

Leonard Techel

Low-Connectivity State Space Exploration using Swarm Model Checking on the GPU

September 19, 2021

supervised by:

Prof. Dr. Sibylle Schupp
Sascha Lehmann

Hamburg University of Technology (TUHH)
Technische Universität Hamburg
Institute for Software Systems
21073 Hamburg

STS
Software
Technology
Systems

Abstract

Using small, independent verification tests, model checking large models with billions of states can be parallelized on GPUs. This approach works great on models with high connectivity. However, it fails when a model has only few edges between states or large portions of the state space are hidden behind bottleneck structures.

Past work on the *Grapple* model checker has tried different approaches including depth-limiting and alternating between breadth-first and depth-first search.

The main goal of this thesis is to: (1) create a systematic way of classifying a model as *low-connectivity* (2) minimize the amount of verification tests needed to maximize the state space coverage of said models. To do that, we provide an implementation of the *Grapple* model checker.

Contents

1	Introduction	1
2	Related Work	3
3	Background	5
3.1	Model Checking	5
3.2	Parallelized Model Checking	6
3.2.1	Swarm Verification	6
3.2.2	Grapple Framework	7
3.2.3	Low-Connectivity Models	7
3.3	CUDA	8
4	Theory	11
4.1	Model Definition	11
4.2	Grapple Model Checker	11
4.2.1	Model Definition and State Generation	11
4.2.2	Queues	12
4.2.3	Hash Table	12
4.3	State Space Diversification	14
4.3.1	Diversification Techniques	14
4.3.2	State Pruning and Start Overs	14
4.4	Waypoints Model	15
4.5	Low-Connectivity Models	15
4.6	Counting Unique States Visited	15
5	Implementation	17
5.1	Source code	17
5.1.1	Main Program	17
5.1.2	CUDA Kernel	17
5.1.3	Queues	18
5.1.4	Model and Successor Generation	19
5.1.5	Hash Tables	19
5.1.6	HyperLogLog++	20
5.2	Usage	21
6	Evaluation	23
6.1	Correctness and Comparison with the Paper	23
6.1.1	Comparison of Waypoints and HLL as Estimators of State Space Coverage	23
6.1.2	Impact of Hash Table Size and Sequential VT Execution	24
6.1.3	Increasing State Space Coverage using Start Overs	25

6.2	Low-Connectivity Model Evaluation	26
6.2.1	Dining Philosophers Problem with Different Numbers of Processes	26
6.2.2	Comparison of Unique States Visited, Total States Visited and State Space Coverage	27
6.2.3	BFS Frontiers of Low-Connectivity Models	27
6.2.4	Start Over Strategy on Low-Connectivity Models	30
7	Conclusion	33
7.1	Future Work	33
7.2	Discussion	33
	Bibliography	35

List of Algorithms

1	Fundamental State Space Exploration Loop	6
2	Grapple state space exploration loop of a single worker	13
3	Waypoints model	15
4	Branching-Free Ternary Operator	19
5	On-The-Fly State Generation on the GPU	20
6	Bitstate hashing	21

List of Figures

3.1	State pruning using hash collisions	7
3.2	Examples of low-connectivity properties	8
5.1	A single VT of the 4D array, mapped to 1D memory	18
6.1	Waypoints model, comparison between reference and our implementation	23
6.2	Comparison between Waypoints and HyperLogLog state space coverage estimation	24
6.3	State space exploration with start overs	26
6.4	BFS frontier visualization of the waypoints model	28
6.5	BFS frontier visualization of low-connectivity models	29
6.6	Start over strategy on low-connectivity models	31

List of Tables

6.1	Exploration of the waypoints model using different hash table capacities in 80 runs à 250 VTs	25
6.2	Parallel and sequential exploration of the waypoints model	25
6.3	Exploration of the dining philosophers problem with 11, 12 and 13 processes	27
6.4	Exploration of low-connectivity models, each after 125 000 VTs	27

1 Introduction

In an explicit-state model checker, state space exploration of large models with billions of states is a time-consuming problem.

Using swarm verification, the problem can be split into small, independent verification tests. Past work has shown that by executing these verification tests in parallel on GPUs, a high-speed model checker can be implemented [3].

GPUs are interesting for model checking because they allow a significant performance increase on parallel algorithms, are widely available and relatively cheap in comparison to specialized hardware like FPGAs/ASICs.

To create the small, independent verification tests, diversification is used. Each verification test only covers a subset of the state space. Together, the verification tests nearly achieve full state space coverage. This approach works great on models with high connectivity. However, it fails when a model has only few edges between states or large portions of the state space are hidden behind bottleneck structures.

Contributions:

- Breakdown of the state space exploration loop using subroutines
- Implementation of a Grapple model checker
- A method of estimating unique states visited with a fixed error
- The *start over* search strategy that can reach deeper states
- A visualization of Grapple BFS frontiers

2 Related Work

- Parallel Breadth-First Search (BFS) [7]
- Swarm Verification [8]
- Swarm Verification Techniques [9]
- Grapple Framework [3]
- On-The-Fly State Generation on the GPU [1]

3 Background

This chapter gives a brief summary on the main topics that this thesis is based on. First, we define the scope of model checking that we are looking at in Section 3.1. We continue with an overview of Swarm Verification on the GPU and the specific problem that we are going to investigate, low-connectivity state space exploration, in Section 3.2. Last but not least, we give an introduction to the CUDA GPU programming framework in Section 3.3.

3.1 Model Checking

Model checking is a formal verification method that checks whether a state machine satisfies a specification. For example, the state machine of an elevator may be verified to meet the safety property of not opening the doors between floors. To do so, a *model checker* searches the state space for counterexamples, also called violations. When a violation is found, the path of state transitions that leads to the violation is reported back.

Model checking can be divided into three main problems: Description of models through state machines, definition of the specification through temporal logic, and algorithms that verify whether a state machine models a specification.

There are two main branches of model checking: Explicit-State Model Checking and Symbolic Model Checking. Explicit-State Model Checking can only verify finite state machines. In particular, the model has to have finite states, each state needs to be representable by a finite-size tuple containing its atomic propositions, and the model changes state through execution of state transitions. To overcome these limitations and verify potentially infinite-size state machines or systems of unknown structure, Symbolic Model Checking uses the abstraction of *symbols*, each representing a set of states and transitions. Within this thesis, we are only considering explicit-state model checking.

Each state in a model is labelled with atomic propositions that hold true while the state is active. An example for such propositions are the current values of variables in a program at a given state. In a specification, different types of properties can then be expressed onto these propositions. Three common properties are reachability, safety and liveness: Reachability means that an atomic proposition holds true at some state in the future. Safety means that an atomic proposition holds true at all states in the future. Liveness means that an atomic proposition holds true infinitely often in the future, meaning that it does not happen that the atomic proposition never holds true. Within this thesis, we are only considering reachability and safety properties.

A challenge all model checking algorithms have to face is the *state explosion problem*. In an asynchronous model of n processes, each consisting of m states, the number of states grows exponentially by the number of processes, namely m^n . This means that even for small models, it is often not possible to fit all reachable states of the system into a computer's memory. Therefore, every model checking algorithm needs to reduce the state space in some sense. However, even then, a non-parallel algorithm may need a lot

of physical time for exhaustive verification of the state space. In exhaustive verification, all states are visited and checked for a violation. [2, 6]

3.2 Parallelized Model Checking

The basic operation in an explicit-state model checker is a state space exploration loop, as defined in Algorithm 1. Model checking can be speed up by parallelizing the state space exploration, for example using a parallel breadth-first search (BFS) [7].

A major challenge in parallelized BFS is the communication overhead between threads: Shared memory does not allow to easily split the work onto a cluster of heterogeneous processors or the massively parallel architecture of a GPU on which thousands of threads can exist simultaneously.

Algorithm 1 Fundamental State Space Exploration Loop

```

while there are unvisited states do
    mark state as visited
    if state violates spec then
        report path to state

```

3.2.1 Swarm Verification

Swarm Verification is a new approach on parallelized model checking [8]. It solves the challenge of parallelizing state-space search by splitting the state space exploration into many small, independent, memory-limited tasks called *Verification Tests* (VTs). Each VT only covers a small subset of the total state space and, using *diversification techniques*, uses a different search path. By executing all VTs and collecting their results, we still achieve nearly full state space coverage.

As VTs are independent of each other, we can easily execute them on heterogeneous computers. Even further, as VTs are also memory-limited, we can massively parallelize them on devices with very limited resources like GPUs.

Multiple diversification techniques can be applied, e.g. randomizing the order of nondeterministic choice on states with multiple outgoing transitions or reversing the search order. The most powerful diversification technique is state pruning using hash collisions: During state exploration, states get marked as visited in a limited-size hash table. By using a different hash function for each VT, the state space gets automatically pruned through hash collisions. The process is illustrated in Figure 3.1, that shows an example state graph on which BFS is performed by two VTs. On the left side, state B and C cause a hash collision, resulting in the subgraph originating in state C being removed. On the right side, state E and F cause a hash collision, resulting in the subgraph originating in state F being removed. Each state exploration on its own does not cover the whole state space, however, by combining the results from both searches, all states are covered again.

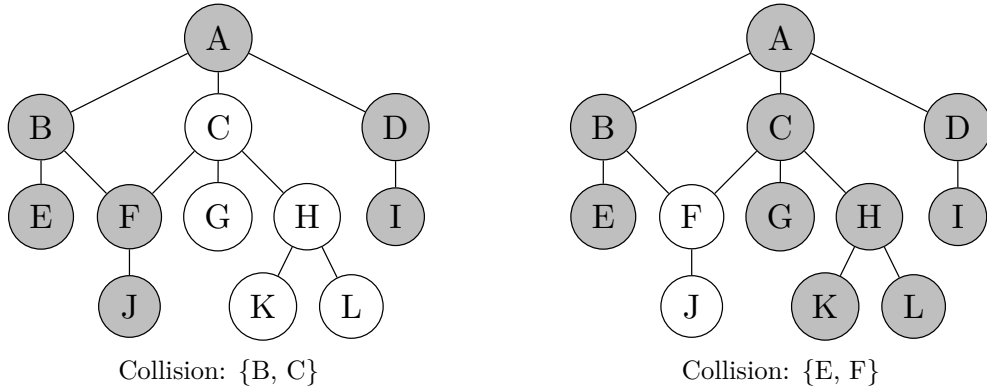


Figure 3.1: State pruning using hash collisions

3.2.2 Grapple Framework

The paper “Swarm model checking on the GPU” by DeFrancisco et al. [3] which this thesis is based on contributes the *Grapple* framework for parallel swarm verification on the GPU using CUDA. In their implementation of Swarm Verification, each VT’s state space exploration runs in a parallel BFS, using CUDA’s built-in synchronization primitives to coordinate threads.

3.2.3 Low-Connectivity Models

A key observation of [3] is that their algorithm’s state space coverage in relation to the number of executed VTs slows down on models with *low connectivity*. A model has *low connectivity* if at least one of the following properties is satisfied:

- **Generally Linear:** The average number of edges per state is close to two: One inbound, one outbound
- **Bottleneck Structures:** A single state or group of states other than the initial state that needs to be passed to reach most of the state space

See Figure 3.2 for examples of these properties.

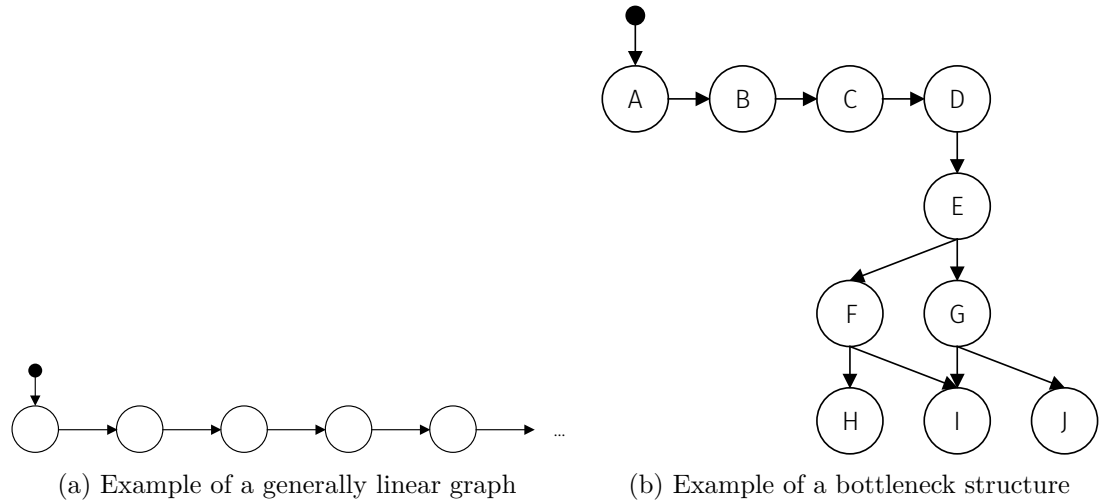


Figure 3.2: Examples of low-connectivity properties

3.3 CUDA

CUDA is NVIDIA's proprietary GPU programming framework that provides automatic massive scalability.

The CUDA model defines a three-level abstraction. At the lowest level, a program function called *kernel* is defined using C/C++. The kernel is executed by defining the amount of *threads* and *blocks* that run in a *grid*: Each grid contains multiple blocks. Each block contains multiple threads. Each thread in a block executes the kernel function using the single instruction, multiple thread (SIMT) execution model. Thus, each thread should work on an independent piece of memory. Only threads within a block can cooperate, for example using shared memory or synchronization barriers.

Internally, the GPU, also called *device*, maps the blocks of threads onto its *streaming multiprocessors*. By doing so, tens of thousands of threads can be executed in parallel. The architecture allows scaling even further onto multiple GPUs.

Furthermore, CUDA defines multiple memory levels, each with different size and access speed.

- *Registers* are used for thread-local variables. There can be at max 255 registers à 32 bit per thread.
- *Shared memory* can only be accessed by the threads of the block in which it got allocated. As it is on-chip, it is faster than global memory. By default, there is 48 KB of shared memory.
- *Constant memory* is part of the device memory, with the constant cache in front of it. Cache hits are served at higher speed than just using global memory. There is 64 KB of constant memory.

- *Global memory* is the largest memory with usually multiple gigabytes of size. It is placed in device memory and thus also the slowest.

CUDA requires the problem to be solved by a program to be split into many small sub-problems. This is a good fit with Swarm Verification, which turns the model checking problem into many small, independent VTs.

4 Theory

This chapter gives a detailed explanation of the algorithms used by our model checker. It starts with an in-depth description of the Grapple framework’s state space exploration loop and its components in Section 4.2. We then continue with a list of known state space diversification techniques in Section 4.3. The chapter finishes with a solution to the problem of counting the number of unique states visited in Swarm Verification in Section 4.6.

4.1 Model Definition

The model is defined using a *Kripke Structure* $M = (S; S_0; R; L)$ where S is the finite set of states, $S_0 \subseteq S$ is the finite set of initial states, $R \subseteq S \times S$ is the left-total transition relation and $L : S \mapsto 2^{AP}$ is the labelling function that assigns each state its atomic propositions.

4.2 Grapple Model Checker

In Grapple, each VT executes an internally parallel state space exploration. We call it *internally parallel* to differentiate between the parallelism *within* a VT and the parallelism *between* multiple VTs.

Each thread in a VT executes Algorithm 2, which is a more abstract version of the Grapple state space exploration loop (Algorithm 1 in [3]). We are going to call this algorithm the *Grapple algorithm*. The Grapple algorithm is based on Holzmann’s parallel BFS [7].

4.2.1 Model Definition and State Generation

State generation strategies of model checkers can be divided into *a priori* and *on-the-fly* generation. In *a priori* generation, the state space graph is completely known before running the verification algorithm, for example as an adjacency matrix. In *on-the-fly* generation, the successors of each state are created on the fly during verification. Grapple uses *on-the-fly* generation on the GPU. Furthermore, Grapple can only verify models with the same number of nondeterministic choices

$$M = ((S; S_0; R; L); P; NDC)$$

The states of a model support two operations: `state_successor($p, ndc, state$)` *on-the-fly* returns a successor to a state and `state_violates($state$)` returns whether a state violates.

$$\text{state_successor} : S \times NDC \times P \mapsto S$$

$$\text{state_violates} : S \mapsto \{0, 1\}$$

4.2.2 Queues

In summary, the parallel BFS is implemented using lock-free communication between threads via a shared, multidimensional queue array. Being lock-free removes the overhead of waiting for locks to be released and allows the algorithm to be split up onto many threads, limited only by the quadratic queue size.

For each pair of threads, there is only one piece of memory to communicate over. Writing and reading to this piece of memory is coordinated by splitting the algorithm into two alternating phases:

In the first phase $t = 0$, each thread i processes all states from its input queues $Q[t][0 \dots N][i]$. Here, N denotes the amount of threads executing a VT, with $N = 32$ being the default. For each input state, successors are generated on-the-fly. For each successor, we check whether it is already visited in a hash table. If a successor is priorly unvisited, we check whether it violates. Violating successors are reported to the host using a buffer. Non-violating successors are written to a random worker j 's input queue $Q[1 - t][i][j]$. When all input states are processed and threads are synchronized, the phases get swapped, i.e. $t = 1 - t$. This process repeats until no more unvisited successors are discovered.

The queues support three operations: `queue_push($Q, state$)` adds a state to the back, `queue_pop(Q)` removes and returns the first state from the front and `queue_empty(Q)` tells whether a queue is empty.

4.2.3 Hash Table

The hash tables support one operation: `mark_visited($V, state$)` marks a state as visited and returns whether it was already visited before.

Completeness A big challenge in Swarm Verification is the question on how much of the state space is actually covered, i.e. whether all violations are found. Due to overlap in state space between VTs and the sheer amount of states being checked by each VT, computing the exact number of unique states visited and thus the exact state space coverage can hardly be done. Instead, we can estimate the state space coverage, as described in Section 4.6.

Termination The algorithm terminates when there are no more unvisited state successors according to the hash table. In this case, all input queues of the next phase remain empty, i.e. $Q[1 - t][0 \dots N][0 \dots N] = \emptyset$.

There are two reasons why at some point no more unvisited states get discovered: Either, we did an exhaustive verification of the model, meaning that all states actually got visited. Or, which happens most of the time in Grapple, the hash table is full, causing

a collision of every new successor with an already visited state. Collisions are intended, as mentioned in Section 3.2.1.

Resulting from this, the algorithm at worst terminates after having visited the number of states being the minimum of queue size and state space size.

Algorithm 2 Grapple state space exploration loop of a single worker

```

 $t \in \{0, 1\} \leftarrow 0$   $\triangleright$  Current algorithm phase
 $visited \subseteq S \leftarrow \emptyset$   $\triangleright$  Hash table of visited states
 $Q[2][N][N][I] \leftarrow \emptyset$   $\triangleright$  Queues

 $\triangleright$  Initial state  $\triangleleft$ 
 $queue\_push(Q[t][0][0], S_0[0])$ 
 $mark\_visited(S_0[0])$ 

 $\_syncthread()$ 

 $done \leftarrow false$ 
while not  $done$  do
  for  $i = 0, \dots, N$  do
    while not  $queue\_empty(Q[t][threadIdx.x][i])$  do
       $state \leftarrow queue\_pop(Q[t][threadIdx.x][i])$ 
      for  $p \in \text{processes}$  do
        for  $ndc \in \text{nondeterministic choices within } p$  do
           $succ \leftarrow state\_successor(p, ndc, state)$ 
           $visited \leftarrow mark\_visited(succ)$ 
          if not  $visited$  then
            if  $state\_violates(succ)$  then
              report path to state
            else
               $next \leftarrow \text{random output queue } n \in N$ 
               $queue\_push(Q[1 - t][next][threadIdx.x], succ)$ 

       $\_syncthread()$ 

   $done \leftarrow queue\_empty(Q[1 - t][0 \dots N][threadIdx.x])$ 
   $t \leftarrow 1 - t$ 

   $\_syncthread()$ 

```

4.3 State Space Diversification

4.3.1 Diversification Techniques

The following diversification techniques are known: [9, 3]

- **State Pruning using Hash Collisions**
Each VT marks its visited states in a limited-size hash table, each using a different hash function. Hash collisions cause different partitions of the state space. Hash table exhaustion causes termination of the exploration when all new states are supposedly already visited.
- **Reverse order of nondeterministic choice**
Iterate through nondeterministic choices in reverse order, i.e. $NDC \dots 0$.
- **Reverse processes**
Iterate through processes in reverse order, i.e. $P \dots 0$.
- **Random order of nondeterministic choice**
Instead of iterating through nondeterministic choices from $0 \dots NDC$, choosing one random order for each VT. The random order remains the same through all searches in a VT.
- **Random process order**
Instead of iterating through processes from $0 \dots P$, choosing one random order for each VT. The random order remains the same through all searches in a VT.
- **Parallel Deep Search (PDS)**
Select one random nondeterministic choice per VT, discard others. The selected nondeterministic choice remains the same through all searches in a VT.
- **process-PDS**
Select one random process per VT, discard others. The selected process remains the same through all searches in a VT.

4.3.2 State Pruning and Start Overs

The basic BFS search terminates as soon as all successors are already visited according to the hash table. Thus, the maximum BFS depth is determined by the hash table's capacity. Also, each VT again starts at the initial state, resulting in all VTs potentially only exploring the tip of a model.

By continuing the BFS within a VT using the set of last unvisited successors, we can reach deeper states of the model. By exploring deeper states, we expect to discover more violations and thus a faster growing state space coverage.

4.4 Waypoints Model

The Waypoints (WP) model is a benchmark for Swarm Verification model checkers. It is first introduced in [9] and is used as primary benchmark in [3]. Our implementation is defined in Algorithm 3.

The model creates a state space of the size of all 32-bit integers, i.e. $2^{32} = 4\,294\,967\,296$ states. To do so, it uses eight processes, each in control of four bits. At successor generation, each process will nondeterministically set one of its bits.

There is a single reachability property: We predefine a set of 100 uniform randomly chosen 32-bit integers, called *waypoints*. A state violates if it is part of the set. As the waypoints are uniformly distributed, the number of unique discovered waypoints is equal to the percentage of state space coverage. For example, 24 discovered waypoints means that approximately 24% of the state space is covered.

Algorithm 3 Waypoints model

```

violations  $\leftarrow$  {100 random 32-bit integers}

function STATE_SUCCESSION(p, ndc, state)
|   return state |  $1 < (4 \cdot \text{process}) + \text{ndc}$ 

function STATE_VIOLATES(state)
|   return state  $\in$  violations

```

4.5 Low-Connectivity Models

4.6 Counting Unique States Visited

In Swarm Verification, we achieve a much faster verification by only covering *nearly* 100% of the state space. Resulting from this, it is important to know how much of the state space is actually covered. To find this out, we need to count unique states visited.

To calculate the percentage of state space covered, we have to know the total state space size. As in this thesis we only consider models with known state space size, we can provide it as pre-calculated constant.

The executed VTs may overlap in explored state space, meaning that across all VTs, a state may be visited multiple times. In order to calculate the exact number of unique states visited, we have to identify distinct states in the stream of all visited states. This is called the *count-distinct problem*.

Intuitively, one could collect all visited states in a set and then take its cardinality. However, this approach requires an amount of memory proportional to the amount of visited states which, as of the state space explosion problem, results in exponential memory usage.

A solution to this problem are *probabilistic cardinality estimators* that approximate the number of unique states within a fixed error using significantly less memory.

We choose the algorithm introduced in “HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm” by Flajolet et al. [4], as it is commonly used and relatively easy to implement.

For each VT, we create a HyperLogLog. Inside each VTs state space exploration, we call the **add** operation for each newly discovered state. When a VT has finished, we **merge** its HyperLogLog with those of all other already finished VTs. We then can execute the **estimate** operation on the global HyperLogLog to get the estimation of unique visited states across all finished VTs.

Cardinality estimation with HyperLogLog yields good results on the Dining Philosophers problem with 15 processes and $3^{15} - 1$ states. As expected, the number of unique states visited grows logarithmic until reaching around 100 %. However, on the Waypoints model with 2^{32} states, the estimation fails.

The paper “HyperLogLog in Practice: Algorithmic Engineering of a State of the Art Cardinality Estimation Algorithm” by Heule et al. [5] finds out that the estimation of cardinalities beyond one billion fails on the original HyperLogLog algorithm. They solve this problem by using a 64-bit hash instead, calling their improved algorithm HLL++.

Using HLL++, we can count unique states visited.

In our implementation, we are using a HLL with 2^{14} registers. This means that there may be a relative error of $\sigma = 1.04/\sqrt{2^{14}} = 0.008125$. The 65 % error bound for each value x is then $[x \cdot (1 - \sigma), x \cdot (1 + \sigma)]$. 95 % and 99 % can be reached using 2σ and 3σ .

5 Implementation

Our implementation serves two purposes: To reproduce the results from the Grapple paper [3]. And to run our experiment series on the state space exploration of low-connectivity models.

This chapter documents our implementation of a Grapple model checker. The general architecture of the model checker is described in Section 5.1. Usage instructions are described in Section 5.2.

5.1 Source code

We chose C++17 as programming language, so we can make full use of the CUDA Toolkit, which is provided as C header. The CUDA Toolkit is used in version 11.4. As build system, CMake 3.16 is used. Code is written using the Object-Oriented Programming paradigm. Our implementation consists of six components:

- The main program, running the CUDA kernels and collecting their results.
- The CUDA kernel, implementing a parallel state space exploration loop.
- The queues, providing lock-free communication between threads in a VT.
- The hash tables in which states are marked as visited.
- The model, providing successor generation and violation checking.
- The HyperLogLog, counting unique states visited across all VTs.

In the following, we are going to describe the implementation-specific details and challenges of each component.

5.1.1 Main Program

The main program runs in a single thread on the host. On startup, it seeds a global pseudorandom number generator (PRNG) and creates a random value for each CUDA thread. Then, for each run, it executes the CUDA kernel in a grid of $K = 250$ blocks, each consisting of $N = 32$ threads. Each block represents a VT, meaning that VTs are executed in batches of 250. Each thread represents a worker of a VT's parallel state space exploration. After each run, the main program collects the discovered violations and number of unique states visited, accumulates them and prints them to the standard output line-by-line as CSV.

5.1.2 CUDA Kernel

The CUDA kernel executes the state space exploration loop, as described in Section 4.2.

5.1.3 Queues

The queues provide a data structure for lock-free communication between threads in a VT, as described in Section 4.2.2. By default, each queue has a capacity of $I = 4$ states. They are only used on-device within the CUDA kernel and stored in global memory due to the limited size of shared memory. Resulting from this, we have to map $K \times 2 \times N \times N \times I$ queues into memory. This proposes two challenges: 1. We have to implement fixed-capacity queues. 2. We have to map multidimensional queues into a one-dimensional, linear memory allocation.

Usually, fixed-capacity queues are implemented as ring buffer. However, due to the design of the two-phase parallel state space exploration, it is sufficient to implement the queue as a singly linked list. In addition to that, we have to keep a pointer to the **head** and **tail** of the linked list. Then, in the first phase, each queue is filled through the **queue_push** operation by adding an element to its linked list and updating the **tail**, until the **tail** points to the last element in memory. In the second phase, each queue is completely emptied through the **queue_pop** operation by retrieving the **head**, then incrementing it until **head** equals **tail**. A queue is empty if both **head** and **tail** are null pointers.

The memory address of thread i 's input queue, which is an output queue of thread j , in algorithm phase t , of VT v , is calculated using the formula:

$$v \cdot (2 \cdot N \cdot N) + t \cdot (N \cdot N) + j \cdot N + i$$

Figure 5.1 illustrates the mapping of a single example VT with $K = 1$ and $N = 4$, meaning that we allocate memory for $1 \times 2 \times 4 \times 4 = 32$ queues. On the top row, each slot represents the memory allocation of a single queue. The remaining rows illustrate the constants enclosed in parentheses in our formula: Each set of output queues has a width of N , which we have to multiply with the index of the current thread to get at its first memory slot. Each phase has a width of $(N \cdot N)$, which we have to multiply with the index of the current phase. Each VT has a width of $(2 \cdot N \cdot N)$, which we again have to multiply with the index of the current VT to get to its first memory position. By adding up all products, and adding the index of the current input queue, we get the address of the exact slot.

0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3
To: 0				To: 1				To: 2				To: 3				To: 0				To: 1				To: 2				To: 3			
Phase 0																Phase 1															
VT 0																															

Figure 5.1: A single VT of the 4D array, mapped to 1D memory

5.1.4 Model and Successor Generation

Models are implemented using a data structure that represents a current state. Each state instance then provides the two operations `state_successor` and `state_violates`. The state instances are stored in the queues.

State generation is usually heavily based on branching through `switch`- and `if`-conditions. However, CUDA's SIMT execution model, in which warps of 32 threads operate on same instruction, causes threads to pause when their execution paths diverge. This means that state generation through branching results in a slower CUDA execution. As a countermeasure, Bartocci et al. propose to calculate all possible state transitions every time, preserving the SIMT execution and at the same time creating different states in each thread [1, Algorithm 3]. To do so, it makes use of the custom ternary operator described in Algorithm 4 that exploits the fact that boolean values evaluate to either zero or one.

Algorithm 4 Branching-Free Ternary Operator

```
function BRANCHFREETERNARY(bool c, int t, int f)
|   return  $c \cdot t + (1 - c) \cdot f$ 
```

A key difference of our implementation is that we do not fit a state's variables into a single integer using bit shifts. Instead, we use regular C++ class data members. The increased memory demand by regular data members can be neglected, as the queues in which we store state instances are in global memory which is usually gigabytes in size. Also, the same behavior as in the paper can be achieved using C++ bit field members.

Using our custom ternary, we can re-use DVE models from the discontinued BEEM model database¹. DVE is a file format used by the DIVINE 3 model checker. We translate DVE models into CUDA compatible C++ with the template described in Algorithm 5.

5.1.5 Hash Tables

The hash tables provide a data structure in which states are marked as visited. They are only used on-device within the CUDA kernel and stored in shared memory, resulting in a highly constrained size. Each VT operates on a single hash table, meaning that all threads in a block share one hash table. For memory efficiency, we use bitstate hashing [3]. Our hash tables can store $2^{18} = 262\,144$ states, taking up 32\,768 bytes of shared memory.

In bitstate hashing, a single bit represents whether a state is visited. Per clock cycle, CUDA can transmit 32 bit of shared memory, which is stored in 32-bit words. Resulting from this, we use a bucket size of 32 bit. The bit representing whether a state is visited is then computed using Algorithm 6.

The hash function's goal is to map random data with unknown distribution onto a nearly uniform distribution, so each bucket of the hash table is used equally likely. In our

¹<https://paradise.fi.muni.cz/beem/>

Algorithm 5 On-The-Fly State Generation on the GPU

```

▷ Global variables                                ◁
glob ← 0
▷ Process-Local variables use arrays. Here, N is the number of processes ◁
state[N] ← {0, ..., 0}

function STATE_SUCCESSION(p, ndc, state)
  ▷ Evaluate all guards before calculating transitions ◁
  guard1 ← state[p] = 0
  guard2 ← state[p] = 1 ∧ glob ≤ 2
  guard3 ← state[p] = 1 ∧ glob > 2
  ▷ Transition from state 0 to 1 ◁
  next.state ← BranchFreeTernary(guard1, 1, state.state)
  ▷ Transition from state 1 to 1 ◁
  next.glob ← BranchFreeTernary(guard2, state.glob + 1, state.glob)
  next.state ← BranchFreeTernary(guard2, 1, next.state)
  ▷ Transition from state 1 to 0 ◁
  next.glob ← BranchFreeTernary(guard3, 0, next.glob)
  next.state ← BranchFreeTernary(guard3, 0, next.state)

  return next

```

implementation, hashing is done using the Jenkins Hash function². The MurMurHash3³, which is used by our HyperLogLog, yields similar results.

Hash collisions resulting in false-positives are intended, as explained in Section 4.3.2. To achieve that in every VT different states cause a collision, each is including a different seed into the hash. The rate of collisions can then be controlled by the hash table's size.

5.1.6 HyperLogLog++

The HyperLogLog (HLL) provides an algorithm for distributed counting of distinct elements, as described in Section 4.6. In our implementation, the grid of VTs of each run is operating on a single HLL that is stored in global memory. Visited states are added to the HLL only on-device. The host collects the HLLs after the kernels have finished, merge them into one big HLL and does the estimation of unique states visited. Merging and estimating the amount of unique states visited is done only on the host.

²<https://www.burtleburtle.net/bob/hash/doobs.html>

³<https://github.com/aappleby/smhasher>

Algorithm 6 Bitstate hashing

```

▷ Initialize hash table buckets. Here,  $N$  is the hash table size





```

5.2 Usage

To verify a model using our model checker, a user has to complete the following steps:

1. Describe their model in a C++ class using Algorithm 5
2. Compile our model checker including their model
3. Execute the single output binary. On execution, the user may set the PRNG seed and number of runs using command-line options. The number of runs needed to achieve a certain state space coverage needs to be evaluated experimentally.
4. Interpret the CSV output

6 Evaluation

This chapter presents our series of empirical experiments conducted on our implementation. It starts with experiments on the correctness of our implementation and a comparison with the Grapple paper’s results in Section 6.1. We then continue with a series of experiments that compare the Waypoints model, which is known to be high-connectivity, with three low-connectivity models in Section 6.2.

All experiments are performed on an NVIDIA GeForce RTX 2080 Ti GPU. The calculation of visited state space is done using a HLL with 2^{14} registers.

6.1 Correctness and Comparison with the Paper

We start the evaluation with whether our implementation reaches the same number of discovered waypoints in relation to the amount of executed VTs as presented in the Grapple paper. We would expect both implementation to behave similarly, discovering at least 80 waypoints in 20 000 VTs and all 100 waypoints in about 100 000 VTs, as shown in Figure 6.1a. However, that is clearly not the case, as presented in Figure 6.1b. Our implementation found significantly fewer waypoints, even though the amount of VTs was 2.5x larger.

To further investigate differences between the implementations, we conducted additional experiments.

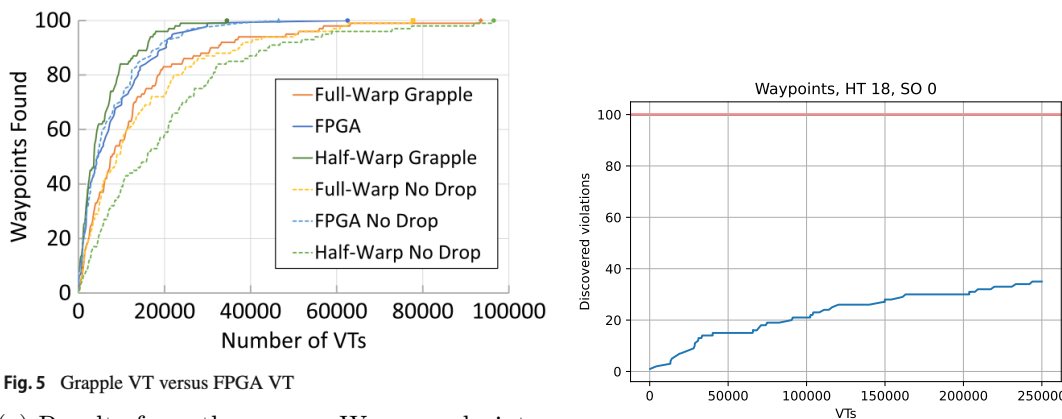


Fig. 5 Grapple VT versus FPGA VT

(a) Results from the paper. We are only interested in *Full-Warp Grapple*

(b) Results of our implementation

Figure 6.1: Waypoints model, comparison between reference and our implementation

6.1.1 Comparison of Waypoints and HLL as Estimators of State Space Coverage

As explained in Section 4.4, counting waypoints is a method for estimating the achieved state space coverage. To make sure that the discovered waypoints of our implementation

correctly indicate the state space coverage, we compared them to the state space coverage estimation of the HyperLogLog introduced in Section 4.6. We would expect both estimators to show similar results. For example, when 20 waypoints are discovered, the HyperLogLog estimation should show about 20 % state space coverage.

Figure 6.2 presents results for both estimators. Clearly, both resemble each other closely. Thus, we assume them to work correct. This also confirmed the results from our first experiment, underlining our implementation’s deficient state space exploration.

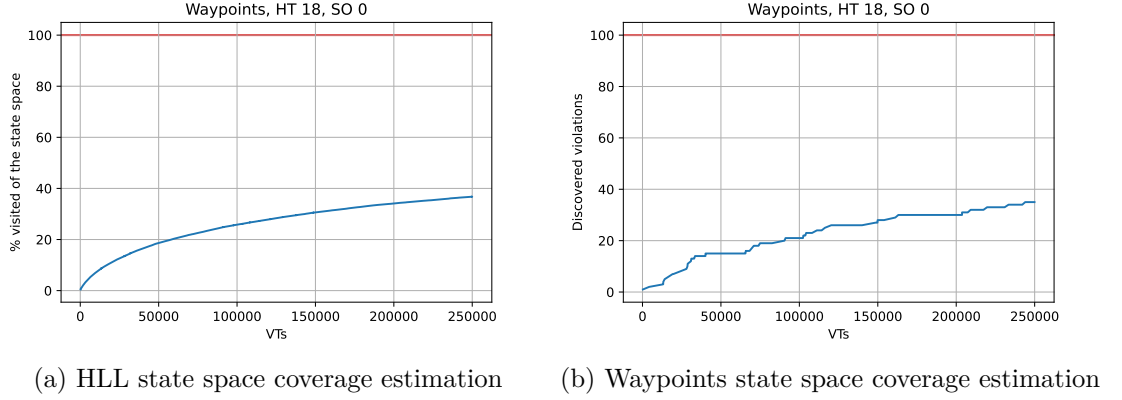


Figure 6.2: Comparison between Waypoints and HyperLogLog state space coverage estimation

6.1.2 Impact of Hash Table Size and Sequential VT Execution

Both implementations share a queue size of $32 \times 32 \times 4$. Beside of that, we identified two key differences of our implementation:

- Our hash tables are 33 % smaller than in the paper. This is due to our bit shift implementation presented in Section 5.1.5 which cannot address the full 48 KiB of shared memory.
- We batch VTs in a CUDA kernel, each executing 250 at once. In the paper, there seems to be only a single VT per CUDA kernel.

To find out whether increasing the hash table size to full 48 KiB increases waypoint discovery such that it is similar to the results from the Grapple paper, we experimentally changed our implementation’s addressing. As the hash table capacity should significantly affect the state space coverage of each VT, we would expect results close to those from the paper. For each experiment, we executed 80 runs à 250 VTs (= 20 000 VTs in total). Results are shown in Table 6.1.

The larger hash table yielded 2.1x as much state space coverage at only 1.45x as many total states visited. Thus, a larger hash table seems beneficial on this particular model.

However, in the paper, a Grapple swarm verification with 392 800 hash table slots per VT discovered 85 waypoints in 20 000 VTs, equalling 85 % state space coverage, which is about 3.7x as much as our large-size hash table implementation. Resulting from this, the large-size hash table is still far from reproducing the paper’s results.

Table 6.1: Exploration of the waypoints model using different hash table capacities in 80 runs à 250 VTs

HT Capacity	State Space Coverage	Total States Visited
262 144	11.04 %	63 568 607
393 056	23.14 %	91 937 680

In the paper, supposedly one VT is executed per *grid*, resulting in one CUDA kernel per VT. To compare our parallelization of VTs in a single CUDA kernel with the sequential execution of multiple kernels, we switched the compiled number of VTs and the executed runs. We would expect both variants to create similar results. As the parallel execution only needs a single kernel setup, we also expect a shorter execution time. Results for both variants are shown in Table 6.2. The state space coverage and total visited states of both experiments are highly similar. The execution time of parallel execution was more than 35x faster. Concluding this, the experiment showed that batch execution of VTs highly speeds up execution time at equal results.

Table 6.2: Parallel and sequential exploration of the waypoints model

VTs per Kernel	Runs	State Space Coverage	Total States Visited	Execution Time
1	250	0.456 %	63 568 343	26.691 s
250	1	0.450 %	63 568 498	0.754 s

The last two experiments showed that there seems to be a lack in our implementation beside the hash table capacity, queue size and execution order.

6.1.3 Increasing State Space Coverage using Start Overs

We suspect that part of the bad exploration performance is that each VT can only reach as deep into the state space as its hash table can find unvisited states. As a countermeasure, we introduced the start over strategy in Section 4.3.2.

The next experiment evaluates whether start overs actually increase the growth of state space coverage and thus the model checking performance. By reaching deeper states, we expect to discover more violations and thus a faster growing state space coverage. To do so, we compared 1000 exploration runs of the waypoints model with an arbitrary, much lower hash table capacity of $2^{14} = 16\,384$ slots and 15 start overs with the default configuration of $2^{18} = 262\,144$ hash table slots and 0 start overs. We chose 15 start overs to mitigate the smaller hash table, which is $\frac{1}{16}$ the size, resulting in 16 total searches

in the first exploration, including the initial search. Results for both explorations are presented in Figure 6.3a. With start overs, the exploration took 243s. Without, it took 442s.

As the number of total visited states by a VT varies between the two explorations due to for example different hash collisions, we also compared the cumulated number of total visited states of both explorations in Figure 6.3b.

Clearly, the start over strategy resulted in a significant increase of discovered violations on a similar number of visited states and VTs. Also, the execution time of the start over strategy was 1.82x faster. We conclude that start overs create a significant performance increase on the waypoints model.

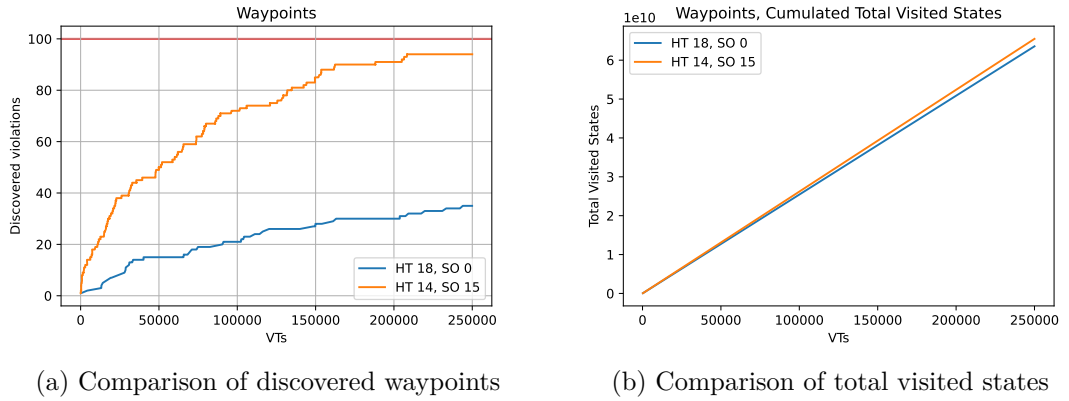


Figure 6.3: State space exploration with start overs

6.2 Low-Connectivity Model Evaluation

In this section, we try to make key differences of low-connectivity state space exploration tangible. The paper suspects three models of being low-connectivity:

- Dining philosophers problem
- Anderson queue lock mutual exclusion algorithm
- Peterson mutual exclusion protocol

6.2.1 Dining Philosophers Problem with Different Numbers of Processes

The paper suggests that a state space size surpassing the hash table capacity significantly slows down the state space exploration, resulting in a higher number of VTs needed to achieve $> 99\%$ state space coverage. To find out how much the state space size affects the exploration performance, we compared three variants with 11, 12 and 13 processes of the dining philosophers problem. We expect a slow-down starting with 12 processes using our the default 2^{18} hash table capacity. As presented in Table 6.3, exploring DP-13 took

213x the VTs as for DP-11, even though it is only 9x its state space size. This confirmed our hypothesis of a significant slow-down for this particular model.

Table 6.3: Exploration of the dining philosophers problem with 11, 12 and 13 processes

Model	State Space Size	Hash Table Utilization	# of VTs to explore
DP-11	177 146	68 %	100 % in 250 VTs
DP-12	531 440	203 %	100 % in 250 VTs
DP-13	1 594 322	608 %	100 % in 53 250 VTs

6.2.2 Comparison of Unique States Visited, Total States Visited and State Space Coverage

In the paper, the models yielded very different results concerning the unique and total states visited. Thus, the next experiment compares these performance indicators for our set of models on a fixed amount of 125 000 VTs. We would expect clear differentiation between the models and eventually a clear distinction between high- and low-connectivity. Our results are displayed in Table 6.4. Clearly, the waypoints model achieved the highest number of total states visited. We further observed an interesting property by taking the unique states visited in relation to the total states visited, presented in the *Rate* column: All three low-connectivity models achieved a rate < 0.1 %. In conclusion, the performance deficiency of the three presented low-connectivity models can be observed through unique states visited and total states visited.

Table 6.4: Exploration of low-connectivity models, each after 125 000 VTs

Model	Unique States Visited	Total States Visited	Rate
Anderson 3	$1.227\,79 \times 10^6$	15 477 969 183	0.008 %
Peterson 5	$4.085\,78 \times 10^6$	24 670 432 336	0.017 %
DP 15	$1.382\,13 \times 10^7$	30 436 723 433	0.045 %
Waypoints	$1.217\,59 \times 10^9$	35 725 590 331	3.408 %

6.2.3 BFS Frontiers of Low-Connectivity Models

The algorithm introduced in Section 4.2 is operating in alternating phases that produce so-called *frontiers*. Each frontier is an intermediate state of the BFS, representing the deepest states found until now. A new frontier is built by checking whether generated successors were visited before. For this experiment, we called already visited states *failed* and states added to the new frontier *visited*.

To find out whether different models affect the BFS frontiers in terms of visited and failed states, we executed a single VT and counted visited and failed states for each

frontier generated.

In Grapple, each BFS starts with a single initial state. As of this, we would expect the number of visited states to quickly build up until the number of failed states takes over due to the hash tables filling up.

For visualization of BFS frontiers, we provide two different methods:

1. Plotting the absolute number of visited and failed states.
2. Plotting visited and failed states in relation to each other, both adding up to 100 %.

Figure 6.4 presents both visualizations for the waypoints model.

- WP, DP and Peterson showed the expected curve in both representations
- WP, DP and Peterson have an upper limit of generated successor states per frontier
- The amount of BFS rounds and number of states per round highly differed between models
- The Anderson model showed two interesting differences: A big notch in the absolute representation (bottleneck?) and a near constant amount of visited states throughout all BFS rounds in the relative representation

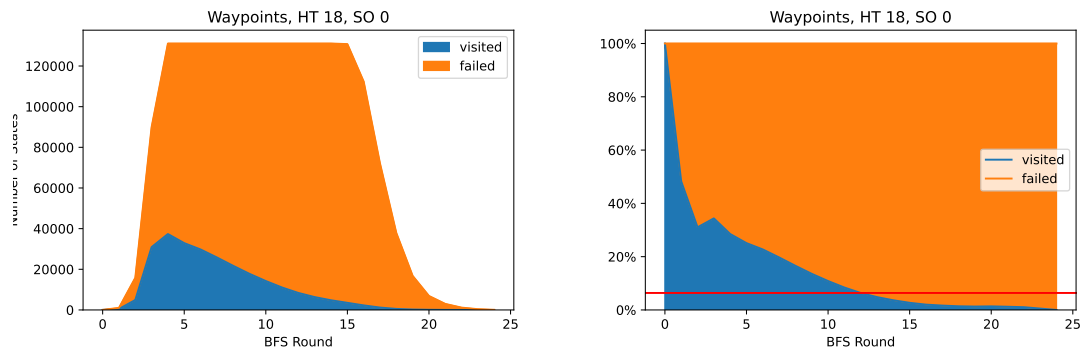


Figure 6.4: BFS frontier visualization of the waypoints model

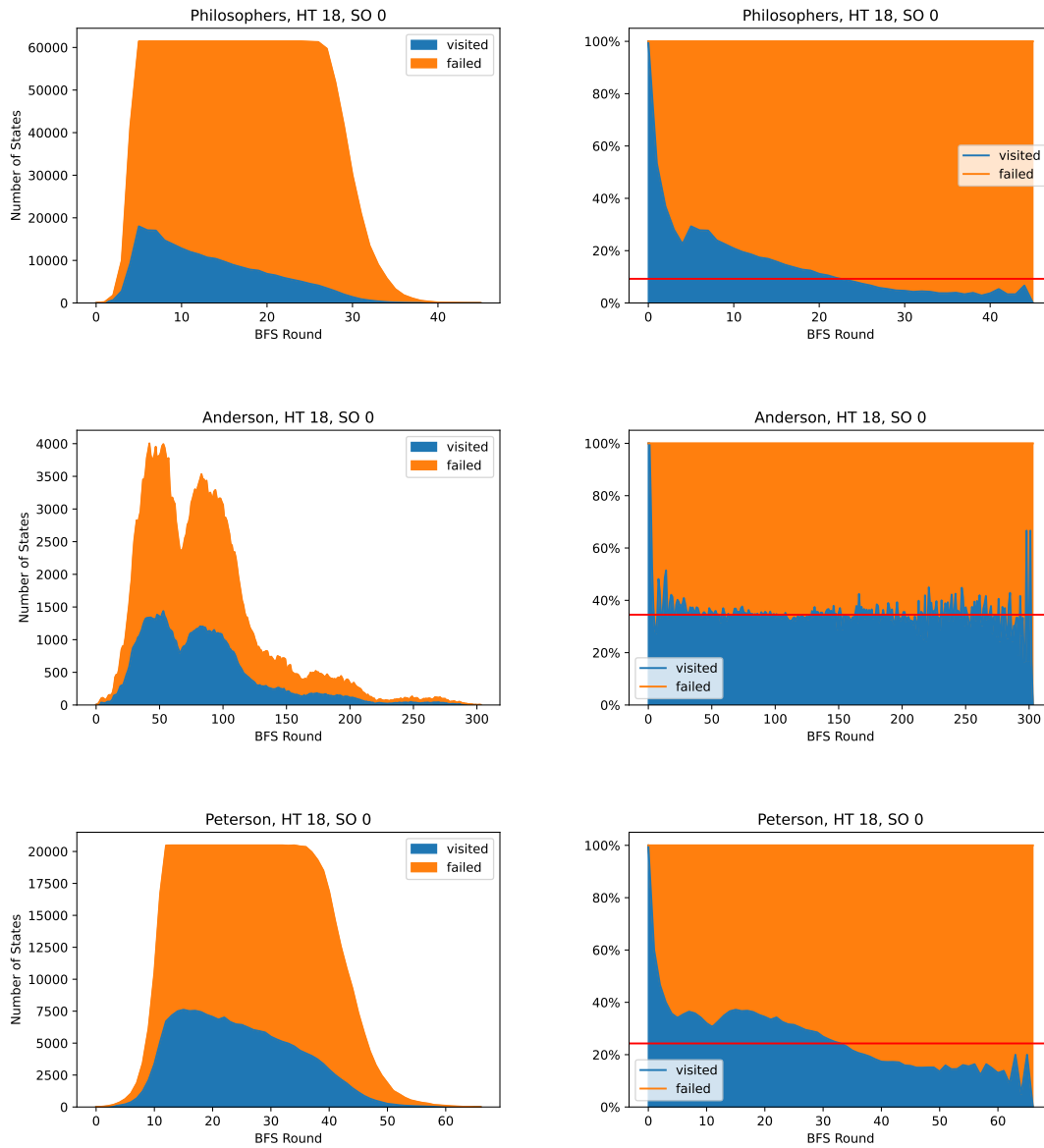


Figure 6.5: BFS frontier visualization of low-connectivity models

6.2.4 Start Over Strategy on Low-Connectivity Models

Our last experiment evaluates whether the start over strategy yields the same performance increase shown in Section 6.1.3 on low-connectivity models. For each of our three low-connectivity models, we conducted two explorations of 500 runs à 250 VTs. One with 2^{18} hash table slots and 0 start overs, the other one with 2^{14} hash table slots and 15 start overs (same configuration as in Section 6.1.3).

- Most interesting result: On the Anderson model, start overs significantly worsen the exploration performance
- DP and Peterson achieved similar results as for the waypoints model
- For all three models, the execution time of the start over strategy was considerably lower

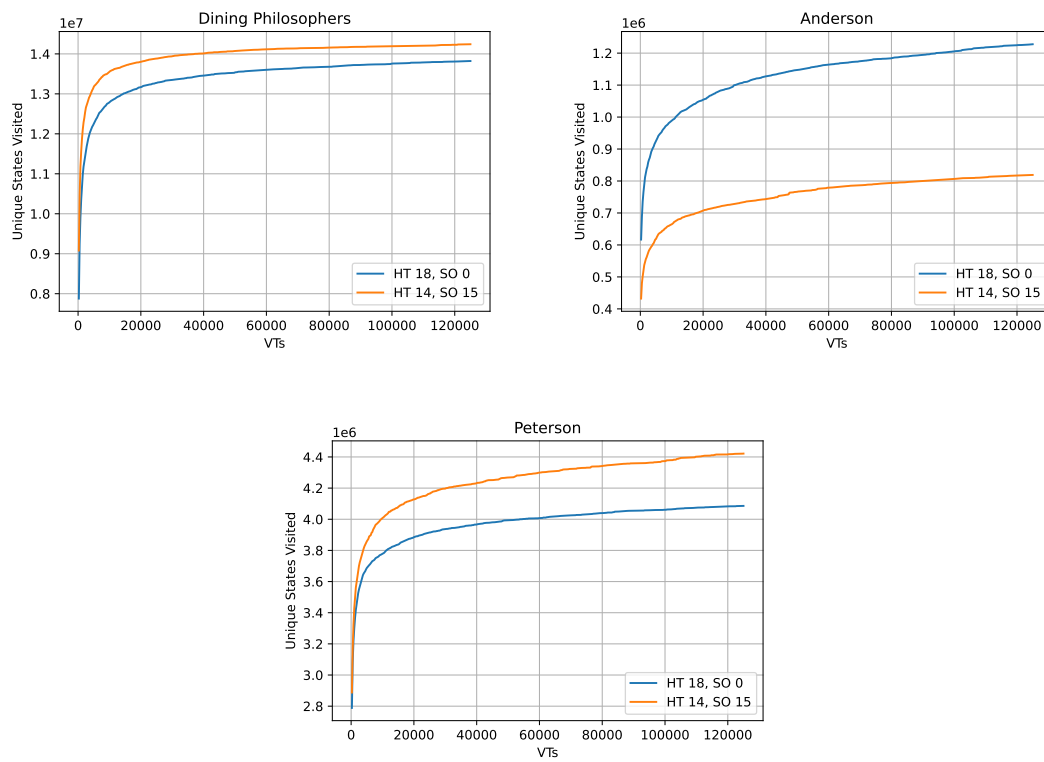


Figure 6.6: Start over strategy on low-connectivity models

7 Conclusion

7.1 Future Work

- Estimate state space size for models of unknown size
- Automatically determine the optimal number of start overs
- Implement usage of multiple GPU devices through CUDA streaming API
- Analyze State Space Diversification, i.e. "how good is a partition?"

7.2 Discussion

- As in the evaluation: The Anderson model draws special interest: Its BFS frontiers show a unique behavior and using the start over strategy, it performs contrary to all other models.

Bibliography

- [1] Ezio Bartocci, Richard DeFrancisco, and Scott A. Smolka. “Towards a GPGPU-Parallel SPIN Model Checker”. In: *Proceedings of the 2014 International SPIN Symposium on Model Checking of Software*. SPIN 2014. San Jose, CA, USA: Association for Computing Machinery, 2014, pp. 87–96. ISBN: 9781450324526. DOI: 10.1145/2632362.2632379.
- [2] Edmund M. Clarke, Thomas A. Henzinger, and Helmut Veith. “Introduction to Model Checking”. In: *Handbook of Model Checking*. Ed. by Edmund M. Clarke et al. Cham: Springer International Publishing, 2018, pp. 1–26. ISBN: 978-3-319-10575-8. DOI: 10.1007/978-3-319-10575-8_1.
- [3] Richard DeFrancisco et al. “Swarm model checking on the GPU”. In: *International Journal on Software Tools for Technology Transfer* 22.5 (Oct. 2020), pp. 583–599. ISSN: 1433-2787. DOI: 10.1007/s10009-020-00576-x.
- [4] Philippe Flajolet et al. “HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm”. In: *Discrete Mathematics & Theoretical Computer Science DMTCS Proceedings* vol. AH, 2007 Conference on Analysis of Algorithms (AofA 07) (Jan. 2007). DOI: 10.46298/dmtcs.3545.
- [5] Stefan Heule, Marc Nunkesser, and Alexander Hall. “HyperLogLog in Practice: Algorithmic Engineering of a State of the Art Cardinality Estimation Algorithm”. In: *Proceedings of the 16th International Conference on Extending Database Technology*. EDBT ’13. Genoa, Italy: Association for Computing Machinery, 2013, pp. 683–692. ISBN: 9781450315975. DOI: 10.1145/2452376.2452456.
- [6] Gerard J. Holzmann. “Explicit-State Model Checking”. In: *Handbook of Model Checking*. Ed. by Edmund M. Clarke et al. Cham: Springer International Publishing, 2018, pp. 153–171. ISBN: 978-3-319-10575-8. DOI: 10.1007/978-3-319-10575-8_5.
- [7] Gerard J. Holzmann. “Parallelizing the Spin Model Checker”. In: *Model Checking Software*. Ed. by Alastair Donaldson and David Parker. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 155–171. ISBN: 978-3-642-31759-0. DOI: 10.1007/978-3-642-31759-0_12.
- [8] Gerard J. Holzmann, Rajeev Joshi, and Alex Groce. “Swarm Verification”. In: *2008 23rd IEEE/ACM International Conference on Automated Software Engineering*. 2008 23rd IEEE/ACM International Conference on Automated Software Engineering. 2008, pp. 1–6. DOI: 10.1109/ASE.2008.9.
- [9] Gerard J. Holzmann, Rajeev Joshi, and Alex Groce. “Swarm Verification Techniques”. In: *IEEE Transactions on Software Engineering* 37.6 (2011), pp. 845–857. DOI: 10.1109/TSE.2010.110.