

1 Checking Data-Race Freedom of GPU Kernels, 2 Compositionally

3 Tiago Cogumbreiro¹, Julien Lange²,
4 Dennis Liew Zhen Rong¹, and Hannah Zicarelli¹

5 ¹ University of Massachusetts Boston
6 ² Royal Holloway, University of London

7 **Abstract.** GPUs offer parallelism as a commodity, but they are diffi-
8 cult to program correctly. Static analyzers that guarantee data-race free-
9 dom (DRF) are essential to help programmers establish the correctness
10 of their programs (kernels). However, existing approaches produce too
11 many false alarms and struggle to handle larger programs. To address
12 these limitations we formalize a novel compositional analysis for DRF,
13 based on access memory protocols. These protocols are behavioral types
14 that codify the way threads interact over shared memory.

15 Our work includes fully mechanized proofs of our theoretical results, the
16 first mechanized proofs in the field of DRF analysis for GPU kernels. Our
17 theory is implemented in `Faial` (anonymized name), a tool that outper-
18 forms the state-of-the-art. Notably, it can correctly verify at least $1.41 \times$
19 more real-world kernels, and it exhibits a linear growth in 4 out of 5
20 experiments, while others grow exponentially in all 5 experiments.

21 1 Introduction

22 GPUs are massively parallel devices that promise a great return on investment
23 at a cost: they are notably difficult to program. In GPU programming, hundreds
24 of lightweight threads share portions of arrays in parallel (without locks) —
25 very different from the programming model of multithreaded programs written
26 in C or Java with heavy-weight heterogeneous threads. Data-race freedom (DRF)
27 analysis aims to guarantee that for all possible executions, every array cell being
28 written by one thread cannot be concurrently accessed by another thread.

29 In the field of static analysis of DRF in GPU programs, there is a tension
30 between efficiency and correctness (no missed data-races and no false alarms)
31 that thus far is unresolved. Bug finding tools [25, 26, 33] favor correctness over
32 efficiency: they provide correct results at small scales, by simulating the program
33 execution. Such tools are incapable of handling certain parameters symbolically
34 (*e.g.*, array size) and can easily exhaust users’ resources (*e.g.*, loops with long
35 iteration spaces or unknown bounds). Approaches based on Hoare logic [5, 7, 21]
36 can cope with medium-sized programs, do not miss data-races, and do not require
37 array size information; however, they suffer from a high-rate of false alarms and
38 require code annotations written by concurrency experts. Finally, tools that can



Fig. 1: Work-flow of the verification.

39 cope with larger programs and do not require array size information either miss
40 data-races [23] or overwhelm the user with false alarms [38].

41 To appease this tension, we introduce a novel static DRF analysis that can
42 handle larger programs and produce fewer false alarms than related work, with-
43 out missing data-races. Additionally our analysis does not require code anno-
44 tations or array size information. Our verification framework hinges on *access*
45 *memory protocols*, a new family of behavioral types [1] that codify the way
46 threads interact through shared memory. Our behavioral types also make evi-
47 dent two aspects of the analysis that can be made separate: concurrency analysis
48 (*i.e.*, could these two expressions run in parallel?) and data-race conflict detec-
49 tion (*i.e.*, do these array indices match?).

50 *Contributions and synopsis* This paper includes the following contributions.

51 (1) In §3, we formalize the syntax, semantics, and well-formedness conditions
52 for access memory protocols, which are behavioral types for GPU programs.
53 This behavioral abstraction results in a simpler yet more expressive theory than
54 previous works, *e.g.*, it does not require user-provided loop invariants.

55 (2) In §4, we show that our DRF analysis of access memory protocols can be
56 soundly and completely reduced to the satisfiability of an SMT formula, see
57 Theorems 1 and 3. Our theory and results on access memory protocols are fully
58 mechanized in Coq. To the best of our knowledge, this is the first mechanized
59 proof of correctness of a DRF analysis for GPU programs.

60 (3) We show that our DRF analysis of access memory protocols is compositional
61 when protocols satisfy a structural property, see Theorem 2. Additionally, we
62 show how to transform protocols when they do not meet this property.

63 (4) In §5 we present *Faial*, which infers access memory protocols from CUDA
64 kernels and implements our theory. Our experiments show that *Faial* is more
65 precise and scales better than existing tools.

66 (5) In §6, we present a thorough experimental evaluation of *Faial* against related
67 work [5, 23, 25, 26], the largest comparative study of GPU verification (5 tools
68 in 260 kernels, 3 tools compared in 487 kernels). *Faial* verified 217 out of 227
69 real-world kernels (at least $1.42\times$ more than other tools) and correctly verified
70 more handcrafted tests than other tools (4 out of 5). In a synthetic benchmark
71 suite (250 kernels), *Faial* is the only tool to exhibit linear growth in 4 out of 5
72 experiments, while others grow exponentially in all 5 experiments.

73 Our paper is accompanied by an implementation (*Faial*, see § A), an evaluation
74 framework (inc. datasets), and proof scripts (in Coq) for each theorem. Should
75 the paper be accepted, these will be submitted for artifact evaluation.

Listing 2.1: Examples of racy kernels, l.h.s. is from [34] and r.h.s. simplifies l.h.s. for clarity, with one-dimensional array and thread identifier, and 1-stride loops.

```

1 for (int r = 0; r < N; r++) {
2   for (int i = 0; i < TILE_DIM; i += BLOCK_ROWS)
3     { tile[tid.y+i][tid.x] = idata[index_in+i*width]; }
4   syncthreads();
5   for (int j = 0; j < TILE_DIM; j += BLOCK_ROWS)
6     { odata[index_out+j*height] = tile[tid.x][tid.y+j]; }

```

```

1 for (int r = 0; r < N; r++) {
2   for (int i = 0; i < M; i++)
3     { tile[tid] = ...; }
4   syncthreads();
5   for (int j = 0; j < M; j++)
6     { ... = tile[tid+j]; }

```

76 2 Overview

77 This section gives an overview of our approach by examining a data-race we
78 found in published work [16] and [34]. We discuss the challenges that such ex-
79 amples pose to the state-of-the-art of DRF analysis. Then we introduce a veri-
80 fication framework based on *access memory protocols*: behavioral types [1] that
81 codify the way threads interact via shared memory. Figure 1 gives an overview
82 of the verification pipeline. We start from CUDA kernels, from which we infer
83 access memory protocols. Protocols are then checked for well-formedness and
84 transformed in three steps into formulas that are verified by an SMT solver.

85 2.1 Challenges of GPU Programming

86 **GPU programming model** The key component of GPU programming is the
87 *kernel* program, or just kernel, that runs according to the Single-Instruction-
88 Multiple-Thread (SIMT) execution model, where multiple threads run a single
89 instruction concurrently. A kernel is parameterized by a special variable that
90 holds a thread identifier, henceforth named *tid*. In parallel, each member of a
91 group of threads runs an instantiated copy of the kernel by supplying its identifier
92 as an argument. Threads communicate via shared memory (arrays) and mediate
93 communication via barrier synchronization (an execution point where all threads
94 must wait for each other before advancing further). Writes are only visible to
95 other threads after a barrier synchronization, *i.e.*, there is no guarantee that a
96 write of a thread can be read by another before a barrier synchronization.

97 GPU programming platforms usually group threads hierarchically in multi-
98 ple levels, across which no inter-groups synchronization is possible. In both the
99 literature [6, 23] and this work, the focus is on intra-group communication as
100 inter-group errors can be seen as a special case of intra-group errors.

101 **Challenges** We motivate the difficulty of analyzing data-races by studying a
102 programming error found in the wild, reported in Listing 2.1 (left). This excerpt
103 comes from a tutorial [34] on optimizing numeric algorithms for GPUs. The code
104 listing transposes a matrix *N*-times with an outer loop indexed by variable *r*.

105 Remarkably, the tutorial [34] does not inform the readers that Listing 2.1
106 contains a subtle *data-race*: one transpose-operation starts (the writes to tile
107 in line 3) without awaiting the termination of the previous transpose-operation

Listing 2.2: Minimal representative example of an access memory protocol highlighting the data-race in Listing 2.1.

```

1 // r = 0
2 forU j in 0..M // for (int j = 0; j<M; j++)
3   {rd[tid+j]}; // _ = tile [tid+i];
4 // r = 1
5 forU i in 0..M // for (int i = 0; i<M; i++)
6   {wr[tid]} // tile [tid] = _;
```

108 (the reads from `tile` in line 6), thus corrupting the data over time and possibly
 109 skewing the timing of the optimization to appear faster than it should be. We
 110 found a related data-race in [16], which reuses code from [34].

111 Our tool, `Faial`, successfully identifies the program state that triggers the
 112 data-race in Listing 2.1: when $r=1$ and $N=2$. However, state-of-the-art tools
 113 struggle to accurately analyze Listing 2.1, as evaluated in Section 6 (Claim 1:
 114 Test 1). Symbolic execution tools, such as [25, 26], timeout for $N > 1$, and, in
 115 general, cannot handle symbolic (unknown) bounds. `GPUVerify` [6], a tool based
 116 on Hoare logic, reports a false alarm: a spurious data-race when $r=0$ and $N=1$.
 117 And `PUG` [23] incorrectly identifies the example as DRF, as its analysis appears
 118 to ignore data-races originating from different iterations of a loop.

119 2.2 Memory Access Protocols by Example

120 We now investigate the data-race in Listing 2.1 with an access memory proto-
 121 col. For presentation purposes, we focus our discussion on Listing 2.1 (r.h.s.),
 122 that simplifies the l.h.s. whilst retaining the root cause of its data-race, which
 123 stems from the interaction between both loops. We discuss how we support
 124 multi-dimensional arrays, multi-dimensional thread identifiers, and arbitrary
 125 loop strides in Section 5. In our Coq formalism the notion of “accesses” (and
 126 their dimensions) is a parameter of the theory, thus orthogonal to the theory
 127 presented here.

128 Consider the execution of the end of the first iteration ($r=0$) and the beginning
 129 of the second ($r=1$) iteration of the outer-loop. In this case, the execution of the
 130 `j`-loop when $r=0$ is not synchronized with the execution of the `i`-loop when $r=1$ as
 131 there is no call to `__syncthreads()` in between.

132 The access memory protocol in Listing 2.2 captures this *partial* execution
 133 from the viewpoint of array `tile`. By design access memory protocols over ap-
 134 proximate kernels by abstracting away *what* data is being written to/read from
 135 an array, to focus on *where* data is being written. The protocol models the two
 136 problematic loops of Listing 2.1, *i.e.*, the `j`-loop when $r=0$ and the `i`-loop when $r=1$.
 137 The first loop reads (`rd[tid+j]`) from the array, while the second writes (`wr[tid]`)
 138 to it. Evaluation of a protocol follows the SIMT model: each thread evaluates
 139 `wr[tid]` by instantiating `tid` with their unique identifier (hereafter, an integer),
 140 *e.g.*, thread 0 yields `wr[0]` and thread 1 yields `wr[1]`.

141 **Analysis of unsynchronized protocols** We say that a protocol is DRF when
 142 all concurrent accesses are pair-wise DRF, *i.e.*, when issued by different threads
 143 on the same index, then neither access is a write. For instance the respective
 144 sets of concurrent accesses of threads 0 and 1 in Listing 2.2 is given below

$$\begin{array}{ccc} \text{tid} = 0 & & \text{tid} = 1 \\ \{\text{rd}[j] \mid 0 \leq j < M\} \cup \{\text{wr}[0]\} & \text{DRF with?} & \{\text{rd}[1+j] \mid 0 \leq j < M\} \cup \{\text{wr}[1]\} \end{array}$$

145 When $M > 1$, thread 0 (l.h.s) accesses `rd[1]` and thread 1 (r.h.s) accesses `wr[1]`.
 146 Thus, there is a data-race on index 1 of the array.

147 A fundamental challenge of static DRF verification is how to handle loops.
 148 Symbolic execution approaches that unroll loops, *e.g.*, [25, 26], cannot handle
 149 large nor symbolic iteration spaces. Static approaches that use Hoare logic,
 150 *e.g.*, [5, 7, 21], require user-provided loop invariants. Another approach is to re-
 151 duce loops to verifying the satisfiability of a corresponding universally quantified
 152 formula, *e.g.*, [24, 30]. This has the advantage of being fast and not requiring in-
 153 variants. However, its previous application to GPU programming, *i.e.*, PUG,
 154 is unsound due to the interaction between barrier synchronizations and loops,
 155 *e.g.*, PUG misses the data-race in Listing 2.1. We give more details in Section 6.

156 *Our Approach* A key contribution of our work is to identify conditions that allow
 157 a kernel to be reduced to a first-order logic formula, by precisely characterizing
 158 the effect of barrier synchronization in loops. To this end, the language of access
 159 memory protocols distinguishes syntactically between protocol fragments that
 160 synchronize from those that do not. For instance, the protocol in Listing 2.2 is
 161 identified as *unsynchronized*, as it does not include any synchronization.

162 In Section 4, we show that the DRF analysis of unsynchronized protocols can
 163 be precisely reduced to a first-order logic formula, where universally quantified
 164 formulae represent loops, thus obviating the need to unroll them explicitly. For
 165 instance, we reduce the verification of Listing 2.2 to asking whether for all M ,
 166 t_1 , and t_2 , where $t_1 \neq t_2$ are thread identifiers, the following holds:

$$\begin{array}{c} \forall j_1, i_1, j_2, i_2: 0 \leq j_1 < M \wedge 0 \leq i_1 < M \wedge 0 \leq j_2 < M \wedge 0 \leq i_2 < M \implies \\ \{\text{rd}[t_1 + j_1]\} \cup \{\text{wr}[t_1]\} \quad \text{DRF with?} \quad \{\text{rd}[t_2 + j_2]\} \cup \{\text{wr}[t_2]\} \end{array}$$

167 This formula is *unprovable* since `rd[$t_1 + j_1$]` races with `wr[t_2]` when, *e.g.*, $t_1 = 0$,
 168 $t_2 = 1$, $j_1 = 1$, and $M > 1$. Hence, our technique flags Listing 2.2 as racy.

169 **Analysis of synchronized protocols** The protocol in Listing 2.3 (left) models
 170 *all* the interactions over the shared array tile from Listing 2.1. This protocol
 171 consists of one outer loop `r` that contains two inner loops separated by a barrier
 172 synchronization (`sync`). The first inner loop writes (`wr[tid]`) to the array, while
 173 the second reads (`rd[tid + j]`) from the array.

174 This protocol illustrates how our language syntactically differentiates be-
 175 tween protocols fragments that synchronize from those that do not. Concretely,
 176 our language precludes an unsynchronized loop (`forU x ∈ n..m {u}`) from calling
 177 `sync` anywhere in u , and it requires that a synchronized loop (`forS x ∈ n..m {p}`)

Listing 2.3: Access memory protocol (left) of array tile from Listing 2.1 and its aligned version (right).

<pre> 1 for^S r in 0..N { 2 for^U i in 0..M { wr[tid] } 3 sync; 4 for^U j in 0..M { rd[tid + j] } 5 }</pre>	aligns to	<pre> 1 for^U i in 0..M { wr[tid] } 2 sync; 3 for^S r in 1..N { 4 for^U j in 0..M { rd[tid + j] } 5 for^U i in 0..M { wr[tid] } 6 sync; } 7 for^U j in 0..M { rd[tid + j] }</pre>
---	-----------	---

178 includes at least one occurrence of `sync`. The superscript `U` (resp. `S`) stands for
 179 synchronized (resp. *unsynchronized*). This distinction can be inferred automat-
 180 ically and yields a compositional analysis, as we explain below.

181 The behavior of synchronized loops is difficult to analyse because they may
 182 contain data-races that span more than one iteration. For instance an instruction
 183 of iteration r in Listing 2.3 may race with an instruction of iteration $r+1$.

184 *Our Approach* In this work we show that the DRF analysis of synchronized
 185 protocols can safely be reduced to a first-order logic formula when such loops
 186 are *aligned*, *i.e.*, when there is a synchronization exactly before the loop and at
 187 the end of its body. In Section 4.1 we show how to transform an arbitrary access
 188 memory protocol into an aligned protocol using a syntax-driven transformation
 189 technique called *barrier aligning*. Intuitively, barrier aligning normalizes loops
 190 so that they do not “leak” accesses between iterations. The right-hand side of
 191 Listing 2.3 shows the result of applying *barrier aligning* on the protocol from
 192 Listing 2.3 (left). Observe that the fragment before the aligned loop (line 1)
 193 corresponds to the unsynchronized part of the original loop (before `sync`). The
 194 original loop itself is rearranged so that the part succeeding `sync` is moved to
 195 the beginning of the aligned loop (lines 3–6). The fragment following the aligned
 196 loop (line 7) corresponds to the unsynchronized loop that appears after the `sync`
 197 in the original protocol.

198 In Section 4.1 we show that aligned protocols enable *compositional* verifica-
 199 tion: protocol fragments between two barriers can be analyzed independently.
 200 This compositional analysis is possible because (i) there is no causality between
 201 instructions, except through `sync` and (ii) aligned protocols syntactically delimit
 202 the causality induced by `sync`. For instance, the aligned protocol in Listing 2.3
 203 can be reduced to analyzing the following three protocol fragment independently:

$$\begin{array}{l}
 \text{for}^U i \in 0..M \{ \text{wr}[\text{tid}] \} \quad \text{for}^U j \in 0..M \{ \text{rd}[\text{tid} + j] \} \\
 \text{for}^S r \in 1..N \{ \text{for}^U j \in 0..M \{ \text{rd}[\text{tid} + j] \}; \text{for}^U i \in 0..M \{ \text{wr}[\text{tid}] \}; \text{sync} \}
 \end{array}$$

204 The first two protocols are handled like Listing 2.2 because they are unsynchro-
 205 nized. Representing a synchronized loop as a formula becomes possible when
 206 the protocol is *aligned*: both threads must share the same value for r at each
 207 iteration. Hence, we reduce the verification to asking whether for all $N, M, t_1,$

208 and t_2 where $t_1 \neq t_2$ and the following holds:

$$\forall r, j_1, i_1, j_2, i_2: 1 \leq r < N \wedge 0 \leq j_1 < M \wedge 0 \leq i_1 < M \wedge 0 \leq j_2 < M \wedge 0 \leq i_2 < M \\ \implies \{\text{rd}[t_1 + j_1]\} \cup \{\text{wr}[t_1]\} \text{ DRF with? } \{\text{rd}[t_2 + j_2]\} \cup \{\text{wr}[t_2]\}$$

209 Our technique identifies Listing 2.3 as racy since this formula is *unprovable*, i.e.,
210 $\text{rd}[t_1 + j_1]$ races with $\text{wr}[t_2]$ when $r = 1, t_1 = 0, t_2 = 1, j_1 = 1, N > 1$ and $M > 1$.

211 3 Access Memory Protocols

212 An access memory protocol describes the interaction between a group of threads
213 and a single shared-memory location. A protocol records *where* in memory ac-
214 cesses take place, but abstracts away from *what* data is read from/written to
215 memory. The language of protocols distinguishes between an unsynchronized
216 protocol fragment $u \in \mathcal{U}$, that disallows synchronization, from a synchronized
217 fragment $p \in \mathcal{S}$ that must include a synchronization. The syntax and semantics
218 of access memory protocols is given in Figure 2. Our operational semantics is in-
219 spired by the synchronous, delayed semantics (SVD) from Betts et al. [6], where
220 threads execute independently and communicate upon reaching a barrier.

221 Hereafter, i, j, k are metavariables over non-negative integers picked from the
222 set \mathbb{N} . An arithmetic expression n is either: an integer variable x , an integer i ,
223 or a binary operation on integers that yields an integer. A boolean expression b
224 is either a boolean literal, an arithmetic comparison \diamond , or a propositional logic
225 connective \circ . We write $n \downarrow i$ when expression n evaluates to integer i , where
226 evaluation is defined in the natural way. We overload the notation for Boolean
227 expressions, e.g., $b \downarrow \text{true}$ means that expression b evaluates to **true**.

228 *Unsynchronized fragment* A protocol $u \in \mathcal{U}$ either does nothing (**skip**), accesses
229 a shared memory location $o[i]$ (reads from/writes to index i), performs sequential
230 composition, or loops. Figure 2 gives the semantics of unsynchronized protocols,
231 which is parameterized by a set of thread identifiers $\mathcal{T} \subseteq \mathbb{N}$, where $|\mathcal{T}| \geq 2$.

232 Evaluation of an unsynchronized protocol u by a thread identifier i , written
233 $u \downarrow_i P$, yields a *phase*, i.e., a set $P \in \mathcal{P}$ of *access values* $\alpha \in \mathbb{A}$. Each access
234 value, or just access, notation $i:o[j]$, consists of its issuing thread identifier i ,
235 an access mode o (read/write), and an index j . Protocol **skip** produces no ac-
236 cesses. A memory access $o[n]$ evaluates the index and creates a singleton phase.
237 Sequencing, branching, and looping are standard. Similarly to SVD, Rule \mathcal{U} -par
238 executes a copy of the unsynchronized code u for each thread $i \in \mathcal{T}$ by replacing
239 the special variable `tid` by the thread identifier, $u[\text{tid} := i]$, which results in the
240 union of the accesses of all threads.

241 *Synchronized fragment* A protocol $p \in \mathcal{S}$ may perform barrier synchroniza-
242 tion **sync**, run unsynchronized code u , perform sequential composition, and loop.
243 Figure 2 gives the semantics of a protocol, notation $p \downarrow H$. Evaluation of a pro-
244 tocol p yields a *history* (ranged over by H), i.e., a list of phases (P) that records
245 how memory was accessed. We use $::$ as list constructor and \cdot for the usual list
246 concatenation operator. Histories are concatenated using the special \odot -operator.

Syntax

$$\begin{array}{ll}
\mathbb{N} \ni i ::= 0 \mid 1 \mid \dots & o ::= \text{wr} \mid \text{rd} \\
n ::= x \mid i \mid n \star n & \mathbb{A} \ni \alpha ::= i:o[i] \\
\mathcal{B} \ni b ::= \text{true} \mid \text{false} \mid n \diamond n \mid b \circ b & \mathcal{P} \ni P ::= \{\alpha_1, \dots, \alpha_n\} \\
\mathcal{U} \ni u ::= \text{skip} \mid o[n] \mid u; u \mid \text{for}^u x \in n..m \{u\} & H ::= [] \mid P :: H \\
\mathcal{S} \ni p ::= \text{sync} \mid u \mid p; p \mid \text{for}^s x \in n..m \{p\} &
\end{array}$$

Big-step semantics for \mathcal{U}

$$\begin{array}{c}
\boxed{u \downarrow_i P} \quad \boxed{u \downarrow_{\mathcal{T}} S} \\
\mathcal{U}\text{-SKIP} \quad \mathcal{U}\text{-ACC} \quad \mathcal{U}\text{-SEQ} \quad \mathcal{U}\text{-FOR-1} \\
\frac{}{\text{skip} \downarrow_i \emptyset} \quad \frac{n \downarrow j}{o[n] \downarrow_i \{i:o[j]\}} \quad \frac{u_1 \downarrow_i P_1 \quad u_2 \downarrow_i P_2}{u_1; u_2 \downarrow_i P_1 \cup P_2} \quad \frac{(n \geq m) \downarrow \text{true}}{\text{for}^u x \in n..m \{u\} \downarrow_i \emptyset} \\
\mathcal{U}\text{-FOR-2} \\
\frac{(n < m) \downarrow \text{true} \quad u[x := n] \downarrow_i P_1 \quad \text{for}^u x \in n + 1..m \{u\} \downarrow_i P_2}{\text{for}^u x \in n..m \{u\} \downarrow_i P_1 \cup P_2} \\
\mathcal{U}\text{-PAR} \\
\frac{S = \bigcup \{u[\text{tid} := i] \downarrow_i P_i \mid i \in \mathcal{T}\}}{u \downarrow_{\mathcal{T}} S}
\end{array}$$

History concatenation and serialization

$$\boxed{H \cdot H} \quad \boxed{H \odot H} \\
[P_1 \dots P_n] \cdot [P_{n+1} \dots P_{n+k}] = [P_1 \dots P_{n+k}] \quad (H \cdot [P]) \odot ([P'] \cdot H') = H \cdot [P \cup P'] \cdot H'$$

Big-step semantics for \mathcal{S}

$$\boxed{p \downarrow H} \\
\mathcal{S}\text{-SYNC} \quad \mathcal{S}\text{-PAR} \quad \mathcal{S}\text{-SEQ} \quad \mathcal{S}\text{-FOR-1} \\
\frac{}{\text{sync} \downarrow [\emptyset, \emptyset]} \quad \frac{u \downarrow_{\mathcal{T}} P}{u \downarrow [P]} \quad \frac{p \downarrow H \quad q \downarrow H'}{p; q \downarrow H \odot H'} \quad \frac{(n + 1 = m) \downarrow \text{true} \quad p[x := n] \downarrow H}{\text{for}^s x \in n..m \{p\} \downarrow H} \\
\mathcal{S}\text{-FOR-2} \\
\frac{(n < m) \downarrow \text{true} \quad p[x := n] \downarrow H \quad \text{for}^s x \in n + 1..m \{p\} \downarrow H'}{\text{for}^s x \in n..m \{p\} \downarrow H \odot H'}$$

Well-formed protocols

$$\boxed{p \in \mathcal{W}} \\
u; \text{sync} \in \mathcal{W} \quad \frac{p \in \mathcal{W} \quad q \in \mathcal{W}}{p; q \in \mathcal{W}} \quad \frac{p \in \mathcal{W} \quad \text{tid} \notin \text{fv}(n) \cup \text{fv}(m)}{u_1; \text{for}^s x \in n..m \{p; u_2\} \in \mathcal{W}}$$

Data-race, safe phase, and safe history

$$\boxed{\alpha \# \beta} \quad \boxed{\text{safe}(P)} \quad \boxed{\text{safe}(H)} \\
\frac{\text{wr} \in \{o, o'\} \quad i \neq j}{i:o[k] \# j:o'[k]} \quad \frac{\forall \alpha, \beta \in P: \neg(\alpha \# \beta)}{\text{safe}(P)} \quad \frac{\forall P \in H: \text{safe}(P)}{\text{safe}(H)}$$

Fig. 2: Syntax, semantics, and properties of access memory protocols.

247 A barrier synchronization creates two empty phases, corresponding to phases
 248 before and after synchronization. Running an unsynchronized protocol yields a
 249 single phase containing all accesses performed by that protocol. Sequencing two
 250 synchronized protocols p with q merges the last phase of the former with the first
 251 phase of the latter, as these two phases run concurrently. Running one iteration
 252 of a synchronized loop sequences the execution of the first iteration with the
 253 rest of the loop, by merging the last phase of the first iteration with the first
 254 phase of the rest of the loop. Synchronized loops in access memory protocols are
 255 nonempty, hence the base case is when there is one iteration left. This additional
 256 requirement helps with the presentation of our theory as it implies that every
 257 synchronized loop always executes at least one synchronization.

258 A protocol is well-formed, written $p \in \mathcal{W}$, if every unsynchronized fragment is
 259 followed by a barrier synchronization, every synchronized loop includes a barrier
 260 and is not branching on thread-local variables, *i.e.*, `tid`. We write $fv(p)$ (resp.
 261 $fv(n)$) for the free variables of p (resp. n). We discuss how well-formedness is
 262 enforced in Section 5.

263 DRF is formalized at the bottom of Figure 2. Two accesses are in a data-race
 264 (or racy) when there exist two different threads that access the same index k ,
 265 and one of these accesses is a write. Phase P is *safe* iff each pair of access it
 266 contains is not racy. History P is *safe* when all of its phases are safe.

267 4 DRF-Preserving Transformations of Protocols

268 This section presents the main steps of the DRF analysis summarized in Figure 1:
 269 barrier aligning (Section 4.1) and splitting (Section 4.2).

270 This section also includes our key theoretical results. We establish that these
 271 steps preserve and reflect data-races (*i.e.*, any and all data-races are found), see
 272 Theorem 1 and Theorem 3. We make precise the notion of compositionality that
 273 makes our approach scalable in Theorem 2.

274 4.1 Aligning Protocols

275 The first transformation step normalizes protocols by aligning synchronized
 276 loops, which in turn enables a form of compositional verification. The goal of the
 277 transformation is to produce protocols which belong to \mathcal{A} , see top of Figure 3.

278 *Barrier aligning* (or just aligning) is performed by function *align*, given in
 279 the bottom half of Figure 3. The function returns a pair whose first element is an
 280 aligned and synchronized protocol, and whose second element is an unsynchro-
 281 nized protocol. Intuitively, the pair represents a sequence which we delimitate
 282 syntactically. We note that the output of *align*, say (q, u) , can be trivially made
 283 into an aligned protocol: $q; u; \text{sync}$. The case for synchronization is simple, *align*
 284 returns the input protocol as the first component of the pair and *skip* as the
 285 second component (the input protocol is already fully aligned). The case for
 286 sequence consists of the sequential composition of the pair aligned with unsynch-
 287 ronized code using operator \circledast . Sequencing two pairs $(p, u) \circledast (q, u')$ amounts
 288 to sequencing u to the outer-most piece of unsynchronized code present in q .

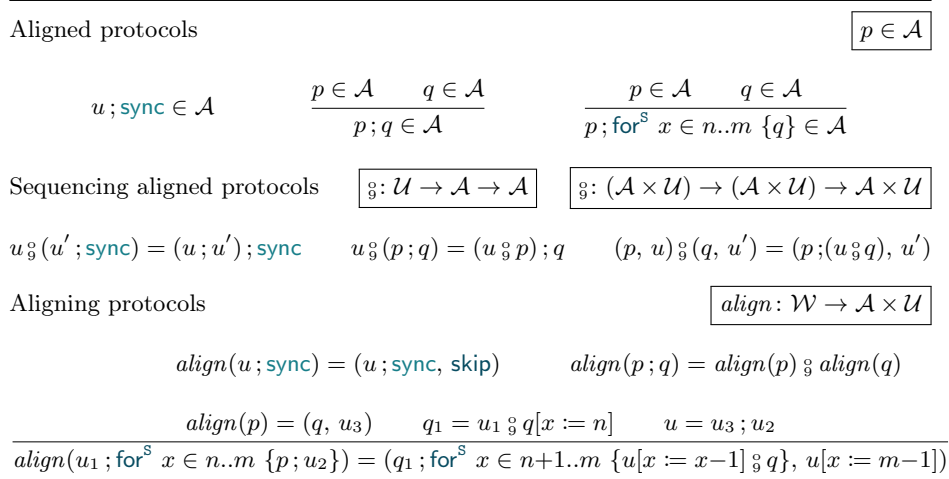


Fig. 3: Aligning protocols.

289 Dealing with synchronized loops is more involved. Given a loop $u_1; \text{for}^s x \in$
 290 $n..m \{p; u_2\}$, we produce a protocol consisting of the fragment preceding the
 291 loop and the synchronized part of its first iteration (q_1), an aligned loop starting
 292 at $n+1$, and the unsynchronized part of its last iteration ($u[x := m-1]$). See
 293 Listing 2.3 for an example of protocol aligning. We note that we can always
 294 unroll the loop because the analysis only considers nonempty synchronized loops;
 295 we discuss how to enforce this assumption in Section 5.

296 We now state two fundamental properties of barrier aligning: preserving and
 297 reflecting DRF (Theorem 1), and enabling compositional verification (Theorem
 298 2). Theorem 1 states that verifying DRF of a well-formed protocol p is
 299 equivalent to verifying DRF of its aligned counterpart.

300 **Theorem 1 (Correctness of Align).** *Let $\text{align}(p) = (q, u)$ and $p \in \mathcal{W}$. If*
 301 *$p \downarrow H_1$ and $q; u \downarrow H_2$, then $\text{safe}(H_1)$ if and only if $\text{safe}(H_2)$.*

302 To state our compositionality result, we introduce a language of contexts:

$$\mathcal{C} ::= [_] \mid u; \text{sync} \mid p; \mathcal{C} \mid \mathcal{C}; p \mid \mathcal{C}; \text{for}^s x \in n..m \{p\} \mid p; \text{for}^s x \in n..m \{\mathcal{C}\}$$

303 The base cases correspond to a hole $[_]$ or an unsynchronized protocol (followed
 304 by sync). The other cases follow the structure of access memory protocols.

305 **Theorem 2 (Compositionality).** *Let \mathcal{C} be a context, s.t. $\mathcal{C}[\text{skip}; \text{sync}]$ is*
 306 *DRF, and $\mathcal{C}[\text{skip}; \text{sync}] \downarrow H$. For all $p \in \mathcal{A}$, if p is DRF, $p \downarrow H'$, and $\text{fv}(p) \subseteq \{\text{tid}\}$,*
 307 *then $\mathcal{C}[p] \in \mathcal{A}$ and $\mathcal{C}[p]$ is also DRF.*

 Syntax

$$\mathcal{L} \ni h ::= \text{skip} \mid n:o[m] \mid h;h \mid \text{var } x \text{ in } n..m;h$$

Product of histories

$$H \otimes H$$

$$H_1 \otimes H_2 = [P_1 \cup P_2 \mid (P_1, P_2) \in H_1 \times H_2]$$

Big-step semantics

$$h \Downarrow H$$

$$\frac{}{\text{skip} \Downarrow [\emptyset]} \quad \frac{n \downarrow i \quad m \downarrow j}{n:o[m] \Downarrow [\{i:o[j]\}]} \quad \frac{h_1 \Downarrow H_1 \quad h_2 \Downarrow H_2}{h_1;h_2 \Downarrow H_1 \otimes H_2} \quad \frac{(n \geq m) \downarrow \text{true}}{\text{var } x \text{ in } n..m;h \Downarrow [\emptyset]}$$

$$\frac{(n < m) \downarrow \text{true} \quad h[x := n] \Downarrow H_1 \quad \text{var } x \text{ in } n+1..m;h \Downarrow H_2}{\text{var } x \text{ in } n..m;h \Downarrow H_1 \cdot H_2}$$

Projection

$$\text{trace}: \mathcal{U} \rightarrow \mathcal{L}$$

$$\text{trace}(o[n]) = \text{tid}:o[n] \quad \text{trace}(\text{for}^u x \in n..m \{u\}) = \text{var } x \text{ in } n..m; \text{trace}(u)$$

$$\text{trace}(u_1; u_2) = \text{trace}(u_1); \text{trace}(u_2) \quad \text{trace}(\text{skip}) = \text{skip}$$

Splitting protocols

$$\text{split}: \mathcal{A} \rightarrow [\mathcal{L}]$$

$$\text{split}(p; q) = \text{split}(p) \cdot \text{split}(q)$$

$$\frac{t_1, t_2 \text{ fresh} \quad h_1 = \text{trace}(u)[\text{tid} := t_1] \quad h_2 = \text{trace}(u)[\text{tid} := t_2]}{\text{split}(u; \text{sync}) = [\text{var } t_1 \text{ in } 1..|\mathcal{T}|; \text{var } t_2 \text{ in } 0..t_1; h_1; h_2]}$$

$$\text{split}(p; \text{for}^s x \in n..m \{q\}) = \text{split}(p) \cdot [\text{var } x \text{ in } n..m; h \mid h \in \text{split}(q)]$$

Fig. 4: Syntax and semantics of symbolic traces, and splitting of protocols.

308 4.2 Splitting Protocols into Symbolic Traces

309 The second verification step, *splitting*, consists in transforming an aligned proto-
 310 col into *symbolic traces*, *i.e.*, symbolic representations of sets of memory accesses
 311 which occur between two synchronizations.

312 *Symbolic traces* Intuitively, symbolic traces are a thin abstraction over an SMT
 313 formula. We describe how to translate a symbolic trace to a formula in Section 5.

314 We give the syntax and semantics of symbolic traces in Figure 4. Express-
 315 ion **skip** terminates a trace. Expression $n:o[m]$ states that thread n accesses
 316 index m with mode o . Expression $h_1;h_2$ composes two symbolic traces using
 317 operator \otimes , also given in Figure 4. Expression $\text{var } x \text{ in } n..m;h$ binds variable x
 318 in h , where variable x is an integer in the range induced from n and m . The
 319 semantics of a symbolic trace yields a history with a phase for each possible vari-

320 able assignment. Expression `skip` yields a single empty phase. Expression $n:o[m]$
 321 evaluates to a singleton set that contains the access value that results from eval-
 322 uating the thread-identifier expression n and the index expression m . Sequencing
 323 histories $h_1; h_2$; h_1 consists of performing the product of phases (operator \otimes), which
 324 consists of merging every phase of H_1 with every phase of H_2 . A variable binder
 325 behaves like a `skip` when the range of values is empty. Otherwise, we fork two histo-
 326 ries H_1 and H_2 . We assign the lower bound of the set in H_1 , and we recursively
 327 evaluate a variable binder where we increment its lower bound in H_2 .

328 *Barrier splitting* is the transformation from aligned protocols to symbolic traces,
 329 performed via functions *trace* and *split*, defined in Figure 4. Function *trace*
 330 extracts the symbolic trace of an unsynchronized program for a single thread.
 331 Memory accesses are tagged with the owner thread `tid`, and unsynchronized loops
 332 are converted into variable bindings. Function *split* returns a list of symbolic
 333 traces. The case for $p; q$ is trivial (operator \cdot stands for list concatenation). The
 334 base case of *split* is for unsynchronized protocol fragment u , which produces a
 335 list containing a single symbolic trace. It introduces fresh variables t_1 and t_2
 336 that represent two (distinct) symbolic thread identifiers. The rest of the trace
 337 consists of the trace of u instantiated to the first thread identifier t_1 followed
 338 by its instantiation to the second thread identifier t_2 . The case for synchronized
 339 loops simply reinterprets the loop as a variable binder. Function *split* leads to an
 340 exponential blow up wrt. nesting of synchronized loops, but this has not posed
 341 problems in practice, *c.f.*, Claim 2.

342 *Example 1.* Let $\hat{p} = \text{wr}[\text{tid} + 1]; \text{rd}[\text{tid} + 2]; \text{sync}$. We have that *split*(\hat{p}) returns:

$$\text{var } t_1 \text{ in } 1..|\mathcal{T}|; \text{var } t_2 \text{ in } 0..t_1; t_1:\text{wr}[t_1+1]; t_1:\text{rd}[t_1+2]; t_2:\text{wr}[t_2+1]; t_2:\text{rd}[t_2+2]$$

343 We show that barrier splitting preserves and reflects DRF.

344 **Theorem 3.** *Let $p \in \mathcal{A}$, such that $p \downarrow H_1$, and $H_2 = [H \mid h \in \text{split}(p) \wedge h \downarrow H]$,
 345 then $\text{safe}(H_1)$ if and only if $\text{safe}(H_2)$.*

346 Hence we have established that aligning (Theorem 1) and splitting (Theorem 3)
 347 preserve and reflect data-races, *i.e.*, any and all data-races are found. Thus,
 348 the only source of approximation in our analysis stems from the inference of
 349 protocols from CUDA kernels, which we discuss in the next section. Theorem 3
 350 highlights the compositionality of our analysis, as each symbolic trace resulting
 351 from function *split* can be analyzed independently.

352 5 Implementation

353 In this section we present our tool, *Faial*, that implements the steps described
 354 in Figure 1. *Faial* takes a CUDA kernel as input and produces results that ei-
 355 ther identify the kernel as DRF or list specific data-races. In this section, we
 356 describe the implementation of the protocol inference, well-formedness checks,
 357 and transformation to SMT.

358 *Inference* This step transforms a CUDA kernel into access memory protocols
 359 (one for each shared array). It uses `libclang` [22] to parse the kernel, a standard
 360 single static assignment (SSA) transformation to simplify the analysis of indices
 361 and arrays, and code slicing to only retain code related to *shared* array accesses.
 362 We note that `Faial` supports constructs of the CUDA programming model that
 363 are not directly modeled by access memory protocols, *e.g.*, unstructured loops,
 364 conditionals, function calls, and multi-dimensional arrays. To support multi-
 365 dimensional thread identifiers, we extend the language of protocols to support
 366 multiple thread identifiers, and adapt function *split* accordingly. The main chal-
 367 lenges are related to loops and function calls.

368 Whenever possible loops are transformed to loops with a stride of 1 following
 369 ideas from loop normalization [23] and abstraction [30]. For instance, in `for(int`
 370 `i=lb;i<ub;i+=s){S}` we change the stride from `s` into 1 by executing the loop body `S`
 371 when the loop variable `i` is divisible by stride, *i.e.*, the loop becomes `for(int`
 372 `i=lb;i<ub;i++)if((i+lb)%s==0){S}`. Similarly, a loop ranging over powers of `n`, *e.g.*,
 373 `for(int i=lb;i<ub;i*=s)`, becomes `for(int i=lb;i<ub;i++)if(powerof(i,s)){S}`, where func-
 374 tion `powerof(i,s)` tests whether `i` is a power of base `s`. We approximate `while`s as
 375 a structured loop with an unknown upper bound.

376 Function calls that manipulate shared memory are uncommon in GPU pro-
 377 gramming. Additionally auxiliary functions that manipulate shared memory
 378 have a compiler annotation to inline their bodies, hence we can inline such calls
 379 easily. `Faial` cannot handle recursive functions, but these rarely occur in practice.
 380 Function calls that do not access shared memory are simply discarded.

381 *Well-formedness* This step ensures that kernels `Faial` analyzes meet the well-
 382 formedness conditions ($p \in \mathcal{W}$) defined in Figure 2, as well as the assumptions
 383 that synchronized loops iterate at least once. First, `Faial` annotates loops with a
 384 synchronized/unsynchronized tag according to the presence of `sync` in the loop
 385 body, then adjusts the precedence of sequencing to group all unsynchronized code
 386 preceding a `sync` or a synchronized loops. Synchronized loops of well-formed pro-
 387 tocols cannot manipulate thread-local variables (*i.e.*, `tid`), an assumption shared
 388 by the CUDA programming model. Hence, `Faial` flags such kernels as erroneous.
 389 Next, `Faial` adds assertions before/after synchronized loops to check that the
 390 loop range is nonempty, *i.e.*, loops execute at least once. Similarly to loops,
 391 conditionals are tagged as synchronized or unsynchronized. Then, `Faial` inlines
 392 synchronized conditionals, *i.e.*, when a synchronized conditional is found, two
 393 copies of the input program are created and each copy is prefixed by a global
 394 assertion corresponding to the condition. `Faial` does not support synchronized
 395 conditionals that appear within synchronized loops. We have not found real-
 396 world kernels that include such a construction.

397 *Quantification* This step transforms each symbolic trace (Figure 4) into an SMT
 398 formula, to check for *safety*, *c.f.*, Figure 2. Essentially, the generated formula
 399 guarantees that the indices of array accesses are distinct when there is at least
 400 one write. We illustrate this straightforward transformation with Example 2.

401 *Example 2.* The formula generated from the trace in Example 1 is given below:

$$\begin{aligned} \forall t_1, t_2: 1 \leq t_1 < 3 \wedge 0 \leq t_2 < t_1 \wedge (m_1 = \text{wr} \vee m_2 = \text{wr}) \implies \\ & ((\text{id}x_1 = t_1 + 1 \wedge m_1 = \text{wr}) \vee (\text{id}x_1 = t_1 + 2 \wedge m_1 = \text{rd})) \\ & \wedge ((\text{id}x_2 = t_2 + 1 \wedge m_2 = \text{wr}) \vee (\text{id}x_2 = t_2 + 2 \wedge m_2 = \text{rd})) \wedge \text{id}x_1 \neq \text{id}x_2 \end{aligned}$$

402 where each symbolic access is translated to a conjunction representing its index
403 ($\text{id}x$) and access mode (m). Observe that the formula enforces that indices $\text{id}x_1$
404 and $\text{id}x_2$ (executed by distinct threads) are different.

405 For multi-dimensional arrays, we generate one pair of indices per dimension, and
406 check that at least one pair is distinct.

407 6 Experimental Evaluation

408 We evaluate *Faial* over several datasets and show how it fares against existing
409 approaches. We structure this evaluation in three claims.

410 **Claim 1: Correctness.** We claim that our approach finds more bugs and raises
411 fewer false alarms than existing tools. To evaluate this claim, we compare *Faial*
412 against four state-of-the-art kernel verification tools over 10 kernels that are
413 known to be tricky to analyze.

414 **Claim 2: Scalability.** We claim that our approach scales better to larger pro-
415 grams. To evaluate this claim, we compare *Faial* against other tools over a set
416 of synthetic benchmarks designed to test the limits of each tool, in terms of run
417 time and memory usage.

418 **Claim 3: Real-world usability.** We claim that our approach is more usable
419 than existing static verification tools on real-world CUDA programs. To evaluate
420 this claim, we use a varied dataset of real-world DRF kernels and measure the
421 false alarm rate, run time, and memory usage of *Faial*, GPUVerify, and PUG.

422 *Benchmarking environment* To make our evaluation reproducible, we developed
423 a benchmarking framework to automate our experiments over the different tools
424 and datasets. For Claim 1 and Claim 3, we designed a tool-agnostic file format for
425 kernel functions and associated metadata (*e.g.*, expected result of DRF analysis,
426 grid and block dimensions, and include directives). And for Claim 2, we created
427 a tool that generates kernels according to given templates, *e.g.*, see Figure 7.

428 We evaluate *Faial* against the following verification tools: GPUVerify [5] v2018-
429 03-22; PUG [23] v0.2; and, GKLEE [25] and SESA [26] v3.0. Experiments for
430 Claim 1 use an Intel i5-6500 CPU, 7.7GiB RAM, and Fedora 33 OS, while
431 Claim 2 and Claim 3 use an Intel i7-10510U CPU, 16GiB RAM, and Pop! OS.

432 *Excluded tools* We excluded ESBMC-GPU [33] and Simulee [38] from the evalu-
433 ation because we were unable to get them to run satisfactorily. Both tools have
434 rudimentary support for verifying arbitrary CUDA kernels. ESBMC-GPU did not
435 find a single data-race in our benchmarks, while Simulee produced false alarms
436 for every DRF-kernel given.

Table 1: Results for Claim 1. DRF indicates that a (static analysis) tool reported a test case as DRF. NRR indicates that a (symbolic execution) tool did not report any data-race. Label x/y indicates that the tool reported x data-races, y of which are actual races. Label *timeout* indicates that the tool did not terminate within 90s. A test passes if the tool returns the expected result and all reported races are valid.

Test	Expected	Faial	GPUVerify	PUG	GKLEE	SESA
1 transposeDiagonal	Racy	1/1	<i>0/2</i>	<i>DRF</i>	<i>timeout</i>	<i>timeout</i>
	DRF	DRF	<i>0/1</i>	DRF	<i>timeout</i>	<i>timeout</i>
2 first-iter	Racy	1/1	<i>0/1</i>	1/1	<i>timeout</i>	<i>timeout</i>
	DRF	DRF	<i>0/1</i>	<i>0/1</i>	<i>timeout</i>	<i>timeout</i>
3 last-iter	Racy	1/1	1/1	<i>0/1</i>	<i>timeout</i>	<i>timeout</i>
	DRF	DRF	<i>0/1</i>	DRF	<i>timeout</i>	<i>timeout</i>
4 last-iter-first-iter	Racy	1/1	<i>0/1</i>	<i>0/1</i>	<i>timeout</i>	<i>timeout</i>
	DRF	DRF	<i>0/1</i>	<i>0/1</i>	<i>timeout</i>	<i>timeout</i>
5 read-index	Racy	<i>0/1</i>	1/1	<i>0/1</i>	<i>NRR</i>	<i>NRR</i>
	DRF	<i>0/1</i>	DRF	<i>0/1</i>	NRR	NRR
Number of tests passed (of 5):		4	1	0	0	0

437 Claim 1: Correctness

438 We have selected a set of tricky kernels to expose false alarms and missed data-
 439 races in Faial, GPUVerify, PUG, GKLEE, and SESA. Our results are reported
 440 in Table 1. The dataset consists of 5 tests, each consisting of two variations
 441 of the same kernel: one racy and one DRF. The racy version of Test 1 (*c.f.*,
 442 Listing 2.1) contains an inter-iteration data-races. The DRF version adds a `sync`
 443 after the second inner loop. Tests 2 to 4 expose various loop-related data-races.
 444 Their protocols are given in Figure 5. In the racy version of Test 2 `wr[tid + 1]`
 445 conflicts with `wr[tid]` of the first iteration. Similarly, in the racy version of Test 3,
 446 `wr[tid + 1]` of the last iteration races with `wr[tid]`. In the racy version of Test 4 the
 447 last iteration of a nested loop races with the first iteration of the following loop.
 448 Test 5 exposes the abstraction gap between kernel and access memory protocols
 449 (which abstract away array elements), see Figure 6.

450 Faial passes more tests than any other tool. Failed Test 5 (two false alarms)
 451 is caused by access memory protocols abstracting away from *what* data is being
 452 read from/written to arrays, *i.e.*, array elements. We report on performance
 453 trade-offs wrt. tracking array elements in Claim 2.

454 GPUVerify passes Test 5 because it tracks array elements, but fails the re-
 455 maining 4 tests. Some reported false alarms are ill-formed, *e.g.*, on the racy
 456 component of Test 2, the report $(0 : \text{wr}[\text{tid}]; 16 : \text{wr}[\text{tid}])$ has disjoint indices.

457 PUG obtains the worst score amongst static tools. Notably, the tool misses a
 458 data-race in Test 1, demonstrating its unsoundness, *c.f.*, Section 2.1.

459 GKLEE and SESA timeout for tests that include loops, as the loop bounds
 460 are unknown. Both tools miss the data-race in Test 5. Symbolic tools may be
 461 able to report data-races when the bound is known, *e.g.*, timeouts start in Test 1
 462 when the bound is at least 2, in Test 2 when the bound is at least 23,000.

<pre> // first-iter wr[tid+1]; for^s x in 0..N { if (x > 0) { wr[tid] }; sync} </pre>	<pre> // last-iter for^s x in 0..N { sync; if (tid < T -1) { wr[tid+1] }; wr[tid + T] } </pre>	<pre> // last-iter-first-iter for^s x in 1..N+1 { for^s y in 1..x+1 { sync; wr[tid+x+y]}; for^s z in N*2..N*3 { wr[tid+z+1]; sync} } </pre>
--	--	---

Fig. 5: Protocols for Tests 2 to 4, *c.f.*, Claim 1, where N is a free thread-global variable. Yellow shaded code only appears in the DRF version of first-iter and last-iter. Red shaded code only appears in the racy version of last-iter-first-iter.

<pre> // Racy kernel A[tid] = tid; int x = A[tid]; A[x+1] = 0; </pre>	<pre> // Protocol A wr[tid]; rd[tid]; wr[x+1] </pre>	<pre> // DRF kernel A[tid] = tid; int x = A[tid]; A[x] = 0; </pre>	<pre> // Protocol A wr[tid]; rd[tid]; wr[x] </pre>
---	--	--	--

Fig. 6: Kernels and protocols for Test 5 (read-index), *c.f.*, Claim 1; x becomes a free thread-local variable as protocols do not model array elements.

463 Claim 2: Scalability

464 We evaluate the scalability of our approach with a synthetic dataset that aims
465 at demonstrating how different kernel constructs affect run time and memory
466 usage of Faial, GKLEE, GPUVerify, PUG, and SESA. Our dataset is divided into
467 five categories, one per syntactical construct in the language of access mem-
468 ory protocols, as well as conditionals, which are supported by our inference step,
469 *c.f.*, Section 5. Figure 7 shows the protocols of the kernel patterns we generate in
470 each category: (i) repeated accesses (read then write), (ii) repeated barrier syn-
471 chronizations separated by writes, (iii) repeated conditionals, (iv) increasingly
472 nested unsynchronized loops, and (v) increasingly nested synchronized loops. In
473 each category, we vary the problem size by repeating a pattern from 1 to 50
474 times. Note that all kernels generated this way are DRF.

475 Figure 8 shows the average run time and memory usage over five runs on
476 logarithmic and linear scales, respectively. For each run, we set a timeout of 90s
477 and we exclude any run that times out or reports a false alarm. Cutoffs in the
478 memory plots are determined by the cutoffs in the run time plots.

479 Overall Faial is the most scalable tool. In 4 out of 5 categories, Faial has
480 the slowest growth for all experiments, and verifies all tests within 0.46s. In the
481 largest problem sizes, our tool is the fastest in 3 categories (access, conditional,
482 unsynchronized loop), 2nd for barriers, and 3rd for synchronized loops. Overall,
483 the memory usage of Faial is competitive with other tools. Faial is the only tool
484 with a near constant time/memory for up to 50 unsynchronized loops, indicating
485 the scalability of reducing unsynchronized loops to universally quantified formu-
486 las. Faial only times out for kernels which consists of >17 nested synchronized
487 loops. However such kernels are uncommon, *e.g.*, the levels of nested synchron-
488 ized loops in the real-word kernels studied in Claim 3 are at most 3.

<pre>// accesses rd[tid + n1* T]; wr[tid + 1* T]; rd[tid + n2* T]; wr[tid + 2* T]; // ...</pre>	<pre>// barriers wr[tid]; sync; wr[tid]; sync; // ...</pre>	<pre>// conditionals if tid==0 {wr[tid]}; if tid==1 {wr[tid]}; // ...</pre>	<pre>// unsynchronized loops for^u i1 in 0..N { wr[tid]; for^u i2 in 0..N { wr[tid]; // ... }}</pre>	<pre>// synchronized loops for^s i1 in 0..N { wr[tid]; sync; for^s i2 in 0..N { wr[tid]; sync; // ... }}</pre>
---	---	---	--	--

Fig. 7: Synthetic protocols generated for Claim 2. N is a free thread-global variable, and n_1, n_2, \dots are positive integer literals.

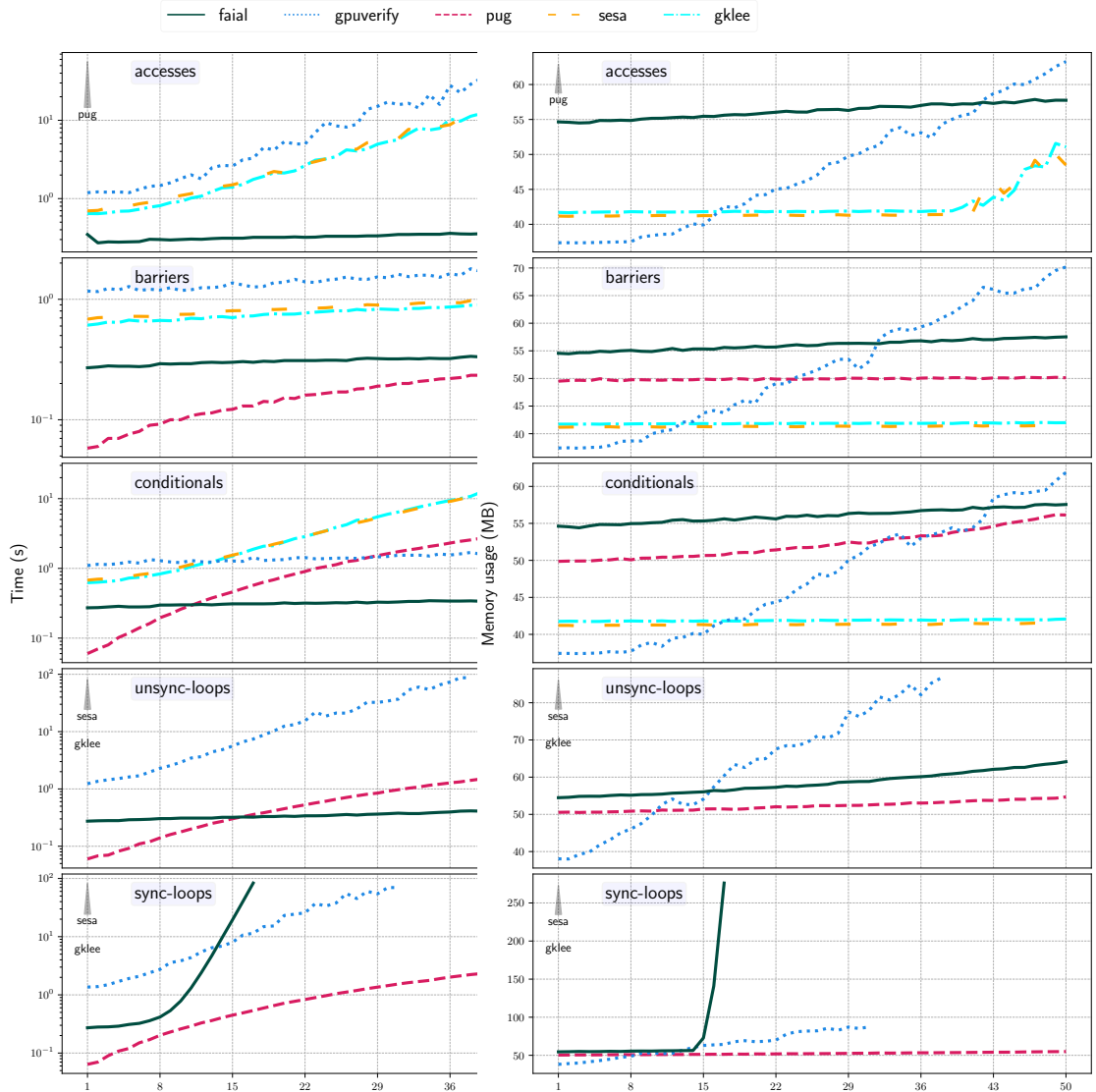


Fig. 8: Results for Claim 2. Run time (left plots) are given on a logarithmic scale, and memory (right plots) are given on a linear scale. Flatter and lower curve is better. Tools annotated with a triangle are excluded due to timeouts or errors.

GPUVerify remains stable in the barrier and conditional categories but is affected negatively by loops and accesses. Loops are a known bottleneck in GPUVerify [2]. In the access category there is an exponential slowdown due to GPUVerify keeping track of what data is being written to/read from array.

PUG tool remains stable with the number of barrier synchronizations but is affected negatively by the number of conditionals and loops. PUG is the fastest tool with smaller inputs, but it raises false alarms in the access category, hence these measurements are omitted from the corresponding plots.

We discuss GKLEE and SESA together since SESA processes GKLEE’s NVCC byte code output by concretizing variables, before passing it to GKLEE itself. There are two main factors that affect negatively these symbolic execution tools: (i) the number of loops, since they unroll each loop; and (ii) the amount of book-keeping required to keep track of what is read from/written to memory. Figure 8 shows clear exponential curves for the access and barrier synchronization categories. Observe that these tools timeout immediately in the loop categories.

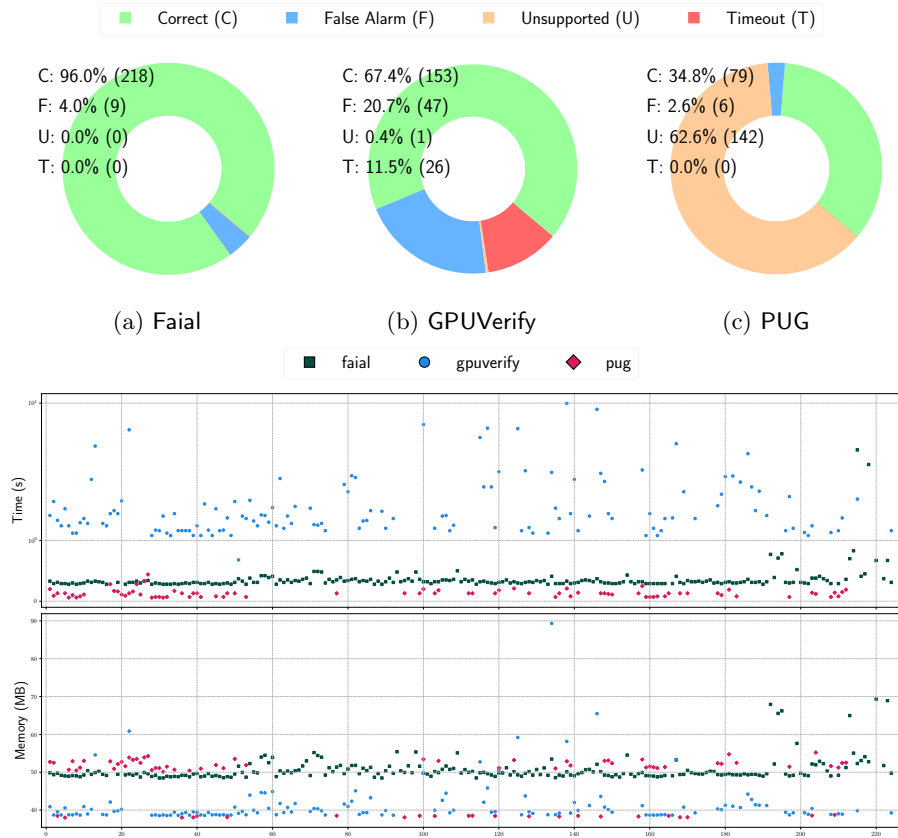
Claim 3: Real-World Usability

TCnoteR1: what is the unsupported kernel by our tool? We evaluate the usability of our approach by comparing Faial with other static verification tools (GPUVerify and PUG) on real-world kernels wrt. rate of false alarm and run time. We curated a set of CUDA kernels from [2], which consists of 3 benchmark suites (totaling 227 CUDA kernels): NVIDIA GPU Computing SDK v2.0 (8 CUDA kernels); NVIDIA GPU Computing SDK v5.0 (166 CUDA kernels); Microsoft C++ AMP Sample Projects (20 kernels); gpgpu-sim benchmarks (33 kernels). All kernels are DRF and have been pre-processed by the authors of [2] to facilitate verification. Each kernel is in a distinct file, all dependencies are available, and kernels are annotated with minimal pre-conditions to allow for automatic analysis (*e.g.*, thread count is given).

As we aim to evaluate fully automatic verification of three tools, we removed code annotations (pre-conditions and loop invariants) specific to GPUVerify. Additionally, we made minor changes to some kernels to meet the limitations of the front-end of Faial and PUG. For instance we converted nested array lookups to use temporary variables and inlined functions calls that operate on arrays in 22 kernels. Another 8 kernels were modified to simplify their control flows. Our curated dataset will be included in our artifact submission.

Figures 9a, 9b, and 9c give the correctness results of Faial, GPUVerify, and PUG, respectively. **Correct** refers to the true-positive rate, *i.e.*, when the tool correctly identifies the kernel as DRF. **False Alarm** refers to the false alarm rate, *i.e.*, when the tool incorrectly identifies the kernel as racy. A kernel is **Unsupported** if it makes the tool crash. A **Timeout** occurs when the tool exceeds the limit of 60s to verify a kernel. The values shown are an average calculated over five runs. Figure 9d shows the average run time and memory usage of every true-positive report (we omit invalid reports) across the three tools.

Overall Faial has the highest rate of true-positives at 96%. Our tool is second in terms of run time and memory usage, showing a good compromise w.r.t. time



(d) Run time (top) and memory usage (bottom) of true-positives. Time (resp. memory) is cropped at 10s (resp. 100MB) and plotted on a logarithmic (resp. linear) scale.

Fig. 9: Results for Claim 3, on a set of 227 DRF CUDA kernels.

533 and space. *Faial* verifies most kernels within 1s, and all kernels that need more
 534 time are only verified by *Faial*. *GPUVerify* shows lower memory usage at the
 535 cost of a higher verification run time. *PUG* verifies the lowest number of kernels
 536 (34.8%), as most kernels are unsupported (62.6%).

537 7 Related Work

538 *SMT-based DRF analyses* Li and Gopalakrishnan propose a direct encoding of
 539 DRF analysis of GPU programs in SMT, with *PUG* [23, 24]. Both *PUG* and *Faial*
 540 follow a similar approach of barrier splitting: having a symbolic representation
 541 of a canonical interleaving, and dividing up the analysis over barrier intervals.
 542 The two major distinctions are that (1) *PUG* misses inter-thread data-races in

543 synchronized loops, *e.g.*, Listing 2.1, and (2) the algorithms of PUG are unspeci-
 544 fied and lack soundness proofs. In [23, §6.3] the authors identify the challenge of
 545 detecting inter-thread data-races, but do not elaborate a solution. Ma *et al.* [30]
 546 present a similar technique to detect data-races and deadlocks in OpenMP pro-
 547 grams (CPU-based parallelism). However, their work does not guarantee DRF,
 548 and they do not formalize their algorithms. In [8], Prasanth *et al.* propose a
 549 polyhedral encoding of DRF for OpenMP programs, which is only applicable to
 550 programs with affine array accesses. However the prevalence of linearized array
 551 expressions in GPU kernels is known to stump polyhedral analysis [15].

552 *Hoare-logic-based DRF analyses* The main drawback of Hoare-logic based tools
 553 is their high rate of false alarms. They also require code annotations from a
 554 concurrency expert to handle loops. GPUVerify [2, 3, 5, 6, 11] can verify CUDA
 555 and OpenCL kernels using Boogie [4] as a backend. GPUVerify also relies on a
 556 two-thread abstraction (pen and paper proof) — in this paper, we present the
 557 first *machine-checked* proof of the two-thread abstraction idea. VeriCUDA [19, 20]
 558 focuses on reasoning about the functional correctness of GPU programs using
 559 Hoare-logic. In [21] the authors extend VeriCUDA to proving DRF. In a sim-
 560 ilar vein, VerCors [7] uses separation logic to prove the functional correctness
 561 and DRF of GPU kernels. Both VeriCUDA and VerCors expect a tool-specific
 562 language, hence cannot handle real-world kernels directly.

563 *Data-race finders* include: dynamic data-race detection, symbolic-execution, and
 564 model-checking. Such techniques are better suited for highly detailed analysis
 565 in smaller kernels, and typically are unable to prove DRF. Dynamic data-race
 566 detection executes a kernel to find data-races on a fixed input, *e.g.*, [13, 17, 18,
 567 27, 28, 32, 39, 40]. This technique only reports real data-races, but suffers from
 568 a slowdown of at least $10\times$ and requires the kernel input data, which might
 569 be unavailable or unknown. Symbolic execution and model checking have been
 570 extended to detect data-races [9, 10, 25, 33, 38]. These techniques do without the
 571 kernel input and can detect more data-races than dynamic data-race detection.

572 *Miscellaneous* Ferrel *et al.* introduce a machine-checked formalism to reason
 573 about the semantics of CUDA assembly [14]. Dabrowski *et al.* mechanize the
 574 DRF-analysis of multithreaded programs [12]. Muller and Hoffmann present a
 575 logic to reason about the evaluation cost of CUDA kernels [31].

576 Session types [35] are a family of behavioral types that codify the message
 577 exchanges that take place over a given *session*, in terms of sends and receives. Ac-
 578 cess memory protocols are similar in that they codify the interactions that take
 579 place over a given shared *array*, in terms of reads and writes. Other behavioral
 580 types have been used to verify parallel and multithreaded systems that com-
 581 municate via message-passing [29, 36, 37]. However these do not capture shared
 582 memory (only message-passing), thus cannot address data-races.

583 8 Conclusion

584 We tackle the problem of statically checking DRF in GPU kernels, with a new
 585 family of behavioral types, *i.e.*, access memory protocols. We provide a novel

586 compositional analysis of access memory protocols, along with fully mechanized
 587 proofs and an implementation. Our evaluation explores challenging and diverse
 588 benchmarks (229 real-world and 258 synthetic kernels) to demonstrate that our
 589 approach is more precise (false alarms and missed alarms), scalable (time/mem-
 590 ory growth), and usable (real-world kernels correctly verified) than other tools.

591 References

- 592 1. Ancona, D., Bono, V., Bravetti, M., Campos, J., Castagna, G., Deniélou, P.M.,
 593 Gay, S.J., Gesbert, N., Giachino, E., Hu, R., Johnsen, E.B., Martins, F., Mascardi,
 594 V., Montesi, F., Neykova, R., Ng, N., Padovani, L., Vasconcelos, V.T., Yoshida,
 595 N.: Behavioral types in programming languages. *Foundations and Trends in Pro-*
 596 *gramming Languages* **3**(2-3), 95–230 (2016). <https://doi.org/10.1561/25000000031>
- 597 2. Bardsley, E., Betts, A., Chong, N., Collingbourne, P., Deligiannis, P., Donaldson,
 598 A.F., Ketema, J., Liew, D., Qadeer, S.: Engineering a static verification tool for
 599 GPU kernels. In: *Proceedings of CAV*. vol. 8559, pp. 226–242. Springer (2014).
 600 https://doi.org/10.1007/978-3-319-08867-9_15
- 601 3. Bardsley, E., Donaldson, A.F., Wickerson, J.: KernelInterceptor: Automating GPU
 602 kernel verification by intercepting kernels and their parameters. In: *Proceedings of*
 603 *IWOCL*. pp. 1–5 (5 2014). <https://doi.org/10.1145/2664666.2664673>
- 604 4. Barnett, M., Chang, B.Y.E., DeLine, R., Jacobs, B., Leino, K.R.M.: Boogie: A
 605 modular reusable verifier for object-oriented programs. In: *Proceedings of FMCO*.
 606 p. 364–387. Springer (2005). https://doi.org/10.1007/11804192_17
- 607 5. Betts, A., Chong, N., Donaldson, A.F., Ketema, J., Qadeer, S., Thomson, P., Wick-
 608 erson, J.: The design and implementation of a verification technique for GPU ker-
 609 nels. *Transactions on Programming Languages and Systems* **37**(3), 1–49 (2015).
 610 <https://doi.org/10.1145/2743017>
- 611 6. Betts, A., Chong, N., Donaldson, A.F., Qadeer, S., Thomson, P.: GPUVerify: a
 612 verifier for GPU kernels. In: *Proceedings of OOPSLA*. pp. 113–132. ACM (2012).
 613 <https://doi.org/10.1145/2384616.2384625>
- 614 7. Blom, S., Huisman, M., Mihelčić, M.: Specification and verification of
 615 GPGPU programs. *Science of Computer Programming* **95**(P3), 376–388 (2014).
 616 <https://doi.org/10.1016/j.scico.2014.03.013>
- 617 8. Chatarasi, P., Shirako, J., Kong, M., Sarkar, V.: An extended polyhedral model
 618 for SPMD programs and its use in static data race detection. In: *Proceedings of*
 619 *LCPC’16*. pp. 106–120. Springer (2017). [https://doi.org/10.1007/978-3-319-52709-](https://doi.org/10.1007/978-3-319-52709-3_10)
 620 [3_10](https://doi.org/10.1007/978-3-319-52709-3_10)
- 621 9. Collingbourne, P., Cadar, C., Kelly, P.H.J.: Symbolic testing of OpenCL code. In:
 622 *Proceedings of HVC*. pp. 203–218. Springer (2012). [https://doi.org/10.1007/978-](https://doi.org/10.1007/978-3-642-34188-5_18)
 623 [3-642-34188-5_18](https://doi.org/10.1007/978-3-642-34188-5_18)
- 624 10. Collingbourne, P., Cadar, C., Kelly, P.H.: Symbolic crosschecking of floating-
 625 point and SIMD code. In: *Proceedings of EuroSys*. pp. 315–328. ACM (2011).
 626 <https://doi.org/10.1145/1966445.1966475>
- 627 11. Collingbourne, P., Donaldson, A.F., Ketema, J., Qadeer, S.: Interleaving and lock-
 628 step semantics for analysis and verification of GPU kernels. In: *Proceedings of*
 629 *ESOP*. pp. 270–289. Springer (2013). [https://doi.org/10.1007/978-3-642-37036-](https://doi.org/10.1007/978-3-642-37036-6_16)
 630 [6_16](https://doi.org/10.1007/978-3-642-37036-6_16)
- 631 12. Dabrowski, F., Pichardie, D.: A certified data race analysis for a Java-
 632 like language. In: *Proceedings of TPHOL*, pp. 212–227. Springer (2009).
 633 https://doi.org/10.1007/978-3-642-03359-9_16

- 634 13. Eizenberg, A., Peng, Y., Pigli, T., Mansky, W., Devietti, J.: BARRACUDA:
635 Binary-level Analysis of Runtime RAces in CUDA programs. In: Proceedings of
636 PLDI. pp. 126–140. ACM (2017). <https://doi.org/10.1145/3062341.3062342>
- 637 14. Ferrell, B., Duan, J., Hamlen, K.W.: CUDA au Coq: A framework for machine-
638 validating GPU assembly programs. In: Proceedings of DATE. pp. 474–479 (2019).
639 <https://doi.org/10.23919/DATE.2019.8715160>
- 640 15. Grosser, T., Ramanujam, J., Pouchet, L.N., Sadayappan, P., Pop, S.: Optimistic
641 delinearization of parametrically sized arrays. In: Proceedings of ICS. pp. 351–360.
642 ACM (2015). <https://doi.org/10.1145/2751205.2751248>
- 643 16. ul Hassan Khan Khan, A., Al-Mouhamed, M., Fatayer, A., Almousa, A.,
644 Baqais, A., Assayony, M.: Padding free bank conflict resolution for CUDA-
645 based matrix transpose algorithm. In: Proceedings of SNPD. pp. 1–6 (2014).
646 <https://doi.org/10.1109/SNPD.2014.6888709>
- 647 17. Holey, A., Mekkat, V., Zhai, A.: HAccRG: Hardware-accelerated data
648 race detection in GPUs. In: Proceedings of ICPP. pp. 60–69 (2013).
649 <https://doi.org/10.1109/ICPP.2013.15>
- 650 18. Kamath, A.K., George, A.A., Basu, A.: ScoRD: A scoped race detec-
651 tor for GPUs. In: Proceedings of ISCA. pp. 1036–1049. IEEE (2020).
652 <https://doi.org/10.1109/ISCA45697.2020.00088>
- 653 19. Kojima, K., Igarashi, A.: A Hoare logic for SIMT programs. In: Proceedings of
654 APLAS. vol. 8301, pp. 58–73. Springer (2013). https://doi.org/10.1007/978-3-319-03542-0_5
- 655 20. Kojima, K., Igarashi, A.: A Hoare logic for GPU kernels. *Transactions on Compu-*
656 *tational Logic* **18**(1), 1–43 (2017). <https://doi.org/10.1145/3001834>
- 657 21. Kojima, K., Imanishi, A., Igarashi, A.: Automated verification of functional correct-
658 ness of race-free GPU programs. *Journal of Automated Reasoning* **60**(3), 279–298
659 (2018). <https://doi.org/10.1007/s10817-017-9428-2>
- 660 22. Lattner, C., Adve, V.: LLVM: A compilation framework for lifelong program
661 analysis & transformation. In: Proceedings of CGO. pp. 75–88. IEEE (2004).
662 <https://doi.org/10.1109/CGO.2004.1281665>
- 663 23. Li, G., Gopalakrishnan, G.: Scalable SMT-based verification of GPU
664 kernel functions. In: Proceedings of FSE. pp. 187–196. ACM (2010).
665 <https://doi.org/10.1145/1882291.1882320>
- 666 24. Li, G., Gopalakrishnan, G.: Parameterized verification of GPU ker-
667 nel programs. In: Proceedings of IPDPSW. pp. 2450–2459 (2012).
668 <https://doi.org/10.1109/IPDPSW.2012.302>
- 669 25. Li, G., Li, P., Sawaya, G., Gopalakrishnan, G., Ghosh, I., Rajan, S.P.: GKLEE:
670 Concolic verification and test generation for GPUs. In: Proceedings of PPOPP.
671 vol. 47, pp. 215–224. ACM (2012). <https://doi.org/10.1145/2370036.2145844>
- 672 26. Li, P., Li, G., Gopalakrishnan, G.: Practical symbolic race checking
673 of GPU programs. In: Proceedings of SC. pp. 179–190. IEEE (2014).
674 <https://doi.org/10.1109/SC.2014.20>
- 675 27. Li, P., Ding, C., Hu, X., Soyata, T.: LDetecter: A low overhead race detector for
676 GPU programs. In: Proceedings of WoDet (2014), [http://wodet.cs.washington.](http://wodet.cs.washington.edu/wp-content/uploads/2014/02/wodet2014-final14.pdf)
677 [edu/wp-content/uploads/2014/02/wodet2014-final14.pdf](http://wodet.cs.washington.edu/wp-content/uploads/2014/02/wodet2014-final14.pdf)
- 678 28. Li, P., Hu, X., Chen, D., Brock, J., Luo, H., Zhang, E.Z., Ding, C.: LD: Low-
679 overhead GPU race detection without access monitoring. *Transactions on Architec-*
680 *ture and Code Optimization* **14**(1), 1–25 (2017). <https://doi.org/10.1145/3046678>
- 681 29. López, H.A., Marques, E.R.B., Martins, F., Ng, N., Santos, C., Vasconcelos,
682 V.T., Yoshida, N.: Protocol-based verification of message-passing par-
683

- 684 allel programs. In: Proceedings of OOPSLA. pp. 280–298. ACM (2015).
685 <https://doi.org/10.1145/2814270.2814302>
- 686 30. Ma, H., Diersen, S.R., Wang, L., Liao, C., Quinlan, D., Yang, Z.: Symbolic analysis
687 of concurrency errors in OpenMP programs. In: Proceedings of ICPP. pp. 510–516.
688 IEEE (2013). <https://doi.org/10.1109/ICPP.2013.63>
- 689 31. Muller, S.K., Hoffmann, J.: Modeling and analyzing evaluation cost of CUDA
690 kernels. Proceedings of the ACM on Programming Languages **5**(POPL) (2021).
691 <https://doi.org/10.1145/3434306>
- 692 32. Peng, Y., Grover, V., Devietti, J.: CURD: A dynamic CUDA
693 race detector. In: Proceedings of PLDI. pp. 390–403. ACM (2018).
694 <https://doi.org/10.1145/3192366.3192368>
- 695 33. Pereira, P., Albuquerque, H., Marques, H., Silva, I., Carvalho, C., Cordeiro, L.,
696 Santos, V., Ferreira, R.: Verifying CUDA programs using SMT-based context-
697 bounded model checking. In: Proceedings of SAC. pp. 1648–1653. ACM (2016).
698 <https://doi.org/10.1145/2851613.2851830>
- 699 34. Ruetsch, G., Micikevicius, P.: Optimizing matrix transpose in CUDA. NVIDIA
700 CUDA SDK Application Note **18** (2009), [https://www.cs.colostate.edu/
701 ~cs675/MatrixTranspose.pdf](https://www.cs.colostate.edu/~cs675/MatrixTranspose.pdf)
- 702 35. Takeuchi, K., Honda, K., Kubo, M.: An interaction-based language and its typing
703 system. In: Proceedings of PARLE. LNCS, vol. 817, pp. 398–413. Springer (1994).
704 https://doi.org/10.1007/3-540-58184-7_118
- 705 36. Vasconcelos, V.T.: Session types for linear multithreaded functional
706 programming. In: Proceedings of PDP. pp. 1–6. ACM (2009).
707 <https://doi.org/10.1145/1599410.1599411>
- 708 37. Vasconcelos, V.T., Ravara, A., Gay, S.: Session types for functional mul-
709 tithreading. In: Proceedings of CONCUR. pp. 497–511. Springer (2004).
710 https://doi.org/10.1007/978-3-540-28644-8_32
- 711 38. Wu, M., Ouyang, Y., Zhou, H., Zhang, L., Liu, C., Zhang, Y.: Simulee: Detecting
712 CUDA synchronization bugs via memory-access modeling. In: Proceedings of ICSE.
713 pp. 937–948. ACM (2020). <https://doi.org/10.1145/3377811.3380358>
- 714 39. Zheng, M., Ravi, V.T., Qin, F., Agrawal, G.: GRace: A low-overhead mechanism
715 for detecting data races in GPU programs. In: Proceedings of PPOPP. pp. 135–146.
716 ACM (2011). <https://doi.org/10.1145/1941553.1941574>
- 717 40. Zheng, M., Ravi, V.T., Qin, F., Agrawal, G.: GMRace: Detecting data races in GPU
718 programs via a low-overhead scheme. Transactions on Parallel and Distributed
719 Systems **25**(1), 104–115 (2014). <https://doi.org/10.1109/TPDS.2013.44>

720 **A Demonstration of Faial**

721 Here we present a demonstration of Faial in three examples.

722 *Example 3.* The following is command-line usage of Faial on Listing 2.3 (l.h.s).
723 In addition to CUDA kernels, Faial allows for the evaluation of access memory
724 protocols stored in the proto file type as shown. Barrier aligning is also shown.

```

725 // Source protocol
726 $ cat inter-iteration.proto
727 shared tile ;
728 const global N,
729       global M,
730       local tid,
731       where distinct [tid] &&
732       N > 0 && M > 0;
733
734 foreach (r in 0..N) {
735   foreach (i in 0..M) { rw tile [tid]; }
736   sync;
737   foreach (j in 0..M) { ro tile [tid + j]; }
738 }
739
740 // Step 3: aligned protocol
741 $ faial -A --steps 3 inter-iteration.proto
742 ; a-lang
743 ; a-prog 1
744 locations : tile ;
745 globals : N, M;
746 locals : tid;
747 invariant : (proj($T1, tid) != proj($T2, tid) && N > 0) && M > 0;
748
749 code {
750   sync;
751   foreach (i in 0..M) {
752     rw tile [tid];
753   }
754   sync;
755   foreach* (r in 1..N) {
756     foreach (j in 0..M) {
757       ro tile [tid + j];
758     }
759     foreach (i in 0..M) {
760       rw tile [tid];
761     }
762     sync;
763   }
764   foreach (j in 0..M) {
765     ro tile [tid + j];
766   }
767   sync;
768 }
769 ; end of a-lang
770
771 // Final analysis
772 $ faial inter-iteration.proto
773 *** DATA RACE ERROR ***
774
775 Array: tile [1]
776 T1 mode: R
777 T2 mode: W
778
779 -----
780 Globals Value
781 -----
782 M          2

```



```

783 -----
784 N      2
785 -----
786 r      1
787 -----
788
789 -----
790 Locals T1 T2
791 -----
792 i      0  0
793 -----
794 j      1  0
795 -----
796 tid    0  1
797 -----
    
```

798 *Example 4.* To show intermediate languages within Faial, we now run Example 1
 799 from Section 4.2. Here we show the inlined protocol, aligned protocol, flattened
 800 phases, and generated booleans.

```

801 // Source protocol
802 $ cat example-1.proto
803 shared A;
804 const local tid
805     where distinct [tid];
806
807 rw A[tid + 1];
808 ro A[tid + 2];
809 sync;
810
811 // Step 1: inline assignments and replace key-values
812 $ faial -A --steps 1 example-1.proto
813 locations : A;
814 globals : ;
815 locals : tid;
816 invariant : proj($T1, tid) != proj($T2, tid);
817
818 code {
819     rw A[1 + tid];
820     ro A[2 + tid];
821     sync;
822 }
823
824 // Step 3: aligned protocol
825 $ faial -A --steps 3 example-1.proto
826 ; a-lang
827 ; a-prog 1
828 locations : A;
829 globals : ;
830 locals : tid;
831 invariant : proj($T1, tid) != proj($T2, tid);
832
833 code {
834     sync;
835     rw A[1 + tid];
836     ro A[2 + tid];
837     sync;
838     sync;
839 }
840 ; end of a-lang
841
842 // Step 6: flatten phases
843 $ faial -A --steps 6 example-1.proto
844 ; flatacc
845 ; acc 1
846 location : A;
    
```

```

847  locals : tid;
848  pre: true;
849  {
850    rw[1 + tid] if proj($T1, tid) != proj($T2, tid);
851    ro[2 + tid] if proj($T1, tid) != proj($T2, tid);
852  }
853  ; end of flatacc
854
855  // Step 7: generate booleans
856  $ faial -A ---steps 7 example-1.proto
857  ; symbexp
858  ; bool 1
859  array: A
860  predicates : ;
861  decls: tid$T2, tid$T1, $T2$mode, $T2$idx$0, $T1$mode, $T1$idx$0;
862  goal: ((tid$T1 != tid$T2 && ($T1$mode ==1 && $T1$idx$0 ==1 +tid$T1)) || (tid$T1 != tid$T2 &&
863    (($T1$mode ==0 && $T2$mode ==1) && $T1$idx$0 ==2 +tid$T1)) && ((tid$T1 != tid$T2
864    && ($T2$mode ==1 && $T2$idx$0 ==1 +tid$T2)) || (tid$T1 != tid$T2 && (($T2$mode ==0 &&
865    $T1$mode ==1) && $T2$idx$0 ==2 +tid$T2))) && $T1$idx$0 ==$T2$idx$0);
866  ; end of symbexp
867
868  // Final analysis
869  $ faial example-1.proto
870  *** DATA RACE ERROR ***
871
872  Array: A[2]
873  T1 mode: W
874  T2 mode: R
875
876  -----
877  Locals  T1 T2
878  -----
879  tid    1  0
880  -----

```

881 *Example 5.* To demonstrate usage on a real-world CUDA kernel, Faial is run on
882 a matrix transpose from [34], *c.f.*, Listing 2.1 (l.h.s.). The inferred protocol is
883 also shown. Note that Faial handles the representation of multiple arrays in the
884 same protocol, but each array is analyzed independently.

```

885  // CUDA kernel source
886  $ cat transposeCoalesced.cu
887  #include <cuda.h>
888
889  #define TILE_DIM 16
890  #define BLOCK_ROWS 16
891
892  __global__ void kernel (float* odata, float * idata, int width, int height, int nreps) {
893
894    __requires(height == 2048);
895    __requires(width == 2048);
896
897    __shared__ float tile[TILE_DIM][TILE_DIM];
898
899    int xIndex = blockIdx.x * TILE_DIM + threadIdx.x;
900    int yIndex = blockIdx.y * TILE_DIM + threadIdx.y;
901    int index_in = xIndex + (yIndex)*width;
902
903    xIndex = blockIdx.y * TILE_DIM + threadIdx.x;
904    yIndex = blockIdx.x * TILE_DIM + threadIdx.y;
905    int index_out = xIndex + (yIndex)*height;
906
907    for (int r=0; r < nreps; r++) {
908      for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
909        tile[threadIdx.y+i][threadIdx.x] = idata[index_in+i*width];
910      }

```

```

911     __syncthreads();
912
913
914     for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
915         odata[index_out+i*height] = tile[threadIdx.x][threadIdx.y+i];
916     }
917 }
918 }
919
920 // Inferred protocol
921 $ faial -b 16,16 -g 64,64 -A ---steps 0 transposeCoalesced.cu
922 arrays : tile , odata, idata;
923 scalars : width, nreps, height, gridDim.y, gridDim.x, blockDim.y, blockDim.x, blockDim.y,
924          blockDim.x;
925 pre: (((proj($T1, threadIdx.x) != proj($T2, threadIdx.x) || proj($T1, threadIdx.y) != proj($T2,
926        threadIdx.y)) && blockDim.x < gridDim.x) && threadIdx.x < blockDim.x) && blockDim.y <
927        gridDim.y) && threadIdx.y < blockDim.y;
928
929 code {
930     local threadIdx.y;
931     local threadIdx.x;
932     assert (height == 2048)
933     assert (width == 2048)
934     local xIndex = (16 * blockDim.x) + threadIdx.x;
935     local yIndex = (16 * blockDim.y) + threadIdx.y;
936     local index_in = xIndex + (yIndex * width);
937     local xIndex = (16 * blockDim.y) + threadIdx.x;
938     local yIndex = (16 * blockDim.x) + threadIdx.y;
939     local index_out = xIndex + (yIndex * height);
940     foreach (r in 0 .. nreps) {
941         foreach (i in 0 .. 16; i + 16) {
942             rw tile[threadIdx.y + i, threadIdx.x];
943             ro idata[index_in + (i * width)];
944         }
945         sync;
946         foreach (i in 0 .. 16; i + 16) {
947             rw odata[index_out + (i * height)];
948             ro tile[threadIdx.x, threadIdx.y + i];
949         }
950     }
951 }
952
953 // Final analysis
954 $ faial -b 16,16 -g 64,64 transposeCoalesced.cu
955 *** DATA RACE ERROR ***
956
957 Array: tile [15, 14]
958 T1 mode: W
959 T2 mode: R
960
961 -----
962 | Globals | Value |
963 |-----|-----|
964 | blockDim.x | 0 |
965 |-----|-----|
966 | blockDim.y | 0 |
967 |-----|-----|
968 | nreps | 2 |
969 |-----|-----|
970 | r | 1 |
971 |-----|-----|
972
973 -----
974 | Locals | T1 | T2 |
975 |-----|-----|
976 | i | 0 | 0 |
977 |-----|-----|
978 | i1 | 0 | 0 |

```

```
979 -----  
980 threadldx.x 15 14  
981 -----  
982 threadldx.y 14 15  
983 -----
```