P Systems Implementation on GPUs

²6

Δ

6.1 Introduction

The development of P system simulators is usually driven by the importance 5 of certain models. Usually, these simulators are implemented in a flexible way, 6 allowing not only to simulate a wide variety of P systems but also to help 7 construct simulators for other models. An example of this flexibility is P-Lingua 8 framework[13]. However, for certain applications and models, efficient simulation 9 tools are required. For instance, the simulation of population dynamics P systems 10 is crucial for model validation of real ecosystems and for virtual experimentation. 11 In this case, the faster the simulation tool, the shorter the time to construct a 12 valid model. Another interest behind the development of efficient tools is also for 13 analyzing theoretical aspects of P systems (parallelism, non-determinism, etc.) and 14 how to bridge them with today, *in-silico* technology [48].

There are several ways to accelerate the simulation of P systems: changing the ¹⁶ technology where to implement the simulators (e.g., from interpreted languages like ¹⁷ Java to compiled ones like C++), increasing the power of the processors where to ¹⁸ run the simulations (e.g., increasing the clock frequency, the memory bandwidth and ¹⁹ clock, etc.), or using high performance computing (HPC) technologies to implement ²⁰ real parallelism. The main trend when developing efficient simulators has been the ²¹ last one: taking advantage of the inherent parallelism of the models and mapping it ²² into parallel platforms such as clusters, supercomputers, accelerators, etc. ²³

According to [15], we can define high performance computing (HPC) as "the ²⁴ practice of aggregating computing power in a way that delivers much higher ²⁵ performance than one could get out of a typical desktop computer or workstation ²⁶ in order to solve large problem instances in science, engineering, or business". ²⁷ This is usually accomplished by means of parallelism, since it is the basis for ²⁸ the acceleration of large and complex real-world applications. The maximum ²⁹ exponent of HPC is known as supercomputing, where the computing power of ³⁰ current technology is being continuously pushed. A ranking of the most powerful ³¹

supercomputers can be consulted at Top500 website [46]. To the time of writing, ³² most of the top 10 supercomputers are based on nodes extended with accelerators. ³³

HPC accelerators are dedicated chips that serve as co-processors, extending in ³⁴ this way the computing power. Examples of accelerators are FPGAs and GPUs. ³⁵ The latter refers to the processors inside the graphics cards, which take over the ³⁶ task of graphics generation in computers. However, with the increasing demand ³⁷ in 3D rendering for gaming and video, and as foreseen by Elster [12] and others, ³⁸ GPUs (Graphics Processing Units) have evolved to a massively parallel processor ³⁹ that is suitable for parallel computing. Today, GPUs are the enabling technology for ⁴⁰ trending areas such as Deep Learning, Data Science, physics simulation, real-time ⁴¹ ray tracing graphics, etc. ⁴²

Concerning P systems and their applications, GPUs have been shown to be an 43 alternative for accelerating simulations. In [4], the double parallelism of P systems 44 were mapped over the double parallelism inside GPUs. This idea has been refined 45 over the time, in such a way that the simulators are now better adapted to GPU 46 parallelism. We can identify three types of simulators: those developed for very 47 specific P systems or family of P systems (specific simulators) and developed for a 48 wide range of P systems inside a variant (generic simulators) and a hybrid simulator 49 that receives high-level information to be better adapted (adaptive simulators). 50

In this chapter, we will introduce all the concepts related with GPU computing 51 and its applications. Later, we will go through some P system simulators depending 52 on their type: specific, generic, or hybrid. Finally, we will provide some guidelines 53 on how to develop new simulators for P systems on GPUs. 54

6.2 GPU Computing

In this section, we will introduce the main concepts of GPU computing, including 56 CUDA and modern GPU architecture. This will provide the required background to 57 understand the design of P system simulators on GPUs. 58

6.2.1 The Graphics Processing Unit

The first Graphics Processing Units (GPUs) were introduced back in 1999 [11], in 60 order to overcome the bottleneck created by the CPU when generating real-time 61 graphics, such as in videogames and in 3D rendering. Since then, every graphics 62 card has integrated such kind of specific processors. Usually, we refer to GPU and 63 graphics card as synonyms. It is important to remark the place where the GPU is 64 located on a computer. The GPU is connected with the CPU through a data bridge 65 (Northbridge), which is also used to access the main memory (RAM). Currently, 66 modern GPUs are connected through the so called PCI Express bus, which runs at 67 more than 64GB/s (when using 16-lane configuration). 68

GPUs, since they were born, install processor cores that are specialized for 69 graphics (pixel colorization, etc.). Thus, these processors are able to include more 70

55

cores than usual CPUs because they are more specific and, so, lightweight. This 71 GPU architecture has evolved over the time, being both more parallel and flexible. 72 The former means that it included more and more computing cores, and the 73 latter means that these cores were more programmable. In fact, since 2002, the 74 graphics pipeline implemented in GPU hardware became programmable through 75 small programs called shaders. Programming languages such as Cg, GLSL, DirectX 76 Shading languages, etc. are employed for shaders. There were two types of shaders 77 (for vertices and fragments), but in 2007 they were unified (e.g., the U of CUDA 78 stands for Unified). 79

The evolution of shaders led to a new area called GPGPU (general purpose ⁸⁰ computing on the GPU), whose name was coined by Mark Harris in 2002. The ⁸¹ target of this research area is to develop parallel methodologies to program GPUs ⁸² for other purposes rather than graphics, such as scientific computing. Today, this ⁸³ is more known as GPU computing and has been settled as an alternative within ⁸⁴ HPC. That is, GPUs are nowadays an HPC accelerator that can be found in the most ⁸⁵ powerful supercomputers. GPUs are good at data parallelism. More specifically, ⁸⁶ they are based on SPMD programming model (Single Program Multiple Data): the ⁸⁷ GPU processes many independent elements in parallel using the same program [35]. ⁸⁸

Currently, GPU computing is also a heterogeneous computing system [17], 89 where the GPU is known as device and the CPU is known as host. The host has 90 a role of a master, which takes over the execution and manages the different devices 91 that can be in the system. Devices are co-processors that help to accelerate the 92 algorithms by executing code in a parallel fashion, reducing in this way the overhead 93 on the host. This trend is being consolidated with OpenCL [31], the first free, 94 open standard for multi-platform, heterogeneous parallel programming of modern 95 processors found in PCs, servers, and embedded devices. OpenCL is being used not 96 only for GPUs but also for FPGAs, multicore CPUs, etc. However, the drivers and 97 compilers developed by each manufacturer of chips (NVIDIA, Intel, AMD...) are 98 not up to date and lack full support. This is why a new standard, called SYCL, is 99 being conceived, but it is still experimental (to date).

GPUs can be programmed with both OpenCL and SYCL [44]. Moreover, 101 NVIDIA GPUs can be also programmed with CUDA [32], which is a proprietary 102 technology that is very mature and has lot of functionality. By having a quick look 103 to the literature, it is possible to see that CUDA is the most used platform for 104 GPU computing. On the other side, AMD GPUs can be also programmed with a 105 CUDA-like environment called HIP [34], which is based on RoCm. This allows 106 programmers to translate easily CUDA code to AMD technology. Other standard 107 languages and platforms to program GPUs are based on the graphics pipeline, such 108 as OpenGL, GLSL, and Vulkan [47]. They can also be used for GPU computing in 109 a not very complex way. 110

In short, GPU computing poses a highly parallel architecture with thousands of 111 lightweight processor cores and high-bandwidth memory that can be programmed 112 with several standard languages. The most evolved one is CUDA, but it works 113 only for NVIDIA GPUs. Since the introduction of CUDA in 2007, many scientific 114 applications have used GPU computing. Their low cost compared to the perfor- 115

mance offered has made GPUs an attractive alternative. In fact, currently they are 116 the enabling technology (i.e., if you want to tackle these problems, you should be 117 using GPUs) for Deep Learning [18], nanopore DNA sequencing, and high-energy 118 particle trajectory reconstruction in LHCb HLT 1 TODO:CITE PAPER DANIEL 119 CAMPORA. 120

6.2.2 CUDA Programming Model

In this chapter, we will be focusing on CUDA [32, 33], since, as mentioned, it is the 122 most widely technology for GPU computing, and on top of that, the majority of P 123 system simulators have been developed with CUDA. Let us recall that CUDA, as 124 usual in GPU computing, offers a heterogeneous system to the user, where the CPU 125 is known as the host, and the GPU is the device. The execution flow in a CUDA 126 program is like in any common program; it starts with the CPU main function. At 127 some points, the CPU asks the GPU to allocate memory, transfer memory, launch 128 computation, retrieve results, etc. 129

CUDA devices take advantage of data-parallel program sections and accelerate 130 their execution. A CUDA program therefore consists of one or more phases that 131 are executed either in the host or in device. Sequential and control phases are 132 implemented in the host code, while phases which exhibit a large amount of data 133 parallelism are implemented in the device code. A CUDA program is also a unified 134 source code covering both sides. 135

We call kernels to those functions executed by the device (GPU). When they are 136 requested to be executed by the host, they allocate an execution grid on device. A 137 grid typically populates a large number of execution threads that work in SPMD 138 fashion: they execute the same piece of code (the kernel function) to probably 139 different portions of data. Actually, a kernel is written as a usual function in a 140 programming language (so far, only C++, Fortran or Python) but using special 141 keywords given by the CUDA API. These keywords, such as the thread identifier, 142 might take different values at different threads in run time, so that they can have 143 index different data elements, or even take different execution path (although this 144 is not optimal). CUDA threads are much lighter than CPU threads. A CUDA 145 programmer can assume that these threads take a few cycles to be generated and 146 scheduled. This contrasts with the threads of the CPU, which normally require 147 thousands of clock cycles to be managed.

Threads within a grid are arranged in a two-level hierarchy. At the higher level, 149 each grid consists of a two-dimensional array of *thread blocks*. At the lower level, 150 each block is organized as a three-dimensional array of threads. All blocks in a 151 grid have the same number and organization of threads. Moreover, each block is 152 identified by a two-dimensional identifier and each thread within its block by a threedimensional identifier. To date, a thread block can contain, at most, 1024 threads. 154 Threads within a block can easily cooperate through a special fast memory (see 155 below) and special warp-wide operations (see next section) and be synchronized 156 with a barrier operation. 157

One kernel is executed by just one grid that, has mentioned, arranges an array 158 of thread blocks, each with the same configuration of threads. It is also possible 159 to launch several kernels at the same time on a GPU with different grids. This is 160 allowed by the so-called CUDA streams. They are concepts similar to lanes, where 161 kernels get executed. As long as there are available resources on the GPU, the 162 kernel executions can be done simultaneously. CUDA streams also allow to overlap 163 execution with memory transfers, designing in this way full computing pipelines. 164

CUDA programmers also have to explicitly manage the memory layout and 165 hierarchy. GPUs offer different memory spaces arranged in a hierarchy where to 166 store the data of a parallel program. In CUDA, the host and the devices have separate 167 memory spaces (as it is in the real hardware). In order to run a kernel on the device, 168 the data has to be there. Hence, enough memory should be first allocated, and 169 later the relevant data has to be transferred from host to device. Similarly, after 170 the execution of the kernels on the device, the resulting data has to be copied from 171 device to host memory, and finally the allocated device memory should be released. 172 From this point of view, we can assume that CUDA uses static memory allocations. 173 Dynamic memory is already supported at the kernels, but there are restrictions, and 174 it drastically downgrades the performance.

The memory that serves as communication channel between the host and the 176 device is called global (or device) memory. The host can allocate memory and copy 177 memory, and CUDA threads can access it and make modifications. However, this 178 memory is the slowest in the GPU, but the largest one. The best performance is 179 achieved when contiguous threads (according to their identifiers) access contiguous 180 positions of data, in what is called coalesced memory access. This helps to maximize 181 the utilization of the memory bandwidth.

Threads can use a common memory space when they are in the same thread 183 block, which is called shared memory. Accesses to shared memory are very 184 fast when done in a coalesced way. However, it is a small space of up to just 185 kilobytes. On the other side, we can find cached memories. They are memories 186 that automatically speedup the access to repeated data through a cache. Examples 187 of them are constant and texture memories. They are read-only memories for the 188 GPU, and the CPU can just copy data in there, under certain restrictions. Moreover, 189 modern GPUs have two levels of cache memory for accesses to global memory, but 190 this is completely transparent to the code. 191

In summary, algorithms implemented in CUDA are structured as follows [17]: 192

- 1. The host initializes the program, reading the input data.
- The host allocate enough memory space in global memory for input and auxiliary 194 data.
- 3. The host copies the input and auxiliary data to the device.
- 4. The host launches a kernel to the device, with the following 197 syntax: kernelName <<< numBlocks, numThreads, streamId, 198 sharedMemory >>> (param1, param2,...) 199

193

- 5. When a kernel is executed, at the device side:
 - (a) The threads of each block read its corresponding data portion from global 201 memory to shared memory or to internal variables (also called registers, see 202 next section).
 - (b) Threads work with the data directly on the shared memory or with their 204 registers. 205
 - (c) Threads copy these data back to global memory.
- 6. The host might call more kernels, copy more data, retrieve data, etc.
- 7. When the algorithm is done, the host copies back the results from global memory 208 of the device. 209

As mentioned before, threads from different blocks cannot cooperate directly, 210 but only through the global memory and using a special set of atomic operations. 211 These operations are implemented by implicit locks, so that accesses to desired data 212 elements can be efficiently synchronized through them. However, this is restricted 213 to the use of a small set of operations. 214

6.2.3 GPU Architecture

A modern GPU architecture [19] consists of a processor array which has hundreds 216 (even thousands) of *SP* (streaming processor) cores organized in *SM* (streaming 217 multiprocessor). In this sense, every SM contains the following units: SP arithmetic 218 cores, SFU single-precision floating point units (for specific operations such as 219 sine, cosine, reciprocal square root, etc.), double precision units, instruction cache, 220 read only constant cache, read/write shared memory and L1 cache memory, a set 221 of 32-bit registers, and access to the off-chip memory (device/local memory). The 222 arithmetic units are capable to execute several instructions per clock cycle, and they 223 are fully pipelined, running at frequencies around 1 GHz (depending on the GPU). 224 The amount of cores, floating point units, shared memory, etc. depends on the GPU 225 itself.

SMs is able to manage and execute thousands of threads in hardware with zero 227 scheduling overhead. Each thread has its own thread execution state and can execute 228 an independent code path. This execution is done in a *SIMT (Single-Instruction 229 Multiple-Thread)* fashion [19], where threads execute the same instruction on 230 different piece of data. SMs create, manage, schedule, and execute threads in groups 231 of 32 threads (of the same thread block). This set of 32 threads is called *warp*. Each 232 SM can handle several warps. Individual threads of the same warp must be of the 233 same type and start together at the same program address in order to be scheduled 234 simultaneously, but they are free to branch (e.g., an if then else clause) and 235 execute independently at the cost of serialization.

When a grid is created to execute a kernel, the thread blocks are created and ²³⁷ assigned to SMs. An SM can handle several thread blocks, but a thread block is ²³⁸ assigned only to one SM. Then, the thread block is split into warps and they are ²³⁹ scheduled. When a warp is selected, its threads are executed on SPs as long as the ²⁴⁰

215

200

206

threads are synchronized in the same execution flow of the code. If the threads of a 241 warp diverge, the warp serially executes each taken branch path, disabling threads 242 that are not on that path. When all the paths complete, the threads re-converge to the 243 original execution path. 244

As shown by the CUDA programming model [17], the GPU contains several 245 memory spaces. First of all, both the GPU and the CPU have separated memory 246 spaces. They are connected through a bus that can be PCI Express × 16 bus standard. 247 Global memory is the largest (up to several gigabytes) but the slowest one. Although 248 there is a two-level cache memory system to speedup repeated access to same data, 249 an access to global memory is around 400 times slower than accessing on-chip 250 memory spaces (such as shared memory or registers). Moreover, in order to fulfill 251 the memory bandwidth, threads should make coalesced accesses. 252

As mentioned, inside each SM we can find a memory space called shared ²⁵³ memory. Its size is measured in kilobytes, but its access is very fast, even close ²⁵⁴ to the accesses to registers. This memory space is also split into shared memory ²⁵⁵ and L1 cache memory. The latter is transparent to CUDA programs, while shared ²⁵⁶ memory is manually managed (one can allocate space and let threads to copy and ²⁵⁷ modify data). There are also many other units for cached memory which is read only ²⁵⁸ for the cores. Finally, SMs incorporate a large amount of registers, whose access is ²⁵⁹ the fastest since they are next to the cores. They are used to allocate the values of ²⁶⁰ single variables declared in the code (e.g., iterators, auxiliary variables, etc.). ²⁶¹

6.2.4 Good Practices

CUDA is supported by a wide range of tools [33], including a compiler (called 263 nvcc), the driver for the GPU, libraries, and examples. They are freely available at 264 their website. There is also a vast amount of documentation, books, and literature 265 in this respect. CUDA is not only the most mature platform for GPU computing but 266 also the one with the largest community and support. It is important to know the 267 compiler options for automatic optimizations (like -O3) and to understand and use 268 the libraries (e.g., CuRAND for random number generation, CuSPARSE for sparse 269 matrix representations and operations, etc.).

It is also a good practice to start developing a reference program in sequential 271 C/C++ before starting implementing in CUDA. This is critical in order to first 272 understand the algorithm, secondly to validate the parallel version, and also to run 273 benchmarks and performance analysis. 274

Finally, let us introduce four ways to accelerate the execution of a program on a 275 GPU with CUDA [35]. They will help to understand the designs of GPU-based P 276 system simulators: 277

• *Emphasize parallelism*: GPUs prefer to run thousands of lightweight threads. 278 Thus, the algorithms should permit dividing the computation into many inde- 279 pendent pieces by decreasing the resources assigned to each thread and avoiding 280 synchronization. 281

- *Minimize branch divergence*: if a warp is broken because divergence in the path 282 executed by the threads, then there is no real parallelism. 283
- Maximize arithmetic intensity: computation is relatively cheap for today's GPUs, 284 but bandwidth is precious. It is better to maximize the computational operations 285 per memory transaction. Shared memory or registers can help for this purpose. 286
- *Exploit streaming bandwidth*: on the other side, GPUs and their on-board 287 memory have a peak bandwidth much faster than in CPUs. It is achieved 288 by streaming memory access patterns: coalesced access to aligned memory 289 positions. A good way to maximize the bandwidth in an algorithm is by the 290 scatter/gather strategy. 291

6.3 Generic Simulations

In this section, we will introduce a type of simulation of P systems, which is 293 called generic simulation [22, 25]. We will describe how to implement this kind 294 of simulators in CUDA and provide two illustrative examples. 295

6.3.1 Definition

When implementing a P system simulator, it is important to understand what type of ²⁹⁷ P systems we want to simulate before starting the development. We will say that a ²⁹⁸ generic simulator is a simulator developed for a wide range of P systems belonging ²⁹⁹ to the same variant. If the simulator is able to handle a large variety of P systems ³⁰⁰ (with very distinct rules and alphabets, even designed to solve different problems), ³⁰¹ then it is generic. Sometimes the types of P systems are restricted somehow for the ³⁰² sake of simplicity, but as long as the simulator accepts P systems from different ³⁰³ families (but for the same variant), we will say it is generic and not specific. ³⁰⁴

In this scenario, it is not possible to know what can happen in the computation 305 at a certain transition step. Therefore, it has to be prepared for any situation, so we 306 need to cover worst-case scenarios when developing such kind of simulators. For 307 example, we need to provide an upper bound of existing objects at a certain step, 308 in order to avoid memory overflows. Furthermore, in principle, all rules might be 309 selected for execution at a certain step (until their applicability is checked). The 310 rules must be stored in memory since we do not know them until the P system 311 model is parsed. As mentioned in Chap. 2, simulators are usually defined by three 312 and be designed to reproduce either a single computation or all computations of the 314 input P system. 315

292

The memory layout is also an essential part of a simulator, since P system ³¹⁶ simulators have been demonstrated to be memory and memory bandwidth bound. ³¹⁷ When storing the information of P system configuration, we can use either [25]: ³¹⁸

- *Sparse representation*: using a large array to store multiplicities, with a position ³¹⁹ per each possible object (all objects defined in the alphabet). The access is direct ³²⁰ since the object identifier is the index where to access the array. However, if many ³²¹ objects are not present at a certain moment, the array will be full of zeroes. ³²²
- *Dense representation*: using a double array with a component for the object 323 identifier and the other for multiplicity. We need to search for the object, unless 324 we track them and we know exactly where they are store at any moment. This 325 can help to drastically reduce the size because objects with multiplicity zero can 326 be discarded. 327

Generic simulators usually use sparse representations, since, in this way, they 328 can identify objects very efficiently, in O(1). The object identifier is employed as 329 the index to access the array where storing the multisets, and by representing the 330 whole alphabet, we make sure of an upper bound for the worst-case scenario. 331

6.3.2 Simulating P Systems with Active Membranes

In this section we will depict the very first GPU simulator for P systems ever 333 developed. It was a generic simulator and helped to understand how to better map 334 the parallelism of P systems on the parallelism on GPUs. The simulated models 335 were of the variant P systems with active membranes and elementary division. 336

The original work is published in [4, 21, 22]. The full framework of simulators ³³⁷ for P systems with active membranes, including the sequential, fast sequential, and ³³⁸ CUDA parallel simulators, is called *PCUDA*. It is a subproject of the PMCGPU ³³⁹ project and can be downloaded from the official website http://sourceforge.net/p/ ³⁴⁰ pmcgpu [45] or the repository https://github.com/RGNC/pcuda. ³⁴¹

6.3.2.1 Recognizer P Systems with Active Membranes

Families of cell-like P systems whose membrane structures does not grow, that ³⁴³ is, there is no rules producing new membranes in the system, only can solve in ³⁴⁴ polynomial time and uniform way, problems in class **P**. Therefore, new ingredients ³⁴⁵ are needed in order to be able to provide efficient solutions of computationally ³⁴⁶ hard problems by making use of an exponential workspace, expressed in terms of ³⁴⁷ number of membranes and number of objects, created in linear time. In [36], a new ³⁴⁸ computing model, called *P system with active membranes*, was introduced. In these ³⁴⁹ systems, the membranes have associated with electrical charges and make use of ³⁵⁰ *division rules*, inspired from the *mitosis* and *meiosis* processes, as a mechanism to ³⁵¹ generate in linear time, an exponential workspace. Polynomial time and uniform ³⁵² solutions to **NP**-complete problems were given by using families of the new ³⁵³ computing model. ³⁵⁴

332

Next, P systems with active membranes and division rules only for elementary 355 division are formally defined. 356

Definition 6.1. A P system with active membranes of degree q > 1 is a tuple 357 $\Pi = (\Gamma, H, \mu, \mathcal{M}_1, \dots, \mathcal{M}_q, \mathcal{R}, i_{out})$, where: 358

- 1. Γ and *H* are finite alphabets such that $\Gamma \cap H = \emptyset$;
- 2. μ is a rooted tree with q nodes labeled by elements from H (the root is labeled 360 by $i_s \in H$; 361
- 3. $\mathcal{M}_1, \ldots, \mathcal{M}_q$ belongs to $\mathcal{M}(\Gamma)$, that is, all of them are multisets over Γ ;
- 4. \mathcal{R} is a finite set of rules, of the following forms:
 - (a) $[a \to u]_{h}^{\alpha}$, for $h \in H, \alpha \in \{+, -, 0\}, a \in \Gamma, u \in M(\Gamma)$ (object evolution 364 rules); 365
 - (b) $a \begin{bmatrix} 1 \\ 1 \\ h \end{bmatrix}_{h}^{\alpha_{1}} \rightarrow \begin{bmatrix} b \end{bmatrix}_{h}^{\alpha_{2}}$, for $h \in H \setminus \{i_{s}\}, \alpha_{1}, \alpha_{2} \in \{+, -, 0\}, a, b \in \Gamma$ (send-in 366 rules) 367
 - 368
 - (c) $[a]_{h}^{\alpha_{1}} \rightarrow b[]_{h}^{\alpha_{2}}$, for $h \in H, \alpha_{1}, \alpha_{2} \in \{+, -, 0\}, a, b \in \Gamma$ (send-out rules); (d) $[a]_{h}^{\alpha} \rightarrow b$, for $h \in H \setminus \{i_{s}, i_{out}\}, \alpha \in \{+, -, 0\}, a, b \in \Gamma$ (dissolution rules); 369
 - (e) $[a]_{h}^{\alpha_{1}} \rightarrow [b]_{h}^{\alpha_{2}} [c]_{h}^{\alpha_{3}}, \text{ for } h \in H \setminus \{i_{s}, i_{out}\}, \alpha_{1}, \alpha_{2}, \alpha_{3} \in \{+, -, 0\}, a, b, c \in [a]_{h}^{\alpha_{1}} \rightarrow [b]_{h}^{\alpha_{2}} [c]_{h}^{\alpha_{3}}, \text{ for } h \in H \setminus \{i_{s}, i_{out}\}, \alpha_{1}, \alpha_{2}, \alpha_{3} \in \{+, -, 0\}, a, b, c \in [a]_{h}^{\alpha_{1}} \rightarrow [b]_{h}^{\alpha_{2}} [c]_{h}^{\alpha_{3}}, \text{ for } h \in H \setminus \{i_{s}, i_{out}\}, \alpha_{1}, \alpha_{2}, \alpha_{3} \in \{+, -, 0\}, a, b, c \in [a]_{h}^{\alpha_{1}} \rightarrow [b]_{h}^{\alpha_{2}} [c]_{h}^{\alpha_{3}}, \text{ for } h \in H \setminus \{i_{s}, i_{out}\}, \alpha_{1}, \alpha_{2}, \alpha_{3} \in \{+, -, 0\}, a, b, c \in [a]_{h}^{\alpha_{1}} \rightarrow [b]_{h}^{\alpha_{2}} [c]_{h}^{\alpha_{3}} [c$ 370 Γ (division rules for elementary membranes); 371
- 5. $i_{out} \in H \cup \{env\}$, where $env \notin \Gamma \cup H$.

A P system with active membranes of degree $q \ge 1, \Pi = (\Gamma, H, \mu, \mathcal{M}_1, \dots, \mathcal{M}_q, 373)$ \mathcal{R} , i_{out}), can be viewed as a set of q membranes, injectively labeled by elements 374 of H, arranged in a hierarchical structure μ given by a labeled rooted tree (called 375 membrane structure) whose root is called the *skin membrane* (labeled by i_s), such 376 that (a) each membrane has associated with an electrical charge from the set 377 $\{+, -, 0\}$; (b) $\mathcal{M}_1, \ldots, \mathcal{M}_p$ represent the finite multisets of *objects* (symbols of 378) the working alphabet Γ) initially placed in the q membranes of the system; (c) \mathcal{R}_{379} is a finite set of rules over Γ associated with the labels; and (d) $i_{out} \in H \cup \{env\}$ 380 indicates the output zone. We use the term zone i to refer to membrane i in the case 381 $i \in H$ and to refer to the "environment" of the system in the case i = env. The 382 leaves of μ are called *elementary membranes*. 383

Next, the semantics of the new computing model is described. A *configuration* 384 (or instantaneous description) C_t at an instant t of a P system with active membranes 385 is described by the following elements: (a) the membrane structure at instant t and $_{386}$ (b) all multisets of objects over Γ associated with all the membranes present in the 387 system at that moment. 388

An object evolution rule $[a \rightarrow u]_h^{\alpha}$ is applicable to a configuration C_t at an 389 instant t, if there exists a membrane labeled by h with electrical charge α , in C_t , 390 such that contains object a. When applying such a rule to such a membrane, one 391 object a is consumed and objects from the multiset u is produced in that membrane. 392

 C_t at an instant t, if there exists a membrane labeled by h with electrical charge α_1 , 394 in C_t such that h is not the label of the root of μ and its parent membrane contains 395 object a. When applying such a rule to such a membrane, one object a is consumed $_{396}$

359

362

363

from the parent membrane, and object *b* is produced in the corresponding membrane ³⁹⁷ labeled by *h*. Besides, the charge α_1 of that membrane *h* is replaced by α_2 . ³⁹⁸

A send-out communication rule $[a]_{h}^{\alpha_{1}} \rightarrow b []_{h}^{\alpha_{2}}$ is applicable to a configuration 399 C_{t} at an instant t, if there exists a membrane labeled by h with electrical charge α_{1} , 400 in C_{t} such that it contains object a. When applying such a rule to such a membrane, 401 one object a is consumed from such membrane h, and object b is produced in the 402 parent of such membrane (in the case that such membrane is the skin, then object 403 b is produced in the environment). Besides, the charge α_{1} of that membrane h is 404 replaced by α_{2} .

A dissolution rule $[a]_{h}^{\alpha} \rightarrow b$ is applicable to a configuration C_{t} at an instant 406 t, if there exists a membrane labeled by h with electrical charge α , in C_{t} , different 407 from the skin membrane and the output zone, such that it contains object a. When 408 applying such a rule to such a membrane, one object a is consumed, membrane h is 409 dissolved, and one object b and the remaining objects of the membrane where the 410 rule is applied are sent to its parent (or the first ancestor that has not been dissolved). 411

A division rule for elementary membrane $[a]_{h}^{\alpha_{1}} \rightarrow [b]_{h}^{\alpha_{2}}$ $[c]_{h}^{\alpha_{3}}$ is applicable to 412 a configuration C_{t} at an instant t, if there exists an elementary membrane labeled by 413 h with electrical charge α_{1} , in C_{t} , different from the skin membrane and the output 414 zone, such that it contains object a. When applying such a rule to such a membrane, 415 the membrane with label h is divided into two membranes with the same label; in 416 the first copy, one object a is replaced by one object b; in the second one, one object 417 a is replaced by one object c; all the other objects are replicated and copies of them 418 are placed in the two new membranes. Besides, the charge α_{1} of the first created 419 membrane h is replaced by α_{2} , and the charge α_{1} of the second created membrane 420 h is replaced by α_{3} .

In P systems with active membranes, the rules are applied according to the 422 following principles: 423

- At one transition step: (i) one object and one membrane can be used by only one 424 rule, selected in a non-deterministic way, and (ii) at most a rule of types (b)–(e), 425 selected in a non-deterministic way, can be applied to a membrane, and then it is 426 applied once. 427
- Object evolution rules can be simultaneously applied to a membrane with one 428 rule of types (b)–(e). If it is the case, object evolution rules will be applied in a 429 maximally parallel manner. 430
- If an object evolution rule and a division rule are applied to a membrane at the 431 same transition step, then we suppose that first the evolution rule is applied, and 432 then the division is produced. Of course, this process takes only one transition 433 step. 434
- The skin membrane and the output membrane, if any, can never get divided nor 435 dissolved.
 436

Given a P system with active membranes, $\Pi = (\Gamma, H, \mu, M_1, \dots, M_q, \mathcal{R}, i_{out})$, 437 the *initial configuration* of Π is $C_0 = (\mathcal{M}_1, \dots, \mathcal{M}_q)$. A configuration is a *halting* 438 *configuration* if no rule of the system is applicable to it. We say that configuration 439

443

 C_1 yields configuration C_2 in one *transition step*, denoted $C_1 \Rightarrow_{\Pi} C_2$, if we can 440 pass from C_1 to C_2 by applying the rules from \mathcal{R} following the previous remarks. A 441 *computation* of Π is a (finite or infinite) sequence of configurations such that:

- 1. The first term of the sequence is the initial configuration of the system.
- 2. Each non-initial configuration of the sequence is obtained from the previous 444 configuration in one transition step 445
- 3. If the sequence is finite (called *halting computation*), then the last term of the 446 sequence is a halting configuration. 447

All computations start from an initial configuration and proceed as stated above; 448 only halting computations give a result, which is encoded by the objects present in 449 the output zone i_{out} in the halting configuration. 450

Let us notice that these P systems have some important features: (a) They use 451 three electrical charges; (b) the polarization of a membrane can be modified by 452 the application of a rule; (c) the label of a membrane cannot be modified by the 453 application of a rule; and (d) they do not use cooperation neither priorities. 454

Decision problems are associated with languages in such manner that *solving* a 455 decision problem is defined by *recognizing* the language associated with it. For that, 456 *recognizer membrane systems* were introduced in [41] (called *decision P systems*), 457 and complexity classes associated with these systems were introduced in [40]. 458 Over the last few years, the previous methodology for addressing the **P** *versus* **NP** 459 problem has been applied in the framework of *Membrane Computing*. 460

A computing model in the paradigm of Membrane Computing (generically called 461 *membrane system*) is said to be a *recognizer system* if it has the following syntactic 462 and semantic peculiarities: (a) the working alphabet has two distinguished objects 463 (yes and no); (b) there exist an input alphabet strictly contained in the working 464 alphabet and an *input membrane*; (c) the initial content of each compartment is a 465 multiset of objects from the working alphabet not belonging to the input alphabet; 466 (d) all computations of the system are halt; and (e) for each computation, either 467 object yes or object no (but not both) must have been released to the environment 468 and only at its last step. Recognizer membrane systems have the ability to accept or 469 reject multisets over the input alphabet. Specifically, given a recognizer membrane 470 system Π , for each multiset *m* over the input alphabet, a new initial configuration 471 is obtained by adding the multiset m to the content of the input compartment at 472 the initial configuration of Π (the system Π with this new initial configuration 473 associated with m is denoted by $\Pi + m$). Then, we say that system Π accepts 474 (respectively, reject) the input multiset m if and only if all computations of the 475 system $\Pi + m$ answer yes (resp. no). 476

Unlike a Turing machine where there is an infinite tape, all the elements that 477 make up a recognizer membrane system have a finite description. Therefore, while 478 a decision problem (with an infinite set of instances) can be solved by a single Turing 479 machine, an infinite family of recognizer membrane systems is necessary to solve it. 480

Following [40], we say that a family $\Pi = \{\Pi(n) \mid n \in \mathbb{N}\}$ of recognizer 481 membrane systems solves a decision problem X in *polynomial time and uniform* 482

way if the following holds: (i) the family Π can be generated by a deterministic 483 Turing machine working in polynomial time; and (ii) there exists a pair (*cod*, *s*) of 484 polynomial-time computable functions (over the set of instances of *X*) such that (a) 485 for each instance $u \in I_X$, s(u) is a natural number and cod(u) is an input multiset 486 of the system $\Pi(s(u))$; (b) for each $n \in \mathbb{N}$, the set $s^{-1}(\{n\})$ is a finite set; and (c) the 487 family Π is *polynomially bounded*, *sound*, and *complete* with regard to (*X*, *cod*, *s*) 488 (see [40] for details).

Given a computing model \mathcal{R} of recognizer membrane systems, $\mathbf{PMC}_{\mathcal{R}}$ denotes 490 the set of decision problems solvable by families from \mathcal{R} in polynomial time 491 and uniform way. This complexity class is closed under complement and under 492 polynomial-time reduction [40]. Hence, if X is a complete problem for a complexity 493 class \mathcal{K} and $X \in \mathbf{PMC}_{\mathcal{R}}$, then we deduce that $\mathcal{K} \cup \mathbf{co} - \mathcal{K} \subseteq \mathbf{PMC}_{\mathcal{R}}$.

P systems with active membranes (without dissolution and using division rules 495 only for elementary membranes) have been successfully used to design polynomial 496 time solutions to (weak and strong) **NP**-complete problems (e.g., SAT [42], 497 Subset Sum [38], Knapsack [39], Partition [14], etc.). It is important to 498 note that some of these solutions only make use of two polarizations in their design. 499

6.3.2.2 Simulation Algorithm

The simulator is based on the sequential simulator for P systems with active 501 membranes provided in pLinguaCore [13]. In this design, the simulation process 502 is a loop divided into two stages: *selection stage* and *execution stage*. The selection 503 stage consists in the search for rules to be executed in each membrane of a given 504 configuration. The selected rules are executed at the execution stage, what finalizes 505 the simulation of a computation step (or transition). 506

The input data for the selection stage contains the description of the membranes 507 with their multisets (strings over the working alphabet of objects, labels associated 508 with the membrane, etc.) and the set of defined rules. The output data of this stage 509 are the multisets of selected rules. Only the execution stage changes the information 510 of the configuration. It is the reason why execution stage needs synchronization 511 when accessing to the membrane structure and the multisets. 512

At the end of the execution stage, the simulation process restarts the selection 513 stage in an iterative way until a halting configuration is reached. This stop condition 514 is twofold: a certain number of iterations or a final configuration is reached. On one 515 hand, we define a maximum number of iterations at the beginning of the simulation. 516 On the other hand, a halting configuration is obtained when there are no more rules 517 to select at selection stage. As previously explained, the halting configuration is 518 always reached since it is a simulator for recognizer P systems. 519

Non-determinism affects the selection stage, since it is possible to have more than 520 one selectable rule but only one can be executed. For example, two evolution rules 521 can be executed using the same object, a division rule and a send-in rule that can be 522 selected in the same membrane at the same time. In order to avoid non-determinism 523 somehow, the simulator assumes only confluent P systems. Thus, instead of working 524 with the entire tree of possible computations, the simulator selects and simulates 525 only one computation path, since all paths are guaranteed to give the same answer. 526

553

We can take advantage of this property by selecting path using the lowest cost rules. 527 We will measure this cost in number of membranes and synchronization operations. 528 These are the conditions that could damage the simulation performance the most. 529 In this context, we introduce the following priorities among rules in the selection 530 stage: 531

- 1. *Dissolution rules*: they decrease the number of membranes (highest priority). 532
- 2. *Evolution rules*: they do not need any communication among membranes (which 533 avoids synchronization). 534
- 3. *Send-out rules*: they do need communication between the given membrane and 535 its parent (adding one object to its parent). 536
- 4. *Send-in rules*: they do need communication between the given membrane and its ⁵³⁷ parent (reserving one object from its parent and adding the object to itself). ⁵³⁸
- 5. *Division rules*: they increase the number of membranes (lowest priority). 539

During the execution stage, the information of the system can vary by including 540 new objects inside membranes, dissolving membranes, dividing membranes, etc., 541 obtaining a new configuration. This new configuration will be the input data for the 542 selection stage of the next iteration. 543

Finally, note that this two-staged algorithm allows to keep a coherence in the 544 simulation. If we perform selection and execution of rules, one by one, it would be 545 difficult to ensure the semantic constraints of the system. Moreover, the selected and 546 executed rules in a step of the simulator may not correspond to the rules applied in 547 a computing step of the theoretical model. An alternative solution might be to take 548 two copies of the configuration, one to be updated with the right-hand sides of the 549 rules and another to select rules (subtracting the left-hand sides of rules). As this 550 involves a bigger use of memory, the simulator uses the two stages and a temporary 551 data structure to store information of the selection of rules.

6.3.2.3 Sequential Simulator

As previously mentioned, CUDA programming model [33] is based on the C/C++ 554 language [16]. Therefore, the first recommended step when developing applications 555 in CUDA is to start from a baseline algorithm written in C++, identifying the parts 556 that can be susceptible to be parallelized on the GPU. In this work, we have based 557 on the simulator for P systems with active membranes developed in pLinguaCore 558 [13]. This sequential (or single-threaded) simulator is programmed in JAVA, so the 559 first step was to translate the code to C++. 560

The first version of the sequential simulator implements the structure of membranes by using C++ pointers and dynamic memory allocations. Each membrane stores a pointer to its parent, a pointer to the first of its children, another pointer to one of its brothers (having the same parent membrane), the charge, and the multiset of objects. The multiset of objects is also implemented by a (dynamic) linked list based on pointers. Each object in the multiset stores its multiplicity (if zero, it is deleted to save memory space) and a pointer to the next object. Therefore, memory spaces for membranes and objects are created and deleted "on demand." The rules of 568 the system are statically stored, so that we can easily access to the rules associated to each membrane, by using its label and charge. Furthermore, the multiset of selected rules is also implemented using a dynamic linked list. However, we found that this drastically slowdown the simulation, since objects get created and consumed at every step, and hence, we are continually allocating and destroying memory, what is very time consuming.

Therefore, a simulator using static structures that get allocated at the beginning 575 of the simulator was developed and shown to be 160 times faster than the first 576 version [21]. These structures are the same also for the parallel simulator, so they 577 are replicated at both sides in order to achieve fair comparisons. In summary, the 578 memory layout to represent the P system is based on the following data structures: 579

- *Multisets*: an array storing the multisets of the objects using a sparse representation. Since, for simplicity, it is assumed that the simulated P system can contain only two levels in the memory hierarchy (a skin membrane and elementary membranes), the representation of the environment, skin, and elementary membranes are separated. The amount of elementary membranes is set initially by the user.
- Charges: an array storing the charge of each membrane.
- *Rule sets*: an array storing rules information. It is indexed by using a membrane 586 label, a charge, and finally an object index. Given that it is possible to have more 587 than one rule associated to the same object, and assuming that the P system is 588 confluent, only one rule of each type is stored. 589

One major problem to overcome is the competition for objects between different 590 membranes. In this case, internal membranes applying send-in rules are competing 591 for the objects in the parent. We loop the tree from the top to the bottom, so the top 592 level membranes have more priority using its objects than internal membranes using 593 send-in rules. 594

The input of the simulator (the P system with active membranes to simulate) is 595 given by a binary file. It is a file whose information is encoded in Bytes and bits (not 596 understandable by humans like plain text), which is suitable for compressing data. 597 This binary file contains all the information of the P system (alphabet, labels, rules, 598 etc.) which is the input of the simulator. The format is depicted in [21]. pLinguaCore 599 2.0 [13] is able to translate a P system written in P-Lingua language into a binary file. 600 First, we define the P system into P-Lingua. pLinguaCore translates it to a binary file, which is used as the input of the simulator. The output is a plain text generated 602 with a format similar to the one provided in pLinguaCore. 603

6.3.2.4 Parallel Simulation on CUDA

Whenever we design algorithms in the CUDA programming model, the main effort 605 is dividing the required work into processing pieces, which have to be processed 606 by *TB* thread blocks of *T* threads each. Using a thread block size of T = 256, it is 607 empirically determined to obtain the overall best performance on the Tesla C1060 608 [43]. Each thread block accesses to one different set of input data and assigns a 609 single or small constant number of input elements to each thread. 610

585



Fig. 6.1 Basic design of the parallel simulator on the GPU. From [4,21]

Each thread block can be considered independent to the other, and it is at this 611 level at which internal communication (among threads) is cheap using explicit 612 barriers to synchronize, and external communication (among blocks) becomes 613 expensive, since global synchronization can only be achieved by the barrier implicit 614 between successive kernel calls. The need of global synchronization in the designs 615 requires successive kernel calls even to the same kernel. 616

Figure 6.1 shows the overall design of the simulator on the GPU [4]. Thread 617 blocks and threads are distributed as follows. Each membrane of the simulated 618 P system is attributed to each thread block. In this way, the parallelism between 619 membranes by using the parallelism between thread blocks is identified. However, 620 this is tricky. Membranes can communicate accordingly to the hierarchical tree 621 structure, while thread blocks are all independent. Communication through send-out 622 and dissolution rules (down-up direction) is controlled by globally synchronizing 623 the selection and execution stages. This is implemented by using different kernels. 624 However, send-in rules (up-down direction in the tree) are more complicated to 625 control. In this case, different membranes can compete for single objects. The 626 sequential simulator controls this issue by looping the tree from the top to the 627 bottom. However, the parallel simulator has to run all the membranes in parallel. 628 Therefore, for the sake of simplicity, the parallel simulator can handle only two 629 levels of membrane hierarchy: the skin (controlled by the host) and the rest of 630 elementary membranes (controlled by the thread blocks in device). This is the 631 tree structure we can find in the literature for the majority of solutions based on 632 P systems with active membranes (note that division rules enlarge the tree width-633 wise) [40]. 634

Furthermore, each individual thread is assigned to each object within a membrane 635 (corresponding to its thread block). It is responsible for identifying the rules that can 636 be executed using the corresponding object, that is, rules that have that object in their 637 left-hand sides. Since all blocks must have the same number of threads, and each 638 membrane can contain a different multiset of objects in every time step, we identify 639 as common for all membranes the whole alphabet. Note that threads can work with 640 many objects that do not really exist in the membrane, as all the alphabet of objects 641 is usually not present within a membrane at a given instant. In fact, the simulator 642 assigns multiple objects to the same thread for not restricting the number of objects 643 in the alphabet. However, the number of objects in the alphabet must be divisible by 644 a number smaller than 512 (the maximum number of threads, per thread block), in 645 order to equally distribute the objects among the threads.

The simulator contains five kernels to implement the selection and execution 647 stages [21]. The first kernel implements the selection stage and also the execution 648 stage for evolution rules. The other four kernels implement the other execution rules 649 (dissolution, division, send-out, and send-in rules). All the kernels follow this basic 650 design. The selection kernel starts with the selection stage. After the selection stage, 651 we also execute in this kernel the evolution rules. These rules are executed inside this 652 kernel for three main reasons: the evolution rules do not imply communication (and 653 therefore, synchronization) among membranes; they are executed in a maximal way, 654 and this decision allows us to use less global memory because it is not necessary to 655 store the selected evolution rules for the execution stage. The rest of the rules to 656 be applied are executed in four different kernels, one kernel per each kind of rule 657 (dissolution, division, send-out, and send-in). 658

Algorithm 2 shows the pseudo-code of the simulator. First of all, the data needed 659 for the computation is moved to the GPU. Then, the code calls the selection kernel 660 which returns the selected rules for the current configuration of the P system. Among 661 the possible selected rules, there will be different kinds of rules to be executed. 662 Therefore, the type of those rules is identified for launching only the required kernels 663 to accomplish the execution stage. As explained before, this process iterates until the 664 maximum number of steps is reached or the system returns an answer. Finally, the 665 result data is copied back to the CPU. 666

6.3.2.5 Performance Comparative Analysis

In this section, the performance of the developed simulators is compared. This 668 is done by a very simple example, with the aim of studying the behavior of the 669 CUDA kernels. In order to evaluate the performance of the simulator, a family of 670 P systems was designed, named test P system, where it is easy to vary the number 671 of membranes as well as the number of objects [4]. This test P system also fits the 672 behavior of the GPU since only evolution and division rules are defined (without 673 communication and dissolution rules), and every object in every membrane will 674

| Algorithm 2: Parallel simulator of P systems on the GPU, from [21] | | |
|--|--|--|
| 1: configuration \leftarrow initialConfiguration | | |
| 2: selectedRules $\leftarrow \emptyset$ | | |
| 3: step $\leftarrow 0$ | | |
| 4: is Final Configuration \leftarrow false | | |
| 5: CopyDataFromCPUtoGPU(configuration) | | |
| 6: CopyDataFromCPUtoGPU(rules) | | |
| 7: while step $<$ maxStep \land NOT isFinalConfiguration do | | |
| 8: kernelSelection(rules,configuration,selectedRules) | | |
| 9: if DISSOLUTION \in selectedRules then | | |
| 10: kernelDissolution(rules,configuration,selectedRules) | | |
| 11: end if | | |
| 12: if DIVISION \in selectedRules then | | |
| 13: kernelDivision(rules,configuration,selectedRules) | | |
| 14: end if | | |
| 15: if SEND-OUT \in selectedRules then | | |
| 16: kernelSendOut(rules,configuration,selectedRules) | | |
| 17: end if | | |
| 18: if SEND-IN \in selectedRules then | | |
| 19: kernelSendIn(rules,configuration,selectedRules) | | |
| 20: end if | | |
| 21: step \leftarrow step + 1 | | |
| 22: isFinalConfiguration ← checkFinalConfiguration(configuration) | | |
| 23: end while | | |
| 24: CopyDataFromGPUtoCPU(configuration) | | |

evolve according to a given rule. The defined P system is of the following form ${}^{675}\Pi = (O, H, \mu, \omega_1, \omega_2, R)$, where: 676

| • | $O = \{d, o_i / 0 \le i \le n\},$ | 677 |
|---|---------------------------------------|-----|
| • | H = 1, 2, | 678 |
| • | $\mu = [[]_2]_2,$ | 679 |
| • | $\omega_1 = \emptyset, \omega_2 = O,$ | 680 |
| • | R = | 681 |
| | | |

| (i) Evolution rules: $[o_i \rightarrow o_i]_2^0$, $0 \le i \le n$ | 682 |
|--|-----|
| (ii) Division rule: $[d]_2^0 \rightarrow [d]_2^0 [d]_2^0$ | 683 |

Thus, the test P system allows us to take control of the number of objects in the 684 system by modifying the *n* parameter. Furthermore, the number of rules changes 685 along with the number of objects, and the number of membranes in every step of 686 the computation is equal to 2^s , where *s* is the step number. Lastly, the number of 687 evolution rules selected and executed per membrane in every step is invariable, since 688 they are defined one per object and all the objects of the alphabet are presented in 689 every membrane labeled with 2.



Fig. 6.2 Comparing the execution time for one step of the fast sequential and parallel simulators, by increasing the number of membranes in the system and using a total of 2560 objects in the alphabet. From [4,21]



Fig. 6.3 Comparing the speedup for one step of the fast sequential and parallel simulators, by increasing the number of membranes in the system and using a total of 2560 objects in the alphabet. From [4, 21]

Figures 6.2 and 6.3 show the results obtained for the parallel simulator versus 691 the sequential version. Notice that in both graphs the Y-axis is also represented 692 in a logarithmic form. The benchmark covers the parallelism between membranes 693 by exponentially increasing the number of membranes. It can be seen that the 694 CPU simulator also increases its time exponentially from the beginning (with four 695 membranes) until reaching the final configuration (with 32768 membranes). The 696 CUDA simulator, which assigns 256 threads per block (each thread handles 10 697 elements per membrane), also increases its execution time in a near exponential

way, but the performance difference is about $5.7 \times$, and this difference enlarges with 698 the number of membranes (from 1024), because the resources of the GPU are fully 699 utilized. 700

6.3.3 Simulating Population Dynamics P Systems

In this section, a simulator for Population Dynamics P (PDP) systems is revisited. 702 It is a generic simulator implementing the DCBA algorithm for PDP systems. 703

The original work is published in [21, 22, 24, 28]. The framework of generic 704 simulators for PDP systems on GPU is called *ABCDGPU*. It is a subproject of 705 the PMCGPU project and can be downloaded from the official website http:// 706 sourceforge.net/p/pmcgpu [45] or the repository https://github.com/RGNC/abcd-707 gpu.

6.3.3.1 Population Dynamics P Systems

Population Dynamics P systems are a variant of *multienvironment P systems with* 710 *active membranes* [6–8]. The model consists of a directed graph of environments, 711 each of them containing a P system where electrical charges are associated with 712 membranes. All P systems share the same *skeleton*, in the sense that they have the 713 same working alphabet, the same membrane structure, and the same set of rules. 714 Nevertheless, in this framework each rule has associated a probability function 715 which can vary for each environment. 716

Definition 6.2. A Population Dynamics P system (PDP) of degree (q, m), $q, m \ge 717$ 1, taking $T \ge 1$ time units, is a tuple 718

$$\Pi = (G, \Gamma, \Sigma, T, \mathcal{R}_E, \mu, \mathcal{R}, \{f_{r,j} : r \in \mathcal{R}, 1 \le j \le m\}, \{\mathcal{M}_{i,j} : 1 \le i \le q, 1 \le j \le m\})$$
719

where:

- G = (V, S) is a directed graph. Let $V = \{e_1, \ldots, e_m\}$.
- Γ and Σ are alphabets such that $\Sigma \subsetneq \Gamma$.
- *T* is a natural number.
- \mathcal{R}_E is a finite set of rules of the form $(x)_{e_j} \overrightarrow{p_r} (y_1)_{e_{j_1}} \cdots (y_h)_{e_{j_h}}$, where 724 $x, y_1, \ldots, y_h \in \Sigma$, $(e_j, e_{j_l}) \in S$, $1 \leq l \leq h$, and $p_r : \{1, \ldots, T\} \longrightarrow [0, 1]$ 725 is a computable function such that for each $e_j \in V$ and $x \in \Sigma$, the sum of 726 functions associated with the rules of the type $(x)_{e_j} \overrightarrow{p_r} (y_1)_{e_{j_1}} \cdots (y_h)_{e_{j_h}}$ is the 727 constant function 1.
- μ is a rooted tree labeled by $1 \le i \le q$, and by symbols from the set EC = 729{0, +, -}.
- \mathcal{R} is a finite set of rules of the form $u[v]_i^{\alpha} \to u'[v']_i^{\alpha'}$, where $u, v, u', v' \in 731$ $M_f(\Gamma), u + v \neq \emptyset, 1 \leq i \leq q \text{ and } \alpha, \alpha' \in \{0, +, -\}$, such that there is no 732 rules $(x)_{e_j} \stackrel{\rightarrow}{p_r} (y_1)_{e_{j_1}} \cdots (y_h)_{e_{j_h}}$ and $u[v]_i^{\alpha} \to u'[v']_i^{\alpha'}$ having $x \in u$. 733

701

709

720

721

722

- For each $r \in \mathcal{R}$ and $1 \leq j \leq m$, $f_{r,j} : \{1, \ldots, T\} \longrightarrow [0, 1]$ is a computable 734 function such that for each $u, v \in M_f(\Gamma), 1 \leq i \leq q, \alpha, \alpha' \in \{0, +, -\}$ and 735 $1 \leq j \leq m$, the sum of functions $f_{r,j}$ with $r \equiv u[v]_i^{\alpha} \rightarrow u'[v']_i^{\alpha'}$, is the constant 736 function 1.
- For each $i, j \ (1 \le i \le q, 1 \le j \le m), \mathcal{M}_{i,j}$ is a finite multiset over Γ . 738

A Population Dynamics P system defined as above can be viewed as a set of m 739 environments e_1, \ldots, e_m interlinked by the edges from the directed graph G. Each 740 environment e_j only can contain symbols from the alphabet Σ , and all of them also 741 contain a P system skeleton, $\Pi_j = (\Gamma, \mu, \mathcal{M}_{1,j}, \ldots, \mathcal{M}_{q,j}, \mathcal{R})$, of degree q, where: 742

- (a) Γ is the working alphabet whose elements are called objects.
- (b) μ is a rooted tree which describes a membrane structure consisting of q 744 membranes injectively labeled by 1, ..., q. The skin membrane (the root of the 745 tree) is labeled by 1. We also associate electrical charges from the set {0, +, -} 746 with membranes.
- (c) $\mathcal{M}_{1,j}, \ldots, \mathcal{M}_{q,j}$ are finite multisets over Γ , describing the objects initially 748 placed in the *q* regions of μ , within the environment e_j . 749
- (d) \mathcal{R} is the set of evolution rules of each P system. Every rule $r \in \mathcal{R}$ in Π_j has a 750 computable function $f_{r,j}$ associated with it. For each environment e_j , we denote 751 by \mathcal{R}_{Π_j} the set of rules with probabilities obtained by coupling each $r \in \mathcal{R}$ with 752 the corresponding function $f_{r,j}$.

Therefore, there is a set \mathcal{R}_E of communication rules between environments, and 754 the natural number *T* represents the simulation time of the system. The set of rules 755 of the whole system is $\bigcup_{i=1}^{m} \mathcal{R}_{\prod_i} \cup \mathcal{R}_E$. 756

The *semantics* of Population Dynamics P systems is defined through a nondeterministic and synchronous model (in the sense that a global clock is assumed). 758 Next, we describe some semantics aspects of these systems. 759

An evolution rule $r \in \mathcal{R}$, of the form $u[v]_i^{\alpha} \to u'[v']_i^{\alpha'}$, is applicable to 760 each membrane labeled by *i*, whose electrical charge is α , and it contains the 761 multiset *v*, and its parent contains the multiset *u*. When such rule is applied, the 762 objects of the multisets *v* and *u* are removed from membrane *i* and from its parent 763 membrane, respectively. Simultaneously, the objects of the multiset *u'* are added to 764 the parent membrane *i*, and objects of multiset *v'* are introduced in membrane *i*. 765 The application also replaces the charge of membrane *i* to α' . In each environment 766 e_j , the rule *r* has associated a probability function $f_{r,j}$ that provides an index of the 767 applicability when several rules compete for objects. In this model, the cooperation 768 degree is given by |u| + |v|.

A rule $r \in \mathcal{R}_E$, of the form $(x)_{e_j} \overrightarrow{p_r} (y_1)_{e_{j_1}} \dots (y_h)_{e_{j_h}}$, is applicable to the 770 environment e_j if it contains object x. When such rule is applied, object x passes 771 from e_j to e_{j_1}, \dots, e_{j_h} possibly modified into objects y_1, \dots, y_h respectively. At 772 any moment t $(1 \le t \le T)$ for each object x in environment e_j , if there exist 773 communication rules of the type $(x)_{e_j} \overrightarrow{p_r} (y_1)_{e_{j_1}} \dots (y_h)_{e_{j_h}}$, then one of these rules 774

will be applied. If more than one such a rule can be applied to an object at a given 775 instant, the system selects one randomly, according to their probability which is 776 given by $p_r(t)$. 777

For each j $(1 \le j \le m)$, there is just one further restriction, concerning the 778 consistency of charges: in order to simultaneously apply several rules of \mathcal{R}_{Π_j} to 779 the same membrane, all the rules must produce the same electrical charge in the 780 membrane in which to be applied. Thus, we will say that the rules of the system, 781 in this computational framework, are applied in a *non-deterministic*, *maximally* 782 *consistent*, and *parallel* way. 783

An *instantaneous description* or *configuration* of the system at any instant *t* is 784 a tuple of multisets of objects present in the *m* environments and at each of the 785 regions of each Π_j , together with the polarizations of the membranes in each P 786 system. We assume that all environments are initially empty and that all membranes 787 initially have a neutral polarization. We assume a global clock exists, synchronizing 788 all membranes and the application of all the rules (from \mathcal{R}_E and from \mathcal{R}_{Π_j} in all 789 environments). 790

In each time unit, we can transform a given configuration in another configuration 791 by using the rules from the whole system as follows: at each transition step, the rules 792 to be applied are selected in a non-deterministic way according to the probabilities 793 assigned to them, and all applicable rules are simultaneously applied in a maximal 794 way. In this way, we get *transitions* from one configuration of the system to the next 795 one. 796

A *computation* is a sequence of configurations such that the first term of the 797 sequence is the initial configuration of the system, and each non-initial configuration 798 of the sequence is obtained from the previous configuration by applying rules of 799 the system in a maximally consistent and parallel manner with the restrictions 800 previously mentioned. 801

6.3.3.2 Simulation Algorithm

The simulation algorithms for PDP systems called BBB and DCBA [21, 27] are 803 based on the grouping of rules into blocks. These groups are constructed by 804 looking the left-hand side. Note that rules having the same left-hand side must 805 have associated probabilities summing 1. Specifically, DCBA works using a refined 806 definition of block, called consistent block [21, 27], as shown in Definition 6.3. 807 DNDP [21, 30] does not use the concept of blocks, but it selects rules by a random 808 loop instead. 809

Definition 6.3. Rules from R and R_E are classified into consistent blocks by either 810 of the following: 811

- (a) The rule block associated with $(i, \alpha, \alpha', u, v)$ is $B_{i,\alpha,\alpha',u,v} = \{r \in R : 812 LHS(r) = (i, \alpha, u, v) \land charge(RHS(r)) = \alpha'\}$
- (b) The rule block associated with (e_j, x) is $B_{e_j, x} = \{r \in R_E : LHS(r) = (e_j, x)\}$. 814

The selection of rules in BBB and DCBA relies always first on selecting blocks, 815 calculating a multinomial random variate, and therefore obtaining a selection of 816 rules within each block. In this sense, we can say that rules within a block will 817 not compete among objects when using BBB and DCBA, because they are selected 818 altogether. This, again, does not hold in DNDP, where rules are selected individually 819 according to the probabilities. Block competition takes place whenever two blocks 820 have distinct but overlapping left-hand sides. 821

DCBA tackles the resource competition issue by performing a proportional s22 distribution of objects among competing blocks. This is done by using the *distribution table*, which is a system-wide time having blocks per columns and pairs s24 (object, region) per rows. Algorithm 3 shows a summary of the algorithm, which s25 can be depicted in [27]. It can be seen that, as usual, each loop iteration is made by s26 two stages: selection and execution. Selection stage consists of three phases: phase s27 1 distributes objects to the blocks in a certain proportional way, phase 2 ensures s28 *maximality* by checking the maximal number of applications of each block, and s29 phase 3 translates from block to rule applications by calculating random numbers s30 using a multinomial distribution. Finally, execution stage (or phase 4) generates the s31 right-hand side of rules.

| Algorithm 3: Sketch of DCBA algorithm for PDP systems |
|---|
| 1: Initialization of the algorithm: static distribution table (columns: blocks, Rows: |
| (objects,membrane)) |
| 2: for $t \leftarrow 0 \dots T$ do |
| 3: Selection stage: |
| 4: Phase 1 (Distribution of objects among rule blocks) |
| 5: Phase 2 (Maximality selection of rule blocks) |
| 6: Phase 3 (Probabilistic distribution, blocks to rules) |
| 7: Execution stage |
| 8: end for |
| |

The proportional distribution of objects along the blocks is carried out through 833 a table which implements the relations between blocks (columns) and objects in 834 membranes (rows). We always start with a static (general) table, and depending on 835 the current configuration of the PDP system, the table is dynamically modified by 836 deleting columns related to non-applicable blocks. Note that after phase 1, we have 837 to assure that the maximality condition still holds. This is normally conveyed by a 838 random loop over the remaining blocks. 839

Finally, DCBA also handles the consistency of rules by defining the concept of 840 consistent blocks [21, 27]: rules within a block have the same left-hand side and 841 the same charge in the right-hand side. There is a further restriction within phase 842 1: if two non-consistent blocks (having different associated right-hand charge) can 843 be selected in a configuration, the simulation algorithm will return an error, or 844 optionally non-deterministically choose a subset of consistent blocks. 845

6.3.3.3 Design of the Parallel Simulator

Normally, the end user (i.e., ecological experts and model designers) runs many 847 simulations on each set of parameters to extract statistical information of the 848 probabilistic model. This can be automated by adding an outermost loop for 849 simulations in the main procedure of the DCBA, which is easily parallelized. 850

At first glance, these two levels of parallelism (simulations and environments 851 [23]) could fit the double parallelism of the CUDA architecture (thread blocks and 852 threads). For example, we could assign each simulation to a block of threads and 853 each environment to a thread (since they require synchronization at each time step). 854 However, the number of environments depends inherently on the model. Typically, 855 2 to 20 environments are considered, which is not enough for fulfilling the GPU 856 resources. Number of simulations typically range from 50 to 100, which is sufficient 857 for thread blocks, but still a poor number compared to the several hundred cores 858 available on modern GPUs.

Thus, the selection of rule blocks (phase 1) and rules (phase 2 and 3) is further parallelized. Hence, the simulator can utilize a huge number of thread blocks by distributing simulations (parallel simulations, as memory can store them) and environments in each one and process each rule block by each thread. Since there are normally more rule blocks (thousand of them) than threads per thread block (up to 512), 256 threads are created, which iterate over the rule blocks in tiles. This design is graphically shown on Fig. 6.4. Each phase of the algorithm has been designed following the general CUDA design explained above and implemented separately as individual kernels. Thus, simulations and environments are synchronized by the successive calls to the kernels.



Fig. 6.4 General design of the CUDA-based simulator: 2D grid and 1D thread blocks. Threads loop the rule blocks in tiles. From [21,28]

6.3.3.4 GPU Implementation of the DCBA Phases

The main challenge at phase 1 is the expanded static table \mathcal{T}_j construction. The 871 size of this table is of order $O(|B| \cdot |\Gamma| \cdot (q + 1))$, where |B| is the number of 872 rule blocks, $|\Gamma|$ is the size of the alphabet (total amount of different objects), and 873 q + 1 corresponds to the number of membranes plus the space for the environment. 874 A full implementation of \mathcal{T}_j can be expensive for large systems and very sparse: 875 competitions for one object appears for a relatively small number of blocks. Thus 876 the expanded static table is implemented by a *virtual table*, which is similar but 877 based on the information of the rules: 878

- Operations over columns: they can be transformed to operations for each rule 879
 block and their left-hand sides.
 880
- Operations over rows: they can be transformed to operations over the left-hand 881 sides of rule blocks and storing the partial results into a global array (one position 882 per row).

Further auxiliary data structures are used to virtually simulate the table [23]: 884

| • | activationVector: the information of filtered blocks is stored here as Boolean | 885 |
|---|---|-----|
| | values. | 886 |
| • | addition: the total sums of the rows are stored using this global vector, one per | 887 |
| | each pair object and region. | 888 |
| • | MinN: the minimum numbers per column are stored here. | 889 |
| • | BlockSel: the total number of applications for each block is stored here. | 890 |
| • | RuleSel: the total number of applications for each rule is stored here. | 891 |
| | | |

The implementation of phase 1 is actually done by means of three kernels, 892 executing one after the other and using the same grid configuration as mentioned 893 in Fig. 6.4. The second and third kernels are executed several times according to 894 parameter A (accuracy) of DCBA [28]: 895

- 1. Kernel for Filters: FILTERS 1 and 2 are implemented here.
- 2. Kernel for Normalization: the two parts (row additions and minimum calculations) of the normalization step is implemented in a kernel. The two parts are synchronized by *synchtreads* CUDA instruction. The work assigned to threads is divergent, that is, each thread works with one rule block, but writes information for each object appearing in the LHS. Therefore, the writes to *addition* are carried out by atomic operations.
- 3. Kernel for Updating and FILTER 2. As before, the work of each thread is 903 divergent. Thus, the update of the configuration is also implemented with atomic 904 operations. Moreover, the *BlockSel* gets updated with the new distribution of 905 selection. 906

870

Phase 2 is the most challenging part when parallelizing by blocks. The selection 907 of blocks is performed in an inherently sequential way: we need to know how much 908 a block can consume before checking the next one. At least, phase 2 can be run 909 simultaneously to each environment and simulation. For this phase, a special version 910 of this kernel was designed. This kernel dynamically calculates the competition of 911 blocks, so that the dependencies of blocks are pre-calculated in order to know which 912 blocks can be selected independently of each other, and everything is done in shared 913 memory. *BlockSel* gets updated with the last selections, and the configuration is 914 also updated to prevent other blocks from being selected. 915

Phase 3 requires a random number generation system for multinomial distributions that was not existing for CUDA. A dedicated implementation was done, 917 called *CURAND_BINOMIAL*, and it is based on the accelerated uniform and normal 918 random variate generation in CUDA with *CuRAND* and the BINV algorithm. This 919 is therefore used to calculate multinomial distributions per rule block and write to 920 *RuleSel*. 921

Finally, for phase 4, the rule selection *RuleSel* is used to generate the right-hand 922 side, by using atomic operations over the configuration. The parallelization is done 923 by using a similar grid configuration as shown in Fig. 6.4, but looping over rules 924 instead of rule blocks. 925

6.3.3.5 Performance Results of the Simulator

In order to test the performance of the simulators, a random generator of PDP 927 systems was designed (designated *pdps-rand*). These randomly created PDP systems have no biological meaning. The purpose is to stress the simulator in 929 order to analyze the implemented designs with different topologies. *pdps-rand* is 930 parametrized in such a way that it can create PDP systems of a desired size. 931

The parallel simulator on the GPU (*pdp-gpu-sim*) and a parallel simulator on 932 multicore CPUs (*pdp-omp-sim*, for 1 (sequential), 2 and 4 cores) are compared. All 933 experiments were run on a GPU server: Linux 64-bit server, with a 4-core (2 GHz) 934 dual socket Intel i5 Xeon Nehalem processor, 12 GBytes of DDR3 RAM, and two 935 NVIDIA Tesla C1060 graphics cards (240 cores at 1.30 GHz, 4 GBytes of memory). 936 GPU cores are typically slower than CPU cores. 937

The test analyses the performance when increasing the parallelism level of the 938 CUDA threads within thread blocks, that is, the number of rule blocks. The speedup 939 achieved by *pdp-gpu-sim* versus *pdp-omp-sim* is shown in Fig. 6.5. The number of 940 simulations is fixed to 50 and the environments to 20 (hence, a total of 1000 thread 941 blocks). The number of objects is proportionally increased together with the number 942 of rule blocks, in such a way that the ratio for number of rule blocks and number of 943 objects is always 2. The mean LHS length is 1.5 (this is normal value for many real 944 ecosystem models, as seen in the literature). The speedup gets stable to around 7× 945 on the number of rule blocks for the GPU versus CPU. For the multicore versions 946 with 2 and 4 CPUs, the speedups are maintained to $4.3 \times$ and $3 \times$, respectively. In 947 the experiments, this number is also achieved when running with 10⁶ rule blocks. 948

As stated in [23], parallelizing by simulations yields the largest speedups 949 on multicore platforms. In order to test the efficiency of the simulator when 950



Fig. 6.5 Scalability of the simulators when increasing the number of rule blocks, from [21, 28]



Fig. 6.6 Scalability of the simulators when increasing the number of simulations, from [21,28]

increasing the number of simulations, rule blocks are fixed to 50000, environments $_{951}$ to 20, objects to 5000, and mean LHS length to 1.5. As shown in Fig. 6.6, the $_{952}$ GPU achieves better runtime than the multicore implementations. The speedup is $_{953}$ maintained to $4.5 \times$ using one core, $3.5 \times$ for 2 cores, and $2.7 \times$ for 4 cores. $_{954}$

6.4 Specific Simulations

In this section, we will introduce a type of simulation of P systems, which is called 956 specific simulation [25]. We will describe how to implement this kind of simulators 957 in CUDA and provide two major examples. 958

6.4.1 Definition

When implementing a P system simulator, it is important to understand what type 960 of P systems we want to simulate before starting the development. We will say that 961 a specific simulator is any simulator developed for just a certain P system or family 962 of P systems. In other words, if we focus on just one P system or a parametrized 963 definition of P systems forming a family, then we need a specific simulator. On the 964 one hand, this is a restricted version of a simulator, since it can handle a reduced 965 variety of P systems, but it also helps to adjust better the simulator to parallel 966 architectures. 967

Basically, restricting the P systems to be simulated helps us to take a better 968 control of the algorithm by predicting when certain things will happen. On the one 969 hand, by just focusing on a P system variant, we will know which kind of rules 970 and what semantics apply. On the other hand, by knowing exactly the P systems to 971 simulate, we can develop tailored code to certain parts of the model. For example, 972 we can see in the literature that when developing solutions to certain problems, it 973 is very useful to design it by making a scheme of the computation. This means 974 that the computation tree of the designed P system is usually bound and can be 975 divided into stages. For example, in SAT solutions, there is usually a stage where an 976 exponential amount of membranes is generated by applying division rules, and it is 977 known at which moment the stage starts and ends, because it is part of the design of 978 the solution. By making an exhaustive analysis, it would be even possible to predict 979 which objects can appear in which membranes and when.

Therefore, specific simulators can take advantage of that information in order to adapt the code and the data structures. Specific functions and kernels can be written for each stage, and the memory layout to store the P system information can be drastically reduced. In fact, the information of rules (left and right-hand sides) can be encoded in the source code, instead of storing them in memory, because we know the rules.

The memory layout is also an essential part of a simulator, since P system 987 simulators have been demonstrated to be memory and memory bandwidth bound. 988 When storing the information of P system configuration, we can use either of the 989 following[25]: 990

- *Sparse representation*: using a large array to store multiplicities, with a position 991 per each possible object (all objects defined in the alphabet). The access is direct 992 since the object identifier is the index where to access the array. However, if many 993 objects are not present at a certain moment, the array will be full of zeroes. 994
- *Dense representation*: using a double array with a component for the object 995 identifier and the other for multiplicity. We need to search for the object, unless 996 we track them and we know exactly where they are store at any moment. This 997 can help to drastically reduce the size because objects with multiplicity zero can 998 be discarded. 999

Specific simulators can use dense representation, if we can know where each 1000 object is at every transition step, so that the rules can access them directly. Moreover, 1001 there are many objects that work as counters, for example o_i , $0 \le i \le n$. They are 1002 distinct objects but at the end they are related. If we know that o_i and o_j , with $i \ne j$ 1003 will not be present at the same time, then we can use just one position, or even use 1004 a variable (register) to store the subindex. 1005

Finally, it is important to remark that specific simulators must keep being full 1006 simulators; that is, if the simulator goes beyond the P system model and skips 1007 representing the P system features, then we are simulating something else. In other 1008 words, we will say that a program is simulating a P system if we can ask the 1009 program, at any transition step, any piece of information of the state of the P system 1010 (configuration, rules applied, etc.). Thus, a specific simulator should be developed 1011 in such a way that we could extract from it the configuration of the P system or the 1012 rules that have been applied.

6.4.2 Simulating a SAT Solution with Active Membrane P Systems 1014

The first specific simulator implemented in CUDA was a family of recognizer 1015 P systems with active membranes designed to solve the SAT problem in linear time 1016 (but with exponential workspace). In this section, we will discuss the design of this 1017 simulator and its performance achieved with CUDA. 1018

The original work is published in [3, 5, 21, 22]. The framework of all these 1019 simulators is named *PCUDASAT*, and it can be downloaded from the official website 1020 http://sourceforge.net/p/pmcgpu [45] or the repository https://github.com/RGNC/ 1021 pcudasat.

6.4.2.1 SAT Solution with Active Membranes

Let $\varphi = C_1 \land \cdots \land C_m$ be a propositional formula in **CNF** such that the set of 1024 variables of the formula is $Var(\varphi) = \{x_1, \dots, x_n\}$, consisting of *m* clauses $C_i = 1025$ $y_{i,1} \lor \cdots \lor y_{i,k_i}, 1 \le i \le m$, where $y_{i,i'} \in \{x_j, \neg x_j : 1 \le j \le n\}$ are the literals of φ . 1026 We can assume that the formula is in simplified expression, i.e., no clause contains 1027 two occurrences of the same literal, and no clause can contain, simultaneously, a 1028 literal and its negation. The SAT problem is defined as follows: given a Boolean 1029 formula in conjunctive normal form (CNF), to determine whether or not there exists 1030 a truth assignment to its variables on which the formula evaluates true. 1031

The solution to SAT based on recognizer P system with active membranes is 1032 defined as $\Pi_{am-SAT}(\langle m, n \rangle) = (\Gamma, \Sigma, \mu, \mathcal{M}_1, \mathcal{M}_2, \mathcal{R}, 2)$ of degree 2, for each pair 1033 of natural numbers $m, n \in \mathcal{N}$. Specifically: 1034

- The input alphabet is $\Sigma = \{x_{i,j}, \overline{x}_{i,j} | 1 \le i \le m, 1 \le j \le n\}.$ 1035
- · The working alphabet is

$$\Gamma = \Sigma \cup \{c_k | 1 \le k \le m+2\} \cup \{d_k | 1 \le k \le 3n+2m+3\} \cup \\ \cup \{r_{i,k} | 0 \le i \le m, 1 \le k \le 2n\} \cup \{e,t\} \cup n\{Yes, No\}$$
¹⁰³⁷

1023

| • | The set of labels is $\{1, 2\}$. | 1038 |
|---|--|------|
| • | The initial structure of membranes is $\mu = [[]_2]_1$. | 1039 |
| • | The initial multisets associated with the membranes are $\mathcal{M}_1 = \emptyset \setminus \mathcal{M}_2 = \{d_1\}$. | 1040 |
| • | The input membrane is the one labeled by 2. | 1041 |
| • | The set \mathcal{R} consists of the following rules: | 1042 |
| | 6 | |
| | (a) $[d_k]_2^0 \to [d_k]_2^+ [d_k]_2^-$, for $1 \le k \le n$. | 1043 |
| | (b) $[x_{i,1} \to r_{i,1}]_2^+$, $[\overline{x}_{i,1} \to r_{i,1}]_2^-$, for $1 \le i \le m$. | 1044 |
| | $[x_{i,1} \rightarrow \lambda]_2^-, \ [\overline{x}_{i,1} \rightarrow \lambda]_2^+, \text{ for } 1 \le i \le m.$ | 1045 |
| | (c) $[x_{i,j} \to x_{i,j-1}]_2^+$, $[x_{i,j} \to x_{i,j-1}]_2^-$, for $1 \le i \le m$, $2 \le j \le n$. $[\overline{x}_{i,j} \to x_{i,j-1}]_2^-$ | 1046 |
| | $\overline{x}_{i,j-1}]_2^+$, $[\overline{x}_{i,j} \rightarrow \overline{x}_{i,j-1}]_2^-$, for $1 \le i \le m, \ 2 \le j \le n$. | 1047 |
| | (d) $[d_k]_2^+ \to []_2^0 d_k$, $[d_k]_2^- \to []_2^0 d_k$, for $1 \le k \le n$. | 1048 |
| | $d_k[]_2^0 \to [d_{k+1}]_2^0$, for $1 \le k \le n-1$ }. | 1049 |
| | (e) $[r_{i,k} \to r_{i,k+1}]_2^0$, for $1 \le i \le m, \ 1 \le k \le 2n-1$. | 1050 |
| | (f) $[d_k \to d_{k+1}]_1^0$, for $n \le k \le 3n-3$; $[d_{3n-2} \to d_{3n-1}e]_1^0$. | 1051 |
| | (g) $e[]_2^0 \to [c_1]_2^+; [d_{3n-1} \to d_{3n}]_1^0.$ | 1052 |
| | (h) $[d_k \to d_{k+1}]_1^0$, for $3n \le k \le 3n + 2m + 2$. | 1053 |
| | (i) $[r_{1,2n}]_2^+ \to []_2^- r_{1,2n}$. | 1054 |
| | (j) $[r_{i,2n} \to r_{i-1,2n}]_2^-$, for $1 \le i \le m$. | 1055 |
| | (k) $r_{1,2n}[]_2^- \to [r_{0,2n}]_2^+$. | 1056 |
| | (l) $[c_k \to c_{k+1}]_2^-$, for $1 \le k \le m$. | 1057 |
| | (m) $[c_{m+1}]_2^+ \to []_2^+ c_{m+1}.$ | 1058 |
| | (n) $[c_{m+1} \to c_{m+2}t]_1^0$. | 1059 |
| | (<i>o</i>) $[t]_{1}^{0} \to []_{1}^{+}t.$ | 1060 |
| | (p) $[c_{m+2}]_1^+ \to []_1^- Yes.$ | 1061 |
| | $(q) \ [d_{3n+2m+3}]_1^0 \to []_1^+ No.$ | 1062 |
| | | |

We also consider a polynomial encoding (cod, s) of the SAT problem in the 1063 family $\Pi_{am-SAT} = \{\Pi_{am-SAT}(t) \mid t \in N\}$. The function *cod* associates to the 1064 previously described propositional formula φ (an instance of SAT with *n* variables 1065 and *m* clauses), the following multiset of objects 1066

$$cod(\varphi) = \bigcup_{i=1}^{m} \{x_{i,j} | x_j \in C_i\} \cup \{\overline{x}_{i,j} | \neg x_j \in C_i\}$$
 1067

In this case, object $x_{i,j}$ represents that variable x_j in clause C_i .

The *size* function, *s*, is defined as follows $s(\varphi) = \langle m, n \rangle = \frac{(m+n)\cdot(m+n+1)}{2} + m$. 1069 Then, $cod(\varphi)$ is an input multiset of the system $\prod_{am-SAT}(s(\varphi))$ and the pair 1070 (cod, s) is therefore a polynomial encoding of the SAT problem in the family 1071 \prod_{am-SAT} . Thus, the system of the family $\prod_{am-SAT} processing the instance \varphi$ will 1072 be the P system with active membranes $\prod_{am-SAT}(s(\varphi))$ with input multiset $cod(\varphi)$. 1073

The system $\Pi_{am-SAT}(s(\varphi))$ with input $cod(\varphi)$ is confluent, and its computation 1074 is structured in four phases as follows: 1075

- *Generation phase*: all possible relevant truth assignment is generated for the set 1076 of variables of the formula $\{x_1, \ldots, x_n\}$. It is achieved by using division rules 1077 in the internal membranes (labeled by 2). This will allow the generation of 2^n 1078 membranes that will encode all possible assignments. Nevertheless, in this phase, 1079 while the valuations are being generated, the clauses that are true by the encoded 1080 valuation in each internal membrane are checked. This idea is implemented 1081 through a very sophisticated process by which only the truth values 1 and 0 are 1082 given to the variable 1. This variable 1 corresponds to the variable x_1 in the first 1083 loop step, but by a set of indices, the variable 1 corresponds to the variable x_2 in 1084 the second loop step, and so on. This phase is executed in 3n 1 computation 1085 steps, and only the rules (a), (b), (c), (d) and (e) are applied.
- Synchronization phase: it prepares the system for the checking phase synchronizing the execution of the system by unifying certain sub-indices of some objects. 1088 The execution of this phase consumes 2n computation steps, and only rules 1089 (e), (f) and (g) are executed. 1090
- *Check-out phase*: in this phase, it is determined how many clauses are true for 1091 each truth assignment encoded by the internal membranes. This is done using the 1092 objects c_k (k > 1), whose appearance in a membrane means that exactly k 1 1093 clauses are made true by the encoded valuation in that membrane. This phase is 1094 executed in 2m steps, and rules (h), (i), (j), (k) and (l) are applied. 1095
- *Output phase*: in this phase the system provides the corresponding output 1096 depending on the analysis of the check-out phase. That is, this step performs 1097 a search of the internal membranes encoding a solution (i.e., containing object 1098 c_{m+1}). If a membrane satisfies the above condition, the object *Yes* is sent to 1099 the environment, and the system stops. Otherwise, the object *No* is sent to the 1100 environment and the system stops. The execution of this phase is done in 4 steps 1101 and the used rules are (m), (n), (o), (p) and (q).

6.4.2.2 Sequential Simulator and Data Structures

The sequential simulator design is based on the four main phases of a P system 1104 computation from Π_{am-SAT} : generation, synchronization, check-out, and output. 1105 Thus, the computation of the P system to simulate (from the family Π_{am-SAT}) is 1106 reproduced by sequentially executing these phases. Firstly, the generation phase 1107 is executed, generating 2^n membranes by dividing each one in *n* steps, where *n* 1108 is the number of variables of the input CNF formula. Since we know the value 1109 of *n*, the simulator knows the amount of membranes to generate before starting 1110 the simulation. After that, the simulator executes the synchronization phase which 1111 evolves the objects following the rules previously explained. The check-out phase the output phase sends out the correct answer to the environment. 1114

It is important to remark that the semantics of the P system is reproduced by 1115 the simulation algorithm, so the simulator is specific for this solution. Thus, the 1116

only input for the simulator is the CNF formula provided in DIMACS CNF format. 1117 We can assume therefore that the simulator behaves as a SAT solver, receiving a 1118 propositional formula and giving the corresponding answer. However, this solver is 1119 implemented following a solution by means of P systems. 1120

The first challenge is to decrease the sparsity of the data structures storing the 1121 configuration of the P system. After an exhaustive analysis of the computation, the 1122 upper bound of number of distinct objects appearing in a membrane can be fixed 1123 to the size of the input multiset (the number of literals in the input propositional 1124 formula). Indeed, one can observe that the size of the right-hand side of evolution 1125 rules is always 1. Thus, every object in the input multiset always evolves to either 1126 another or disappears. By definition, send-in, send-out, and division rules do not 1127 generate more than one object in the right-hand side.

Hence, the representation of the P system is made by an array storing the 1129 multisets of objects for every membrane labeled by 2. The amount of elements 1130 per membrane equals to, as mentioned above, the size of the input multiset (total 1131 number of literals in the formula, $|cod(\varphi)|$). This array is initially allocated for the 1132 maximum amount of membranes 2 that the P system will create, which is 2^n (note 1133 that *n* is defined in the input file). Only the first one is initialized by storing the full 1134 input multiset. Division rules will initialize each membrane later on.

The encoding of objects for the input multiset can be made at a bit-level within 1136 integers of 32 bits. Each integer stores the following (8 bits for each field): 1137

| 1. ' | The name of the object (x or \overline{x}) | 11 | 38 |
|------|---|-----|----|
| 2. | Reserved space. | 11: | 39 |
| 3. | Variable (subindex <i>i</i>). | 11. | 40 |
| 4. | Clause (subindex <i>j</i>). | 11/ | 41 |

It is noteworthy that the membrane charges are not stored, since we can observe 1142 from the computation that a partition of membranes having positive and negative 1143 charges can be done over the array. In other words, the first half of membranes 1144 are positive, and the other half (new ones) negative. The skin membrane is not 1145 represented, since its purpose is to store objects sent out from membranes, those 1146 which are sent in to the same membranes in the next step. This process is therefore 1147 simulated within each membrane, avoiding to store the information for the skin 1148 membrane. Other objects, such as *yes*, *no*, and *c* counter, are also placed as variables 1149 encoded directly in source code. 1150

6.4.2.3 Design of the GPU Simulator

The parallel simulator on the GPU was designed to take over the most demanding 1152 phases on the computation of Π_{am-SAT} , which are the first three phases. The last 1153 one (output phase) is developed on the CPU. In order to map the parallelism into the 1154 GPU, the simulator assigns a thread block to each membrane, as shown in Fig. 6.7. 1155 In this way, the parallelism among membranes is represented. Moreover, each thread 1156 is assigned to each object of the input multiset, which is a literal of the input formula 1157 (with the exception of object d_1). This mapping is common to all the defined kernels. 1158



Fig. 6.7 Design of the parallel simulator for Π_{am-SAT} . Each generated membrane is assigned to a thread block, and each object (initially, the input multiset) is assigned to a thread. From [3,21]

Algorithm 4 shows the pseudocode for the host side the simulator. The generation 1159 phase is simulated by using three kernels which execute the rules in this phase. This 1160 is an iterative process of n steps where the kernels are called n times. A tailored 1161 kernel for division is designed to make copies of all membranes executing division 1162 rules, whose behavior is shown in Fig. 6.8. There is a double parallelism in the 1163 division, as shown in the figure; threads within a block are in charge of making 1164 the copies for the new membrane and changing just the corresponding object. But 1165 this is also repeated for each thread block. In each iteration, the simulator adjusts 1166 the number of thread blocks before calling the kernel, since new membranes are 1167 created. That is, the membranes are distributed along the two-dimensional grid of 1168 thread blocks. 1169

When the exponential amount of membranes is created, synchronization and 1170 check-out phases are executed. This is simulated within just one kernel for both 1171 phases, in parallel for each membrane. Global synchronization is not necessary 1172 because there is no communication among the internal membranes at these phases. 1173 Finally, the output phase is developed on the CPU, checking the conditions and 1174 launching the result of the computation. 1175



Fig. 6.8 Parallel division on GPU at generation phase in Π_{am-SAT} . Each membrane is divided by a thread block, which copies the objects also in parallel by using several threads. If the number of objects is larger than 256, the thread block repeats the process until covering all the objects. From [5]

Algorithm 4: Parallel Simulator of Π_{am-SAT} , at host side. From [3,21]

- 1: {Initialization}
- 2: *Threads* $\leftarrow |cod(\varphi)|$ {The number of literals in the CNF formula}
- 3: Blocks \leftarrow (1, 1) {One block in the 2-dimensional grid}
- 4: $d \leftarrow 0$ {A counter}
- 5: *numMembranes* \leftarrow 1 {Number of membranes}
- 6: *psystem* ← *allocateGPUMemory*(2^{*n*}) {Allocate enough memory to represent the P system}
- 7: {Generation phase}
- 8: repeat
- 9: Division_kernel <<< Blocks, Threads >>> (psystem, numMembranes)
- 10: $numMembranes \leftarrow numMembranes \times 2$
- 11: Blocks ← AdjustBlocks(psystem, numMembranes) {Distribute membranes among blocks}
- 12: Send_out_kernel <<< Blocks, Threads >>> (psystem, numMembranes)
- 13: Send_in_kernel <<< Blocks, Threads >>> (psystem, numMembranes)
- 14: $d \leftarrow d + 1$
- 15: **until** d < n {Repeat n times (number of variables)}
- 16: {Synchronization and Check-out phases}
- 17: Syn_Check_kernel <<< Blocks, Threads >>> (psystem, numMembranes)
- 18: {Output phase (executed on the CPU)}
- 19: *Output(psystem, numMembranes)*



Fig. 6.9 Simulation performance for *am-sat-gpu* vs *am-sat-seq* when increasing the number of membranes (x-axis). From [3, 21]

6.4.2.4 Performance Analysis

Next, we analyze the performance of the simulators above described for Π_{am-SAT} : 1177 the sequential simulator developed in C++ (from now, *am-sat-seq*) and the GPU 1178 parallel simulator on CUDA (*am-sat-gpu*). The experimental results were obtained 1179 using a Tesla C1060 GPU. 1180

Figure 6.9 shows the experimental performance of the cell-like simulators (in a 1181 log scale) when increasing the number of membranes in the P system (and hence, 1182 the number of blocks in the GPU and also the variables in the CNF formula) until 1183 reaching 2^{12} membranes. The number of simulated membranes is restricted by the 1184 available memory of the system. The number of literals in the formula is fixed to 1185 256, which means 256 threads per block.

It can be seen that once the GPU resources have been fully occupied, the 1187 execution time increases linearly with the number of blocks. In this case, we report 1188 up to $94 \times$ of speedup between *am-sat-seq* and *am-sat-gpu*. However, Fig. 6.9 shows 1189 the speedup becomes a constant number of $100 \times$ when the number of membranes is 1190 greater than 128 K.¹ This is the number of blocks launched in the grid of the GPU. 1191

We finalize the performance analysis by also considering the data management 1192 (allocation and transfer) time of the GPU. This is also very important, because it 1193 is part of the solution. Figure 6.10 shows the speedup achieved by comparing *am*-1194 *sat-gpu* (with data management) and *am-sat-seq*. We can see that for small amounts 1195 of membranes, the speedup is below 1, what means a worst performance. However, 1196 after 32 K membranes, the speedup is $1.23 \times$, and it is increased along with the 1197 number of membranes until $64 \times$ for 4 M membranes. This is caused by the decrease 1198 in the kernels time, and the time of handling the data is almost constant for any 1199

¹Note that we use here "K" and "M" for binary prefixes "kilo" and "mega", respectively. Therefore, $128 \text{ K} = 2^{17} = 131072$.



Fig. 6.10 Achieved speedup of *am-sat-gpu* against *am-sat-seq* considering also the GPU data management, when increasing the number of membranes (x-axis). From [3,21]

size. Note that the data management performed by *am-sat-gpu* is the following: data 1200 allocation, initial configuration (only 1 membrane) transfer, and answer (object yes 1201 or not) transfer. The information of the P system during the computation is always 1202 kept on the GPU memory. 1203

6.4.3 Simulating a SAT Solution with Tissue P Systems

In this section, we depict a specific simulator for the family of recognizer tissue P 1205 systems with cell division. We first explain the data structures and the phases that 1206 compound the simulation algorithm, then the sequential version, and after that the 1207 parallel one based on CUDA, describing the different optimizations taken for each 1208 phase of the simulator. 1209

The original work is published in [21, 22, 26]. This simulation framework 1210 is named *TSPCUDASAT*, and it can be downloaded from the official website 1211 http://sourceforge.net/p/pmcgpu [45] or the repository https://github.com/RGNC/ 1212 tspcudasat.

6.4.3.1 Recognizer Tissue P System with Cell Division

In the paradigm of membrane computing, a new computing model (*tissue P system* 1215 *with cell division*) is introduced by using the biological membranes placed in the 1216 nodes of a directed graph, inspired from the cell inter-communication in tissues 1217 [20]. Besides, cell division is an elegant process that enables organisms to grow 1218 and reproduce. Mitosis is a process of cell division which results in the production 1219 of two daughter cells from a single parent cell. Daughter cells are identical to one 1220 another and to the original parent cell. Through a sequence of steps, the replicated 1221 genetic material in a parent cell is equally distributed to two daughter cells. While 1222 there are some subtle differences, mitosis is remarkably similar across organisms. 1223

1204

Definition 6.4. A tissue P system with cell division of degree q > 1, is a tuple 1224 $\Pi = (\Gamma, \mathcal{E}, \mathcal{M}_1, \dots, \mathcal{M}_q, \mathcal{R}, i_{out}),$ where: 1225

- 1. Γ is a finite *alphabet*:
- 2. $\mathcal{E} \subseteq \Gamma$;
- 3. $\mathcal{M}_1, \ldots, \mathcal{M}_a$ are finite multisets over Γ ;
- 4. \mathcal{R} is a finite set of communication rules of the following forms:
 - (a) Communication rules: (i, u/v, j), for $i, j \in \{0, 1, 2, \dots, q\}, i \neq j, u, v \in \{0, 1, 2, \dots, q\}$ 1230 $M(\Gamma), |u| + |v| > 0;$ 1231
 - (b) Division rules: $[a]_i \rightarrow [b]_i [c]_i$, where $i \in \{1, 2, \dots, q\}, i \neq i_{out}$ and 1232 $a, b, c \in \Gamma$. 1233
- 5. $i_{out} \in \{0, 1, \dots, q\}.$

A tissue P system with cell division $\Pi = (\Gamma, \mathcal{E}, \mathcal{M}_1, \dots, \mathcal{M}_q, \mathcal{R}, i_{out})$ of degree $q \geq 1$ 1235 1, can be viewed as a set of q cells, labeled by $1, \ldots, q$, with an environment labeled 1236 by 0 such that (a) $\mathcal{M}_1, \ldots, \mathcal{M}_q$ represent the finite multisets of *objects* (symbols of 1237) the working alphabet Γ) initially placed in the q cells of the system; (b) \mathcal{E} is the set 1238 of objects initially located in the environment of the system, all of them available in 1239 an arbitrary number of copies; (c) \mathcal{R} is a finite set of rules over Γ associated with 1240 the cells and the environment; and (d) $i_{out} \in \{0, 1, \dots, q\}$ indicates the output zone. 1241 We use the term *zone* i to refer to cell i, in the case $1 \le i \le q$ and to refer the 1242 environment in the case i = 0. 1243

A communication rule (i, u/v, j) is called a symport rule if $u = \lambda$ or v = 1244 λ . A symport rule $(i, u/\lambda, j)$ provides a virtual arc from zone i to zone j. A 1245 communication rule (i, u/v, j) is called an *antiport rule* if $u \neq \lambda$ and $v \neq \lambda$. An 1246 antiport rule (i, u/v, j) provides two arcs: one from zone i to zone j and another one 1247 from zone j to zone i. Therefore, every tissue P system has an underlying directed 1248 graph whose nodes are the zones (cells and the environment) of the system and the 1249 arcs are obtained from communication rules. 1250

A configuration (or instantaneous description) C_t at an instant t of a tissue P 1251 system Π is a tuple whose components are the multisets over Γ associated with 1252 each cell present in the system at moment t and the multiset over $\Gamma \setminus \mathcal{E}$ associated 1253 with the environment at moment t. 1254

A communication rule (i, u/v, j) is applicable to zones i, j to a configuration 1255 C_t at instant t, if in that configuration the multiset u is contained in one zone i and 1256 multiset v is contained in one zone j. When applying such a communication rule to 1257such zones, the objects of the multiset represented by u are sent from zone i to zone 1258 *j*, and simultaneously, the objects of multiset v are sent from zone *j* to zone *i*. The $_{1259}$ *length* of communication rule (i, u/v, j) is defined as |u| + |v|. 1260

A division rule $[a]_i \rightarrow [b]_i[c]_i$ is applicable to a configuration at an instant 1261 t, if one cell i belongs to that configuration containing object a. When applying a 1262division rule $[a]_i \rightarrow [b]_i[c]_i$ to such a cell *i*, under the influence of object *a*, that cell 1263 is divided into two cells with the same label i; in the first copy, object a is replaced 1264 by object b, in the second one, object a is replaced by object c; all the other objects

1226 1227

1228

1229

residing in such a cell *i* are replicated, and copies of them are placed in the two new $_{1265}$ cells. The output cell i_{out} and any cell with input degree equal to zero cannot be $_{1266}$ divided. $_{1267}$

The rules of a tissue P system with cell division are applied as follows: 1268 communication rules will be applied in a non-deterministic maximally parallel 1269 manner as it is customary in membrane computing but with the following important 1270 remark: if a cell divides, then the division rule is the only one which is applied 1271 for that cell at that step; the objects inside that cell do not evolve by means 1272 of communication rules. In other words, before division a cell interrupts all its 1273 communication channels with the other cells and with the environment. The new 1274 cells resulting from division will interact with other cells or with the environment 1275 only at the next step—providing that they do not divide once again. The label of a 1276 cell precisely identifies the rules which can be applied to it.

Given a tissue P system with cell division, $\Pi = (\Gamma, \mathcal{E}, \mathcal{M}_1, \dots, \mathcal{M}_q, \mathcal{R}, i_{out})$, the 1278 initial configuration of Π is $C_0 = (\mathcal{M}_1, \dots, \mathcal{M}_q; \emptyset)$. A configuration is a halting 1279 configuration if no rule of the system is applicable to it. We say that configuration 1280 C_1 yields configuration C_2 in one transition step, denoted $C_1 \Rightarrow_{\Pi} C_2$, if we can 1281 pass from C_1 to C_2 by applying the rules from \mathcal{R} following the previous remarks. In 1282 tissue P systems with cell division, the concepts of computation, recognizer system, 1283 and polynomial time and uniform solution to a decision problem are introduced in a 1284 similar way than in P systems with active membranes.

6.4.3.2 SAT Solution with Tissue P Systems

This section presents an efficient solution to SAT problem by means of family 1287 of recognizer tissue P systems with cell division. Let $\varphi = C_1 \land \cdots \land C_m$ be 1288 a propositional formula in **CNF** such that the set of variables of the formula is 1289 $Var(\varphi) = \{x_1, \ldots, x_n\}$ and consists of *m* clauses $C_j = y_{j,1} \lor \cdots \lor y_{j,k_j}, 1 \le i \le m$, 1290 where $y_{j,j'} \in \{x_i, \neg x_i : 1 \le i \le n\}$ are the literals of φ . Without loss of generality, 1291 we can assume that the formula is in simplified expression.

For each pair of natural numbers $m, n \in \mathbf{N}$, we will consider the recognizer 1293 tissue P system with cell division $\Pi_{tsp-SAT}(\langle m, n \rangle) = (\Gamma, \Sigma, \mu, \mathcal{M}_1, \mathcal{M}_2, R, 2)$ of 1294 degree 2, defined as follows: 1295

- The input alphabet is $\Sigma = \{x_{i,j}, \overline{x}_{i,j} | 1 \le i \le n, 1 \le j \le m\}$
- The working alphabet is

$$\Gamma = \Sigma \cup \{a_i, t_i, f_i \mid 1 \le i \le n\} \cup \{r_i \mid 1 \le i \le m\} \cup \cup \{T_i, F_i \mid 1 \le i \le n\} \cup \{T_{i,j}, F_{i,j} \mid 1 \le i \le n, 1 \le j \le m+1\} \cup \cup \{b_i \mid 1 \le i \le 2n+m+1\} \cup \{c_i \mid 1 \le i \le n+1\} \cup \cup \{d_i \mid 1 \le i \le 2n+2m+nm+1\} \cup \cup \{e_i \mid 1 \le i \le 2n+2m+nm+3\} \cup \{f, g, yes, no\}$$

$$1298$$

- The environment alphabet is $\mathcal{E} = \Gamma \{ yes, no \}$.
- The set of labels is $\{1, 2\}$.

1286

1296

1297

1299

| • | The initial multisets associated with the cells are $\mathcal{M}_1 = \{\text{yes}, \text{no}, b_1, c_1, d_1, e_1\}$ | 1301 |
|----|---|------|
| | and $\mathcal{M}_2 = \{f, g, a_1, a_2, \dots, a_n\}.$ | 1302 |
| • | The input cell is the one labeled by 2, and the output region is the environment. | 1303 |
| • | The set \mathcal{R} is formed by the following rules: | 1304 |
| | | |
| 1. | Division rule: | 1305 |
| | (a) $[a_i]_2 \to [T_i]_2 [F_i]_2$, for $i = 1, 2,, n$. | 1306 |
| 2. | Communication rules: | 1307 |
| | (b) $(1, b_i/b_{i+1}^2, 0)$, for $i = 1,, n$. | 1308 |
| | (c) $(1, c_i/c_{i+1}^2, 0)$, for $i = 1,, n$. | 1309 |
| | (d) $(1, d_i/d_{i+1}^2, 0)$, for $i = 1,, n$. | 1310 |
| | (e) $(1, e_i/e_{i+1}, 0)$, for $i = 1,, 2n + 2m + nm + 2$. | 1311 |
| | (f) $(1, b_{n+1}c_{n+1}/f, 2)$. | 1312 |
| | (g) $(1, d_{n+1}/g, 2)$. | 1313 |
| | (h^*) $(1, f^2/f, 0).$ | 1314 |
| | (h) $(2, c_{n+1}T_i/c_{n+1}T_{i,1}, 0)$, for $i = 1,, n$. | 1315 |
| | (i) $(2, c_{n+1}F_i/c_{n+1}F_{i,1}, 0)$, for $i = 1,, n$. | 1316 |
| | (j) $(2, T_{i,j}/t_i T_{i,j+1}, 0)$, for $i = 1,, n$ and $j = 1,, m$. | 1317 |
| | (k) $(2, F_{i,j}/f_i F_{i,j+1}, 0)$, for $i = 1,, n$ and $j = 1,, m$. | 1318 |
| | $(l) (2, b_i/b_{i+1}, 0).$ | 1319 |
| | (m) $(2, d_i/d_{i+1}, 0)$, for $i = n + 1,, 2n + m$. | 1320 |
| | (n) $(2, b_{2n+m+1} t_i x_{i,j}/b_{2n+m+1} r_j, 0).$ | 1321 |
| | (<i>o</i>) $(2, b_{2n+m+1} f_i \overline{x}_{i,j}/b_{2n+m+1} r_j, 0)$, for $1 \le i \le n$ and $1 \le j \le m$. | 1322 |
| | (p) $(2, d_i/d_{i+1}, 0)$, for $i = 2n + m + 1, \dots, 2n + m + nm$. | 1323 |
| | (q) $(2, d_{2n+m+nm+j} r_j/d_{2n+m+nm+j+1}, 0)$, for $j = 1,, m$. | 1324 |
| | (r) $(2, d_{2n+2m+nm+1}/f \text{ yes}, 1).$ | 1325 |
| | (s) $(2, yes/\lambda, 0)$. | 1326 |
| | (t) $(1, e_{2n+2m+nm+3} f no/\lambda, 0)$. | 1327 |

Next, we consider a polynomial encoding (cod, s) of the SAT problem in the 1328 family $\Pi_{tsp-SAT} = {\Pi_{tsp-SAT}(t) | t \in \mathbb{N}}$. The function *cod* associates to 1329 the previously described propositional formula φ that is an instance of SAT with 1330 parameters *n* (number of variables) and *m* (number of clauses), with the following 1331 multiset of objects 1332

$$cod(\varphi) = \bigcup_{i=1}^{m} \{x_{i,j} | x_i \in C_j\} \cup \{\overline{x}_{i,j} | \neg x_i \in C_j\}$$
1333

In this case, object $x_{i,j}$ represents that variable x_i belongs to clause C_j . The *size* 1334 function, *s*, is defined as follows $s(\varphi) = \langle m, n \rangle = \frac{(m+n) \cdot (m+n+1)}{2} + m$. The 1335 system of the family $\Pi_{tsp-SAT}$ to process the instance φ will be the tissue P system 1336 $\Pi_{tsp-SAT}(s(\varphi))$ with input multiset $cod(\varphi)$.

The execution of the system $\Pi_{tsp-SAT}(s(\varphi))$ with input $cod(\varphi)$ is structured in 1338 six phases: 1339

- Valuations generation phase: in this phase all the possible relevant truth valuations are generated for the set of variables of the formula $\{x_1, \ldots, x_n\}$. This is 1341 accomplished by using division rules (*a*), whereby each object x_i produces two 1342 new cells, one having the object T_i that codifies the true value of the variable x_i , 1343 y and the other having the object T_i that codifies the false value of the variable 1344 x_i . Thus, 2^n cells are obtained in *n* computation steps. These cells are labeled 1345 by 2, and each one codifies each possible truth valuation of the set of variables 1346 $\{x_1, \ldots, x_n\}$. Meanwhile, the objects *f*, *g* are replicated in each created cell. This 1347 phase spends *n* computation steps. 1348
- Counter generation phase: simultaneously, and using the rules (b), (c), (d), and 1349 (e), the counters b_i , c_i , d_i , e_i of the cell labeled by 1 are evolving such that in 1350 each computation step the number of objects in each one is doubling. Thereby, 1351 through this process and after *n* steps, we get 2^n copies of the objects b_{n+1} , c_{n+1} , 1352 and d_{n+1} . Objects b's will be used to check which clauses are satisfied for each 1353 truth valuation. Objects c's are used to obtain a sufficient number of copies of 1354 t_i , f_i (namely, *m*). Objects d's will be used to check if there is at least one 1355 valuation satisfying all clauses. Finally, objects e's will be used to produced, 1356 in its case, the object no at the end of the computation. 1357
- *Checking preparation phase*: this phase aims at preparing the system for checking clauses. For this, at step n + 1 of the computation, and by the application of the rules (f) and (g), the counters b_{n+1} , c_{n+1} , d_{n+1} of the cell 1 are exchanged for the objects f and g of the 2^n cell 2. Thus, after this step, each cell labeled by two has a copy of the objects b_{n+1} , c_{n+1} , d_{n+1} , while cell 1 has 2 copies of the objects f and g.

Subsequently, the presence of an object c_{n+1} in each one of the 2^n cells labeled 1364 by 2 allows to generate the objects $T_{i,1}$ and $F_{i,1}$. By the application of rules (j) 1365 and (k), these objects allow the emergence of m copies of t_i and m copies of 1366 f_i , according to the values of truth or falsity that a cell 2 assigns to a variable 1367 x_i . This process spends n + m steps since there is only one object c_{n+1} in each 1368 cell 2, and moreover, for each $i = 1, \ldots, n$, the rules (j) and (k) are applied 1369 exactly m consecutively times. Simultaneously, in the first steps of this process, 1370 the application of the rule (h^*) makes the cell labeled by 1 to appear only one 1371 copy of the object yes. Simultaneously in this phase, the counters b_i , d_i and e_i 1372 are evolving by the applications of the corresponding rules. 1373

• *Checking clauses phase*: in this phase it is determined which clauses are true 1374 for every truth valuation encoded by a cell labeled by 2. This phase starts at 1375 the computation step (n + 1) + (n + m) + 1 = 2n + m + 2. Using the rules 1376 (*n*) and (*o*), the true clauses are checked for each valuation encoded by a cell, 1377 so that the appearance of an object r_j in a cell 2 means that the corresponding 1378 valuation makes true the clause C_j . Bearing in mind that a single copy of the 1379 object b_{2n+m+1} is in each cell, the phase takes *nm* computation steps. 1380

Thus, the configuration $C_{2n+m+nm+1}$ is characterized by the following:

- It contains exactly 2^n cells labeled by 2. Each one contains the object 1382 $d_{2n+m+nm+1}$, and copies of objects r_j for each clause C_j are made true by 1383 the encoded valuation in the cell. 1384
- It contains a unique cell labeled by 1, containing a copy of objects 1385 yes, no, f, g and the counter $e_{2n+m+nm+2}$. 1386

This phase consumes m computation steps.

- Formula checking phase: in this phase it is determined if there exists any 1388 valuation making true the *m* clauses of the formula. For this, the rules of 1389 type (*q*) are used, analyzing in an ordered way (first the clause C_1 , after that 1390 clause C_2 , and so on) if the clauses of the formula are being satisfied by the 1391 represented valuation in the corresponding cell labeled by 2. For example, from 1392 counter $d_{2n+m+nm+1}$ appearing in every cell 2, the appearance of the object r_1 1393 (the valuation makes true clause C_1) permits to generate in that cell the object 1394 $d_{2n+m+nm+2}$. This object, in turn, permits to evolve object $d_{2n+m+nm+3}$ if in that 1395 cell appears the object r_2 . In this manner, a valuation represented by a cell labeled 1396 by 2 makes true the formula φ if and only if the object $d_{2n+m+nm+n+1}$ appears 1397 in the content of that cell in the configuration $C_{2n+m+nm+m+1}$.
- *Output phase*: in this phase, the system will provide the corresponding output, 1399 depending on the analysis in the formula checking phase. 1400

If the formula φ is satisfiable, then there is some cell in the configuration 1401 $C_{2n+m+nm+m+1}$ that contains an object $d_{2n+m+nm+m+1}$. In this case, the applica-1402 tion of rule (r) sends an object f and the object yes to the cell 1. The object yes 1403 therefore disappears from cell 1, and consequently, rule (t) cannot be applied. In 1404 the next computation step, the application of the rule (s) produces an object yes 1405 in the environment (for the first time during the whole computation) and the 1406 process ends. 1407

If the formula φ is not satisfiable, then there no exist any cell in the 1408 configuration $C_{2n+m+nm+m+1}$ containing an object $d_{2n+m+nm+m+1}$. In this case, 1409 the rule (r) is not applicable, and in the next computation step, the counter e_i 1410 evolves, providing an object $e_{2n+m+nm+m+3}$ in cell 1. This object permits the 1411 application of rule (t), since the object no and f remain in cell 1. In this way, 1412 the object no is sent in the next computation step, and the computation finalizes. 1413

6.4.3.3 Sequential Simulation and Data Structure

For an easier implementation, the simulation algorithm has been divided into five 1415 (simulation) phases, instead of the six phases in $\Pi_{tsp-SAT}$ since we merge some 1416 of them. Each of these simulation phases are implemented in code as separated 1417 functions whenever is possible. They corresponds to the application of certain rules, 1418 as explained below: 1419

• *Generation phase*: it performs the application of rules from (*a*) to (*e*) of systems 1420 from $\Pi_{tsp-SAT}$. Therefore, it comprises the two first phases of the theoretical 1421 model: valuations generation phase and counters generation phase. 1422

1381

1387

- *Exchange phase*: it simulates the application of rules (f) and (g). It comprises 1423 the first part of the checking preparation phase. 1424
- *Synchronization phase*: it applies the rules from (*h*) to (*m*), so comprising the 1425 second part of the checking preparation phase. 1426
- *Checking phase*: it performs the application of rules from (*n*) to (*p*). Thus, it is 1427 the checking clauses phase we identified in the theoretical model. 1428
- *Output phase*: it applies rules from (q) to (t). It then performs both the formula the phase identified in the theoretical model.

The sequential simulator implements these five simulation phases directly in 1431 code. The input of the simulator is the same than the one used in the simulator for 1432 the cell-like solution Π_{am-SAT} . A DIMACS CNF file is provided, and the simulator 1433 outputs the response of the computation. Therefore, it acts merely as a SAT solver, 1434 but the implementation follows the computation of the systems from the family 1435 $\Pi_{tsp-SAT}$.

Furthermore, we have adopted a set of optimizations to improve the performance 1437 of the sequential simulator. After several tests, we show that the best optimizations 1438 are as follows [21]: 1439

- As the exchange phase is very simple, it is then implemented after the generation 1440 phase loop, within the same function. 1441
- The full synchronization phase is applied to one cell before going to the next one. 1442 This allows to exploit data locality in cache memories. 1443
- In the checking phase, the objects r_j , for $1 \le j \le m$, are inserted in order in 1444 the corresponding array whenever they are created. Thus, the output phase can 1445 be easily performed, in such a way that it is not necessary to loop all the objects 1446 coming from the input multiset (literals). Now it is enough to check if there exists 1447 the *m* objects r_j . 1448

For this solution, the memory layout for the representation of the tissue P 1449 system differentiates between cells labeled 1 and 2, having a different data structure 1450 representing each type of cell in the system. 1451

First, cell 1 is represented as an array having a maximum dimension of five 1452 elements. That is, the multiset for cell 1 has the maximum amount of five objects. 1453 These five objects are the three counters, *b*, *c*, and *d* (which are initially in this cell), 1454 and the two objects *yes* and *no* (that will final answer to the problem). Note that 1455 the size of the array for cell 1 is always constant, as it is independent of the input 1456 parameters of the simulator. 1457

Second, the cells labeled by 2 are also represented by a one-dimensional array. 1458 All of them are stored inside this large array, since it is initially allocated to store 1459 the maximum amount of cells (2^n) . By studying a computation of the systems $\Pi_{\text{tsn-SAT}}$, we conclude that the maximum number of objects appearing in a cell 1460 2 is $(2n) + 4 + |cod(\varphi)|$, where: 1461

- $|cod(\varphi)|$ elements for the initial multiset,
- *n* elements for objects $T_{i,j}$ and $F_{i,j}$, for $1 \le i \le n$ and $1 \le j \le m$. Note that an 1463 object $T_{i,i}$ and an object $F_{i,i}$, for any *i*, cannot be simultaneously placed within 1464 a cell 2. Moreover, the index j is used sequentially in the computation steps of 1465 the system, i.e., replacing objects in the evolution process of incrementing the 1466 second index. For all of this, *n* elements are enough to store those objects. 1467
- *n* elements for objects t_i and f_i , for $1 \le i \le n$. Note that objects f_i and t_i , for 1468 i = j, cannot be simultaneously placed within a cell 2, so *n* elements are enough 1469 to store those objects. 1470
- 4 elements for counter objects a, b, c, and d. They will be replaced for counter 1471 objects f and g. 1472

The objects are represented similarly to the simulator for Π_{am-SAT} . In this 1473 case, we recover the reserved space utilized to store the multiplicity of the object, 1474 inasmuch as it exceeds 1. In summary, they are encoded at bit-level within integers 1475 of 32 bits that store the following (8 bits for each field): 1476

- 1. The name of the object (x or \overline{x}) 1477 2. Multiplicity of the object. As there are objects whose multiplicity can exceed 2^8 , 1478 this field can eventually be joined to the next one (variable). 1479 3. Variable (subindex *i*). 1480
- 4. Clause (subindex *j*).

6.4.3.4 Design of the Parallel Simulator

The design of this parallel simulator is driven by the same structure of phases we 1483 have used for the sequential one. Separated CUDA kernels are utilized to speedup 1484 the execution of each phase. 1485

The general assignment of work for threads and thread blocks is summarized in 1486 Fig. 6.11. Each thread block corresponds to each cell labeled by 2 created in the 1487 system. However, unlike the previous simulator for the cell-like solution, we do not 1488 assign a thread per literal. The assignment of each thread, this time, is different for 1489 each simulation phase. The work mapping per phase is therefore as follows: 1490

- Generation phase: the number of thread blocks is iteratively increased together 1491 with the amount of cells created in each computation step. We distribute cells 1492 along the two-dimensional grid through successive kernel calls. Each thread 1493 block contains $(2n) + 4 + |cod(\varphi)|$ threads. That is, the amount of elements 1494 assigned to each cell in the global array storing multisets. Threads are then used 1495 to copy each individual elements of the corresponding cell when it is divided. 1496
- *Exchange phase*: it is executed at the kernel for generation phase, using the 1497 same amount of thread blocks, but only the corresponding threads perform the 1498 exchange. 1499

1462

1481



Fig. 6.11 General design of the parallel simulator for $\Pi_{tsp-SAT}$. From [21, 26]

- Synchronization phase: the thread blocks are assigned to the cells labeled by 15002, like in the last step of the generation phase. For this phase, the number of 1501threads is *n* (number of variables). If we use the same amount of threads than 1502in generation phase, most of them will be idle. So it is preferred to launch less 1503threads, but performing effective work. We have experimentally corroborated this 1504fact. 1505
- *Checking phase*: the number of thread blocks is again assigned to be the number 1506 of cells labeled by 2. However, for this phase we use a block size of $|cod(\varphi)|$. 1507 That is, each thread is used to execute, in parallel, rules of type (n) and 1508 (o). The result at the SAT problem resolution level, each thread checks if the 1509 corresponding literal makes true its clause, depending on the truth assignment 1510 encoded by the cell assigned to the thread block. 1511
- *Output phase*: rules of type (q) are sequentially executed in a separate kernel, 1512 again using $|cod(\varphi)|$ threads per block and 2^n thread blocks $(2^n$ is the number of 1513 cells labeled by 2). 1514

For this solution, we have applied a small set of optimizations, focused on 1515 the GPU implementation, to improve the performance of the parallel simulator. 1516 We identify that the simulator runs twice faster than the simulator without these 1517 optimizations. We will use the optimized version of the parallel simulator to perform 1518 the comparisons. These optimizations are oriented to improve two performance 1519 aspects of GPU computing, what leads us to consider two kind of optimizations. 1520 The first one is to *emphasize the parallelism*. This optimization aims to increase 1521 the number of threads per block (to the recommended amount from 64 to 256), 1522 so it allows to fulfill warps and hide latency. The second is to *exploit streaming* 1523 *bandwidth*. To do this, the data is loaded first to the shared memory, and operated 1524 there, avoiding global memory (expensive) accesses. Next, we show the specific 1525 optimizations we have carried out for each phase: 1526

- *Generation phase*: no optimizations were implemented here, since the implementation already satisfies the first optimization type. The second type will require a more sophisticated implementation, like the one presented in Sect. 6.4.2.
- *Exchange phase*: this phase, as it is joined with the generation phase, has no 1530 optimizations. 1531
- Synchronization phase: the two optimization types are implemented here. The 1532 second optimization type is carried out by using shared memory to avoid global 1533 memory accesses. The first type is performed by increasing the number of threads 1534 per block. For the simulator, we can assume that *n* (number of variables, and the 1535 number of threads per block) is a small number, since the number of cells grows 1536 exponentially with respect to it. For example, let n = 32. Then, 2^{32} cells will be 1537 created, what require $2^{32}(68 + |cod(\varphi)|)$ bytes (in gigabytes: $272 + 4|cod(\varphi)|$). 1538 This number obviously exceeds the amount of available device memory. We 1539 therefore need to increase the number of threads per block, since n < 32 means 1540 to not fulfil a CUDA warp. A solution here is to assign more than one cell to 1541 each thread block. This amount is $\frac{256}{n}$, being 256 the optimum number of threads 1542 per block. It allows us to reach a number of threads close to the optimum one. 1543 However, we have to take care also of having enough shared memory to load the 1544 data of every assigned cell.
- *Checking phase*: since $|cod(\varphi)|$ can be greater than 32, we then keep this number 1546 as the number of threads per block. However, we use shared memory to speedup 1547 the accesses to the elements of the array. 1548
- *Output phase*: as in the previous phase, we also use shared memory, and the 1549 number of threads per block is kept to $|cod(\varphi)|$. 1550

6.4.3.5 Performance Analysis

In this subsection, we analyze the performance of the two simulators developed for 1552 the family of tissue-like P systems $\Pi_{tsp-SAT}$: the sequential simulator developed in 1553 C++ (from now, *tsp-sat-seq*) and the parallel simulator on the GPU (*tsp-sat-gpu*). 1554

Figure 6.12 shows the results for both simulators when increasing the number of 1555 cells (by increasing the number of variables in the input CNF formulas), considering 1556 only kernel runtime for *tsp-sat-gpu*. For this case, we can observe that again the 1557 kernels of *tsp-sat-gpu* run faster than *tsp-sat-seq*. However, the performance gain is 1558 increased with the amount of cell 2 created by the system. For 64 membranes, the 1559 speedup is of $2\times$, but for 2 M cells it is of $8.3\times$.



Fig. 6.12 Simulation performance for *tsp-sat-seq* and *tsp-sat-gpu* when increasing the number of membranes (x-axis). From [21, 26]



Fig. 6.13 Speedup achieved running Test 2 (256 Objects/Cell) for *tsp-sat-gpu* and *tsp-sat-seq* considering also the GPU data management, when increasing the number of membranes (x-axis). From [21, 26]

Finally, we show the speedup achieved by the simulator *tsp-sat-gpu*, taking into 1561 account also the amount of time consumed by the data management (allocation and 1562 transfer). It is observed that, since the data management time is fixed for all the sizes 1563 (copy the initial multiset and retrieve the final answer), the speedup exceeds 1 only 1564 after 128 K membranes. Systems with smaller number of cells are executed slower 1565 than in the CPU, because of the data management. However, for very large systems, 1566 the speedup is as large as with only kernels. The maximum speedup we report for 1567 this simulator is given for 4 M cells, up to $10 \times (Fig. 6.13)$.

AQ1

6.5 Adaptive Simulations

In this section, we will introduce a third type of simulation of P systems, which is 1570 called adaptive simulation. 1571

6.5.1 Definition

We have discussed the difference between generic and specific simulators. In 1573 this section, we will discuss a hybrid type, which is called adaptative, or simply 1574 adaptive, simulation. A simulator of this kind is initially a generic simulator, which 1575 is designed to simulate a wide range of P systems within a variant. However, the 1576 simulator is provided with high-level information that can be either discarded (then 1577 remaining as generic) or used to adapt the simulation to improve its efficiency. 1578

In this sense, an adaptative simulator has the goal of getting closer to specific 1579 simulators without losing generality; that is, they are generic simulators with 1580 improved performance by taking advantage of extra information provided directly 1581 by designers (e.g., modules). For example, if the algorithm scheme of the computation is known by the designer (as it is, as discussed for the specific simulators), then 1583 it can be given to the simulator in order to be able to discard rules at selection stage 1584 (because the algorithm scheme is known). 1585

Next, we will overview the first adaptive simulator for P systems implemented 1586 so far, which is published in [29]. This simulation framework is implemented 1587 within *ABCD-GPU*, and it can be downloaded from the official website http:// 1588 sourceforge.net/p/pmcgpu [45] or the repository https://github.com/RGNC/abcd-1589 gpu/tree/adaptative. 1590

6.5.2 Simulating Population Dynamics P Systems

The idea of adaptative simulators was introduced and analyzed in [29]. It is inspired 1592 in the way directives work in common programming languages. They are special 1593 syntactic elements that tell extra information to the compiler, allowing to better 1594 adapt the code for some purposes if the compiler accepts it (e.g., in OpenMP, one 1595 call can easily ask to parallelize the iterations of a loop). This way, a P system model 1596 designer can also provide very useful information to the simulator, rather than just 1597 the syntactic and/or semantic elements of the P system to simulate, such as the algorithmic scheme of the computation. 1599

Specifically in PDP systems, ecosystem modelers often use algorithmic schemes 1600 for their models [6]. This is given as cycle that is repeated (per year, per season, etc.). 1601 A cycle in the model is a fixed amount of transition steps where a sequence of modules take place. These modules reproduce certain processes such as reproduction of 1603 species, feeding, migration, etc. Moreover, these modules consist of certain rules 1604 that are carefully designed to model the corresponding process. Therefore, we can 1605

1569

1572

say that somehow the model designer already knows which rules can be executed in 1606 each time step. Thus, if they are able to provide that information, the simulator can 1607 take advantage of this to dismiss rules automatically at each step. 1608

The PDP system simulator was turned into adaptative. First, the model designer is 1609 able to provide the information of the modules they are defining by using the new P-Lingua 5 software [37]. This new version now includes new syntax elements called 1611 *features*. They are written as @featureName = featureValue and can be 1612 defined globally (for the whole system) or locally (for individual rules). ABCD-1613 GPU takes this information to organize the rules by modules. If the simulator does 1614 not recognize the information provided by the features, it can proceed and simulate 1615 the system without problems. 1616

In summary, there are two main pieces of information that has to be declared in 1617 order to define modules: 1618

- 1. Information about the modular structure of the model. This includes module 1619 names and their temporal relation. The latter indicates when a module starts 1620 inside a cycle and which modules will follow a given one. 1621
- 2. Information about distribution of rules in modules. That is, which module each 1622 rule belongs to. 1623

The simulator precomputes which modules are active in each step within the 1624 cycle before starting the simulation. In this way, this information can be used to 1625 easily identify the rules that might be applicable at each transition step. For this 1626 purpose, the rule blocks and the rules are sorted in order to compact them into 1627 modules; rules belonging to the same module are put one after the other. The kernels 1628 of ABCD-GPU are expanded to accept extra indexes indicating the modules and 1629 where the rules of the modules are. In this way, the threads as distributed in Fig. 6.4 1630 will have a shorter loop, because the rule blocks (and rules) are just those from 1631 the module being active. Furthermore, if the solution has parallel modules in a 1632 cycle, then they can also run in parallel thanks to CUDA streams. We can launch 1633 the kernels for phase 1 also in parallel at different streams, one per module. As for 1634 environments and simulations, the behavior remains as before.

6.5.2.1 Analysis of Performance Results

Next, the behavior and performance of the adaptative PDP system simulators for 1637 GPU and OpenMP are analyzed. The model employed as benchmark is based on 1638 the tritrophic interactions presented in [9, 10]. This is a virtual ecosystem that 1639 was defined to illustrate PDP systems as a modeling framework. In this model, 1640 three trophic levels are represented: grass, herbivores, and carnivores. These species 1641 interact with each other, reproduce, and move along the 10 environments when no 1642 food is encountered. Rule block competitions take place. For instance, all herbivores 1643 compete for grass that is represented by a single object, G.

For benchmarking purposes, the model has been generalized so that the number 1645 of species can be changed. The corresponding parameters (probabilities, amount of 1646 copies eaten per species, etc.) are generated randomly. This was possible thanks to 1647

the ability of P-Lingua 5 to incorporate calls from the model to random number 1648 generation functions. Moreover, the modules of the model are identified by P- 1649 Lingua 5 features. 1650

In this section, the benchmark carried out to the adaptative PDP systems 1651 simulator is analyzed. The two versions of the simulator are compared: generic and 1652 adaptive versions of ABCD-GPU. The extended tritrophic model is used as input. In 1653 all experiments, 20 years of the virtual ecosystem are simulated (corresponding to 1654 180 transition steps of the PDP systems). The A parameter of DCBA is set to 2. No 1655 output was asked, so only the simulation runtimes were measured. The scalability 1656 of the simulators is analyzed by increasing the number of species. Specifically, 7 1657 will be used to denote the *base* model, which has in fact 7 species. In order to 1658 have an idea of the dimensions of the model, the ratio of rule blocks per species is 1659 approximately 22: 21985 rule blocks are generated for 1000 species, being 9990 1660 communication rule blocks and 11,995 skeleton rule blocks. Another parameter 1661 affecting scalability is the amount of simulations running in parallel. For this reason, 1662 50 simulations were launched for the tests. The following two configurations of CPU 1663 and GPU hardware were used to run the simulations (short names are provided in 1664 bold): 1665

- (i7) Intel i7-8700 CPU at 3.20 GHz, having 12 logical cores (6 physical) 1666
- (P100) Tesla P100 GPU, having 3584 cores at 1.33 GHz

A cross comparison of runtimes and speedups achieved by GPU compared to 1668 CPU is shown in Fig. 6.14, which corresponds to the speedups reached by the above 1669 simulation times. The GPU is faster, in both *adaptive* and *generic* versions, than 1670 the multicore counterparts when handling middle and large models. Only for the 1671



Fig. 6.14 Comparison of P100 versus i7 with 8 threads, for both generic and adaptive versions of *abcd-gpu* tested for different number of species in the model. It shows the corresponding speedups of P100 against i7 for both version. 50 simulations were run. Bar plots use logarithmic scale for y-axis. From [29]

small base model, the GPU is a bit slower (above $0.9 \times$). Speedups are higher with 1672 larger models, being around $30 \times$ and $50 \times$ for *adaptive* and *generic* simulators, 1673 respectively, and for 2000 species. When simulating hundreds of species, $6 \times$ and 1674 $10 \times$ accelerations were obtained for *adaptive* and *generic* versions. Finally, the 1675 speedup of the GPU is lower when using the *adaptive* version, given that the impact 1676 of the *modular* scheme is better for the CPU than for the GPU. 1677

We can conclude that this design helped to improve the performance by $2.5 \times$ 1678 extra when using a P100 GPU [29].

6.6 Conclusions

GPUs have been established as a massively parallel processor and an enabling 1681 technology where programmers currently accelerate scientific applications. They 1682 provide a good parallel platform to simulate P systems due to the double parallel 1683 nature that both GPUs and P systems present. Their shared memory system also 1684 helps to efficiently synchronize the simulation of the models. Moreover, they are a 1685 cheap and scalable parallel architecture that can be seen in current HPC solutions. 1686

However, the results in the literature [22, 48] show that P systems simulations 1687 are memory bandwidth bound: they spend more time accessing and updating data 1688 (multisets) than executing computation. The main cause is that simulating P systems 1689 requires a high synchronization degree (e.g., the global clock of the models, rule 1690 cooperation, rules competition, etc.), and the number of operations to execute per 1691 memory access is very small (P systems execute rewriting rules). This restricts the 1692 design of parallel simulators. A parallel simulator designer has to be careful with 1693 the representation and management of each P system ingredient. A bad step taken 1694 on GPU programming can easily break parallelism and, so, performance. 1695

We can identify a taxonomy of simulators developed so far. Generic parallel 1696 simulators are intended to be flexible enough to simulate a wide range of P systems 1697 within a variant. They also take advantage of P systems parallelism to speedup the 1698 simulation. However, when working with highly flexible simulators, the P systems 1699 design has to be reconsidered to achieve performance, in such a way that they 1700 execute as many rules as possible in each computation step. Some variants simulated 1701 by generic simulators are P systems with active membranes and elementary division 1702 and Population Dynamic P systems. Other related works also include spiking neural 1703 P systems variants [1, 2].

On the other hand, specific simulators are designed for just certain P systems 1705 within a solution or family. This way, the simulator can be designed adapting 1706 all parts to the P systems, since their scheme is known at developing time. The 1707 performance achieved in these simulators are much higher, but it comes at restricting 1708 the P systems to simulate. For example, one cannot define new rules to simulate, 1709 since they are already predetermined. In the middle term, we have a new type of 1710 simulation called adaptive. Basically, it is a generic simulator but that includes high 1711 level information that can be either discarded by the simulator (going generic) or 1712 used to adapt the simulation and achieve better performance (adaptive). 1713

We can identify some challenges for the future. For example, concerning the 1714 memory bandwidth limit, one challenge is to design P system variants where the 1715 model contains a higher computational intensity. Moreover, memory accesses can 1716 be partially reduced by improving data structures using a compacted, dense, and 1717 well-ordered memory representation of P systems. A challenge is to use a dense 1718 representation in an effective way in generic simulators. Finally, a P system model 1719 with cooperation in the LHS usually leads to this issue, making it more difficult 1720 when the cooperation is larger. 1721

References

AQ2

| 1. | J.P. Carandang, F.G.C. Cabarle, H.N. Adorna, N.H.S. Hernandez, M.A. Martínez-del-Amor, | 1723 |
|----|---|------|
| | Handling non-determinism in spiking neural P systems: algorithms and simulations. Fundam. | 1724 |
| | Inf. 164(2–3), 139–155 (2019). https://doi.org/10.3233/FI-2019-1759 | 1725 |
| 2. | J.P. Carandang, J.M.B. Villaflores, F.G.C. Cabarle, H.N. Adorna, M.A. Martínez-del-Amor, | 1726 |
| | CuSNP: spiking neural P systems simulators in CUDA. Rom. J Inf Sci Technol. 20(1), 57-70 | 1727 |
| | (2017) | 1728 |
| 3. | J.M. Cecilia, J.M. García, G.D. Guerrero, M.A. Martínez-del-Amor, I. Pérez-Hurtado, M.J. | 1729 |
| | Pérez-Jiménez, Simulating a P system based efficient solution to SAT by using GPUs. J Logic | 1730 |
| | Algebraic Program. 79(6), 317–325 (2010). https://doi.org/10.1016/j.jlap.2010.03.008 | 1731 |
| 4. | J.M. Cecilia, J.M. García, G.D. Guerrero, M.A. Martínez-del-Amor, I. Pérez-Hurtado, M.J. | 1732 |
| | Pérez-Jiménez, Simulation of P systems with active membranes on CUDA. Briefings in Bioinf. | 1733 |
| | 11(3), 313–322 (2010). https://doi.org/10.1093/bib/bbp064 | 1734 |
| 5. | J.M. Cecilia, J.M. García, G.D. Guerrero, M.A. Martínez-del-Amor, M.J. Pérez-Jiménez, M. | 1735 |
| | Ujaldón, The GPU on the simulation of cellular computing models. Soft Comput. 16(2), 231- | 1736 |
| | 246 (2012). https://doi.org/10.1007/s00500-011-0716-1 | 1737 |
| 6. | M.A. Colomer, A. Margalida, M.J. Pérez-Jiménez, Population Dynamics P system (PDP) | 1738 |
| | models: a standardized protocol for describing and applying novel bio-inspired computing | 1739 |
| | tools. PLOS ONE 8(5), e60698 (2013). https://doi.org/10.1371/journal.pone.0060698 | 1740 |
| 7. | M.A. Colomer, A. Margalida, D. Sanuy, M.J. Pérez-Jiménez, A bio-inspired computing model | 1741 |
| | as a new tool for modeling ecosystems: the avian scavengers as a case study. Ecol. Model. | 1742 |
| | 222(1), 33-47 (2011). https://doi.org/10.1016/j.ecolmodel.2010.09.012 | 1743 |
| 8. | M.A. Colomer, A. Margalida, L. Valencia, A. Palau, Application of a computational model for | 1744 |
| | complex fluvial ecosystems: the population dynamics of zebra mussel Dreissena polymorpha | 1745 |
| | as a case study. Ecol. Complexity 20, 116–126 (2014). https://doi.org/10.1016/j.ecocom.2014. | 1746 |
| | 09.006 | 1747 |
| 9. | M.A. Colomer, I. Pérez-Hurtado, M.J. Pérez-Jiménez, A. Riscos-Núñez. Comparing simulation | 1748 |
| | algorithms for multienvironment probabilistic P systems over a standard virtual ecosystem, in | 1749 |
| | IEEE Fifth International Conference on Bio-Inspired Computing: Theories and Applications | 1750 |
| | (BIC-TA 2010), vol. 1 (2010). https://doi.org/10.1109/BICTA.2010.5645258 | 1751 |
| 0. | M.A. Colomer, I. Pérez-Hurtado, M.J. Pérez-Jiménez, A. Riscos-Núñez, Comparing simulation | 1752 |
| | algorithms for multienvironment probabilistic P systems over a standard virtual ecosystem. | 1753 |
| | Nat. Comput. 11(3), 369–379 (2012). https://doi.org/10.1007/s11047-011-9289-2 | 1754 |
| 1. | T.S. Crow, Evolution of the graphical processing unit. Master's thesis, University of Nevada | 1755 |
| | Reno (2004). http://www.cse.unr.edu/~fredh/papers/thesis/023-crow/GPUFinal.pdf | 1756 |
| 2. | A.C. Elster, High-Performance Computing: past, present, and future, in Applied Parallel | 1757 |
| | Computing, Lecture Notes in Computer Science, ed. by J. Fagerholm, J. Haataja, J. Järvinen, | 1758 |
| | M. Lyly, P. Raback, V. Savolainen, vol. 2367 (2006), pp. 433-444. https://doi.org/10.1007/3- | 1759 |
| | 540-48051-X 43 | 1760 |

- M. García-Quismondo, R. Gutiérrez-Escudero, M.A. Martínez-del-Amor, E. Orejuela-Pinedo, 1761
 I. Pérez-Hurtado, P-Lingua 2.0: A software framework for cell-like P systems. Int. J. Comput. 1762
- Commun. Control **4**(3), 234–243 (2009). https://doi.org/10.15837/ijccc.2009.3.2431 1763 14. M.A. Gutiérrez-Naranjo, M.J. Pérez-Jiménez, A. Riscos-Núñez, A fast P system for finding a balanced 2-partition. Soft Comput. **9**, 673–678 (2005). https://doi.org/10.1007/s00500-004-0397-0 1766
- 15. Inside HPC blog. http://insidehpc.org
- B.W. Kernighan, D. Ritchie. *The C Programming Language*, 2nd edn. (Prentice Hall, Englewood Cliffs 1988)
- 17. D.B. Kirk, W.W. Hwu, Programming Massively Parallel Processors: A Hands-on Approach, 1770
 3rd edn. (Morgan Kaufmann, San Francisco, 2016). https://www.sciencedirect.com/science/ 1771
 book/9780128119860 1772
- A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional 1773 neural networks. Adv. Neural Inf. Process. Syst. 25(2) (2012). https://doi.org/10.1145/3065386 1774
- 19. E. Lindholm, J. Nickolls, S. Oberman, J. Montrym, NVIDIA Tesla: a unified graphics and 1775 computing architecture. IEEE Micro 28(2), 39–55 (2008). https://doi.org/10.1109/MM.2008. 1776 31
- 20. C. Martín-Vide, Gh. Păun, J. Pazos, A. Rodríguez-Patón, Tissue P systems. Theor. Comput. 1778
 Sci. 296(2), 295–326 (2003). https://doi.org/10.1016/S0304-3975(02)00659-X
 1779
- 21. M.A. Martínez-del-Amor, Accelerating Membrane Systems Simulators Using High Performance Computing with GPU. Ph.D. Thesis, Universidad de Sevilla, 2013. http://hdl.handle.
 1781
 net/11441/15644
 1782
- M.A. Martínez-del-Amor, M. García-Quismondo, L.F. Macías-Ramos, L. Valencia-Cabrera, A. 1783 Riscos-Núñez, M.J. Pérez-Jiménez, Simulating P systems on GPU devices: a survey. Fundam. 1784 Inf. 136(3), 269–284 (2015). https://doi.org/10.3233/FI-2015-1157 1785
- M.A. Martínez-del-Amor, I. Karlin, R.E. Jensen, M.J. Pérez-Jiménez, A.C. Elster, Parallel 1786 simulation of probabilistic P systems on multicore platforms, in *Proceedings of the Tenth* 1787 *Brainstorming Week on Membrane Computing*, ed. by M. García-Quismondo, L.F. Macías-Ramos, Gh. Păun, L. Valencia-Cabrera, vol. II (Fénix Editora, 2012), pp. 17–26 1789
- 24. M.A. Martínez-del-Amor, L.F. Macías-Ramos, L. Valencia-Cabrera, M.J. Pérez-Jiménez, 1790 Parallel simulation of Population Dynamics P systems: updates and roadmap. Nat. Comput. 1791 15(4), 565–573 (2015). https://doi.org/10.1007/s11047-016-9566-1 1792
- M.A. Martínez-del Amor, D. Orellana-Martín, I. Pérez-Hurtado, L. Valencia-Cabrera, A. 1793 Riscos-Núñez, M.J. Pérez-Jiménez, Design of specific P systems simulators on GPUs, in 1794 *Membrane Computing. CMC 2018*, ed. by T. Hinze, G. Rozenberg, A. Salomaa, C. Zandron. 1795 Lecture Notes in Computer Science, vol. 11399 (2019), pp. 202–207. https://doi.org/10.1007/ 1796 978-3-030-12797-8_14
- M.A. Martínez-del-Amor, J. Pérez-Carrasco, M.J. Pérez-Jiménez, Characterizing the parallel 1798 simulation of P systems on the GPU. Int. J. Unconv. Comput. 9(5–6), 405–424 (2013) 1799
- M.A. Martínez-del-Amor, I. Pérez-Hurtado, M. García-Quismondo, L.F. Macías-Ramos, L. 1800 Valencia-Cabrera, A. Romero-Jiménez, C. Graciani-Díaz, A. Riscos-Núñez., M.A. Colomer, 1801 M.J. Pérez-Jiménez, DCBA: simulating Population Dynamics P Systems with proportional 1802 object distribution, in *Membrane Computing. CMC 2012*, ed. by E. Csuhaj-Varjú, M. Gheorghe, G. Rozenberg, A. Salomaa, G. Vaszil. Lecture Notes in Computer Science, vol. 7762 1804 (2012), pp. 291–310. https://doi.org/10.1007/978-3-642-36751-9_18
- M.A. Martínez-del-Amor, I. Pérez-Hurtado, A. Gastalver-Rubio, A.C. Elster, M.J. Pérez-Jiménez, Population Dynamics P Systems on CUDA, in *Computational Methods in Systems* 1807 *Biology*, ed. by D. Gilbert, M. Heiner. Lecture Notes in Computer Science, vol. 7605 (2012), 1808 pp. 247–266. https://doi.org/10.1007/978-3-642-33636-2_15 1809
- 29. M.A. Martínez-del-Amor, I. Pérez-Hurtado, D. Orellana-Martín, M.J. Pérez-Jiménez, Adaptative parallel simulators for bioinspired computing models. Future Gener. Comput. Syst. 107, 1811 469–484 (2020). https://doi.org/10.1016/j.future.2020.02.012

| 30. | M.A. Martínez-del-Amor, I. Pérez-Hurtado, M.J. Pérez-Jiménez, A. Riscos-Núñez, M.A. | 1813 |
|-----|--|------|
| | Colomer, A new simulation algorithm for multienvironment probabilistic P systems, in 2010 | 1814 |
| | IEEE Fifth International Conference on Bio-Inspired Computing: Theories and Applications | 1815 |
| | (BIC-TA), Changsha, 2010, vol. 1 (2010), pp. 59–68. https://doi.org/10.1109/BICTA.2010. | 1816 |
| | 5645352 | 1817 |
| 31 | A Munshi BR Gaster TG Mattson I Fung D Ginsburg OpenCI Programming Guide | 1818 |
| 51. | A. Mulishi, D.K. Gastel, 1.O. Matison, J. Fung, D. Ghisourg, OpenCL Frogramming Guide, 1st edn. (Addison Weeley, Reading, 2011) | 1010 |
| 22 | I Nickella I Buck M Carland K Skedron Scalable nerallal programming with CUDA, is | 1019 |
| 52. | J. Mickons, I. Buck, M. Garland, K. Skaufon, Scalable parallel programming with CODA. Is | 1820 |
| | CUDA the parallel programming model that application developers have been waiting for? | 1821 |
| | Queue $6(2)$, 40–53 (2008). https://doi.org/10.1145/1365490.1365500 | 1822 |
| 33. | NVIDIA CUDA C Programming Guide. https://docs.nvidia.com/cuda/cuda-c-programming- | 1823 |
| | guide/index.html. Accessed June 2019 | 1824 |
| 34. | N. Otterness, J. Anderson, AMD GPUs as an alternative to NVIDIA for supporting real-time | 1825 |
| | workloads, in Proceedings of the 32nd Euromicro Conference on Real-Time Systems, (2020), | 1826 |
| | pp. 10:1–10:23. https://doi.org/10.4230/LIPIcs.ECRTS.2020.10 | 1827 |
| 35. | J.D. Owens, M. Houston, D. Luebke, S. Green, J.E. Stone, J.C. Phillips, GPU computing. Proc. | 1828 |
| | IEEE 96(5), 879–899 (2008). https://doi.org/10.1109/JPROC.2008.917757 | 1829 |
| 36. | Gh. Păun, P systems with active membranes: attacking NP-complete problems. J. Autom. Lang. | 1830 |
| | Comb. 6, 75–90 (1999) | 1831 |
| 37 | I Pérez-Hurtado D Orellana-Martín G Zhang M I Pérez-Iiménez P-lingua in two steps: | 1832 |
| 57. | flexibility and efficiency I Membr Comput 1(2) 93-102 (2019) https://doi.org/10.1007/ | 1833 |
| | c/1965_010_00014_1 | 1000 |
| 20 | MI Dávaz limánaz A Disass Núñaz Salving the Subast Sum problem by D systems | 1034 |
| 50. | with active membranes. N. Canon Commut. 22 (4), 220, 256 (2005), https://doi.org/10.1007/ | 1835 |
| | with active memoranes. N. Gener. Comput. $23(4)$, $339-336$ (2003). https://doi.org/10.100// | 1836 |
| 20 | BF0303/03/ | 1837 |
| 39. | M.J. Perez-Jimenez, A. Riscos-Nunez, A linear-time solution for the knapsack problem with | 1838 |
| | active membranes, in Membrane Computing. WMC 2003, ed. by C. Martín-Vide, G. Mauri, | 1839 |
| | Gh. Păun, G. Rozenberg, A. Salomaa. Lecture Notes in Computer Science, vol. 2933 (2004), | 1840 |
| | pp. 250–268. https://doi.org/10.1007/b95207 | 1841 |
| 40. | M.J. Pérez-Jiménez, A. Romero-Jiménez, F. Sancho-Caparrini, Complexity classes in models | 1842 |
| | of cellular computing with membranes. Nat. Comput. 2(3), 265–285 (2003). https://doi.org/10. | 1843 |
| | 1023/A:1025449224520 | 1844 |
| 41. | M.J. Pérez-Jiménez, A. Romero-Jiménez, F. Sancho-Caparrini, Decision P systems and the | 1845 |
| | $\mathbf{P} \neq \mathbf{NP}$ conjecture, in <i>Membrane Computing</i> . WMC 2002, ed. by Gh. Păun, G. Rozenberg, | 1846 |
| | A. Salomaa, C. Zandron. Lecture Notes in Computer Science, vol. 2597 (2003), pp. 388–399. | 1847 |
| | https://doi.org/10.1007/3-540-36490-0_27 | 1848 |
| 42 | M I Pérez-Iiménez A Romero-Iiménez E Sancho-Caparrini A polynomial complexity class | 1849 |
| 12. | in P systems using membrane division I Autom Lang Comb 11 423–434 (2006) | 1850 |
| 43 | N Satish M Harris M Garland Designing efficient sorting algorithms for manycore GPUs in | 1951 |
| 45. | Proceedings of the 2000 IEEE International Symposium on Parallel & Distributed Processing | 1051 |
| | (IDDBS '00) Demo 2000 (IEEE Computer Society, Silver Series, 2000), nr. 1, 10, https://doi. | 1052 |
| | (<i>IPDPS 09</i>), Rome, 2009 (IEEE Computer Society, Silver Spring, 2009), pp. 1–10. https://doi. | 1853 |
| | org/10.1109/1PDPS.2009.5101005 | 1854 |
| 44. | W. Shin, K.H. Yoo, N. Baek, Large-scale data computing performance comparisons on SYCL | 1855 |
| | heterogeneous parallel processing layer implementations. Appl. Sci. 10, 1656 (2020). https:// | 1856 |
| | doi.org/10.3390/app10051656 | 1857 |
| 45. | The PMCGPU (Parallel simulators for Membrane Computing on the GPU) project website. | 1858 |
| | http://sourceforge.net/p/pmcgpu. Accessed June 2019 | 1859 |
| 46. | The top 500 supercomputer site. http://www.top500.org | 1860 |
| 47. | A. Torres-Moríño, M.A. Martínez-del-Amor, F. Sancho-Caparrini, GPU-parallel crowd sim- | 1861 |
| | ulation with Vulkan, in Proceedings of the 18th High Performance Computing & Simulation | 1862 |
| | (<i>HPCS 2020</i>). (2020, in press) | 1863 |
| 48. | G. Zhang, Z. Shang, S. Verlan, M.A, Martínez-del-Amor, C. Yuan, L. Valencia-Cabrera, M.J. | 1864 |
| | Pérez-Jiménez. An overview of hardware implementation of Membrane Computing models | 1865 |
| | ACM Comput. Surv. 53(4), 90 (2020), https://doi.org/10.1145/3402456 | 1866 |
| | · · · · · · · · · · · · · · · · · · · | |

AUTHOR QUERIES

- AQ1. Missing citation for "Fig. 6.13" was inserted. Please check if appropriate. Otherwise, please provide citation for "Fig. 6.13". Note that the order of main citations of figures in the text must be sequential.
- AQ2. Please note that the references in the list have been arranged alphabetically and corresponding citations have been changed accordingly. Please check.

uncorrected