KOSTAS FERLES\*, The University of Texas at Austin, USA

BENJAMIN SEPANSKI\*, The University of Texas at Austin, USA

RAHUL KRISHNAN, The University of Texas at Austin, USA

JAMES BORNHOLT, The University of Texas at Austin, USA

ISIL DILLIG, The University of Texas at Austin, USA

A monitor is a widely-used concurrent programming abstraction that encapsulates all shared state between threads. Monitors can be classified as being either *implicit* or *explicit* depending on the primitives they provide. Implicit monitors are much easier to program but typically not as efficient. To address this gap, there has been recent research on automatically synthesizing explicit-signal monitors from an implicit specification [Ferles et al. 2018], but prior work does not exploit all paralellization opportunities due to the use of a *single* lock for the entire monitor. This paper presents a new technique for synthesizing *fine-grained* explicit-synchronization protocols from implicit monitors. Our method is based on two key innovations: First, we present a new static analysis for inferring *safe interleavings* that allow violating mutual exclusion of monitor operations *without* changing its semantics. Second, we use the results of this static analysis to generate a MaxSAT instance whose models correspond to correct-by-construction synchronization protocols. We have implemented our approach in a tool called CORTADO and evaluate it on monitors that contain parallelization opportunities. Our evaluation shows that CORTADO can synthesize synchronization policies that are competitive with, or even better than, expert-written ones on these benchmarks.

### 1 INTRODUCTION

Concurrent programming is difficult because it requires developers to consider interactions between multiple threads of execution and mediate access to shared resources and data. Programming languages can offer higher-level abstractions to reduce this complexity by making concurrent programming more declarative. One such abstraction is the *monitor* [Hansen 1973; Hoare 1974], which is an object that encapsulates shared state and allows threads access to it only through a set of *operations*, between which the monitor enforces mutual exclusion.

Ideally, developers would implement monitors using *implicit synchronization*, wherein the *only* synchronization primitive is a waituntil(P) operation that blocks threads until condition P is satisfied. The compiler or runtime can then automatically generate the necessary *explicit synchronization* operations (locks, condition variables, etc.) to implement the monitor in a way that respects the semantics of the implicit monitor. However, automatically deriving an efficient explicit monitor from its implicit specification is a challenging problem, and there have been several recent research efforts, including both run-time techniques like AutoSynch [Hung and Garg 2013] and compile-time tools like Expresso [Ferles et al. 2018], to support implicit-synchronization monitors.

While these state-of-the-art approaches make it possible to program using implicit monitors, they still achieve sub-optimal performance because they adhere closely to the monitor's mutual exclusion

<sup>41</sup> \*Both authors contributed equally to the paper.

Authors' addresses: Kostas Ferles, Computer Science Department, The University of Texas at Austin, USA, kferles@cs.utexas.
 edu; Benjamin Sepanski, Computer Science Department, The University of Texas at Austin, USA, ben\_sepanski@utexas.edu;
 Rahul Krishnan, Computer Science Department, The University of Texas at Austin, USA, rahulk@cs.utexas.edu; James
 Bornholt, Computer Science Department, The University of Texas at Austin, USA, bornholt@cs.utexas.edu; Isil Dillig,
 Computer Science Department, The University of Texas at Austin, USA, bornholt@cs.utexas.edu; Isil Dillig,

48 https://doi.org/

1 2

3

5

6

7 8

9

10

11

12

13

14

15

16

17

18

19

20

21

22 23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

<sup>47 2018. 2475-1421/2018/1-</sup>ART1 \$15.00

requirement. They generally use a single lock for the entire monitor and allow access by at most 50 one thread at a time across all monitor operations. In practice, however, many monitors can admit 51 52 additional concurrency while still preserving the *appearance* of mutual exclusion. For example, consider a FIFO queue monitor that provides take and put operations. These two operations can 53 safely run concurrently unless the queue is empty or full, as they will not access the same slot in 54 the queue. Today, realizing this *fine-grained* concurrency requires expert developers to fall back 55 to hand-written explicit synchronization. These implementations are subtle and error-prone, and 56 57 there is no easy way for developers to determine when they have extracted the maximum possible concurrency from such an implementation. 58

This paper presents a new technique to automatically synthesize fine-grained explicit-synchronization 59 monitors. Our technique takes as input an implicit monitor that specifies the desired operations and 60 automatically generates an implementation that allows as much concurrency as possible between 61 those operations while still preserving the appearance of mutual exclusion. The key idea is to de-62 compose each monitor operation into a set of *fragments* and allocate a set of locks to each fragment 63 to enforce the mutual exclusion requirement while allowing as many fragments as possible to run 64 concurrently. The resulting implementation selectively acquires and releases locks at fragment 65 boundaries within each operation and signals condition variables as needed. 66

At a high level, our approach operates in three phases to generate a high-performance explicit synchronization monitor from its implicit version:

- **Signal placement:** First, we use an off-the-shelf technique [Ferles et al. 2018] to infer a signaling regime which determines where to insert signaling operations on condition variables. While the output of this tool is sufficient to synthesize a single-lock implementation, it does not admit any additional concurrency wherein different threads can perform monitor operations simultaneously.
- **Static analysis:** Second, we perform static analysis to infer sufficient conditions for correctness. That is, the output of the static analysis is a set of conditions such that if the synthesized monitor obeys them, it is guaranteed to be correct-by-construction. A key challenge for this static analysis is to determine which fragments can safely execute concurrently without creating a potential violation of the monitor semantics. The analysis simulates *interleaving* each fragment between the fragments of other operations and determines which possible interleavings are safe.
  - Synchronization protocol synthesis via MaxSAT: Finally, we reduce the synthesis problem to a maximum satisfiability (MaxSAT) instance from whose solution an explicit sychronization protocol can be extracted. The hard constraints in the MaxSAT problem enforce the correctness requirements extracted by the static analysis, while the soft constraints encode two competing objective functions: minimizing the total number of locks used, while maximizing the number of pairs of fragments that can run concurrently.

We have implemented our proposed approach in a tool called CORTADO that operates on Java monitors and evaluated it on a collection of monitor implementations that are (1) drawn from popular open-source projects and (2) contain parallelization opportunities that can be achieved via fine-grained locking. Given only the implicit monitor as input, CORTADO synthesizes explicitsynchronization monitors that perform as well as, or better than, hand-written explicit implementations by expert developers. Compared to state-of-the-art automated tools for synthesizing explicit monitors [Ferles et al. 2018], CORTADO-synthesized monitors extract more concurrency and therefore perform much better (up to 39.1×) on heavily contended workloads.

In summary, this paper makes four main contributions:

- A new technique for automatically synthesizing fine-grained monitor implementations that admit the maximum possible concurrency.
- A novel static static analysis for inferring safe interleaving opportunities between threads.

67

68 69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

```
1
                                                                            class ArravBlockingOueue {
              class ArrayBlockingQueue {
                                                                       2
                                                                              int first = 0, last = 0; Object[] queue;
          1
99
                int first = \emptyset, last = \emptyset, count = \emptyset;
                                                                              AtomicInteger count = new AtomicInteger(0);
         2
                                                                      3
                Object[] queue;
          3
                                                                       4
100
                                                                              Lock putLock = new Lock(), takeLock = new Lock();
                                                                       5
          4
101
                ArrayBlockingQueue(int capacity) {
                                                                              Condition notFull = putLock.newCondition();
          5
                                                                       6
                  if (capacity < 1)</pre>
                                                                              Condition notEmpty = takeLock.newCondition();
          6
                                                                       7
102
                                                                              // Constructor is the same as the implicit version.
          7
                    throw new IllegalArgumentException();
                                                                       8
103
          8
                  this.queue = new Object[capacity];
                                                                       9
          9
                                                                      10
                                                                              void put(Object o) {
104
         10
                                                                      11
                                                                                putLock.lock()
105
                void put(Object o) {
                                                                                while (count.get() == queue.length)
         11
                                                                      12
                  // Fragment 1
                                                                                 notFull.await();
         12
                                                                      13
106
         13
                  waituntil(count < queue.length);</pre>
                                                                      14
                                                                                queue[last] = o;
107
                                                                                last = (last + 1) % queue.length;
                  // Fragment 2
         14
                                                                      15
         15
                  queue[last] = o;
                                                                      16
                                                                                int c = count.getAndIncrement();
108
         16
                  // Fragment 3
                                                                      17
                                                                                putLock.unlock();
109
                  last = (last + 1) % queue.length;
                                                                                if (c == 0) {
         17
                                                                      18
         18
                  // Fragment 4
                                                                      19
                                                                                 takeLock.lock():
110
         19
                  count++;
                                                                      20
                                                                                  notEmpty.signalAll();
111
                                                                                  takeLock.unlock();}}
                3
         20
                                                                      21
         21
                                                                      22
112
                Object take() {
         22
                                                                      23
                                                                              Object take() {
113
         23
                  // Fragment 5
                                                                                takeLock.lock():
                                                                      24
         24
                  waituntil(count > 0):
                                                                      25
                                                                                while (count.get() == 0)
114
                                                                                 notEmpty.await();
         25
                  // Fragment 6
                                                                      26
115
                  Object r = queue[first];
                                                                                Object r = queue[first];
         26
                                                                      27
                                                                                queue[first] = null;
         27
                  queue[first] = null;
                                                                      28
116
                                                                                first = (first + 1) % queue.length;
         28
                  // Fragment 7
                                                                      29
117
                  first = (first + 1) % queue.length;
                                                                                int c = count.getAndDecrement();
         29
                                                                      30
         30
                  // Fragment 8
                                                                      31
                                                                                takeLock.unlock();
118
                                                                                if (c == queue.length) {
         31
                  count--;
                                                                      32
119
                                                                                  putLock.lock();
         32
                  return r:
                                                                      33
         33
                }
                                                                      34
                                                                                  notFull.signalAll();
120
              }
                                                                                  putLock.unlock();}
         34
                                                                      35
121
                                                                                return r;}}
                                                                      36
```

(a) Implicit-synchronization ArrayBlockingQueue.

(b) Explicit-synchronization ArrayBlockingQueue.

Fig. 1. Motivating example.

- A MaxSAT encoding to automate reasoning about *both* the correctness and performance of the synthesized explicit-synchronization monitor.
- An implementation of our technique, CORTADO, that outperforms both state-of-the-art automated tools and expert-written code on benchmarks that can be parallelized via fine-grained locking.

# 2 OVERVIEW

In this section, we give an overview of our approach through a motivating example. Given the implicit-synchronization monitor shown in Figure 1a, our goal is to automatically synthesize an efficient and semantically equivalent explicit-synchronization monitor like the one presented in 134 Figure 1b. In what follows, we walk through this example and describe how our technique is able to automatically generate the code in Figure 1b.

#### Implicit-synchronization monitor 2.1 138

Our technique takes as input an *implicit-synchronization monitor* that specifies which operations 139 should execute atomically and when certain operations are allowed to proceed but does not fix 140 a specific synchronization protocol for realizing that behavior. For example, Figure 1a shows an 141 implicit monitor that implements a limited capacity blocking queue via a bounded circular array 142 buffer. This monitor defines two operations, put and take, that execute atomically (i.e., the body of 143 each method must appear to execute as one indivisible unit). The put operation adds an object if the 144 queue is not full, and take removes an object if the queue is not empty. If one of these method calls 145 cannot proceed (i.e., queue is full or empty), the monitor blocks the calling thread's execution using 146

147

122

123

124 125

126

127

128

129 130

131

132

133

135

a waituntil statement until the operation can be executed. For example, the waituntil statement
 at line 13 in take blocks execution until there is at least one object in the queue.

As Figure 1a illustrates, implicit-synchronization monitors make concurrent programming simpler because they are *declarative*: they merely state which operations are atomic and when operations can proceed, but they do *not* specify a particular synchronization protocol for realizing that desired behavior. However, most programming languages do not offer implicit synchronization facilities; so, concurrent programs must instead be implemented in terms of *explicit* synchronization constructs such as locks and condition variables, as we discuss next.

#### 157 2.2 Explicit-synchronization monitor

158 Figure 1b shows an explicit-synchronization implementation of the bounded queue from Figure 1a 159 that is written by an expert. This implementation uses two distinct locks, putLock and takeLock, 160 to protect the put and take methods respectively. The explicit-synchronization monitor also uses 161 an atomic integer for the count field, transforming reads into get() calls (e.g., line 12) and writes 162 into the appropriate atomic method (e.g., count.getAndIncrement() on line 16). The expert-written 163 monitor performs explicit signaling via condition variables notFull and notEmpty that are associated 164 with putLock and takeLock respectively. When a thread cannot execute one of these operations, it 165 calls await on the appropriate condition variable to block its execution (lines 13 and 26). A thread blocked in put can only be unblocked by a corresponding take that frees up space in the queue. To 166 167 do so, the take must acquire putLock and perform a signal operation on condition variable notFull 168 (lines 33–35). The logic for take is symmetric (lines 19–21).

Although the expert-written version has more locks than a single global-lock implementation, its 169 170 performance will often be better: Introducing two locks allows put and take to execute concurrently, 171 although multiple concurrent puts are still serialized, as are multiple takes. A single global lock 172 would admit no concurrency in this case and would still incur the same synchronization overhead 173 of acquiring and releasing a lock on every method call. The expert implementation mitigates the 174 overhead of having two locks by acquiring locks selectively: take only acquires the putLock if 175 it is possible for there to be a put operation currently blocked waiting for space in the queue, 176 which happens only if the queue was full when take ran (the put/takeLock case is symmetric). This 177 example demonstrates the intricacy of synthesizing fine-grained locking protocols: instead of only 178 minimizing the total number of locks, we must also try to maximize the available concurrency. 179

#### 2.3 Our Approach

Our tool CORTADO automatically synthesizes the efficient explicit-synchronization monitor in Figure 1b given the implicit version from Figure 1a. It does so in three phases: First, it infers when and how signaling operations should take place. Second, it performs static analysis to infer sufficient conditions for the synthesized monitor to be correct. Third, it encodes the synchronization protocol synthesis problem as a MaxSAT instance and uses a model of the MaxSAT problem to generate an explicit-sychronization monitor. Since prior work can already handle the first phase, we only focus on the the latter two phases in the following discussion.

**Granularity**. The granularity of our synthesized locking protocol is at the level of *code fragments*, where each fragment is a single-entry region of code within a single method. For example, the fragments chosen for the blocking queue example are indicated by comments in Figure 1a. Fragments are the indivisible unit of concurrency in our approach: we aim to maximize the number of fragments that can run concurrently, but we do not modify the code within a fragment to introduce extra concurrency (e.g., by removing data races). Hence, the explicit monitor synthesized by our approach acquires and releases locks only at fragment boundaries.

1:4

156

215

216

217

218

219

220

221

222

223 224

225

226

227

228

229

230

231

232

240

241

242 243

245

Static analysis. To ensure correctness of the synthesized monitor, our technique needs to 197 enforce the following three key requirements: 198

- 199 (1) **Data-race freedom:** Fragments that involve a data race must not be able to run concurrently.
- 200 (2) **Deadlock freedom:** Locks must be acquired and released in an order that prevents deadlocks.
- 201 (3) Atomicity: Each monitor operation should *appear* to take place as one indivisible unit. That is, 202 even though the implementation can allow thread interleavings inside monitor operations, the 203 resulting state should be equivalent to one where each method executes truly atomically. 204

Here, the second requirement (i.e., deadlock freedom) does not necessitate any static analysis, as 205 we can prevent deadlocks by imposing a static total order  $\leq$  on locks [Birrell 1989] and ensuring 206 that locks are acquired and released in a manner that is consistent with  $\leq$ . However, in order to 207 ensure data-race freedom and atomicity, we need to perform static analysis of the source code 208 to identify (1) code fragments that have a data race, and (2) interleaving opportunities between 209 code fragments. Since detection of data races is a well-studied problem, the novelty of our static 210 analysis lies in identifying safe interleaving opportunities. Hence, the key question addressed by 211 our analysis is the following: Given a code fragment f executed by thread t, and two consecutive 212 code fragments  $f_1$ ,  $f_2$  executed by a different thread t', is it safe to interleave the execution of f in 213 between  $f_1$  and  $f_2$  while ensuring that monitor operations appear to take place atomically? 214

To answer this question, our method performs a novel static analysis to identify a set of such safe interleavings. For instance, going back to the running example, our analysis determines that it is safe to interleave the execution of fragment 4 in Figure 1a in between fragments 5 and 6 by checking a number of commutativity relations between code fragments. In this instance, since our analysis proves that fragment 4 left-commutes [Lipton 1975] with fragment 5 and right-commutes [Lipton 1975] with 6 and all of its successors, we identify this as a safe interleaving opportunity. On the other hand, our analysis concludes that interleaving fragment 4 in between 1 and 2 is not safe because fragment 4 does not left-commute with fragment 1 - intuitively, this is because fragment 4 can falsify predicate count < queue.length that appears in the waituntil statement of fragment 1.

MaxSAT overview. Once we identify possible data races and safe interleavings via static analysis, we use this information to generate a MaxSAT instance whose solution corresponds to a finegrained synchronization protocol. Specifically, our MaxSAT encoding uses a variable  $h_{f_i}^{l_j}$  to indicate that code fragment  $f_i$  must hold lock  $l_i$  and generates both hard constraints (for correctness) and soft constraints (for efficiency) over these variables. Thus, if the MaxSAT solver returns a model in which variable  $h_{f_i}^{l_j}$  is assigned to true, this means that the synthesized code must acquire lock  $l_j$ prior to executing fragment f<sub>i</sub>. Similarly, our MaxSAT encoding introduces a variable a<sub>fld</sub> indicating that field *fld* should be implemented using an atomic type.

The hard constraints in our MaxSAT encoding correspond to the three correctness requirement 233 mentioned earlier, namely (1) data race prevention, (2) deadlock freedom, and (3) atomicity. On the 234 other hand, soft constraints encode our optimization objective. In what follows, we give a brief 235 overview of the different types of constraints in our encoding, focusing only on constraints that 236 involve lock acquisition variables  $h_{f_i}^{l_j}$ . However, it is worth noting that our technique also generates constraints on atomic variables  $a_{fld}$  and can automatically convert fields to atomic types whenever 237 238 239 doing so is safe and more efficient than introducing a lock.

**Data-race freedom.** Given a pair of code fragments  $(f_i, f_j)$  that have a potential data race according to the static analysis, our MaxSAT encoding introduces hard constraints of the form  $\bigvee_k (h_{f_i}^{l_k} \wedge h_{f_i}^{l_k})$  stating that  $f_i$  and  $f_j$  must share at least one common lock. For example, in Figure 1a, our analysis determines that fragments 4 and 8 cannot run in parallel since they both write to the 244

same memory location count. Thus, the MaxSAT instance contains boolean constraints to make sure that two different threads cannot execute count-- and count++ at the same time.

**Deadlock freedom.** Our approach precludes deadlocks by imposing a total order  $\leq$  on locks. In particular, it enforces that a thread *t* can only acquire lock *l* if *t* does *not* already hold any lock *l'* where l' < l. For example, in Figure 1a, suppose the locking protocol determines that fragments 1 and 2 must hold all locks in sets  $L_1$  and  $L_2$  respectively. Between executing the two fragments, the code will need to acquire all locks in  $L_2 \setminus L_1$ . Hence, we add constraints i < j for every pair of locks  $l_j \in L_2 \setminus L_1$  and  $l_i \in L_1 \cap L_2$  so that those locks can be acquired while respecting the order  $\leq$ .

Atomicity. Our MaxSAT encoding also includes constraints to ensure that monitor operations appear to execute atomically. Suppose that our static analysis determines that a thread cannot safely execute code fragment f in between some other thread's execution of code fragments  $f_1$ and  $f_2$ . To prevent such an unsafe interleaving, we add hard constraints to ensure that fragments f,  $f_1$ , and  $f_2$  all share at least one common lock. For example, since our analysis determines that fragment 4 (count++) cannot be interleaved with any other pair of fragments in the same method put (running concurrently on a different thread), our MaxSAT encoding includes a hard constraint asserting that fragment 4 must share a lock with all other fragments in the put method.

**Soft constraints.** Because the efficiency of the synthesized code depends on both the allowed parallelization opportunities as well as the number of locks, our optimization objective tries to *minimize* the number of locks and *maximize* the number of fragments that can run in parallel. To encode the latter objective, our MaxSAT encoding includes soft contraints asserting that any two parallelizable fragments must *not* share a lock. On the other hand, to encode the former objective, we add a soft constraint stating that no fragment in *m* is holding lock *l*.

**Monitor generation.** A solution of the generated MaxSAT instance determines (a) which fragments should hold which locks, (b) which fields should be implemented using atomic types, and (c) which locks should be associated with which condition variables. Thus, together with the output of the signal placement technique [Ferles et al. 2018], a model of the MaxSAT problem can be automatically translated into the target monitor implementation. For our running example, CORTADO synthesizes precisely the implementation in Figure 1b given the implicit monitor of Figure 1a.

#### 3 PRELIMINARIES

In this section, we describe our source and target languages and define what it means for an explicit synchronization monitor to correctly implement an implicit one.

## 3.1 Background on Monitors

In this work, we assume that all shared resources between threads are handled by a *monitor class* M which consists of fields F and set of operations (methods) O. The fields F constitute the *only* shared state between threads, which can only access shared state by performing one of the monitor operations  $o \in O$ . These operations can be performed by an arbitrary, yet fixed, number of threads, and locations reachable through arguments are assumed to be thread-local. We represent each thread by a unique identifier from set  $T \subseteq \mathbb{N}$ , and we model memory locations using *access paths* (AP) [Landi and Ryder 1992] of the form  $\pi = v(.f)*$ , consisting of a base variable v optionally followed by a finite sequence of field accesses. We also assume that a special this variable stores the memory location of the monitor object.

Definition 3.1. (Monitor State). A monitor state  $\sigma : T \times AP \to \mathbb{N}$  is a mapping from pairs  $(t, \pi)$  (where *t* is a thread identifier and  $\pi$  an access path) to a value.

			Monitor $M$	::= monitor $M \{(fld \mid sync \mid m)^*\}$
295 296	Monitor M	$u = monitor M \left( \left( f d \mid m \right)^* \right)$	Field <i>fld</i>	$::= \tau f := e$
297 298 299 300	Field <i>fld</i> Method <i>m</i>	$::= \pi (\vec{v}) \{ ccr^* \}$	Sync sync	<pre>::= Lock l := new Lock()   CondVar cv := l.newCondVar()   Atomic[τ] a := e</pre>
301	CCR ccr	<pre>::= waituntil(p);s</pre>	Method <i>m</i>	$::= m(\vec{v}) \{ ccr^* \}$
302 303	Stmt s	::= skip   v := e   v.f := e	CCR ccr	$::= (ls)^*$
304 305		$  v.m(\vec{e})   [if (e)]? goto 1   ls_1; ls_2$	Stmt s	$ ::= skip   v := e   v.f := e    v.m(\vec{e})   [if (v)]? goto 1 $
306	LStmt ls	::= 1:? s		$  ls_1; ls_2$
307 308	(a) Implicit	-synchronization monitor language.		$a_{pre} := a.update(\chi \chi.e)$
309			LStmt <i>ls</i>	::= 1:? s
310			(b) Explicit-	synchronization monitor language.

Fig. 2. Source & target languages. We use e and p for expressions and predicates respectively.

#### Source Language 3.2

311

312

313

322

323

324

325

326

327 328

329

330

331

332

333

334 335

336

Our source language, presented in Figure 2a, corresponds to implicit synchronization monitors 314 without explicit locking or signaling. The body of each monitor operation consists of a sequence of 315 so-called Conditional Critical Regions (CCRs) [Hoare 1971], which in turn consist of a waituntil 316 statement followed by one or more regular non-blocking statements. We refer to the predicate of 317 the waituntil statement of a CCR as its guard and to the rest of the statements as its body. A thread 318 executes the body of the CCR atomically if its guard evaluates to true; otherwise it suspends its 319 execution and exits the monitor until the predicate becomes true. More formally, the semantics of 320 our source language are defined via the notion of an *implicit monitor history*: 321

Definition 3.2. (Implicit monitor history). Given a set of threads interacting with each other through monitor  $M_s = (F, O)$ , an *implicit monitor history*  $h_i$  is a sequence  $(ccr_1, t_1) \dots (crr_n, t_n)$ where each  $ccr_i$  is a CCR of  $M_s$  and  $t_i$  is a thread identifier.

Given history  $h_i$ , we define an *argument mapping*  $v_i$  to be a list whose *i*'th element maps formal parameters of  $Method(ccr_i)$  to their actual value for each event  $(ccr_i, t_i)$  in  $h_i$ .

Definition 3.3. (Implicit monitor semantics). Given a monitor  $M_s$ , initial state  $\sigma$ , and monitor history  $h_i$  with argument mapping  $v_i$ , the operational semantics of M is defined using a judgment  $M_{s} \vdash (h_{i}, v_{i}, \sigma) \Downarrow \sigma'$  indicating that the new monitor state is  $\sigma'$  after executing  $h_{i}$  on state  $\sigma$ .

Because our source language is very similar to the one used in Ferles et al. [2018], we omit a formal definition of the operational semantics. Following that work, we also consider an implicit history to be valid only if it respects the program order of the input monitor.

## 3.3 Target Language

Figure 2b presents the language of explicit-synchronization monitors. The overall structure of 337 this target language is similar to the source language but with a few important differences. First, 338 an explicit monitor contains locks, conditional variables, and atomic fields, collectively referred 339 to as synchronization variables. Second, CCRs in the target language do not contain waituntil 340 statements; instead, the logic of a waituntil statement is implemented by calling methods on the 341 appropriate condition variable. We assume that synchronization variables support all the standard 342

class M { class M {  $\text{ int } x \ = \ \emptyset, \ y \ = \ \emptyset, \ z \ = \ \emptyset;$ int x = 0, y = 0, z = 0; void foo() { x++; y++; } Lock 11 = new Lock(), 12 = new Lock(); 344 void bar() { z++; } } void foo() { l1.lock(); x++; y++; l1.unlock(); } 345 void bar() { 12.lock(); z++; 12.unlock(); } } 346 (a) A simple implicit monitor. (b) An explicit monitor implementation of Figure 3a. 347  $h_i = (foo, t_1)(bar, t_2)$ 348  $h_e = (11.1ock(), t_1)(x++, t_1)(y++, t_1)(11.unlock(), t_1)(12.1ock(), t_2)(z++, t_2)(12.unlock(), t_2)(z++, t_2)(z$ 349  $h'_{e} = (11.1ock(), t_{1})(x++, t_{1})(12.1ock(), t_{2})(y++, t_{1})(z++, t_{2})(11.unlock(), t_{1})(12.unlock(), t_{2})$ 350 (c) Examples of implicit and explicit histories. 351 352 Fig. 3. A simple implicit monitor and its explicit implementation. 353 synchronization operations present in modern concurrent languages (e.g., await, signal, signalAll, 354 etc.). Finally, our target language contains a special update statement for performing updates on 355 atomic fields: it takes as argument an atomic field a and a unary function f and updates the value of 356 a atomically as f(a). For instance, the statement  $c_{pre} := c.update(\lambda \chi, \chi + 1)$  atomically increments 357 c by one and stores the value of c before the update in  $c_{pre}$ . 358 359 Definition 3.4. (Explicit monitor history). Given a set of threads executing in monitor  $M_t$  = 360 (F, O), an explicit monitor history  $h_e$  is a sequence  $(s_1, t_1) \dots (s_n, t_n)$  where each  $s_i$  is a (non-361 composite) statement of a monitor operation  $o \in O$  and  $t_i$  is a thread identifier. 362 Leveraging the same notion of *argument mappings* defined in Section 3.2, we define explicit 363 monitor semantics as follows: 364 365 Definition 3.5. (Explicit monitor semantics). Given a monitor  $M_t$ , initial state  $\sigma$ , and monitor 366 history  $h_e$  with argument mapping  $v_e$ , the operational semantics of  $M_t$  is defined using a judgment 367  $M_t \vdash (h_e, v_e, \sigma) \downarrow \sigma'$  indicating that the new state is  $\sigma'$  after executing  $h_e$  on initial state  $\sigma$ . 368 The full operational semantics of our target language is given in Appendix C. 369 370 3.4 **Relating Implicit and Explicit Histories** 371 In order to formalize the correctness of our approach, we need to relate an implicit history  $h_i$  of a 372 source monitor  $M_s$  with an explicit history  $h_e$  of its corresponding target version  $M_t$ . Because every 373 history of an implicit monitor  $M_s$  induces a corresponding history of its explicit version  $M_t$ , we 374 define an operation called that Expand that "translates" an implicit history to an explicit one. That 375 is, given an implicit history  $h_i$  with argument mapping  $v_i$  and state  $\sigma$ , Expand<sub>M</sub>  $(h_i, v_i, \sigma)$  returns a 376 pair  $(h_e, v_e)$ , where  $h_e$  is a history of  $M_t$  containing all statements executed by  $h_i$  under initial state 377  $\sigma$  and  $v_e$  is the argument mapping for  $h_e$ . 378 *Example 3.6.* Consider the implicit monitor of Figure 3a and its explicit counterpart in Figure 3b. 379 For histories  $h_i$  and  $h_e$  from Figure 3c we have Expand<sub>*M*<sub>e</sub></sub> $(h_i, v_i, \sigma) = (h_e, v_e)$  for some  $v_i, v_e$ . 380 381 Using this Expand operation, we can classify explicit histories as being sequential or interleaved: 382 Definition 3.7. (Sequential history) Let  $M_t$  be an explicit monitor implementation of  $M_s$ . We 383 say that an explicit history  $h_e$  of monitor  $M_t$  with argument mapping  $v_e$  is sequential iff there exist 384 a history  $h_i$  of  $M_s$ , argument mapping  $v_i$ , and initial state  $\sigma$  such that Expand<sub>M<sub>i</sub></sub> $(h_i, v_i, \sigma) = (h_e, v_e)$ . 385 386 In other words, a sequential history corresponds to an execution in which statements of the 387 explicit monitor are not interleaved between threads. 388 *Example 3.8.* Going back to Figure 3c, history  $h_e$  is sequential but  $h'_e$  is not. 389 Next, we introduce the notion of *well-formed histories*, which, intuitively, respect the program 390 order of the original implicit monitor: 391

Kostas Ferles, Benjamin Sepanski, Rahul Krishnan, James Bornholt, and Isil Dillig

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

392

1:8

1:9

Definition 3.9. (Well-formed history) Let  $\Pi(h, t)$  be the projection of h onto thread t (i.e., it filters out all elements of h not involving thread t). We say that a history  $h_e$  of  $M_t$  is well-formed iff, for every thread t, there exists sequential histories  $h_e^1, \ldots, h_e^n$  such that  $\Pi(h_e, t) = h_e^1 \cdots h_e^n$ .

Intuitively, well-formed histories respect program dependence in the original monitor for every thread. By definition, every sequential history is also well-formed. In the remainder of this paper, we implicitly mean *well-formed* explicit history whenever we refer to an explicit history.

*Example 3.10.* Histories  $h_e$ ,  $h'_e$  from Figure 3c are both well-formed. However, the following history is not well-formed because it does not respect program order: (12.unlock(), t)(12.lock(), t)

Definition 3.11. (Interleaved history) We say that a history  $h_e$  of  $M_t$  is interleaved iff it is (1) well-formed and (2) not sequential.

*Example 3.12.* History  $h'_e$  from Figure 3c is interleaved.

Next, we define what it means for an explicit history to *simulate* an implicit history.

Definition 3.13. (Simulation relation). Let  $M_t$  be an explicit version of implicit monitor  $M_s$ . We say that an explicit history  $h_e$  of  $M_t$  with argument mapping  $v_e$  simulates  $(h_i, v_i)$  of  $M_s$  on input  $\sigma$ , denoted  $(h_e, v_e) \sim (h_i, v_i)$ , if there exist sequential history  $h'_e$  and  $v'_e$  such that:

(1)  $\forall t. \Pi(h_e, t) = \Pi(h'_e, t)$  and (2) Expand<sub>M<sub>e</sub></sub> $(h_i, v_i, \sigma) = (h'_e, v'_e)$ .

In other words,  $h_e$  simulates a history of the original monitor if it is a (well-formed) permutation of some sequential history  $h'_e$  of the explicit monitor  $M_t$ .

*Example 3.14.* Going back to Figure 3c, we have  $(h'_e, v') \sim (h_i, v)$  for some v, v'.

### 3.5 Correctness of Explicit-Synchronization Monitors

Using the concepts introduced in the previous section, we now formalize what it means for an explicit monitor to *correctly implement* an implicit one.

Definition 3.15. (State equivalence) Let  $\sigma$  be a program state of an implicit monitor  $M_s$  and  $\sigma'$  that of an explicit monitor  $M_t$ . We say that  $\sigma$  and  $\sigma'$  are equivalent modulo  $M_s$ , denoted  $\sigma \equiv_{M_s} \sigma'$ , iff for all  $(t, \pi)$  in the domain of  $\sigma$ , we have  $\sigma(t, \pi) = \sigma'(t, \pi)$ 

Intuitively, this notion of equivalence between two monitor states ignores any additional synchronization fields and local variables introduced by translating M to an explicit-synchronization monitor. Finally, we can define the correctness of an explicit monitor as follows:

Definition 3.16. (Correctness) We say that an explicit monitor  $M_t$  correctly implements an implicit monitor  $M_s$ , denoted as  $M_s \sim M_t$ , iff for all input states  $\sigma_s$ ,  $\sigma_t$  s.t.  $\sigma_s \equiv_{M_s} \sigma_t$ , we have:

(1) 
$$\forall h_i, v_i. \ M_s \vdash (h_i, v_i, \sigma_s) \Downarrow \sigma'_s \Longrightarrow \left( M_t \vdash (\text{Expand}_{M_t}(h_i, v_i, \sigma_s), \sigma_t) \downarrow \sigma'_t \land \sigma'_s \equiv_{M_s} \sigma'_t \right)$$

 $(2) \ \forall h_e, v_e. \ M_t \vdash (h_e, v_e, \sigma_t) \downarrow \sigma'_t \Longrightarrow (\exists h_i, v_i. \ (h_e, v_e) \backsim (h_i, v_i) \land M_s \vdash (h_i, v_i, \sigma_s) \Downarrow \sigma'_s \land \sigma'_s \equiv_{M_s} \sigma'_t)$ 

The first correctness condition simply states that  $M_t$  does not eliminate any feasible behaviors of  $M_s$ . The second condition, on the other hand, states that every feasible history of  $M_t$  simulates *some* implicit history that results in the same state. Intuitively, this means that all statement interleavings allowed by  $M_t$  provide the illusion that all operations of  $M_s$  are executed atomically.

#### 4 MAIN ALGORITHM

In this section, we present our main synthesis algorithm. Specifically, Section 4.1 introduces some
 preliminary definitions and proves an NP-completeness result to justify the reduction to MaxSAT.
 Then, Section 4.2 presents the high-level algorithm, Section 4.3 presents the static analysis for
 inferring safe interleavings, and Sections 4.4 presents the details of the MaxSAT encoding.

#### **Fragment Dependency Graphs and NP-Completeness** 4.1 442

443 Our main synthesis algorithm is parametrized over a partitioning of the input monitor into code 444 fragments, where each code fragment defines a unit of computation that we need to assign locks 445 to. In this section, we clarify our assumptions about these code fragments and prove the NP-446 completeness of the problem for a given choice of partition.

First, to define what we mean by a valid partition, we represent each method of the monitor as a 448 standard control-flow graph (CFG), where each atomic statement belongs to its own basic block. 449 Given a control-flow graph G and node n, we write Preds(G, n) to indicate the predecessor nodes 450 of n in G and Succs(G, n) to indicate its successors. Then, a valid partition of a method into code 451 fragments is defined as follows: 452

Definition 4.1. (Partition) Let G = (V, E) be the CFG representation of a method. Then, a partition of this method is a set of CFGs  $\{G_1, \ldots, G_n\}$  with  $G_i = (V_i, E_i)$  such that:

- (1)  $V = \bigoplus_{i=1}^{n} V_i$  and  $E_i = E \cap (V_i \times V_i)$
- (2) For every  $G_i$ , there is at most one node  $n \in V_i$  such that  $Preds(G, n) \nsubseteq V_i$
- (3) Every waituntil(p) statement must belong to its own  $G_i$  i.e., if a node  $n \in V$  is a waituntil statement, then there exists a  $G_i = (\{n\}, \emptyset)$

Intuitively, a partition is a set of sub-CFGs such that (1) these sub-CFGs cover all nodes of the original CFG, (2) each sub-CFG has a unique entry node, and (3) waituntil statements belong to their own sub-CFG. We refer to the code snippet represented by each sub-CFG as a code fragment and define a notion of *fragment dependency graph (FDG)* as follows:

Definition 4.2. (FDG) Given a method m with CFG G = (V, E) and a partition of G into  $\{G_1, \ldots, G_n\}$ , a fragment dependency graph (FDG) is a directed acyclic graph (V', E') such that (1) every  $f_i \in V'$  is the code fragment associated with  $G_i$ ; (2) there is an edge  $(f_i, f_i) \in E'$  iff there is an edge in *G* from any exit node of  $G_i$  to the entry node of  $G_i$ .

Example 4.3. Figure 4 presents the FDG of method take for the partition in Figure 1a,

469 Observe that we require the FDG to be 470 acyclic, so some partitions do not give rise to valid FDGs. In the rest of this paper, we assume that partitions obey this restriction so that all 473 cycles are contained within individual nodes of 474 the FDG. We also lift this notion of FDG from





individual methods to entire monitors in the obvious way (i.e., union of all FDGs). As we will see in the next section, our synthesis algorithm operates over FDG representations of monitors. Next, we state the following NP-completeness result to justify our MaxSAT encoding:

THEOREM 4.4. (NP-Completeness) Let  $\mathcal{G} = (V, E)$  be an FDG of monitor M, and let  $\Pi \subseteq V \times V$  be a set of fragment pairs that can run in parallel. Then, deciding whether there exists a synchronization protocol with at most k locks and that allows all pairs in  $\Pi$  to run in parallel is NP-Complete.

PROOF. By reduction from the edge clique cover problem [Michael and Quint 2006]. The proof can be found in Appendix B. 

## 4.2 Synthesis Algorithm

In this section, we describe our core synthesis procedure, which is summarized in Figure 5. At a high 486 level, the SYNTHESIZEMONITOR algorithm consists of the following steps. First, it uses the technique 487 of Ferles et al. [2018] to infer signaling operations (line 4). This yields a partially concretized 488 monitor M' with signaling operations but no locking. Next, it constructs an FDG representation 489

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

447

453

454

455

456

457

458 459

460

461

462 463

464

465

466

467

468

471

472

475

476

477

478

479

480

481

482

483

484

485

490

1:10

1:11

492	1:	procedure SynthesizeMonitor(M)
493	2:	<b>input:</b> <i>M</i> : an implicit-synchronization monitor.
494	3:	output: a semantically equivalent explicit-synchronization monitor.
495	4:	$M' \leftarrow PlaceSignals(M)$ $\triangleright$ Use technique of Ferles et al. to infer signaling operations
496	5:	$\mathcal{G} \leftarrow \text{ConstructFDG}(M')$
497	6:	$\mathcal{N}_u \leftarrow \text{ComputeMaxLocks}(\mathcal{G})$
498	7:	$\mathcal{S} \leftarrow \text{StaticAnalyze}(\mathcal{G})$
499	8:	$opt \leftarrow -1$
500	9:	for $i \in [1, \mathcal{N}_u]$ do
501	10:	$(\mathcal{H}, \mathcal{S}) \leftarrow \text{MaxSatEncoding}(M, \mathcal{G}, \mathcal{S}, i)$
502	11:	$(p, v, timeout) \leftarrow Solve(\mathcal{H}, S)$
503	12:	if timeout $\lor (v \le opt)$ then break
504	13:	$(best, opt) \leftarrow (p, v)$
505	14.	return lntrumont(hest C M')
506	14:	return intrument(best, g, M)
507		Fig. 5. Main Synthesis Algorithm.
508		

of the resulting monitor M' as defined in Section 4.1 (line 5). Third, it infers an *upper bound*  $N_u$ on the maximum number of locks that the synthesized code should use (line 6). Then, it statically analyzes the FDG to infer requirements that the synthesized code needs to obey (line 7) and uses the results of the previous steps to generate the MaxSAT encoding (line 11). Finally, it instruments M' (line 14) using the synchronization protocol inferred using MaxSAT. Since the most involved aspects of this algorithm are the MaxSAT encoding and inference of safe interleavings, we defer a detailed discussion of these to the next two subsections and focus on the rest.

523 *Iterative exploration of lock count.* As mentioned above, our synthesis algorithm conceptually 524 reduces the protocol synthesis problem to MaxSAT and uses an off-the-shelf solver to maximize 525 our optimization objective. To achieve this goal, one option is to generate the MaxSAT encoding 526 based on the maximum possible locks (obtained via the call to ComputeMaxLocks) and then let 527 the solver figure out the optimal number of locks to use. However, in practice, such an approach 528 does not scale because the size of the encoding increases with respect to the maximum number of 529 locks allowed. That is, for many realistic problems, the MaxSAT solver fails to terminate within 530 a reasonable time limit if we generate the encoding based on the maximum possible locks. Thus, 531 instead of directly generating a very large MaxSAT formula up front, our SynthesizeMonitor 532 procedure enters a loop (lines 9-13) wherein it gradually increases the maximum number of locks 533 allowed (and hence the size of the MaxSAT encoding). If we get to a point where the MaxSAT solver 534 starts timing out (indicated by boolean variable called *timeout*) or we fail to increase the objective 535 value despite using a larger upper bound on locks (see line 12), then the procedure terminates 536 with the best sychronization policy found so far. While this strategy does not guarantee global 537 optimality, it is much more practical than the alternative. 538

539

491

540 *Signaling operations.* Our synthesis algorithm uses 541 an auxiliary procedure called PlaceSignals [Ferles et al. 542 2018] which yields a monitor M' that belongs to an in-543 termediate language that is identical to our source lan-544 guage (Figure 2a) except that it contains explicit sig-545 naling operations. Specifically, this intermediate lan-546 guage contains two additional signaling directives: (1)

```
void put(Object o) {
 waituntil(count < queue.length):
 boolean wasEmpty = count == 0;
 queue[last] = o;
 last = (last + 1) % queue.length;
 count++:
 broadcast(count == 0, wasEmpty);
```

Fig. 6. Method put with explicit signals.

547 signal (p, c) which notifies a single thread that is blocked on predicate p if condition c holds, and 548 (2) broadcast (p, c) which notifies all threads blocked on p if c holds. Figure 6 shows the result of 549 calling PlaceSignals on the put procedure from Figure 1a. 550

FDG construction. Recall that an FDG is a generalized version of a control-flow graph where 551 nodes are code fragments rather than basic blocks, and each code fragment is a unit of computation 552 that our algorithm should assign locks to. Since there can be many ways to partition a given 553 method into code fragments, the ConstructFDG procedure invoked at line 5 of Figure 5 implements 554 a particular heuristic for partitioning a method into code fragments. In particular, the more the 555 number of code fragments, the more the parallelization opportunities; thus, our ConstructFDG procedure tries to maximize the number of code fragments while maintaining the invariant that the FDG is acyclic and that each code fragment must have a unique entry point (see Section 5). 558

Computing upper bound on locks. Because the MaxSAT encoding assumes a fixed number 559 560 of locks, our synthesis algorithm calls the ComputeMaxLocks procedure at line 6 to compute an upper bound on the number of locks needed. Given an FDG  $\mathcal{G} = (V, G)$ , the key idea behind this 562 procedure is to construct a so-called *conflict graph*  $G_C = (V, E_C)$  where (f, f') is in  $E_C$  iff fragments f and f' have a data race. Since it can be shown that the optimal solution to our problem is an *edge* 563 clique cover [Michael and Quint 2006] of this conflict graph (see Appendix), we can use known 564 theorems (e.g., Mantel's theorem, Alon [1986] etc.) to obtain an upper bound on the number of 565 locks needed without having to solve an NP-complete problem.<sup>1</sup>

Static analysis. Recall from Section 2 that the constraints in our MaxSAT encoding utilize information obtained via static analysis. Thus, line 7 of Figure 5 statically analyzes the input monitor to obtain the following three pieces of information:

- Atomic fields  $\mathcal{F}$ : One of the goals of the analysis is to infer a set of fields that could *potentially* be implemented using Atomic types. Thus, our static analysis checks whether (a) a field of type T has a corresponding AtomicT version, and (b) whether all updates to this field can be implemented using one of the methods provided by AtomicT.
- Data races R: The second goal of our static analysis is to identify pairs of fragments that would 575 have a data race if they do not use a shared lock. Thus, given a pair of fragments (f, f'), our 576 static analysis checks whether f writes to a memory location l that is accessed in f'. 577
  - Interleaving opportunities *I*: Finally, a third key goal of the analysis is to identify safe interleaving opportunities between fragments. Since this aspect of the analysis is quite involved, we discuss it in the next subsection.

581 MaxSAT encoding. As mentioned in Section 2, our MaxSAT encoding uses two types of boolean 582 variables, namely (1)  $h_{f_i}^{l_j}$  indicating that fragment  $f_i$  must hold lock  $l_j$  and (2)  $a_f$  indicating that field 583 f should be converted to atomic. Hence, a model of the MaxSAT problem can be easily converted 584 to a so-called *locking protocol*  $(\mathcal{L}, \mathcal{A}, \mathcal{P})$ , where  $\mathcal{L}$  is an assignment from fragments to a set of 585 locks,  $\mathcal{A}$  is a set of fields that should be implemented using atomic types, and  $\mathcal{P}$  is a mapping from 586

556

557

561

566 567

568

569

570

571

572

573

574

578

579

<sup>&</sup>lt;sup>1</sup>In our implementation, we use multiple upper bounds using known theorems and return the best one.

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

waituntil guards to locks. In particular, we have  $l_j \in \mathcal{L}(f_i)$  if and only if  $h_{f_i}^{l_j}$  is assigned to true in 589 590 the model returned by the MaxSAT solver, and we have  $fld \in \mathcal{A}$  if  $a_{fld}$  is assigned to true. Due to 591 the constraints in our MaxSAT encoding, it is similarly easy to derive  $\mathcal{P}$ : because our encoding 592 ensures that every occurrence of a waituntil(p) statement is protected by the same set of locks S, 593 we associate one of the locks l in S with the condition variable introduced for predicate p.<sup>2</sup>

Instrumentation. The last step of our algorithm is to synthesize the explicit-synchronization monitor via the Instrument procedure invoked at line 14. Given a synchronization protocol  $(\mathcal{L}, \mathcal{A}, \mathcal{P})$ , the Instrument procedure performs the following steps:

- (1) First, it introduces all the synchronization fields (locks, condition variables, and atomic fields) 598 599 that appear in the protocol.
- (2) It converts every update to an atomic field to the corresponding atomic update statement. 600
- (3) Finally, it introduces all the necessary locking and signaling operations to implement the 601 synthesized synchronization protocol. 602

We refer the interested reader to Appendix D for more details on the instrumentation.

THEOREM 4.5. (Correctness) Given an implicit-synchronization monitor M in the language of Figure 2a, if SynthesizeMonitor(M) returns M', then we have  $M \sim M'$ .

PROOF. Can be found in Appendix E.

#### Analysis to Identify Safe Interleavings 4.3

We now describe how to infer safe interleaving opportunities between threads while ensuring that monitor operations *appear* to take place atomically. Given a fragment dependency graph  $\mathcal{G} = (V, E)$ for a monitor M, an interleaving opportunity (or interleaving for short) is a pair (v, e) where  $v \in V$ is a code fragment of M and  $e = (v_1, v_2) \in E$  is an edge of the FDG. Intuitively, such an interleaving is safe if some thread can execute v in between some other thread's execution of  $v_1$  and  $v_2$  without violating atomicity. The goal of our static analysis is to identify a set I of such safe interleavings. In what follows, we formalize *safe interleavings* and describe an analysis for identifying them.

*Formalizing safe interleavings.* To formalize the notion of safe interleaving, we need to keep 619 track of which fragments of the monitor were executed in what order. For this purpose, given an FDG  $\mathcal{G} = (V, E)$  of M, we define a fragmented monitor  $M_{\mathcal{G}}$  to be the same as M except that every fragment in  $\mathcal{G}$  is placed in its own method. Observe that histories of  $M_{\mathcal{G}}$  encode all possible interleavings of fragments in  $\mathcal{G}$ . In this sense, histories of  $M_{\mathcal{G}}$  are similar to explicit monitor histories but are slightly higher level in that they allow interleavings between fragments rather than atomic statements. Thus, we adapt the same notions of sequential, well-formed, and interleaved histories from Section 3.3 to fragmented monitors, as illustrated by the following examples.

*Example 4.6.* Given monitor M from Figure 1a, its fragmented version  $M_G$  splits put and take into four different methods, each named put<sub>i</sub> and take<sub>i</sub>. Given history h = (take, t) and initial state  $\sigma$  with a non-empty queue, we have

 $\mathsf{Expand}_{M_G}(h, v, \sigma) = ((take_1, t)(take_2, t)(take_3, t)(take_4, t), v')$ 

where  $take_1, \ldots, take_4$  denote fragments 5-8 in Figure 1a and v, v' are empty argument mappings. 632

594

595

596

597

603

604

605

606 607

608 609

610

611

612

613

614

615

616

617 618

620

621

622

623

624

625

626

627

628

629 630

631

633

<sup>&</sup>lt;sup>2</sup>Specifically, when choosing which lock l in set S to designate as the representative, we choose the smallest lock in S634 according to the total order. Because all locks held by a thread must be released before it blocks on a condition variable and 635 must be acquired after it gets notified (with method await releasing and acquiring l internally), choosing the smallest lock prevents deadlocks. 636

*Example 4.7.* In the example above,  $\text{Expand}_{M_{\mathcal{G}}}(h, v, \sigma)$  is both sequential and well-formed. However,  $(take_1, t)$ ,  $(take_2, t)$  is not well-formed because it does not involve all four methods, and  $(take_1, t)$ ,  $(take_3, t)$ ,  $(take_2, t)$ ,  $(take_4, t)$  is also not well-formed because it executes  $take_3$  before  $take_2$ . Finally, the following history is an interleaved (and, by definition, well-formed) history where threads t and t' execute method take and put respectively:

$$h_{\mathcal{G}} = (put_1, t)(put_2, t)(put_3, t)(take_1, t')(put_4, t)(take_2, t')(take_3, t')(take_4, t')$$
(1)

Furthermore, for this history we have  $(h_{\mathcal{G}}, v_{\mathcal{G}}) \sim ((put, t)(take, t'), v)$  for some  $v_{\mathcal{G}}$  and v. That is,  $h_{\mathcal{G}}$  simulates a history of M where thread t executes method put and t' executes take.

Definition 4.8. (Interleaving) Given an FDG  $\mathcal{G} = (V, E)$  for monitor M, an *interleaving* is a pair (v, e) where  $v \in V$  and  $e \in E$ . Furthermore, given a history h of fragmented monitor  $M_{\mathcal{G}}$ , we write  $X(h_{\mathcal{G}})$  to denote the set of all interleavings that occur in h.

*Example 4.9.* For the history  $h_G$  from Eq. 1, we have:

 $X(h_{\mathcal{G}}) = \{(take_1, (put_3, put_4)), (put_4, (take_1, take_2))\}$ 

This is the case because this history executes  $take_1$  in between consecutive fragments  $put_3$  and  $put_4$  of some other thread. Similarly, we have  $(put_4, (take_1, take_2)) \in \chi(h_G)$  because it executes  $put_4$  in between  $take_1$  and  $take_2$ .

Definition 4.10. (Safe interleavings). Let  $\mathcal{G}$  be an FDG of monitor M. We say that a set of interleavings S is safe, if for every input state  $\sigma$  and every interleaved history  $h_{\mathcal{G}}$  of  $M_{\mathcal{G}}$  we have:

If  $X(h_{\mathcal{G}}) \subseteq S$  and  $M_{\mathcal{G}} \vdash (h_{\mathcal{G}}, v_{\mathcal{G}}, \sigma) \Downarrow \sigma'$  then  $\exists h, v. (h_{\mathcal{G}}, v_{\mathcal{G}}) \backsim (h, v)$  and  $M \vdash (h, v, \sigma) \Downarrow \sigma'$ 

In other words, a set of interleavings *S* is safe if for every interleaved history of  $h_{\mathcal{G}}$  whose interleavings are a subset of *S* we can prove that  $h_{\mathcal{G}}$  leads to the same final state as *some* history *h* of *M* where *h* simulates  $h_{\mathcal{G}}$ . This definition essentially lifts the second correctness criterion of Definition 3.16 to a fragmented monitor.

**Inferring Safe Interleavings.** We now turn our attention to the problem of *inferring* safe interleavings. Given a monitor M and its FDG  $\mathcal{G} = (V, E)$ , our goal is to find a set  $I \subseteq V \times E$  such that all interleavings in I are safe. However, a key challenge is that the space of all safe interleavings is exponential (i.e., the power set of  $V \times E$ ), so, even if we had a procedure for checking whether some set I is safe, enumerating all candidates would be computationally intractable.

To overcome this challenge, we introduce the notion of *strong safety* that allows us to build I iteratively. In particular, note that if  $S_1$  and  $S_2$  are both safe interleaving sets according to Definition 4.10, it may *not* be the case that  $S_1 \cup S_2$  is also a safe interleaving. However, to build I incrementally, we need a notion of safe interleaving that is closed under union. For this purpose, we introduce a notion of *strong safety* for a single interleaving (v, e). Since strongly safe interleavings enjoy the property of being closed under union, this notion lends itself to a computationally feasible technique for computing safe interleaving sets. In the remainder of this section, we define strong safety and present our static analysis for computing safe interleaving sets. Towards this goal, we first introduce the notions of *left* and *right commutativity* for our context:

Definition 4.11. (Left/Right Commutativity). Given fragments v and v', we say that v left commutes with v', denoted LeftCommute(v, v'), iff, whenever  $M_{\mathcal{G}} \vdash ((v', t')(v, t), v, \sigma) \Downarrow \sigma'$  holds, so does  $M_{\mathcal{G}} \vdash ((v, t)(v', t'), v, \sigma) \Downarrow \sigma'$ . Conversely, v right commutes with v', denoted RightCommute(v, v'), iff  $M_{\mathcal{G}} \vdash ((v, t)(v', t'), v, \sigma) \Downarrow \sigma'$  implies  $M_{\mathcal{G}} \vdash ((v', t')(v, t), v, \sigma) \Downarrow \sigma'$ .

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:14

In other words, a fragment v left commutes with v' if, whenever v executes just after v', the 687 resulting state is the same as if v had executed just before v'. For example,  $f_4$  (i.e. count++) in Figure 1a 688 left-commutes with  $f_5$  since increasing count right after waituntil(count > 0) is equivalent to 689 increasing count just before waituntil(count >0). That is, assuming that waituntil(count>0) was 690 not blocked before executing count++, then it will still not be blocked after executing count++. 691 However,  $f_4$  does not left-commute with  $f_1$ : when count equals queue.length - 1, incrementing 692 count just after waituntil (count < queue.length) is *not* equivalent to incrementing queue.length 693 before the waituntil statement. That is, if waituntil(count < queue.length) did not block before 694 executing count++, we cannot guarantee that it also does not block after executing count++. 695 Next, we use this notion of left and right commutativity to define strong safety: 696

<sup>697</sup> Definition 4.12. (Strong safety). Let  $\mathcal{G} = (V, E)$  be an FDG for monitor M, and let  $E^*$  denote the <sup>698</sup> reflexive transitive closure of E. We say that an interleaving (v, e), where  $e = (v_s, v_t)$ , is strongly <sup>699</sup> safe if the following conditions are satisfied:

- (1)  $\forall v^-.(v^-, v_s) \in E^* \Longrightarrow LeftCommute(v, v^-)$
- (2)  $\forall v^+.(v_t, v^+) \in E^* \Longrightarrow RightCommute(v, v^+)$

That is, an interleaving (v, e) is said to be *strongly safe* if we can prove that fragment v left 703 704 commutes with every possible predecessor of  $v_s$  and that it right commutes with every possible successor of  $v_t$ . To see why these conditions imply safety, recall that a set of interleavings S is safe 705 706 if, for any history  $h_G$  whose interleavings are a subset of S, we can find some (sequential) history 707 of the original monitor that simulates  $h_G$ . Assuming S contains only strongly safe interleavings, we can create such a sequential history by "removing" interleavings one at a time from  $h_G$ . For instance, 708 709 let  $\chi = (v, (v_s, v_t)) \in S$  be an interleaving that occurs in  $h_G$ . Since  $\chi$  is strongly safe, we can always obtain an equivalent history  $h'_G$  that has strictly less interleavings than  $h_G$  by commuting v past 710 711 either every successor of  $v_t$  or every predecessor of  $v_s$  that appears in  $h_G$ . 712

*Example 4.13.* For the monitor from Figure 1a, we can show that every interleaving (v, e) where vbelongs to method take and edge e belongs to method put (and vice versa) is strongly safe. However, none of the interleavings where v and e belong to the same method are strongly safe. Finally, because both of the interleavings of the history  $h_{\mathcal{G}}$  from Eq. 1 are strongly safe, we can derive a sequential history that simulates history  $h_{\mathcal{G}}$  by swapping  $(take_1, t')$  with  $(put_4, t)$ .

We now state a key theorem that underlies the correctness of our approach:

THEOREM 4.14. Let G be an FDG and let  $\chi_1, \ldots, \chi_n$  be strongly safe interleavings. Then,  $\{\chi_1, \ldots, \chi_n\}$  satisfies Definition 4.10 (i.e., is a safe interleaving set for G).

PROOF. Can be found in Appendix E.

Static analysis algorithm. Finally, we conclude this section by presenting our static analysis 724 algorithm (shown in Figure 7) for computing a set I of safe interleavings. At a high level, this 725 algorithm identifies which (v, e) pairs are strongly safe and then returns their union, which by 726 Theorem 4.14, corresponds to a safe interleaving set. To check whether an interleaving (v, e) (for 727  $e = (v_s, v_t)$  is strongly safe, we must check if v left commutes with each predecessor of  $v_s$  and 728 right commutes with each successor of  $v_t$ . As shown in the LEFTCOMMUTE procedure, we reduce 729 the verification of left commutativity to the problem of verifying a Hoare triple. In particular, 730 given fragments v, v', we generate a code snippet  $S_L$ ;  $S_R$  where (1)  $S_L$  is an alpha-renamed version 731 of v'; v with waituntil's replaced by assume statements, and (2)  $S_R$  is an alpha-renamed version 732 of v, v' with waituntil's replaced by assert statements. Note that we turn waituntil's in  $S_L$  into 733 assumes because the definition of left commutativity assumes that v'; v has terminated. On the 734

735

701

702

718

719

720

721

722

723

736		
737	1:	procedure FindSafeInterleavings( $\mathcal{G}$ )
738	2:	<b>input:</b> An FDG representation $\mathcal{G} = (V, E)$ of monitor <i>M</i>
739	3:	<b>output:</b> A set $I$ of all safe interleavings
740	4:	$I \leftarrow \emptyset$
741	5:	for $v \in V$ , $e = (v_s, v_t) \in E$ do
742	6:	$V_s^* \leftarrow \{ v' \mid (v', v_s) \in E^* \} $ > All predecessor vertices that reach $v_s$ .
745	7:	$V_t^* \leftarrow \{ v' \mid (v_t, v') \in E^* \} $ > All successor vertices of $v_t$ .
745	8:	<b>if</b> $(\forall v_s^* \in V_s^*$ . LeftCommute $(v, v_s^*)) \land (\forall v_t^* \in V_t^*$ . LeftCommute $(v_t^*, v))$ <b>then</b>
746	9:	$I \leftarrow I \cup \{(v, e)\}$
747 748	10:	return I
749	11:	<b>function</b> LeftCommute( $v, v'$ )
750	12:	<b>input:</b> Two fragments <i>v</i> , <i>v</i> '
751	13:	<b>output:</b> true iff $v$ left commutes with $v'$
752	14:	$X \leftarrow \{x \mid x \text{ is a variable in } v \text{ or } v'\}.$
754	15:	$X_L \leftarrow \{x_l \text{ fresh name} \mid x \in X\}$ $X_R \leftarrow \{x_r \text{ fresh name} \mid x \in X\}$
755	16:	$S_L \leftarrow (v'; v)$ [assume/waituntil, $X_L/X$ ] $S_R \leftarrow (v; v')$ [assert/waituntil, $X_R/X$ ]
756	17:	<b>return</b> Verify $({X_L = X_R} S_L; S_R {X_L = X_R})$

Fig. 7. Algorithm to find all safe interleavings.

other hand, we need to show that  $S_R$  does not block; thus, we assert that the predicates in the 761 762 waituntil statement evaluate to true under the assumption that they also evaluate to true in  $S_L$ . 763 Finally, in addition to showing that waituntil's are not blocked, we also need to establish that 764 the monitor state is the same in  $S_L$  and  $S_R$ . Thus, the Hoare triple we construct checks that the 765 values of variables are the same at the end, assuming that they are the same in the beginning. Note that the implementation of right commutativity is the same with v and v' swapped; thus, 766 767 RIGHTCOMMUTE(v, v') can be checked by directly calling LEFTCOMMUTE(v', v).

#### 4.4 MaxSAT Encoding

In this section, we describe our MaxSAT encoding which is formalized as inference rules in Figure 8. Recall that the encoding procedure takes as input (a) an FDG representation of the monitor, (b) the results of the static analysis, and (c) an upper bound on the maximum number of locks, and it produces a set of hard constraints  $\mathcal H$  and a set of soft constraints  $\mathcal S$ . In the remainder of this section, we describe the inference rules in Figure 8 for generating these constraints in more detail.

Variables. Our MaxSAT encoding uses two types of variables. First, we introduce variables of 776 the form  $h_{v_i}^{l_j}$  indicating that fragment  $v_i$  needs to hold lock  $l_j$ . Thus, given an FDG with *n* vertices and an upper bound N on the number of locks, our encoding contains  $n \times N$  such variables. 779 The second type of variable used in our encoding is of the form  $a_{fld}$  indicating that fld should be implemented using an atomic type.

*Mutex encoding.* Given a set of fragments F and an upper bound N on the number of locks, we often need to enforce that all fragments in F share at least one of the N possible locks. We write

757 758

759 760

768 769

770

771

772

773

774

775

777 778

780

781

782

$$\begin{array}{c|c} \hline \text{Race-1} & \frac{IsFrag(v_1) \quad IsFrag(v_2) \quad Races = \mathcal{R}(v_1, v_2) \quad Races \subseteq \mathcal{F} \quad Races = \{ \text{ this.} f \} \\ \hline \text{Mutex}(\{v_1, v_2\}, \mathcal{N}) \lor a_f \in \mathcal{H} \\ \hline \text{Race-2} & \frac{IsFrag(v_1) \quad IsFrag(v_2) \quad Races = \mathcal{R}(v_1, v_2) \quad Races \neq \emptyset \quad (|Races| > 1 \lor Races \notin \mathcal{F}) \\ \hline \text{Mutex}(\{v_1, v_2\}, \mathcal{N}) \in \mathcal{H} \\ \hline \text{I-Leave} & \frac{IsFrag(v) \quad IsEdge(e) \quad e = (v_s, v_t) \quad \neg SafeInterleaving(v, e) \\ \hline \text{Mutex}(\{v_v, v_s, v_t\}, \mathcal{N}) \in \mathcal{H} \\ \hline \text{Wart} & \frac{F = \{ f \mid IsFrag(f), f \equiv \texttt{waituntil}(p) \} \\ \hline \text{Mutex}(F, \mathcal{N}) \in \mathcal{H} \quad \bigwedge_{i=1}^{\mathcal{N}} \wedge_{v_1, v_2 \in F} \left( h_{v_1}^{\mu_i} \leftrightarrow h_{v_2}^{\mu_i} \right) \in \mathcal{H} \\ \hline \text{Min-Lock} & \frac{MF = \{ v \mid IsFrag(v), Method(v) = m \} \\ \hline \bigcup_{i=1}^{\mathcal{N}} f_{\in MF} \neg h_f^{i_i} \} \subseteq S \\ \hline \text{Max-Par} & \frac{IsFrag(v_1) \quad IsFrag(v_2) \quad \mathcal{R}(v_1, v_2) = \emptyset}{\neg Mutex(\{v_1, v_2\}, \mathcal{N}) \in S} \end{array}$$

$$\boxed{\text{Aux-Defs}} \quad Mutex(F, \mathcal{N}) = \bigvee_{i=1}^{\mathcal{N}} \bigwedge_{f \in F} h_f^{l_i} \quad LockOrder((v_s, v_t), \mathcal{N}) = \bigwedge_{1 \le \ell < u \le \mathcal{N}} \neg \left(h_{v_s}^u \land h_{v_t}^u \land \neg h_{v_s}^\ell \land h_{v_t}^\ell\right)$$

Fig. 8. Inference rules for MAXSATENCODING(M, G, S, N) procedure. G = (V, E) is an FDG of monitor M,  $S = (\mathcal{F}, \mathcal{R}, \mathcal{I})$  are the results of the static analysis, and N is an upper bound on the number of locks. Predicate IsFrag(v) is true if  $v \in V$ , IsEdge(e) if  $e \in E$ , and SafeInterleaving(v, e) if  $(v, e) \in \mathcal{I}$ . Relations Methods(M) and Preds(M) return all methods of monitor M and all predicates that appear as an argument of a waituntil statement in M respectively.

 $Mutex(F, \mathcal{N})$  to denote this requirement. In particular, as shown at the bottom of Figure 8, this is defined as  $Mutex(F, \mathcal{N}) = \bigvee_{i=1}^{\mathcal{N}} \bigwedge_{f \in F} h_f^{l_i}$ .

*Hard constraints.* Next, we describe the hard constraints generated by our MaxSAT encoding. These hard constraints  $\mathcal{H}$  correspond to correctness requirements on the synthesized protocol and include (1) data race freedom, (2) correct signaling and deadlock freedom and (3) atomicity. Specifically, the first two rules in Figure 8 deal with data race freedom, the next rule deals with atomicity, and the last two rules deal with deadlock freedom and correct signaling.

*RACE-1.* The first rule, labeled RACE-1, deals with data race freedom of two fragments that have a data race on a *single* monitor field f. The premises of this rule stipulate that  $v_1, v_2$  are fragments that race *only* on field f which can be converted to atomic (i.e., this.f  $\in \mathcal{F}$ ). In this case, we prevent data races by either (1) enforcing that  $v_1, v_2$  share a lock (the *Mutex* constraint) or (2) ensuring that field f is converted to an atomic field.

*RACE-2.* The next RACE-2 rule prevents data races between fragments where the data race cannot be resolved by making one of the fields atomic. In particular, given two fragments  $v_1$ ,  $v_2$  that have a data race, this rule simply enforces that they share a common lock via the *Mutex* function.

*I-LEAVE.* The next rule generates constraints to ensure that monitor operations appear to take place atomically. In particular, if the static analysis cannot prove (v, e) to be a strongly safe interleaving (recall Definition 4.12), then we need to ensure that a thread cannot execute v when some other thread is executing e. To do so, we ensure that  $v, v_s, v_t$  all share a common lock by generating a *Mutex* hard constraint for these three fragments.

*L-ORDER.* The next rule, labeled L-ORDER, ensures that the resulting synchronization protocol is deadlock-free. Specifically, for every edge  $e = (v_s, v_t)$  in the input FDG, this rule generates a hard constraint, via *LockOrder*(e, N) (defined at the bottom of Figure 8), that ensures that every lock acquisition respects the total order on locks. In particular, for every pair of locks l, u such that l < u, *LockOrder*(e, N) prevents the synchronization protocol from violating the global order on locks. Recall that a protocol violates this global order if it acquires the "smaller" lock l between  $v_s$  and  $v_t$  while both  $v_s$  and  $v_t$  hold lock u. Thus, the hard constraint generated by *LockOrder*(e, N) prevents this from happening.

Example 4.15. Assuming 
$$\mathcal{N} = 2$$
, this rule generates  $\neg \left(h_{v_s}^{l_2} \wedge h_{v_t}^{l_2} \wedge \neg h_{v_s}^{l_1} \wedge h_{v_t}^{l_1}\right)$  for edge  $(v_s, v_t)$ .

*WAIT*. The last hard constraint rule, called WAIT, is used for associating a single lock with each condition variable. In particular, since all fragments of the form waituntil(p) must hold the same set of locks, this rule generates two hard constraints for every waituntil predicate p of the input monitor: (1) a mutex constraint for all waituntil(p) fragments and (2) a constraint that enforces that all waituntil(p) fragments must share *all* common locks.

**Soft Constraints.** As discussed earlier, our goal is to generate a synchronization protocol that is not only correct-by-construction but one that also results in efficient code. Hence, as a proxy metric for efficiency, we want to (1) minimize the number of locks and atomic fields that are introduced, and (2) maximize the number of fragments that can run in parallel. The remaining three rules in Figure 8 introduce soft constraints to encode this optimization objective.

MIN-LOCK. The rule labeled MIN-LOCK is used for minimizing the number of locks. However, instead of simply minimizing the total number of locks used by the protocol, the soft constraints generated by this rule minimize the number of locks used *per method*. Even though this is not equivalent to minimizing the number of locks used by the entire protocol, we have found this approach to synthesize protocols with a more even distribution of locks among the monitor methods. In practice, such protocols are more desirable because they avoid scenarios where a subset of the methods incur a higher synchronization cost than others. Specifically, this rule generates a soft constraint for every lock  $l \in \{l_1...l_N\}$  and every method m of M and asserts that none of the fragments in *m* hold lock *l*.

MIN-ATOM. The MIN-ATOM rule generates soft constrains to minimize the number of fields that are made atomic by asserting that  $a_{fld}$  is assigned to false.

*MAX-PAR*. The last rule called MAX-PAR generates soft constraints to maximize parallelism. Specifically, for every pair of fragments (v, v') that do not have data races, we add a soft constraint stating that v and v' do not share any locks.

We conclude this Section with a theorem that states the correctness of our MaxSAT encoding.

THEOREM 4.16. Let *m* be a model of the generated MaxSAT instance and  $(\mathcal{L}, \mathcal{A}, \mathcal{P})$  be the synchronization protocol constructed as follows:

$$\mathcal{L} = \left\{ v \mapsto \left\{ l \mid m[h_v^l] \right\} \right\} \ \mathcal{R} = \left\{ \mathsf{fld} \mid m[a_{fld}] \right\} \ \mathcal{P} = \left\{ p \mapsto l_i \mid \mathsf{IsWait}(v, p), i = \min(\{j \mid m[h_v^l]\}) \right\}$$

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

where, IsWait(v, p) is true if v is a waituntil statement on p. Then,  $(\mathcal{L}, \mathcal{A}, \mathcal{P})$  is a correct synchro-883 nization protocol. 884

PROOF. Can be found in Appendix E.

#### 5 IMPLEMENTATION

885

886 887

888

889

890

891

892

893

894

895

896 897

898

899

900

901

902

903

904

905

906

907

908

909

911

912

913

914

915

916

917

918

919 920

921

922

923

We have implemented our approach in a tool called CORTADO that emits explicit-synchronization monitors in Java. CORTADO is based on the Soot program analysis infrastructure [Vallée-Rai et al. 1999] and the Z3 SMT solver [De Moura and Bjørner 2008]. In particular, we use Soot to perform various kinds of static analyses needed by our method (e.g., pointer analysis) and to translate the input monitor to an explicit-synchronization monitor in Java. Furthermore, we leverage Z3 for solving MaxSAT instances and discharging the validity queries that arise when checking commutativity between fragments. In the remainder of this section, we discuss several design choices and optimizations that were not discussed previously.

*Weights of soft constraints.* As expected, the quality of the synthesized protocol depends on the model returned by the MaxSAT solver. In practice, we have observed certain types of soft constraints to be more important than others for efficiency. Thus, our implementation assigns different weights for different classes of soft constraints. For instance, because it is always preferable to use an atomic field instead of a lock, CORTADO assigns a higher weight to soft constraints generated by rule Мил-Атом from Figure 8 than the ones generated by rule Мил-Lock.<sup>3</sup>

Constructing FDGs. As mentioned in Section 4, CORTADO uses a heuristic to partition the input CFG into fragments. The goal of this heuristic is to maximize parallelization opportunities while ensuring that the partition results in a valid FDG according to Definition 4.2. Our heuristic places every loop in its own fragment (to make sure that the FDG is well-formed) and, for code outside loops, CORTADO creates a new fragment whenever it detects an update to monitor state (i.e., this.fld = \*). In practice, we found this heuristic to achieve a good balance between the number of parallelization opportunities and the size of the resulting FDG.<sup>4</sup> 910

Static analysis optimization. Our approach uses an off-the-shelf pointer analysis to detect which pairs of FDG fragments do not have a data race (and, so, can run in parallel). However, such an approach, based on pointer analysis alone, often leads to imprecision. For example, Soot's pointer analysis cannot prove that fragments 2 and 6 in Figure 1a do not contain any races, as it does not reason about individual array elements. Therefore, in order to increase the precision of the static analysis, CORTADO implements an SMT-based static analysis on top of Soot's built-in pointer analysis and generates appropriate verification conditions (similar to the ones generated by Gurfinkel et al. [2015]) to prove that memory accesses of two fragments are disjoint.

## **EVALUATION**

We evaluated CORTADO's ability to generate fine-grained locking protocols by performing a set of experiments that are designed to answer the following research questions:

- **RO1** How does the code generated by CORTADO compare against explicit-synchronization monitors 924 written by experts? 925
- RQ2 How does the technique implemented in CORTADO compare against other compile-time 926 state-of-the-art approaches targeting implicit-synchronization monitors? 927

928

<sup>929</sup> <sup>3</sup>An ablation study that demonstrates the need for adjusting the weights of soft constraints can be found in Appendix A.3.

<sup>&</sup>lt;sup>4</sup>An ablation study that justifies the design of this heuristic can be found in Appendix A.2. 930

<sup>931</sup> 

**RQ3** How does the static analysis for inferring safe interleavings impact the quality of the code generated by CORTADO?

RQ4 How long does CORTADO take to synthesize code and how complex are the resulting protocols?

To answer these research questions, we conducted experiments on ten explicit-synchronization monitors from popular open source repositories. Aside from CORTADO, we consider two additional baselines, described below, that aid us in answering our second and third research questions.

**Benchmarks.** The benchmarks used in our evaluation are collected from popular open source GitHub repositories. We wrote a crawler (based on GitHub's REST API [GitHub 2022]) to automatically identify candidate explicit-synchronization monitors implemented in Java by searching for keywords like lock, unlock, await, etc. We then manually inspected class files returned by the crawler in decreasing order of GitHub popularity (stars and forks) and identified self-contained monitor-style classes that encapsulate shared state accessed by multiple threads. We included such a monitor in our benchmarks only if it satisfies the following conditions: (1) the class has a welldefined API for client threads and (2) it contains *parallelization opportunities that can be realized via fine-grained locking*.<sup>5</sup> We manually isolated the shared state and monitor methods of the class file to obtain a standalone explicit-synchronization monitor and then manually translated it to an equivalent implicit monitor. To convert a benchmark to an implicit monitor, we simply removed all synchronization code (i.e., locking and signaling operations) and introduced appropriate waitunti1 statements. In total, we collected 10 monitors from popular repositories such as Spring Framework (a Java-based framework for creating enterprise applications), Java JDK, Apache Spark (an analytics engine for large-scale data processing), etc.<sup>6</sup>

**Baselines.** As mentioned above, our evaluation uses two additional baselines in order to answer RQ2 and RQ3. To compare against other compile-time techniques (RQ2), we evaluate EXPRESSO [Ferles et al. 2018], a tool that addresses the same problem as this paper but generates an explicit signal monitor using a *single global lock* shared by all monitor methods. To evaluate the importance of our static analysis (RQ3), we created an ablated version of CORTADO, called ABLATED, which uses a very coarse analysis to infer safe interleavings. This ablated version considers (v, ( $v_s$ ,  $v_t$ )) a safe interleaving only if v does not have any data races with any predecessor (resp. successor) of  $v_s$  (resp.  $v_t$ ). This is a sufficient condition for strong safety, but it only requires checking data races rather than discharging a set of Hoare triples.

**Evaluating performance.** Following prior work [Ferles et al. 2018; Hung and Garg 2013], we evaluate the performance of each monitor implementation by performing *saturation tests* [Cherem et al. 2008] wherein threads perform monitor operations without doing any additional work. We collect our performance measurements using the Java Microbenchmark Harness (JMH) [Shipilev et al. 2021]. All measurements are conducted on a 112-way (56-core × 2 SMT) Intel Xeon CPU W-3275 2.50GHz with 256 GB of memory using JDK 1.8.0\_272. In this section, we present results for each benchmark for up to 128 threads, chosen as an arbitrary stopping point past the total number of hyper-threads. Results for up to 256 threads can be found in Appendix A.

1:20

<sup>&</sup>lt;sup>5</sup>If a monitor does not contain parallelization opportunities, our technique generates code equivalent to that synthesized by Ferles et al. [2018]. Since the goal of our evaluation is to evaluate CORTADO's ability to generate fine-grained locking protocols, we did not include benchmarks from prior work [Ferles et al. 2018; Hung and Garg 2013] that do not contain such parallelization opportunities.

<sup>&</sup>lt;sup>6</sup>All benchmarks are publicly available here: https://github.com/utopia-group/cortado

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

## 981 6.1 Performance Results

982

983

984

985

986

987

988

989

990

991 992

993

994

995

Figure 9 plots the average time taken per monitor method invocation (i.e., milliseconds/operation) against the number of threads. In what follows, we analyze the plots in more detail and present several conclusions drawn from these results. Because our benchmarks only contain monitors where fine-grained locking is beneficial, we emphasize that our conclusions only apply to such monitors.

**Comparison against hand-written implementations (RQ1).** For every benchmark, the explicit synchronization monitor generated by CORTADO performs *better* than the expert-written implementation as the number of threads increases. In particular, CORTADO-synthesized code performs on average  $3.7 \times^7$  and up to  $39.1 \times$  times faster than the original code.

**Comparison against Expresso (RQ2).** CORTADO-generated explicit monitors perform better than Expresso explicit monitors generated from the same implicit specification on all benchmarks. CORTADO-synthesized code outperforms Expresso by 4.0× on average (and up to 48.7×).

996 Comparison against ABLATED (RQ3). Finally, we analyze how CORTADO compares to its 997 simplified version, ABLATED, which does not use the results of the safe interleavings analysis from 998 Section 4.3. In five cases, the code generated by ABLATED is equivalent to the code generated by 999 EXPRESSO and therefore worse than CORTADO. In two other cases (PausableThreadPoolExecutor 1000 and ProgressTracker), ABLATED generates code different from both EXPRESSO and CORTADO. For 1001 PausableThreadPoolExecutor, the code generated by ABLATED is slower than that of EXPRESSO 1002 because many of the operations it parallelizes are very cheap, so the overhead of extra locks 1003 outweighs their benefit. On the other hand, our static analysis detects several safe interleavings 1004 which enable CORTADO to synthesize a protocol with cheaper synchronization operations. Finally, 1005 for the remaining three cases, the code generated by ABLATED matches the one generated by 1006 CORTADO. This ablation study demonstrates that the safe interleaving analysis from Section 4.3 1007 helps extract additional concurrency on five of our benchmarks. 1008

# 10091010 6.2 Synthesis Time & Protocol Complexity

To evaluate the cost and complexity of synthesizing code with CORTADO (RQ4), Table 1 summarizes its running time and presents some statistics about the synthesized protocols. For each benchmark, we report the running time for the various phases of the tool: pointer analysis with Soot, signal placement with EXPRESSO, and synthesis with CORTADO. We also report the number of locks and atomic fields in the synthesized protocol.

Table 1 shows that CORTADO terminates in under one minute for all but two benchmarks. For these two outliers, the synthesis time is dominated by EXPRESSO's monitor invariant inference, which is necessary for signal placement. Overall, CORTADO is able to extract better performance than EXPRESSO alone with only a small additional compile-time cost.

The last three columns in Table 1 provide statistics about the synthesized explicit monitors. Most monitors benefit from CORTADO's ability to introduce atomic fields, which reduces the overhead of operations on monitor state that would otherwise require a lock. The lines of code (LOC) results show that CORTADO synthesizes explicit monitors that are on average 1.7× larger than their implicit specifications.

1025 1026

 <sup>&</sup>lt;sup>7</sup>In order to handle outliers such as in JobWrapper, for all reported aggregate speedups (max, mean, etc.) we throw out data
 points with a z-score greater than two.



Fig. 9. Performance Results For All Tools. The y-axis is in log scale and time measurements are in milliseconds. The shadowed regions surrounding each line present 99.9% confidence interval of each measurement.

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1077 1078

1075

1076

1:22

		Synthesis Time (secs)				Synthesized Protocol		
Benchmark	LOC	Soot	Expresso	Cortado	Total	#Lock/#Op	#Atomic/#Op	LOC
ArrayBlockingQueue	287	16.7	1897.7	372.7	2302.5	2 / 18	1 / 25	514
ConcurrencyThrottleSupport	33	17.0	1.4	0.2	18.7	1/1	1 / 4	68
CountableThreadPool	54	20.4	0.8	0.2	21.5	1/1	1/5	85
JobWrapper	33	17.0	0.6	0.1	17.7	1/1	1/3	63
PausableThreadPoolExecutor	79	20.7	0.8	0.9	22.5	3/4	2/9	122
ProgressTracker	65	0.7	6.8	0.2	8.0	2 / 7	1 / 4	119
RealmThreadPoolExecutor	34	18.4	0.3	0.1	18.7	1/1	1/3	61
RoundTripWorker	62	16.7	3.1	0.4	20.4	2/4	1/4	103
SinkQueue	75	15.7	981.4	34.6	1031.8	2/4	1/8	131
WSDataListener	158	18.3	5.2	1.0	25.1	4 / 11	0 / 0	222

Table 1. Synthesis time for each phase of CORTADO and summaries of the synthesized protocols. LOC is lines
 of code. Soot indicates pointer analysis time, EXPRESSO is the time for monitor invariant generation and signal
 placement, and CORTADO shows the additional time on top of Soot and EXPRESSO.

#### 7 RELATED WORK

1093 1094 1095

1096

Monitor abstractions. The notion of monitors as an organizing abstraction for concurrent program ming originates with Hoare [1974] and Hansen [1973]. Monitors offer the same synchronization
 facilities as semaphores—the ability to coordinate multiple threads and enforce mutual exclusion—
 but encapsulate the state protected by those facilities and automate mutual exclusion when entering
 and exiting the monitor's operations. Lampson and Redell [1980] further extended the monitor
 abstraction in the Mesa programming language to handle spurious wake ups.

These early monitor abstractions are *explicit-signal* monitors in the taxonomy of Buhr et al. 1103 [1995] because they require the programmer to explicitly insert condition variables and signalling 1104 operations to coordinate threads within the monitor. This requirement places both a safety and a 1105 liveness burden on the programmer: they must place signals correctly to preserve invariants about 1106 the monitor's state, but must also insert enough signals to avoid deadlock. An alternative is to use 1107 an implicit-signal (or automatic-signal) monitor, in which signals are inserted automatically by the 1108 compiler, language runtime, or operating system. Hoare [1971] proposed the notion of conditional 1109 critical regions (CCRs), which allow for monitor operations to block until a guard predicate over 1110 the monitor state is satisfied by some other thread. A CCR implementation would automatically 1111 block and signal threads in a fashion consistent with this guard semantics. 1112

Implicit-signal monitors simplify concurrent programming, but come at a steep performance 1113 cost–Buhr et al. [1995] estimate that implicit-signal monitors are  $10-50\times$  slower than explicit ones. 1114 More recent work has tried to lower the cost of implicit-signal monitors. AutoSynch [Hung and Garg 1115 2013] uses a combination of compile-time instrumentation and run-time evaluation to efficiently 1116 compute which threads should be woken when monitor state changes. This approach lowers the 1117 cost of implicit monitors to be close to, or sometimes better than, explicit ones. Expresso [Ferles et al. 1118 2018] takes a different approach, using compile-time static analysis to synthesize an explicit-signal 1119 monitor equivalent to an implicit-signal version given as input. In this way, Expresso is able to 1120 erase the dynamic cost of implicit-signal monitors, and in most cases is comparable to hand-written 1121 explicit monitors. However, Expresso uses a single lock for the entire monitor and does not allow 1122 concurrent execution of threads within the monitor even when safe. Our work expands on this 1123 direction by using a richer static analysis to infer additional concurrency opportunities and uses 1124 MaxSAT to synthesize a safe and efficient locking protocol. Hence, our key contribution is to 1125 synthesize an explicit monitor that *appears* to match the semantics of the implicit one, but runs 1126 1127

monitor operations concurrently when possible and efficient. As we show in the evaluation, our proposed approach can often make the synthesized monitor *faster* than a hand-written equivalent.

1130

1:24

Automatic synchronization. An appealing approach to lower the difficulty of concurrent pro-1131 gramming is to deploy program analysis and synthesis techniques for automation. The common 1132 abstraction for much of this work is for the programmer to annotate *atomic sections* that should 1133 be executed atomically. Emmi et al. [2007] present a technique for lock allocation to an annotated 1134 program. They reduce the problem to integer linear programming and deploy the resulting tool 1135 on large-scale C and Java programs. Other approaches [Halpert et al. 2007; Hicks et al. 2006; Mc-1136 Closkey et al. 2006], on the other hand, take a purely static analysis route and attempt to maximize 1137 parallelism based solely on the results of the analysis. Cherem et al. [2008] present an alternative 1138 technique that uses runtime support to enable finer-grained concurrency. Compared to these efforts, 1139 CORTADO applies to the more limited domain of monitors, but in exchange for this limitation is 1140 able to reason about conditional signalling and can allow atomic sections to run concurrently so 1141 long as the illusion of atomicity is maintained. 1142

Other approaches start from a sequential program and automatically generate an equivalent concurrent program. The closest work to ours in this space is that of Golan-Gueta et al. [2011] which generates concurrent data structures given their sequential implementation. Compared to our method, their approach is applicable only to data structures that satisfy certain shape properties and all synthesized programs adhere to the same locking protocol, whereas CORTADO generates a synchronization protocol specialized to the input monitor.

1149 Concurrency verification. CORTADO reasons about concurrent program executions by building on 1150 work in concurrent program analysis and verification. Our notion of left- and right-commutativity 1151 (Definition 4.11) comes from Lipton's work on reduction as a concurrency proof technique [Lipton 1152 1975]. Reduction translates interleaved program executions to simpler, equivalent sequential execu-1153 tions by exploiting the commutativity properties of individual program steps. We use the same 1154 idea but in reverse: starting with a sequential history (Definition 3.7), we use a static analysis of 1155 commutativity to determine how to safely introduce interleavings into that history, and use that 1156 information to determine how to assign locks to program fragments. 1157

# 1159 8 CONCLUSION

1160 We presented a technique for synthesizing fine-grained synchronization protocols for implicitly 1161 synchronized monitors. Our approach first employs a novel static analysis to identify safe inter-1162 leavings opportunities between code fragments and uses the results of this analysis to generate a 1163 MaxSAT encoding whose solution can be used to synthesize an efficient and correct-by-construction 1164 explicit-synchronization monitor. We have implemented our method in a tool called CORTADO and 1165 evaluated its effectiveness eight monitors collected from popular open source applications. The 1166 results of our experimental evaluation demonstrate that CORTADO is able to generate non-trivial 1167 synchronization protocols that are  $3.7 \times$  times faster than the original implementation on average 1168 (and up to 39.1× times for some outliers). 1169

# ACKNOWLEDGMENTS

This material is based upon work supported by the National Science Foundation under Grant No. nnnnnn and Grant No. mmmmmm. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation.

1176

#### 1177 **REFERENCES**

- Noga Alon. 1986. Covering graphs by the minimum number of equivalence relations. *Combinatorica* 6, 3 (Sep 1986), 201–206.
   https://doi.org/10.1007/BF02579381
- 1180 Andrew D Birrell. 1989. An introduction to programming with threads. Digital Systems Research Center.
- Peter A. Buhr, Michael Fortier, and Michael H. Coffin. 1995. Monitor Classification. ACM Comput. Surv. 27, 1 (1995), 63–107.
- Sigmund Cherem, Trishul M. Chilimbi, and Sumit Gulwani. 2008. Inferring locks for atomic sections. In *Proceedings of the* ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI). Tucson, AZ, USA, 304–315.
- Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An efficient SMT solver. *Tools and Algorithms for the Construction and Analysis of Systems* (2008), 337–340.
- Michael Emmi, Jeffrey S. Fischer, Ranjit Jhala, and Rupak Majumdar. 2007. Lock allocation. In *Proceedings of the 34th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2007, Nice, France, January 17-19, 2007.* 291–296.
- Kostas Ferles, Jacob Van Geffen, Isil Dillig, and Yannis Smaragdakis. 2018. Symbolic Reasoning for Automatic Signal
   Placement. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Philadelphia, PA, USA) (*PLDI 2018*). Association for Computing Machinery, New York, NY, USA, 120–134. https:
   //doi.org/10.1145/3192366.3192395
- GitHub. 2022. GitHub REST API. https://docs.github.com/en/rest
- Guy Golan-Gueta, Nathan Bronson, Alex Aiken, G Ramalingam, Mooly Sagiv, and Eran Yahav. 2011. Automatic fine-grain locking using shape properties. *ACM SIGPLAN Notices* 46, 10 (2011), 225–242.
- Arie Gurfinkel, Temesghen Kahsai, and Jorge A Navas. 2015. SeaHorn: A framework for verifying C programs (competition contribution). In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer, 447–450.
- 1196Richard L Halpert, Christopher JF Pickett, and Clark Verbrugge. 2007. Component-based lock allocation. In 16th International<br/>Conference on Parallel Architecture and Compilation Techniques (PACT 2007). IEEE, 353–364.
- <sup>1197</sup> Per Brinch Hansen. 1973. Operating System Principles. Prentice-Hall.
- Michael Hicks, Jeffrey S Foster, and Polyvios Pratikakis. 2006. Lock inference for atomic sections. In Proceedings of the First
   ACM SIGPLAN Workshop on Languages, Compilers, and Hardware Support for Transactional Computing.
- C. A. R. Hoare. 1971. Towards a theory of parallel programming. In *Operating Systems Techniques, Proceedings of a Seminar* at *Queen's University, Belfast*. Belfast, Northern Ireland.
- C. A. R. Hoare. 1974. Monitors: An Operating System Structuring Concept. Commun. ACM 17, 10 (1974), 549–557.
- Wei-Lun Hung and Vijay K. Garg. 2013. AutoSynch: an automatic-signal monitor based on predicate tagging. In ACM
   SIGPLAN Conference on Programming Language Design and Implementation PLDI. Seattle, WA, USA, 253–262.
- Butler W. Lampson and David D. Redell. 1980. Experience with Processes and Monitors in Mesa. Commun. ACM 23, 2 (1980), 105–117.
- 1206 William Landi and Barbara G. Ryder. 1992. A Safe Approximate Algorithm for Interprocedural Aliasing. SIGPLAN Not. 27, 7 (July 1992), 235–248. https://doi.org/10.1145/143103.143137
- Richard J. Lipton. 1975. Reduction: A Method of Proving Properties of Parallel Programs. Commun. ACM 18, 12 (1975), 717-721.
- 1209Bill McCloskey, Feng Zhou, David Gay, and Eric Brewer. 2006. Autolocker: synchronization inference for atomic sections. In1210Conference record of the 33rd ACM SIGPLAN-SIGACT symposium on Principles of programming languages. 346–358.
- TS Michael and Thomas Quint. 2006. Sphericity, cubicity, and edge clique covers of graphs. *Discrete Applied Mathematics* 154, 8 (2006), 1309–1313.
- 1212
   Aleksey Shipilev, Sergey Kuksenko, Astrand Astrand, Staffan Freiberg, and Henrik Loef. 2021. OpenJDK: jmh. http:

   1213
   //openjdk.java.net/projects/code-tools/jmh/
- Raja Vallée-Rai, Phong Co, Etienne Gagnon, Laurie Hendren, Patrick Lam, and Vijay Sundaresan. 1999. Soot-a Java bytecode
   optimization framework. In *Proceedings of the 1999 conference of the Centre for Advanced Studies on Collaborative research*.
   IBM Press, 13.
- 1217

- 1219
- 1220
- 1221
- 1222
- 1223
- 1224 1225

# 1226 A ADDITIONAL EXPERIMENTAL EVALUATIONS

This section presents additional experimental data for the evaluation of Section 6 as well as ablations
 related to several design decisions in the implementation of CORTADO. In particular, Section A.1
 presents additional data points for the evaluation presented in Section 6, Section A.2 presents an
 ablation study that justifies the design of our heuristic for constructing an FDG, and Section A.3
 presents an ablation study that demonstrates the need for adjusting the weights of soft constraints
 in our MaxSAT encoding.

# A.1 Additional Data Points per Benchmark

In this Section we present additional data points for all the experiments presented in Section 6. As mentioned in Section 6, we chose 128 threads as a stopping point past the number of total hyper-threads in the machine used in our evaluation. In this Section, we provide data points up to 256 threads.

Figure 10 presents the results for all benchmarks in our evaluation up to 256 threads. As demonstrated by Figure 10, the general trend for all benchmarks is the same as the date presented in Section 6. Note that for benchmarks where the code generated by CORTADO exhibits similar run-time performance with the other three implementation we consider in our evaluation (e.g., RoundTripWorker), context-switches seem to dominate the running time as the number of threads increases.

# 1248 1249 A.2 Ablation Study for FDG Construction Heuristic

This section presents an ablation study for the design of our FDG construction heuristic described in 1250 Section 5. To justify the decisions behind the design of our heuristic, we implemented two additional 1251 versions of CORTADO that differ in the way they construct the FDG. In particular, we implemented 1252 a version that creates finer-grained FDGs than CORTADO and one that creates coarser-grained ones. 1253 The finer-grained version, named STMT-ABLATION, puts every non-composite statement outside of 1254 a loop in its own fragment. Statements inside a loop are grouped together in the same fragment 1255 since FDGs are acyclic. The coarser version of our tool, named CCR-ABLATION, simply puts each 1256 CCR in its own fragment. 1257

Figure 11 presents the results of this ablation study. As demonstrated by the results, there are several cases where CORTADO performs better than at least one of its two modified versions. The cases where all three versions perform similarly are benchmarks where the synthesized synchronization protocol mainly exploits data-level parallelism among different CCRs in the monitor, which our tool can exploit given any FDG. Overall, this study demonstrates the need for a customized heuristic for constructing an FDG suitable for maximizing parallelism of implicit synchronisation monitors.

1265 1266 1267

## A.3 Ablation Study for MaxSAT Soft Constraint Weights

Finally, this Section presents an ablation of the soft constraints weights in our MaxSAT encoding.
 As mentioned in Section 5, CORTADO assigns different weights to different classes of soft constraints
 because some of them are more important for synthesizing the optimal synchronization protocol.
 To demonstrate this, we have created a modified version of our tool, named WEIGHT-ABLATION,
 that assigns the same weight to all soft constraints.

1274

1234 1235

1236

1237

1238

1239



Fig. 10. Performance results for all tools up to 256 threads. The y-axis is in log scale and time measurements
 are in milliseconds. The shadowed regions surrounding each line present 99.9% confidence interval of each
 measurement.



Fig. 11. Performance results for the FDG ablation study. The y-axis is in log scale and time measurements are in milliseconds. The shadowed regions surrounding each line present 99.9% confidence interval of each measurement.

1371 1372

1369

1370

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:28

The results of this ablation are presented in Figure 13. As this figure demonstrates, there are several cases where the ablated version of the tool performs significantly worse than CORTADO. To give a concrete example of why this is the case, consider the monitor of Figure 12 that contains two methods both of which modify field x. Be-cause the body of method bar can be interleaved between the waituntil statement and the increment statement of method foo, the optimal synchronization protocol would convert field *x* to an atomic integer and introduce a lock that would only be held in method foo. However, an 

```
class M {
    int x = 0;
    void foo() {
        waituntil(x < 10);
        x++;
    }
    void bar() {
        x--;
    }
}</pre>
```

Fig. 12. A simple implicit monitor.

equivalent protocol would be to simply protect both foo and bar with the same global lock. So, if all constraints have an equal weight, CORTADO could generate both of these protocols, since they would have the same optimum objective value. As mentioned in Section 5, CORTADO's MaxSAT encoding prefers assignments where a race between two fragments (like the one on field x) are resolved via an atomic field rather than a lock. This forces CORTADO to generate the optimal solution for the monitor of Figure 13. The adjusted weights for other classes of soft constraints try to steer CORTADO to better performing synchronization protocols in a similar way.



Fig. 13. Performance results for the soft constraints weights ablation. The y-axis is in log scale and time measurements are in milliseconds. The shadowed regions surrounding each line present 99.9% confidence interval of each measurement.

1469 1470

1467

1468

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:30

#### 1471 B PROOF OF NP-COMPLETENESS

<sup>1472</sup> To aid the reader, we restate Theorem 4.4.

THEOREM B.1. (NP-Completeness) Let  $\mathcal{G} = (V, E)$  be the FDG representation of a monitor M and let  $\Pi \subseteq V \times V$  be a set of fragment pairs that can safely run in parallel. Then, deciding whether there exists a synchronization protocol with at most k locks and atomic fields that allows all pairs in  $\Pi$ to run in parallel is an NP-Complete problem.

<sup>1478</sup> PROOF OF THEOREM 4.4. We prove the theorem by reduction to the edge clique cover prob-<sup>1479</sup> lem [Michael and Quint 2006]. Let G = (V, E) be an undirected graph. For each  $v \in V$ , let E(v) be <sup>1480</sup> the set of edges incident to v.

<sup>1481</sup> Define monitor M as follows: for each edge  $e \in E$  define a new field  $f_e$  of the monitor, initially <sup>1482</sup> set to zero. For each vertex  $v \in V$ , define a CCR  $inc_v()$  which increments each  $f_e$  for  $e \in E(V)$ , has <sup>1483</sup> a guard of  $\top$ , and returns nothing.

Let  $G_M$  be the control-flow graph of monitor M. Note that there is one waituntil( $\top$ ) statement for each  $inc_v()$  and |E(v)| = degree(v) increment statements in  $inc_v()$ , so there are |V| + 2|E|total nodes in  $G_M$ . Define  $\{G_{1,M}, \ldots, G_{|V|+2|E|,M}\}$  to be the partition of  $G_M$  into singletons. Then, let  $\mathcal{G}_M = (V_M, E_M)$  to be the fragment dependency graph obtained from this partition, and let  $\Pi \subseteq V_M \times V_M$  be the set of fragment pairs that can safely run in parallel.

We write  $frag_{v,e}$  for the fragment which increments  $f_e$  in method  $inc_v()$ , and  $waituntil_v$  for the  $waituntil(\top)$  statement at the beginning of  $inc_v()$ . Observe that

$$\Pi = \{ (f_1, f_2) \mid \exists v \in V \text{ such that } f_1 = waituntil_v, \text{ or } f_2 = waituntil_v \}$$
(2)  
$$\cup \{ (frag_{v_1, e_1}, frag_{v_2, e_2}) \mid e_1 \neq e_2 \}.$$

Note that any synchronization protocol which implements  $f_e$  for some  $e = (u, v) \in E$  as an atomic variable is equivalent to one which wraps a unique lock around each  $frag_{u,e}$  and  $frag_{v,e}$ . Therefore, we only need to consider synchronization protocols which use only locks.

<sup>1497</sup> Suppose we are given a synchronization protocol which allows all pairs in  $\Pi$  to run in parallel and uses exactly *k* locks,  $\{\ell_1, \ldots, \ell_k\}$ , for some  $k \in \mathbb{N}$ . Define the vertices holding each lock to be

$$C_i = \left\{ v \in V \mid \exists e \in E(v) \text{ such that } frag_{v,e} \text{ holds lock } \ell_i \right\}.$$
(3)

1501 We claim that  $\{C_1, \ldots, C_k\}$  is a clique edge cover of *G*. First, observe that every edge  $e = (u, v) \in E$ corresponds to two fragments in  $\mathcal{G}_{M}$ :  $frag_{u,e}$  and  $frag_{v,e}$ . Since these fragments must not run in 1502 1503 parallel (due to a data race), they must share some lock. Let  $\ell_i$  be that lock. By the definition of  $C_i$ in Equation 3,  $u \in C_i$  and  $v \in C_i$ . Therefore, every edge appears in  $C_i$  for some  $1 \le i \le k$ . Second, 1504 suppose that  $u \neq v \in V$  are both contained in  $C_i$  for some *i*. By Equation 3, there must be some 1505 1506  $e_u \in E(u)$  and  $e_v \in E(v)$  such that both  $frag_{v,e_n}$  and  $frag_{u,e_n}$  hold lock *i*. Since the two fragments share a lock, we know  $(frag_{v,e_v}, frag_{u,e_u}) \notin \Pi$ . Therefore,  $e_v = e_u$  (by Equation 2). Hence, there is 1507 1508 an edge  $e_u = e_v = (u, v)$  in *E*. Consequently, any two distinct vertices in  $C_i$  are incident, so  $C_i$  is a 1509 clique.

A symmetric argument shows how to construct a synchronization protocol using exactly *k* locks from any edge clique-cover of *G* which has *k* cliques.

<sup>1512</sup> We have shown that, given an arbitrary graph G, in polynomial time we may compute a monitor <sup>1513</sup> M such that M has a synchronization protocol using at most k locks and atomic variables if and <sup>1514</sup> only if G has an edge clique-cover with at most k cliques.

1515

1490

1495

1496

1500

- 1516
- 1517
- 1518 1519

Kostas Ferles, Benjamin Sepanski, Rahul Krishnan, James Bornholt, and Isil Dillig

1520  
1521 (1) 
$$\frac{e = (s, t) \quad \text{CHECKNOTIF}(\mathcal{T}, t) \quad \neg \text{Synch}(s) \quad (s, v), \sigma \Downarrow \sigma' \quad \mathcal{T}' = \text{UPDATESTATE}(\mathcal{T}, t) \\ (\sigma, e, v, \mathcal{T}) \Rightarrow (\sigma', \epsilon, \epsilon, \mathcal{T}')$$

(

(2) 
$$e = (s, t) \quad \text{CHECKNOTIF}(\mathcal{T}, t) \quad s = 1. \text{lock}()$$
$$\underbrace{\text{LockHeld}(\mathcal{T}, t, \sigma[l]) \quad \mathcal{T}' = \text{BlockThreadOnLock}(\mathcal{T}, t, \sigma[l]) \quad \text{NoDeadLocks}(\mathcal{T}')}_{(\sigma, e, v, \mathcal{T}) \Rightarrow (\sigma', \epsilon, \epsilon, \mathcal{T}')}$$

(3) 
$$\begin{array}{c} e = (s,t) \quad \text{CHECKNOTIF}(\mathcal{T},t) \quad s = 1.\operatorname{lock}() \\ \neg \operatorname{LockHeld}(\mathcal{T},t,\sigma[l]) \quad \mathcal{T}' = \operatorname{AcqLock}(\mathcal{T},t,\sigma[l]) \\ \hline (\sigma,e,v,\mathcal{T}) \Rightarrow (\sigma',\epsilon,\epsilon,\mathcal{T}') \end{array}$$

(4) 
$$\frac{e = (s, t) \quad \text{CHECKNOTIF}(\mathcal{T}, t) \quad s = 1. \text{unlock}() \quad \mathcal{T}' = \text{RelLock}(\mathcal{T}, t, \sigma[l])}{(\sigma, e, v, \mathcal{T}) \Rightarrow (\sigma', \epsilon, \epsilon, \mathcal{T}')}$$

(5) 
$$\frac{e = (s, t) \quad \text{CHECKNOTIF}(\mathcal{T}, t) \quad s = c. \text{await}() \quad \mathcal{T}' = \text{BLOCKONCVAR}(\mathcal{T}, t, \sigma[c])}{(\sigma, e, v, \mathcal{T}) \Rightarrow (\sigma', \epsilon, \epsilon, \mathcal{T}')}$$

6) 
$$\frac{e = (s, t) \quad \text{CHECKNOTIF}(\mathcal{T}, t) \quad s = \text{c.signal}() \quad \mathcal{T}' = \text{SIGCVAR}(\mathcal{T}, \sigma[c])}{(\sigma, e, v, \mathcal{T}) \Rightarrow (\sigma', \epsilon, \epsilon, \mathcal{T}')}$$

(7) 
$$\frac{e = (s, t) \quad \text{CHECKNOTIF}(\mathcal{T}, t) \quad s = \text{c.signalAll}() \quad \mathcal{T}' = \text{BCASTCVAR}(\mathcal{T}, \sigma[c])}{(\sigma, e, v, \mathcal{T}) \Rightarrow (\sigma', \epsilon, \epsilon, \mathcal{T}')}$$

$$(\sigma, e, \nu, \mathcal{T}) \Rightarrow (\sigma', \epsilon, \epsilon, \mathcal{T}')$$

(8) 
$$\overline{(\sigma, e :: h, v :: v', \mathcal{T})} \Rightarrow (\sigma', h, v', \mathcal{T}')$$

Fig. 14. Semantics for our target language 2b. Here, his a history, v a list of arguments for every element in h, and  $\mathcal{T}$  a tuple of sets and mappings that keep track of all pending signaling and locking operations. Methods and predicates that appear in small caps are defined the text.

## <sup>1551</sup> C TARGET LANGUAGE OPERATIONAL SEMANTICS

This section presents the semantics of our target language presented in Figure 2b. As mentioned in Section 3, given an explicit monitor  $M_t$ , initial state  $\sigma$ , and monitor history  $h_e$  with argument mapping  $v_e$ , the operational semantics of  $M_t$  is defined using a judgment  $M_t \vdash (h_e, v_e, \sigma) \downarrow \sigma'$ indicating that the new state is  $\sigma'$  after executing  $h_e$  on initial state  $\sigma$ . The semantics of such a monitor are implemented using the inference rules of Figure 14 that use judgements of the form  $(\sigma, h_e, v_e, \mathcal{T}) \Rightarrow (\sigma', h'_e, v'_e, \mathcal{T}')$ . Here,  $\mathcal{T}$  and  $\mathcal{T}'$  are tuples of the form  $(\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L)$  and their role is to keep track all signaling and locking operations of the history. Specifically, each element of the tuple is defined below: 

- $\mathcal{B}_S \subseteq T \times CVar$ : a set of thread-condition variable pairs,  $(t, c) \in \mathcal{B}_S$  means that thread *t* is blocked on condition variable *c*.
- $N_S \subseteq T \times CVar$ : a set of thread-condition variable pairs,  $(t, c) \in N_S$  means that thread t has been notified on condition c.
- $\mathcal{H}_L \subseteq T \times L$ : a set of thread-lock pairs,  $(t, l) \in \mathcal{H}_L$  means that thread *t* holds lock *l*.
- 1566  $\mathcal{B}_L \subseteq T \times L$ : a set of thread-lock pairs,  $(t, l) \in \mathcal{B}_L$  means that thread t is blocked waiting to 1567 acquire lock l.

1:32

1569

1586 1587

1588

1589

1605

1606

1607

1608

1570	1: <b>p</b> 1	rocedure UpdateState( $\mathcal{T}, t$ )
1571	2:	<b>input:</b> $\mathcal{T} = (\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L)$
1572	3:	<b>input:</b> <i>t</i> , thread executing a non-synchronization statement.
1573	4:	output: updated mappings.
1574	5:	if $(t,c) \in \mathcal{B}_S$ then
1575	6:	$\mathcal{B}'_{S} \leftarrow \mathcal{B}_{S} \setminus \{(t,c)\} \ \mathcal{N}'_{S} \leftarrow \mathcal{N}_{S} \setminus \{(t,c)\}$
1576	7:	if $(t, l) \in \mathcal{B}_l$ then
1577	8:	$\mathcal{H}'_{l} \leftarrow \mathcal{H}_{L} \cup \{(t,l)\}$
1578	9:	$\mathcal{N}_{l}^{L} \leftarrow \mathcal{N}_{L} \setminus \{(t, l)\}$
1579	10:	$\mathcal{B}_{L}^{'} \leftarrow \mathcal{B}_{L} \setminus \{(t, l)\}$
1580	11.	return $(\mathcal{B}'_{\mathcal{A}}, \mathcal{N}'_{\mathcal{A}}, \mathcal{H}'_{\mathcal{A}}, \mathcal{B}'_{\mathcal{A}})$
1581	11.	$\mathcal{L}_{\mathcal{S}},\mathcal{L},\mathcal{L},\mathcal{L},\mathcal{L},\mathcal{L},\mathcal{L},\mathcal{L},$
1582		
1583		
1584		
1585		Fig. 15. Procedure UpdateState

Fig. 15. Procedure UpdateState

•  $N_L \subseteq T \times L$ : a set of thread-lock pairs,  $(t, l) \in N_L$  means that thread t can acquire a previously held lock l

We say that  $M_t \vdash (h_e, v_e, \sigma) \downarrow \sigma'$  if and only if  $(\sigma, h_e, v_e, \mathcal{T}) \Rightarrow^* (\sigma', \epsilon, \epsilon, \mathcal{T}')$ , where  $\Rightarrow^*$  is the 1590 reflexive transitive closure of relation  $\Rightarrow$ . In other words, a  $h_e$  is a valid explicit history according 1591 to our operational semantics only if the rules of Figure 14 can "consume" the entire history. If none 1592 1593 of the rules of Figure 14 apply to a history, then we consider the computation of relation  $\Rightarrow$  stuck and thus the history is not valid. 1594

1595 On a high level, the rules of our operational semantics iterate over all statements of input history 1596 h and updates the sets inside  $\mathcal{T}$  accordingly. Because during the execution of a h a thread t might 1597 perform a blocking operation (e.g., call 1.lock on a state where 1 is being held), the rules require 1598 every statement to be executed in a state where a thread is not blocked. To ensure this, every rule in Figure 14 requires predicate  $M_t \vdash (h_e, v_e, \sigma) \downarrow \sigma'$ , defined below, to hold for the executing thread 1599 1600 t.

CheckNotif
$$((\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L), t) = (t, c) \in \mathcal{B}_S \leftrightarrow (t, c) \in \mathcal{N}_S \land (t, l) \in \mathcal{B}_L \leftrightarrow (t, l) \in \mathcal{N}_L$$

Essentially, CHECKNOTIF( $(\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L), t$ ) requires every thread t that was previously blocked by some operation  $((t, \_) \in \mathcal{B}_S \text{ or } (t, l) \in \mathcal{B}_L)$  to be first notified  $((t, \_) \in \mathcal{N}_S \text{ or } (t, l) \in \mathcal{N}_L)$ in order to execute a statement.

In what follows, we explain each of the rules of Figure 14 in more detail.

*Rule (1).* This rule applies for all statements that are not a synchronization statement (i.e., lock 1609 or signal operation). Because the operational semantics for non-synchronization statements are 1610 well-studied, we assume the existence of an oracle  $\downarrow$  that give a statement s and its argument v, 1611 it returns the resulting monitor state  $\sigma'$ . Furthermore, because thread t could be blocked before 1612 executing statement s, this rule uses procedure UPDATESTATE (defined in Figure 15) to update the 1613 sets inside  $\mathcal{T}$  accordingly. Specifically, if thread t was blocked in some condition variable c, then 1614 procedure UPDATESTATE removes pair (t, c) from both  $\mathcal{B}_S$  and  $\mathcal{N}_S$  (recall that if t was blocked then 1615 it is guaranteed to be notified). Similarly, if thread t was blocked on some lock l, then procedure 1616 1617

1618 UPDATESTATE add the pair (t, l) to  $\mathcal{H}_L$  (i.e., now t holds lock l) and removes it from  $\mathcal{N}_L$  and  $\mathcal{B}_L$ 1619 (same as in the condition variable case).

*Rule (2).* This rule applies to all statements where a thread *t* attempts to acquire lock 1 that is currently held by another thread. In order to determine whether a lock is held by another thread, this rule makes use of predicate LOCKHELD defined as follows:

$$LockHeld((\mathcal{B}_{S}, \mathcal{N}_{S}, \mathcal{H}_{L}, \mathcal{B}_{L}, \mathcal{N}_{L}), t, l) = l \in \mathcal{H}_{L}[t']. t \neq t'$$

Then, the rule marks thread *t* as blocked on lock *l* by using the following procedure that updates map  $\mathcal{B}_L$  by adding pair (t, l):

 $BLOCKTHREADONLOCK((\mathcal{B}_{S}, \mathcal{N}_{S}, \mathcal{H}_{L}, \mathcal{B}_{L}, \mathcal{N}_{L}), t, l) = (\mathcal{B}_{S}, \mathcal{N}_{S}, \mathcal{H}_{L}, \mathcal{B}_{L}', \mathcal{N}_{L})$ where  $\mathcal{B}_{L}' = \mathcal{B}_{L} \setminus \{(t, l)\}$ 

Finally, the rule requires that the new attempt to acquire lock l does not introduce any deadlocks by invoking oracle NoDeadLocks. This oracle detects any cycles in the lock acquisition by examining maps  $\mathcal{B}_L$  and  $\mathcal{H}_L$ .

*Rule (3).* Conversely, the third rule applies to all cases where thread thread *t* attempts to acquire a lock not currently held by some other thread. In this case, the rule simply adds pair (t, l) in map  $\mathcal{H}_L$  as follows:

AcqLock
$$((\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L), t, l) = (\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}'_L, \mathcal{B}_L, \mathcal{N}_L)$$
  
where  $\mathcal{H}'_I = \mathcal{H}_L \cup \{(t, l)\}$ 

*Rule* (4). This rule is triggered when a thread *t* releases lock *l*. The rule performs the following two updates to maps  $\mathcal{H}_L$  and  $\mathcal{N}_L$ :

$$\operatorname{RelLock}((\mathcal{B}_{S}, \mathcal{N}_{S}, \mathcal{H}_{L}, \mathcal{B}_{L}, \mathcal{N}_{L}), t, l) = (\mathcal{B}_{S}, \mathcal{N}_{S}, \mathcal{H}_{L}', \mathcal{B}_{L}, \mathcal{N}_{L}')$$
  
where  $\mathcal{H}_{L}' = \mathcal{H}_{L} \setminus \{(t, l)\}, \ \mathcal{N}_{L}' = \mathcal{N}_{L} \cup \{(t', l)\} \text{ s.t. } (t', l) \in \mathcal{B}_{L}$ 

Specifically, it removes pair (t, l) from  $\mathcal{H}_{L}$  and notifies some thread t' currently blocked on lock l.

*Rule* (5). This rule applies when a thread t calls method await on a condition variable c. The rule simply adds pair (t, c) in set  $\mathcal{B}_S$ .

$$BLOCKONCVAR((\mathcal{B}_{S}, \mathcal{N}_{S}, \mathcal{H}_{L}, \mathcal{B}_{L}, \mathcal{N}_{L}), t, c) = (\mathcal{B}_{S}', \mathcal{N}_{S}, \mathcal{H}_{L}, \mathcal{B}_{L}, \mathcal{N}_{L})$$
  
where  $\mathcal{B}_{S}' = \mathcal{B}_{S} \setminus \{(t, c)\}$ 

*Rules (6) and (7).* These two rules are used when a thread signals or broadcasts a condition variable *c*. They simply update set  $N_S$  as follows:

SIGCVAR(
$$(\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L), c$$
) =  $(\mathcal{B}_S, \mathcal{N}'_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L)$   
where  $\mathcal{N}'_S = \mathcal{N}_S \cup \{(t, c)\}$  s.t.  $(t, c) \in \mathcal{B}_S$ 

BCASTCVAR(
$$(\mathcal{B}_S, \mathcal{N}_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L), c$$
) =  $(\mathcal{B}_S, \mathcal{N}'_S, \mathcal{H}_L, \mathcal{B}_L, \mathcal{N}_L)$   
where  $\mathcal{N}'_S = \mathcal{N}_S \cup \{(t, c) \mid (t, c) \in \mathcal{B}_S\}$ 

Specifically, rule 6 adds *a single* thread *t* currently blocked on condition variable *c* in  $N_S$ , whereas rule 7 adds *all* such threads in  $N_S$ .

*Rule (8).* Finally, rule 8 recursively applies the procedure to the whole input history *h*.

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:34

# 1667 D MONITOR INSTRUMENTATION.

In this section, we describe procedure Instrument which given an implicit-synchronization monitor M, it corresponding FGD  $\mathcal{G} = (V, E)$ , and a synchronization protocol  $\mathcal{S} = (\mathcal{L}, \mathcal{A}, \mathcal{P})$ , it instrument protocol  $\mathcal{S}$  into M yielding an explicit-synchronization monitor M' equivalent to M. This is achieved by first introducing all the necessary synchronization fields (locks, condition variables, and atomic fields) in the input class and then instrumenting locking and signaling operations in all methods as follows:

- Lock acquisition and release: The synthesized code must ensure that all the locks in  $\mathcal{L}(f)$  are held when executing fragment f. Thus, for every edge (f, f') in the FDG, we instrument the code to acquire locks  $\mathcal{L}(f') \setminus \mathcal{L}(f)$  and release locks  $\mathcal{L}(f) \setminus \mathcal{L}(f')$ . Furthermore, as mentioned in Section 2, we acquire and release these locks according to a static total order to prevent deadlocks.
- Blocking on predicates: Our instrumentation must also convert every waituntil statement to a sequence of operations on locks and condition variables. Specifically, we instrument a waituntil(p) statement as follows:

while(!p) { ln.unlock(); ...; l2.unlock(); c.await(); l2.lock(); ...; ln.lock(); }

where c is the condition variable associated with p; 11, ... In are the locks associated with this fragment, and 11 is the lock associated with condition variable c.

- Signaling operations: Finally, we instrument signaling operations introduced by PlaceSignals to acquire and release the appropriate locks. In particular, given a statement signal(p,c) (similarly for broadcast(p,c)), our instrumentation generates the following code:
- 1687 1688

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1681

1682

1683

if (c) { lp.lock(); cp.signal(); lp.release(); }

where cp is the condition variable for predicate p and 1p is the corresponding lock for cp.

Procedure Instrument is presented in Figure 16 in the form of inference rules that use the following two judgements:

- ν ⊢ Δ → Δ', where ν is a subset of the arguments of procedure Instrument (we overload operator → depending on the arguments) and Δ is one of the following: the input monitor, a field, a method, a CCR or a statement.
  - $\mathcal{L}, \mathcal{G} \vdash v \hookrightarrow v'$ , where  $\mathcal{L}$  is the lock map of the input synchronization protocol  $\mathcal{S}$ , FDG is the input  $\mathcal{G}$ , and v is a fragment in  $\mathcal{G}$ .

The meaning of each judgement is that whenever procedure Instrument is applied to an element that appears on the left-hand side of an arrow ( $\rightsquigarrow$ ,  $\hookrightarrow$ ) it generates the element on the right-hand side.

**Overall Structure.** The core logic of this procedure is to recursively iterate every element of the input monitor and use the inferred synchronization protocol in order to convert each element to an equivalent element of the target language. At a top level, the procedure begins by transforming every field and method of the input monitor. For every method, the procedure recursively visits every CCR using operator  $\rightsquigarrow$ . Then, for every CCR, it collects all its fragments and uses operator  $\hookrightarrow$  to instrument all the lock operations dictated by the input protocol. In what follows, we give a brief description of every rule presented in Figure 16.

MTR. This is the top-level rule called by procedure Instrument and performs the following
 tasks: 1. it introduces all the synchronization fields (locks and condition variables) needed by the
 synchronization protocol and initializes them accordingly and 2. it recursively calls itself for every
 field, and method of *M*.

 $F_{LD-1} & F_{LD-2}$ . These two rules are used to translate fields of *M*, with the first one being applicable to fields that must be converted to atomic fields and the second one to fields that should remain the

same. Only the first rule alters the original field by converting to an atomic field with the samename as the original.

 $\begin{array}{ll} & & \\ \hline & METHOD \ \& \ CCR. \end{array} \ These two rules simply recursively apply operator <math>\rightsquigarrow$  to their constituent elements. \\ \end{array}

1726 FRAG-STMT. This rule applies to all statements *s* that are a fragment in the input  $\mathcal{G}$ . It first uses 1727 operator  $\hookrightarrow$  to instrument all necessary lock operations and then uses a special oracle  $\rightarrow_{\mathcal{A}}$  that 1728 converts all operations that involve a field converted to atomic to the equivalent update statement 1729 in the target language.<sup>8</sup>

1730 WAIT. Rule labeled WAIT is a special case of the above rule because, by definition, every waituntil 1731 statement defines its own fragment in an FDG. Similar to the rule above, this rule also uses operator 1732  $\hookrightarrow$  to instrument the appropriate lock operation in the fragment but it additionally translates 1733 the waituntil statement into an equivalent statement in the target language that uses condition 1734 variables. As mentioned in Section 4, each waituntil statement is translated into a while loop that 1735 waits on the appropriate condition variable and properly releases and acquires all locks before and 1736 after the call to method await. Additionally, it acquires all locks needed by its successor statement 1737 v in the FDG and releases all locks held by it but not needed by v (similar to the logic described 1738 below). 1739

1740 SIG. This rule applies to all fragments that are a signalling directive of the monitor's intermediate 1741 representation.<sup>9</sup> In a similar manner as the rule for waituntil statements, this rule first uses operator 1742  $\hookrightarrow$  to instrument all lock operations needed to implement the synchronization protocol. Then, it 1743 consults the predicate map  $\mathcal{P}$  of the synchronization protocol to acquire the appropriate lock and 1744 perform the signaling operation on the associated condition variable.

17451746Instrumenting Fragments With Lock Operations. Finally, we describe operator  $\hookrightarrow$  which1747given a fragment v, the lock map  $\mathcal{L}$  of the input synchronization protocol S, and the FDG  $\mathcal{G}$ , it1748instruments all the necessary lock operations. The logic of this operator is split between two groups1749of rules, described in more detail below:

- Rules for entry & exit fragments (i.e., fragments without predecessors and successors respectively), which are handled by rules ENTRY-FRAG and EXIT-FRAG respectively. These rules simply lookup the entry or exit fragment in  $\mathcal{L}$  and acquire or release the locks returned by the  $\mathcal{L}$  accordingly.
- Rules for fragments with successors. Fragments that contain some successor in the graph are handled by rules BRANCH-FRAG, REG-FRAG-1 and REG-FRAG-2. The logic for each of these rule is similar, i.e., for any successor fragment  $v_s$  of fragment v, the instrumentation releases all locks required by v but not by  $v_s$  ( $\mathcal{L}[v] \setminus \mathcal{L}[v_s]$ ) and acquires all locks required by  $v_s$  but not v ( $\mathcal{L}[v_s] \setminus \mathcal{L}[v]$ ). All these operation are operation in accordance to the global lock order to prevent deadlocks. Last, it is worth mentioning that the main difference of these three rules is how they instrument the edge between v and  $v_s$ . That is, if v ends with a goto statement

1764

1761

1750

1751

1752

<sup>&</sup>lt;sup>1762</sup> <sup>8</sup>Due to its simplicity, we omit a formal description of oracle  $\rightarrow_{\mathcal{A}}$ .

<sup>&</sup>lt;sup>1763</sup> <sup>9</sup>For simplicity, we assume that every signaling operation defines its own fragment.

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:37

(conditional or not), then the instrumentation redirects the control-flow appropriately so all lock operations occur along edge  $(v, v_s)$ .

Finally, we conclude with the following theorem that states the correctness of our instrumentation
 phase.

THEOREM D.1. Let  $S = (\mathcal{L}, \mathcal{A}, \mathcal{P})$  be a synchronization protocol inferred over FGD  $\mathcal{G} = (V, E)$ of input monitor M and M' be the result of procedure Instrument for M. Then, the following three conditions hold:

- (1) For every fragment  $v \in V$ ,  $l_i \in \mathcal{L}[v]$  iff fragment v holds lock  $l_i$  in M'
- (2) If i < j, then  $l_i$  is never acquired whenever  $l_j$  is held.

(3) Field  $f \in \mathcal{A}$  iff all its occurrences in M have been replaced with an atomic operation in M'.

PROOF. Proof can be find in Appendix E.

## E CORRECTNESS-RELATED PROOFS

This section contains all the proofs related to the correctness of our approach. Section E.1 presents the proof of Theorem 4.5, Section E.2 presents the proof of Theorem 4.14, Section E.3 presents the proof of Theorem 4.16, and Section E.4 presents the proof of Theorem D.1.

#### 1785 E.1 Proof of Theorem 4.5

Theorem 4.5 states the following:

(Correctness). We say that an explicit monitor  $M_t$  correctly implements an implicit monitor  $M_s$ , denoted as  $M_s \sim M_t$ , iff for all input states  $\sigma_s, \sigma_t$  s.t.  $\sigma_s \equiv_{M_s} \sigma_t$ , we have:

(1) $\forall h_i, v_i. \ M_s \vdash (h_i, v_i, \sigma_s) \Downarrow \sigma'_s \Longrightarrow$	$\left(M_t \vdash (Expand_{M_t}(h_i, \nu_i, \sigma_s), \sigma_t) \downarrow \sigma'_t \land \sigma'_s \equiv_{M_s} \sigma'_t\right)$
--	--

 $(2) \ \forall h_e, v_e. \ M_t \vdash (h_e, v_e, \sigma_t) \downarrow \sigma'_t \Longrightarrow (\exists h_i, v_i. \ (h_e, v_e) \backsim (h_i, v_i) \land M_s \vdash (h_i, v_i, \sigma_s) \Downarrow \sigma'_s \land \sigma'_s \equiv_{M_s} \sigma'_t)$ 

PROOF. For all proofs in this Section, we assume the correctness of procedure PLACESINGALS (proved in previous work [Ferles et al. 2018]).

The proof of condition (1) above follows directly from Theorems 4.16, D.1, and the correctness of procedure PLACESINGALS. The proof of condition (2) follows directly from Theorem 4.14, Theorem 4.16, Theorem D.1, and correctness of PLACESIGNALS.

#### E.2 Proof of Theorem 4.14

In this section, we present the proof of Theorem 4.14 which we reiterate here for convenience.

Before presenting the actual proof, we first introduce some auxiliary notation, relations, and 1802 lemmas. First, given a history h, we define a predicate  $h[[(v_1, t_1)_{i_1}, \dots, (v_k, t_k)_{i_k}]]$  that evaluates to 1803 true iff each event  $(v_i, t_i)_i$  is *i*-th element in h and  $i_1 < \ldots < i_k$ . That is, this predicate encodes 1804 that these event occur in this particular order within h. Second, we define Next(h, i, t) as follows: 1805  $min(\{j \mid j > i, h_G[[(, t)_i]]\})$ . In other words, Next(h, i, t) returns the first element in h after 1806 index *i* whose thread identifier is *t*. Additionally, we use h[i] to denote the *i*-th element in *h* and 1807 h[i:j], i < j to denote the "sub-history" of *h* between its *i*-th element (inclusive) and *j*-th element 1808 (exclusive). Finally, we extend the definition of a history projection to filter out elements that do 1809 *not* involve a thread , e.g.,  $\Pi(h, \neg t)$  filters out all events of *h* that involve thread *t*. 1810

Next, using the notation above we define some relations that identify interleavings inside a history of a fragmented monitor  $M_{\mathcal{G}}$ .

1777

1778 1779

1780

1784

1786 1787

1788

1789 1790 1791

1792 1793

1794

1795

1796

1797

1798 1799

1800

Kostas Ferles, Benjamin Sepanski, Rahul Krishnan, James Bornholt, and Isil Dillig

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:38

Definition E.1. (History Interleaving). Given history  $h_{fdg} = (V, E)$  and interleaving  $\chi = (v, e = (v_s, v_t))$ . We define the occurrences of  $\chi$  as follows:

$$Interleavings(\chi, h_{\mathcal{G}}) = [(j, (i, k)) \mid h_{\mathcal{G}} \llbracket (v_s, t)_i, (v, t')_j, (v_t, t)_k \rrbracket, t \neq t', k = Next(h_{\mathcal{G}}, i, t)]$$

Also, we write  $X(h_G)$  to be the set of all interleavings that occur in  $h_G$ , i.e.

$$X(h_G) = \{ \chi \mid X = Interleavings(\chi, h_G), |X| > 0 \}$$

Finally, we write  $X_{\#}(h_{\mathcal{G}})$  to denote the *number* of interleavings inside *h*. Formally:

$$\mathcal{X}_{\#}(h_{\mathcal{G}}) = \sum_{\chi \in V \times E} |Interleavings(\chi, h_{\mathcal{G}})|$$

Next, given a fragment *v* we assume the existence of two predicates, namely, *EntryFrag* and *ExitFrag*, that hold only if *v* is the entry fragment or the exit fragment of its CCR respectively. Based on these relations, we define the next relation that partitions a fragment history into CCR sub-histories.

*CCR History Partition.* Let  $h_{\mathcal{G}}$  be a history of fragments in FDG  $\mathcal{G}$ . We define the CCR partition of  $h_{\mathcal{G}}$  that returns a list of potentially overlapping sub-histories of  $h_{\mathcal{G}}$  as follows:

$$CCRPart(h_{\mathcal{G}}) = \begin{bmatrix} h_{\mathcal{G}}[i:j] \mid & h_{\mathcal{G}}[i] = (v_{in}, t, \_), \ EntryFrag(v_{in}), \\ j = min(\{k \mid k > i, \ h_{\mathcal{G}}[k+1] = (v_{out}, t, \_), ExitFrag(v_{out})\}) \end{bmatrix}$$

Let  $P = CCRPart(h_G)$ , then we use P[i] to refer to the *i*-th sub-history in *P*. Note we assume that partitions returned by *CCRPart* are ordered according to the index of the first element in the sub-history. That is, if *ccr*<sub>1</sub> began its execution before *ccr*<sub>2</sub> in  $h_G$ , then the partition of *ccr*<sub>1</sub> appears before the partition of *ccr*<sub>2</sub> in *P*. Furthermore, given a CCR partition P[i], we write *Thread*(P[i]) to represent the thread of the first element in sub-history P[i].

**Removing Interleavings from Histories**. Before we prove our main theorem, we define some transformations on interleaved histories that helps us remove interleavings.

First, given a CCR sub-history that is interleaved, we define its sequential history as follows:

Definition E.2. Sequential CCR Sub-history Let  $h_{\mathcal{G}}$  be an interleaved history,  $P = CCRPart(h_{\mathcal{G}})$ , and  $h'_{\mathcal{G}} = P[i]$  be an interleaved CCR partition. We define that the sequential history of  $h'_{\mathcal{G}}$ , denoted as  $Seq_{|CCR}(h'_{\mathcal{G}})$  to be the following history:  $\Pi(h'_{\mathcal{G}}, t)\Pi(h'_{\mathcal{G}}, \neg t)$ , where t = Thread(P[i]). Given an argument mapping  $\nu$  for history  $h'_{\mathcal{G}}$ , we write  $Seq(\nu)$  to denote the corresponding argument mapping for  $Seq_{|CCR}(h'_{\mathcal{G}})$ .

We now prove the following useful lemmas about sequential CCR sub-histories.

LEMMA E.3. Let  $h_{\mathcal{G}}$  be an interleaved sub-history and  $h'_{\mathcal{G}} = Seq_{|CCR}(h_{\mathcal{G}})$ , then the following two things hold:

$$(1) X(h'_{\mathcal{G}}) \subseteq X(h_{\mathcal{G}})$$
$$(2) X_{\#}(h'_{\mathcal{G}}) < X_{\#}(h_{\mathcal{G}})$$

PROOF. Both of the properties logically follow from the construction of  $h'_{\mathcal{G}}$ . That is, a sequential CCR sub-history of the form  $(v_1, t) \dots (v_i, t), (v_{i+1}, t'), \dots, (v_j, t'')$ , where all elements before the *i*-th position are from thread *t* and all elements past that are from some thread *t'* s.t.  $t \neq t'$ . On the other hand, the original history  $h_{\mathcal{G}}$  is of the form:

$$(v_1, t) \dots (v_k, t) (v_{k+1}, t') \dots (v_j, t)$$

Where  $h_{\mathcal{G}}[1:k+1] = h'_{\mathcal{G}}[1:k+1]$  (i.e.,  $h_{\mathcal{G}}$  and  $h'_{\mathcal{G}}$  have a common prefix). Therefore, since by its construction  $h'_{\mathcal{G}}$  does not move the relative order of element is  $h_{\mathcal{G}}$  that do not involve thread t,

1:40

if an interleaving  $\chi \in \mathcal{X}(h'_{\mathcal{G}})$  then we also have  $\chi \in \mathcal{X}(h_{\mathcal{G}})$ . Conversely, any thread interleaving that involved thread t in  $h_{\mathcal{G}}$  does not appear in  $h_{\mathcal{G}}$  (by construction). Since by its definition the interleaved history  $h_{\mathcal{G}}$  contains at least one interleaving that involves an edge executed by t, we can conclude that  $\mathcal{X}_{\#}(h'_{\mathcal{G}}) < \mathcal{X}_{\#}(h_{\mathcal{G}})$ .

LEMMA E.4. Let  $h_{\mathcal{G}}$  be an interleaved sub-history and  $h'_{\mathcal{G}} = Seq_{|CCR}(h_{\mathcal{G}})$ , then if  $X(h_{\mathcal{G}})$  is a set of strongly safe interleavings we have that:  $\forall \sigma, v.M_{\mathcal{G}} \vdash (h_{\mathcal{G}}, v, \sigma) \Downarrow \sigma' \Rightarrow M_{\mathcal{G}} \vdash (h'_{\mathcal{G}}, Seq(v), \sigma) \Downarrow \sigma'$ 

**PROOF.** We prove this by induction on the number of distinct interleavings of history  $h_G(X_{\#}h_G)$ .

*Base Case:*  $X_{\#}(h_{\mathcal{G}}) = 1$ . If there is a single interleaving in  $h_{\mathcal{G}}$ , this implies that  $h_{\mathcal{G}}$  is of the form:

$$(v_1, t) \dots (v_i, t') \dots (v_j, t)$$

where  $t = Thread(h_{\mathcal{G}})$  and  $(v_i, t')$  is the only element in  $h_{\mathcal{G}}$  not executed by t. Because, the interleaving of  $h_{\mathcal{G}}$  is strongly safe, we have that fragment  $v_i$  executed by t', right commutes with any possible successor of the edge it interleaves. Also, by definition,  $h'_{\mathcal{G}}$  is  $(v_1, t) \dots (v_j, t)(v_i, t')$ . Combining this two facts with lemma E.3, we can prove the theorem for our base case:

$$\forall \sigma, v.M_{\mathcal{G}} \vdash (h_{\mathcal{G}}, v, \sigma) \Downarrow \sigma' \to M_{\mathcal{G}} \vdash (h'_{\mathcal{G}}, Seq(v), \sigma) \Downarrow \sigma'$$

Inductive Step. In our inductive step, we assume that our lemma holds for  $X_{\#}(h_{\mathcal{G}}) = n$  and we are going to prove it for n + 1. The logic is similar to the base case, specifically, we get the right-most interleaved fragment in  $h_{\mathcal{G}}$  and right-commute to the end of the history while obtaining a semantically equivalent history  $h'_{\mathcal{G}}$ . After that, we can apply our inductive hypothesis on  $h'_{\mathcal{G}}$ , which again proves our goal.

Finally, we prove our main theorem, which we re-iterate below for convenience.

**Theorem 4.14.** Let  $\mathcal{G}$  be an FDG and let  $\chi_1, \ldots, \chi_n$  be strongly safe interleavings. Then,  $S = \{\chi_1, \ldots, \chi_n\}$  is a safe interleaving set for  $\mathcal{G}$ .

PROOF. By definition of safe set of interleavings, we have to prove the following for every interleaved history  $h_G$  of monitor  $M_G$ 

If 
$$X(h_G) \subseteq S$$
 and  $M_G \vdash (h_G, v_G, \sigma) \Downarrow \sigma'$  then  $\exists h, v. (h_G, v_G) \backsim (h, v)$  and  $M \vdash (h, v, \sigma) \Downarrow \sigma'$ 

In order to prove that, we have to prove that for every interleaved history  $h_{\mathcal{G}}$  that only allows interleavings in *S* and argument mapping  $v_{\mathcal{G}}$  we can find a history of the original monitor with corresponding argument mapping s.t.,  $((h_{\mathcal{G}}, v_{\mathcal{G}}) \sim (h, v))$ . Which in turn means that we have to find a sequential history of  $M_{\mathcal{G}} h'_{\mathcal{G}}$  s.t.

(1) 
$$\forall t. \pi(h_G, t) = \pi(h'_G, t)$$
 and (2)  $\mathsf{Expand}_{M_G}(h, v, \sigma) = (h'_G, v_G)$ 

To prove the goal above, we start with an arbitrary interleaved history  $h_{\mathcal{G}}$  s.t.  $X(h_{\mathcal{G}}) \subseteq S$  and convert it to a sequential history  $h'_{\mathcal{G}}$  with the above properties. We perform this proof, by first creating the CCR partition of  $h_{\mathcal{G}}$ ,  $P = CCRPart(h_{\mathcal{G}})$ , and then induct on the number of partitions in *P* that *are* interleaved.

**Base Case: One interleaved CCR in P.** Let  $h_{\mathcal{G}}^{ccr} = P[i]$  be the interleaved history in  $h_{\mathcal{G}}$ . Now, let  $h'_{\mathcal{G}} = h_{\mathcal{G}}[Seq_{|CCR}(h_{\mathcal{G}}^{ccr})/h_{\mathcal{G}}^{ccr}]$ . Because  $h_{\mathcal{G}}^{ccr}$  is the only interleaved sub-history in P and because of lemma E.3, we have that  $h'_{\mathcal{G}}$  is a sequential history s.t.  $\forall t, v_{\mathcal{G}}. \pi(h_{\mathcal{G}}, t) = \pi(h'_{\mathcal{G}}, t)$ . Furthermore, because of lemma E.4 we have  $\forall \sigma, v.M_{\mathcal{G}} \vdash (h_{\mathcal{G}}, v_{\mathcal{G}}, \sigma) \Downarrow \sigma' \Rightarrow M_{\mathcal{G}} \vdash (h'_{\mathcal{G}}, Seq(v_{\mathcal{G}}), \sigma) \Downarrow \sigma'$ . Finally, because we have  $M_{\mathcal{G}} \vdash (h_{\mathcal{G}}, v_{\mathcal{G}}, \sigma) \Downarrow \sigma'$  for some  $\sigma$ , this implies that  $\mathsf{Expand}_{M_{\mathcal{G}}}(h, v_{\mathcal{G}}, \sigma) = (h'_{\mathcal{G}}, Seq(v_{\mathcal{G}}))$ , which in turns implies  $(h'_{\mathcal{G}}, v_{\mathcal{G}}) \sim (h, Seq(v_{\mathcal{G}}))$ .

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1935

1936

1937

1938 1939

1940

1941 1942

1943 1944

1945

1946 1947 1948

**Inductive Step.** Next, we assume that our theorem holds for up to *n* interleaved CCRs in *P*, and will prove it for n + 1. Similarly as above, we find the smallest *i* s.t.  $P[i] = h_{\mathcal{G}}^{ccr}$  is an interleaved history. Again, we construct  $h'_{\mathcal{G}} = h_{\mathcal{G}}[Seq_{|CCR}(h_{\mathcal{G}}^{ccr})/h_{\mathcal{G}}^{ccr}]$ . Because of lemma E.3, we have that the number of interleaved histories in  $CCRPart(h'_{\mathcal{G}})$  has strictly fewer number of interleaved sub-histories than *P*. Therefore, by our inductive hypothesis, we have that  $(h'_{\mathcal{G}}, v'_{\mathcal{G}} \sim (h, Seq(v_{\mathcal{G}})))$ for some history of *h* of *M*. This, combined with lemma E.4, proves that  $(h_{\mathcal{G}}, v_{\mathcal{G}}) \sim (h, Seq(v_{\mathcal{G}}))$ .

### 1969 E.3 Proof of Theorem 4.16

1972

1973 1974 1975

1978

1979

1980

1981

1982

1983

1984

1985

1987 1988

1989

1990

1991

1992 1993

1994

1995 1996

1997

1998

1999

2000

2001

2002

2005

2006

2007

We now prove theorem 4.16 which states the correctness of our MaxSAT encoding.

**Theorem 4.16.** Let *m* be a model of the generated MaxSAT instance and  $(\mathcal{L}, \mathcal{A}, \mathcal{P})$  be the synchronization protocol constructed as follows:

$$\mathcal{L} = \left\{ v \mapsto \left\{ l \mid m[h_v^l] \right\} \right\} \ \mathcal{A} = \left\{ \mathsf{fld} \mid m[a_{fld}] \right\} \mathcal{P} = \left\{ p \mapsto l_i \mid \mathsf{IsWait}(v, p), i = \min(\{j \mid m[h_v^{l_j}]\}) \right\}$$

where, IsWait(v, p) is true if v is a waituntil statement on p. Then,  $(\mathcal{L}, \mathcal{A}, \mathcal{P})$  is a correct synchronization protocol.

PROOF. As mentioned earlier, a synchronization protocol must meet the following correctness criteria:

- (1) If two fragments  $v_1, v_2$  have a race (i.e.,  $\mathcal{R}(v_1, v_2) \neq \emptyset$ ), then the protocol must prevent this race with a lock or an atomic field.
- (2) If a fragment interleaving  $\chi = (v, e)$  is not safe, then the synchronization protocol must not allow fragment *v* to execute in between edge *e*.
  - (3) The protocol must be deadlock-free.

We show that, by construction, a model m returned by a MaxSAT solver always satisfies the above conditions.

- (1) Model *m* prevents any races between two fragments because it must satisfy *all* hard constraints generated by rules RACE-1 and RACE-2 from Figure 8. Therefore, *m* will force two racy fragments to either share a lock or, when possible, convert all operations involving the racy field to equivalent atomic ones.
- (2) Similarly, because model *m* must satisfy the hard constraints generated by rule I-LEAVE, any interleaving that was deemed unsafe by our static analysis is guaranteed to be infeasible in the resulting synchronization monitor.
  - (3) Finally, because of rules WAIT and L-ORDER, the resulting synchronization monitor is guaranteed to be deadlock-free. Specifically, the hard constraints generated by rule L-ORDER enforce the invariant that all lock acquisitions respect the global lock order. Whereas, the hard constraints of rule WAIT, enforce the same invariant for the translation of a waituntil statement into an equivalent statement in the target language (see Figure 16).

#### E.4 Proof of Theorem D.1

Finally, we prove the correctness of our monitor instrumentation procedure (Fig. 16).

THEOREM E.5. Let  $S = (\mathcal{L}, \mathcal{A}, \mathcal{P})$  be a synchronization protocol inferred over FGD  $\mathcal{G} = (V, E)$ of input monitor M and M' be the result of procedure Instrument for M. Then, the following three conditions hold:

(1) For every fragment  $v \in V$ ,  $l_i \in \mathcal{L}[v]$  iff fragment v holds lock  $l_i$  in M'

- 2010 (2) If i < j, then  $l_i$  is never acquired whenever  $l_j$  is held.
  - (3) Field  $f \in \mathcal{A}$  iff all its occurrences in M have been replaced with an atomic operation in M'.

**PROOF.** All three conditions can be proved by providing certain guarantees for a subset of the rules of Figure 16. Note that operator  $\rightsquigarrow$  (Figure 16) is guaranteed to visit every code fragment  $v \in V$  in the FDG, since it recursively visits every element of the input monitor until it discovers *all* fragments of the given FDG. Next, we prove all three conditions.

2017 Condition (1): For this condition, we need to prove that both fragments will only hold the locks 2018 required by the synthesized protocol S. The logic of this proof depends on the number and type of 2019 predecessors of fragment v. We now present a case analysis:

Zero predecessors. This is the case of an entry fragment of a method in *G*. Due to the structure of our input language and the definition of an FDG, this fragment *must* be a fragment that contains a single waituntil statement. The instrumentation of such a fragments is handled by rules WAIT and ENTRY-FRAG. Note, that rule WAIT first calls ENTRY-FRAG which acquires all locks needed by the fragmented defined by the waituntil fragment.

At least one predecessor. These types of fragments are handled by rules WAIT, BRANCH-FRAG, REG-FRAG-1, and REG-FRAG-2. All these rules maintain the following invariant for the fragment vthat triggers them: before transferring control to any of v's successor, they release all locks needed by v but not needed by the successor (locksets of the form  $R_i$ ) and acquire all locks needed by the successor but not held by v (locksets of the form  $A_i$ ). This invariant combined with the fact that these are the only ways to transfer control flow in our input language, ensure that before executing a fragment in the output monitor all necessary locks (and only those) will be acquired.

*Condition (2):* This directly follows from:

- (1) That procedure Instrument uses auxiliary relation Acq to instrument lock acquisitions, which as shown in Figure 16 does so in increasing order of lock indices.
- (2) The guarantee provided by Theorem 4.16 that the synthesized protocol acquires locks in increasing order along every control-flow edge.

Condition (3): This condition is ensured by rule FRAG-STMT of Figure 16 that ensures oracle  $\rightarrow_{\mathcal{R}}$  is called on every fragment of  $\mathcal{G}$ .