

Using GPU for Multi-agent Soil Simulation

Guillaume Laville, Kamel Mazouzi,
Christophe Lang and Laurent Philippe
Institut FEMTO-ST/DISC
Université de Franche-Comté, France
Email: guillaume.laville@univ-fcomte.fr

Nicolas Marilleau
UMI 209 UMMISCO IRD-UPMC
Bondy, France
Email: nicolas.marilleau@ird.fr

Abstract—Multi-Agent Systems (MAS) can be used to model systems where the global behavior cannot be uniformly represented by standard techniques such as partial differential equations or linear systems because the system elements have their own independent behavior. This is, for instance, the case in complex systems such as daily mobility in a city for example. Depending on the system size the computing power needs for the MAS may be as big as for more traditional linear numerical systems and may need to be parallelized to fully represent real systems. Graphical Processing Units (GPU) have already proven to be an efficient support to execute large linear programs. In this paper we present the use of GPU for the execution of Sworm, a multi-scale MAS system. We show that GPU computing can be efficient in that less regular case and when the agent behavior is simple. We advocate for a wider use of the GPU in Agent Based Models in particular for multi-scale systems with work distribution between the CPU and GPU.

Keywords-multi-agent systems; GPU; parallel computing

I. INTRODUCTION

Physical systems are often modeled by mathematical representations and their behavior simulated by dynamic differential equations. This is possible when the behavior of the elements of the system is uniformly driven by the same law, but when these elements have their own behavior, the modeling process is too complex to rely on this approach. Multi-Agent Systems (MAS) is a recognized approach to model and simulate systems where individuals have an autonomous behavior that cannot be simulated by the evolution of a set of variables driven by mathematical laws. MAS are often used to simulate natural or collective phenomena whose actors are too numerous or various to provide a unified algorithm describing the system evolution. The agent-based approach is to divide these complex systems into individual, self-contained entities with their smaller set of attributes and functions. But, as for mathematical simulations, when the size of the Multi-Agent System increases the need of computing power and memory also increases. For this reason, multi-agent systems should benefit from the use of distributed computing architectures. Clusters and Grids are often identified as the main solution to increase simulation performance but Graphical Processing Units (GPU) are also a promising technology with an attractive performance/cost ratio.

Conceptually an Agent Based Model (ABM) is a distributed system as it favors the definition and description of large sets of individuals, the agents, that can be run in parallel. Most

of the agent-based simulators are however designed with a sequential scheme in mind and these simulators seldom use more than one core for their execution. Due to simulation scheduling constraints, data sharing and exchange between agents and the huge amount of interactions between agents and their environment, it is indeed rather difficult to distribute an agent based simulator to take advantage of new multi-threaded computer architectures. Thus, guidelines and tools dedicated to MAS paradigm and HPC is now a need for other complex system communities. Note that, from the described structure (large number of agents sharing data), we can conclude that MAS would more easily make good use of many-cores computing than from other kinds of parallelism.

Another fact that advocates for the use of many-core in MAS is the growing need for multi-scale simulations. Multi-scale simulations explore problems with interactions between several scales. The different scales use different granularities of the structure and potentially different models. Most of the time the lower scale simulations provide results to higher scale simulations. In that case the execution of the simulations can easily be distributed between the local cores and a many-core architecture, i.e. a GPU device.

We explore in this paper the use of many-core architectures to execute many-core simulations. Our case study is the Sworm simulator that aims at reproducing effects of earthworms on bacteria dynamics in a bulked soil. In Section ?? of the paper we present the work related to MAS and parallelization with a special focus on many-core use. We give an overview of the Sworm simulator in Section ?? and we detail its GPU implementation in Section ?. We present the experiments and results in Section ?? then we conclude on the possible generalization of our work.

II. GPU AND MULTI-AGENT SYSTEMS

A. Agent-Based Simulations

Agent-Based systems are often used to simulate natural or collective phenomena whose actors are too numerous or various to provide a simple unified algorithm describing the studied system dynamic [?]. The implementation of an agent based simulation usually starts by designing the underlying Agent Based Model (ABM). Most ABM are based around a few types of entities: Agent, Environment and Interaction Organization [?]. In the complex system domain, the environment often describes a real space, its structure (e.g. matter

organization in soils) and its dynamics (e.g. organic matter decomposition). It is a virtual world in which agents represent studied entities (e.g. biotic organisms) evolving. The actual agent is animated by a behavior that can range between reactivity (only react to external stimuli) and cognition (lives its own process alongside other individuals). Interaction and Organization define functions, types and patterns of communications of their member agents in the system [?]. Note that, depending on the MAS, agents can communicate either directly through special primitives or indirectly through the information stored in the environment.

Agent based simulations have been used for more than one decade to simulate complex systems and have proved their interest in various scientific communities. Nowadays generic agent based frameworks are promoted such as Repast [?] or NetLogo [?] to implement simulations. Many ABM such as the crown model representing a city-wide scale [?] tend however to require a large number of agents to provide a realistic behavior and reliable global statistics. Moreover, an achieved model analysis needs to run many simulations identified into an experiment plan to obtain enough confidence in a simulation. In this case the available computing power often limits the simulation size and the result range thus requiring the use of parallelism to explore bigger configurations.

For that, three major approaches can be identified:

- 1) parallelizing experiments plan on a cluster or a grid (one or few simulations are submitted to each core) [?], [?]
- 2) parallelizing the simulator on a cluster (the environment of the MAS is split and run on several distributed nodes) [?], [?]
- 3) optimizing simulator by taking advantage of computer resources (multi-threading, GPU, and so on)[?]

In the first case, the experiments are run independently of one another and only the simulation parameters change between two runs. A simple version of an existing simulator can thus be used. It does not imply model code changes except for Graphical Unit Interface extracting as the experiments does not run interactively but in batch. In the second and the third case, model and code modifications are necessary. Only few frameworks however introduce distribution in agent simulation (Madkit [?] or MASON [?]) and parallel implementations are often based on the explicit use of threads using shared memory [?] or cluster libraries such as MPI [?].

Parallelizing a multi-agent simulation is however complex due to space and time constraints. Multi-agent simulations are usually based on a synchronous execution of time steps by the agents that share the same environment. At each time step, numerous events (space data modification, agent motion) and interactions between agents happen. Distributing the simulation on several computers or grid nodes thus implies to guaranty a distributed synchronous execution and coherency, which often leads to poor performances or complex synchronization problems. Multi-cores execution or delegating part of this execution to others processors as GPUs [?] are usually easier to implement since all the threads share the data and the local clock.

Different agent patterns can be adopted in an ABMs such as cognitive and reactive ones [?]. Cognitive agents act on the environment and interact with other agents according to a complex behavior. This behavior takes a local perception of the virtual world and the agent past (a memory characterized by an internal state and belief, imperfect knowledge about the world) into account. Reactive agents have a much systematic pattern of action based on stimuli response schemes (no or few knowledge and state conservation in agent). The evolution of the ABM environment, in particular, is often represented with this last kind of agents. As their behavior is usually simple, we propose in this paper to delegate part of the environment and of the reactive agents execution to the graphical processing unit of the computer. This way we can balance the load between both CPU and GPU execution resources.

In the particular case of multi-scale simulations such as the Swarm simulation [?] the environment may be used at different levels. Since the representation of the whole simulated soil area would be costly, the environment is organized as a multi-level tree of small soil cubes which can be lazily instantiated during the simulation. This allows to gradually refine distribution details in units of soil as agents progress and need those informations, by using a fractal process based on the bigger-grained already instantiated levels. This characteristic, especially for a fractal model, could be the key of the distribution. For instance, each branch of a fractal environment could be identified as an independent area and parallelized. In addition Fractal is a famous approach to describe multi-scale environment (such as soil) and its organization [?]. In that case the lower scale simulations can also be delegated to the GPU card to limit the load of the main (upper scale) simulation.

B. Introduction to GPU computing

Graphics Processing Unit devices are based on a massively-parallel hardware architecture, which allows the execution of a big amount of floating point instructions at the same time and their computing power is often measured in hundred of GFlops. These units were originally designed to be used by the graphic driver for multimedia [?] or graphic processes. The last few years saw however the apparition of new generations of graphic cards based on more general purpose execution units. These units can be programmed using GPGPU platforms to solve various problems in HPC applications such as linear algebra resolutions. They can be used both for dense problems, such as CuBLAS, a GPU implementation of the BLAS library, MAGMA [?] and sparse ones such as [?][?]. These applications are based on matrix manipulations where a set of operations is applied multiple times on each input data. The same kind of treatments are used in geometric computations which also suit the GPU particularly and yield to big performance gains.

These matrix-based data representation and SIMD computations are not so straightforward in Multi-Agent Systems, where data structures and algorithms are tightly coupled to the described simulation. However, works from existing literature show that MAS can benefit from these performance gains on

various simulation types, such as traffic simulation [?], cellular automata [?], mobile-agent based path-finding [?] or genetic algorithms [?].

An application-specific adaptation process was required in the case of these MAS: some of the previous examples are driven by mathematical laws (path-finding) or use a natural mapping between a discrete environment (cellular automaton) and GPU cores. Unfortunately, this mapping often requires algorithmic adaptations in other models but the experience shows that the more reactive a MAS is the more adapted its implementation is to GPU.

C. MAS Implementation on GPU

The first step in the adaptation of an ABM to GPU platforms is the choice of the language. On the one hand, the Java programming language is often used for the implementation of MAS due to its availability on numerous platforms or frameworks and its focus on high-level, object-oriented programming. On the other hand, GPU platforms can only run specific languages as OpenCL or CUDA. OpenCL (supported on AMD, Intel and NVIDIA hardware) better suits the portability concerns across a wide range of hardware needed the agent simulators, as opposed to CUDA which is a NVIDIA-specific library.

OpenCL is a C library which provides access to the underlying CPUs or GPUs using an asynchronous interface. Various OpenCL functions allow the compilation and the execution of programs on these execution resources, the copy of data buffers between devices, or the collection of profiling information.

This library is based around three main concepts:

- the *kernel*, which represents a runnable program containing instructions to be executed on the GPU.
- the *work-item* (or task), which is analogous to the concept of thread on GPU, in that it represents one running instance of a GPU kernel.
- the *work-group* (or task-group) which is a set of work-items sharing some memory to speed up data accesses and computations. Synchronization operations such as barrier can only be used across the same work-group.

Running an OpenCL computation consists in launching numerous work-items that execute the same kernel. The work-items are submitted to a submission queue to optimize the available cores usage. A calculus is achieved once all these kernel instances have terminated.

The number of work-items used in each work-group is an important implementation choice which determines how many tasks will share the same cache memory. Data used by the work-items can be stored as N-dimensions matrices in local or global GPU memory. Since the size of this memory is often limited to a few hundred of kilobytes, choosing this number often implies a compromise between the model synchronization or data requirements and the available resources.

In the case of agent-based simulations, each agent can be naturally mapped to a work-item. Work-groups can then be used to represent groups of agents or simulations sharing

common data (such as the environment) or algorithms (such as the background evolution process).

III. THE MIOR MODEL

The MIOR [?] model was developed by the GEODES Team of IRD (ex UMMISCO) to simulate local interactions in a soil between microbial colonies and organic matters. It reproduces each small cubic units (0.002) of soil as a MAS.

Multiple implementations of the MIOR model have already been realized, in Smalltalk and Netlogo, in 2 or 3 dimensions. The last implementation, used in our work and referenced as MIOR in the rest of the paper, is freely accessible online as WebSimMior¹.

MIOR is based around two types of agents: (i) the Meta-Mior (MM), which represents microbial colonies consuming carbon and (ii) the Organic Matter (OM) which represents carbon deposits occurring in soil.

The Meta-Mior agents are characterized by two distinct behaviors:

- *breath*: the action converts mineral carbon from the soil to carbon dioxide CO_2 that is released in the soil.
- *growth*: by this action each microbial colony fixes the carbon present in the environment to reproduce itself (augments its size). This action is only possible if the colony breathing needs were covered, i.e. enough mineral carbon is available.

These behaviors are described in the following algorithm ??.

Algorithm 1 Evolution step of each Meta-Mior (microbial colony) agent

```

Input: mmList List of MM
Input: omList List of OM
Input: world Environment MIOR
breathNeed  $\leftarrow$  world.respirationRate  $\times$  mm.carbon
growthNeed  $\leftarrow$  world.growthRate  $\times$  mm.carbon
availableCarbon  $\leftarrow$  totalAccessibleCarbon(mm)
if availableCarbon > breathNeed then
  // Breath
  mm.active  $\leftarrow$  true
  availableCarbon  $\leftarrow$  availableCarbon -
  consumCarbon(mm, breathNeed)
  world.CO2  $\leftarrow$  world.CO2 + breathNeed
  if availableCarbon > 0 then
    // Growth
    growthConsum  $\leftarrow$  max(totalAccessCarbon(mm),
    growthNeed)
    consumCarbon(mm, growthConsum)
    mm.carbon  $\leftarrow$  mm.carbon + growthConsum
  end if
else
  mm.active  $\leftarrow$  false
end if

```

Since this simulation takes place at a microscopic scale, a large number of these simulations must be executed at each macroscopic simulation step to model realistic-sized unit of soil, despite the small computation cost of each individual simulation.

¹<http://www.IRD.fr/websimmior/>

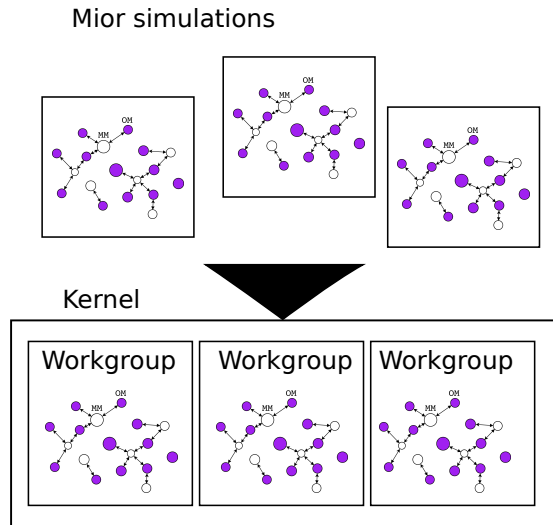


Fig. 1. Execution distribution retained on GPU

IV. MIOR IMPLEMENTATION ON GPU

As a starting point, we realized a simple GPU implementation of the MIOR simulator, with only minimal changes to the CPU algorithm. Execution times showed the inefficiency of this approach, and highlighted the necessity of adapting the simulator to take advantage of the GPU execution capabilities. In this part, we show the main changes which were realized to adapt the MIOR simulator on GPU architectures.

A. GPU Architecture of the MIOR model

Each MIOR simulation is represented by a work-group, and each agent by a work-item. A kernel is responsible of the life cycle process of each agent of the model. It is executed by all work-items of the simulation on their own gpu core.

The usage of a work-group for each simulation allows to easily execute multiple simulations in parallel, as shown on figure ???. Thus it becomes possible to exploit all the cores at the same time by taking advantage of execution overlap possibilities provided by OpenCL, even if an unique simulation is too small to use all the available GPU cores. However, the maximum size of a work-group is limited (512), which only allows use to execute one simulation per work-group when using 310 threads (number of OM in the reference model) to execute the simulation.

The usage of the GPU to execute multiple simulations is initiated by the CPU. The latter keeps total control of the simulator execution flow. Thus, optimized scheduling policies (such as submitting kernels in batch, or limiting the number of kernels, or asynchronously retrieving the simulation results) can be defined to minimize the cost related to data transfers between CPU and GPUs.

B. Data structures translation

The adaptation of the MIOR model to GPU requires the mapping of the data model to OpenCL data structures. The

environment and the agents are represented by arrays of structures, where each structure describes the state of one entity. The behavior of these entities are implemented as OpenCL functions to be called from the kernels during execution.

Four main data structures are defined: (i) an array of MM agents, representing the state of the microbial colonies. (ii) an array of OM agents, representing the state of the carbon deposits. (iii) a topology matrix, which stores accessibility information between the two types of agents of the model (iv) a world structure, which contains all the global input data (metabolism rate, numbers of agents) and output data (quantity of CO_2 produced) of the simulation.

These data structures are initialized by the CPU and then copied on the GPU.

The world topology is stored as a two-dimension matrix which contains OM indexes on the abscissa and the MM index on the ordinate. Each agent walks through its line/column of the matrix at each iteration to determine which agents can be accessed during the simulation. Since many agents are not connected this matrix is sparse, which introduces a big number of useless memory accesses. To reduce the impact of these memory accesses we propose a compacted, optimized representation of this matrix based on [?], as illustrated on the figure ???. Since memory allocations are impossible yet in OpenCL and only provided in the latest revisions of the CUDA standard, these matrices are statically allocated to handle the worst-case scenario where all OM and MM are linked, since the actual occupation of the matrix cells cannot be deduced without some kind of preprocessing computations.

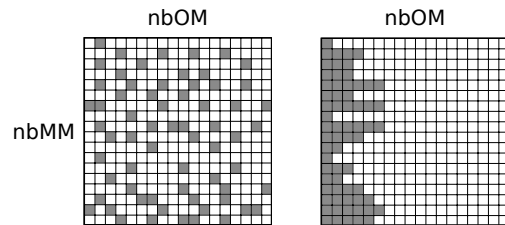


Fig. 2. Compact representation of the topology of a MIOR simulation

The compact representation considers each line of the matrix as an index list, and only stores accessible agents continuously, to reduce the size of the list and the number of non productive accesses.

C. Critical resources access management

One of the main issue in a MIOR model is to ensure that all the microbial colonies will have an equitable access to carbon resources, when multiple colonies share the same deposits. Access synchronizations are mandatory in these cases to prevent conflicting updates on the same data that may lead to calculus error (e.g. loss of matter).

On a massively parallel architecture such as GPUs, this kind of synchronization conflicts can however lead to an inefficient implementation by enforcing a quasi-sequential execution. It is necessary, in the case of MIOR as well as for other ABM, to ensure that each work-item is not too constrained in its execution.

From the sequential algorithm 1 where all the agents share the same data, we have developed a parallel algorithm composed of three sequential stages separated by synchronization barriers. This new algorithm is based on the distribution of the available OM carbon deposits into parts at the beginning of each execution step. The three stages, illustrated on figure ??, are the following:

- 1) *distribution*: The available carbon in each carbon deposit (OM) is equitably dispatched between all accessible MM in the form of parts.
- 2) *metabolism*: Each MM consumes carbon in its allocated parts for its breathing and growing processes.
- 3) *gathering*: Unconsumed carbon in parts is gathered back into the carbon deposits.

This solution suppresses the data synchronization needed by the first algorithm and thus the need for synchronization barriers.

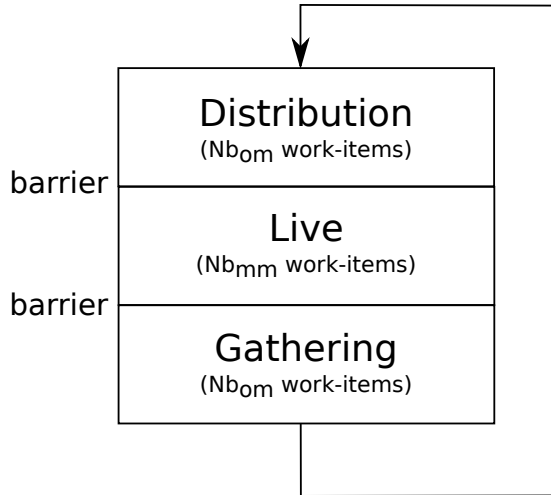


Fig. 3. Decomposition in OpenCL kernels of one MIOR parallel execution cycle.

D. Termination detection

The termination of a MIOR simulation is reached when the model stabilizes, and no more CO_2 is produced. This detection can be done either on the CPU or the GPU.

In the first case, it is the CPU which controls the GPU simulation process: (i) the CPU starts the execution of a simulation step on the GPU. (ii) Once this step is finished, the CPU retrieves the GPU data, and determines if another iteration must be launched or if the simulation has terminated. This approach allows a fine-grain control over the GPU execution, but it requires many costly transfers as each iteration results must be sent from the GPU to the CPU. In the case of the

MIOR model, these costs are mainly due to the inherent PCI-express port latencies, rather than to bandwidth limitation, since data sizes remains rather small, in the order of few dozens of Megabytes.

In the second case the termination detection is directly implemented on the GPU: The CPU does not have any feedback while the simulation is running, but retrieves the results once the kernel execution is finished. This approach minimizes the number of transfers between the CPU and the GPU.

V. EXPERIMENTS

In this section we present several MIOR GPU implementations using the distribution/gathering process described in Section ?? and compare their performance on two distinct hardware platform; i.e two different GPU devices. Five incremental MIOR implementations were realized with an increasing level of adaptation for the algorithm: In all cases, we choose the average time over 50 executions as a performance indicator.

- the **GPU 1.0** implementation is a direct implementation of the existing algorithm and its data structures where data dependencies were removed and using the non-compact topology representation described in Section ??
- the **GPU 2.0** implementation uses the previously described compact representation of the topology and remains otherwise identical to the GPU 1.0 implementation.
- the **GPU 3.0** implementation introduces the manual copy into local (private) memory of often-accessed global data, such as carbon parts or topology information.
- the **GPU 4.0** implementation is also a variant of the GPU 1.0 implementation which allows the execution of multiple simulations for each kernel execution.
- the **GPU 5.0** implementation is a multi-simulation of the GPU 2.0 implementation, using the execution of multiple simulations for each kernel execution as for GPU 4.0

The two last implementations – **GPU 4.0** and **GPU 5.0** – illustrate the gain provided by a better usage of the hardware resources, thanks to the driver execution overlap capabilities. A sequential version of the MIOR algorithm, labeled as **CPU**, is provided for comparison purpose. This sequential version is developed in Java, the same language used for GPU implementations and the Swarm model.

For these performance evaluations, two platforms are used. The first one is representative of the kind of hardware which is available on HPC clusters. It is a cluster node dedicated to GPU computations with two Intel X5550 processors running at 2.67GHz and two Tesla C1060 GPU devices running at 1.3GHz and composed of 240 cores (30 multi-processors). Only one of these GPU is used in these experiments, at the moment. The second platform illustrates what can be expected from a personal desktop computer built a few years ago. It uses a Intel Q9300 CPU, running at 2.5GHz, and a Geforce 8800GT GPU card running at 1.5GHz and composed of 112 cores (14 multi-processors). The purpose of these two platforms is to assess the benefit that could be obtained either when a scientist

has access to specialized hardware as a cluster or tries to take benefit from its own personal computer.

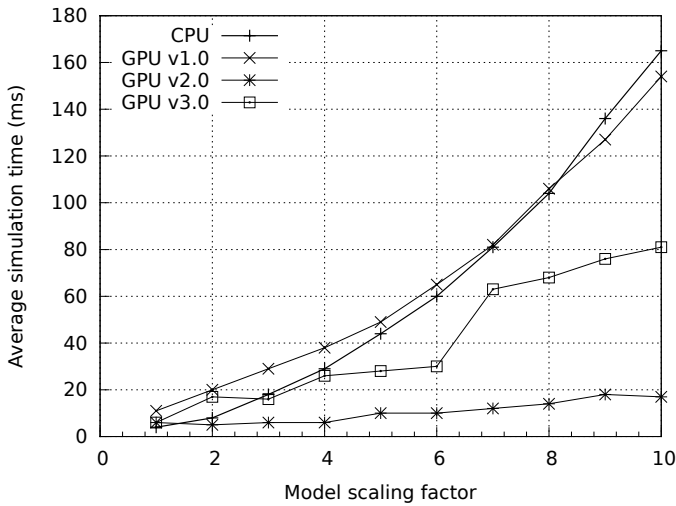


Fig. 4. CPU and GPU performance on a Tesla C1060 node

Figures ?? and ?? show the execution time for 50 simulations on the two hardware platforms. A size factor is applied to the problem: at scale 1, the model contains 38 MM and 310 OM, while at the scale 6 these numbers are multiplied by six. The size of the environment is modified as well to maintain the same average agent density in the model. This scaling factor display the impact of the chosen size of simulation on performance.

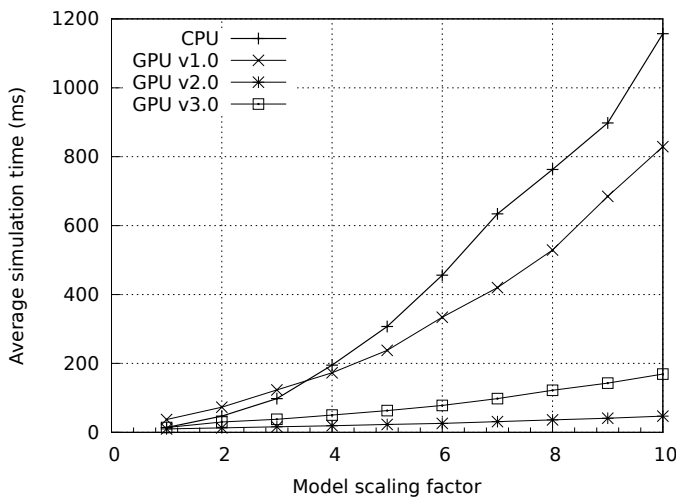


Fig. 5. CPU and GPU performance on a personal computer with a Geforce 8800GT

The charts show that for small problems execution times of all implementations are very close. This is because the GPU implementation does not have enough threads (representing agents) for an optimal usage of GPU resources. This trend changes after scale 5 where GPU 2.0 and GPU 3.0 take the advantage on the GPU 1.0 and CPU implementations. This advantage continues to grow with the scaling factor, and

reaches a speedup of 10 at the scale 10 between the fastest single-simulation GPU implementation and the first, naive one GPU 1.0.

Multiple trends can be observed in these results. First, optimizations for the GPU hardware show a big, positive impact on performances, illustrating the strong requirements on the algorithm properties to reach execution efficiency. These charts also show that despite the vast difference in numbers of cores between the two GPUs platforms the same trends can be observed in both case. We can therefore expect similar results on other GPU cards, without the need for more adaptations.

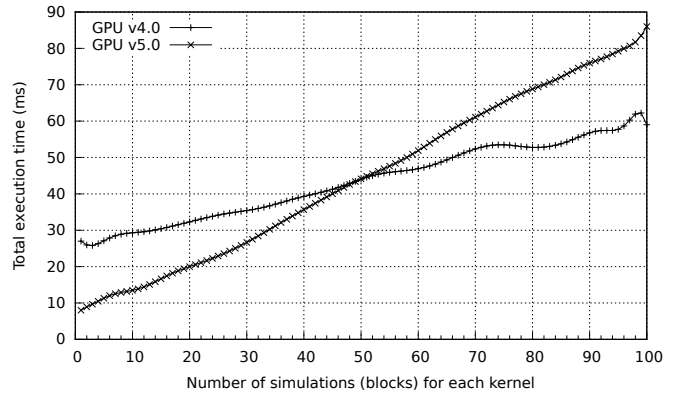


Fig. 6. Execution time of one multi-simulation kernel on the Tesla platform

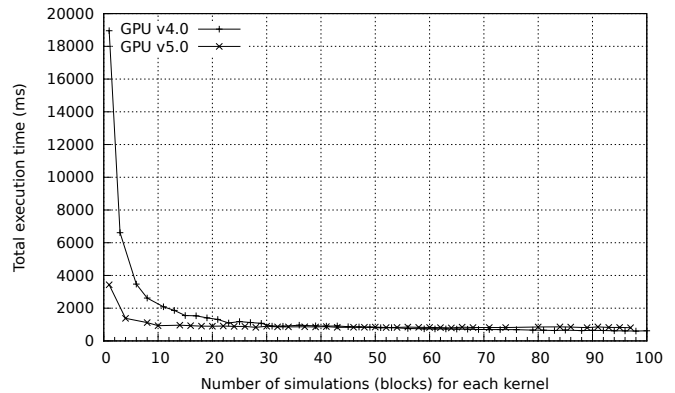


Fig. 7. Total execution time for 1000 simulations on the Tesla platform, while varying the number of simulations for each kernel

There are two ways to measure multi-simulations performance: (i) by executing only one kernel, and varying its size (the number of simulations executed in parallel), as seen on Figure ??), to test the costs linked to the parallelization process or (ii) by executing a fixed number of simulations (Figure ??) and varying the size of each kernel.

Figure ?? illustrates the execution time for only one kernel. It shows that for small numbers of simulations run in parallel, the compact implementation of the model topology is faster than the two-dimension matrix representation. This trends reverse with more than 50 simulations in parallel, which

can be explained either by the non linear progression of the synchronization costs or by the additional memory required for the access-efficient representation.

Figure ?? illustrates the execution time of a fixed number of simulations. It shows that for a small number of simulations run in parallel, the costs resulting of program setup, data copies and launch on GPU are determinant and very detrimental to performance. Once the number of simulations executed for each kernel grows, these costs are counterbalanced by computation costs. This trend is more marked in the case of the sparse implementation (GPU 4.0) than the compact one but appears on both curves. With more than 30 simulations for each kernel, execution times stalls, since hardware limits are reached. This indicates that the cost of preparing and launching kernels become negligible compared to the computing time once a good GPU occupation rate is achieved.

used for the computations described in the paper.

VI. CONCLUSION

This paper addresses the issue of complex system simulation by using agent based paradigm and GPU hardware. From the experiment on an existing Agent Based Model of soil science (Mior model) we intend to provide useful information on the architecture, the algorithm design and, the data management to run Agent based simulations on GPU. The first result of this work is that adapting the algorithm to a GPU architecture is possible and suitable to speed up agent based simulations. Coupling CPU with GPU seems even to be an interesting way to better take advantage of the power given by computers and clusters: cognitive agents can be run on the CPU and reactive agents or environment be run on the GPU. Note that simulations of reactive agents, with a simple behavior, give interesting and efficient results whereas the running of cognitive agent (such as social agents) on GPU are too complicated and less adapted to this kind of execution. From our point of view this adaptation process is lighter than a traditional parallelization on distributed nodes and not much difficult than a standard multi-threaded parallelization, since all the data remains on the same host and can be shared in central memory. The OCL adaptation also enables a portable simulator that can be run on different graphical units. Even using a mainstream card as the GPU card of a standard computer can lead to significant performance improvements. This is an interesting result as it opens the field of inexpensive HPC to the ABM community.

In this perspective, we are working on a generalization of these tools to other multi-agent models based around reactive computations. One first challenge is the definition of common, efficient, reusable data structures, such as grids or lists. Another goal of this solution is to provide easier means to control the distribution of specific processes on CPU or GPU, to allow the easy exploitation of the strengths of each platform in the same multi-agent simulation.

ACKNOWLEDGMENT

The authors would like to thanks the Mésocentre de calcul de Franche-Comté, which provided the computer facilities