

Swarm Model Checking on the GPU

Richard DeFrancisco¹, Shengsun Cho¹,
Michael Ferdman¹, and Scott A. Smolka¹

Stony Brook University, Stony Brook NY 11794-2424, USA

Abstract. We present *Grapple*, a new and powerful framework for explicit-state model checking on GPUs. Grapple is based on *swarm verification* (SV), a model-checking technique wherein a collection or *swarm* of small, memory- and time-bounded *verification tests* (VTs) are run in parallel to perform state-space exploration. SV achieves high state-space coverage via diversification of the search strategies used by constituent VTs. Grapple represents a swarm implementation for the GPU. In particular, it runs a parallel swarm of internally-parallel VTs, which are implemented in a manner that specifically targets the GPU architecture and the SIMD parallelism its computing cores offer. Grapple also makes effective use of the GPU shared memory, eliminating costly inter-block communication overhead. We conducted a comprehensive performance analysis of Grapple focused on the various design parameters, including the size of the queue structure, implementation of guard statements, and non-deterministic exploration order. Tests are run with multiple hardware configurations, including on the Amazon cloud. Our results show that Grapple performs favorably compared to the SPIN swarm and a prior non-swarm GPU implementation. Although a recently debuted FPGA swarm is faster, the deployment process to the FPGA is much more complex than Grapple’s.

Keywords: GPU · model checking · swarm verification · Grapple.

1 Introduction

Modern computing exists in a space that is increasingly parallel, distributed, and heterogeneous. High-performance co-processors such as GPUs (Graphics Processing Units) are utilized in many super-computing applications due to their high computational throughput, energy efficiency, and low cost [5]. GPGPU (General-Purpose Computing on a GPU) is achieved through the use of GPU programming languages such as the Open Computing Language (OpenCL) [6] and the Compute Unified Device Architecture (CUDA) [1].

In 2014, we adapted the multicore SPIN algorithm of [24] to the GPU [12]. While our approach achieved speedups up to 7.26x over traditional SPIN, and 1.26x over multicore SPIN, it was severely limited by the memory footprint of the GPU, and by an explicit limit on state-vector size set by the hash function [8].

The introduction of *Swarm Verification* (SV) in [27] represented an entirely new approach to parallel MC. In SV, a large number of MC instances are executed in parallel, each with a restricted memory footprint and a different search

path. Each instance is called a *verification test* (VT), because it does not seek to cover the full state space as a model checker would. Through the use of *diversification techniques*, VTs are independent of one other in terms of the portions of the model’s state space they cover. By executing a sufficiently large number of VTs, one is therefore statistically guaranteed to achieve nearly complete, if not complete coverage of the entire state space.

In this paper, we present *Grapple*, bringing the light-weight yet powerful nature of SV to the massively parallel GPU architecture. While other swarm implementations run internally sequential VTs in parallel, Grapple VTs are internally parallel and evolved from our previous GPU-based MC design [12]. Each VT runs on a single block of the GPU, with a bitstate hash table in shared memory, compacting per-state storage *by a factor of 64* compared to the cuckoo tables used in [12]. These tables use the hash function of [29], eliminating the hard 64-bit state vector limit of our previous model checker.

Grapple VTs run in parallel on all available GPU streaming multiprocessors (SMs), and make efficient use of the GPU scheduler to quickly replace jobs the instant an SM becomes available. As VTs are independent of each other and each one is tightly bound to a single chip on hardware, there is no need for inter-block communication or additional synchronization primitives.

To assess Grapple’s performance, we used a benchmark specifically designed for SV-based model checkers [27, 16]: a model that can randomly generate more than 4 billion states. Exploration progress in the benchmark is captured by the visitation of 100 randomly distributed states, or *waypoints*, with 100 waypoints approaching complete state-space exploration. Our experiments, which we ran on multiple hardware configurations, including the Amazon cloud [2], evaluate the impact of variations in queue size, guard-statement implementation, and nondeterministic exploration order.

We also compared Grapple’s performance with the FPGA swarm implementation of [16], the CPU swarm of [27], and the original (non-swarm) GPU implementation of [12]. Grapple easily outperforms the GPU implementation and the CPU swarm, and reaches all waypoints in a number of VTs comparable to that required by the FPGA implementation. While it cannot compete in raw speed with the hardware-level FPGA implementation, it offers much easier deployment, with VTs that complete in under a second. We additionally evaluated Grapple using multiple configurations of the Dining Philosophers problem, a small model with a known state-space size and deadlock violation.

In summary, our main contributions are as follows. (i) We introduce Grapple, a GPU-based swarm verification model checker with internally parallel verification tasks. (ii) We analyze structural elements of VTs (e.g., search strategy, queue size, guard logic, number of threads per VT) to determine how they impact the rate of exploration. (iii) We compare Grapple’s performance to previous SV implementations on the CPU [27] and FPGA [16], as well as to our non-swarm GPU-based model checker [12].

The rest of the paper is organized as follows. Section 2 provides background on GPU hardware, the CUDA programming model, the SPIN model checker,

and swarm verification. Section 3 presents our Grapple model checker. Section 4 presents our various experimental results. Section 5 considers related work. Section 6 interprets our findings and offers directions for future work.

2 Background

To motivate our design decisions for Grapple, we first explain the intricacies of GPU hardware and the associated CUDA programming model, and provide an overview of the SPIN model checker [7], on which Grapple is based. Further details on the GPU hardware and CUDA are available in the CUDA C Programming Guide [4].

2.1 GPU Hardware Model

The GPU is a high-performance co-processor designed to efficiently render 3D graphics in real time. GPUs are well-suited for linear algebra, matrix arithmetic, and other computations frequently used in graphical applications. As illustrated in Figure 1, the GPU architecture consists of a scalable array of N multithreaded *Streaming Multiprocessors* (SMs), each of which is made up of M *Stream Processor* (SP) cores. Each core is equipped with a fully pipelined integer-arithmetic logic unit (ALU) and a floating-point unit (FPU) that execute one integer or floating-point instruction per clock cycle. Each SM controls a *warp* of 32 threads, executing the same instructions in lock-step for all threads.

The GPU features a number of memory types, differing in access speed, capacity, and read/write availability. Global memory is large (order of gigabytes), available device-wide, but relatively slow. Constant memory is a cached, read-only memory intended for storing constant values that are not updated during execution. Finally, each SM has a shared memory region (16–48 KB). In practice, accessing shared memory can be *up to 100 times faster* than using global memory for the same transaction.

Devices connect to the host machine using the the PCIe bus. Communication between the host and device are extremely costly compared to on-board memory accesses, including those that use global memory.

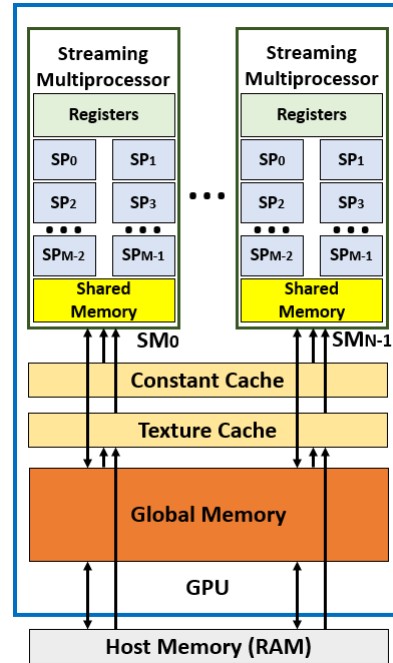


Fig. 1: GPU hardware model. SP = Stream Processor

2.2 CUDA Programming Model

CUDA is the proprietary NVIDIA programming model for general-purpose computing on their GPU architecture. While the alternative model, OpenCL, is universally compatible with all GPU architectures, the high-performance of CUDA has led to wide adoption. We decided to write Grapple in CUDA for this reason, but an OpenCL implementation would be very similar.

The CUDA parallel computing model uses tens of thousands of lightweight *threads* assembled into one- to three-dimensional *thread blocks*. A thread executes a function called a *kernel*, which contains the computations to be run in parallel. Each thread uses different parameters. Threads located in the same thread block can work together in several ways. They can insert a synchronization point into the kernel, which requires all threads in the block to reach that point before execution can continue. They can also share data during execution. In contrast, threads located in different thread blocks cannot communicate in such ways and essentially operate independently.

Shared-memory transactions are typically parallel to some number n distinct banks, but if two or more address requests fall in the same bank, the collision causes a serialization of the access. It is therefore important to understand addressing patterns when utilizing shared memory. Register management is also critically important. Use of registers is partitioned among all threads and, as such, using a large number of registers within a CUDA kernel will limit the number of threads that can run concurrently. Double and long long variables, use of shared memory, and unoptimized block/warp geometry all lead to increased register use. If available registers are exhausted, the contents will spill over into local memory- a special type of device memory with the same high-latency and low-bandwidth as global memory.

The SIMD nature of warps on SPs has a great impact on code structure for the GPU. As warps act in lock-step, any *branching logic* encountered by a warp must have all branches explored by all threads. The data created during the additional branch exploration is simply discarded. This phenomena is referred to as *branch divergence* and is warp-local; other warps continue to perform independently of the divergent warp. This can lead to scheduling conflicts where non-branching warps must wait for the divergent warps to complete. It is also generally a performance loss within a warp, especially for cases where one or more branches is long but uncommonly taken.

Finally, kernels can be launched in parallel on a single device, as long as that device has the capacity to do so. Streams are command sequences that execute in order internally, but can be concurrent with each other. The number of concurrent streams is device dependent, and additional streams will queue until the device has availability. Streams are unnecessary to run parallel commands on multiple devices, and are not needed for pipelining data transfers with kernel execution. Two commands from multiple streams cannot run concurrently if the host specifies memory manipulation or kernel launches on stream 0 (default) between them. Synchronization, where necessary, can be invoked within a stream or across streams with provided CUDA sync statements.

2.3 SPIN Model Checker

SPIN [7] is a widely used model checker designed to verify multi-threaded software. SPIN has an ever-growing list of features and options, including optimization techniques, property specification types, and hardware support. State spaces can be pruned using partial order reduction, speed can be increased by changing search strategies or disabling certain checks, and memory footprint can be reduced through bitstate hashing. SPIN can handle safety and liveness properties, any LTL specification, Büchi automata, never claims, and invariant assertions. Multicore support was added in 2007 [23], improved in 2012 [24], and extended to liveness properties in 2015 [20].

A central feature of the 2012 algorithm is the structure holding the frontier of newly discovered states. In order to assign these states to the N worker threads, SPIN uses two sets of $N \times N$ queues. By splitting each frontier queue into an $N \times N$ structure all threads can communicate without the need for mutex locks. Of these two queue sets, one (output) fills with a new frontier as the other (input) empties the current frontier. When the input queue is empty, all threads synchronize and the two swap labels. This process continues until both the input and output are empty or a violation is found. We adopt this structure for Grapple.

2.4 Swarm Verification

Recently, support for large-scale parallel model checking on CPU-based systems was added to SPIN in the form of swarm verification (SV) [26, 27, 25]. SV is a technique wherein a large number of small *verification tasks* (VTs) are run in parallel on many independent processors, including multiple CPUs, multicore CPUs, and distributed systems [25]. The term verification test is used in place of model checker or verifier because these tests are not guaranteed to complete. Instead, each test is given a set amount of time and memory to explore whatever portion of the state space it can. VTs can be as small as a number of KBs.

Each VT is independent, and the state space it covers is differentiated through the use of various *diversification techniques*. These techniques include reversing search direction or search order, randomizing nondeterministic choice order of transitions, and other perturbations of the original search algorithm. VTs do not share resources nor need to live on the same physical machine. Given enough parallel hardware, all VTs can run concurrently. When these resources are more limited, VTs will be scheduled like any other batch of independent programs.

The most potent diversification technique is the use of *statistically independent hash functions*. With up to 10^8 suitable unique 32-bit hash polynomials, in addition to other search diversification methods, the potential number of distinct concurrent searches is easily in the billions [25]. Hash functions reduce the state space graph via collisions; as each hash table is much smaller than the total number of states, collisions are frequent. If we treat each collision as valid (consider them the same state, even if that is not the case), the state space will be quickly, and naturally pruned.

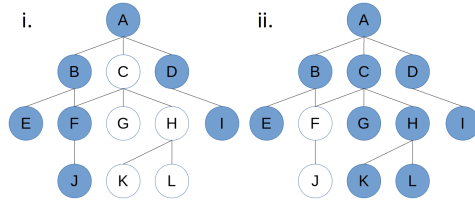


Fig. 2: Pruning States via Hash Collision. (i) Hash Collision $\{B, C\}$ on trace $ABEFJDI$. (ii) Hash Collision $\{E, F\}$ on trace $ABECEGHKLDI$.

Figure 2 depicts state-space pruning via collision. In both searches, a left-favoring Depth-First Search strategy is used, but their hash tables use different hash polynomials to store states. In the left graph, nodes B and C have the same hashed value, so C appears to be the same state and will not be expanded. In the right graph, E and F have the same value, preventing the expansion of

F . While this method of pruning all but assures that individual VT will not reach all states in the state space, with a sufficient number of diverse VTs, the swarm as a whole will achieve full coverage.

3 Swarm Verification via the Grapple Model Checker

The Grapple model checker brings the power of GPU computing to the model-checking problem via swarm verification. For simplicity of presentation, we discuss Grapple’s design in terms of a *Waypoints* (WPs) benchmark specifically designed for SV-based model checkers [27, 16]. The WP benchmark involves a model that can randomly generate more than 4 billion states. Said model is comprised of 8 processes each in control of 4 bits. At successor generation, the current process will nondeterministically set one of its bits to 1. Exploration progress in the benchmark is captured by the visitation of 100 randomly distributed states, or *waypoints*, with 100 waypoints suggesting a nearly complete state-space exploration. This style of presentation does not in any way imply that Grapple is limited to this one benchmark; it is still a general-purpose model checker. And indeed, we present results from an additional model in Section 4.

Although traditionally each VT is a small, sequential version of SPIN, this is not the case for Grapple VTs, which run on the GPU. As discussed in Section 2.1, the GPU has a SIMD/SIMT programming model: a single instruction or set of instructions is given to a group of threads operating on different data. Warps of 32 threads execute in lock-step, and all branches in logic must be fully explored by the entire warp. Mimicking SPIN by running a completely sequential VT on an entire warp would waste massive amounts of resources. Instead, we use a modified version of the 2014 GPU MC algorithm [12] to run a single, internally-parallel VT per warp. VTs execute independently in parallel outside of the warp, but internally (i.e., within a given VT), all data structures are shared among the threads and there is a single state space to explore.

While the queue structure and general search algorithm remain the same as the 2014 MC, there are a number of alterations made to the GPU VT to take advantage of the new swarm environment. First and foremost, the hash table is a bitstate implementation moved to shared memory. This hash table is only shared

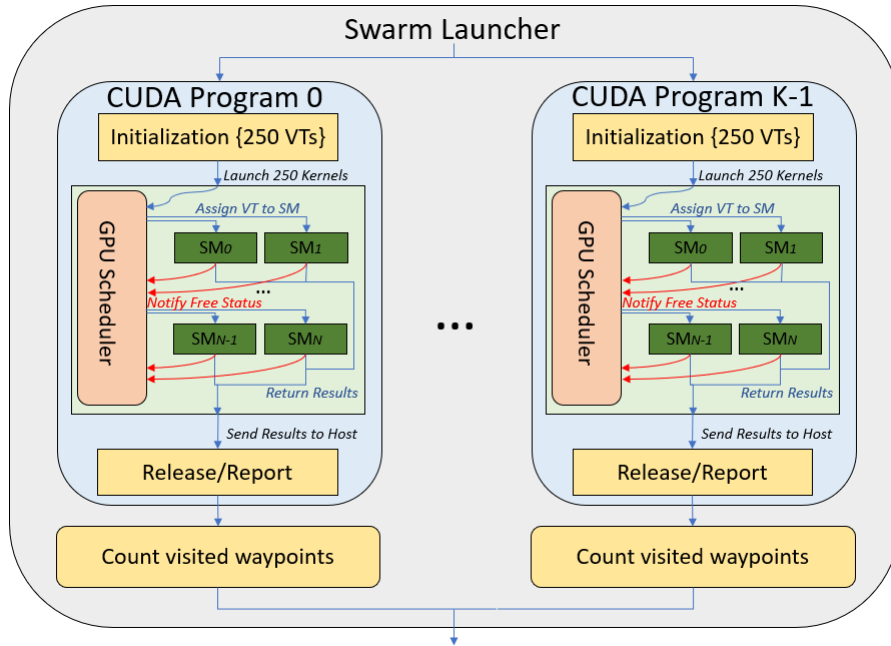
among threads within a VT and not between VTs. Factors typically considered weaknesses of a shared-memory approach are its locality to an SM and its small size (48KB maximum). With each SIMD-parallel VT limited to a single warp, all threads within the VT are guaranteed to be on the same SP within the same SM, and therefore all have access to this structure. The 48 KB limit is not an issue for VTs utilizing bitstate hashing, as such a table can hold nearly 400,000 entries. This is on the low end of the scale for a VT compared to those in other SV implementations [27, 25], but VTs of this size were shown to work well in a recent FPGA implementation [16].

Also as in the FPGA implementation, cuckoo hashing [8] has been replaced with an AB mix function based on the Bob Jenkins Linear Feedback Shift Register (LFSR) [29]. For this purpose, two random integers, A and B , are generated on the host machine for each VT and included as parameters in the VT’s kernel launch. This hash function change is motivated by a desire to better align with the FPGA implementation, as well as the elimination of the multiple-function schema used in the cuckoo algorithm. The random variables are reused on the GPU in some search strategies as quick random-digit generators, as on-device random generation tends to be convoluted and this method is more efficient. Since each VT is relegated to a single warp, the fast-barrier synchronization [41] used in the previous GPU MC implementation has also been removed. Instead, the on-board CUDA `__syncthreads()` function is used at the required synchronization points.

Grapple, like the FPGA swarm [16], runs multiple VTs within a single program, with additional copies of that program launched by script if necessary. In contrast, the SPIN swarm [27] is coordinated by a script that simply launches every VT as an independent thread. A Grapple program running on the GPU initiates multiple VTs, each a CUDA kernel, and utilizes streams to run these kernels in parallel whenever possible. The number of VTs that a core program can launch is dependent upon the hardware of the device(s) available, the memory footprint of each VT, and how initialization and memory transfers are handled. In the current design, all variables and structures are initialized, transferred to the GPU before kernel launch, transferred back to the host after kernel completion, and then freed in a single batch. Theoretically, more VTs could be launched within a program and additional efficiency squeezed out if the transfers were pipelined with some VT execution, but the current arrangement also has benefits.

Since the primary diversification techniques in Grapple are alterations in hash polynomial, search structure, and nondeterminism order, most of the host-level set-up is common across VTs. Overall, these common elements reduce the cost of this process to be nearly negligible when compared to time spent on the device. In this case, pipelining would increase overall complexity of the core code with minimal benefit. On the theme of common initialization, structures are placed in constant memory whenever possible so all VTs gain fast read-only access.

Figure 3 illustrates the control flow of Grapple. Upon start-up, a swarm script launches a CUDA program on all available hardware devices (GPUs). When

Fig. 3: Control flow for Grapple with $250 \cdot K$ VTs.

there is only a single device, these K programs must sequentialize with each other, with one program launching after the execution of the previous program and its sort instance (the Linux sort utility is used to count WPs) has terminated. Internally, each CUDA program initializes a number of VTs, in this case 250, sharing common data wherever possible to minimize overhead. Examples of this includes setting up the initial state and sending WP identifiers to GPU constant memory. This initialization/pre-launch procedure runs on the CPU (host).

Each VT is assigned to a single stream, and as many of them as possible will be launched in parallel to the N streaming multiprocessors (SMs) available on the device. The number (250) of VTs maintained by a given GPU program is a function of the global memory footprint of all structures associated with those VTs. While the hash tables are assigned to the 48KB of on-chip shared memory, frontier queues and other support structures must still hold the full-length global state vectors and combine to reach the upper limits of the GPU global memory. Despite sitting on global memory, these structures are still access-limited to a single VT, maintaining VT independence.

Once launched, a VT executes its complete search until its frontier queues are empty, and there are no more states to be explored. This exhaustion process is driven by the limited size of the hash table, and the collision-based pruning mentioned described in Section 2.4.

To achieve maximal utilization of SMs and therefore maximal parallelism at the SM-level, VTs are assigned to SMs using pipelining: as soon as a VT completes its execution on a SM, the GPU scheduler replaces it with a new VT,

until all VTs within the CUDA program have been executed on some SM. At this point, the host collects the discovered WPs from all 250 VTs and appends this information to a single output file. All data structures on both the GPU and CPU are released, and the program terminates. The output file is read by a sort utility, and current progress reported by the swarm script. The next GPU program is launched, and the process continues until all GPU programs in the swarm are exhausted.

Note that for a single GPU system, a swarm of size 50,000 VTs requires 200 sequentially launched CUDA programs. One of the benefits of Grapple, and SV in general, is that if additional GPUs are available, even on different machines in different locations, these 200 CUDA programs can run in parallel with each other without additional modification. These other GPUs may have more memory or more SMs, allowing more VTs per program or more concurrent execution of VTs, respectively.

Due to the abridged nature of VT searches, minute changes in control flow can have a major impact on the set of visited states for each VT. As hash collisions are resolved by dropping the new entry, even differences in the order of constituent operations change the results. To better understand a VT’s behavior, we offer in Algorithm 1 a comprehensive breakdown of a VT’s main control loop. Furthermore, in Section 4, we conduct a series of tests that illuminate the effects of making even minor changes to the code.

The nondeterministic choice (NDC) has a variety of different implementation options. Traditionally, all nondeterministic options would be accessed in order as in standard *BFS* (parallel BFS in this case) behavior. With minor modification, all nondeterministic options can be visited in random order. To minimize the amount of branching logic, all NDC order possibilities are enumerated in constant memory, and the selection of order is completely random for each step in the loop.

As described in Section 2.3, Grapple VTs use a set of $N \times N$ queue structures to allow lock-free communication between threads. Each thread has a set of N input queues and N output queues, with I slots in each queue. We call an $N \times N \times I$ set of queues a *queue structure*. In Section 4, we consider a queue structure in Grapple to be the same as a queue in SPIN and FPGA VTs. For this to hold, I will often be as small as four or five slots.

In Grapple, the input and output queue structures are sets of pointers to a single array in GPU global memory. To avoid illegal memory access, a VT must first check that there are slots available when attempting to insert a new state. In Algorithm 1, this check happens after a state is marked visited. If there are no queue slots available, the state is dropped and its successors potentially lost. If instead the queue check happens before the state is marked visited, the same state (or a state with the same hash value) can be visited later. This second location is used in FPGA and Grapple VTs.

The logic employed with this check also plays a factor in Grapple’s performance. If the check prevents writing outside the bounds of the underlying array structure, hence referred to as the *old guard*, it will still allow threads to write

to unintended targets. A stricter boundary check, the *new guard*, enforces the local limitation of I . In practice, illustrated in Section 4, VTs with the old guard have better performance.

Algorithm 1 State-Space Exploration Loop
executed by each VT thread

Each of a given VT's N parallel threads does the following:

```

while none of thread  $i$ 's output queues are empty do
  for all  $N$  of thread  $i$ 's input queues do
    while input queue  $j$  is not empty do
      for all processes in the model do
        for all nondeterministic choices NDC
          within a process do
            successor = successor_generation(process,
              NDC, state);
            selection = (mix(a, b, state));
            hashed_value = (selection/8) % table_size;
            sel = selection%8;
            visited_state = table[hashed_value];
            table[hashed_value] |= (1<<sel);
            if (visited_state &(1 <<sel)) == 0 then
              Report state back to CPU for check
              against 100 WPs
              Pick random thread  $i' \in N$  to
              output to
              if  $i'$  has slots then
                Insert the new state into queue  $i'$ 
              end if//implicit else drop the state
            end if
          end for//close for (NDC)
        end for//close for (process)
      end while
    end for
  _synctreads();
  Check output queues for emptiness
end while

```

dropped states to this point has been of random drops or partial-match drops (hash collisions). It is also possible to do complete explicit-state drops for specific state-vector matches. The default behavior of the FPGA swarm is to consider WPs to be violations. When one of these states is encountered, it is reported and dropped without generating successors. While for other models this behavior may lead to unreachable portions of the state space, it is not the case for the WP model. Our Grapple tests include variants with and without this WP dropping behavior.

The reason for the better behavior of the old guard is as follows. When a thread n attempts to write to another thread q , the new guard would make sure n is not writing to $q + 1$ instead. If n is attempting to write to the (non-existent) $I + 1$ slot of q , it instead overwrites slot 0 of $q + 1$. In practice, this is a random state-drop that replaces a *shallow* state in the queue structure with a *deeper* one. Both guards lead to a state-drop, but the old guard favors keeping deep states while the new guard favors shallow states. In general, the Grapple implementation uses the old, deep-state-favoring, guard logic.

All discussion of

4 Experimental Results

In this section, we present experimental results for Grapple. The first set of experiments use the WP benchmark to test variants of the Grapple VT design, and allow us to compare performance with the SPIN [27] and FPGA [16] swarms, as well as with our non-swarm GPU implementation [12]. All of these tests use the same 100 WPs, selected from a random distribution over the 32-bit integer space. The GPU used in these experiments is an Nvidia Geforce 660Ti GPU with 2GB GPU global memory, and 7 SMs. This is an older, inexpensive GPU model but still allows for Grapple to show sufficient performance. SPIN experiments run Swarm 3.2 with SPIN 6.4.7, using an Intel dual-socket server that has two Xeon E5-2670v3 CPUs (24 cores total) running at 2.3GHz and Hyper-Threading enabled (48 hardware threads total), with 128GB of RAM. FPGA experiments are done with cycle-accurate SystemC simulations using Xilinx Vivado HLS 2017.4, targeting a Xilinx Virtex-7 XC7V690T FFG1761-3 FPGA. The test environments for the SPIN and FPGA experiments are the same as in [16]. Additionally, we include experiments using the Dining Philosopher’s problem in order to demonstrate Grapple’s ability to discover a known deadlock violation. Finally, we show Grapple’s potential in a high-performance environment by running WP benchmark tests on Amazon’s EC2 GPU cloud platform [2].

4.1 WP Benchmark

FPGA experiments use internally sequential VTs with 48KB of storage each. The FPGA runs in batches of 44 concurrent VTs, starting a new batch when the previous one finishes. Unlike the general-purpose VT designs of the GPU and CPU swarms, which can be applied to any Promela model, the FPGA swarm is currently limited (hardwired) to the 32-bit random number generator. Fortunately, there are still some variants of this WP benchmark to test against.

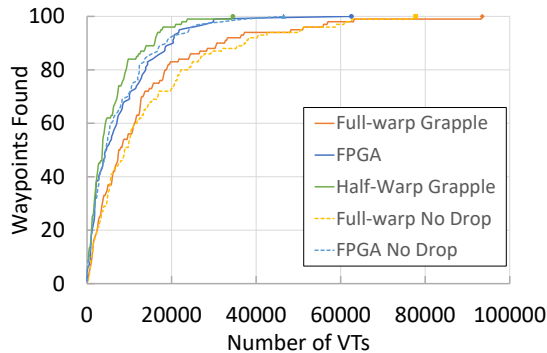


Fig. 4: Grapple VT vs FPGA VT.

Figure 4 shows combined results of two FPGA swarm variants and three Grapple variants running the WP benchmark. In the standard configuration, WPs are recorded upon discovery, considered a violation, and the state is dropped. Non-WP states first check the queue, and are marked visited and propagate if there are slots available or drop and remain unvisited if the queue is full. This allows the state (or a colliding state) to potentially be visited

later by the same VT. In later Grapple tests, we refer to this control flow as “FPGA-style”, as it matches the behavior of VTs in [16].

Half-warp (16 threads per VT) Grapple leads the FPGA in number of WPs from the very beginning and reaches the 100th WP in 34,500 VTs, over 28,000 fewer VTs than its counterpart. The full-warp (32 threads per VT) Grapple implementation, however, is outpaced by the FPGA. The FPGA completes the WP benchmark in 30,947 fewer VTs. While these three versions share the same control flow and queue structure size (4,096 entries), the half-warp Grapple implementation has much better performance when using the WP/VT metric. In terms of raw speed, however, the half-warp version is slower, with VTs lasting 650ms compared to the full-warp’s average of 451ms. Both Grapple versions cannot match the hardware-level speed of the FPGA implementation, but Grapple offers fast VTs with a much easier deployment process than the FPGA swarm.

There is also an alternate control flow, wherein the 100 WPs are reported but otherwise treated like any other state. In this case, all 100 are discovered by the FPGA in 46,515 VTs or roughly 74.4% the number of VTs as the previous iteration. Full-warp Grapple also sees improvement, completing in 77,750 VTs. This is not significant enough to catch up with the FPGA or half-warp Grapple. A no-drop version of half-warp Grapple was not included in these tests.

On FPGA hardware, the swarms from Figure 4 complete in an extremely fast 12.5s for the original and 9.3 seconds for no-drop, with individual VTs lasting only ~ 0.2 ms. These swarms, however, were run on a cycle-accurate FPGA simulator, where one second of simulated time takes approximately one hour of wall-clock time. The simulation allows for more useful data collection without harming FPGA performance, and is cheaper and faster than deploying to a physical FPGA.

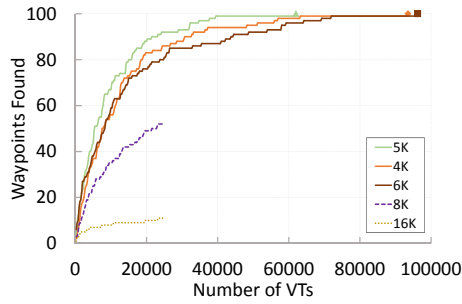


Fig. 5: Impact of frontier size on Grapple search.

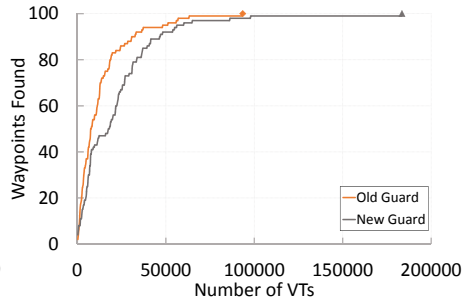


Fig. 6: Impact of guard logic change on Grapple search.

Figure 5 shows the impact of the queue structure size on Grapple’s performance. This test was inspired by the WP/VT difference between earlier half-warp vs full-warp tests. For the same size queue structure ($N \times N \times I$), a Grapple half-warp VT has more slots per thread (a smaller N value means a larger I value). Since the number of slots can impact state-drops (see Section 3), we ran a series of tests expanding the queue structure size (and thus the I value) for

full-warp Grapple. When $I = 16$ or $I = 8$ (16,384 or 8,192 total queue structure size), by 25,000 VTs we determined that these versions would not outperform the $I = 4$ control and terminated the swarms. $I = 6$ performs just slightly worse than the control. Grapple achieved peak performance with $I = 5$ (5,120 queue structure size), reaching 100 WPs in 62,000 VTs. This is better than the 93,500 VTs of the control, but still worse than the 34,500 of half-warp Grapple. Since half-warp Grapple uses a queue structure of 4,096 elements ($I = 16$ with $N = 16$), but outperforms all full-warp versions in WP/VT, the difference in performance requires further study. It is likely due to a low-level bottleneck, such as register access patterns or to differences in exploration order arising from the fewer random thread options.

We also tested the impact of altering the guard logic for full-warp Grapple’s queue structure, as explained in Section 3. Both versions use a queue structure with 4,096 entries, and otherwise identical control flow. Figure 6 shows the old guard logic maintaining a WP lead throughout the lifetime of the swarm, reaching the 100th WP in 93,500 VTs. The new guard logic takes an additional **90,000 VTs** to find all 100 WPs, with $\sim 47\%$ of the search spent looking for the final WP.

Unlike in Grapple and FPGA tests, where the hash table size is always 48KB per VT, SPIN swarm experiments run on a variety of different hash table sizes. While a 48KB hash table would be ideal for comparison purposes, a SPIN swarm requires hash tables to be multiples of 32. With the table size set to 32KB, SPIN ran for over a week without discovering all 100 WPs, after which we terminated the search. The next step up, with 64KB-hash-table VTs, managed to find 90 WPs in 263,220 seconds (just over three days). As in [16], the optimal configuration for the SPIN swarm seems to be a 256MB table per VT. This version uncovers all 100 WPs in 10,890 seconds, $\sim 3.4x$ as long as half-warp Grapple or $\sim 1.8x$ as long as full-warp Grapple.

The optimal setting for SPIN VTs requires over 5000x the amount of memory per VT as Grapple and the FPGA. A larger memory footprint for each VT lets a VT cover a greater portion of the state-space, but at the cost of longer execution time per VT. The SPIN results suggest that either the overhead for creating many small SPIN VTs hinders their effectiveness, or that SPIN’s implementation of diversification techniques favor larger VTs. While SPIN could run more concurrent VTs if more machines were available, improving performance, the same could be said for the Grapple and FPGA versions.

Non-swarm GPU tests were difficult for this model. Our original implementation in [12] called for four full explicit-state cuckoo hash tables to contain every possible state vector. Although the WP benchmark uses randomly generated 32-bit states, the states are still wrapped in a 64-bit unsigned long integer. Following the original MC design, the total hash storage alone would be 128GB, much larger than the 2GB of global memory on this GPU. Converting this checker to bitstate hashing allows us to cut the hash storage to a more-reasonable 500MB. However, this does not account for the other support structures that still use full 64-bit state vectors. The simplest solution is to run a

version that is 250x the size of a single Grapple VT, since we know 250 Grapple VTs can be allocated in one CUDA program without exhausting memory. A table this size can hold just over 98 million states, a fraction of the statespace generated by the WP benchmark. Our non-swarm checker explores this space in 352 seconds, reaching 10 WPs. As a standalone program, the GPU MC clearly cannot compete with the full state-space exploration of Grapple.

4.2 Dining Philosophers Model

Table 1 contains results for Dining Philosophers, where each philosopher picks up the left stick, then the right, releases the left and then the right. There is a violating state (deadlock) when all philosophers pick up their respective left stick concurrently. The minimum number of VTs tested is 7, since less than 7 would take the same amount of time to run on this GPU. For versions with more processes, we use sets of 451 VTs (an arbitrary large number that fits within the GPU memory footprint), but for DP10 and DP11, we determined that more precision would be better than just saying $x \leq 451$. The number of VTs needed to fully explore the state space increases dramatically when increasing the number of processes to 12. This is as expected, as DP11 has 177,146 states to fit into 392,800 slots per VT ($\sim 45\%$ occupancy), while DP12 has 531,440 states to fit into the same number of slots ($\sim 135\%$ occupancy). Beyond 12, we prematurely terminate the search due to the low rate of new state discovery.

Table 1: Dining Philosophers Model in Grapple.

Number of processes	% of VTs finding violation	Average VT execution time	State space size	# of VTs to explore	% of state space covered by first 451VTs
10	67.72	195ms	59048	100% in 7	100
11	46.65	366ms	177146	100% in 14	100
12	25.55	677ms	531440	100% in 3157	99.99
13	13.75	832ms	1594322	99.21% in 24,805	98.65
14	11.18	882ms	4782968	97.76% in 13,530	92.72
15	11.35	902ms	14348906	93.19% in 50,061	76.56

The final column of Table 1 shows the percentage of the state space covered in the first 451 VTs. Due to search overlap, the number of unique states visited grows logarithmically with the number of VTs. The effect is more pronounced in a deterministic model like Dining Philosophers, since the only source of diversification in Grapple for such models is the VT’s hash polynomial.

4.3 Large-scale Results

For our large-scale experiments, we used two Amazon EC2 nodes [2], one with 4 and one with 8 Tesla V100 devices. Each device features 16GB global memory and 80 SMs. All devices for each configuration run concurrently and their reported WPs are collected by a script on the host. As in the previous tests, each VT is independent and features data structures private to said VT. There is no

inter-GPU communication other than WP counting by the script. Each CUDA program runs 2,000 VTs between reports to the host.

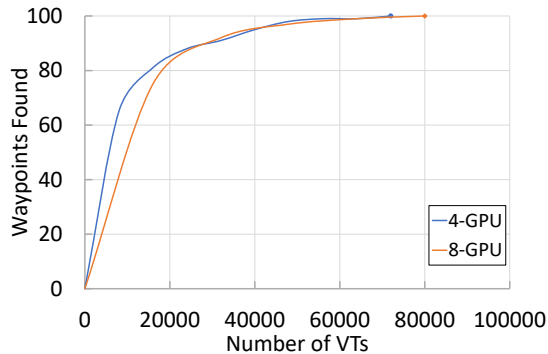


Fig. 7: Grapple with 16 threads/VT on Amazon EC2

As in Figure 7, the 4-GPU node reaches all 100WPs in 72,000 VTs (18,000 per GPU). The 8-GPU node reaches all 100 in 80,000 VTs (10,000 per GPU). Even with state-recording overhead they complete in 42 minutes and 21 minutes, respectively. This is faster than our previous results with such recording disabled. Turning off state-recording results in a reduction of average VT time from 1.02250s to 203.51ms. This is a significant reduction of 80.1%.

5 Related Work

In [23], SPIN was extended to support dual-core processors, using nested DFS to check safety and liveness properties. This work was extended to multicore systems for safety properties in [24] and liveness properties in [20]. Despite the earlier debut of a distributed model checker [11], the dual-core version of SPIN was the first parallel MC to reach wide adoption. Other work sought to avoid the naturally sequential depth-first post-order found in dual-core SPIN’s nested DFS algorithm by leveraging the parallelism in breadth-first reachability analysis on both distributed [35] and multicore systems [10]. This was mainly accomplished using two algorithms: One Way Catch Them Young (OWCTY) and Maximal Accepting Predecessors (MAP). Both algorithms perform parallel reachability analysis, but differ in the way they detect cycles in the state-space graph.

Early GPU-based MC efforts focused on *a priori* graph exploration, as opposed to generating new states *on-the-fly* [9, 17, 22, 28, 32]. The first on-the-fly GPU approach used the GPU to generate new states with enabled transitions, and the CPU for duplicate detection [18]. This is not unlike waypoint counting in Grapple, but their system makes less efficient use of the GPU hardware and is not based on SV. GPUexplore [38] was introduced in 2014 along with our own GPU-based model checker [12]. While we tried to redesign SPIN to take advantage of the GPU architecture, GPUexplore worked on Labeled Transition Systems (LTSs) and followed a symbolic approach. Grapple uses VTs based on our 2014 design, so it is still very different than GPUexplore. A GPU-based on-the-fly reachability checking system for LTSs that achieved 50–100x performance over sequential search was presented in [40].

In [36], GPUs were used for strong and branching bisimilarity checking. A GPU-based method for liveness checking for finite-state concurrent system ap-

peared in [37]. Three partial-order reduction algorithms were implemented on the GPU in [33], bringing GPUexplore closer to parity with existing CPU-based checkers. A second version of GPUexplore was released that same year, with improvements made to lock-less hashing and thread synchronization [39]. Unlike [37], this version does not include support for liveness properties. Scalability tests for GPUexplore were carried out in [13], achieving 5.5 million states/second on a 61.9 million state model. Additionally, they used GPUexplore to pit the 2015 Maxwell Architecture Nvidia Titan X GPU against the 2016 Pascal Titan X GPU, averaging a 1.73x improvement on the new device. A more in-depth comparison between cuckoo hashing and the GPUexplore table was carried out in [14], concluding that cuckoo hashing is 3x faster for random data and up to 9x faster for non-random data.

A GPU-based parameter-synthesis tool for stochastic systems was presented in [15]. Utilizing a single GPU, it achieves up to 31x the performance of sequential approaches. A multi-core version of the LTSMIN model checker [31] outperformed the 2005 multi-core SPIN and the 2008 multi-core DiVinE model checkers. In [19], a new multi-core DFS algorithm called CNDFS with better performance than the OWCTY algorithm was presented. This technique uses a swarm approach with state coloring to perform cycle detection concurrently with state-space exploration. LTSMIN saw further improvements in 2015 including support for new modeling languages [30].

In [21, 34], an FPGA was used to accelerate the exploration of a relatively small 10,000-state model, achieving a 50x speed-up compared to its software equivalent. The FPGA swarm of [16], to which this work is compared, achieved a 900x improvement over a SPIN swarm for a model of a much more substantial size (4B+ states). While this scale of improvement is unlikely for a single GPU device, the process of deployment to the FPGA is much more complex compared to the GPU. Additionally, their FPGA swarm was designed specifically for the 32-bit WP model, while Grapple can handle arbitrary Promela models.

6 Conclusions

We have presented Grapple, a new framework for highly efficient explicit-state model checking on the GPU. Grapple is based on swarm verification (SV), and its features include: a parallel swarm of internally parallel verification tasks (VTs); GPU-optimized implementations of hash functions and bitstate representation of visited states; and optimal use of GPU shared memory, thereby eliminating inter-block communication/synchronization overhead. Our experimental results show that Grapple outperforms multicore SV [25] and GPU non-SV [12] approaches, and that it uses a number of VTs similar to that required by an FPGA swarm [16].

Future work includes adding support for larger state vectors, allowing us to test Grapple with larger-scale model instances from the BEEM database [3]. We will also investigate new diversification techniques, including randomized process order and alternative NDC search strategies.

References

1. About CUDA | NVIDIA developer zone. <https://developer.nvidia.com/about-cuda>.
2. Amazon EC2 P3 instances. <https://aws.amazon.com/ec2/instance-types/p3/>.
3. BEEM: BEncmarks for EXplicit Model Checkers- ParaDiSe. <http://paradise.fi.muni.cz/beem/>.
4. CUDA C programming guide. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>.
5. Green 500 — TOP500 supercomputer sites. <https://www.top500.org/green500/>.
6. OpenCL technology™- intel.com. <http://software.intel.com/OpenCL>.
7. Spin- formal verification. <https://http://spinroot.com/>.
8. D. A. F. Alcantara. Efficient hash tables on the GPU, 2011. Copyright - Copyright ProQuest, UMI Dissertations Publishing 2011; Last updated - 2014-01-23; First page - n/a; M3: Ph.D.
9. J. Barnat, P. Bauch, L. Brim, and M. Češka. Designing fast LTL model checking algorithms for many-core GPUs. *Journal of Parallel and Distributed Computing*, 72(9):1083–1097, 2012. Accelerators for High-Performance Computing.
10. J. Barnat, L. Brim, and P. Ročkal. Scalable multi-core LTL model-checking. In D. Bošnački and S. Edelkamp, editors, *Proc. of SPIN 2007: the 14th international SPIN conference on Model checking software*, volume 4595 of *Lecture Notes in Computer Science*, pages 187–203. Springer Berlin Heidelberg, 2007.
11. J. Barnat, L. Brim, and J. Stříbrná. Distributed LTL model-checking in SPIN. In M. Dwyer, editor, *Proc. of SPIN 2001: the 8th International SPIN Workshop on Model Checking of Software*, volume 2057 of *Lecture Notes in Computer Science*, pages 200–216. Springer Berlin Heidelberg, 2001.
12. E. Bartocci, R. DeFrancisco, and S. A. Smolka. Towards a GPGPU-parallel SPIN model checker. In *Proceedings of the 2014 International SPIN Symposium on Model Checking of Software*, pages 87–96. ACM, 2014.
13. N. Cassee, T. Neele, and A. Wijs. On the scalability of the GPUexplore explicit-state model checker. In *Proceedings Third Workshop on Graphs as Models (GaM 2017)*, Uppsala, Sweden, 2017.
14. N. Cassee and A. Wijs. Analysing the performance of GPU hash tables for state space exploration. *Electronic Proceedings in Theoretical Computer Science, EPTCS*, 263:1–15, 2017.
15. M. Češka, P. Pilař, N. Paoletti, L. Brim, and M. Kwiatkowska. PRISM-PSY: Precise gpu-accelerated parameter synthesis for stochastic systems. In M. Chechik and J.-F. Raskin, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, pages 367–384, Berlin, Heidelberg, 2016. Springer Berlin Heidelberg.
16. S. Cho, M. Ferdman, and P. Milder. FPGASwarm: High throughput model checking using FPGAs. In *28th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2018.
17. Y. Deng, B. D. Wang, and S. Mu. Taming irregular EDA applications on GPUs. In *Proceedings of the ICCAD '09: the International Conference on Computer-Aided Design*, ICCAD '09, pages 539–546, New York, NY, USA, 2009. ACM.
18. S. Edelkamp and D. Sulewski. Efficient explicit-state model checking on general purpose graphics processors. In J. Pol and M. Weber, editors, *Proc. of SPIN'10: the 17th International SPIN Conference on Model Checking Software*, volume 6349 of *Lecture Notes in Computer Science*, pages 106–123. Springer Berlin Heidelberg, 2010.

19. S. Evangelista, A. Laarman, L. Petrucci, and J. van de Pol. Improved multi-core nested depth-first search. In S. Chakraborty and M. Mukund, editors, *Automated Technology for Verification and Analysis*, Lecture Notes in Computer Science, pages 269–283. Springer Berlin Heidelberg, 2012.
20. I. Filippidis and G. J. Holzmann. An improvement of the piggyback algorithm for parallel model checking. In *Proceedings of the 2014 International SPIN Symposium on Model Checking of Software*, pages 48–57. ACM, 2014.
21. M. E. Fuess, M. Leeser, and T. Leonard. An FPGA implementation of explicit-state model checking. In *Proceedings of the 2008 16th International Symposium on Field-Programmable Custom Computing Machines*, FCCM '08, pages 119–126, Washington, DC, USA, 2008. IEEE Computer Society.
22. P. Harish and P. Narayanan. Accelerating large graph algorithms on the GPU using CUDA. In S. Aluru, M. Parashar, R. Badrinath, and V. Prasanna, editors, *Proc. of HiPC'07: the 14th international conference on High performance computing*, volume 4873 of *Lecture Notes in Computer Science*, pages 197–208. Springer Berlin Heidelberg, 2007.
23. G. Holzmann and D. Bošnački. The design of a multicore extension of the SPIN model checker. *IEEE Transactions on Software Engineering*, 33(10):659–674, Oct 2007.
24. G. J. Holzmann. Parallelizing the SPIN model checker. In A. Donaldson and D. Parker, editors, *Proc. of SPIN 2012: the 19th International Workshop on SPIN Model Checking Software*, volume 7385 of *Lecture Notes in Computer Science*, pages 155–171. Springer Berlin Heidelberg, 2012.
25. G. J. Holzmann. Cloud-based verification of concurrent software. In *Proceedings of the 2016 International Conference on Verification, Model Checking, and Abstract Interpretation*. Springer-Verlag, 2016.
26. G. J. Holzmann, R. Joshi, and A. Groce. Swarm verification. In *Proceedings of the 2008 23rd IEEE/ACM International Conference on Automated Software Engineering*, ASE '08, pages 1–6, Washington, DC, USA, 2008. IEEE Computer Society.
27. G. J. Holzmann, R. Joshi, and A. Groce. Swarm verification techniques. *Software Engineering, IEEE Transactions on*, 37(6):845–857, 2011.
28. S. Hong, S. K. Kim, T. Oguntebi, and K. Olukotun. Accelerating CUDA graph algorithms at maximum warp. In *Proc. of PPOPP '11: the 16th ACM Symposium on Principles and Practice of Parallel Programming*, pages 267–276, 2011.
29. B. Jenkins. A hash function for hash table lookup. <https://burtleburtle.net/bob/hash/doobs.html>.
30. G. Kant, A. Laarman, J. Meijer, J. van de Pol, S. Blom, and T. van Dijk. LTSmin: High-performance language-independent model checking. In *Tools and Algorithms for the Construction and Analysis of Systems*, pages 692–707, Berlin, Heidelberg, 2015. Springer Berlin Heidelberg.
31. A. Laarman, J. van de Pol, and M. Weber. Multi-core LTSmin: Marrying modularity and scalability. In M. Bobaru, K. Havelund, G. Holzmann, and R. Joshi, editors, *NASA Formal Methods*, volume 6617 of *Lecture Notes in Computer Science*, pages 506–511. Springer Berlin Heidelberg, 2011.
32. L. Luo, M. Wong, and W. Hwu. An effective GPU implementation of breadth-first search. In *Proc. of DAC '10: the 47th Design Automation Conference*, DAC '10, pages 52–55, 2010.
33. T. Neele, A. Wijs, D. Bošnački, and J. van de Pol. Partial-order reduction for GPU model checking. In *Automated Technology for Verification and Analysis: 14th*

- International Symposium, ATVA 2016, Chiba, Japan, Proceedings*, pages 357–374. Springer International Publishing, Oct. 2016.
34. M. E. Tie. Accelerating explicit state model checking on an FPGA: PHAST. Master’s thesis, Northeastern University, 2012.
 35. K. Verstoep, H. Bal, J. Barnat, and L. Brim. Efficient large-scale model checking. In *Parallel Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on*, pages 1–12, May 2009.
 36. A. Wijs. GPU accelerated strong and branching bisimilarity checking. In C. Baier and C. Tinelli, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, pages 368–383, Berlin, Heidelberg, 2015. Springer Berlin Heidelberg.
 37. A. Wijs. BFS-based model checking of linear-time properties with an application on GPUs. In *Computer Aided Verification: 28th International Conference, CAV 2016, Toronto, ON, Canada, Proceedings, Part II*, pages 472–493. Springer International Publishing, July 2016.
 38. A. Wijs and D. Bořnački. GPUexplore: Many-core on-the-fly state space exploration using GPUs. In E. Ábráham and K. Havelund, editors, *Proc. of TACAS 2014: the 20th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, volume 8413 of *Lecture Notes in Computer Science*, pages 233–247. Springer Berlin Heidelberg, 2014.
 39. A. Wijs, T. Neele, and D. Bořnački. GPUexplore 2.0: Unleashing GPU explicit-state model checking. In J. Fitzgerald, C. Heitmeyer, S. Gnesi, and A. Philippou, editors, *FM 2016: Formal Methods*, pages 694–701, Cham, 2016. Springer International Publishing.
 40. Z. Wu, Y. Liu, J. Sun, J. Shi, and S. Qin. GPU accelerated on-the-fly reachability checking. In *2015 20th International Conference on Engineering of Complex Computer Systems (ICECCS)*, pages 100–109, 2015.
 41. S. Xiao and W. chun Feng. Inter-block GPU communication via fast barrier synchronization. In *Proc. of IPDPS 2010: the IEEE International Symposium on Parallel Distributed Processing*, pages 1–12, April 2010.