

# The Design and Implementation of a Verification Technique for GPU Kernels

ADAM BETTS, Imperial College London  
NATHAN CHONG, Imperial College London  
ALASTAIR F. DONALDSON, Imperial College London  
JEROEN KETEMA, Imperial College London  
SHAZ QADEER, Microsoft Research Redmond  
PAUL THOMSON, Imperial College London  
JOHN WICKERSON, Imperial College London

We present a technique for the formal verification of GPU kernels, addressing two classes of correctness properties: data races and barrier divergence. Our approach is founded on a novel formal operational semantics for GPU kernels termed *synchronous, delayed visibility (SDV)* semantics, which captures the execution of a GPU kernel by multiple groups of threads. The SDV semantics provides operational definitions for barrier divergence and for both inter- and intra-group data races. We build on the semantics to develop a method for reducing the task of verifying a massively parallel GPU kernel to that of verifying a sequential program. This completely avoids the need to reason about thread interleavings, and allows existing techniques for sequential program verification to be leveraged. We describe an efficient encoding of data race detection and propose a method for automatically inferring the loop invariants that are required for verification. We have implemented these techniques as a practical verification tool, GPUVerify, that can be applied directly to OpenCL and CUDA source code. We evaluate GPUVerify with respect to a set of 162 kernels drawn from public and commercial sources. Our evaluation demonstrates that GPUVerify is capable of efficient, automatic verification of a large number of real-world kernels.

Categories and Subject Descriptors: F.3.1 [**Logics and Meanings of Programs**]: Specifying and Verifying and Reasoning about Programs

General Terms: Theory, Verification

Additional Key Words and Phrases: Verification, GPUs, concurrency, data races, barrier synchronization

## ACM Reference Format:

Betts, A., Chong, N., Donaldson, A. F., Ketema, J., Qadeer, S., Thomson, P., Wickerson, J. 2014. The Design and Implementation of a Verification Technique for GPU Kernels. *ACM Trans. Program. Lang. Syst.* 0, 0, Article 0 (December 2014), 49 pages.

DOI : <http://dx.doi.org/10.1145/0000000.0000000>

## 1. INTRODUCTION

In recent years, massively parallel *accelerator* processors, primarily graphics processing units (GPUs) from companies such as AMD, Nvidia and ARM, have become widely

---

This work was supported by the EU FP7 STREP project CARP (project number 287767), by EPSRC project EP/G051100/2, and by two EPSRC-funded PhD studentships. Part of the work was carried out while Alastair Donaldson was a Visiting Researcher at Microsoft Research Redmond. This paper is a revised and extended version of [Betts et al. 2012].

Author's addresses: A. Betts, N. Chong, A. F. Donaldson, J. Ketema, P. Thomson, J. Wickerson, Department of Computing, Imperial College London, UK; S. Qadeer, Microsoft Research, Redmond, USA.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2014 ACM 0164-0925/2014/12-ART0 \$15.00

DOI : <http://dx.doi.org/10.1145/0000000.0000000>

available to end-users. Accelerators offer significant compute power at a low cost, and algorithms in domains such as medical imaging [Cates et al. 2004], computer vision [Salas-Moreno et al. 2013], computational fluid dynamics [Harris 2004] and DNA sequence alignment [Li and Durbin 2009; Klus et al. 2012] can be accelerated to beat CPU performance, sometimes by orders of magnitude. Importantly, GPUs enable performance improvements not only in terms of run-time, but also energy efficiency [Lokhmotov 2011].

GPUs present a serious challenge for software developers. A system may contain one or more of the plethora of devices on the market. Applications must exhibit *portable correctness*, operating correctly on *any* GPU accelerator. Software bugs in media, entertainment and business domains can have serious financial implications, and bugs in safety critical domains, such as medical applications, can be more drastic still. Because GPUs are being used increasingly in such areas, there is an urgent need for verification techniques to aid the construction of correct GPU software.

In this work we address the problem of static verification of GPU kernels written using mainstream programming models such as OpenCL [Khronos OpenCL Working Group 2012], CUDA [Nvidia 2012a] and C++ AMP [Microsoft Corporation 2012]. We focus on two classes of bugs that make writing correct GPU kernels harder than writing correct sequential code: *data races* and *barrier divergence*.

In contrast to the well-understood notion of data races, there does not appear to be a formal definition of barrier divergence for GPU programming. Our work gives a precise characterization of barrier divergence via an operational semantics based on predicated execution, which we call *synchronous, delayed visibility (SDV)* semantics. The SDV semantics models the execution of a kernel by multiple groups of threads, and provides operational definitions for inter- and intra-group data races as well as barrier divergence. While predicated execution has been used for code generation by GPU kernel compilers, our work is the first to use predicated operational semantics for the purpose of specification and verification.

Founded on the SDV semantics, we present a verification method which transforms a massively parallel GPU kernel into a sequential program such that correctness of the sequential program implies data race- and barrier divergence-freedom of the kernel. Our method achieves scalability by exploiting the fact that data races and barrier divergence are *pairwise* properties, allowing a reduction to an arbitrary pair of threads. For this reduction to be sound, it is necessary to use abstraction to model the possible effects of other threads. We consider two practical abstraction methods, the *adversarial* abstraction and the more refined *equality* abstraction, and prove formally that the reduction to a pair of threads is sound when combined with these abstractions.

Reducing the analysis of a highly concurrent program to the verification of a sequential program completely avoids the need to reason about thread interleavings, and enables reuse of existing verification techniques for sequential programs. We present novel heuristics for automatically inferring loop invariants required for verification.

We discuss the design and implementation of GPUVerify, a practical tool for verifying GPU kernels using our method, which can be applied directly to OpenCL and CUDA source code. GPUVerify builds on the Clang/LLVM compiler infrastructure<sup>1</sup>, the Boogie verification system [Barnett et al. 2005] and the Z3 theorem prover [de Moura and Bjørner 2008].

We have used GPUVerify to analyze a set of 162 OpenCL and CUDA kernels. We summarize an experiment with the initial implementation of GPUVerify (which was limited to intra-group race checking [Betts et al. 2012]) where we divided these kernels into separate *training* and *evaluation* sets, tuned the capabilities of the tool to the

<sup>1</sup><http://clang.llvm.org/>

training set and then ran the tool *blindly* on the evaluation set (none of our team were familiar with the kernels in the evaluation set). We found that fully automatic verification was achieved for 49 out of 70 kernels from the evaluation set (69%). We also present a comparison of the initial GPUVerify implementation with PUG, a formal analysis tool for CUDA kernels by Li and Gopalakrishnan [2010]. GPUVerify performs competitively with PUG when verifying correct kernels and rejects defective kernels in several cases where PUG reports false negatives. Additionally, the equality abstraction offered by GPUVerify is more refined than the shared state abstraction supported by PUG, allowing verification of real-world kernels for which PUG reports false positives.

We also report on a more recent version of GPUVerify, which supports inter-group race checking, to analyze our benchmarks, using a combination of improved invariant inference and manual invariant annotations to achieve verification across the entire benchmark set. We report on the performance of GPUVerify, assess the overhead associated with performing inter-group race checking on top of intra-group race checking, and compare the performance and precision of GPUVerify when the adversarial vs. equality abstractions are used.

In summary, our paper makes the following contributions:

- *Synchronous, delayed visibility* (SDV) semantics: a formal operational semantics for GPU kernels based on predicated execution, inter- and intra-group data-race freedom, and barrier divergence freedom.
- The *two-thread reduction*: an abstract version of the SDV semantics that models a pair of threads, and a soundness theorem showing that data race- and barrier divergence-freedom of a massively parallel kernel according to the SDV semantics can be established by proving data race- and barrier divergence-freedom for each pair of threads using the abstract semantics.
- A verification method for GPU kernels which exploits the SDV semantics and two-thread reduction to transform analysis to a sequential program verification task, allowing reuse of existing analysis techniques based on verification condition generation and automated theorem proving.
- A method for automatically inferring the invariants needed for our verification method.
- An extensive evaluation of our verifier on a collection of 162 publicly available and commercial GPU kernels.

### Contributions in Relation to Prior Work

We make explicit how the current article significantly extends a previous conference version [Betts et al. 2012] and complements other papers that have arisen from the GPUVerify project [Collingbourne et al. 2013; Chong et al. 2013; Bardsley and Donaldson 2014; Bardsley et al. 2014a; Chong et al. 2014].

*Treatment of Multi-Group GPU Kernels.* The original presentation of our verification method was restricted to the analysis of single-group GPU kernels. Here we treat the multi-group case by extending the formal semantics for kernels, implementing inter-group race checking in the GPUVerify tool, and conducting an experimental evaluation to assess the overhead associated with handling multiple groups vs. just a single group.

*A Full Soundness Proof for the Two-Thread Reduction.* The scalability of our verification method hinges on reducing the analysis to consider an arbitrary pair of threads. Only a short proof sketch for the soundness of this “two-thread reduction” was presented in our prior work. In the current article we present, for the first time, a detailed formal soundness proof. This is an important contribution, because a similar reduction

is also employed by other methods, such as the PUG verifier of Li and Gopalakrishnan [2010].

*Mechanical Checking of the Language and Semantics Definitions.* We have used the Isabelle proof assistant to type check all the definitions, semantic rules and theorem statements appearing in the paper. The result is published in the Archive of Formal Proofs [Wickerson 2014]. We have also conducted a partial mechanization of the soundness proof for the two-thread reduction.

*New Experimental Evaluation.* Rather than repeating the extensive experimental evaluation presented in our prior work [Betts et al. 2012], we instead summarize the results of this evaluation and present a new evaluation in which we assess the verification performance of a more recent version of the GPUVerify tool, as well as evaluating the overhead associated with inter-group race checking and performance characteristics of two shared state abstractions used by the tool.

*Relation to Other Papers Arising from the GPUVerify Project.* In addition to the original version of the paper [Betts et al. 2012], we have published three technical articles that extend the GPUVerify verification method, each in a distinct manner [Collingbourne et al. 2013; Chong et al. 2013; Bardsley and Donaldson 2014]. Each of these extensions is separate from, and complementary to, the new contributions made in the current article.

Following our prior work [Betts et al. 2012], the results presented in this paper apply to structured GPU kernels expressed using `if` and `while` constructs. In related work, we have generalized the lock-step execution semantics used by our method to apply to unstructured control flow graphs, allowing the technique to be applied at the level of the intermediate representation used by a compiler [Collingbourne et al. 2013].

We have proposed *barrier invariants* [Chong et al. 2013] as a mechanism to allow users to improve on the precision offered by the *adversarial* and *equality* abstractions used in this work. The shared state abstractions discussed here are straightforward to apply automatically, but are inflexible. Barrier invariants are instead provided by a user; they offer a tunable, but heavy-weight and manually-driven, approach to GPU kernel verification.

These works, like the current contribution, are restricted to a barrier synchronizing programming model. We have conducted a preliminary investigation into lifting the method beyond this model to cater for kernels in which threads communicate using atomic operations and warp-based synchronization [Bardsley and Donaldson 2014].

As well as technical articles that extend the core verification method, a recent tool paper discusses engineering experience and insight gained during the GPUVerify project [Bardsley et al. 2014a], and an invited article provides a tutorial overview of the technique [Donaldson 2014]. We have also applied GPUVerify as a component in an automatic verification method for parallel prefix sum kernels [Chong et al. 2014].

## Reproducibility

The GPUVerify tool is open source and available online:

<http://multicore.doc.ic.ac.uk/tools/GPUVerify>

On the tool web page we make available the versions of GPUVerify used for the experiments reported in this paper, together with the non-commercial benchmarks used for our evaluation.

## 2. GPU KERNEL PROGRAMMING

A typical GPU (see Figure 1) consists of a large number of simple *processing elements* (*PEs*), sometimes referred to as *cores*. Subsets of the PEs are grouped together into

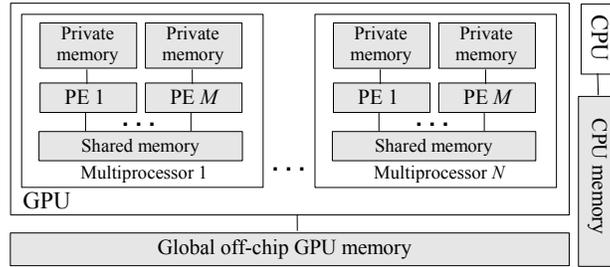


Fig. 1. Schematic overview of a typical GPU architecture

Table I. Equivalent terms for what we shall call *private*, *shared* and *global* memory in CUDA, OpenCL and C++ AMP

Term	CUDA	OpenCL	C++ AMP
private	local	private	local
shared	shared	local	tile static
global	global	global	global

Table II. Equivalent terms for *thread*, *group* and (where applicable) *sub-group* in CUDA, OpenCL and C++ AMP

Term	CUDA	OpenCL	C++ AMP
thread	thread	work-item	thread
group	thread block	work-group	tile
sub-group	warp	N/A	N/A

*multiprocessors*, such that all PEs within a multiprocessor execute in lock-step, in single instruction multiple data (SIMD) fashion. Distinct multiprocessors on a GPU can execute independently. Each PE is equipped with a small private memory, and PEs located on the same multiprocessor can access a portion of *shared* memory dedicated to that multiprocessor. All PEs on the GPU have access to a large amount of off-chip memory known as *global* memory, which is usually separate from main CPU memory.

Today, there are three major GPU programming models: OpenCL, an industry standard proposed by the Khronos OpenCL Working Group [2012] and widely supported (in particular, OpenCL is AMD's primary high-level GPU programming model); CUDA, from Nvidia [2012a]; and C++ AMP, from Microsoft [2012]. Table I summarizes the terms used in the programming models to refer to *private*, *shared* and *global* memory as described above.

*Threads and Groups.* All three programming models provide a similar high-level abstraction for mapping computation across GPU hardware, centered around the notion of a *kernel* program being executed by many parallel *threads*, together with a specification of how these threads should be partitioned into *groups*. The kernel is a template specifying the behavior of an arbitrary thread, parameterized by thread and group id variables. Expressions over these ids allow distinct threads to operate on separate data and follow differing execution paths through the kernel. Threads in the same group can synchronize during kernel execution, while threads in distinct groups execute completely independently.

The runtime environment associated with a GPU programming model must interface with the driver of the available GPU to schedule execution of kernel threads across processing elements. Typically each group of threads is assigned to one of the multiprocessors of the GPU, so that distinct groups execute in parallel on different multiprocessors. If the number of threads in a group is  $N$  and the number of PEs in a multipro-

Table III. Predicated forms for conditionals and loops

Program fragment	Predicated form
<pre> <b>if</b> (lid &gt; N)   x = 0; <b>else</b>   x = 1; </pre>	<pre> p = (lid &gt; N); p <math>\Rightarrow</math> x = 0; !p <math>\Rightarrow</math> x = 1; </pre>
<pre> <b>while</b> (i &lt; x) {   i++; } </pre>	<pre> p = (i &lt; x); <b>while</b> (<math>\exists t :: t.p</math>) {   p <math>\Rightarrow</math> i++;   p <math>\Rightarrow</math> p = (i &lt; x); } </pre>

cessor is  $M$ , then a group is divided into  $\lceil \frac{N}{M} \rceil$  *sub-groups*, each consisting of up to  $M$  threads. Execution of a single group on a multiprocessor then proceeds by scheduling the sub-groups in an interleaving manner. Each thread in a given sub-group is pinned to a distinct PE, and all threads in the same sub-group execute together in lock-step, following exactly the same control path. Distinct sub-groups may follow different control paths.

Table II summarizes the specific terms used by the three main GPU programming models to refer to *threads*, *groups* and (in the case of CUDA) *sub-groups*. OpenCL and C++ AMP aim for portability across GPUs from multiple vendors, so do not allow a kernel to query the device-specific size or structure of thread sub-groups.<sup>2</sup> As CUDA is Nvidia-specific, CUDA programmers *can* write kernels that make assumptions about the division of threads into sub-groups. However, such kernels will not easily port to general-purpose GPU programming languages, and may break when executed on future generations of Nvidia hardware using a different sub-group size. Thus GPU kernels that do not make assumptions about the size of thread sub-groups are preferable.

*Predicated Execution.* Recall that the PEs in a GPU multiprocessor execute in lock-step, as a SIMD processor array. Threads within a sub-group occupy the PEs of a multiprocessor, and thus must also execute in lock-step. Conditional statements and loops through which distinct threads in the same sub-group should take different paths must therefore be simulated, and this is achieved using *predicated execution*.

Consider the conditional statement in the top-left of Table III, where `lid` denotes the local id of a thread within its group and `x` is a local variable stored in private memory. This conditional can be transformed into the straight-line code shown in the top-right of the figure, which can be executed by a sub-group in lock-step. The meaning of a statement *predicate*  $\Rightarrow$  *command* is that a thread should execute *command* if *predicate* holds for that thread, otherwise the thread should execute a no-op. All threads evaluate the condition `lid > N` into a local Boolean variable `p`, then execute both the *then* and *else* branches of the conditional, predicated by `p` and `!p`, respectively.

Loops are turned into predicated form by dictating that all threads in a sub-group continue to execute the loop body until the loop condition is false for *all* threads in the sub-group, with threads for whom the condition does not hold becoming disabled. This is illustrated for the loop in the bottom-left of Table III (where `i` and `x` are local variables) by the code fragment shown in the bottom-right of the figure. First, the condition `i < x` is evaluated into local variable `p`. Then the sub-group loops while `p` remains true for some thread in the sub-group, indicated by  $\exists t : t.p$ . The loop body is predicated by `p`, and thus has an effect only for enabled threads.

We present precise operational semantics for predicated execution in Section 3.

<sup>2</sup>The recent OpenCL 2.0 extension specification [Khronos OpenCL Working Group 2014b, Section 9.17, p. 133] contains an optional extension to support sub-groups of threads.

*Barrier Synchronization.* When a thread  $t_1$  writes to an address in shared or global memory, the result of this write is not guaranteed to become visible to another thread  $t_2$  unless  $t_1$  and  $t_2$  *synchronize*. As noted above, there is *no* mechanism for threads in distinct groups to synchronize during kernel execution.<sup>3</sup> Threads in the same group can synchronize via *barriers*. Intuitively, a kernel thread belonging to group  $g$  waits at a *barrier* statement until every thread in  $g$  has reached the barrier. Passing the barrier guarantees that all writes to shared and global memory by threads in  $g$  occurring before execution of the barrier have been committed.

Our analysis through writing GPU kernels and talking to GPU developers is that there are two specific classes of bugs that make writing correct GPU kernels more difficult than writing correct sequential code: *data races* and *barrier divergence*.

## 2.1. Data Races

We distinguish between two kinds of data races in GPU kernels. An *inter-group data race* occurs if there are two threads  $t_1$  and  $t_2$  from *different* groups such that  $t_1$  writes to a location in global memory and  $t_2$  writes to or reads from this location. An *intra-group data race* occurs if there are two threads  $t_1$  and  $t_2$  from the same group such that  $t_1$  writes to a location in global or shared memory,  $t_2$  writes to or reads from this location, and no *barrier* statement is executed between these accesses. Races can lead to non-deterministic kernel behavior, and computation of incorrect results.

## 2.2. Barrier Divergence

If threads in the same group *diverge*, reaching different barriers as in the following kernel fragment:

```
if ((lid%2) == 0)
  barrier(); // Even threads hit first barrier
else
  barrier(); // Odd threads hit second barrier
```

then kernel behavior is undefined. According to CUDA [Nvidia 2012a]:

*“execution is likely to hang or produce unintended side effects.”*

While there is clarity across all programming models as to the meaning of barrier divergence in loop-free code, the situation is far from clear for code with loops. Consider the example kernel shown on the left of Figure 2. This kernel is intended to be executed by one group of four threads, and declares an array `A` of two shared buffers, each of size four. Private variable `buf` is an index into `A`, representing the *current* buffer. The threads execute a nest of loops. On each inner loop iteration a thread reads the value of the current buffer at index `lid+1` modulo 4 and writes the result into the non-current buffer at index `lid`. A barrier is used to avoid data races on `A`. Notice that local variables `x` and `y` are set to 4 and 1, respectively, for thread 0, and to 1 and 4, respectively, for all other threads. As a result, we expect thread 0 to perform four outer loop iterations, each involving one inner loop iteration, while other threads will perform a single outer loop iteration, consisting of four inner loop iterations.

According to the guidance in the CUDA documentation such a kernel appears to be valid: all threads will hit the barrier statement four times. Taking a snapshot of the array `A` at each barrier and at the end of the kernel, we might expect to see the

<sup>3</sup>In OpenCL 2.0, limited forms of inter-work-group synchronization can be implemented using atomic operations [Khronos OpenCL Working Group 2014a]. However, an inter-work-group barrier still cannot be reliably implemented due to lack of progress guarantees between groups.

<pre> <b>shared int</b> A[2][4];  <b>void</b> kernel() {   <b>int</b> buf, x, y, i, j;   x = (lid == 0 ? 4 : 1);   y = (lid == 0 ? 1 : 4);   buf = i = 0;   <b>while</b> (i &lt; x) {     j = 0;     <b>while</b> (j &lt; y) {       <b>barrier</b>();       A[1-buf][lid]         = A[buf][(lid+1)%4];       buf = 1 - buf;       j++;     }     i++;   } } </pre>	<pre> ... p = (i &lt; x); <b>while</b> (<math>\exists t::t.p</math>) {   p <math>\Rightarrow</math> j = 0;   q = p &amp;&amp; (j &lt; y);   <b>while</b> (<math>\exists t::t.q</math>) {     q <math>\Rightarrow</math> <b>barrier</b>();     q <math>\Rightarrow</math> A[1-buf][lid]       = A[buf][(lid+1)%4];     q <math>\Rightarrow</math> buf = 1 - buf;     q <math>\Rightarrow</math> j++;     q <math>\Rightarrow</math> q = p &amp;&amp; (j &lt; y);   }   p <math>\Rightarrow</math> i++;   p <math>\Rightarrow</math> p = (i &lt; x); } </pre>
---	---

Fig. 2. Illustration of the subtleties of barriers in nested loops

following:

$$\begin{aligned}
 A &= \{\{0, 1, 2, 3\}, \{-, -, -, -\}\} \\
 &\rightarrow \{\{0, 1, 2, 3\}, \{1, 2, 3, 0\}\} \rightarrow \{\{2, 3, 0, 1\}, \{1, 2, 3, 0\}\} \\
 &\rightarrow \{\{2, 3, 0, 1\}, \{3, 0, 1, 2\}\} \rightarrow \{\{0, 1, 2, 3\}, \{3, 0, 1, 2\}\},
 \end{aligned}$$

under the assumption that, during each inner loop iteration, threads 1, 2 and 3 wait at the barrier until thread 0 reaches the barrier again by proceeding to the next iteration of the outer loop.

However, consider the predicated version of the kernel shown in part on the right of Figure 2. This is the form in which the kernel executes on an Nvidia GPU. The four threads comprise a single sub-group. All threads will enter the outer loop and execute the first inner loop iteration. Then thread 0 will become disabled ( $q$  becomes *false*) for the inner loop. Thus the barrier will be executed with some, but not all, threads in the sub-group enabled. On Nvidia hardware, a barrier is compiled to a `bar.sync` instruction in the PTX (Parallel Thread Execution) assembly language. According to the PTX documentation [Nvidia 2012b],

*“if any thread in a [sub-group] executes a `bar` instruction, it is as if all the threads in the [sub-group] have executed the `bar` instruction.”*

Thus threads 1, 2 and 3 will *not* wait at the barrier until thread 0 returns to the inner loop: they will simply continue to execute past the barrier, performing three more inner loop iterations. In effect, thread 0 is disabled while threads 1, 2 and 3 perform these additional loop iterations, because in practice the threads execute in lock step. This yields the following sequence of state-changes to  $A$ :

$$\begin{aligned}
 A &= \{\{0, 1, 2, 3\}, \{-, -, -, -\}\} \\
 &\rightarrow \{\{0, 1, 2, 3\}, \{1, 2, 3, 0\}\} \rightarrow \{\{0, 3, 0, 1\}, \{1, 2, 3, 0\}\} \\
 &\rightarrow \{\{0, 3, 0, 1\}, \{1, 0, 1, 0\}\} \rightarrow \{\{0, 1, 0, 1\}, \{1, 0, 1, 0\}\}.
 \end{aligned}$$

Table IV. The litmus test of Figure 2 yields a range of results across varying platforms

Architecture	Final state of $A$
Nvidia Tesla C2050	$\{\{0, 1, 0, 1\}, \{1, 0, 1, 0\}\}$
AMD Tahiti	$\{\{0, 1, 2, 3\}, \{1, 2, 3, 0\}\}$
ARM Mali-T600	$\{\{0, 1, 2, 3\}, \{3, 0, 1, 2\}\}$
Intel Xeon X5650	$\{\{*, *, *, 1\}, \{3, 0, 1, 2\}\}$

After the inner loop exits, thread 0 becomes enabled, but all other threads become disabled, for a further three outer loop iterations, during each of which thread 0 executes a single inner loop iteration. The state of  $A$  thus remains  $\{\{0, 1, 0, 1\}, \{1, 0, 1, 0\}\}$ .

The OpenCL standard [Khronos OpenCL Working Group 2012] gives a better, though still informal definition, stating:

*“If barrier is inside a loop, all [threads] must execute the barrier for each iteration of the loop before any are allowed to continue execution beyond the barrier.”*

which at least can be interpreted as rejecting the example of Figure 2.

To investigate this issue in practice, we implemented the litmus test of Figure 2 in both CUDA and OpenCL and (with help from contacts in the GPU industry; see Acknowledgments) ran the test on GPU architectures from Nvidia, AMD and ARM, and on an Intel Xeon CPU (for which there is an OpenCL implementation). Our findings are reported in Table IV. Observe that the test result does not agree between any two vendors. The Nvidia results match our above prediction. The AMD result also appears to stem from predicated execution. ARM’s Mali architecture does *not* use predicated execution [Lokhmotov 2011], so perhaps unsurprisingly gives the ‘intuitive’ result we might expect. For Intel Xeon, we found that different threads reported different values for certain array elements in the final shared state, indicated by asterisks in Table IV, which we attribute to cache effects.

The example of Figure 2 is contrived in order to be small enough to explain concisely and examine exhaustively. It does, however, illustrate that barrier divergence is a subtle issue, and that non-obvious misuse of barriers can compromise correctness and lead to implementation-dependent results. Clearly a more rigorous notion of barrier divergence is required than the informal descriptions found in the CUDA and OpenCL documentation.

We give a precise, operational definition for barrier divergence in Section 3 which clears up this ambiguity. In essence, our definition states that if a barrier is encountered by a group of threads executing in lock-step under a predicate, the predicate must hold *uniformly* across the group, i.e., the predicate must be *true* for all threads, or *false* for all threads. This precise definition facilitates formal verification of barrier divergence-freedom and is incorporated in the GPUVerify tool we present in Section 5. GPUVerify can accurately analyze the kernel of Figure 2, flagging the barrier divergence issue.

### 3. OPERATIONAL SEMANTICS FOR GPU KERNELS

Our aim is to verify data race- and barrier divergence-freedom for GPU kernels. In order to do this, we require an operational semantics for GPU kernels that specifies precisely the conditions under which data races and barrier divergence occur.

For checking barrier divergence, the least conservative assumption we can safely make is that a thread group consists of a *single* sub-group, so that all threads in a group execute in lock-step. Thus, when analyzing barrier divergence-freedom for a kernel in which each thread group consists of  $N$  threads, we assume for the purposes of analysis that the sub-group size is also  $N$  (so that each thread group trivially contains

just one sub-group), regardless of the sub-group size associated with any particular architecture. We can then define barrier divergence to occur if and only if the lock-step thread group executes a barrier and the threads are not uniformly enabled: the current predicate of execution holds for some threads but not others. Clearly if we can prove barrier divergence-freedom for a kernel under this strong assumption of lock-step behavior, the kernel will be free from barrier divergence if thread groups are divided into sub-groups at a finer level of granularity.

For example, when analyzing barrier divergence for a kernel in which thread blocks have size 1024, we assume that the threads execute in lock step as a single sub-group size of 1024. Recent GPUs from AMD, Nvidia and ARM exhibit sub-group sizes of 64, 32 and 1, respectively [AMD 2013; Nvidia 2012a; Lokhmotov 2011]. Our worst-case assumption thus caters for the AMD scenario with the largest sub-group size, as well as the smaller sub-group sizes exhibited by Nvidia and ARM GPUs. Furthermore, the conservative assumption also caters for possible sub-group sizes associated with any other (present or future) architectures.

For race checking, the scenario is reversed: the least conservative safe assumption is that threads in the same group interleave completely asynchronously between pairs of barriers, with no guarantees as to the relative order of statement execution between threads (so that essentially every sub-group consists of just a single thread). This is the case, for example, on ARM's Mali GPU architecture [Lokhmotov 2011]. Clearly if race-freedom can be proved under this most general condition, then a kernel will remain race-free if, in practice, certain threads in a group execute synchronously.

We propose a semantics which we call *synchronous, delayed visibility (SDV)*. Under SDV, group execution is *synchronous*, allowing precise barrier divergence checking. The shared memory accesses of all threads are logged, and when a thread writes to shared memory, the visibility of this write to other threads in the group is *delayed* until the group reaches a barrier. Delaying the visibility of writes ensures that threads do *not* see a synchronized view of shared and global memory between barriers, catering for the fact that execution might not *really* be fully synchronous. Logging accessed locations allows racing accesses to be detected. The SDV semantics does *not* attempt to directly capture the manner in which threads execute on real GPU architectures which exhibit a mixture of synchronous and interleaving execution in practice; instead the semantics provides a basis for proving data race- and barrier divergence-freedom.

To describe the SDV semantics formally, we define the *Kernel Programming Language (KPL)*, which captures the essential features of mainstream languages for writing GPU kernels such as CUDA and OpenCL.

### 3.1. Syntax

The syntax for KPL is shown in Figure 3. A KPL kernel declares the number of groups (groups: number) and the number of threads in each group (threads: number) that will execute the kernel. Thus, the total number of threads that will execute the kernel is equal to the number of declared groups multiplied by the number of declared threads. The group and thread declarations are followed by a list of procedure declarations and a main statement. Each procedure has a name, a single parameter and a body; for brevity we do not model multiple parameters or return values.

For simplicity, but without loss of generality, threads have access to a single shared array *sh* which we refer to as *global memory*. We make no distinction between shared and global memory (cf. Figure 1) in our programming language; from a verification perspective, it suffices to view shared memory as being part of global memory (by employing the ids of groups in the calculation of memory addresses). We assume that every local variable and each indexable element of global memory has type *Word*, the type of memory words. Any value of type *Word* can be interpreted as an integer or a

```

kernel      ::=  groups: number
                threads: number
                proc*
                main: stmt
proc        ::=  procedure name var stmt
stmt        ::=  basic_stmt | stmt; stmt
                | local name stmt
                | if local_expr stmt else stmt
                | while local_expr stmt
                | while local_expr stmt
                | call name(local_expr)
                | barrier
                | break | continue | return
basic_stmt  ::=  loc := local_expr
                | loc := sh[local_expr]
                | sh[local_expr] := local_expr
local_expr  ::=  lid | gid | loc
                | constant literal of type Word
                | local_expr op local_expr
loc         ::=  name | V
name        ::=  any valid C name
number     ::=  any positive number

```

Fig. 3. Syntax for Kernel Programming Language

Boolean. In practice, Word will also represent floating point numbers; structured data will be represented by sequences of values of type Word.

Local expressions are built from constant literals, locations, and the built-in variables *lid* and *gid* which refer to the local id of a thread within its group and the id of the group in which the thread occurs, respectively. Compound expressions are built using an unspecified set of binary operators; the syntax could easily be extended to accommodate operators of other arities.

A thread may update one of its local variables by performing a local computation ( $v := e$ , where  $e$  is a local expression), or by reading from global memory ( $v := sh[e]$ , where  $e$  is a local expression determining which index to read from). A thread may also update global memory ( $sh[e_1] := e_2$ , where  $e_1, e_2$  are expressions over local variables, with  $e_1$  determining which index to write to, and  $e_2$  the value to be written). For simplicity, we assume that all local variables are scalar.

Compound statements are constructed via sequencing, conditional branches, local variable introduction, loops, and procedure calls in the standard way. KPL provides a number of other statements: **barrier**, which is used to synchronize threads in a group; **break**, which causes execution to break out of the closest enclosing loop; **continue**, which causes execution to jump to the head of the closest enclosing loop; and **return**, which causes execution to return from the closest enclosing procedure call.

Figure 3 specifies two syntactic elements which should *not* appear directly in a KPL program; they are used in the semantic rules of Figure 6 which we will explain in Section 3.2. These are: a special **while** statement, used to model the dynamic semantics of while loops in which we have to distinguish between the first and subsequent iterations of a loop, and a set  $V$  of locations from which storage for local variables is allocated as they come into scope dynamically.

We assume that programs are well-formed according to the usual rules, e.g., statements should only refer to declared variables and variable introduction should not

hide a variable introduced earlier in an enclosing scope. In addition, we require the main statement not to contain a return; this non-consequential restriction simplifies our presentation somewhat (it avoids having to apply the function *relim* of Figure 6, discussed in Section 3.2 below, to the main statement).

We do not formalize features of GPU kernels such as multi-dimensional groups and arrays. However, our verification method and implementation, described in Section 5, handles both.

### 3.2. Semantics

Given a function  $f : A \rightarrow B$  and elements  $a \in A, b \in B$ , we write  $f[a \mapsto b]$  to denote the function  $g : A \rightarrow B$  such that  $g(x) = f(x)$  for all  $x \in A \setminus \{a\}$ , and  $g(a) = b$  otherwise. We abbreviate  $f[a \mapsto b][c \mapsto d]$  with  $a \neq c$  to  $f[a \mapsto b, c \mapsto d]$ . Viewing a tuple with named components as a function mapping names to element values, we also use this notation to specify updates to tuples. We use  $\langle s_1, s_2, \dots, s_k \rangle$  to denote a sequence of length  $k$ , and write  $\langle \rangle$  for the empty sequence. We write  $s : ss$  for a sequence whose first element is  $s$ , and whose remaining elements form the sequence  $ss$ ; the operator ‘:’ associates to the right.

In our semantics, each thread is equipped with a *shadow* copy of global memory. At the start of kernel execution, the shadow memory of all threads is identical. During execution, a thread reads and modifies its shadow memory locally, and maintains a read set and a write set recording those addresses in global memory that the thread has accessed. During each access the read and write sets are checked for data races. If a race has occurred, execution aborts. When a barrier statement is reached with all threads in a group enabled, the write sets are used to build a consistent view of global memory between the threads in a group, the shadow memories are all updated to agree with this view, and the read and write sets are cleared with regard to the accesses of threads within groups.

In what follows let  $P$  be a KPL kernel.

**3.2.1. Thread, Group and Kernel States.** Let  $gs$  denote the number of groups executing  $P$ , specified via groups:  $gs$  in the definition of  $P$ . Similarly, let  $ts$  denote the number of threads in each group, specified via threads:  $ts$ . Let  $\mathcal{G}$  denote the set  $\{0, 1, \dots, gs - 1\}$ , and for each  $i \in \mathcal{G}$  let  $\mathcal{T}_i$  denote the set  $\{0, 1, \dots, ts - 1\}$ .<sup>4</sup> Together,  $\mathcal{G}$  and  $\mathcal{T}$  specify the identifiers of all the threads running the kernel, via the function

$$tids(\mathcal{G}, \mathcal{T}) \triangleq \{(i, j) \mid i \in \mathcal{G} \wedge j \in \mathcal{T}_i\}.$$

We now define thread, group and kernel states.

**Thread States.** A *thread state* is a tuple  $(l, sh, R, W)$  where:

- $l : V \rightarrow \text{Word}$  is the local variable store of the thread (recall that  $V$  is a set of locations); included are locations *gid* and *lid* specifying, respectively, the id of the group to which the thread belongs, and the id of the thread within this group.
- $sh : \mathbb{N} \rightarrow \text{Word}$  is the shadow copy of global memory owned by the thread.
- $R, W \subseteq \mathbb{N}$  are the read and write sets of the thread recording the global addresses the thread has accessed since the last barrier; these sets are used to detect intra-group data races.

We use  $\tau$  to denote a thread state, and  $\tau.l, \tau.sh, \tau.R, \tau.W$  to refer to the components of  $\tau$ . The set of all thread states is denoted by `ThreadStates`. Given a local expression

<sup>4</sup>Observe that it is not strictly necessary to parameterize  $\mathcal{T}$  in  $i$ , as  $\mathcal{T}_i = \{0, 1, \dots, ts - 1\}$  for all  $i \in \mathcal{G}$ . However, the parameterization will simplify the presentation of the kernel abstractions presented below.

$e$  and a thread state  $\tau$ , we write  $e^\tau$  for the result of evaluating  $e$  according to  $\tau.l$ . We do not provide a concrete definition of  $e^\tau$ , which depends on the nature of the base type *Word* and the available operators, except we specify that, for a storage location  $v$ ,  $v^\tau = \tau.l(v)$ .

**Group States.** A *group state* of a group  $i \in \mathcal{G}$  is a tuple  $(\gamma_{ts}, R, W)$  where:

- $\gamma_{ts} : \mathcal{T}_i \rightarrow \text{ThreadStates}$  records a thread state for each thread in group  $i$ .
- $R, W \subseteq \mathbb{N} \times \mathcal{T}_i$  are the read and write sets of group  $i$ . These are used to detect *inter-group* data races. Each element  $(x, j)$  records a global address  $x$  accessed by thread  $j$  in the group since the start of the kernel. It is necessary to store the local id  $j$  so that when we later reduce executions to just two threads, accesses by threads other than the two selected ones can be identified and removed from  $R$  and  $W$ .

We use  $\gamma$  (without a subscript) to denote a group state. Given  $\gamma = (\gamma_{ts}, R, W)$ , we use  $\gamma(j)$  to refer to  $\gamma_{ts}(j)$  and we refer to  $R$  and  $W$  by  $\gamma.R$  and  $\gamma.W$ , respectively. A group state  $\gamma$  of a group  $i$  is *valid* if  $\gamma(j).l(\text{gid}) = i$  and  $\gamma(j).l(\text{lid}) = j$  for all  $j \in \mathcal{T}_i$ , i.e., knowledge of the group id  $i$  is consistent across the threads in the group and each thread has a unique local id corresponding to its location in the group.

The set of all valid group states for a group  $i$  is denoted  $\text{GroupStates}_i$ , and we use  $\text{GroupStates}$  to denote  $\bigcup_{i \in \mathcal{G}} \text{GroupStates}_i$ , the set of all valid group states.

**Kernel States.** A *predicated statement* is a pair  $(s, p)$ , with  $s \in \text{stmt}$  and  $p \in \text{local\_expr}$ . Intuitively,  $(s, p)$  denotes a statement  $s$  that should be executed if  $p$  holds, and otherwise should have no effect. The set of all predicated statements is denoted  $\text{PredStmts}$ .

A *kernel state* is a tuple  $(\kappa, ss)$  where:

- $\kappa : \mathcal{G} \rightarrow \text{GroupStates}$  records a group state for each group in the kernel.
- $ss \in \text{PredStmts}^*$  is a sequence of predicated statements to be executed by the kernel.

A kernel state  $(\kappa, ss)$  is *valid* if, for each  $i \in \mathcal{G}$ ,  $\kappa(i) \in \text{GroupStates}_i$ , i.e., each group has a unique id. The set of all valid kernel states is denoted  $\text{KernelStates}$ .

Observe that a kernel state does not include a single, definitive global memory component: global memory is represented via the shadow copies held by the individual threads in groups, which are initially consistent, and are made consistent again at barriers for threads within a group.

A kernel state  $(\kappa, ss) \in \text{KernelStates}$  is a valid *initial state* of  $P$  if:

- $ss = \langle (s, \text{true}) \rangle$ , where  $s$  is declared in  $P$  via  $\text{main} : s$ .
- $\kappa(i).R = \kappa(i).W = \emptyset$  for all  $i \in \mathcal{G}$ .
- $\kappa(i)(j).R = \kappa(i)(j).W = \emptyset$  for all  $i \in \mathcal{G}$  and  $j \in \mathcal{T}_i$ .
- $\kappa(i)(j).sh = \kappa(i')(j').sh$  for all  $i, i' \in \mathcal{G}$  and  $j \in \mathcal{T}_i, j' \in \mathcal{T}_{i'}$ .
- $\kappa(i)(j).l(v) = \text{false}$  for all  $i \in \mathcal{G}, j \in \mathcal{T}_i$  and  $v \in V$ .

The first requirement ensures that the kernel executes the main statement. The second and third requirement guarantee that the read and write sets of both groups and threads are initially empty. The fourth requirement ensures that threads have a consistent but arbitrary initial view of global memory. We single out the final requirement as it impacts on our semantics in a somewhat subtle manner:

**Remark 3.1 (Local Storage Locations are Initially false).** The fifth requirement ensures that local storage locations are initialized to the value of type *Word* corresponding to *false*. Thus whenever fresh variables are introduced by the semantic rules of Figure 6 these variables initially store the value *false*. This requirement simplifies our presentation of *break* and *continue* in the rules of Figure 6.

$$\begin{array}{c}
\frac{\neg p^\tau}{(\tau, \text{basic\_stmt}, p) \rightarrow_t \tau} \text{ (T-DISABLED)} \qquad \frac{p^\tau \quad l' = \tau.l[v \mapsto e^\tau]}{(\tau, v := e, p) \rightarrow_t \tau[l \mapsto l']} \text{ (T-ASSIGN)} \\
\\
\frac{p^\tau \quad l' = \tau.l[v \mapsto \tau.sh(e^\tau)] \quad R' = \tau.R \cup \{e^\tau\}}{(\tau, v := sh[e], p) \rightarrow_t \tau[l \mapsto l', R \mapsto R']} \text{ (T-READ)} \\
\\
\frac{p^\tau \quad sh' = \tau.sh[e_1^\tau \mapsto e_2^\tau] \quad W' = \tau.W \cup \{e_1^\tau\}}{(\tau, sh[e_1] := e_2, p) \rightarrow_t \tau[sh \mapsto sh', W \mapsto W']} \text{ (T-WRITE)}
\end{array}$$

Fig. 4. Rules for predicated execution of basic statements by threads

$$\begin{array}{c}
\frac{\forall j \in \mathcal{T}_i . (\gamma(j), \text{basic\_stmt}, p) \rightarrow_t \gamma'(j) \quad \text{group\_race}_i(\gamma')}{(\gamma, \text{basic\_stmt}, p) \rightarrow_{g(i)} \text{error}} \text{ (G-RACE)} \\
\\
\frac{\forall j \in \mathcal{T}_i . (\gamma(j), \text{basic\_stmt}, p) \rightarrow_t \gamma'(j) \quad \neg \text{group\_race}_i(\gamma') \quad \gamma'.R = \gamma.R \cup \bigcup_{j \in \mathcal{T}_i} (\gamma'(j).R \times \{j\}) \quad \gamma'.W = \gamma.W \cup \bigcup_{j \in \mathcal{T}_i} (\gamma'(j).W \times \{j\})}{(\gamma, \text{basic\_stmt}, p) \rightarrow_{g(i)} \gamma'} \text{ (G-BASIC)}
\end{array}$$

(a) Rules for basic statements

$$\begin{array}{c}
\frac{\forall j \in \mathcal{T}_i . \neg p^{\gamma(j)}}{(\gamma, \text{barrier}, p) \rightarrow_{g(i)} \gamma} \text{ (G-NO-OP)} \qquad \frac{\exists j \neq k \in \mathcal{T}_i . p^{\gamma(j)} \wedge \neg p^{\gamma(k)}}{(\gamma, \text{barrier}, p) \rightarrow_{g(i)} \text{error}} \text{ (G-DIVERGENCE)} \\
\\
\frac{\forall j \in \mathcal{T}_i . p^{\gamma(j)} \quad \forall j \in \mathcal{T}_i . \gamma'(j) = \gamma(j)[sh \mapsto \text{merge}_i(\gamma), R \mapsto \emptyset, W \mapsto \emptyset]}{(\gamma, \text{barrier}, p) \rightarrow_{g(i)} \gamma'} \text{ (G-SYNC)}
\end{array}$$

(b) Rules for barriers

Fig. 5. Rules for lock-step execution of basic statements and barriers by groups

**3.2.2. Predicated Execution.** We now describe predicated execution of threads, groups, and kernels.

*Predicated Execution of Threads.* The rules of Figure 4 define the binary relation

$$\rightarrow_t \subseteq (\text{ThreadStates} \times \text{PredStmts}) \times \text{ThreadStates}$$

describing the evolution of one thread state to another under execution of a predicated basic statement. For readability, given a thread state  $\tau$  and predicated statement  $(s, p)$ , we write  $(\tau, s, p)$  instead of  $(\tau, (s, p))$ .

Rule T-DISABLED ensures that a predicated statement has no effect if the predicate does not hold, indicated by  $\neg p^\tau$  in the premise of the rule; T-ASSIGN updates  $\tau.l$  according to the assignment; T-READ updates the local store of the thread with the appropriate element from the shadow copy of shared memory owned by the thread, and records the address that was read from; rule T-WRITE is analogous.

*Predicated Execution of Groups.* Figure 5 defines the binary relation

$$\rightarrow_{g(i)} \subseteq (\text{GroupStates}_i \times \text{PredStmts}) \times (\text{GroupStates}_i \cup \{\text{error}\})$$

describing how a group state of a group  $i$  mutates into another as a result of executing a predicated basic statement or barrier, where *error* is a designated error state. As for

threads, given a group state  $\gamma$  and predicated statement  $(s, p)$ , we write  $(\gamma, s, p)$  instead of  $(\gamma, (s, p))$ .

Intra-group data races are detected via rule G-RACE. An intra-group race occurs if there is a thread  $j$  in group  $i$  such that the write set of  $j$  intersects with either the read or write set of a thread  $k$  different from  $j$  but in the same group. The predicate  $\text{group\_race}_i(\gamma)$  used to detect intra-group races is formally defined as:

$$\text{group\_race}_i(\gamma) \triangleq \exists j \neq k \in \mathcal{T}_i . \gamma(j).W \cap (\gamma(k).R \cup \gamma(k).W) \neq \emptyset .$$

Rule G-RACE states that if the execution of a basic statement by all threads in a group leads to a group state  $\gamma'$  in which  $\text{group\_race}_i(\gamma')$  holds, then the *error* state is reached.

Lock-step execution of basic statements by a group is achieved by having each thread in the group execute the statement; the order in which they do so is irrelevant due to delayed visibility (G-BASIC). The rule requires the execution of the basic statement to be free from intra-group data races. The rule also records which addresses of global memory were read from or written to by the threads in the group; this is achieved by propagating the addresses from the read and write sets of the individual threads into the read and write set of the group.

Execution of barrier with all threads in a group disabled has no effect (G-NO-OP). If the group is due to execute a barrier statement under predicate  $p$  but not all threads agree on the truth of  $p$ , then the *error* state is reached (G-DIVERGENCE). This *precisely* captures the notion of barrier divergence discussed in Section 2.

Rule G-SYNC captures the effect of barrier synchronization under predicate  $p$  in the case where all threads in the group agree that  $p$  is valid. A new group state  $\gamma'$  is constructed, in which for each thread  $j$ , the read and write sets of  $j$  are empty in  $\gamma'(j)$  and the local component  $\gamma'(j).l$  is identical to the local component for the thread before the barrier. The barrier also enforces a consistent view of shared memory across the group by setting the shadow memories of all threads to the same value. This is achieved by the function  $\text{merge}_i$ . If thread  $j$  has recorded a write to a shared memory location  $z$ , i.e.,  $z \in \gamma(j).W$ , then  $\text{merge}_i(\gamma)$  maps  $z$  to the value at address  $z$  in the shadow memory of thread  $j$ , i.e., to  $\gamma(j).sh(z)$ .

Formally,  $\text{merge}_i(\gamma)$  is a map satisfying the following constraints:

$$\frac{j \in \mathcal{T}_i \quad z \in \gamma(j).W}{\text{merge}_i(\gamma)(z) = \gamma(j).sh(z)} \quad \frac{\forall j \in \mathcal{T}_i . z \notin \gamma.W}{\text{merge}_i(\gamma)(z) = \gamma(0).sh(z)}$$

The value of  $\text{merge}_i(\gamma)$  is guaranteed to be unique, otherwise there exists a  $z$  belonging to  $\gamma(j).W$  and  $\gamma(k).W$  for  $j \neq k$ , in which case execution would have aborted earlier via G-RACE. Observe that, in the second constraint,  $\gamma(0)$  is arbitrary, and could have been chosen to be any  $\gamma(i)$ ; the value of  $sh(z)$  is consistent across threads in this case, since no thread has written to the memory location  $z$ .

*Predicated Execution of Kernels.* Figure 6 defines the binary relation

$$\rightarrow_k \subseteq \text{KernelStates} \times (\text{KernelStates} \cup \{\text{error}\}),$$

where *error* is again a designated error state. This relation describes the evolution of a kernel as it executes a sequence of predicated statements.

Inter-group races are detected via rule K-INTER-GROUP-RACE. An inter-group race occurs if there is a group  $i$  such that the write set of  $i$  intersects with either the read or write set of a group  $j$  different from  $i$  (ignoring the particular thread responsible for each read and write). The predicate  $\text{kernel\_race}(\kappa)$  used to detect inter-group races is formally defined as:

$$\text{kernel\_race}(\kappa) \triangleq \exists i \neq j \in \mathcal{G} . \text{fst}(\kappa(i).W) \cap \text{fst}(\kappa(j).R \cup \kappa(j).W) \neq \emptyset$$

$$\begin{array}{c}
\frac{\forall i \in \mathcal{G} . (\kappa(i), \text{basic\_stmt}, p) \rightarrow_{g(i)} \kappa'(i) \quad \text{kernel\_race}(\kappa')}{(\kappa, (\text{basic\_stmt}, p) : ss) \rightarrow_k \text{error}} \text{ (K-INTER-GROUP-RACE)} \\
\\
\frac{\exists i \in \mathcal{G} . (\kappa(i), \text{basic\_stmt}, p) \rightarrow_{g(i)} \text{error}}{(\kappa, (\text{basic\_stmt}, p) : ss) \rightarrow_k \text{error}} \text{ (K-INTRA-GROUP-RACE)} \\
\\
\frac{\forall i \in \mathcal{G} . (\kappa(i), \text{basic\_stmt}, p) \rightarrow_{g(i)} \kappa'(i) \quad \neg \text{kernel\_race}(\kappa')}{(\kappa, (\text{basic\_stmt}, p) : ss) \rightarrow_k (\kappa', ss)} \text{ (K-BASIC)} \\
\\
\text{(a) Rules for basic statements} \\
\\
\frac{\exists i \in \mathcal{G} . (\kappa(i), \text{barrier}, p) \rightarrow_{g(i)} \text{error}}{(\kappa, (\text{barrier}, p) : ss) \rightarrow_k \text{error}} \text{ (K-DIVERGENCE)} \\
\\
\frac{\forall i \in \mathcal{G} . (\kappa(i), \text{barrier}, p) \rightarrow_{g(i)} \kappa'(i)}{(\kappa, (\text{barrier}, p) : ss) \rightarrow_k (\kappa', ss)} \text{ (K-SYNC)} \\
\\
\text{(b) Rules for barriers} \\
\\
\frac{}{(\kappa, (S_1; S_2, p) : ss) \rightarrow_k (\kappa, (S_1, p) : (S_2, p) : ss)} \text{ (K-SEQ)} \\
\\
\frac{\text{fresh } v \in V}{(\kappa, (\text{local } x S, p) : ss) \rightarrow_k (\kappa, (S[x \mapsto v], p) : ss)} \text{ (K-VAR)} \\
\\
\frac{\text{fresh } v \in V}{(\kappa, (\text{if } e S_1 \text{ else } S_2, p) : ss) \rightarrow_k (\kappa, (v := e, p) : (S_1, p \wedge v) : (S_2, p \wedge \neg v) : ss)} \text{ (K-IF)} \\
\\
\frac{\text{fresh } v \in V}{(\kappa, (\text{while } e S, p) : ss) \rightarrow_k (\kappa, (\underline{\text{while}} e \text{ belim}(S, v), p \wedge \neg v) : ss)} \text{ (K-OPEN)} \\
\\
\frac{\exists i \in \mathcal{G}, j \in \mathcal{T}_i . (p \wedge e)^{\kappa(i)(j)} \quad \text{fresh } u, v \in V}{(\kappa, (\underline{\text{while}} e S, p) : ss) \rightarrow_k (\kappa, (u := e, p) : (\text{celim}(S, v), p \wedge u \wedge \neg v) : (\underline{\text{while}} e S, p) : ss)} \text{ (K-ITER)} \\
\\
\frac{\forall i \in \mathcal{G}, j \in \mathcal{T}_i . \neg(p \wedge e)^{\kappa(i)(j)}}{(\kappa, (\underline{\text{while}} e S, p) : ss) \rightarrow_k (\kappa, ss)} \text{ (K-DONE)} \\
\\
\frac{\text{fresh } u, v \in V \quad S = \text{Body}(f)[\text{Param}(f) \mapsto u]}{(\kappa, (\text{call } f(e), p) : ss) \rightarrow_k (\kappa, (u := e; \text{relim}(S, v), p \wedge \neg v) : ss)} \text{ (K-CALL)} \\
\\
\text{(c) Rules for statements}
\end{array}$$

Fig. 6. Rules for lock-step execution of statements by kernels

where  $\text{fst}(X) = \{x \mid (x, y) \in X\}$ . Rule **K-INTER-GROUP-RACE** states that if collective execution of a predicated basic statement by all groups leads to a state  $\kappa'$  in which  $\text{kernel\_race}(\kappa')$  holds, then the *error* state is reached. Rule **K-INTRA-GROUP-RACE** detects intra-group races by lifting the application of rule **G-RACE** to the kernel level.

Collective execution of a predicated basic statement is achieved by having every group execute the statement; the order in which they do so is irrelevant (K-BASIC). The rule requires that the collective execution was free from inter-group data races (via the  $\neg\text{kernel\_race}(\kappa')$  condition) and intra-group data races (via the condition that each group state in  $\kappa$  transitions to a non-error group state in  $\kappa'$ ).

Rule K-DIVERGENCE lifts application of rule G-DIVERGENCE to the kernel level, i.e., the kernel transitions to the *error* state if barrier divergence is detected in one of the groups. Collective execution of a barrier statement in the absence of barrier divergence is achieved via rule K-SYNC where each group transitions to a new group state either by rule G-SYNC or rule G-NO-OP. Observe that the read and write sets of groups are not cleared, since barriers are used solely for intra-group synchronization and do not achieve inter-group synchronization.

The remaining rules of Figure 6 describe predicated execution of compound statements. Rule K-SEQ is straightforward. Rule K-VAR creates storage for a new local variable  $x$  by allocating a fresh location  $v$  in  $V$  and substituting all occurrences of  $x$  in  $S$  by  $v$ ; we use  $S[x \mapsto v]$  to denote this substitution. Rule K-IF decomposes a conditional statement into a sequence of predicated statements: the guard of the conditional is evaluated into a new location  $v$ ; the *then* branch  $S_1$ , is executed by all threads in all groups under predicate  $p \wedge v$  (where  $p$  is the predicate already in place on entry to the conditional), and the *else* branch  $S_2$ , is executed by all threads in all groups under predicate  $p \wedge \neg v$ .

Rules K-OPEN, K-ITER and K-DONE together model predicated execution of a while loop. In what follows, we say that a break or continue statement is *top-level* in a loop if the statement appears in the loop body but is not nested inside any further loops.

Rule K-OPEN converts a while loop into a **while** loop by creating fresh storage to model break statements. A fresh location  $v$  is selected;  $v$  records whether a thread has executed a break statement associated with the while loop. Like all local storage,  $v$  initially has the value *false* (see Remark 3.1): no thread has executed break on loop entry. The function *belim* is applied to the loop body. This function takes a statement  $S$  and a location  $v$  and replaces each top-level break statement inside  $S$  by the statement  $v := \text{true}$ . The predicate for the execution of the **while** loop becomes  $p \wedge \neg v$  to model that the statements in the loop have no effect after the execution of a break statement. A similar technique to model break statements is used by Habermaier and Knapp [2012].

The K-ITER rule models execution of loop iterations, and handles continue statements. The rule fires each time a loop iteration is executed. For a given loop iteration, two fresh local storage locations  $u$  and  $v$  are selected and (as per Remark 3.1) both are initialized to *false*.<sup>5</sup> Location  $u$  is used to store the valuation of the loop guard. Location  $v$  is used to record whether a thread has executed a top-level continue statement during the current loop iteration. The value of  $v$  is initially *false*, since no thread has executed a continue statement at the beginning of an iteration. The statement  $u := e$  (executed under the enclosing predicate  $p$ ) evaluates the loop guard into  $u$ . The function *celim* is applied to the loop body; this function takes a statement  $S$  and a location  $v$  and replaces each top-level continue statement inside  $S$  by the statement  $v := \text{true}$ . The loop body, after elimination of continue statements, is executed under the predicate  $p \wedge u \wedge \neg v$ : a thread is enabled during the current iteration if the incoming predicate holds ( $p$ ), the loop guard evaluates to *true* at the start of the iteration ( $u$ ) and the thread has not executed a continue statement ( $\neg v$ ). Note that, due to rule K-OPEN, the incoming predicate  $p$  includes a conjunct recording whether the thread

<sup>5</sup>Note that  $u$  and  $v$  are not re-used between loop iterations: a fresh  $u$  and  $v$  are selected at the start of each loop iteration. Hence, there is no need to reset  $u$  and  $v$  at the end of a loop iteration.

has executed a `break` statement. After the loop body, the `while` construct is considered again.

Thus, all threads continuously execute the loop body using K-ITER until, for *every* thread in *every* group, (a) the enclosing predicate  $p$  becomes *false*, either because this predicate was *false* on loop entry or because the thread has executed `break`, or (b) the loop condition no longer holds for the thread.<sup>6</sup> When (a) or (b) holds for all threads, loop exit is handled by rule K-DONE.

The rule K-CALL models the execution of a call to a procedure  $f$ . This involves executing the statement corresponding to the body of the called procedure ( $Body(f)$ ) after replacing all occurrences of its formal parameter ( $Param(f)$ ) with a location storing the evaluation result of the actual parameter expression. All threads execute the entire body of a procedure in lock-step. A fresh storage location  $v$  is used to record whether a thread has executed a return statement. Initially this location is set to *false*, and the function *relim* replaces each return statement in  $Body(f)$  with the statement  $v := true$ . The procedure body is executed under the predicate  $p \wedge \neg v$  (where  $p$  is the existing predicate of execution at the point of the call) so that execution of a return statement by a thread is simulated by the thread becoming disabled for the remainder of the procedure body.

### 3.3. Formalization in Isabelle

The definitions in this section have been formalized in the Isabelle proof assistant [Nipkow et al. 2002], and are available from the Archive of Formal Proofs [Wickerson 2014]. For the most part, the Isabelle definitions closely resemble those presented above, but there are three notable differences. First, we do not model the abuse of notation whereby  $\gamma(j)$  abbreviates  $\gamma_{ts}(j)$ , since the use of such *type coercions* was found to necessitate many additional type annotations. Second, we extend kernel states with a third component, to record all the variables that have been used. This is necessary for calculating the fresh variables required by several of the rules in Figure 6, but is left implicit, in accordance with common practice, in the previous subsections. Third, because Isabelle does not allow *dependent types*, we cannot define functions between arbitrary sets (at least, not without extensive use of `typedefs`, which can complicate proofs) such as  $\gamma_{ts} : \mathcal{T}_i \rightarrow \text{ThreadStates}$ . In such situations, we instead declare  $\gamma_{ts}$  as a partial function from local ids to ThreadStates, and make  $dom(\gamma_{ts}) = \mathcal{T}_i$  available as an assumption when proving theorems.

In Section 4.2 we discuss our experience using these formalized definitions to partially mechanize a proof of soundness for the *two-thread reduction*; this reduction is our next main contribution.

## 4. THE TWO-THREAD REDUCTION

The SDV operational semantics of Section 3 has all threads execute in lock-step. As such, our massively-parallel  $N$ -threaded kernel can be encoded as a single sequential program, where each instruction is replicated  $N$  times, and where assertions are inserted to check for data races and barrier divergence. In principle, existing techniques for sequential program verification can be applied directly to this program and, as discussed in detail in Section 5, we have investigated the use of the Boogie verification system [Barnett et al. 2005] for this purpose. We have found that directly encoding massively-parallel kernels as sequential programs leads to infeasibly large

<sup>6</sup>Execution of `continue` does *not* directly contribute to the conditions under which threads exit a loop. This is because on executing a `continue` statement a thread becomes disabled for the remainder of the current loop iteration, but the thread is *re-enabled* at the start of the next iteration, provided the loop guard still holds for the thread.

verification conditions if an explicit representation of threads is used: in Section 6.4 we illustrate this blow-up in verification condition size, and the corresponding effect this has on verification performance, using a simple example. The explicit representation can be replaced with quantification over all threads, but we have found this to lead to verification conditions that cannot be automatically handled by Z3, a state-of-the-art theorem prover [de Moura and Bjørner 2008] which Boogie uses by default.<sup>7</sup>

To overcome these practical limitations and achieve scalable reasoning, we describe an alternative encoding that involves reducing the number of threads under consideration from  $N$  to just *two*.

Observe that the properties of data race- and barrier divergence-freedom are *pairwise* properties: a race occurs when accesses by two threads conflict, and barrier divergence occurs when a barrier is executed in a state where one thread is enabled and another thread in the same group is disabled. We can therefore consider an operational semantics where the predicated execution of only *two* threads is modeled. If we can use such a semantics to prove a kernel data race- and barrier divergence-free for a pair of distinct but otherwise *arbitrary* threads, we can conclude correctness of the kernel. This is because the verification process considers all possible execution traces so that, because the two threads under consideration are arbitrary, interactions between all possible pairs of threads are implicitly considered.

Formally, if the original execution involves all of the threads in  $tids(\mathcal{G}, \mathcal{T})$  (see Section 3.2.1), then we need to check data race- and barrier divergence-freedom for all restrictions of  $(\mathcal{G}, \mathcal{T})$  to a pair of threads, i.e., for restrictions to  $(\mathcal{G}', \mathcal{T}')$  such that

- $\mathcal{G}' \subseteq \mathcal{G}$  and  $\mathcal{T}'_i \subseteq \mathcal{T}_i$  for all  $i \in \mathcal{G}'$ , and
- $tids(\mathcal{G}', \mathcal{T}')$  consists of exactly two elements.

It is important to clarify that in practice we do not enumerating all possible two-thread reductions. Rather, as discussed in Section 5 when we discuss the implementation of GPUVerify, we use symbolic constants to model the identifiers of two arbitrary threads during analysis.

We now present the reduction to a pair of threads formally and prove that it is sound. Our GPUVerify verification technique and tool, described in Section 5, depend upon this reduction. The design of the PUG verifier for CUDA kernels by Li and Gopalakrishnan [2010] also hinges on a reduction to a pair of threads. However, this is the first formal presentation and proof of soundness for the reduction method.

Observe that the definitions of (valid) group and kernel states carry over to the sets  $\mathcal{G}'$  and  $\mathcal{T}'_i$  as defined above and that the same holds for the kernel level execution rules from Figure 6. However, for the above approach to be sound, we must *approximate* shared state handling, abstracting the values written to the shared state by threads that are not modeled. This can be achieved in multiple ways. We consider the following strategies:

*Equality abstraction:* Threads manipulate a shadow copy of the shared state. At a barrier, the shadow copies are set to be *arbitrary*, but *equal*. Thus on leaving the barrier, the threads have a consistent view of shared memory.

*Adversarial abstraction:* The shared state is made irrelevant; reads from the shared state into local variables are replaced with non-deterministic assignments to local variables.

We have found several example kernels where race-freedom hinges on threads agreeing on the values read from certain shared locations. In these cases, the adversarial

<sup>7</sup>The difficulty of reasoning about quantifiers is, of course, not specific to Z3, but a general obstacle faced by any automated theorem prover.

$$\frac{\forall j \in \mathcal{T}'_i . p^{\gamma(j)} \quad \exists sh' . \forall j \in \mathcal{T}'_i . \gamma'(j) = \gamma(j)[sh \mapsto sh', R \mapsto \emptyset, W \mapsto \emptyset]}{(\gamma, \mathbf{barrier}, p) \rightarrow_{g(i)} \gamma'} \text{ (G-SYNC-EQ)}$$

(a) Updated rule for the equality abstraction

$$\frac{p^\tau \quad \exists v' . l' = \tau.l[v \mapsto v'] \quad R' = \tau.R \cup \{e^\tau\}}{(\tau, v := sh[e], p) \rightarrow_t \tau[l \mapsto l', R \mapsto R']} \text{ (T-READ-ADV)}$$

$$\frac{p^\tau \quad W' = \tau.W \cup \{e_1^\tau\}}{(\tau, sh[e_1] := e_2, p) \rightarrow_t \tau[sh \mapsto sh', W \mapsto W']} \text{ (T-WRITE-ADV)}$$

$$\frac{\forall j \in \mathcal{T}'_i . p^{\gamma(j)} \quad \forall j \in \mathcal{T}'_i . \exists sh' . \gamma'(j) = \gamma(j)[sh \mapsto sh', R \mapsto \emptyset, W \mapsto \emptyset]}{(\gamma, \mathbf{barrier}, p) \rightarrow_{g(i)} \gamma'} \text{ (G-SYNC-ADV)}$$

(b) Updated rules for the adversarial abstraction

Fig. 7. Rules for the adversarial and equality abstractions

abstraction is too coarse for successful verification. However, in many cases, it does not matter what *specific* value is stored in shared memory, only that all threads see the same value. The equality abstraction suffices for such cases. Our use of equality abstraction allows us to improve upon the precision of the work of Li and Gopalakrishnan [2010] which is limited to the adversarial abstraction.

The execution relation for equality abstraction, written  $\rightarrow_{k, \mathcal{E}(\mathcal{G}', \mathcal{T}')}$ , is obtained by replacing rule G-SYNC in Figure 5 with rule G-SYNC-EQ from Figure 7(a), and replacing all references to  $\mathcal{G}$  and  $\mathcal{T}$  in all other semantic rules in Figures 5 and 6 with  $\mathcal{G}'$  and  $\mathcal{T}'$ , respectively. In rule G-SYNC-EQ, instead of merging the shared states of the threads in group  $i$ , an arbitrary value is chosen for the shared state and the shadow copy of each thread in the group is assigned this value. Clearing of read and write sets and checking for barrier divergence is left unchanged.

The execution relation for adversarial abstraction, written  $\rightarrow_{k, \mathcal{A}(\mathcal{G}', \mathcal{T}')}$ , is obtained by replacing rules T-READ and T-WRITE in Figure 4, and G-SYNC in Figure 5, with the rules T-READ-ADV, T-WRITE-ADV and G-SYNC-ADV from Figure 7(b), respectively, as well as replacing all references to  $\mathcal{G}$  and  $\mathcal{T}$  in all other semantic rules in Figures 5 and 6 with  $\mathcal{G}'$  and  $\mathcal{T}'$ , respectively. Rule T-READ-ADV makes the shared state irrelevant by updating the local store with an arbitrary value  $v'$  rather than one obtained from the shared state, although it still logs the read in the read set of the thread. Rule T-WRITE-ADV sets the shared state to an arbitrary value  $sh'$ , but still logs the write in the write set of the thread.<sup>8</sup> Similar to G-SYNC-EQ, rule G-SYNC-ADV does away with merging the shared states of the threads in a group  $i$ ; the rule sets the shared state of each thread to an arbitrary value, clears the read and write sets and checks for barrier divergence.

Comparing rules G-SYNC-EQ and G-SYNC-ADV, the premises are almost identical. The key difference is that in rule G-SYNC-EQ a *single* value is selected for the shared state, common to all threads, while in rule G-SYNC-ADV a value for the shared state

<sup>8</sup>Technically, this rule does not need to change the shared state at all, since the shared state is ignored by the T-READ-ADV. However, the rule in its current form is better able to simulate the concrete T-WRITE rule in the upcoming soundness proof.

is selected separately for each thread. Thus with the adversarial abstraction, threads do *not* see a unified view of shared memory after a barrier synchronization.

#### 4.1. Soundness of the Two-Thread Reduction

We shall first prove the soundness of the two-thread reduction when using the equality abstraction.

Throughout, let a *concrete execution* be a non-empty sequence of kernel states

$$\langle (\kappa_0, ss_0), \dots, (\kappa_{n-1}, ss_{n-1}) \rangle$$

where  $(\kappa_0, ss_0)$  is a valid initial kernel state and where each successive pair is related by the  $\rightarrow_k$  relation of Figure 6. An *execution under the equality abstraction with respect to  $(\mathcal{G}', \mathcal{T}')$*  is defined similarly, except that it uses the  $\rightarrow_{k, \mathcal{E}(\mathcal{G}', \mathcal{T}')}$  relation instead. Recalling that  $\mathcal{G} = \{0, 1, \dots, gs - 1\}$ , and that  $\mathcal{T}_i = \{0, 1, \dots, ts - 1\}$  for each  $i \in \mathcal{G}$ , we now have the following:

**THEOREM 4.1 (SOUNDNESS: EQUALITY ABSTRACTION).** *Let  $P$  be a KPL kernel executed by  $gs$  groups and  $ts$  threads per group. If no execution of  $P$  under the equality abstraction with respect to any  $(\mathcal{G}', \mathcal{T}')$  (representing a pair of threads) leads to error, then no concrete execution of  $P$  leads to error.*

**PROOF.** We prove the contrapositive. Thus, let us assume that  $P$  has a concrete execution leading to *error*. We will show that there also exists an execution for  $P$  under the equality abstraction that leads to *error*.

Let the concrete execution of  $P$  that leads to *error* be:

$$\rho \triangleq \langle (\kappa_0, ss_0), \dots, (\kappa_{n-1}, ss_{n-1}) \rangle.$$

We now construct an execution of  $P$  under the equality abstraction with respect to some  $(\mathcal{G}', \mathcal{T}')$  (representing a pair of threads) that also leads to *error*.

Observe that in the final state  $(\kappa_{n-1}, ss_{n-1})$ , either rule K-INTER-GROUP-RACE, K-INTRA-GROUP-RACE or K-DIVERGENCE applies, because only these rules can lead to *error*. In each case, exactly two threads are responsible for the rule being applicable. Assume that these threads are  $s = (i_1, j_1)$  and  $t = (i_2, j_2)$ , with the first component of each pair being the group id of the thread and the second component being the local id. Define  $\mathcal{G}'$  and  $\mathcal{T}'$  as follows:

- if  $i_1 = i_2$  then  $\mathcal{G}' = \{i_1\}$  and  $\mathcal{T}_{i_1} = \{j_1, j_2\}$ , and
- if  $i_1 \neq i_2$  then  $\mathcal{G}' = \{i_1, i_2\}$ ,  $\mathcal{T}_{i_1} = \{j_1\}$  and  $\mathcal{T}_{i_2} = \{j_2\}$ .

We now define a function  $proj^{s,t}$  that takes a kernel state and filters out all components not pertaining to  $s$  and  $t$ :

$$proj^{s,t}(\kappa, ss) \triangleq (\kappa', ss)$$

where  $dom(\kappa') = \mathcal{G}'$  and  $\kappa'(i) = proj_i^{s,t}(\kappa(i))$  for all  $i \in \mathcal{G}'$ . For group states, we define:

$$proj_i^{s,t}(\gamma_{ts}, R, W) \triangleq (\gamma_{ts} \upharpoonright_{\mathcal{T}'_i}, R \cap (\mathbb{N} \times \mathcal{T}'_i), W \cap (\mathbb{N} \times \mathcal{T}'_i))$$

where  $f \upharpoonright_A$  denotes the restriction of the domain of  $f$  to  $A$ . The intersection with  $\mathbb{N} \times \mathcal{T}'_i$  restricts the read and write sets of the group to only those accesses by the thread(s) in  $\mathcal{T}'_i$ .

We next define a mapping  $\mathcal{E}$  which, given a concrete execution  $\rho$ , will yield an execution  $\mathcal{E}(\rho)$  under the equality abstraction. For the most part, it simply maps each successive kernel state  $(\kappa_i, ss_i)$  to the state  $proj^{s,t}(\kappa_i, ss_i)$ . Complexity arises when we reach a state that satisfies the following condition:

Rule K-ITER applies, but the predicate  $p \wedge e$  holds for neither  $s$  nor  $t$ . (\*)

The above situation occurs when, during the execution of a loop in a concrete execution, the guard no longer holds for threads  $s$  and  $t$ , but does hold for at least one other thread. The appropriate rule in the concrete execution is K-ITER, but in the abstracted execution, we must invoke the K-DONE rule instead to exit the loop. Accordingly, when  $\mathcal{E}$  reaches a state  $(\kappa_i, ss_i)$  where the condition (\*) holds, it discards all of the following states up to and including the first state  $(\kappa_j, ss_j)$  where K-DONE applies and  $ss_j = ss_i$ . Here, the condition  $ss_j = ss_i$  ensures that this is the end of the *current* loop, not the end of a loop nested in the current one.

The definition of  $\mathcal{E}$  is as follows:

$$\mathcal{E}(\langle \rangle) \triangleq \langle \rangle$$

$$\mathcal{E}((\kappa, ss) : \rho) \triangleq \begin{cases} \text{proj}^{s,t}(\kappa, ss) : \mathcal{E}(\text{dropUntil}(\text{loopEnd}_{\kappa,ss}, \rho)) & \text{if } (\kappa, ss) \text{ satisfies (*)} \\ \text{proj}^{s,t}(\kappa, ss) : \mathcal{E}(\rho) & \text{otherwise} \end{cases}$$

where

$$\text{dropUntil}(p, \langle \rangle) \triangleq \langle \rangle$$

$$\text{dropUntil}(p, x : xs) \triangleq \begin{cases} xs & \text{if } p(x) \text{ holds} \\ \text{dropUntil}(p, xs) & \text{otherwise} \end{cases}$$

and where the predicate  $\text{loopEnd}_{\kappa,ss}(\kappa', ss')$  holds if and only if  $ss = ss'$  and K-DONE applies in state  $(\kappa', ss')$ .

We prove that  $\mathcal{E}(\rho)$  is an execution under the equality abstraction that leads to *error*, proceeding by induction on the length of  $\rho$ . Thus, suppose that  $\rho = (\kappa, ss) : \rho'$ . In the case where  $\rho'$  is empty—the base case—we obtain  $\langle \text{proj}^{s,t}(\kappa, ss) \rangle$ , which is clearly an execution of  $P$  under the equality abstraction. Moreover, since we have that  $(\kappa, ss)$  leads to *error*, it follows by our choice of  $s$  and  $t$  that so does  $\langle \text{proj}^{s,t}(\kappa, ss) \rangle$ .

For the induction step, let us first assume that the condition (\*) does *not* apply to  $(\kappa, ss)$ . We must now show that  $\text{proj}^{s,t}(\kappa, ss) : \mathcal{E}(\rho')$  is an execution of  $P$  under the equality abstraction. Hence, assume that  $\rho' = (\kappa', ss') : \rho''$ . The induction hypothesis guarantees that  $\mathcal{E}((\kappa', ss') : \rho'')$  is an execution under the equality abstraction leading to *error*. It remains to demonstrate that  $\text{proj}^{s,t}(\kappa, ss)$  can take an abstract step to  $\text{proj}^{s,t}(\kappa', ss')$ . Since we know that  $(\kappa, ss)$  takes a concrete step to  $(\kappa', ss')$ , we can complete this part of the proof by a case distinction on the rules from Figure 6, i.e., the abstract step simply invokes the equality-abstracted version of the rule invoked by the concrete step, where in the case of a barrier, the arbitrary shared state  $sh'$  in the rule G-SYNC-EQ is chosen to simulate the concrete kernel state.

If the condition (\*) *does* apply to  $(\kappa, ss)$ , then we must show that

$$\mathcal{E}(\text{dropUntil}(\text{loopEnd}_{\kappa,ss}, \rho'))$$

is an execution of  $P$  under the equality abstraction. In order to apply the induction hypothesis, we must know that the sequence  $\text{dropUntil}(\text{loopEnd}_{\kappa,ss}, \rho')$  is non-empty. Thus, suppose for the moment that the sequence is empty. The condition (\*) now implies that  $(\kappa, ss)$  is a state at the head of a loop, and that threads  $s$  and  $t$  are disabled throughout the remainder of the loop execution. The sequence being empty means that the concrete execution reaches *error* before exiting the loop. Yet, since the transition to *error* is brought about by threads  $s$  and  $t$ , at least one must be enabled: contradiction. Hence,  $\text{dropUntil}(\text{loopEnd}_{\kappa,ss}, \rho')$  is indeed non-empty and it must therefore have a first state  $(\kappa'', ss'')$ . We can now apply the induction hypothesis, to obtain that  $\mathcal{E}(\text{dropUntil}(\text{loopEnd}_{\kappa,ss}, \rho'))$  is an execution of  $P$  under the equality abstraction and leads to *error*. It remains to show that we can take an abstract step from  $\text{proj}^{s,t}(\kappa, ss)$

to  $proj^{s,t}(\kappa'', ss'')$ : this is justified by a single invocation of the K-DONE rule. The invocation is possible, as the shadow copies of the global memory as owned by threads  $s$  and  $t$  did not change between  $(\kappa, ss)$  and  $(\kappa'', ss'')$ . If the shadow copies did change, either the condition (\*) did not hold in  $(\kappa, ss)$ , or rule G-SYNC was applied to a group in  $\mathcal{G}'$  somewhere between  $(\kappa, ss)$  and  $(\kappa'', ss'')$ ; both of these contradict the assumption that the condition (\*) holds. In particular, the rule G-SYNC can only apply if all threads in the group are enabled, which is not the case according to condition (\*).

Finally, we note that since  $(\kappa_0, ss_0)$  is a valid initial kernel state, so is  $proj^{s,t}(\kappa_0, ss_0)$  (the first state of  $\mathcal{E}(\rho)$ ). Hence,  $\mathcal{E}(\rho)$  is an execution of  $P$  under the equality abstraction leading to *error* and the result now follows.  $\square$

We shall now prove the soundness of the two-thread reduction when using the adversarial abstraction. Here, *execution under the adversarial abstraction with respect to*  $(\mathcal{G}', \mathcal{T}')$  is defined to be a non-empty sequence of kernel states where each successive pair is related by the  $\rightarrow_{k, \mathcal{A}(\mathcal{G}', \mathcal{T}')}$  relation.

**THEOREM 4.2 (SOUNDNESS: ADVERSARIAL ABSTRACTION).** *Let  $P$  be a KPL kernel executed by  $gs$  groups and  $ts$  threads per group. If no execution of  $P$  under the adversarial abstraction with respect to any  $(\mathcal{G}', \mathcal{T}')$  (representing a pair of threads) leads to *error*, then no concrete execution of  $P$  leads to *error*.*

**PROOF.** By replaying the proof of Theorem 4.1, with G-SYNC-EQ replaced by G-SYNC-ADV and where in the case of K-BASIC we observe that the arbitrary values chosen by rules T-READ-ADV and T-WRITE-ADV can be chosen in accordance with the shared state and evaluation of local expressions, respectively.  $\square$

*Remark 4.3.* Observe that as the contents of the shared state is not actually used by the adversarial abstraction, we can completely omit the shared state from the semantics. Our verification method presented in the next section exploits this observation.

*Incompleteness of the Two-Thread Reduction.* No completeness results corresponding to the above two soundness results hold, as the arbitrary choices employed in the abstractions allow for strictly more behavior. Moreover, as already observed by Habermair and Knapp [2012], the abstractions allow for possible detection of races and barrier divergence in code not reachable during real kernel executions due to an infinite loop preceding the offending code. For instance, the following KPL fragment does *not* have a race on the shared location 42 because the failure of thread 0 to exit the loop prevents any other thread from writing to the location, but when only two arbitrary threads are considered, we cannot deduce that the loop does not terminate.

```
while (lid = 0) continue;
sh[42] := lid
```

#### 4.2. Use of the Isabelle Proof Assistant

As explained in Section 3.3, we used the Isabelle proof assistant [Nipkow et al. 2002] to formalize all the definitions presented in Section 3; we have also formalized the further definitions of this section using Isabelle. Even without conducting any proofs, the strict typing discipline imposed by Isabelle allowed us to find and fix several ambiguities in the original definitions.

We had intended to formalize the soundness results in Theorems 4.1 and 4.2, but found the task to be more challenging than anticipated, mainly due to the complexity of the KernelStates type and the sheer number of execution rules. The requirement for ‘fresh’ variables in several of the execution rules proved particularly challenging; indeed, “if one wants to formalize such proofs in a theorem prover, then dealing with

binders, renaming of bound variables, capture-avoiding substitution, etc., is very often a major problem” [Urban and Narboux 2009].

Nonetheless, the use of Isabelle in a ‘top-down’ manner, to explore the high-level structure of the proof rather than to prove each statement from first principles, was very helpful. We were able to significantly improve the quality of an earlier hand proof of Theorem 4.1, and gain increased confidence in its correctness. For instance, to make the abstract semantics successfully simulate the concrete semantics, the  $R$  and  $W$  components of each group state must tag each accessed location with the local id of the accessing thread, so that those locations not accessed by thread  $s$  or  $t$  can be identified and removed. This observation was made as a direct result of our failure to prove a lemma in Isabelle.

Our Isabelle proof script, which comprises the definitions presented in Sections 3 and 4, plus an outline proof of Theorem 4.1, is available online.<sup>9</sup>

## 5. THE DESIGN AND IMPLEMENTATION OF GPUVERIFY

Armed with the SDV semantics and abstractions of Section 3, we now consider the problem of verifying that GPU kernels are data race- and barrier divergence-free. For this purpose, we have designed a tool, GPUVerify, built on top of the Boogie verification system [Barnett et al. 2005]. Boogie takes a program annotated with loop invariants and procedure contracts, and decomposes verification into a set of formulas to be checked automatically by the Z3 theorem prover [de Moura and Bjørner 2008]. We describe the challenges associated with automatically translating OpenCL and CUDA kernels into a Boogie intermediate representation (Section 5.1), a technique for transforming this Boogie representation of the kernel into a standard sequential Boogie program whose correctness implies data race- and barrier divergence-freedom of the original kernel (Section 5.2), and a method for automatically inferring invariants and procedure contracts to enable automatic verification (Section 5.3).

### 5.1. Translating OpenCL and CUDA into Boogie

To allow GPUVerify to be applied directly to source code we have implemented a compiler that translates GPU kernels into an intermediate Boogie form. Our compiler is built on top of the Clang/LLVM infrastructure, which supports both OpenCL and CUDA. We target CUDA 5.0 [Nvidia 2012a], and versions of OpenCL up to 1.2 [Khronos OpenCL Working Group 2012].

There were effectively three challenges with respect to compilation. First, many industrial applications utilize features of OpenCL and CUDA not present in vanilla C, such as declaring variables as vector or image types, or calling intrinsic functions; we therefore invested significant engineering effort into designing equivalent Boogie types and functions. Second, the Boogie language does not support floating point values directly, thus we modeled them abstractly via uninterpreted functions. This sound over-approximation can in principle lead to false positives, but we have not found this to be a problem in practice since floating point operators are not used to compute array indexing expressions or to determine synchronization-related control flow. The third issue, namely handling of pointers, is technically more interesting and is now discussed in depth.

*Modeling Pointers.* Boogie is a deliberately simple intermediate language, and does not support pointer data types natively. We have devised an encoding of pointers in Boogie which we explain using an example. For readability we use C-like syntax rather than the Boogie input language.

<sup>9</sup><http://multicore.doc.ic.ac.uk/tools/GPUVerify/IsabelleFormalisation/>

Table V. Translating pointer usage into Boogie

Source	Generated Boogie
<code>p = A;</code>	<code>p = int_ptr(A_base, 0);</code>
<code>q = p;</code>	<code>q = p;</code>
<code>foo(p);</code>	<code>foo(p);</code>
<code>p = q + 1;</code>	<code>p = int_ptr(q.base, q.offset + 1);</code>
<code>x = p[e];</code>	<pre> <b>if</b> (p.base == A_base)   x = A[p.offset + e]; <b>else if</b> (p.base == B_base)   x = B[p.offset + e]; <b>else</b>   <b>assert</b>(false); </pre>
<code>p[e] = d;</code>	<pre> <b>if</b> (p.base == A_base)   A[p.offset + e] = d; <b>else if</b> (p.base == B_base)   B[p.offset + e] = d; <b>else</b>   <b>assert</b>(false); </pre>

Suppose a kernel declares exactly two character arrays (in any memory space) and two character pointers:

```

char A[1024], B[1024];
char *p, *q;

```

In this case GPUVerify generates the following types:

```

enum int_ptr_base = { A_base, B_base, null };

struct int_ptr {
  int_ptr_base base;
  int offset;
};

```

Thus an integer pointer is modeled as a pair consisting of a base array, or the special value `null` if the pointer is null, and an integer offset from this base. To cater for arbitrary type-casting, the offset is in terms of bytes rather than array elements. GPUVerify incorporates an optimization to transform offsets to work at the level of array elements when it can be deduced using a simple type-based analysis that an array is always accessed at element-level granularity.

Pointers `p` and `q` can be assigned offsets from `A` or `B`, can be assigned `null`, or can be left uninitialized. An uninitialized pointer has an arbitrary value in the generated Boogie code. Table V shows how uses of `p` and `q` are translated into Boogie.

Statement `p = q + 1` demonstrates that pointer arithmetic is straightforward to model using this encoding. Pointer read and writes are modeled by a case split on all the possible bases for the pointer being dereferenced. If no base matches then the pointer is either uninitialized or `null`. These illegal dereferences are captured by an assertion failure. This encoding exploits the fact that in GPU kernels there are a finite, and usually small, number of explicitly declared pointer targets.

We deal with stack-allocated local variables whose addresses are taken by rewriting these variables as arrays of length one, and transforming the corresponding accesses to such variables appropriately. This suffices by the fact that GPU kernel languages do not permit recursion.

*Points-to Analysis.* The case-split associated with pointer dereferences can hamper verification of kernels with pointer-manipulating loops, requiring loop invariants that narrow down the permissible arrays to which pointers can refer. To avoid this in many

<pre> <b>void barrier();</b>  <b>shared bool gr[SZ][SZ];</b>  <b>void kernel() {</b>   <b>int k = 0;</b>   <b>while (k &lt; SZ) {</b>     <b>if (!gr[lidY][lidX]) {</b>       <b>if (gr[lidY][k] &amp;&amp; gr[k][lidX]) {</b>         gr[lidY][lidX] = true;       }     }     <b>barrier();</b>     k++;   } } </pre> <p style="text-align: center;">(a) Example kernel</p>	<pre> <b>void barrier();</b> <b>void LOG_RD_gr(int y, int x);</b> <b>void LOG_WR_gr(int y, int x);</b>  <b>shared bool gr[SZ][SZ];</b>  <b>void kernel() {</b>   <b>int k = 0;</b>   <b>while (k &lt; SZ) {</b>     LOG_RD_gr(lidY, lidX);     <b>if (!gr[lidY][lidX]) {</b>       LOG_RD_gr(lidY, k);       LOG_RD_gr(k, lidX);       <b>if (gr[lidY][k] &amp;&amp; gr[k][lidX]) {</b>         LOG_WR_gr(lidY, lidX);         gr[lidY][lidX] = true;       }     }     <b>barrier();</b>     k++;   } } </pre> <p style="text-align: center;">(b) Kernel after race instrumentation</p>
---	--

Fig. 8. Example illustrating how GPUVerify instruments a kernel to detect races

cases, we have implemented the flow- and context-insensitive pointer analysis algorithm of Steensgaard [1996]. Although this over-approximates the points-to sets, our experience with GPU kernels is that aliasing is scarce and therefore precision is high. Returning to the above example, suppose the points-to analysis determines that  $p$  may only refer to array  $A$  (or be null or uninitialized). In this case, the assignment  $p[e] = d$  is translated to:

```

if (p.base == A_base)
  A[p.offset + e] = d;
else
  assert(false);

```

As well as checking for dereferences of null or uninitialized pointers, the `assert(false)` ensures that potential bugs in our implementation of the points-to analysis do not lead to unsound verification.

## 5.2. Reducing Data Race- and Barrier Divergence-Checking to Sequential Program Verification

Having compiled an OpenCL or CUDA kernel into corresponding Boogie form, GPUVerify attempts to verify the kernel. We describe the verification strategy employed by GPUVerify using a worked example.

Consider the kernel of Figure 8(a), adapted from part of a C++ AMP application that computes the transitive closure of a graph using Warshall's algorithm, and simplified for ease of presentation. The kernel is written for a single, 2-dimensional group of  $SZ \times SZ$  threads. The local id of a thread is 2D, with  $x$  and  $y$  components `lidX` and `lidY`, respectively. The kernel declares a 2D shared array of Booleans, `gr`, representing the adjacency matrix of a graph.

*Access Logging Instrumentation.* A kernel is first instrumented with calls to procedures that will log accesses to shared arrays. Figure 8(b) shows the example kernel of Figure 8(a) after access logging instrumentation. Observe for example that the condition `gr[lidY][k] && gr[k][lidX]` involves two read accesses to `gr`, and is thus preceded by two calls to `LOG_RD_gr`.<sup>10</sup>

*Reduction to a Pair of Threads.* After access logging, the kernel must be translated into a form that models the predicated execution of multiple threads. Initially, we attempted to encode the SDV semantics of Section 3 directly, modeling lock-step execution of *all* threads. Unfortunately, modeling in this way required heavy use of quantifiers, especially for implementing the G-SYNC rule of Figure 5 and associated  $\text{merge}_i$  function. This led to Boogie programs outside the decidable theory supported by the Z3 theorem prover. As a result, verification of small (micro-sized) kernels took in the order of minutes, while verification attempts for large kernels quickly exhausted memory limits.

To achieve better scalability, we instead use the two-thread reduction of Section 4 to transform a kernel into a form where the predicated execution of only *two* threads is modeled. The local and group identities of the two threads are represented by symbolic constants, constrained such that the threads are either in different groups, or in the same group but with different local ids.

Because a two-threaded predicated program with lock-step execution is essentially a *sequential* program consisting of parallel assignments to pairs of variables, reasoning about GPU kernels at this level *completely avoids* the problem of exploring interleavings of concurrent threads, and allows us to leverage existing techniques for reasoning about sequential programs. Furthermore, because the ids of the threads under consideration are symbolic, the sequential program can only be proven correct if correctness holds for *all* thread id combinations. In practice this is achieved by generating verification conditions that are discharged to an SMT solver.

GPUVerify supports both the adversarial and equality abstractions studied in Section 4. The equality abstraction is essential in verification of some kernels (including the kernel of Figure 8(a)) whose race-freedom hinges on threads agreeing on the values read from certain shared locations. We show in Section 6.2.3 that using the adversarial abstraction, when it suffices, typically proves to be more efficient than employing the equality abstraction. GPUVerify chooses the abstraction on array-by-array basis. We have implemented an inter-procedural taint analysis to over-approximate those shared arrays whose values may influence control flow. Arrays that may influence control flow are handled using the equality abstraction and all others using the adversarial abstraction. We study the precision of this heuristic experimentally in Section 6.2.3.

While the equality or adversarial abstractions suffice for verification of the vast majority of kernels we have studied, the equality abstraction is not sufficient when correctness depends upon richer properties of the shared state. For instance, suppose a kernel declares shared arrays `A` and `B`, and includes a statement:

```
A[B[lid]] = ...
```

Write-write race freedom of `A` requires that `B[i] != B[j]` for all distinct `i` and `j`. In practice, we have found that this prohibits verification of a number of kernels that perform a prefix sum operation into an array `B`, and then use `B` to index into an array `A` as shown above. A richer shared state abstraction is investigated by Chong et al. [2013].

<sup>10</sup>For ease of presentation we ignore here the fact that, due to short-circuit evaluation, the read from `gr[k][lidX]` will not be issued if the value read from `gr[lidY][k]` is *false*. Our implementation does handle short-circuit evaluation correctly.

```

void barrier(bool en1, bool en2);
void LOG_RD_gr(bool en1, int y1, int x1, bool en2, int y2, int x2);
void LOG_WR_gr(bool en1, int y1, int x1, bool en2, int y2, int x2);
bool gr1[SZ][SZ], gr2[SZ][SZ];

void kernel() {
  int k1, k2;
  bool LC1, LC2, P1, P2, Q1, Q2; // Predicates
  // Assume the two group ids lie in appropriate range;
  // trivial in this example as there is only one group:
  assume(0 <= gid1 && gid1 < 1 && 0 <= gid2 && gid2 < 1);
  // Assume the thread ids lie in appropriate range:
  assume(0 <= lidX1 && lidX1 < SZ && 0 <= lidX2 && lidX2 < SZ);
  // Assume that if the two threads are in the same group (which, in this
  // single-group example, they must be) their local ids are distinct:
  assume((gid1 == gid2) ==> (lidX1 != lidX2 || lidY1 != lidY2));
  k1, k2 = 0, 0;
  LC1, LC2 = k1 < SZ, k2 < SZ;
  while (LC1 || LC2) {
    LOG_RD_gr(LC1, lidY1, lidX1, LC2, lidY2, lidX2);
    P1, P2 = LC1 && !gr1[lidY1][lidX1],
             LC2 && !gr2[lidY2][lidX2];
    LOG_RD_gr(P1, lidY1, k1, P2, lidY2, k2);
    LOG_RD_gr(P1, k1, lidX1, P2, k2, lidX2);
    Q1, Q2 = P1 && gr1[lidY1][k1] && gr1[k1][lidX1],
             P2 && gr2[lidY2][k2] && gr2[k2][lidX2];
    LOG_WR_gr(Q1, lidY1, lidX1, Q2, lidY2, lidX2);
    gr1[lidY1][lidX1], gr2[lidY2][lidX2] =
      Q1 ? true : gr1[lidY1][lidX1],
      Q2 ? true : gr2[lidY2][lidX2];
    barrier(LC1, LC2);
    k1, k2 = LC1 ? k1 + 1 : k1, LC2 ? k2 + 1 : k2;
    LC1, LC2 = LC1 && k1 < SZ, LC2 && k2 < SZ;
  }
}

```

Fig. 9. The kernel of Figure 8(a) after transformation to two-thread predicated form

Figure 9 shows the result of transforming the access-instrumented version of the kernel (Figure 8(b)) into a form where the predicated execution of a pair of arbitrary, distinct threads is modeled, using the equality abstraction. The transformation using the adversarial abstraction is identical, except that the arrays `gr1` and `gr2` are eliminated, and reads from these arrays are made non-deterministic.

The two-dimensional local id of the first thread is represented by `lidX1, lidY1`, and similarly for the second thread. The one-dimensional group id for the first thread is represented by `gid1`, and similarly for the second thread.

The first `assume` statement constrains the group ids to lie within an appropriate range. As our worked example contains only a single group, the `assume` constrains each of `gid1` and `gid2` to be equal to zero.

The second `assume` statement constrains the local id components to lie within appropriate ranges. The two-dimensional group of threads in our example has dimensions  $SZ \times SZ$ , thus each local id component must lie in the range  $[0..SZ - 1]$ .

The final `assume` statement forces the group and thread ids to be chosen such that if the threads are in the same group they cannot have identical local ids. The left-hand-side of the implication holds in our example because, with only a single group, `gid1` and `gid2` must be equal; thus the `assume` statement forces a difference in at least one of the pairs `lidX1, lidX2` or `lidY1, lidY2`.

Local variable  $k$  is duplicated, and the assignment  $k = 0$  is replaced with a parallel assignment, setting  $k_1$  and  $k_2$  to zero. The kernel declares fresh Boolean variables  $LC$ ,  $P$  and  $Q$  (duplicated for each thread). These are used to model predicated execution of the *while* loop ( $LC$ ) and the outer and inner conditionals ( $P$  and  $Q$  respectively). In the examples of Section 2, and in the operational semantics of Section 3, we specified that under predicated execution, a *while* loop should continue to execute while there exists a thread for which the condition holds. In the presence of just two threads, existential quantification turns into disjunction, hence the loop condition  $LC_1 \vee LC_2$ .

In Figure 9, parameters to the `LOG_RD_gr` and `LOG_WR_gr` procedures are duplicated, with a parameter being passed for each thread. In addition, a *predicate* parameter, `en`, is passed for each thread, recording whether the thread is enabled during the call (cf. the incoming predicate  $p$  in rule K-CALL of Figure 6). If `LOG_RD_gr` is called with *false* as its `en1` parameter, this indicates that the first thread is not enabled, and thus a read should not be logged for this thread. Similarly, `barrier` is equipped with a pair of predicate parameters, `en1` and `en2`.

*Handling Multiple Procedures.* During the transformation to two-threaded form, the parameter list of each user-defined procedure is duplicated, and (as with the `LOG` and `barrier` procedures) enabled predicates are added for each thread. The procedure body is then translated to two-threaded, predicated form, with every statement guarded by the *enabled* predicate parameters. Correspondingly, actual parameters are duplicated at call sites, and the current predicates of the execution passed as *enabled* parameters.

*Checking for Barrier Divergence.* Under the two-thread encoding, inserting a check for barrier divergence is trivial: the `barrier` procedure merely asserts that its arguments `en1` and `en2` are equal if the group ids of the two threads correspond. This two-threaded version of rules G-DIVERGENCE (Figure 5) and K-DIVERGENCE (Figure 6) precisely matches the notion of barrier divergence presented formally in Section 3. We may wish to *only* check barrier divergence-freedom for a kernel, if verifying race-freedom proves too difficult. This is sound under *adversarial* abstraction, where every read from the shared state returns an arbitrary value. A kernel that can be shown to be barrier divergence-free under this most general assumption is guaranteed to be barrier divergence-free under any schedule of shared state modifications. If we prove barrier divergence-freedom for a kernel under the equality abstraction, we can conclude a weaker property than barrier divergence-freedom: that barrier divergence cannot occur unless a data race has occurred. Note that our barrier divergence checking is stricter than that attempted by the PUG verifier of Li and Gopalakrishnan [2010], which merely requires threads which follow different conditional paths through a kernel to pass the same number of barriers.<sup>11</sup> While PUG reports microkernels exhibiting the barrier divergence bugs discussed in Section 2 as successfully verified, such kernels are rejected by GPUVerify.

*Checking for Races.* The `LOG_RD` and `LOG_WR` procedures are responsible for manipulating a read and write set for each thread, and for each of the shared arrays of the kernel. According to the semantics of Section 3 (rules G-RACE and K-RACE of Figure 5 and 6, respectively), race checking involves asserting for each array  $A$  that the read and write sets for  $A$  do not conflict between threads.

We encode read and write sets efficiently by exploiting non-determinism, similar to a method used in prior work by Donaldson et al. [2010; 2011]. For each shared array  $A$

<sup>11</sup>In a subsequent paper on dynamic symbolic execution of CUDA kernels Li et al. [2012b] improve this check to restrict to textually aligned barriers.

with index type  $T$  we introduce the following *instrumentation variables* for each thread  $i$  under consideration (where  $i \in \{1, 2\}$ ):

```
WR_exists_Ai : bool
WR_elem_Ai : T
RD_exists_Ai : bool
RD_elem_Ai : T
```

The Boolean  $WR\_exists\_A_i$  is set to *true* only if the write set of thread  $i$  for  $A$  is non-empty. In this case,  $WR\_elem\_A_i$  represents one element of this write set: an index into  $A$ . The corresponding RD variables for read sets are similar.

Initially  $WR\_exists\_A_i$  and  $RD\_exists\_A_i$  are *false* for each thread because the read/write sets are empty. The  $LOG\_WR\_A$  procedure then works as follows: for each thread  $i$ , if  $i$  is enabled on entry to the procedure (predicate parameter  $eni$  is *true*), then the thread non-deterministically chooses to do *nothing*, or to set  $WR\_exists\_A_i$  to *true* and  $WR\_elem\_A_i$  to the index being logged. Procedure  $LOG\_RD\_A$  operates similarly. This strategy ensures that if  $WR\_exists\_A_i$  holds, then  $WR\_elem\_A_i$  is the index of an *arbitrary* write to  $A$  performed by thread  $i$ . Checking absence of write-write races can then be achieved by placing the following assertion in the  $LOG\_WR\_A$  procedure:

```
assert (!(WR_exists_A1 && WR_exists_A2 && WR_elem_A1 == WR_elem_A2))
```

A similar assertion is used to check read-write races; the procedure  $LOG\_RD\_A$  works analogously.

Because this encoding tracks an *arbitrary* element of each read and write set, if the sets can have a common, conflicting element this will be tracked by both threads along some execution trace, and the generated assertion will fail along this trace. If we can prove for every array that the associated assertions can *never* fail, we can conclude that the kernel is race-free.

*Inter- vs. Intra-Group Race Checking.* A novel contribution of this paper over prior work is a treatment of both intra- and inter-group races; prior work focused only on the intra-group case. We briefly recall how these kinds of races are handled by our semantics (Section 3.2). Intra-group data races are detected by rule G-RACE (Figure 5) which compares the read and write sets of individual threads; an intra-group race detected by G-RACE is detected at the kernel level by rule K-INTRA-GROUP-RACE (Figure 6). The role of barrier operations in avoiding races within a work group is captured at the group level by rule G-SYNC (Figure 5): this rule accumulates the read and write sets of individual threads to form group-wide read and write sets, and then clears the individual thread read and write sets to reflect the fact that synchronization has taken place. The group-level read and write sets capture the fact that data races between threads in distinct groups *cannot* be avoided using barriers; rule K-INTER-GROUP-RACE of Figure 6 detects inter-group races by comparing these sets.

We capture the intent of these semantic rules in our implementation by distinguishing between the cases where the two threads under consideration are in the same or different groups. If the threads are in the *same* group then the instrumentation variables ( $WR/RD\_exists$ ,  $WR/RD\_elem$ , across all arrays) play the role of thread-level read/write sets for checking intra-group data races. These sets must be cleared at barriers, and there is no need to record additional group-level read/write sets (because no further threads are directly considered). If the threads are in *different* groups then the instrumentation variables play the role of group-level read/write sets for checking inter-group data races. These sets must *not* be cleared at barriers, and there is no need to record additional thread-level read/write sets.

Thus, at a barrier, the instrumentation variables are unaffected if the threads under consideration are in different groups. If the threads are in the same group, read and

write sets are cleared by assuming (using Boogie-level `assume` statements) that every `WR_exists` and `RD_exists` is *false*, i.e., by terminating all execution paths along which read or written elements were logged. This matches the thread-level read and write set handling of rule G-SYNC (Figure 5) and K-SYNC (Figure 6).

*Tolerating Benign Write-Write Races.* In practice it is quite common for threads to participate in benign write-write races, where identical values are written to a common location without synchronization. When equality abstraction is used, GPUVerify tolerates this kind of race by adding a conjunct to the above assertion to check that the values written are not equal.

### 5.3. Invariant Inference

GPUVerify produces a Boogie program akin to the transformed kernel of Figure 9, together with implementations of the `barrier` and all `LOG_RD/WR` procedures. This program must be verified in order to prove race- and barrier divergence-freedom of the original kernel. Verification hinges on finding inductive invariants for loops and contracts for procedures.

We have found that invariant generation using abstract interpretation over standard domains (such as intervals or polyhedra) is not effective in verifying GPU kernels. This is partly due to the data access patterns exhibited by GPU kernels, discussed in detail below, where threads do not tend to read or write from contiguous regions of memory, and also due to the predicated form of the programs produced by our verification method.

Instead, we use the Houdini algorithm of Flanagan and Leino [2001] as the basis for inferring invariants and contracts. Houdini is a method to find the largest set of inductive invariants from among a user-supplied pool of candidate invariants. Houdini works as a fixpoint procedure; starting with the entire set of invariants, it tries to prove that the current candidate set is inductive. The invariants that cannot be proved are dropped from the candidate set and the procedure is repeated until a fixpoint is reached. We briefly discuss the relationship between Houdini and other invariant generation techniques in Section 7.

By manually deducing invariants for a set of kernels (the *training* set described in our experimental evaluation, Section 6) we have devised a number of candidate generation rules which we outline below. We emphasize that the candidate invariants generated by GPUVerify are just that: *candidates*. The tool is free to speculatively generate candidates that later turn out to be incorrect: these are simply discarded by Houdini. A consequence is that incorrect or unintended candidates generated due to bugs in GPUVerify cannot compromise the soundness of verification.

Our candidate generation rules are purely heuristic. The only fair way to evaluate these carefully crafted heuristics is to evaluate GPUVerify with respect to a large set of unknown benchmarks. We presented such an evaluation in [Betts et al. 2012], and we summarize the findings of this evaluation in Section 6.1.

*Candidate Invariant Generation Rules.* The following rules (except for the final “variable is zero or a power of two” rule) relate to the manner in which a thread accesses shared data using its thread id. In each case, the rule can be applied where thread id denotes the local id of a thread within its group, or the global id of the thread across the kernel. We use `id` to generically describe both cases, and use `SZ` to denote the size of a thread group in the case that `id` refers to local id, and the total number of threads in the case that `id` refers to global id.

For clarity, we present the essence of each rule; the GPUVerify implementation is more flexible (e.g., being insensitive to the order of operands of commutative operations, and detecting when the id of a thread has been copied into another local vari-

able). For each of the rules associated with shared memory writes, there is an analogous rule for reads. We use  $\implies$  to denote implication.

**Rule: access thread id plus offset.** It is common for a thread to write to an array using its thread id, plus a constant offset (which is often zero) as index; this access pattern is illustrated by the following diagram:



*Observed pattern:*

—  $A[id + C] = \dots$  occurs in a loop

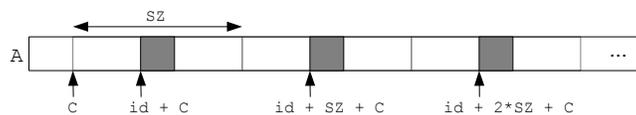
*Generated candidate:*

—  $WR\_exists\_A \implies WR\_elem\_A - C == id$

**Rule: access thread id plus strided offset.** When processing an array on a GPU, it is typically efficient for threads in a group to access data in a coalesced manner as in:

```
for (i = 0; i < 256; i++)
  A[i*SZ + id + C] = ...;
```

This access pattern is illustrated by the following diagram:



*Observed pattern:*

—  $A[id + i*SZ + C] = \dots$  occurs in a loop

—  $i$  is live at the loop head

*Generated candidate:*

—  $WR\_exists\_A \implies (WR\_elem\_A - C) \% SZ == id$

**Rule: access at thread id plus strided offset, with strength reduction.** This is similar to the previous rule. However, GPU programmers commonly apply the *strength reduction* operation manually, rewriting the above code snippet as follows:

```
for (i = id; i < 256*SZ; i += SZ)
  A[i + C] = ...;
```

In this case, the write set candidate invariant will not be inductive in isolation: the invariant  $(i \% SZ) == id$  is required in addition.

*Observed pattern:*

—  $i = id$  appears before a loop

—  $A[i+C] = \dots$  occurs in the loop

—  $i = i + SZ$  appears in the loop

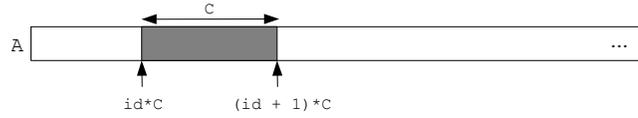
—  $i$  is live at the loop head

*Generated candidates:*

—  $(i \% SZ) == id$

—  $WR\_exists\_A \implies (WR\_elem\_A - C) \% SZ == id$

**Rule: access contiguous range.** It is common for threads to each be assigned a fixed-size chunk of an array to process. This access pattern is illustrated by the following diagram:



*Observed pattern:*

- $A[id \cdot C + i] = \dots$  occurs in a loop
- $i$  is live at the loop head

*Generated candidates:*

- $WR\_exists\_A \implies id \cdot C \leq WR\_elem\_A$
- $WR\_exists\_A \implies WR\_elem\_A < (id + 1) \cdot C$

**Rule: variable  $i$  is zero or a power of two.** GPU kernels frequently perform tree reduction operations on shared memory, as in the code snippet below. Race-freedom is ensured through the use of a barrier, together with a guard ensuring that threads which have dropped out of the reduction computation do not write to shared memory. Verifying race-freedom requires an invariant stating that the loop counter is a power of two, possibly allowing for the value to be zero. The access pattern for such a tree reduction with respect to a group of 8 threads ( $SZ == 8$ ) is illustrated below. A gray square containing a thread id indicates a memory access by the associated thread; dark gray indicates both a read and a write, while light gray indicates a read only.

```

for (i = 1; i < SZ; i *= 2) {
  if ((id % (2*i)) == 0) {
    A[id] += A[id + i];
  }
  barrier();
}

```

The diagram illustrates the access pattern for a tree reduction with 8 threads ( $SZ = 8$ ). It shows three rows of memory access for  $i=1$ ,  $i=2$ , and  $i=4$ . For  $i=1$ , threads 0-7 access elements 0, 0, 2, 2, 4, 4, 6, 6. For  $i=2$ , threads 0-7 access elements 0, 0, 4, 4. For  $i=4$ , threads 0-7 access elements 0. Barriers are shown between the rows.

*Observed pattern:*

- $i = i \cdot 2$  or  $i = i / 2$  occurs in a loop
- $i$  is live at the loop head

*Generated candidates:*

- $i \& (i-1) == 0$
- $i != 0$

The first candidate caters for the case where  $i$  is either a power of two or zero; here,  $\&$  denotes the bitwise-and operator. The second candidate encompasses the stronger condition where  $i$  is *not* zero.

GPUVerify includes a number of additional candidate generation rules that are intimately related to the details of our transformation of a kernel to a predicated sequential program. We omit details of these rules as they are very specific and less intuitive.

We have also designed rules to generate candidate pre- and post-conditions for procedures. We do not discuss these rules: although they allow us to perform modular verification of some GPU kernels, we find that for our current benchmarks (which are representative of the sizes of today's GPU kernels), full procedure inlining yields superior performance to modular analysis.

## 6. EXPERIMENTAL EVALUATION

We have gathered a set of 162 kernels which have been used to drive the development of GPUVerify and to evaluate the capabilities of the tool in practice. The benchmarks come from four suites:

- **AMD SDK**<sup>12</sup>: AMD Accelerated Parallel Processing SDK v2.6 consisting of 71 publicly available OpenCL kernels.
- **CUDA SDK**<sup>13</sup>: Nvidia GPU Computing SDK v2.0 consisting of 20 publicly available CUDA kernels.
- **C++ AMP**<sup>14</sup>: Microsoft C++ AMP Sample Projects consisting of 20 publicly available kernels, translated to CUDA.
- **Basemark CL**<sup>15</sup>: Rightware Basemark CL v1.1 consisting of 51 commercial OpenCL kernels,<sup>16</sup> provided to us under an academic license.

We consider the somewhat out-of-date version 2.0 of the Nvidia SDK to facilitate a direct comparison of GPUVerify with the PUG tool by Li and Gopalakrishnan [2010], the only existing publicly available verifier for CUDA kernels, since PUG is not compatible with more recent versions of the CUDA SDK (the CUDA verification tool described by Leung et al. [2012] is not publicly available). We also restrict our attention to the benchmarks from this SDK which were used to evaluate PUG. Because GPUVerify cannot directly analyze C++ AMP code, we retrieved the set of C++ AMP samples available online on 3 February 2012,<sup>14</sup> and manually extracted and translated the GPU kernel functions into corresponding CUDA kernels. This manual extraction and translation was straightforward.

We scanned each benchmark suite and removed kernels which are immediately beyond the scope of GPUVerify, either because they use atomic operations (7 kernels) or because they involve writes to the shared state using double-indirection as discussed in Section 5.2 (12 kernels). In recent work we have considered prototype support for atomic operations in GPUVerify [Bardsley and Donaldson 2014]. A richer shared state abstraction to handle double-indirection is proposed by Chong et al. [2013], but requires significant manual effort to apply in practice. The number of kernels quoted above for each benchmark suite correspond to the sizes of these suites *after* removal of these kernels.

In Section 6.1 we summarize the results of an experiment using this benchmark set for a comparison of our initial implementation of GPUVerify with PUG, and to evaluate the invariant inference capabilities of the initial implementation.

We then evaluate a more recent version of GPUVerify, equipped with new functionality for inter-group race checking (Section 6.2), assessing the performance of the tool across our benchmark set, the overhead associated with performing inter-group race checks, and the relative benefits in terms of precision and performance of the adversarial and equality abstractions.

Section 6.3 discusses a bug that we detected in an old CUDA SDK benchmark using GPUVerify. Section 6.4 demonstrates the effectiveness of the two-thread reduction.

### 6.1. Summary of Previous Evaluation

In [Betts et al. 2012] we presented a head-to-head comparison of the original GPUVerify implementation with the PUG tool of Li and Gopalakrishnan [2010], and evaluated the automatic invariant inference capabilities of GPUVerify. This evaluation focused on proving intra-group, but not inter-group, race-freedom, as well as barrier divergence-freedom. Instead of replicating this detailed evaluation here we provide a

<sup>12</sup><http://developer.amd.com/tools/heterogeneous-computing/amd-accelerated-parallel-processing-app-sdk/>

<sup>13</sup><https://developer.nvidia.com/cuda-toolkit-archive>

<sup>14</sup><http://blogs.msdn.com/b/nativeconcurrency/archive/2012/01/30/c-amp-sample-projects-for-download.aspx>

<sup>15</sup><http://www.rightware.com/benchmarking-software/product-catalog/>

<sup>16</sup>In [Betts et al. 2012] we reported results for 52 Rightware kernels, but we subsequently discarded one kernel from our benchmark set due to an issue which we have reported to Rightware.

brief summary of findings, preferring to focus on evaluating a more recent version of GPUVerify in Section 6.2.

The experimental evaluation we summarize here was performed on a PC with a 3.4GHz Intel Core i7-2600 CPU, 8GB RAM running Windows 7 (64-bit), with Z3 v3.3. All times reported are averages over 3 runs, using a 5 minute timeout per benchmark.

*Comparing GPUVerify with PUG.* We compared PUG and GPUVerify using the CUDA subset of our benchmarks: the CUDA SDK and C++ AMP suites (40 kernels), configuring PUG to use a 32-bit representation for integers, which is prescribed by the CUDA programming guide [Nvidia 2012a]. We used a synthetic bug injection procedure to randomly mutate each benchmark with an intra-group data race or barrier divergence bug to yield both correct and buggy versions of each benchmark.

Running the tools on correct kernels we found that PUG reported a false positive in three cases, where correctness depended upon threads agreeing on the contents of the shared state. GPUVerify is able to reason about these kernels using the equality abstraction (Section 4). The shared state abstraction of PUG is equivalent to our adversarial abstraction, which is not sufficient for these kernels. Note that GPUVerify decides automatically which shared state abstraction to use; we evaluate the heuristics GPUVerify uses to make this decision in Section 6.2.3. We found that PUG was on average faster than GPUVerify, and in six cases an order of magnitude faster. However, worst-case performance of PUG was significantly worse than that of GPUVerify: our timeout of 5 minutes was reached by PUG for six kernels.

Applying the tools to buggy versions of benchmarks we found that the proof attempts by GPUVerify generally failed within around 5 seconds, whereas the proof attempts by PUG usually failed within half a second: an order of magnitude faster. However, for seven buggy kernels we found that PUG reported *false negatives*: wrongly reporting correctness of the kernel. Of these false negatives, one mutation was an injected barrier divergence while the remaining six were intra-group data races. GPUVerify reported no false negatives.

*Evaluating Automatic Invariant Inference.* We used the following methodology to design and evaluate our invariant inference technique. We divided our benchmarks into two similarly-sized sets: a *training set* and an *evaluation set*, such that details of the evaluation set were unknown to all members of our team. We chose the CUDA SDK, C++ AMP and Basemark CL benchmarks as the training set (92 kernels), and the AMD SDK benchmarks as the evaluation set (70 kernels): members of our team had looked previously at the CUDA SDK and C++ AMP benchmarks, but not at the AMD SDK and Basemark CL benchmarks; however, we wanted to make the evaluation set publicly available, ruling out Basemark CL.

We manually analyzed all benchmarks in the training set, determining invariants sufficient for proving intra-group race- and barrier divergence-freedom. We then distinguished between ‘bespoke’ invariants: complex, kernel-specific invariants required by individual benchmarks; and ‘general’ invariants, conforming to an identifiable pattern that cropped up across multiple benchmarks. The general invariants led us to devise the invariant inference heuristics described in Section 5.3. We implemented these heuristics in GPUVerify and tuned GPUVerify to maximize performance on the training set. We then applied GPUVerify *blindly* to the evaluation set to assess the extent to which our inference technique enabled fully automatic analysis of the AMD SDK kernels. We believe that this approach of applying GPUVerify unassisted to a large, unknown set of benchmarks provides a fair evaluation of the automatic capabilities of the tool.

Using the inference techniques devised with respect to the training set (cf. Section 5.3), GPUVerify was able to verify 49 out of the 70 kernels from the evaluation

set (69%) *fully automatically*. Of these kernels, 48 were verified in 10 seconds or fewer, and the longest verification time was 17 seconds. With race checking disabled, GPUVerify was able to prove barrier divergence-freedom fully automatically for all evaluation benchmarks, in under 10 seconds per kernel. In [Betts et al. 2012] we present a detailed discussion of the reasons why verification failed for 22 of the evaluation kernels.

Many modern static analysis tools achieve low false alarm rates via a careful mixture of deliberately introduced unsoundness in the analysis and ad hoc warning suppression [Bessey et al. 2010]. GPUVerify does *not* follow this approach: the tool attempts to be a ‘real’ verifier, and thus will report verification failure for a kernel unless it was possible to construct a proof of correctness in a sound manner under bit-level accuracy. With this in mind, we believe that being able to verify 49 out of 70 evaluation kernels is a good result.

## 6.2. Evaluation of GPUVerify

Building on the evaluation presented in [Betts et al. 2012], we manually investigated each benchmark for which intra-group race-freedom could not be automatically verified. In each case, we either provided necessary invariants manually, to allow the kernel to verify, or improved the invariant inference capabilities of GPUVerify by adding additional candidate generation rules or generalizing existing rules. These additional and changed candidates all intimately relate to the details of our transformation from a kernel to a predicated sequential program, and details are omitted for that reason. We then turned on the inter-group race checking capabilities of GPUVerify that are novel in this work (see Section 5.2), and again manually supplied invariants and improved invariant inference capabilities until inter-group race-freedom could also be verified across the set of benchmarks.

In three cases, we had to make minor simplifications to allow correct kernels to be verified within our framework. In two binary search kernels from the AMD SDK we had to insert *assume* statements to encode the precondition that the input array is strictly sorted. Another AMD SDK kernel which computes a quasi random sequence exhibited an intricate, intentional benign read-write race which is beyond the scope of the benign tolerance mechanism mentioned in Section 5.2; we inserted barriers to eliminate this race.

We also had to modify a number of examples to remove data races. In a histogram example from the AMD SDK we found an arguably benign data race where threads non-atomically increment the same bucket of a histogram without synchronization. Such races can lead to increments to a bucket being lost non-deterministically, but in many uses of histograms this does not matter. For purposes of analysis we simplified this example so that threads write to disjoint locations; this changes the meaning of the kernel, but retains the features of the kernel that determine whether it is a challenging candidate for proving race-freedom. An AMD SDK kernel implementing a pass of the Floyd-Warshall algorithm adds together input values, which can cause an overflow-induced data race for extremely large inputs. Contacts at AMD regarded this data race report as a false alarm which we suppressed through the use of an addition operator that assumes non-overflowing behavior. Using GPUVerify we found a *non-benign* data race bug in one AMD SDK benchmark, `DwtHaar1D`, arising from an erroneous array indexing expression, which has since been independently fixed by AMD in a more recent version of the SDK.

With respect to this refined benchmark set and tool chain, we now report on:

- The raw performance of GPUVerify across our benchmarks (Section 6.2.1)

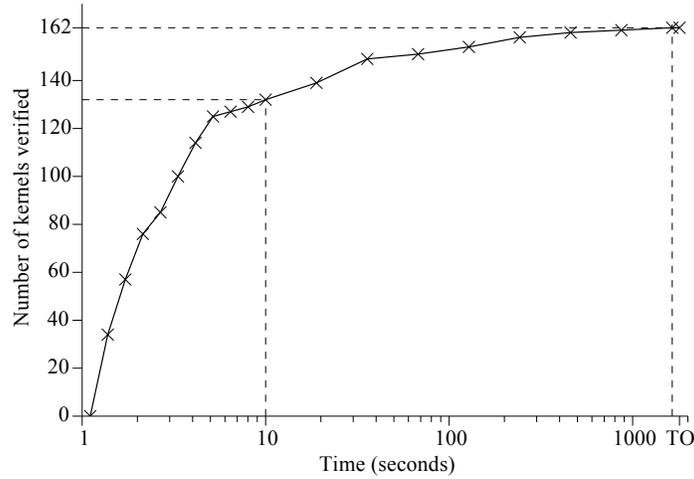


Fig. 10. Cumulative histogram showing the time taken to verify data race- and barrier divergence-freedom using GPUVerify across our benchmark set.

- The overhead associated with performing inter-group race checking in addition to intra-group race checking (Section 6.2.2)
- The performance and precision differences associated with using the equality vs. adversarial shared state abstractions (Section 6.2.3)

Experiments were performed on a PC with a 1.15GHz AMD Phenom CPU, 5.8GB RAM running Ubuntu (64-bit), using the nightly build of GPUVerify from 8 August 2013, and Z3 v4.3.1. All times reported are averages over 10 runs, and a timeout of 1800 seconds (30 minutes) was used in all experiments.

**6.2.1. Performance of GPUVerify.** Figure 10 is a cumulative histogram showing the performance of GPUVerify with respect to the benchmark set. The  $x$ -axis plots the time (in seconds, on a log scale), and the  $y$ -axis plots the number of kernels. A point at position  $(x, y)$  indicates that for  $y$  of the kernels, verification took  $x$  seconds or fewer. The results show that GPUVerify is capable of rapidly analyzing the vast majority of benchmark kernels: in 132 out of 162 cases, verification took less than 10 seconds. The longest verification time was 1643 seconds, for analysis of an FFT kernel in the AMD SDK. Interestingly, the FFT kernel is loop-free, thus the high verification time stems solely from the cost of data race checking, and not from the inference of loop invariants.

**6.2.2. Overhead of Inter-Group Race Checking.** A new addition to GPUVerify over [Betts et al. 2012] is the ability to perform *inter*-group race checking in addition to checking for races between groups. Inter-group race checks are represented by additional assertions in the Boogie program generated by GPUVerify, which in turn lead to larger verification conditions. Our hypothesis was that enabling inter-group race checking would uniformly lead to a drop in the performance of GPUVerify.

Figure 11 compares the performance of GPUVerify with and without inter-group race checking enabled, across our benchmark set. Each point represents a benchmark. The  $x$ -axis plots the time in seconds taken to verify full (i.e., inter- and intra-group) race freedom; the  $y$ -axis the time in seconds taken to verify only intra-group race-freedom. Both axes are plotted using log scales. Thus a point with coordinates  $(x, y)$  corresponds to a kernel for which verification of full race-freedom took  $x$  seconds and verification of only intra-group race-freedom took  $y$  seconds. Points lying below/above

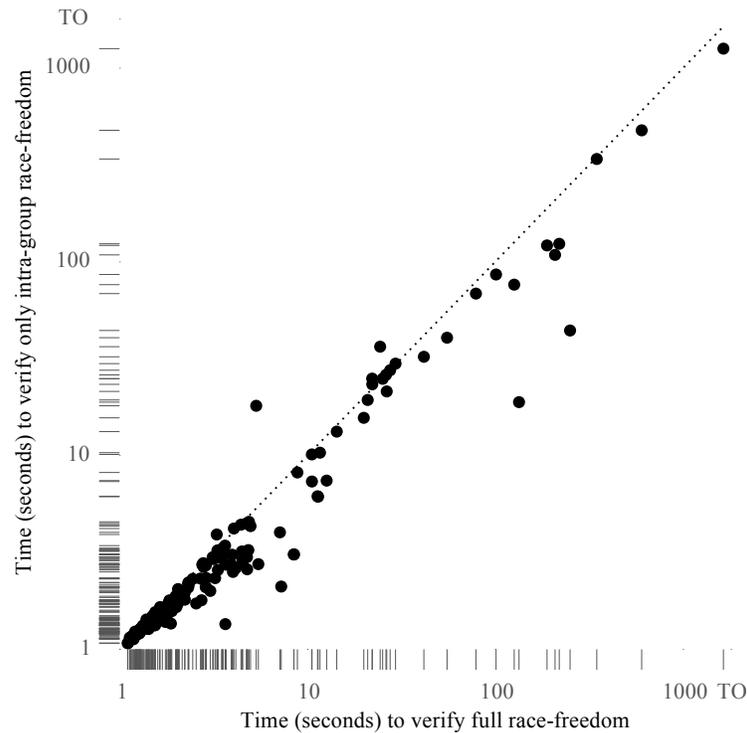


Fig. 11. Scatter plot comparing time taken to verify full race-freedom vs. only intra-group race-freedom across our benchmark set.

the diagonal support/rebut our hypothesis that there is computational overhead associated with adding inter-group race checks.

The data in Figure 11 strongly supports our hypothesis: in all but two cases, verification time either differs negligibly, or slows down when inter-group race checking is enabled.

There are two outliers which do not follow this trend: verification of a histogram kernel from the CUDA SDK takes  $3.4\times$  longer when only intra-group races are analyzed than when full race checking is performed; similarly a bitonic sort kernel from the AMD SDK takes  $1.5\times$  longer. We attribute these outliers to quirks in the implementation of the Z3 theorem prover: we tried the same experiment using a different theorem prover, CVC4 [Barrett et al. 2011], finding that the histogram kernel is not an outlier when CVC4 is used, leaving the bitonic sort kernel as the only outlier, with a slightly reduced slow-down of  $1.3\times$ .

**6.2.3. Adversarial vs. Equality Abstraction.** In Section 4 we discussed two shared state abstractions: the very coarse *adversarial* abstraction, and the slightly more refined *equality* abstraction. Treating an array using the adversarial abstraction has the potential advantage that the array can be *completely* removed from the Boogie program generated by GPUVerify. This reduces the extent to which the theorem prover must reason about arrays, which can be computationally expensive. Having observed this advantage working in practice for specific examples, we implemented the taint anal-

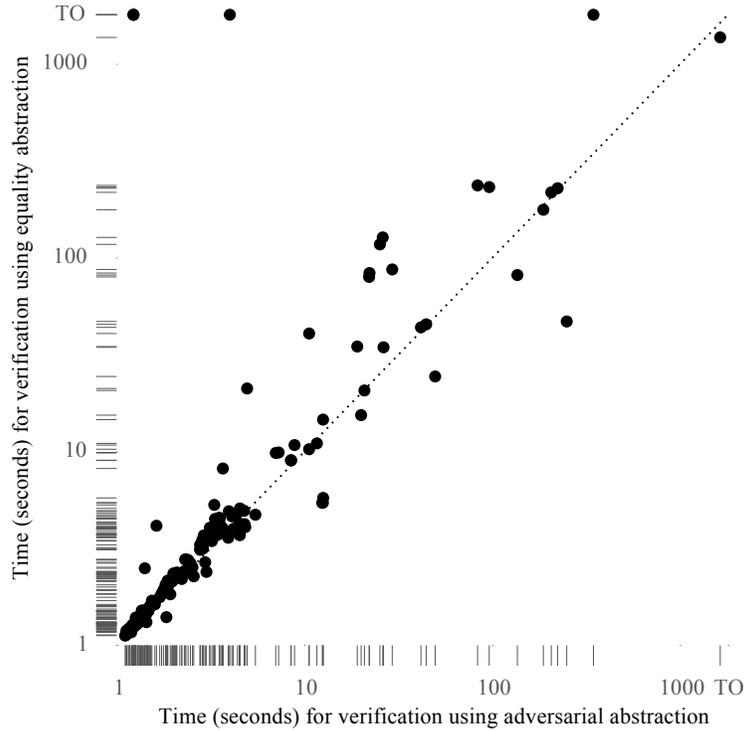


Fig. 12. Scatter plot comparing verification times using the adversarial and equality abstractions. Benchmarks for which verification using the adversarial abstraction fails (7 cases) are not included in the comparison.

ysis mentioned in Section 5.2 which GPUVerify employs to avoid use of the equality abstraction when it does not appear necessary.

We now evaluate both the effectiveness of GPUVerify in automatically deciding when the equality abstraction is necessary, and our hypothesis that it is *desirable* to avoid using the equality abstraction where possible.

To this end, we ran GPUVerify across the benchmark set in two modes, one where the adversarial abstraction is always used for every array, the other where the equality abstraction is always used for every array, adding to the data gathered in Section 6.2.1 where GPUVerify used its default heuristic to choose which abstraction to employ on an array-by-array basis.

*Precision of Taint Analysis for Abstraction Choice.* We found that for 155 out of 162 kernels, verification of full data race- and barrier divergence-freedom succeeded even when the adversarial abstraction was universally forced. In 130 of these cases, our taint analysis also led to the adversarial abstraction being used for all arrays, but in 25 cases GPUVerify decided, unnecessarily, that at least one array should be modeled using the equality abstraction. This is due to the conservative nature of our taint analysis which is neither flow- nor path-sensitive; a more precise taint analysis could allow more aggressive application of the adversarial abstraction.

Among the 7 kernels that required equality abstraction to be used for at least one array, our taint analysis succeeded in 3 cases to apply the equality abstraction sufficiently. In 4 cases, however, our taint analysis failed to choose the equality abstraction

for certain arrays, leading to false positive data race reports. These kernels were: the two binary search kernels discussed in Section 6.2, which depend on properties of the input array that the adversarial abstraction loses; the kernel implementing a pass of the Floyd-Warshall algorithm (discussed in relation to overflow in Section 6.2), for which race-freedom hinges on threads agreeing on the contents of a shared array; a parallel scan kernel exhibiting a benign data race which GPUVerify is only able to tolerate with the knowledge that a shared location is regarded as equal between threads. (In these cases, we used an option to override the taint analysis, using the equality abstraction for all arrays, when gathering data for the results in Section 6.2.1.)

Our evaluation shows that for our benchmark set, GPUVerify applies the equality abstraction relatively judiciously, capturing some but not all of the cases where it is useful.

*Performance Comparison Between Adversarial and Equality Abstractions.* For the 155 kernels for which adversarial abstraction suffices for verification, we investigated our hypothesis that using the equality abstraction would lead to slower verification. We did this by comparing verification times for the extreme cases where either the adversarial abstraction is forced for every array, or the equality abstraction is forced for every array. (By default, GPUVerify chooses which of these abstractions to apply on an array-by-array basis.) Note that the equality abstraction is strictly finer than the adversarial abstraction, so that if a kernel can be verified using the adversarial abstraction for all arrays, it must also be possible to verify the kernel (given sufficient resources) using the equality abstraction for all arrays.

The scatter plot of Figure 12 uses log scales to plot the time in seconds taken for verification using the adversarial abstraction ( $x$ -axis) and equality abstraction ( $y$ -axis). Each point is a benchmark, thus a point at coordinates  $(x, y)$  indicates that verification took  $x$  seconds using the adversarial abstraction and  $y$  seconds using the equality abstraction. Points with  $y$  coordinates at the top of the graph indicate that the timeout of 30 minutes was reached when using the equality abstraction. Points lying above the diagonal support our hypothesis that using the equality abstraction incurs a performance penalty in terms of verification time.

Our hypothesis is mainly supported: ignoring the conglomeration of benchmarks clustered at the bottom left of Figure 12 that verify quickly with either abstraction, most benchmarks verify more efficiently when the adversarial abstraction is used. There are four cases where applying the equality abstraction leads to our timeout of 30 minutes being reached.

There are, however, a number of cases where we found verification to be *faster* when the equality abstraction is used. There are five outliers in Figure 12 for which verification was at least two times faster using the equality abstraction than with adversarial abstraction; in the worst case verification is  $5.3\times$  slower. As in our discussion of outliers when comparing full vs. only intra-group race checking in Section 6.2.2, we attribute this to quirks of Z3. Again, we tried the same experiment using CVC4, obtaining a rather different set of outliers which, although larger, included just two of the five outliers observed with respect to Z3. With sufficient processing power we could guard against abstraction and solver quirks through portfolio solving, running multiple verification instances with different solvers in parallel.

### 6.3. Detection of a Bug in a Previous CUDA SDK Example

Using GPUVerify we discovered a write-write data race in the  $N$ -body example that shipped with version 2.3 of the CUDA SDK. This example uses multiple CUDA kernels to numerically approximate a system of  $N$  interacting bodies [Nyland et al. 2007]. This is an ideal problem for parallelization since interactions between each pair of bodies

```

shared int A[N];

void kernel() {
  for(int d = N/2; d > 0; d = d / 2) {
    if(lid < d) {
      A[lid] += A[lid + d];
    }
    barrier();
  }
}

```

Fig. 13. A reduction kernel used to illustrate the effectiveness of the two-thread reduction

can be calculated independently. The CUDA implementation of this example decomposes the  $N^2$  pair-interactions into smaller  $k \times k$  tiles, each of which is assigned to a one-dimensional group of  $k$  threads. Within each group, every thread is assigned to a distinct body (a row of the tile) and sequentially considers the interactions associated with this body to compute an updated state for the body. The kernel implements an optimization for small values of  $N$  where threads are arranged in two-dimensional groups, and multiple threads within a group are assigned to the same body. Consequently, the interactions calculated by threads assigned to the same body must be summed. A barrier ensures that each thread has completed its sub-calculation, and then a conditional is used to ensure that a single ‘master’ thread performs the summation. However, a data race could occur because a similar condition was not in place to ensure that only this master thread would perform a final update to the position and velocity of the body. As a result, it was possible for the final update by the master thread, using the full summation, to be overwritten by partial results computed by other threads.

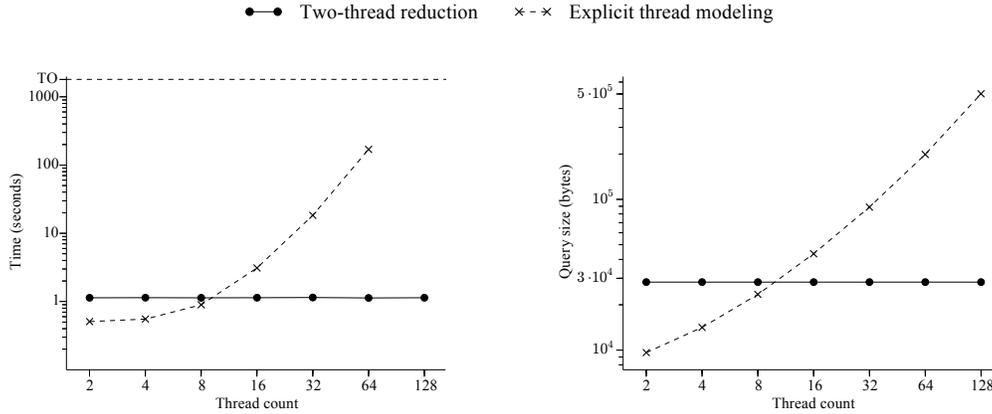
We reported this data race to Lars Nyland [2012] at Nvidia who confirmed that “It was a real bug, and it caused real issues in the results. It took significant debugging time to find the problem.” Nvidia had subsequently fixed this bug in version 3.0 of the CUDA SDK.

#### 6.4. A Demonstration of the Effectiveness of the Two-Thread Reduction

At the start of Section 4 we remarked that reasoning about GPU kernels using the SDV semantics of Section 3 directly, without reducing the number of threads that are modeled, does not scale. To justify this remark we consider the kernel of Figure 13 which computes a reduction over an input array  $A$ . The kernel is designed to be executed by a single group of  $N$  threads, where  $N$  is a power of two, with `lid` denoting the local id of a thread.

We wrote a Python script that takes as input an integer  $N$  and outputs a Boogie program representing the kernel of Figure 13 modeled using SDV semantics, executed explicitly by  $N$  threads (i.e., *without* applying the two-thread reduction). For successively large values of  $N$ , we compared the performance of Boogie on these generated programs vs. the performance of GPUVerify on the original kernel; GPUVerify *does* apply the two-thread reduction, and uses Boogie in its back-end. To ensure a fair comparison we disabled invariant inference in GPUVerify and annotated both the generated Boogie programs and the kernel analyzed by GPUVerify with invariants necessary for verification. We used a 30 minute timeout for each verification task.

Figure 14(a) compares verification times. The  $x$  axis plots the number of threads; the  $y$  axis plots verification time in seconds. Both axes use a log scale. The solid line shows



(a) As the thread count increases, verification time using GPUVerify (solid line) is virtually unaffected; this is due to the two-thread reduction. In contrast, verification time increases dramatically according to thread count when threads are explicitly modeled (dashed line)

(b) The timing results of Figure 14(a) can be linked to the SMT queries associated with each verification task; these queries grow quadratically with thread count when explicit modeling is used, but remain constant when the two-thread reduction is employed

Fig. 14. Comparing verification time and SMT query size when Boogie programs associated with the example kernel of Figure 13 are analyzed, with and without the two-thread reduction, for increasing thread counts

that verification with GPUVerify, using the two-thread reduction, remains virtually constant. On the other hand, the dashed curve shows that verification time grows rapidly (roughly exponentially) with explicit modeling as the thread count increases, exceeding our timeout with 128 threads.

This phenomenon is explained by Figure 14(b) in which the  $x$  axis again plots the number of threads and the  $y$  axis plots the size in bytes of the SMT query generated by Boogie for each verification problem. The formulas generated by GPUVerify (represented by the solid line) have identical size; they differ only in the constant that records thread count, which is a fixed-size bit-vector. In contrast, the formulas associated with explicit modeling (represented by the dashed line) become larger as the thread count increases. The growth is quadratic: a quadratic number of constraints are required to ensure that all the threads under consideration have distinct ids. This growth in formula size places an increasing burden on the SMT solver.

## 7. RELATED WORK

There are numerous existing dynamic and static techniques for data-race detection in programs that use lock-based synchronization or fork-join parallelism; a full discussion of these techniques is beyond the scope of this paper. We note however that this paper is concerned with proving data race- and barrier divergence-freedom in data-parallel programs in which the primary challenges—barrier synchronization and disjoint access patterns based on clever array indexing—are different from those encountered in lock-based and fork-join programs. In the rest of this section, we discuss papers that explicitly handle data-parallel or GPU programs, and the relationship between the verification method employed by GPUVerify and comparable proof techniques used in

protocol verification. We conclude the section with a brief discussion of invariant generation techniques.

*Static Verification of GPU Kernels.* The closest work to GPUVerify is the PUG analyzer for CUDA kernels by Li and Gopalakrishnan [2010]. Although GPUVerify and PUG have a similar goal, scalable verification of GPU kernels, the internal architecture of the two systems is very different. GPUVerify first translates a kernel into a sequential Boogie program that models the lock-step execution of two threads; the correctness of this program implies data race- and barrier divergence-freedom of the original kernel. Next, it infers and uses invariants to prove the correctness of this sequential program. Therefore, we only need to argue soundness for the translation into a sequential program; the soundness of the verification of the sequential program follows directly from the soundness of contract-based verification. On the other hand, PUG performs invariant inference while simultaneously translating the GPU kernel into a logical formula. PUG provides a set of built-in loop summarization rules which replace loops exhibiting certain shared array access patterns with corresponding invariants. Unlike GPUVerify, which must prove or discard all invariants that it generates, the loop invariants inserted by PUG are *assumed* to be correct. While this approach works for simple loop patterns, it has difficulty scaling to general nested loops in a sound way, and this necessitates various restrictions on the input program required by PUG. In contrast, GPUVerify inherits flexible and sound invariant inference from Houdini, regardless of the complexity of the control structure of the GPU kernel.

A recent method for functional verification of GPU kernels by Blom et al. [2014] uses permissions-based separation logic. Here, a kernel must be annotated (currently manually) with an assignment of read or write permissions to memory locations on a per-thread basis. Race-freedom is proven by showing that having permission to write a location excludes other threads from having permission to read or write that location. To reason about communication at barriers, a *barrier specification* is required, to state how permissions are exchanged between threads at a barrier. The approach of Blom et al. [2014] does not employ any thread-reduction abstraction; instead, quantifiers are used to reason about all threads. This method has not yet been automated, so a systematic comparison with GPUVerify for realistic examples is not yet possible.

In related work we have extended the GPU verification method with *barrier invariants* [Chong et al. 2013], which go beyond the equality abstraction described in this work, allowing more sophisticated invariants about the shared state to be established and used in verification. While the emphasis of this paper is on highly automatic verification of lightweight correctness properties, the focus of Chong et al. [2013] is on manual derivation of intricate barrier invariants to establish richer functional properties which can be used in race analysis of complex data-dependent kernels. We have also studied the generalization of lock-step execution from structured programs (as presented here) to unstructured control flow graphs [Collingbourne et al. 2013], and the current implementation of GPUVerify is based on this extension.

*Formal Semantics for GPU Kernels.* Habermaier and Knapp [2012] study the relationship between the single instruction multiple thread (SIMT) execution model of Nvidia GPUs and the standard interleaved semantics of threaded programs, presenting a formal semantics for predicated execution. This semantics shares similarities with the SDV semantics we present in Section 3, but the focus of Habermaier and Knapp [2012] is not on verification of GPU kernels. Recent work by Kojima and Igarashi [2013] provides a formulation of Hoare logic geared towards the SIMT model, in a similar spirit to Habermaier and Knapp [2012].

*Symbolic Execution and Bounded-Depth Verification.* The GKLEE tool by Li et al. [2012b] and KLEE-CL by Collingbourne et al. [2012] perform dynamic symbolic execution of CUDA and OpenCL kernels, respectively, and are both built on top of the KLEE symbolic execution engine [Cadar et al. 2008]. A method for bounded verification of barrier-free GPU kernels via depth-limited unrolling to an SMT formula is discussed by Tripakis et al. [2010]; lack of support for barriers, present in most non-trivial GPU kernels, limits the scope of this method. Symbolic execution and bounded unrolling techniques can be useful for bug-finding—both GKLEE and KLEE-CL have uncovered data race bugs in real-world examples—and these techniques have the advantage of generating concrete bug-inducing tests. The major drawback to these methods is that they cannot verify freedom from defects for non-trivial kernels.

The GKLEE tool specifically targets CUDA kernels, and faithfully models lock-step execution of sub-groups of threads, or *warps* as they are referred to in CUDA (see Table II). This allows precise checking of CUDA kernels that deliberately exploit the warp size of an Nvidia GPU to achieve high performance. By default GPUVerify makes no assumptions about sub-group size, making it useful for checking whether CUDA kernels are portable; GPUVerify includes an option to enable warp-level reasoning [Bardsley and Donaldson 2014].

Both GKLEE and KLEE-CL explicitly represent the number of threads executing a GPU kernel. This allows for precise defect checking, but limits scalability. A recent extension to GKLEE by Li et al. [2012a] uses the notion of *parametric flows* to soundly restrict defect checking to consider only certain pairs of threads. This is similar to the two-thread abstraction employed by GPUVerify and PUG, and leads to scalability improvements over standard GKLEE, at the expense of a loss in precision for kernels that exhibit inter-thread communication (GPUVerify and PUG also suffer from this loss in precision, as discussed under *Incompleteness of the Two-Thread Reduction* in Section 4.1).

The GKLEE tool has also been extended with support for reasoning about atomic operations [Chiang et al. 2013]. Atomic operations are challenging because they allow race-free non-determinism between pairs of barriers, destroying the ability to reason about a single canonical thread schedule, an assumption on which PUG, GPUVerify, GKLEE and KLEE-CL all rest. GKLEE supports bounded reasoning about atomics, geared towards bug-finding, through the use of *delay bounding* [Emmi et al. 2011]. GPUVerify includes prototype support for reasoning about kernels that manipulate data using atomic operations, through a refinement to the adversarial abstraction [Bardsley and Donaldson 2014]. Neither GPUVerify nor GKLEE provides the necessary support for sound reasoning about fine-grained concurrency in GPU kernels which, as demonstrated in recent work, can expose the relaxed memory consistency models of modern GPU architectures [Alglave et al. 2015].

*Dynamic Analysis.* Dynamic analysis of CUDA kernels for data race detection has been proposed by Boyer et al. [2008]. Leung et al. [2012] report on the following technique that combines dynamic and static data race analysis. They first simulate a CUDA kernel while dynamically checking for races. If no races are detected, flow analysis is used to determine whether the control-flow taken during dynamic execution was dependent on input data; if not, the kernel can be deemed race free, otherwise the technique is inconclusive. It appears that this approach can handle kernels that are verifiable using our *adversarial* abstraction. Kernels which GPUVerify can verify only with the *equality* abstraction, due to threads testing input data, are not amenable to analysis using this technique. As discussed in Section 6, the tool associated with [Leung et al. 2012] is not publicly available, hence we do not present an experimental comparison between this work and GPUVerify.

A problem associated with static analysis can be the lack of adequate preconditions for reasoning about the correctness of a procedure. A pragmatic method for overcoming this problem when analyzing OpenCL kernels with GPUVerify is *kernel interception*, where dynamic analysis is used to capture the values of scalar parameters supplied to a kernel at runtime, as well as the specific configuration of threads and groups with respect to which the kernel executes [Bardsley et al. 2014b]. This dynamically collected information provides a very specific precondition under which a kernel can be analyzed. This has the property of lowering the rate of false positives reported by GPUVerify, at the expense of making the verification result specific to the particular runtime configuration that was captured.

*Thread Contracts.* Karmani et al. [2011] present a practical approach for constructing correct parallel programs, based on *thread contracts*. A programmer specifies the coordination and data sharing strategy for their multi-threaded program as a contract. Adequacy of the specification for ensuring race-freedom is then checked statically, while adherence to the specification by the implementation is ascertained via testing. Adapted to the setting of barrier synchronization rather than lock-based coordination, this technique might enable analysis of more complex GPU kernels for which automatic contract inference is infeasible.

*Protocol Verification.* A reduction to two processes, similar to the two-thread reduction employed in GPU kernel verification, is at the heart of a method for verifying cache coherence protocols known as CMP [Chou et al. 2004], which was inspired by the foundational work of McMillan [1999]. With CMP, verification of a protocol for an arbitrary number of processes is performed by model checking a system where a small number of processes are explicitly represented and a highly non-deterministic ‘other’ process over-approximates the possible behaviors of the remaining processes. The unconstrained nature of the ‘other’ process can lead to spurious counterexamples, which must be eliminated either by introducing additional explicit processes, or by adding *non-interference lemmas* such that the actions of the ‘other’ process more precisely reflect the possible actions of processes in the concrete system. The CMP method has been extended and generalized with *message flows* and *message flow invariants* by Talupur and Tuttle [2008]. This extension aids in the automatic derivation of non-interference lemmas by capturing large classes of permissible interactions between processes. Our approach, and that of PUG by Li and Gopalakrishnan [2010], uses the same high-level proof idea as the CMP method: we consider a small number of threads (two), and our default adversarial abstraction models the possible actions of all other threads, analogously to the ‘other’ process. There are, however, many technical differences in the manner by which the high-level proof technique is applied in practice. Principally, we use predicated execution to turn the verification problem into a sequential program analysis task, completely eliminating the need to reason about concurrent threads, while the CMP method involves directly checking the parallel composition of a number of processes using model checking.

*Invariant Generation.* As described in Section 5.3, we use the Houdini algorithm of Flanagan and Leino [2001] to generate loop invariants for verification. Houdini was introduced as an annotation assistant for the Java Extended Static Checker [Leino et al. 2000]. Related template-based invariant generation techniques include [Kahsai et al. 2011; Srivastava and Gulwani 2009; Lahiri and Qadeer 2009]. As discussed in the related work section of Flanagan and Leino [2001], Houdini can be viewed under the framework of abstract interpretation [Cousot and Cousot 1977] where the abstract domain comprises conjunctions of predicates drawn from the set of candidate invariants. Compared with standard predicate abstraction [Graf and Saïdi 1997], which consid-

ers arbitrary Boolean combinations of predicates (and is thus more precise), verification using Houdini requires a linear instead of exponential number of theorem prover calls. In our context, the key advantage of the Houdini approach over traditional abstract interpretation using a fixed abstract domain is *flexibility*. We can easily extend GPUVerify with a richer language of predicates by adding further candidate invariant generation rules; there is no need for careful redesign of an abstract domain.

The main problem with our invariant generation method is that its success is limited by the scope of our candidate invariant generation rules. Interpolation and counterexample-guided abstraction refinement can be used incrementally to generate invariants in response to failed or partial verification attempts [McMillan 2006; Beyer et al. 2007], while the Daikon technique [Ernst et al. 2007] allows program-specific invariants to be speculated through dynamic analysis. It may be possible to draw upon these techniques to improve the invariant inference capabilities of GPUVerify.

## 8. CONCLUSIONS

We have presented the design and implementation of the GPUVerify method for statically verifying data race- and barrier divergence-freedom of GPU kernels. The foundations of the method are our novel SDV semantics and the two-thread reduction which allows scalable verification to be built on the semantic foundations. We have presented a rigorous soundness proof for the two-thread reduction, partially mechanized using Isabelle, and have discussed practical issues associated with implementing the GPUVerify tool. Our large experimental evaluation demonstrates that GPUVerify is effective in verifying and falsifying real-world OpenCL and CUDA GPU kernels.

Our current and planned future activities are mainly focused on improving the automatic invariant inference facilities of GPUVerify to lower the rate of false positives reported by the tool; we believe this is key for industrial uptake of the method.

## ACKNOWLEDGMENTS

Our thanks to Guodong Li and Ganesh Gopalakrishnan for providing us with the latest version of PUG, and for answering numerous questions about the workings of this tool.

Thanks to Matko Botinčan, Mike Dodds, Hristina Palikareva and the anonymous reviewers of the conference version of the paper for their insightful feedback on an earlier draft of this work.

We are also grateful to the anonymous TOPLAS reviewers for their useful feedback on the extended version.

We thank several people working in the GPU industry who provided assistance with this work: Anton Lokhmotov (ARM) provided us with barrier divergence results for ARM's Mali architecture; Lee Howes (AMD) provided such results for AMD's Tahiti architecture; Teemu Uotila and Teemu Virolainen (Rightware) provided access to Basemark CL, and answered various queries about these kernels; Yossi Levanoni (Microsoft) provided early access to C++ AMP samples.

We are grateful to Peter Collingbourne for leading development of Bugle, an LLVM IR-based front-end for GPUVerify which replaced the Clang AST-based front-end described in the conference version of this paper [Betts et al. 2012].

## References

- Jade Alglave, Mark Batty, Alastair F. Donaldson, Ganesh Gopalakrishnan, Jeroen Ketema, Daniel Poetzl, Tyler Sorensen, and John Wickerson. 2015. GPU Concurrency: Weak Behaviours and Programming Assumptions. In *Proceedings of the 20th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'15)*. ACM, New York, 577–591.
- AMD. 2013. OpenCL Programming Guide, Revision 2.7. (Nov. 2013).
- Ethel Bardsley, Adam Betts, Nathan Chong, Peter Collingbourne, Pantazis Deligiannis, Alastair F. Donaldson, Jeroen Ketema, Daniel Liew, and Shaz Qadeer. 2014a. Engineering a Static Verification Tool for GPU Kernels. In *Proceedings of the 26th International Conference on Computer Aided Verification (CAV 2014) (Lecture Notes in Computer Science)*, Vol. 8559. Springer, Berlin, 226–242.

- Ethel Bardsley and Alastair F. Donaldson. 2014. Warps and Atomics: Beyond Barrier Synchronization in the Verification of GPU Kernels. In *Proceedings of the 6th NASA Formal Methods Symposium (NFM 2014) (Lecture Notes in Computer Science)*, Vol. 8430. Springer, Berlin, 230–245.
- Ethel Bardsley, Alastair F. Donaldson, and John Wickerson. 2014b. KernelInterceptor: Automating GPU Kernel Verification by Intercepting Kernels and Their Parameters. In *Proceedings of the 2014 International Workshop on OpenCL (IWOCCL 2014)*. ACM, New York, Article 7, 5 pages.
- Michael Barnett, Bor-Yuh Evan Chang, Robert DeLine, Bart Jacobs, and K. Rustan M. Leino. 2005. Boogie: A Modular Reusable Verifier for Object-Oriented Programs. In *Revised Lectures of the 4th International Symposium on Formal Methods for Components and Objects (FMCO 2005) (Lecture Notes in Computer Science)*, Vol. 4111. Springer, Berlin, 364–387.
- Clark Barrett, Christopher L. Conway, Morgan Deters, Liana Hadarean, Dejan Jovanovic, Tim King, Andrew Reynolds, and Cesare Tinelli. 2011. CVC4. In *Proceedings of the 23rd International Conference in Computer Aided Verification (CAV 2011) (Lecture Notes in Computer Science)*, Vol. 6806. Springer, Berlin, 171–177.
- Al Bessey, Ken Block, Benjamin Chelf, Andy Chou, Bryan Fulton, Seth Hallem, Charles-Henri Gros, Asya Kamsky, Scott McPeak, and Dawson R. Engler. 2010. A Few Billion Lines of Code Later: Using Static Analysis to Find Bugs in the Real World. *Commun. ACM* 53, 2 (2010), 66–75.
- Adam Betts, Nathan Chong, Alastair F. Donaldson, Shaz Qadeer, and Paul Thomson. 2012. GPUVerify: A Verifier for GPU Kernels. In *Proceedings of the 27th ACM International Conference on Object Oriented Programming, Systems, Languages, and Applications (OOPSLA 2012)*. ACM, New York, 113–132.
- Dirk Beyer, Thomas A. Henzinger, Rupak Majumdar, and Andrey Rybalchenko. 2007. Path Invariants. In *Proceedings of the 2007 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2007)*. ACM, New York, 300–309.
- Stefan Blom, Marieke Huisman, and Matej Mihelčić. 2014. Specification and Verification of GPGPU Programs. *Science of Computer Programming* 95, 3 (2014), 376–388.
- Michael Boyer, Kevin Skadron, and Westley Weimer. 2008. Automated Dynamic Analysis of CUDA Programs. In *Proceedings of the 3rd Workshop on Software Tools for MultiCore Systems (STMCS 2008)*. Online proceedings, <http://people.csail.mit.edu/rabbah/conferences/08/cgo/stmcs/>, 6 pages.
- Cristian Cadar, Daniel Dunbar, and Dawson R. Engler. 2008. KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs. In *Proceedings of the 8th USENIX Symposium on Operating Systems Design and Implementation (OSDI 2008)*. USENIX Association, Berkeley, CA, 209–224.
- Joshua E. Cates, Aaron E. Lefohn, and Ross T. Whitaker. 2004. GIST: An Interactive, GPU-Based Level Set Segmentation Tool for 3D Medical Images. *Medical Image Analysis* 8 (2004), 217–231. Issue 3.
- Wei-Fan Chiang, Ganesh Gopalakrishnan, Guodong Li, and Zvonimir Rakamaric. 2013. Formal Analysis of GPU Programs with Atomics via Conflict-Directed Delay-Bounding. In *Proceedings of 5th International NASA Formal Methods Symposium (NFM 2013) (Lecture Notes in Computer Science)*, Vol. 7871. Springer, Berlin, 213–228.
- Nathan Chong, Alastair F. Donaldson, Paul Kelly, Shaz Qadeer, and Jeroen Ketema. 2013. Barrier Invariants: A Shared State Abstraction for the Analysis of Data-Dependent GPU Kernels. In *Proceedings of the 28th ACM International Conference on Object Oriented Programming, Systems, Languages, and Applications (OOPSLA 2013)*. ACM, New York, 605–622.
- Nathan Chong, Alastair F. Donaldson, and Jeroen Ketema. 2014. A Sound and Complete Abstraction for Reasoning About Parallel Prefix Sums. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL 2014)*. ACM, New York, 397–410.
- Ching-Tsun Chou, Phanindra K. Mannava, and Seungjoon Park. 2004. A Simple Method for Parameterized Verification of Cache Coherence Protocols. In *Proceedings of the 5th International Conference on Formal Methods in Computer-Aided Design (FMCAD 2004) (Lecture Notes in Computer Science)*, Vol. 3312. Springer, Berlin, 382–398.
- Peter Collingbourne, Cristian Cadar, and Paul H. J. Kelly. 2012. Symbolic Testing of OpenCL Code. In *Revised Selected Papers of the 7th International Haifa Verification Conference (HVC 2011) (Lecture Notes in Computer Science)*, Vol. 7261. Springer, Berlin, 203–218.
- Peter Collingbourne, Alastair F. Donaldson, Jeroen Ketema, and Shaz Qadeer. 2013. Interleaving and Lock-Step Semantics for Analysis and Verification of GPU Kernels. In *Proceedings of the 22nd European Symposium on Programming (ESOP 2013) (Lecture Notes in Computer Science)*, Vol. 7792. Springer, Berlin, 270–289.
- Patrick Cousot and Radhia Cousot. 1977. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *Proceedings of the 4th ACM SIGACT-*

- SIGPLAN Symposium on Principles of Programming Languages (POPL 1977)*. ACM, New York, 238–252.
- Leonardo Mendonça de Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In *Proceedings of the 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS 2008) (Lecture Notes in Computer Science)*, Vol. 4963. Springer, Berlin, 337–340.
- Alastair F. Donaldson. 2014. The GPUVerify Method: a Tutorial Overview. *Electronic Communications of the EASST 70*, Article 1 (2014), 16 pages. <http://journal.ub.tu-berlin.de/eceasst/article/view/986>
- Alastair F. Donaldson, Daniel Kroening, and Philipp Rümmer. 2010. Automatic Analysis of Scratch-Pad Memory Code for Heterogeneous Multicore Processors. In *Proceedings of the 16th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS 2010) (Lecture Notes in Computer Science)*, Vol. 6015. Springer, Berlin, 280–295.
- Alastair F. Donaldson, Daniel Kroening, and Philipp Rümmer. 2011. Automatic Analysis of DMA Races Using Model Checking and  $k$ -Induction. *Formal Methods in System Design* 39, 1 (2011), 83–113.
- Michael Emmi, Shaz Qadeer, and Zvonimir Rakamaric. 2011. Delay-Bounded Scheduling. In *Proceedings of the 39th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL 2011)*. ACM, New York, 411–422.
- Michael D. Ernst, Jeff H. Perkins, Philip J. Guo, Stephen McCamant, Carlos Pacheco, Matthew S. Tschantz, and Chen Xiao. 2007. The Daikon System for Dynamic Detection of Likely Invariants. *Science of Computer Programming* 69, 1–3 (2007), 35–45.
- Cormac Flanagan and K. Rustan M. Leino. 2001. Houdini, an Annotation Assistant for ESC/Java. In *Proceedings of the International Symposium of Formal Methods Europe (FME 2001) (Lecture Notes in Computer Science)*, Vol. 2021. Springer, Berlin, 500–517.
- Susanne Graf and Hassen Saïdi. 1997. Construction of Abstract State Graphs with PVS. In *Proceedings of the 9th International Conference on Computer Aided Verification (CAV 1997) (Lecture Notes in Computer Science)*, Vol. 1254. Springer, Berlin, 72–83.
- Axel Haberman and Alexander Knapp. 2012. On the Correctness of the SIMT Execution Model of GPUs. In *Proceedings of the 21st European Symposium on Programming (ESOP 2012) (Lecture Notes in Computer Science)*, Vol. 7211. Springer, Berlin, 316–335.
- Mark J. Harris. 2004. Fast Fluid Dynamics Simulation on the GPU. In *GPU Gems*. Addison-Wesley, Boston, MA, Chapter 38, 637–665.
- Temesghen Kahsai, Yeting Ge, and Cesare Tinelli. 2011. Instantiation-Based Invariant Discovery. In *Proceedings of 3rd International NASA Formal Methods Symposium (NFM 2011) (Lecture Notes in Computer Science)*, Vol. 6617. Springer, Berlin, 192–206.
- Rajesh K. Karmani, P. Madhusudan, and Brandon M. Moore. 2011. Thread Contracts for Safe Parallelism. In *Proceedings of the 16th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPOPP 2011)*. ACM, New York, 125–134.
- Khronos OpenCL Working Group. 2012. The OpenCL Specification, Version 1.2. (2012). Document revision 19.
- Khronos OpenCL Working Group. 2014a. The OpenCL C Specification, Version 2.0. (2014). Document revision 26.
- Khronos OpenCL Working Group. 2014b. The OpenCL Extension Specification, Version 2.0. (2014). Document revision 26.
- Petr Klus, Simon Lam, Dag Lyberg, Ming Sin Cheung, Graham Pullan, Ian McFarlane, Giles S. H. Yeo, and Brian Y. H. Lam. 2012. BarraCUDA – A Fast Short Read Sequence Aligner Using Graphics Processing Units. *BMC Research Notes* 5, Article 27 (2012), 7 pages.
- Kensuke Kojima and Atsushi Igarashi. 2013. A Hoare Logic for SIMT Programs. In *Proceeding of the 11th Asian Symposium on Programming Languages and Systems (APLAS 2013) (Lecture Notes in Computer Science)*, Vol. 8301. Springer, Berlin, 58–73.
- Shuvendu K. Lahiri and Shaz Qadeer. 2009. Complexity and Algorithms for Monomial and Clausal Predicate Abstraction. In *Proceedings of the 22nd International Conference on Automated Deduction (CADE-22) (Lecture Notes in Computer Science)*, Vol. 5663. Springer, Berlin, 214–229.
- K. Rustan M. Leino, Greg Nelson, and James B. Saxe. 2000. *ESC/Java User's Manual*. Technical Report SRC Technical Note 2000-002. Compaq Systems Research Center.
- Alan Leung, Manish Gupta, Yuvraj Agarwal, Rajesh Gupta, Ranjit Jhala, and Sorin Lerner. 2012. Verifying GPU kernels by test amplification. In *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'12)*. ACM, New York, 383–394.

- Guodong Li and Ganesh Gopalakrishnan. 2010. Scalable SMT-Based Verification of GPU Kernel Functions. In *Proceedings of the 18th ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2010)*. ACM, New York, 187–196.
- Guodong Li, Peng Li, Geoffrey Sawaya, Ganesh Gopalakrishnan, Indradeep Ghosh, and Sreeranga P. Rajan. 2012b. GKLEE: concolic verification and test generation for GPUs. In *Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPOPP 2012)*. ACM, New York, 215–224.
- Heng Li and Richard Durbin. 2009. Fast and Accurate Short Read Alignment with Burrows-Wheeler Transform. *Bioinformatics* 25, 14 (2009), 1754–1760.
- Peng Li, Guodong Li, and Ganesh Gopalakrishnan. 2012a. Parametric Flows: Automated Behavior Equivalencing for Symbolic Analysis of Races in CUDA Programs. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC 2012)*. IEEE Computer Society Press, Los Alamitos, CA, 10 pages.
- Anton Lokhmotov. 2011. Mobile and Embedded Computing on Mali GPUs. (2011). 2nd UK GPU Computing Conference.
- Kenneth McMillan. 1999. Verification of Infinite State Systems by Compositional Model Checking. In *Proceedings of the 10th Conference on Correct Hardware Design and Verification Methods (CHARME 1999) (Lecture Notes in Computer Science)*, Vol. 1703. Springer, Berlin, 219–234.
- Kenneth L. McMillan. 2006. Lazy Abstraction with Interpolants. In *Proceedings of the 18th International Conference on Computer Aided Verification (CAV 2006) (Lecture Notes in Computer Science)*, Vol. 4144. Springer, Berlin, 123–136.
- Microsoft Corporation. 2012. C++ AMP: Language and Programming Model, Version 1.0. (2012).
- Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. 2002. *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. Lecture Notes in Computer Science, Vol. 2283. Springer, Berlin.
- Nvidia. 2012a. CUDA C Programming Guide, Version 5.0. (2012).
- Nvidia. 2012b. Parallel Thread Execution ISA, Version 3.1. (2012).
- Lars Nyland. 2012. Personal Communication. (April 2012).
- Lars Nyland, Mark Harris, and Jan Prins. 2007. *Fast N-Body Simulation with CUDA*. Addison-Wesley, Upper Saddle River, NJ, Chapter 31, 677–696.
- Renato F. Salas-Moreno, Richard A. Newcombe, Hauke Strasdat, Paul H. J. Kelly, and Andrew J. Davison. 2013. SLAM++: Simultaneous Localisation and Mapping at the Level of Objects. In *Proceedings of the 2013 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2013)*. IEEE Computer Society Press, Los Alamitos, CA, 1352–1359.
- Saurabh Srivastava and Sumit Gulwani. 2009. Program Verification Using Templates over Predicate Abstraction. In *Proceedings of the 2009 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2009)*. ACM, New York, 223–234.
- Bjarne Steensgaard. 1996. Points-to Analysis in Almost Linear Time. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL 1996)*. ACM, New York, 32–41.
- Murali Talupur and Mark R. Tuttle. 2008. Going with the Flow: Parameterized Verification Using Message Flows. In *Proceedings of the 8th International Conference on Formal Methods in Computer-Aided Design (FMCAD 2008)*. IEEE Computer Society Press, Los Alamitos, CA, 8 pages.
- Stavros Tripakis, Christos Stergiou, and Roberto Lublinerman. 2010. Checking Equivalence of SPMD Programs Using Non-Interference. In *Proceedings of the 2nd USENIX Workshop on Hot Topics in Parallelism (HotPar 2010)*. Online proceedings, <http://static.usenix.org/events/hotpar10/>, 5 pages.
- Christian Urban and Julien Narboux. 2009. Formal SOS-Proofs for the Lambda-Calculus. *Electronic Notes in Theoretical Computer Science* 247 (2009), 139–155.
- John Wickerson. 2014. Syntax and Semantics of a GPU kernel Programming Language. *Archive of Formal Proofs* (2014). [http://afp.sourceforge.net/entries/GPU\\_Kernel\\_PL.shtml](http://afp.sourceforge.net/entries/GPU_Kernel_PL.shtml)

Received ?; revised ?; accepted ?