function that refers to the states
53 of their neighbors [5]. In the last decades, CA have been successfully applied to build simulations
54 of complex systems in a wide range of fields, including physics (fluid dynamics, magnetization in
55 solids, reaction-diffusion processes), bio-medicine (viral infections, epidemic spreading), engineering
56 (communication networks, cryptography), environmental science (forest fires, population dynamics),
57 economy (stock exchange markets), theoretical computer science, etc [6–8]. They are currently being
58 very used, in particular, for simulations in geography (specially in urban development planning
59 [9], future development of cities [10], and land use [11]) pedestrian or vehicular traffic [12,13], and
60 bio-medicine (applied to physiological modeling, for example for cancer [14], or epidemic modeling
61 [15]).

62 The application considered in this study is a cellular automata model of laser dynamics introduced
63 by Guisado et. al., capable of reproducing much of the phenomenology of laser systems [16–19].
64 It captures the essence of laser as a complex system in which its macroscopic properties emerge
65 spontaneously due to the self-organization of its basic components. This model is a useful alternative
66 to the standard modeling approach of laser dynamics, based on differential equations, in situations
67 for which the approximations considered for them are not valid, for instance for lasers ruled by stiff
68 differential equations, lasers with difficult boundary conditions, or very small devices. The mesoscopic
69 character of the model also allows to have results impossible to be obtained by the differential equations,
70 such as studying the evolution of its spatio-temporal patterns.

71 In order to reduce the runtime of laser simulations with this model by taking advantage of its
72 parallel nature, a parallel implementation of it for computer clusters (distributed-memory parallel
73 computers), using the message-passing programming paradigm, was introduced in [20,21]. It showed
74 a good performance on dedicated computer clusters [22] and also on heterogeneous non-dedicated
75 clusters with a dynamic load balancing mechanism [23].

76 Due to the excellent ratio performance/price and performance/power of the Graphics Processing
77 Units (GPUs), it is very interesting to implement the model on them. GPUs are massively parallel
78 graphics processors originally designed for running interactive graphics applications, but that can also
79 be used to accelerate arbitrary applications, what is known as GPGPU (General Purpose computation
80 on GPU) [24]. They can run thousands of programming threads in parallel, providing speedups
81 mainly from 10x to 200x compared to CPU (depending on the application and on the optimizations of
82 its implementation), at very affordable prices. Therefore, GPUs have widespread use today in high
83 performance scientific computing. Their architecture is formed by a number of multiprocessors, each

84 of them with a number of cores. All cores of a multiprocessor share one memory unit called shared
85 memory and all multiprocessors share a memory unit called global memory.

86 A first version of a parallel implementation of the model for GPUs was presented in [25]. Even
87 when that first implementation did not explore all the possible optimizations to boost the performance
88 on that platform, it showed that the model could be successfully implemented on GPU. A speedup of
89 14.5 on a NVIDIA GeForce GTX 285 (a consumer-oriented GPU intended for low-end users and gamers)
90 compared to an Intel Core i5 750 with 4 cores at 2.67 GHz was obtained. The GPU implementation
91 described in the present paper differs from that previous one in that this new version has been carefully
92 optimized to extract all possible performance from the GPU and his performance has been evaluated
93 not only on a consumer-oriented GPU, but also on a Tesla scientific high-end GPU.

94 Another interesting parallel platform to implement discrete simulations today are Multi-Core
95 processors. Since around 2005 all general-purpose CPUs implement more than one CPU (or "core") on
96 the processor chip. For a decade, the number of cores in standard Intel x86 processors was modest
97 (mainly from 2 to 8). But in the last years, there are high-end CPUs in the market including up to
98 several dozen cores (now up to 18 cores for Intel Core i9 and up to 56 cores for Intel Xeon Platinum).
99 Therefore, Multi-Core CPUs can start to be competitive with GPUs to implement parallel discrete
100 simulations, specially taking into account that its parallelization with OpenMP is much easier than for
101 GPUs. Therefore, we also present the first parallel implementation of the CA laser dynamics model for
102 Multi-Core architectures and compare its performance on current high-end Multi-Core CPUs to the
103 performance obtained on GPUs.

104 The remainder of the paper is organized as follows: Section 2 reviews the related work in the field
105 of discrete simulations via cellular automata and their parallel implementation on Multi-Core and
106 GPU Architectures. Section 3 describes the methodology employed in this work, trying to give useful
107 indications to researchers interested in parallelizing efficiently their own codes. Section 4 presents the
108 results and discusses their interpretation and significance. Finally, Section 5 summarizes the contents
109 of this paper, the conclusions and indicates interesting future work.

110 2. Related work

111 Most parallel implementations of CA models on Multi-Core processors or GPUs were presented
112 after 2007. In the case of Multi-Core processors, they became generalised only from 2005 onwards, and
113 started to be used for parallel simulations in the following years. As regards GPUs, before 2007 there
114 were few works devoted to the parallel implementation of cellular automata models on GPUs, because
115 they had to adapt somehow their application to a shading language (a special purpose programming
116 language for graphics applications), such as OpenGL. An example is the paper from Gobron et. al.
117 [26], that studies a CA model for a biological retina obtaining a 20x speedup as compared to the CPU
118 implementation. After the introduction in 2007 of CUDA (Compute Unified Device Architecture), a
119 general purpose programming language for GPUs of the NVIDIA manufacturer, followed soon by a
120 multi-platform one called OpenCL, the usage of GPUs in scientific computing exploded.

121 Let us review some relevant parallel implementations of CA models on Multi-Core CPUs and
122 GPUs introduced from 2007 on.

123 Rybacki et. al. [27] presented a study of the performance of seven different very simple cellular
124 automata standard models running on a single core processor, a multi core processor and a GPU. They
125 found that the performance results were strongly dependent on the model to be simulated.

126 Bajzát et. al. [28] obtained an order of magnitude increase in the performance of the GPU
127 implementation of a CA model for an ecological system, compared to a serial execution.

128 Balasalle et. al. [29] studied how to improve the performance of the GPU implementation of one
129 of the simplest two-dimensional CAs—the game of life—by optimizing the memory access patterns.
130 They found that carefull optimizations of the implementation can produce a 65% improvement in
131 runtime from a baseline implementation. However, they did not study other more realistic CA models.

132 Special interest has been devoted to GPU implementations of Lattice Boltzmann methods, a
133 particular class of CA. Some works have been able to obtain spectacular speedups for them. For
134 instance, [30] reported speedups of up to 234x respect to single-core CPU execution without using SSE
135 instructions or multithreading.

136 Gibson et. al. [31] presents the first thorough study of the performance of cellular automata
137 implementations on GPUs and multi-core CPUs with respect to different standard CA parameters such
138 as lattice and neighbourhood sizes, number of states, complexity of the state transition rules, number
139 of generations, etc. They have studied a "toy application", the "game of life" cellular automaton in two
140 dimensions and two multi-state generalizations of it. They employed the OpenCL framework for the
141 parallel implementation on GPUs and OpenMP for multi-core CPUs. That study is very useful for
142 researchers to help them choose the right CA parameters, when that is possible, taking into account
143 their impact in performance. Also to help to explain much of the variation found in reported speedup
144 factors from literature. Our present work is different and complementary to that study in the sense
145 that the game of life is a toy model very useful to study the dependence of performance on general
146 CA parameters, but it is also very interesting to study the parallelization and performance of a real
147 application instead of a toy model such as the game of life, as we do in this work.

148 3. Materials and Methods

149 3.1. Cellular automaton model for laser dynamics simulation

150 We present parallel implementations for Multi-Core CPUs and for GPUs of the cellular automaton
151 model of laser dynamics introduced by Guisado et. al. [16–18].

152 A laser system is represented in this model by a two-dimensional CA which corresponds to a
153 transverse section of the active medium in the laser cavity.

154 Formally the CA model is made of:

155 a) A regular lattice in a two-dimensional space of $L \times L$ cells. Each lattice position is labelled by
156 the indices (i, j) . Also to avoid boundary problems and to best simulate the properties of a macroscopic
157 system we use periodic boundary conditions.

158 b) The state variables associated with each node (i, j) . In the case of a laser system we need two
159 variables: one for the lasing medium $a_{ij}(t)$ and the other for the number of laser photons $c_{ij}(t)$. $a_{ij}(t)$
160 is a boolean variable: 1 represents the excited state of the electron in the lasing medium in cell (i, j)
161 and 0 is the ground state. For the photons $c_{ij}(t)$ is an integer variable in the range $[0, M]$ where M is
162 an upper limit, that represent the number of laser photons in cell (i, j) . The state variables $a_{ij}(t)$ and
163 $c_{ij}(t)$ represent "bunches" of real photons and electron, the values of which are obviously smaller than
164 the real number of photons and electrons in the system and are connected to them by a normalization
165 constant.

166 c) The neighborhood of a cell. In a cellular automata the state variables can change depending on
167 the neighboring cells. In our model the *Moore neighborhood* is employed: the neighborhood of a cell
168 consists of the cell itself and the eight cells around it at positions north, south, east, west, northeast,
169 southeast, northwest and southwest.

170 d) The evolution rules that specify the state variables at time $t + 1$ in function of their state at time
171 t . From a microscopic point of view the physics of a laser can be described by five processes:

172 i) The pumping of the ground state of the laser medium to the excited state. In this way energy
173 is supply to the lasing medium. This process is considered to be probabilistic: If $a_{ij}(t) = 0$ then
174 $a_{ij}(t + 1) = 1$ with a probability λ .

175 ii) The stimulated emission by which a new photon is created when an excited laser medium cell
176 surrounded by one or more photons decays to the ground state: If $a_{ij}(t) = 1$ and the sum of the values
177 of the laser photons states in its neighboring cells is greater than 1, then $c_{ij}(t + 1) = c_{ij}(t) + 1$ and
178 $a_{ij}(t + 1) = 0$.

- 179 iii) The non-radiative decaying of the excited state. After a finite time τ_a a excited laser medium
 180 cell will go to the ground state $a_{ij}(t + 1) = 0$ without the generation of any photon.
- 181 iv) The photon decay. After a given time τ_c , photons will escape and its number will decrease by
 182 one unit $c_{ij}(t + 1) = c_{ij}(t) - 1$.
- 183 v) Thermal noise. In a real laser system there is a thermal noise of photons produced by
 184 spontaneous emissions and they cause the initial start-up of the laser action. Therefore in our CA
 185 model a small number of photons less than 0.01% are added at random positions at each time step.

```

1: Initialize system
2: Input data
3: for time step = 1 to maximum time step do
4:   for each cell in the array do
5:     Apply noise photons creation rule (Fig. 2)
6:     Apply photon and electron decay and evolution of temporal variables
       (Fig. 3 )
7:     Apply pumping and stimulated emission rules (Fig. 4)
8:   end for
9:   Refresh value of  $c$  matrix with contents of  $c'$  matrix
10:  Calculate populations after this time step
11:  Optional additional calculations on intermediate results
12: end for
13: Final calculations
14: Output results
  
```

Figure 1. Pseudo code description of the main program for the CA laser model.

186 3.2. Sequential implementation of the model

187 The algorithmic description of the model using pseudo code is shown in Figs. 1 to 4. The main
 188 program is described in Fig. 1. The structure of the algorithm is based on a time loop, inside of which
 189 there is a data loop to sweep all the CA cells. At each time step, first the state of all the cells of the
 190 lattice is updated by applying the transition rules, and then the total populations of laser photons and
 191 electrons in the upper state are calculated by summing up the values of the state variables a_{ij} and c_{ij} for
 192 all the lattice cells. Because we are emulating a time evolution, the order of the transition rules for each
 193 time step can be switched. Of course, different orders get to slightly different particle quantities, but
 194 on the whole, CA evolution is similar. Fig. 2 defines the implementation of the noise photons creation
 195 rule. The photon and electron decay rules and the evolution of temporal variables are described in Fig.
 196 3. Finally, Fig. 4 describes the implementation of the pumping and stimulated emission rules.

197 In order to simulate a parallel evolution of all the CA cells, we use two copies of the c_{ij} matrix,
 198 called c and c' . In each time step, the new states of c_{ij} are written in c' and the updated values of this
 199 matrix are only copied to c after finishing with all the CA cells. In the algorithmic description of the
 200 implementation of the model we have used two temporal variables, \tilde{a}_{ij} and \tilde{c}_{ij}^k as time counters, where
 201 k distinguishes between the different photons that can occupy the same cell. When a photon is created,
 202 $\tilde{c}_{ij}^k = \tau_c$. After that, 1 is subtracted to \tilde{c}_{ij}^k for every time step and the photon will be destroyed when
 203 $\tilde{c}_{ij}^k = 0$. When an electron is initially excited, $\tilde{a}_{ij} = \tau_a$. After that, 1 is subtracted to \tilde{a}_{ij} for every time
 204 step and the electron will decay to the ground level again when $\tilde{a}_{ij} = 0$.

```

1: {Introduce  $n_n$  number of photons in random positions}
2: for  $n = 0$  to  $n_n - 1$  do
3:   {Generate two random integers in  $(0, size - 1)$  interval}
4:    $i \leftarrow random\_number(0, L_x - 1)$ 
5:    $j \leftarrow random\_number(0, L_y - 1)$ 
6:   {Look for first value of  $k$  for which  $\tilde{c}_{ij}^k = 0$ }
7:   while  $\tilde{c}_{ij}^k \neq 0$  and  $k \leq M$  do
8:      $k \leftarrow k + 1$ 
9:   end while
10:  if  $k \leq M$  then
11:    {Create new photon}
12:     $c'_{ij} \leftarrow c'_{ij} + 1$ 
13:     $\tilde{c}_{ij}^k \leftarrow \tau_c$ 
14:  end if
15: end for

```

Figure 2. Pseudo code diagram for the implementation of the noise photons rule.

205 3.3. Parallel Frameworks for Efficient CA Laser Dynamics Simulation

206 Algorithms described in previous Section 3.2 arise from a direct conversion of the systems of
 207 differential equations that represents the CA laser model. The efficient execution of these algorithms in
 208 parallel platforms to generate fast simulations of a bunch of different input parameters requires many
 209 specific considerations for each hardware platforms. To begin with, modern out-of-order execution
 210 superscalar processors achieves an almost optimal time execution of operations when operands reside
 211 in CPU registers. That is, they reach the so-called *data – flow – limit* of the algorithm, being the most
 212 patent bottlenecks that of the real dependences among operations and the difficult branch predictions.
 213 In fact, taking some simple assumptions around these bottlenecks, some authors have proposed simple
 214 processor performance models that predicts computing times with enough accuracy [32]. Above this,
 215 when many operations cannot be executed over CPU registers, memory model is the other crucial
 216 factor.

217 In relation with our CA model, a simple inspection of the code and of the data evolution brings to
 218 light that memory usage is massive and that an elevated branch misprediction ratio is expected. First
 219 assertion is obvious: matrices that contains c_{ij} , \tilde{c}_{ij}^k , a_{ij} , \tilde{a}_{ij} supposes many megabytes for those lattice
 220 widths that emulates a correct behavior of the laser dynamics. Second assertion comes mainly from
 221 two code features: the use of random values in many decisions representing the particle evolution,
 222 and the chaotic values that particle states take along the life of the simulation. Whereas this paper
 223 concentrates in a laser dynamics model, it is obvious that these features may be present in many CA
 224 simulations, mostly when cooperative phenomena are expected. What is more relevant, the existence
 225 of many branches (some of them in the form of nested conditional structures) in the “hot spot” zones
 226 implies that GPU implementations would suffer from an important deceleration. This is due to the
 227 inherent so-called thread divergence [33] that GPU compilers introduce in these cases, which is one
 228 the main reasons why the performance on these platforms diminishes.

229 Taking into account previous considerations, an accurate timing characterization of main
 230 sequential algorithm pieces was done. This analysis concludes that:

231 - More than 80% of the mean execution time is spent in stimulated emission and pumping rules.
 232 What is more, their execution times have a considerable variance: minimum times are around five
 233 times lower than maximum times. This asserts the effect of random values and the chaotic evolution
 234 of different cell particles.

```

1: for  $j = 0$  to  $L_y - 1$  do
2:   for  $i = 0$  to  $L_x - 1$  do {CA lattice loop}
3:     if  $c_{ij} > 0$  then {Apply photon decay rule}
4:       for  $k = 1$  to  $M$  do
5:         {Subtract 1 to every photon's lifetime}
6:         if  $\tilde{c}_{ij}^k > 0$  then
7:            $\tilde{c}_{ij}^k \leftarrow \tilde{c}_{ij}^k - 1$ 
8:           if  $\tilde{c}_{ij}^k = 0$  then {One photon decays}
9:              $c_{ij} \leftarrow c_{ij} - 1$ 
10:             $c'_{ij} = c_{ij}$ 
11:          end if
12:        end if
13:      end for
14:    end if
15:    if  $a_{ij} = 1$  then {Apply electron decay rule}
16:      {Subtract 1 to time of life of every excited electron}
17:       $\tilde{a}_{ij} \leftarrow \tilde{a}_{ij} - 1$ 
18:      if  $\tilde{a}_{ij} = 0$  then
19:        {One electron decays}
20:         $a_{ij} \leftarrow 0$ 
21:      end if
22:    end if
23:  end for
24: end for

```

Figure 3. Pseudo code diagram for the implementation of the photon and electron decay and evolution of temporal variables rules.

- 235 - Random number computation supposes around the 70% of the pumping rule time.
- 236 - The rest of time resides mainly in photon and electron decay. The oscillatory behavior of particle
- 237 evolution during the stationary phase implies also a considerable variance in these times. This is even
- 238 more exaggerated during transient evolution.
- 239 - Noise photon rule timing is negligible (in fact, its number of iterations is very much inferior than
- 240 the rest of rules).
- 241 Previous facts make necessary the introduction of at least the next changes in both OpenMP and
- 242 GPU code implementations (see <https://github.com/dcagigas/Laser-Cellular-Automata>):
- 243 - Of course, avoiding non re-entrant functions like simple random generators. Even more,
- 244 although generating a seed for each thread should be enough to make random generation independent
- 245 among threads, the deep inner real data dependences that random functions contain lasts in the mean
- 246 longer than the rest of an iteration step. It is preferable an implementation similar to that of the
- 247 *cuRAND* library, that is, a seed for each cell ij , which preserves a good random distribution while
- 248 accelerates each step around a 40%.
- 249 - As only one electron per cell is allowed, suppressing the a_{ij} matrix. Thus, it is considered that if
- 250 \tilde{a}_{ij} is zero the electron is not excited; and it is excited elsewhere.
- 251 - Eliminating the refresh of c matrix (which supposes copying long matrices) with values of \tilde{c} (line
- 252 9 of Fig. 1), by using pointers to these two matrices and swapping these pointers at the end of each
- 253 iteration step.

```

1: for  $j = 0$  to  $L_y - 1$  do
2:   for  $i = 0$  to  $L_x - 1$  do {CA lattice loop}
3:     if  $a_{ij} = 0$  then {Apply pumping rule}
4:       {Generate random number in  $(0, 1)$  interval}
5:        $\xi \leftarrow \text{random\_number}(0, 1)$ 
6:       if  $\xi < \lambda$  then  $\{\lambda: \text{pumping probability}\}$ 
7:         {One electron is pumped}
8:          $a_{ij} \leftarrow 1$ 
9:          $\tilde{a}_{ij} \leftarrow \tau_a$ 
10:        end if
11:      else  $\{(a_{ij} = 1) \rightarrow \text{Apply stimulated emission rule}\}$ 
12:        if  $\text{neighbours}(i, j) > \delta$  then
13:          {Look for first value of  $k$  for which  $\tilde{c}_{ij}^k = 0$ }
14:           $k \leftarrow 1$ 
15:          while  $\tilde{c}_{ij}^k \neq 0$  and  $k \leq M$  do
16:             $k \leftarrow k + 1$ 
17:          end while
18:          if  $k \leq M$  then
19:             $a_{ij} \leftarrow 0$ 
20:             $\tilde{a}_{ij} \leftarrow 0$ 
21:             $c'_{ij} \leftarrow c'_{ij} + 1$ 
22:             $\tilde{c}_{ij}^k \leftarrow \tau_c$ 
23:          end if
24:        end if
25:      end if
26:    end for
27:  end for

```

Figure 4. Pseudo code diagram for the implementation of the pumping and stimulated emission rules.

254 The source code of the different implementations and results achieved are available in
255 <https://github.com/dcagigas/Laser-Cellular-Automata>. The source code is under GPL v3 license.
256 Researchers can download and modify the code freely to run their own particular laser dynamic
257 simulations.

258 3.3.1. OpenMP Framework

259 Previous improvements are quite easy to detect and to implement. However, there are further
260 enhancements that speedup even more a CA simulation when running an OpenMP implementation
261 over multicore platforms. As a result, apart from the OPENMP_NOT_OPTIMIZED version, an
262 optimized one (called simply OPENMP) can be downloaded from the previous *github* page. For the
263 sake of clarity, these further enhancements are grouped and listed in the following points. Moreover,
264 they have been marked in the *github* source code with the symbol @.

265 - After a deeper examination of laser dynamics evolution, it was detected that very few cells
266 contain more than one photon during the stationary evolution. Thus, the habitual matrix arrangement
267 of variable \tilde{c}_{ij}^k , that is, storing consecutively the M values for each cell ij is switched by the following
268 one. All the cells ij are stored consecutively for each of the M possible photons. In terms of the C++
269 code, this three-dimensional matrix is represented by: $\tilde{c}[M][L_x][L_y]$ (see *lifet_f* matrix in the code).

270 The new arrangement implies that all elements of $\tilde{c}[0][][[]]$ are continuously used and then cached, but
271 the rest of elements $\tilde{c}[1 : M - 1][L_x][L_y]$ are scarcely used, so they do not consume precious cache lines.
272 On the contrary, if the habitual matrix arrangement had been used, it would have wasted many cache
273 bytes (almost only 1 of each M elements would have been really utilized during the stationary period).
274 The rest of code pieces where this matrix is manipulated are not decelerated by the new arrangement;
275 e.g. very few cells generate a shifting from $\tilde{c}[k][L_x][L_y]$ to $\tilde{c}[k - 1][L_x][L_y]$, ($k > 0$), when a photon
276 decays.

277 - While previous improvement avoids cache line wasting, memory consumption is another
278 fundamental issue. The analysis of real values of the big code matrices leads to the conclusion
279 that maximum values are small for most physical variables. Thus, instead of 32 bits per element
280 (*unsignedint* variables in C++), real sizes in the optimized version have been reduced to unsigned
281 short int and unsigned char whenever possible. More exactly, this supposed reducing memory size
282 from: $32 \times L_x \times L_y \times (M + 3)$ to $16 \times L_x \times L_y \times (M + 1) + 8 \times 2$ (see *e, f1, f2, lifet_f* matrices in the
283 optimized code).

284 - In order to promote loop vectorization, some conditional branches have been transformed
285 into simple functions. For example, those conditional sentences that increment a counter q when a
286 certain condition p is true have been written like $q += p$. This eases the task of the compiler when
287 introducing SIMD instructions and predicative-like code and prevents many BTB (Branch Target
288 Buffer) misprediction penalties because these conditions are difficult to predict (due to the random
289 nature of particle state evolution).

290 - Loop splitting is another classical technique that reduces execution time when memory
291 consumption is huge, and the loop manages several disjointed data. This occurs in the case of
292 photon and electron decay rules, which have been separated into two different loops in the optimized
293 version. This way, caches are not struggled with several matrixes thus preventing many conflict misses
294 on them.

295 To sum up, previous optimizations achieve around a 2x speedup (see section 4) with respect to the
296 basic one. It is worth to remark that both OpenMP versions give exactly the same particle evolution
297 results.

298 Despite that random number computing has been considerable accelerated by using a seed for
299 each cell (i, j) , it continues to be the most time-consuming piece. A final improvement draws to an
300 approx. additional 3 times speedup of the OpenMP simulation time: instead of computing random
301 numbers during the simulation, generating a list of them previously and using this same list for all the
302 desired simulations (e.g. if different parameters want to be tested like pumping probability, maximum
303 electron and photon lifetimes, etc.).

304 Using a random list as an input for the model eases the checking of results for different platforms,
305 because the output of the simulation must be exactly the same. More precisely, it is needed that a
306 random number is stored in the list for each time step and for each cell.

307 However, this list should be enormous for the pumping rule, even if only a random true/false
308 bit were stored for each time step and for each cell. For example, considering a simulation of 1000
309 steps, a lattice of 4096×4096 would occupy 2 GBytes. Because of this, this improvement has not
310 been considered in the Result section. Nonetheless, the interested reader can test this optimization
311 (note that big lattice sizes would overflow platform memory) simply by defining the constant
312 *RANDOM_VECT_ENABLE* in the *github* OpenMP codes. Defining this constant would generate the
313 random numbers in advance while suppressing its computation during the simulation time.

314 3.3.2. CUDA Framework

315 The CUDA framework has three main kernels (i.e. CUDA functions written in C style code), the
316 same as those of the OpenMP implementation. They are called for each time step sequentially. First
317 the *PhotonNoiseKernel* produces new photons randomly, then the *DecayKernel* performs the electron
318 and photon decay, and last the *PumpingKernel* does the pumping and stimulated emission. This order

319 can be altered but it must be the same as the one used in OpenMP. Otherwise, results could not be
320 exactly the same.

321 There is one last kernel needed: *do_shannon_entropy_time_step*. Most of the variables that are
322 needed to calculate the Shannon Entropy are stored in GPU global memory. Data transfers between
323 computer host memory and GPU memory must be minimized because of its big latency. Despite
324 that the calculations are not parallel, it is more convenient to perform time step Shannon Entropy
325 calculations in GPU memory. There is also a final kernel called *finish_shannon_entropy_time_step* after
326 the time/step loop. However, this last kernel has a low performance impact because it is executed
327 only once.

328 Simulation parameters are defined in a header file; for example, the SIDE constant that determines
329 the grid side of a simulation. In case of CUDA, and GPUs in general, memory size constraints
330 are particularly important when comparing with computer workstations. The GPU global memory
331 available is usually lower than that of a workstation. Thus, data structure types for electrons and
332 photons are important for large grids. Matrix data structures grow by a factor of x4 for each SIDE
333 increment, and x40 in case of the matrix that records photon energy values in each cell (*GPU_tvf*).

334 For example, with a grid SIDE of 8192 (2^{13}) and 4 GB of GPU global memory it is only possible to
335 run the simulation if cells of *GPU_tvf* matrix variable are set to *char* type in C. As mentioned before,
336 this variable is in charge of recording photon life time values in each grid cell. The *char* type is 8 bit
337 size, so only initial photon life time values between 0 and 255 are allowed. The same happens with
338 electron life time values. By default this constant value is set to *short int* (i.e. 16 bits) to allow higher
339 values.

340 The CUDA programming environment and the latest NVIDIA architectures (Pascal, Turing
341 and Volta) also have some restrictions related to integer arithmetic operations. CUDA
342 atomic arithmetic operations only allow the use of *int* data type but not the *short int* or *char*. In
343 the *PhotonNoiseKernel* it is necessary to update new photons in the matrix data structures. Those
344 updates are performed in a parallel way. To avoid race conditions, atomic arithmetic operations are
345 needed when each GPU hardware thread updates a matrix photon cell (two hardware GPU threads
346 could try to update the same cell at the same time). Therefore, it was necessary to use the *int* data type
347 instead of *short int* or *char*, thus increasing the GPU memory size needed by these data structures.

348 CUDA framework has also one extra feature that can be enabled in the source files: the electron
349 evolution output video. A .avi video file showing the electron evolution through the time steps can be
350 produced. This feature involves the transfer of a video frame for each time step from GPU memory
351 to host or computer memory. When activated and for a moderate grid side (1024 or above), the total
352 execution time could be significantly high because of the latency between GPU and host memory
353 transfers. This feature could also be adapted or modified to show photon evolution (nonetheless,
354 electron and photon behaviours are very similar).

355 4. Results

356 We present here the performance evaluation results for the two architectures: Multi-Core and
357 GPU. For each architecture we evaluated the performance on two different hardware platforms that are
358 representative of different target market segments: a high-end chip intended for scientific computing
359 or for servers and a consumer-oriented one.

360 We have tested that the simulation results of both parallel implementations reproduce the output
361 of the original sequential one. As an example we show in Fig. 5 the time evolution of the total
362 number of laser photons and the population inversion in the laser system for values of the parameter
363 corresponding to a laser spiking regime. The results are the same as found in previous publications
364 with a sequential implementation, as [21]. It is shown that when increasing the size of the CA lattice
365 the results are smoother since the model reproduces better the macroscopic behavior of the system
366 with a higher statistics.

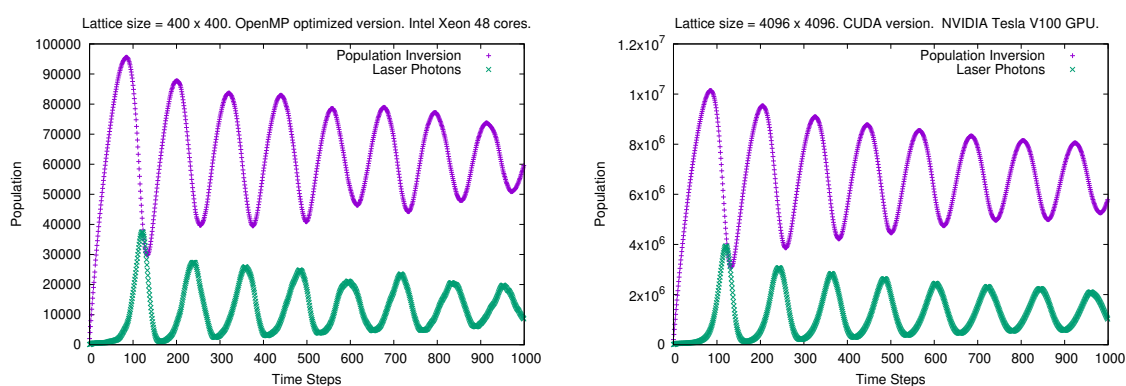


Figure 5. Output of the model for particular values of the system parameters corresponding to a laser spiking behavior. Parameter values: $\lambda = 0.0125$, $\tau_c = 10$, $\tau_a = 180$. The results are smoother for larger lattice sizes. **(Left):** Sequential implementation with a lattice size of 400×400 cells. **(Right):** Sequential implementation with a lattice size of 4096×4096 cells.

367 4.1. Multi-Core architecture

368 The Multi-Core architecture was executed and tested on the following two platforms.

369 4.1.1. High-end Multi-Core CPUs (48 cores)

370 We evaluated the performance on a high-performance server CPU running in the Cloud, using
 371 the Amazon Web Server(AWS) Infrastructure as a Service (IaaS) EC2 service. We run our performance
 372 test on a *m5.24xlarge* instance. It runs on a dual socket server with 2 Intel Xeon Platinum 8259CL
 373 processors with 24 physical cores each (completing 48 physical cores between them), running at a
 374 frequency of 2.50 GHz, with 35,75 MiB of cache memory. The total RAM memory was 373 GiB. Both
 375 processor sockets are linked by Ultra Path Interconnect (UPI), a high speed point-to-point interconnect
 376 link delivering a transfer speed of up to 10.4 GT/s.

377 4.1.2. Consumer-oriented Multi-Core CPU (8 cores)

378 The performance was evaluated on a PC with a Core i9 9900k processor and a total RAM memory
 379 of 16 GiB. The processor frequency was 3.6 GHz and the RAM memory was on a single channel
 380 running at 2400 MHz. This processor has 8 physical cores and each core has 2 hardware threads
 381 (completing a total of 16 Threads).

382 4.2. GPU architecture

383 The following two GPU chips were used to run and test the GPU architecture.

384 4.2.1. High-end GPU chip

385 We evaluated the performance on a *p3.2xlarge* instance of the Amazon Web Server(AWS) IaaS
 386 EC2 Cloud Computing service. It used a NVIDIA Tesla V100 GPU (Volta architecture) with 5,120
 387 CUDA cores and a GPU memory of 16 GiB. The server used an Intel Xeon E5-2686v4 CPU with 4
 388 physical cores running at 2.3 GHz. CPU and GPU were interconnected via PCI-Express Gen. 3, with a
 389 bandwidth of 32 GiB/s.

390 4.2.2. Consumer-oriented GPU chip

391 We used a consumer-oriented NVIDIA GeForce GTX 1050 Ti graphic card (Pascal architecture).
 392 This card has 768 CUDA cores running at 1290 MHz. The total memory is 4 GiB of GDDR5 type, with
 393 a maximum bandwidth of 1120 GiB/s.

394 4.3. Performance and scalability results

395 Fig. 6 shows the speedup of the parallel OpenMP implementation running on the Core i9 9900k
 396 PC (8 cores). In this case, speedups are fairly foreseeable. Firstly, for small lattice sizes most of matrices
 397 reside in CPU caches, thus achieving an excellent speedup versus the sequential (one thread) test until
 398 8 threads. Launching 9 threads implies that a physical CPU must manage two threads (whereas the
 399 rest of CPUs, only one), thus causing the speedup to decrease. However, this problem diminishes for
 400 more threads. Finally, Core i9 simultaneous multi-threading begins to play a role from 9 threads up,
 401 hence speedups above 8 are reached for some tests when launching many threads.

402 On the other hand, for big lattices speedups are stuck according to the maximum RAM bandwidth,
 403 as predicted by the roof-line model [34]. In these cases, RAM memory bandwidth is the bottleneck
 404 that in fact determines program runtimes.

405 In the case of the the high-end server with two 2 Intel Xeon Platinum 8259CL CPUs (see Fig. 7),
 406 our parallel implementation shows an excellent speedup for large lattice sizes. The OpenMP optimized
 407 version reaches a speedup around 30x for 48 threads. The behaviors are similar to those obtained in
 408 the consumer-oriented hardware: a peak on the speedups is reached when launching the same number
 409 of threads than physical CPUs; then accelerations begin to decrease for some more threads, and finally
 410 speedups are recovered when the number of threads doubles the number of physical cores.

411 Nevertheless, for high number of threads, or more precisely, for small number of lattice rows
 412 per thread, the computation-to-communication scale begins to deteriorate speedups. Note that small
 413 lattices (Fig. 7 for lattice width = 512) exhibit patently this problem, whereas big lattice speedups
 414 are almost not deteriorated. This is a well-known effect when scientific applications are massively
 415 distributed [35]. Because the stimulated emission rule uses Moore neighborhood, the more threads in
 416 which we divide the lattice, the more communication between threads are needed, thus degrading the
 417 parallel performance.

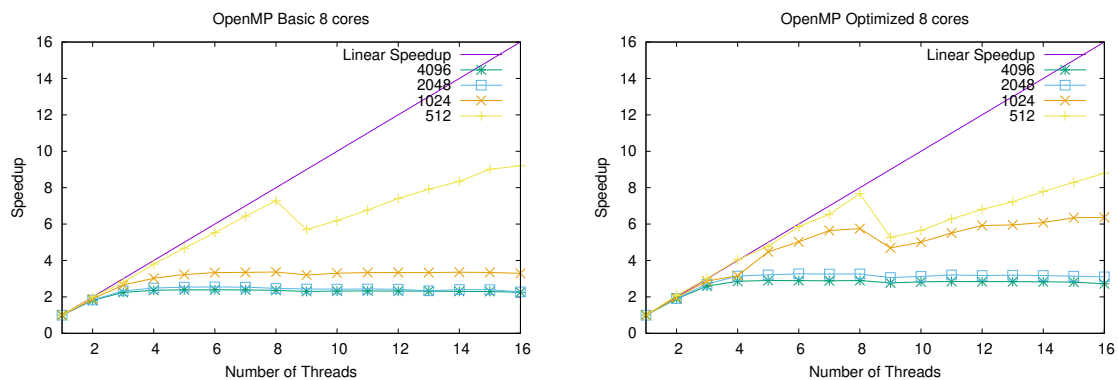


Figure 6. Speedup of the parallel OpenMP implementation running on a consumer-oriented CPU with 8 cores. **(Left:)** Basic implementation without in-depth optimization. **(Right:)** Fully optimized implementation.

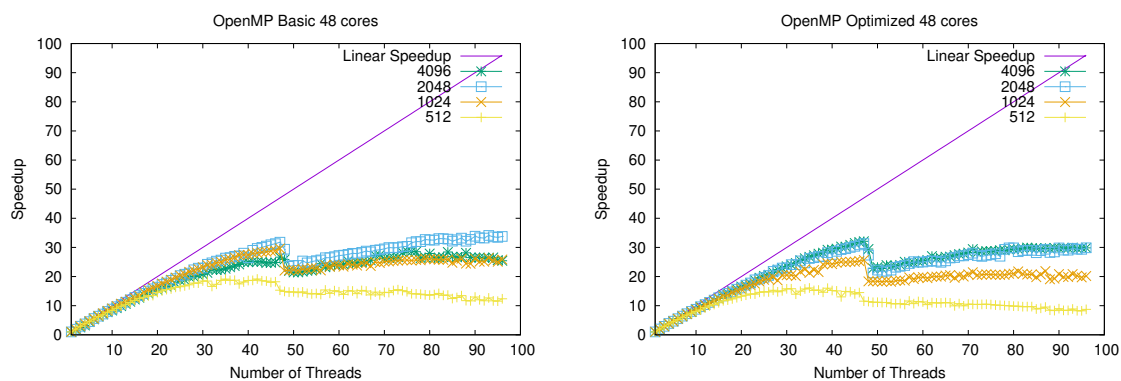


Figure 7. Speedup of the parallel OpenMP implementation running on a high-end dual-socket server with a total of 48 cores. **(Left):** Basic implementation without in-depth optimization. **(Right):** Fully optimized implementation.

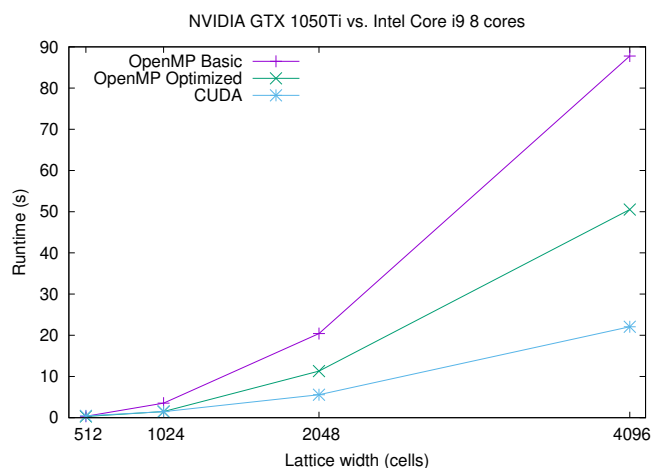


Figure 8. Runtime comparison between NVIDIA GTX 1050 Ti GPU and Intel Core i9 (8 cores) for different CA lattice sizes.

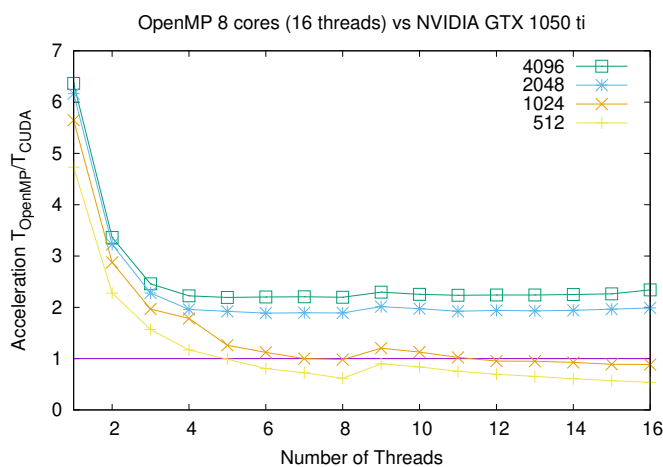


Figure 9. Comparison of OpenMP Optimized and CUDA times when using Intel Core i9 and NVIDIA GeForce GTX 1050 Ti resp.

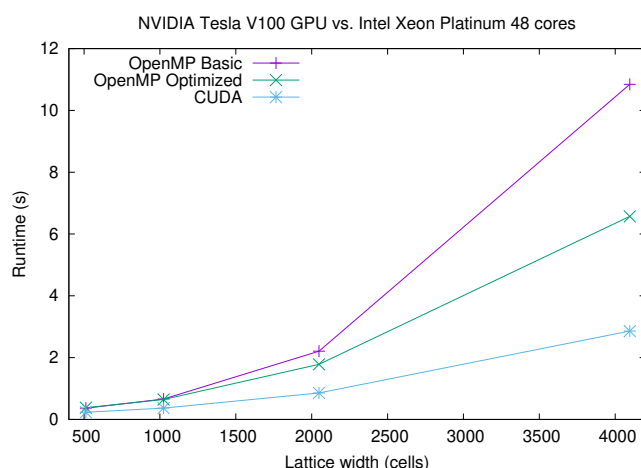


Figure 10. Runtime comparison between NVIDIA Tesla V100 GPU and Intel Xeon Platinum (dual socket with 48 cores in total) for different CA lattice sizes.

418 In Figs. 8 to 10 we show a runtime comparison between both CPU/GPU pairs, for different sizes
 419 and number of Multi-Core Threads. It is interesting to note that for CPU/GPU comparisons shown
 420 in Figs. 8, 9 and Fig. 10) CUDA implementation is, at most, 2.5x faster than the OpenMP Optimized
 421 version. In truth, only for big lattice sizes GPU platforms beat clearly CPU ones: as shown in Fig. 10,
 422 that happens only for a number of threads much smaller than the number of available physical cores
 423 of the CPU and for large system sizes. However, when using the 48 available cores of the CPU, the
 424 CUDA implementation is again only around 2.5x faster than the CPU. Even for small system sizes the
 425 runtime of CPU and GPU is similar.

426 An additional consideration plays in favor of classical CPU platforms. If a previously computed
 427 random list were used as an input for the model, (thus preventing the time spent in computing random
 428 numbers during simulation, see subsection 3.3.1), even OpenMP Optimized simulation times would
 429 be lower than that of CUDA versions. We recall here that, although this alternative is practical only for
 430 medium lattice widths, it speedups OpenMP around three times. This optimization does not favour so
 431 much GPU platforms.

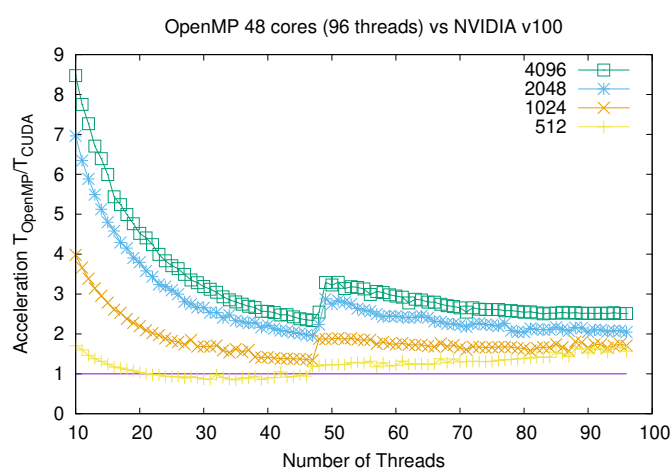


Figure 11. Comparison of OpenMP Optimized and CUDA times: Detail from 10 threads up when using Intel Xeon Platinum and NVIDIA Tesla V100 GPU, resp.

432 5. Conclusions and Future Lines

433 In some previous studies cited in the introduction and section 2, it was pointed out that the
 434 speedups achieved with GPU when comparing to single-core CPUs where above 10X and sometimes

435 above 100x. Even with these ratios, it is evident that a current multicore CPU may approach GPU
436 performance. Moreover, this affirmation must be revised and carefully analyzed for Cellular Automata
437 (CA) applications. In this paper, it is experimentally proved, using almost the same code for a laser
438 dynamics CA (except for the necessary adaptation to each platform), that these distances have been
439 significantly shortened. We conclude that nested conditional structures (in general, many branches in
440 the “hot spot” zones) implies that GPU implementations would suffer from an important deceleration.
441 Even for massive parallel data structures like CA, an approximately 3x speedup is achieved when
442 using high performance computers and GPUs. The other factor that limits CPU and GPU performance
443 for big lattice CAs is obviously the maximum RAM bandwidth, as predicted by the roofline model. If
444 other variables were taken into account like price, TDP, source code maintenance, or easy and rapid
445 software development, it is not clear that GPUs are always the best choice for an efficient parallelization
446 of CA algorithms.

447 Future lines, provided the important conclusions obtained in this paper for the specific case of laser
448 dynamics, include the extension of this analysis to generic Cellular Automata. This will be necessary
449 in order to extract a simple but realistic model that allows to predict the performance on CPU and
450 GPU platforms mainly as a function of the form of its state transition rules, its neighboring relations,
451 the amount of memory, etc. It will be also interesting to investigate the efficient implementation and
452 speedup obtained by implementing this model on Multi-GPUs.

453 **Author Contributions:** Conceptualization, D.C., F.D, and J.L.G.; methodology, D.C., F.D, M.R.L. and J.L.G.;
454 software, D.C., F.D, M.R.L. and J.L.G.; validation, D.C., F.D, and J.L.G.; investigation, D.C., F.D, F.J. and J.L.G.;
455 resources, D.C., F.D, and J.L.G.; writing–original draft preparation, D.C., F.D, F.J. and J.L.G.; writing–review and
456 editing, D.C., F.D, and J.L.G.; visualization, D.C. and M.R.L.; supervision, J.L.G.; funding acquisition, D.C., F.J.
457 and J.L.G.

458 **Funding:** This research was funded by the following research project of Ministerio de Economía, Industria
459 y Competitividad, Gobierno de España (MINECO) and the Agencia Estatal de Investigación (AEI) of Spain,
460 cofinanced by FEDER funds (EU): MABICAP (Bio-inspired machines on High Performance Computing platforms:
461 a multidisciplinary approach, TIN2017-89842P). The work was also partially supported by the computing facilities
462 of Extremadura Research Centre for Advanced Technologies (CETA-CIEMAT), funded by the European Regional
463 Development Fund (ERDF). CETA-CIEMAT belongs to CIEMAT and the Government of Spain.

464 **Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the
465 study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to
466 publish the results.

467 References

- 468 1. Talia, D.; Naumov, N., Parallel cellular programming for emergent computation. In *Simulating Complex*
469 *Systems by Cellular Automata*; Springer, 2010; pp. 357–384.
- 470 2. Bandini, S.; Mauri, G.; Serra, R. Cellular automata: From a theoretical parallel computational model to its
471 application to complex systems. *Parallel Computing* **2001**, *27*, 539–553. doi:10.1016/S0167-8191(00)00076-4.
- 472 3. Wolfram, S. *Cellular automata and complexity*; Addison-Wesley: Reading, MA, 1994.
- 473 4. Ilachinski, A. *Cellular Automata: A Discrete Universe*; World Scientific: Singapore, 2001; p. 808.
- 474 5. Sayama, H. *Introduction to the Modeling and Analysis of Complex Systems*; Open SUNY Textbooks. 14:
475 Geneseo, NY, 2015.
- 476 6. Chopard, B.; Droz, M. *Cellular Automata Modeling of Physical Systems*; Cambridge University Press:
477 Cambridge, MA, USA, 1998.
- 478 7. Sloat, P.; Hoekstra, A. Modeling Dynamic Systems with Cellular Automata. In *Handbook of dynamic system*
479 *modeling*; Fishwick, P., Ed.; Chapman & Hall/CRC, 2007; pp. (21) 1–6.
- 480 8. A.G., H.; Kroc, J.; Sloat, P., Eds. *Simulating Complex Systems by Cellular Automata*; Springer: Berlin,
481 Heidelberg, 2010.
- 482 9. Gounaridis, D.; Chorianopoulos, I.; Koukoulas, S. Exploring prospective urban growth trends under
483 different economic outlooks and land-use planning scenarios: The case of Athens. *Applied Geography* **2018**,
484 *90*, 134–144. doi:10.1016/J.APGEOG.2017.12.001.

- 485 10. Aburas, M.M.; Ho, Y.M.; Ramli, M.F.; Ash'aari, Z.H. The simulation and prediction of spatio-temporal
486 urban growth trends using cellular automata models: A review. *International Journal of Applied Earth*
487 *Observation and Geoinformation* **2016**, *52*, 380–389. doi:10.1016/J.JAG.2016.07.007.
- 488 11. Liu, X.; Liang, X.; Li, X.; Xu, X.; Ou, J.; Chen, Y.; Li, S.; Wang, S.; Pei, F. A future land use simulation model
489 (FLUS) for simulating multiple land use scenarios by coupling human and natural effects. *Landscape and*
490 *Urban Planning* **2017**, *168*, 94–116. doi:10.1016/J.LANDURBPLAN.2017.09.019.
- 491 12. Qiang, S.; Jia, B.; Huang, Q.; Jiang, R. Simulation of free boarding process using a cellular automaton
492 model for passenger dynamics. *Nonlinear Dynamics* **2018**. doi:10.1007/s11071-017-3867-5.
- 493 13. Tang, T.Q.; Luo, X.F.; Zhang, J.; Chen, L. Modeling electric bicycle's lane-changing and retrograde behaviors.
494 *Physica A: Statistical Mechanics and its Applications* **2018**, *490*, 1377–1386. doi:10.1016/J.PHYSA.2017.08.107.
- 495 14. Monteagudo, A.; Santos, J. Treatment Analysis in a Cancer Stem Cell Context Using a Tumor Growth
496 Model Based on Cellular Automata. *PLoS One* **2015**, *10*, e0132306. doi:10.1371/journal.pone.0132306.
- 497 15. Burkhead, E.; Hawkins, J. A cellular automata model of Ebola virus dynamics. *Physica A: Statistical*
498 *Mechanics and its Applications* **2015**, *438*, 424–435. doi:10.1016/J.PHYSA.2015.06.049.
- 499 16. Guisado, J.; Jiménez-Morales, F.; Guerra, J. Cellular automaton model for the simulation of laser dynamics.
500 *Physical Review E* **2003**, *67*, 66708.
- 501 17. Guisado, J.; Jiménez-Morales, F.; Guerra, J. Computational simulation of laser dynamics as a cooperative
502 phenomenon. *Physica Scripta* **2005**, *T118*, 148–152.
- 503 18. Guisado, J.; Jiménez-Morales, F.; Guerra, J. Simulation of the dynamics of pulsed pumped lasers based on
504 cellular automata. *Lecture Notes in Computer Science* **2004**, *3305*, 278–285.
- 505 19. Kroc, J.; Jimenez-Morales, F.; Guisado, J.L.; Lemos, M.C.; Tkac, J. Building Efficient Computational Cellular
506 Automata Models of Complex Systems: Background, Applications, Results, Software, and Pathologies.
507 *Advances in Complex Systems* **2019**, *22*, 1950013.
- 508 20. Guisado, J.; Fernández-de Vega, F.; Jiménez-Morales, F.; Iskra, K. Parallel implementation of a cellular
509 automaton model for the simulation of laser dynamics. *Lecture Notes in Computer Science* **2006**, *3993*, 281–288.
510 doi:10.1007/11758532_39.
- 511 21. Guisado, J.; Jiménez-Morales, F.; Fernández-de Vega, F. Cellular automata and cluster computing: An
512 application to the simulation of laser dynamics. *Advances in Complex Systems* **2007**, *10*, 167–190.
- 513 22. Guisado, J.; Fernandez-de Vega, F.; Iskra, K. Performance analysis of a parallel discrete model for
514 the simulation of laser dynamics. Proceedings of the International Conference on Parallel Processing
515 Workshops. IEEE Computer Society, 2006, pp. 93–99. doi:10.1109/ICPPW.2006.62.
- 516 23. Guisado, J.; Fernández de Vega, F.; Jiménez-Morales, F.; Iskra, K.; Sloat, P. Using cellular automata for
517 parallel simulation of laser dynamics with dynamic load balancing. *International Journal of High Performance*
518 *Systems Architecture* **2008**, *1*, 251–259.
- 519 24. GPGPU. General-Purpose Computation on Graphics Hardware. <http://gpgpu.org>, as available on may
520 2012., 2012.
- 521 25. Lopez-Torres, M.; Guisado, J.; Jimenez-Morales, F.; Diaz-del Rio, F. GPU-based cellular automata
522 simulations of laser dynamics. Proceedings of the XXIII Jornadas de Paralelismo; hgpu.org: Elche,
523 2012; pp. 261–266.
- 524 26. Gobron, S.; Devillard, F.; Heit, B. Retina simulation using cellular automata and GPU programming.
525 *Machine Vision and Applications Journal* **2007**, *18*, 331–342.
- 526 27. Rybacki, S.; Himmelspach, J.; Uhrmacher, A. Experiments With Single Core, Multi Core, and {GPU}-based
527 Computation of Cellular Automata. Advances in System Simulation, 2009. {SIMUL}'09. First International
528 Conference on, 2009.
- 529 28. Bajzát, T.; Hajnal, E. Cell Automaton Modelling Algorithms: Implementation and Testing in {GPU} Systems.
530 {INES} 2011, 15th International Conference on Intelligent Engineering Systems, 2011.
- 531 29. Balasalle, J.; Lopez, M.; Rutherford, M., Optimizing Memory Access Patterns for Cellular Automata on
532 {GPU}s. In *GPU Computing Gems Jade Edition*; Elsevier - Morgan Kaufmann - NVIDIA, 2011; pp. 67–75.
- 533 30. Geist, R.; Westall, J., Lattice-Boltzmann Lighting Models. In *GPU Computing Gems, Emerald Edition*; Elsevier
534 - Morgan Kaufmann - NVIDIA, 2011; pp. 381–399.
- 535 31. Gibson, M.J.; Keedwell, E.C.; Savić, D.A. An investigation of the efficient implementation of cellular
536 automata on multi-core CPU and GPU hardware. *Journal of Parallel and Distributed Computing* **2015**,
537 *77*, 11–25. doi:10.1016/j.jpdc.2014.10.011.

- 538 32. Jongerius, R.; Anghel, A.; Dittmann, G. Analytic Multi-Core Processor Model for Fast Design-Space
539 Exploration. *IEEE TRANSACTIONS ON COMPUTERS* **2018**, *67*.
- 540 33. NVIDIA. CUDA C Best Practices Guide Version. Available in <http://developer.nvidia.com/>, Accessed in
541 Dec. 2019.
- 542 34. Williams, S.; Waterman, A.; Patterson, D. Roofline: an insightful visual performance model for multicore
543 architectures. *Commun. ACM* **2009**, *52*, 65–76.
- 544 35. Hennessy, J.; Patterson, D. *Computer Architecture: A Quantitative Approach, 6th Edition*; Vol. 19,
545 Elsevier-Morgan Kaufmann, 2017.