

ESBMC-GPU

A Context-Bounded Model Checking Tool to Verify CUDA Programs

Felipe R. Monteiro¹, Erickson H. da S. Alves¹, Isabela S. Silva¹,
Hussama A. Ismail¹, Lucas C. Cordeiro^{1,2}, and Eddie B. de Lima Filho^{1,3}

¹*Faculty of Technology, Federal University of Amazonas, Brazil*

²*Department of Computer Science, University of Oxford, United Kingdom*

³*FPF Tech, Brazil*

Abstract

The Compute Unified Device Architecture (CUDA) is a programming model used for exploring the advantages of Graphics Processing Unit (GPU) devices, through parallelization and specialized functions and features. Nonetheless, as in other development platforms, errors may occur, due to traditional software creation processes, which may even compromise the execution of an entire system. In order to address such a problem, ESBMC-GPU was developed, as an extension to the Efficient SMT-Based Context-Bounded Model Checker (ESBMC). In summary, ESBMC processes input code through ESBMC-GPU and an abstract representation of the standard CUDA libraries, with the goal of checking a set of desired properties. Experimental results showed that ESBMC-GPU was able to correctly verify 85% of the chosen benchmarks and it also overcame other existing GPU verifiers.

Keywords: GPU verification, formal verification, model checking, CUDA

1. Introduction

The Compute Unified Device Architecture (CUDA) is a development framework that makes use of the architecture and processing power of Graphics Processing Units (GPUs) [1]. Indeed, CUDA is also an Application Programming Interface (API), through which a GPU's parallelization scheme and tools can be accessed, with the goal of executing kernels. Nonetheless, source code is still written by human programmers, which may result in arithmetic overflow, division by zero, and other violation types. In addition, given that CUDA allows parallelization, problems related to the latter can also occur, due to thread scheduling [2].

11 In order to address the mentioned issues, an extension to the Efficient
 12 SMT-Based Context-Bounded Model Checker (ESBMC) [4] was developed,
 13 named as ESBMC-GPU [5, 6, 7], with the goal of verifying CUDA-based pro-
 14 grams (available online at <http://esbmc.org/gpu>). ESBMC-GPU consists
 15 of an extension for parsing CUDA source code (*i.e.*, a front-end to ESBMC)
 16 and a CUDA operational model (COM), which is an abstract representation
 17 of the standard CUDA libraries (*i.e.*, the native API) that conservatively
 18 approximates their semantics.

19 A distinct feature of ESBMC-GPU, when compared with other approaches
 20 [2, 8, 9, 10], is the use of Bounded Model Checking (BMC) [11] allied to
 21 Satisfiability Modulo Theories (SMT) [12], with explicit state-space explo-
 22 ration [3, 4]. In summary, concurrency problems are tackled, up to an un-
 23 winding bound, while each interleaving itself is symbolically handled; how-
 24 ever, even with BMC, space-state exploration may become a very time-
 25 consuming task, which is alleviated through state hashing and Monotonic
 26 Partial Order Reduction (MPOR) [13]. As a consequence, redundant inter-
 27 leavings are eliminated, without ignoring a program’s behavior.

28 Finally, existing GPU verifiers often ignore some aspects related to mem-
 29 ory leak, data transfer, and overflow, which are normally present in CUDA
 30 programs. The proposed approach, in turn, explicitly addresses them, through
 31 an accurate checking procedure, which even considers data exchange between
 32 main program and kernel. Obviously, it results in higher verification times,
 33 but more errors can then be identified and later corrected, in another devel-
 34 opment cycle.

35 2. Architecture and Implementation

36 ESBMC-GPU builds on top of ESBMC, which is an open source context-
 37 bounded model checker based on SMT solvers for ANSI-C/C++ programs [3,
 38 4], and adds four essential models, as described below.

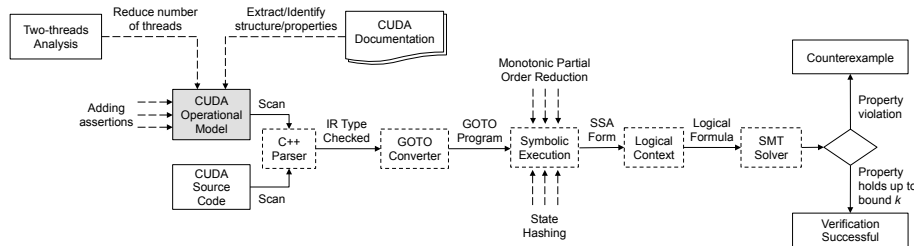


Figure 1: Overview of ESBMC-GPU’s architecture.

39 **1. CUDA Operational Model.** An operational model for CUDA libraries
 40 that provides support to CUDA functionalities, in conjunction with ESBMC,
 41 as shown in Fig. 1. Such an approach, which was previously attempted in the

42 verification of C++ programs [14, 15, 16, 17], consists in an abstract repre-
 43 sentation that reliably approximates the CUDA library’s semantics; however,
 44 COM incorporates pre- and post-conditions into verification processes, which
 45 enables ESBMC-GPU to verify specific properties (cf. Sec. 3). Indeed, COM
 46 allows the necessary control for performing code analysis, where both CUDA
 47 operation and knowledge for model checking its properties are available.

48 ESBMC was designed to handle multi-threaded software, through the use
 49 of an API called Portable Operating System Interface (POSIX – ISO/IEC
 50 9945) [18]. Thus, ESBMC-GPU applies a combination of processing methods
 51 used by Central Processing Units (CPUs) and the POSIX library, where
 52 thread instructions can interleave to create execution paths. Particularly,
 53 COM simulates the behavior of kernel calls using pthread functions (*e.g.*,
 54 `pthread_create`) and combines that with ESBMC, in order to check data
 55 race and specific C/C++ programming language failures (*e.g.*, array out-of-
 56 bounds and pointer safety).

57 **2. Two-threads Analysis.** Similarly to GPUVerify [2] and PUG [8],
 58 ESBMC-GPU also reduces the number of threads (to only two elements),
 59 during the verification of CUDA programs, by considering a NVIDIA Fermi
 60 GPU architecture, in order to improve verification time and avoid the state-
 61 space explosion problem. Besides, in CUDA programs, whilst threads exe-
 62 cute the same parametrized kernel, only two of them are necessary for conflict
 63 check. Thus, such an analysis ensures that errors (*e.g.*, data races) detected
 64 between two threads, in a given subgroup and due to unsynchronized accesses
 65 to shared variables, are enough to justify a property violation [7].

66 **3. State Hashing.** ESBMC-GPU applies state hashing to further eliminate
 67 redundant interleavings and also reduce state space, based on SHA256 hashes
 68 [19]. In particular, its symbolic state hashing approach computes a summary
 69 for a particular state that has already been explored and then indexes the
 70 resulting set, in order to reduce the generation of redundant states. Given
 71 any state computed during the symbolic execution of a specific CUDA kernel,
 72 ESBMC-GPU simply summarizes it and efficiently determines whether it has
 73 been explored before or not, along a different computation path. When this
 74 behavior is confirmed, which happens during the ESBMC-GPU’s symbolic-
 75 execution procedure, then the current computation path does not need to
 76 be further explored in the associated reachability tree (RT). This way, if
 77 ESBMC-GPU reaches such a state, *i.e.*, where a context switch can be taken
 78 (*e.g.*, before a global variable or synchronization primitive) and all shared/lo-
 79 cal variables and program counters are similar to another explored node, then
 80 ESBMC-GPU just considers that an identical node to be further explored,
 81 since reachability subtrees associated to them are also similar [7, 20].

82 **4. Monotonic Partial Order Reduction.** MPOR is used to reduce the
 83 number of thread interleavings, by classifying transitions inside a program

84 as dependent or independent. As a consequence, it is possible to determine
 85 whether interleaving pairs always lead to the same state and then remove du-
 86 plicates in a reachability tree, without ignoring any program’s behavior [20].

87 3. Functionalities

88 Through the integration of COM into ESBMC (*i.e.*, ESBMC-GPU), one
 89 is able to analyze CUDA programs and verify the following properties: data-
 90 race conditions, in order to detect if multiple threads perform unsynchronized
 91 access to shared-memory locations; pointer safety, *i.e.*, whether (i) a pointer
 92 offset does not exceed object bounds and (ii) a pointer is neither NULL nor
 93 invalid; array bounds, in order to ensure that array indices are within known
 94 bounds; arithmetic under- and overflow, which happens when a sum or prod-
 95 uct exceeds the memory limits that a variable can handle; division by zero,
 96 which takes place when denominators, in arithmetic expressions, lead to a
 97 division by zero; and user-specified assertions, *i.e.*, all assertions specified by
 98 users, which is essential to a thorough verification process.

99 In order to check the aforementioned properties, ESBMC-GPU explicitly
 100 explores the possible interleavings (up to the given context bound) and calls
 101 the single-threaded BMC procedure on each one, whenever it reaches an RT
 102 leaf node. Then, the mentioned procedure will stop if it finds a bug or when
 103 all possible RT interleavings has been systematically explored [7].

104 4. Illustrative Example

105 In this part, ESBMC-GPU usage is demonstrated, by using the CUDA
 106 program shown in Fig. 2. First of all, users must replace the default kernel call
 107 (line 16) by an intrinsic function of ESBMC-GPU (line 17). Then, the result-
 108 ing CUDA program can be passed to the command-line version of ESBMC-
 109 GPU, as follows: `esbmc-gpu <file>.cu --unwind <k> --context-switch`
 110 `<c> --state-hashing -I <path-to-CUDA-OM>`, where `<file>.cu` is the
 111 CUDA program, `<k>` is the maximum loop unrolling, `<c>` is a context-switch
 112 bound, `--state-hashing` reduces redundant interleavings, and `<path-to-`
 113 `CUDA-OM>` is the location of the COM library.

114 In the mentioned example, ESBMC-GPU detects an array out-of-bounds
 115 violation. Indeed, this CUDA-based program retrieves a memory region that
 116 has not been previously allocated, *i.e.*, when `threadIdx.x = 1`, the program
 117 tries to access `a[2]`. Importantly, the `cudaMalloc()` function’s operational
 118 model has a precondition that checks if the memory size to be allocated is
 119 greater than zero. In addition, an assertion checks if the result matches to the
 120 expected postcondition (line 19). Therefore, the verification of this program
 121 through ESBMC-GPU produces 54 successful and 3 failed interleavings. For
 122 instance, one possible failed interleaving is represented by the threads exe-
 123 cutions $t_0 : a[1] = 0$; $t_1 : a[2] = 1$, where $a[2] = 1$ represents an incorrect

124 access to the array index a . It is worth noticing that CIVL, ESBMC-GPU,
 125 and GKLEE are also able to detect this array out-of-bounds violation, but
 126 GPUVerify fails, as it reports a true incorrect result (missed bug).

```

1 #include <...>
2 #define BLOCKS 1
3 #define THREADS 2
4 --global-- void kernel(int *A) {
5     A[threadIdx.x + 1] = threadIdx.x;
6 }
7 int main(){
8     int *a;
9     int *dev_a;
10    int size = THREADS*sizeof(int);
11    a = (int*)malloc(size);
12    cudaMalloc((void**)&dev_a, size);
13    for (int i = 0; i < THREADS; i++)
14        a[i] = 0;
15    cudaMemcpy(dev_a, a, size, cudaMemcpyHostToDevice);
16    // kernel<<<BLOCKS,THREADS>>>(dev_a);
17    ESBMC_verify_kernel(kernel, BLOCKS, THREADS, dev_a);
18    for (int i = 0; i < THREADS; i++)
19        assert(a[i]==i);
20    cudaFree(dev_a);
21    return 0;
22 }

```

Figure 2: Illustrative CUDA code example.

127 5. Experimental Evaluation

128 In order to evaluate ESBMC-GPU’s precision and performance, bench-
 129 marks¹ were extracted from the available literature (*i.e.*, NVIDIA GPU
 130 Computing SDK v2.0 [21] and Microsoft C++ AMP Sample Projects [22]),
 131 which covers basic functions commonly used by real CUDA applications.
 132 The present experiments answer two research questions: *(i)* How accurate
 133 is ESBMC-GPU when verifying the chosen benchmarks? *(ii)* How does
 134 ESBMC-GPU’s performance compare to other existing verifiers?

135 In order to answer both questions, all benchmarks were verified with 4
 136 GPU verifiers (ESBMC-GPU v2.0, GKLEE v2012, GPUVerify v1811, and
 137 CIVL v1.7.1), on an otherwise idle Intel Core i7-4790 CPU 3.60 GHz, with
 138 16 GB of RAM, running Ubuntu 14.04 OS. Importantly, all presented ex-
 139 ecution times are actually CPU times, *i.e.*, only the elapsed time periods
 140 spent in the allocated CPUs, which was measured with the `times` system
 141 call (POSIX system). An overview of the experimental results is shown in
 142 Fig. 3, where *True* represents bug-free benchmarks, *False* represents buggy
 143 benchmarks, *Not supported* represents benchmarks that could not be veri-
 144 fied, *Correct* represents the percentage of benchmarks correctly verified, and
 145 *Incorrect* represents the percentage of benchmarks incorrectly verified (*i.e.*,
 146 a verification tool reports an unexpected result). As one may notice, the
 147 present experimental results show that ESBMC-GPU reached a successful

¹A detailed description of all benchmarks is available at <http://esbmc.org/gpu/>

148 verification rate of approximately 85%, while GKLEE, GPUVerify, and CIVL
 149 reported 72%, 50%, and 35%, respectively². More precisely, ESBMC-GPU
 150 supports the verification of benchmarks related to array bounds (3%), as-
 151 sertive statements (5%), data race (11%), NULL pointers (3%), and other
 152 specific CUDA functionalities (63%).

153 **Limitations.** ESBMC-GPU was unable to correctly verify 24 bench-
 154 marks, which are related to constant memory access (2%), CUDA’s specific
 155 libraries (4.5%), and the use of pointers to functions, structures, and `char`
 156 type variables, when passed as kernel call arguments (4.5%). In addition, it
 157 only reported 3% of incorrect true and 1% of incorrect false results.

158 **Performance.** MPOR resulted in a performance improvement of ap-
 159 proximately 80%, by decreasing the verification time from 16 to 3 hours, while
 160 the two-threads analysis further reduced that to 789.6 sec. Although such
 161 techniques have considerably improved the ESBMC-GPU’s performance, it
 162 still takes longer than the other evaluated tools: GPUVerify (98.36 sec),
 163 GKLEE (108.32 sec), and CIVL (708.52 sec). This is due to thread inter-
 164 leavings, which combine symbolic model checking with explicit state-space
 165 exploration [7]. In addition, ESBMC-GPU still presents the highest accuracy,
 166 with less than 6 seconds per benchmark.

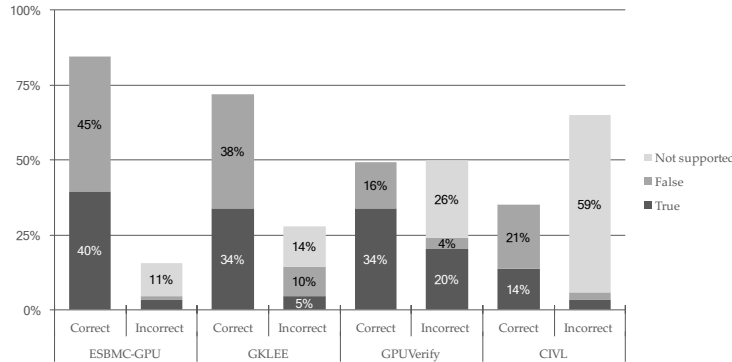


Figure 3: Experimental evaluation of ESBMC-GPU against other verifiers.

167 6. Conclusions and Future Work

168 ESBMC-GPU marks the first application of an SMT-based context-BMC
 169 tool that recognizes CUDA directives [7]. Besides, it further simplifies veri-
 170 fication models and provides fewer incorrect results, if compared with GK-
 171 LEE, GPUVerify, and CIVL. Finally, it presents improved ability to detect
 172 array out-of-bounds and data race violations. Future work aims to support
 173 stream interleaving and implement further techniques to reduce the number
 174 of thread interleavings, by taking into account GPU symmetry.

²All experimental results are available at <http://esbmc.org/gpu/>

175 References

- 176 [1] Cheng, J., Grossman, M., and McKercher, T. *Professional CUDA C*
177 *Programming*. John Wiley and Sons, Inc. 2014.
- 178 [2] Betts, A., Chong, N., Donaldson, A., Qadeer, S., and Thomson, P.
179 *GPUVerify: A Verifier for GPU Kernels*. OOPSLA 2012; 113–132.
- 180 [3] Cordeiro, L. and Fischer, B. *Verifying Multi-threaded Software using*
181 *SMT-based Context-Bounded Model Checking*. ICSE 2011; 331–340.
- 182 [4] Cordeiro, L. and, Fischer, B., and Marques-Silva, J. *SMT-Based*
183 *Bounded Model Checking for Embedded ANSI-C Software*. IEEE Trans.
184 *Software Eng.* 2012; 38(4):957–974.
- 185 [5] Pereira, P., Albuquerque, H., Marques, H., Silva, I., Carvalho, C., San-
186 tos, V., Ferreira, R., and Cordeiro, L. *Verificação de Kernels em Pro-*
187 *gramas CUDA usando Bounded Model Checking*. WSCAD-SSC 2015;
188 24–35.
- 189 [6] Pereira, P., Albuquerque, H., Marques, H., Silva, I., Carvalho, C., San-
190 tos, V., Ferreira, R., and Cordeiro, L. *Verifying CUDA Programs using*
191 *SMT-Based Context-Bounded Model Checking*. SAC SVT track 2016;
192 1648–1653.
- 193 [7] Pereira, P., Albuquerque, H., Silva, I., Marques, H., Monteiro, F., Fer-
194 reira, R., and Cordeiro, L. *SMT-Based Context-Bounded Model Checking*
195 *for CUDA Programs*. Concurrency Computat.: Pract. Exper. 2016 (to
196 appear).
- 197 [8] Li, G. and Gopalakrishnan, G. *Scalable SMT-based Verification of GPU*
198 *Kernel Functions*. FSE 2010; 187–196.
- 199 [9] Li, G., Li, P., Sawaya, G., Gopalakrishnan, G., Ghosh, I., and Rajan, S.
200 *GKLEE: Concolic Verification and Test Generation for GPUs*. PPOPP
201 2012; 215–224.
- 202 [10] Zheng, M., Rogers, M., Luo, Z., Dwyer, M., and Siegel, S. *CIVL: Formal*
203 *Verification of Parallel Programs*. ASE 2015; 830–835.
- 204 [11] Biere, A. *Bounded Model Checking. Handbook of Satisfiability*. IOS Press
205 2009; 457–481.
- 206 [12] Barrett C, Sebastiani R, Seshia S, Tinelli C. *Satisfiability Modulo The-*
207 *ories. Handbook of Satisfiability*. IOS Press 2009; 825–885.

- 208 [13] Kahlon V, Wang C, Gupta A. *Monotonic Partial Order Reduction: An*
 209 *Optimal Symbolic Partial Order Reduction Technique*. CAV 2009; 398–
 210 413.
- 211 [14] Ramalho, M., Freitas, M., Sousa, F., Marques, H., Cordeiro, L., and
 212 Fischer, B. *SMT-Based Bounded Model Checking of C++ Programs*.
 213 ECBS 2013; 147–156.
- 214 [15] Monteiro, F., Cordeiro, L., and de Lima Filho, E. *Bounded Model Check-*
 215 *ing of C++ Programs Based on the Qt Framework*. GCCE 2015; 179–
 216 447.
- 217 [16] Garcia, M., Monteiro, F., Cordeiro, L., and de Lima Filho, E.
 218 *ESBMC^{QtOM}: A Bounded Model Checking Tool to Verify Qt Applica-*
 219 *tions*. SPIN 2016; 97–103.
- 220 [17] Monteiro, F., Garcia, M., Cordeiro, L., and de Lima Filho, E. *Bounded*
 221 *Model Checking of C++ Programs based on the Qt Cross-Platform*
 222 *Framework*. Software Testing, Verification & Reliability. 2017; (to ap-
 223 pear).
- 224 [18] Institute of Electrical and Electronics Engineers, Inc. *IEEE Stan-*
 225 *dard for Information Technology - Portable Operating System Interface*
 226 *(POSIX) Base Specifications, IEEE Std 1003.1-2008 (Revision of IEEE*
 227 *Std 1003.1-2004)*. IEEE 2008.
- 228 [19] Federal Information Processing Standard 180-2. *Secure Hash Standard*.
 229 National Institute of Standards and Technology, 2002.
- 230 [20] Morse, J. *Expressive and Efficient Bounded Model Checking of Concur-*
 231 *rent Software*. University of Southampton, PhD Thesis 2015.
- 232 [21] NVIDIA. *CUDA Toolkit Release Archive*.
 233 <https://developer.nvidia.com/cuda-toolkit-archive>, 2015.
- 234 [22] Microsoft Corporation. *C++ AMP sample projects for download (MSDN blog)*.
 235 [blogs.msdn.com/b/nativeconcurrency/archive/2012/01/30/c-amp-](https://blogs.msdn.com/b/nativeconcurrency/archive/2012/01/30/c-amp-sample-projects-for-download.aspx)
 236 [sample-projects-for-download.aspx](https://blogs.msdn.com/b/nativeconcurrency/archive/2012/01/30/c-amp-sample-projects-for-download.aspx), 2012.

238 **Required Metadata**

239 **Current executable software version**

240 Ancillary data table required for sub version of the executable software.

Nr.	(executable) Software metadata description	Please fill in this column
S1	Current software version	2.0
S2	Permanent link to executables of this version	http://esbmc.org/gpu/
S3	Legal Software License	Apache v2.0
S4	Computing Operating System	Ubuntu Linux OS
S5	Installation requirements & dependencies	GNU Libtool; Automake; Flex & Bison; Boost C++ Libraries; Multi-precision arithmetic library developers tools (libgmp3-dev package); SSL development libraries (libssl-dev package); CLang 3.8; LLDB 3.8; GNU C++ compiler (multilib files); libc6 and libc6-dev packages
S6	Link to user manual	http://esbmc.org/gpu/
S7	Support email for questions	lucas.cordeiro@cs.ox.ac.uk

Table 1: Software metadata (optional)

241 **Current code version**

242 Ancillary data table required for subversion of the codebase.

Nr.	Code metadata description	Please fill in this column
C1	Current code version	<i>v2.0</i>
C2	Permanent link to code/repository used for this code version	https://github.com/ssvlab/esbmc-gpu
C3	Legal Code License	GNU Public License
C4	Code versioning system used	git
C5	Software code languages, tools, and services used	C++
C6	Compilation requirements, operating environments & dependencies	GNU Libtool; Automake; Flex & Bison; Boost C++ Libraries; Multi-precision arithmetic library developers tools (libgmp3-dev package); SSL development libraries (libssl-dev package); CLang 3.8; LLDB 3.8; GNU C++ compiler (multilib files); libc6 and libc6-dev packages
C7	Link to developer documentation	http://esbmc.org/gpu
C8	Support email for questions	lucas.cordeiro@cs.ox.ac.uk

Table 2: Code metadata (mandatory)