# Automated Synthesis of Asynchronizations

SIDI MOHAMED BEILLAHI, Université de Paris, France

AHMED BOUAJJANI, Université de Paris, France

CONSTANTIN ENEA, Université de Paris, France

SHUVENDU LAHIRI, Microsoft Research Lab - Redmond, USA

Asynchronous programming is widely adopted for building responsive and efficient software. Modern languages such as C# provide async/await primitives to simplify the use of asynchrony. However, the use of these primitives remains error-prone because of the non-determinism in their semantics. In this paper, we propose an approach for refactoring a given sequential program into an asynchronous program that uses async/await, called asynchronization. The refactoring process is parametrized by a set of methods to replace with given asynchronous versions, and it is constrained to avoid introducing data races. Since the space of possible solutions is exponential in general, we focus on characterizing the delay complexity that quantifies the delay between two consecutive and distinct outputs. We show that this is polynomial time modulo an oracle for solving reachability (assertion checking) in sequential programs. We also describe a pragmatic approach has been implemented and evaluated on a number of non-trivial C# programs extracted from open-source repositories.

## 1 INTRODUCTION

1 2 3

4

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

40

41

42

43

Asynchronous programming is widely adopted for building responsive and efficient software. Unlike synchronous procedure calls, asynchronous procedure calls may run only partially and return the control to their caller. Later, when the callee finishes execution, a callback procedure registered by the caller is invoked.

As an alternative to the tedious model of asynchronous programming that required explicitly registering callbacks with asynchronous calls, C# 5.0 [Bierman et al. 2012] introduced the async/await primitives. These primitives allow the programmer to write code in a familiar sequential style without explicit callbacks. An asynchronous procedure, marked by the keyword async, returns a task object that the caller uses to "await" it. Awaiting may suspend the execution of the caller, if the awaited task did not finish, but does not block the thread it is running on. The code after the await is the continuation that is automatically called back when the callee result is ready. For instance, on the right of Figure 1, the method ContentLength calls an asynchronous method GetStringAsync that returns a task object t5 used to await it at line 24. Executing this await suspends the execution of ContentLength and returns the control to its caller MainAsync, assuming that GetStringAsync did not finish. Passing the control to the caller is a constraint of the await semantics. When GetStringAsync finishes and the thread is idle, the continuation after line 24 is scheduled. This paradigm has become popular across many languages, eg, C++, JavaScript, Python.

While simplifying the writing of asynchronous programs, the async/await primitives introduce concurrency which is notoriously complex. Depending on the scheduler, the code in between a call and a matching await (referring to the same task) may execute before some part of the awaited task (if the latter passed the control to its caller before finishing), or after the awaited task finished. For instance, on the right of Figure 1, the assignment x=0 in MainAsync between the call and the await of ReadFile may execute before ReadFile finishes, if the await in ReadFile suspends

Authors' addresses: Sidi Mohamed Beillahi, Université de Paris, France, beillahi@irif.fr; Ahmed Bouajjani, Université de
 Paris, France, abou@irif.fr; Constantin Enea, Université de Paris, France, cenea@irif.fr; Shuvendu Lahiri, Microsoft Research
 Lab - Redmond, USA, shuvendu@microsoft.com.

48 https://doi.org/

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

```
50
        1 void Main(string url0) {
                                                                    1 async Task MainAsync(string url0) {
        2 string url1 = ReadFile("url1.txt");
                                                                    2 Task<string> t1 = ReadFile("url1.txt");
51
        3 x = 0:
                                                                    3 \times = 0:
52
        4 int val0 = ContentLength(url0);
                                                                   4 Task<int> t2 = ContentLength(url0);
53
                                                                   5 string url1 = await t1:
                                                                   6 int val0 = await t2;
54
        7 int val1 = ContentLength(url1);
                                                                   7 Task<int> t3 = ContentLength(url1);
55
                                                                   8 int val1 = await t3;
56
        9 int r = x;
                                                                   9 int r = x:
        10 Debug.Assert(r == val0 + val1);
                                                                   10 Debug.Assert(r == val0 + val1);
57
        11 }
                                                                   11 }
58
59
        13 string ReadFile(string fn) {
                                                                   13 async Task<string> ReadFile(string fn) {
        14 StreamReader reader = new StreamReader(fn);
                                                                   14 StreamReader reader = new StreamReader(fn);
60
        15 string content = reader.ReadToEnd();
                                                                   15 Task<string> t4 = reader.ReadToEndAsync();
61
                                                                   16 string content = await t4;
62
        17 return content;
                                                                   17 return content;
        18 }
                                                                   18 }
63
64
        20 int ContentLength(string url) {
                                                                   20 async Task<int> ContentLength(string url) {
65
        21 HttpClient clt = new HttpClient()
                                                                   21 HttpClient clt = new HttpClient();
        22 string urlContents = clt.GetString(url);
                                                                   22 Task<string> t5 = clt.GetStringAsync(url);
66
                                                                   23 int r1 = x;
        23 int r1 = x;
67
                                                                   24 urlContents = await t5;
68
        25 x = r1 + urlContents.Length;
                                                                   25 x = r1 + urlContents.Length;
        26 return urlContents.Length;
                                                                   26 return urlContents.Length;
69
        27 }
                                                                   27 }
70
```

Fig. 1. A synchronous C# program and an asynchronous refactoring (x is a static variable).

its execution and passes the control to MainAsync (when reaching the await at line 16 because ReadToEndAsync did not finish), or after ReadFile finishes, otherwise (the call to ReadToEndAsync finishes before reaching its matching await at line 16). The resemblance with sequential code can be especially deceitful since this non-determinism is opaque. It is common that await instructions are placed immediately after the corresponding call which limits the benefits that one can obtain from executing code in the caller concurrently with code in the callee [Okur et al. 2014].

In this paper, we address the problem of writing efficient asynchronous code that uses async/await primitives. We propose a procedure for automated synthesis of asynchronous programs equivalent to a given synchronous (sequential) program *P*. This can be seen as a way of refactoring synchronous code to asynchronous code. Since solving this problem in its full generality would require checking equivalence between arbitrary programs, which is known to be hard, we consider a restricted space of asynchronous program candidates that are defined by substituting synchronous methods in *P* with asynchronous versions (assumed to be behaviorally equivalent). The substituted methods are supposed to be leaves of the call-tree, i.e., they do not call any other method in *P*. Such programs are called *asynchronizations* of *P*. A practical instantiation of this problem is replacing IO synchronous calls for, e.g., reading/writing files, managing http connections, with asynchronous versions.

For instance, let us consider the sequential C# program on the left of Fig. 1. The Main invokes ReadFile and ContentLength in a synchronous way (using the standard call stack semantics). ReadFile reads and returns the content of a file while ContentLength returns the length of the text in a webpage. The two URLs given as input to ContentLength are given as input to Main or read from some file using ReadFile. The program uses a variable x to aggregate the lengths of all pages accessed by ContentLength. Note that this program passes the assertion at line 10.

97 98

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

1:2

The time-consuming primitives for reading files, StreamReader.ReadToEnd, or the content of a webpage, HttpClient.GetString<sup>1</sup>, are an obvious choice for being replaced with equivalent *asynchronous* counterparts, i.e., StreamReader.ReadToEndAsync and HttpClient.GetStringAsync, respectively. Performing such tasks asynchronously can lead to significant boosts in performance.

The program on the right of Fig. 1 is an example of an asynchronization of the program on the 103 left where the calls to StreamReader.ReadToEnd and HttpClient.GetString are replaced with 104 asynchronous counterparts (assumed to have the same effect). The syntax of async/await imposes 105 that every method that transitively calls one of the substituted methods must also be declared to be 106 asynchronous, e.g., MainAsync and ContentLength. Then, an asynchronous call must be followed 107 by an await statement that specifies the control location where that task should have completed 108 (e.g., the return value should have been computed). For instance, the call to ReadToEndAsync at 109 line 15 is immediately followed by an await since the next instruction (at line 17) uses the value 110 computed by ReadToEndAsync. Therefore, synthesizing such refactoring boils down to finding a 111 correct placement of awaits for every method that transitively calls a substituted method (we do 112 not consider "deeper" refactoring like rewriting conditionals or loops). 113

We consider an equivalence relation between a synchronous program and an asynchronization 114 that corresponds to absence of data races in the asynchronization. Data race free asynchronizations 115 are called sound. Relying on absence of data races instead of a more precise equivalence relation like 116 117 equality of reachable sets of states could prevent enumerating some number of asynchronizations that reach the same set of states as the original synchronous program. However, checking equality 118 of reachable sets of states is known to be hard in general, and relying on absence of data races is 119 a well established compromise. For instance, the asynchronization on the right of Fig. 1 is sound 120 (data-race free) since the two calls to ContentLength that access x do not "overlap" and both finish 121 before the read at line 9. The accesses to x in the asynchronization are performed in the same order 122 as in the original synchronous program. 123

The asynchronization on the right of Fig. 1 is not the only sound (data-race free) asynchronization 124 of the program on the left. For instance, the await at line 24 can be moved one statement up (before 125 the read of x) and the resulting program remains equivalent to the sequential one. Thus, we 126 consider the problem of enumerating *all* sound asynchronizations of a sequential program P w.r.t. 127 substituting a set of methods with asynchronous versions. Enumerating all sound asynchronizations 128 makes it possible to deal separately with the problem of choosing the best asynchronization in 129 terms of performance based on some metric (e.g., performance tests). This problem reduces to 130 finding all possible placements of awaits that do not introduce data races. 131

In general, the number of (sound) asynchronizations is exponential in the number of method 132 calls in the program. Therefore, we focus on the *delay* complexity of this enumeration problem, 133 i.e., the complexity of the delay between outputting two consecutive (distinct) outputs. Note that a 134 trivial enumeration of all asynchronizations and checking equivalence to the original program for 135 each one of them has an exponential delay complexity. We also consider the problem of computing 136 optimal asynchronizations that maximize the distance between a call and a matching await. The 137 code in between these two statements can execute in parallel with the awaited task, and therefore, 138 optimal asynchronizations maximize the amount of parallelism. Note however that it is hard to 139 argue that such maximal parallelism translates always to maximal performance in practice. 140

We show that both the delay complexity of the enumeration problem, and the complexity of computing an optimal asynchronization are polynomial time modulo an oracle for solving reachability (assertion checking) in *sequential* programs (they both reduce to a quadratic number

147

 <sup>&</sup>lt;sup>145</sup> <sup>1</sup>Actually, the .Net platform does not contain such a method. We use it here to simplify the exposition. Reading the content
 <sup>146</sup> of a webpage should pass through WebRequest and HttpWebResponse objects. The explanations would remain valid.

of reachability queries). The former relies on the latter via a rather surprising result, which differs 148 from other concurrency synthesis problems (e.g., insertion of locks), which is that the optimal 149 asynchronization is *unique*. This holds even if the optimality is relative to a given asynchronization 150  $P_a$  which intuitively, imposes an upper bound on the distance between awaits and matching calls. 151 In general, one could expect that avoiding data races could reduce to a choice between moving 152 one await or another closer to the matching call. We show that this is not necessary because 153 essentially, the optimal asynchronization is required to be equivalent to a *sequential* program, 154 155 which is deterministic and executes statements in a fixed order. To show the robustness of these results, we also investigate the related problem of synthesizing sound multi-threaded refactorings, 156 where every method call is executed by a different thread. We show that the techniques used to 157 compute asynchronizations can be extended to this case as well. 158

As a more pragmatic approach, we define a procedure for computing sound asynchronizations 159 which relies on a bottom-up interprocedural data-flow analysis. Intuitively, the placement of awaits 160 is computed by traversing the call graph bottom up, from "base" methods that do not call any other 161 method in the program, to methods that call only base methods, and so on. Each method m is 162 considered only once, and the placement of awaits in *m* is derived based on a data-flow analysis 163 that computes read or write accesses made in the callees. We show that this procedure computes 164 optimal asynchronizations of abstracted programs where every Boolean condition in if-then-else 165 constructs or while loops is replaced with non-deterministic choice. These asynchronizations are 166 sound for the concrete programs as well. This procedure enables a polynomial delay enumeration 167 of the sound asynchronizations of abstracted programs. 168

We implemented the asynchronization enumeration based on data-flow analysis in a prototype tool for C# programs. We evaluated this implementation on a number of non-trivial programs extracted from open source repositories. This evaluation shows that sound asynchronizations can be enumerated efficiently and in some cases, we found asynchronizations that increase the amount of parallelism (the distance between calls and awaits). This demonstrates that our techniques have the potential to become the basis of refactoring tools that allow programmers to improve their usage of async/await primitives.

176 In summary, this paper makes the following contributions:

- We define the problem of data race-free asynchronization synthesis for refactoring sequential code to equivalent asynchronous code.
- We show that the optimization problem of computing a data race-free asynchronization that maximizes the distance between calls and awaits admits a unique solution.
  - We investigate the delay complexity of data race-free asynchronization synthesis.
  - A pragmatic algorithm based on data-flow analysis for computing asynchronizations.
  - A prototype implementation of this algorithm and an evaluation of this prototype on a benchmark of non-trivial C# programs extracted from open-source repositories.

### 2 OVERVIEW

We give an overview of our techniques for synthesizing sound asynchronizations using as example 187 the synchronous program on the left of Fig. 2. We discuss asynchronizations obtained by replacing 188 the calls to IO with equivalent asynchronous counterparts IOAsync. The program on the center of 189 Fig. 2 is the "weakest" asynchronization where the awaits cannot be moved further away from their 190 matching calls because of the use of the return values. The methods MainAsync, F1, and F2 are 191 declared to be asynchronous since they (in)directly call IOAsync. This program is not a solution to 192 our synthesis problem since it is not equivalent to the sequential program. It admits new behaviors, 193 an example being pictured in Fig. 3 (edges represent execution order): the accesses to x and y in 194 MainAsync occur before the write to x in F1 and the write to y in F2, respectively. These statements 195

177

178

179

180

181

182

183

184 185

<sup>196</sup> 

213

214

197	1 void Main() {	1 async Task MainAsync() {	1 async Task MainAsync() {
198	2 F1();	2 Task t1 = F1();	2 Task t1 = F1();
100	3 F2();	3 Task t2 = F2();	3 Task t2 = F2();
199			4 await t1;
200	5 int r1 = x;	5 int r1 = x;	5 int r1 = x;
201			6 await t2;
000	7 int r2 = y;	7 int r2 = y;	7 int r2 = y;
202	8 x = r1 + r2;	8 x = r1 + r2;	8 x = r1 + r2;
203		9 await t1;	
204		10 await t2;	
205	11 }	11 }	11 }
205	12 void F1() {	12 async Task F1() {	12 async Task F1() {
206	<pre>13 int r3 = IO();</pre>	<pre>13 Task<int> t3 = IOAsync();</int></pre>	<pre>13 Task<int> t3 = IOAsync();</int></pre>
207		14 int r3 = await t3;	14 int r3 = await t3;
208	15 x = r3;	15 x = r3;	15 x = r3;
200	16 }	16 }	16 }
209	17 void F2() {	17 async Task F2() {	17 async Task F2() {
210	18 int r4 = IO();	<pre>18 Task<int> t4 = IOAsync();</int></pre>	<pre>18 Task<int> t4 = IOAsync();</int></pre>
011		19	<pre>19 int r4 = await t4;</pre>
211	20 y = r4;	20 y = r4;	20 y = r4;
212	21 }	21 }	21 }

Fig. 2. A synchronous C# program and two asynchronizations (x and y are static variables).

215 were executed in the opposite order in the sequential program. They form three data races because they are not ordered by the control-flow of the asynchronization. There is another execution where 216 they execute as in the original program: if tasks t3 and t4 finish immediately, then the await has 217 no effect, and F1 and F2 finish before returning control to their caller MainAsync. 218



Fig. 3. An execution with three data races on x and y. The blocks of instructions executed in F1 and F2 are marked with red and blue outlines, respectively. Each call is decomposed into two blocks representing what is executed before they are suspended (due to an await of IOAsync) and when the continuation is scheduled.

We define a procedure for computing a data racefree asynchronization, which is optimal in the sense that it maximizes the distance between calls and matching awaits. This is an iterative process that repairs data races starting from the "weakest" asynchronization on the center of Fig. 2. For instance, the data race between the write to x in MainAsync and the write to x in F1 in Fig. 3 can be repaired by moving the await t1 one position up, before the the write to x in MainAsync. This way, the write to x in F1 will always execute first. The call to F1 that matches this await and the write to x in MainAsync are regarded as the root cause of this data race.

For efficiency, the data races are enumerated and repaired in a certain order, that avoids superfluous repair steps. This order prioritizes data races involving statements that would execute first in the original sequential program. For instance, in Fig. 3, the first data race to repair involves the read of x from MainAsync and the write to x in F1, because these statements are the first to execute in the original sequential program among the other statements involved in data races. Repairing this data race consists

in moving await t1 before the read of x from MainAsync, which implies that F1 completes before 242 the read of x. This repair is defined from a notion of *root cause* of a data race, that in this case, 243 contains the call to F1 and the read of x from MainAsync. Interestingly, this repair step removes 244

245

235

236

237

238

239

240

the write-write data race between the write to x in MainAsync and the write to x in F1 as well. If 246 we would have repaired these data races in the opposite order, we would have moved await t1 247 first before the write to x, and then, before the read of x. Similarly, the data race between the read 248 of y from MainAsync and the write to y in F2 is repaired by moving the await t2 before the the 249 read from y in MainAsync (the call to F2 and the read from y in MainAsync are the root cause of this 250 data race). The resulting program is shown on the right of Fig. 2 and it is equivalent to the original 251 sequential program. We show that the problem of computing root-causes of data races which are 252 253 minimal in this order can be reduced in polynomial time to reachability (assertion checking) in sequential programs. 254

Asynchronizations can be partially ordered depending on the distance, i.e., the number of 255 statements from the original program, between a call and a matching await. The left of Fig. 4 256 pictures an excerpt of this partial order where an asynchronization is represented as a vector of 257 258 distances, the first element is the number of statements between the call and the await on t1 (we count only statements that appear in the sequential program as well and exclude awaits), and so 259 on. The asynchronization on the right of Fig. 1 corresponds to (1, 1, 0, 0). The bottom of this order 260 represents an asynchronization where every call is immediately followed by await, and it has the 261 same semantics as the original program. The top element is the "weakest" asynchronization, which 262 was given in the middle of Fig. 1. The right of Fig. 4 gives the set of all sound asynchronizations. 263 The program on the right of Fig. 1 is the biggest element, i.e., moving any await further away from 264 the matching call introduces a data race. 265

To enumerate all sound asynchronizations, we 266 perform a top-down traversal of the partial order on 267 the left of Fig. 4. We first compute the biggest ele-268 ment that is data race free. Although this is a partial 269 order, we show that this element is actually unique. 270 Intuitively, uniqueness is proved by contradiction, 271 showing that the least upper bound of two maxi-272 mal incomparable sound asynchronizations is also a 273 274 sound asynchronization. For instance, the least common ancestor of the two sound asynchronizations 275 with vectors of distances (1, 0, 0, 0) and (0, 1, 0, 0) is 276 the sound asynchronization (1, 1, 0, 0). 277

$$\begin{array}{c}(4,3,0,0)\\\vdots\\(2,1,0,0)&(1,2,0,0)\\(1,1,0,0)&(0,1,0,0)\\(0,1,0,0)&(1,0,0,0)\\(0,0,0,0)&(0,0,0,0)\end{array}$$

Fig. 4. Partially-ordered sets of asynchronizations of the program on the left of Fig. 2. The edges connect comparable elements, smaller elements being below bigger elements.

Then, for each immediate predecessor  $P_a$  of the biggest sound asynchronization, we compute the biggest sound asynchronization which is smaller than  $P_a$  (the first step is a particular case where  $P_a$ is the top element). As an extension of the previous case, this is also unique, and called an optimal asynchronization relative to  $P_a$ . Ensuring that traversals starting in different immediate predecessors explore disjoint parts of the asynchronization space requires some additional constraints on the exploration, which are explained in Section 5. The enumeration finishes when reaching the bottom, which is data race free by definition, on all branches of the recursion.

As a pragmatic alternative, we propose a procedure based on static analysis for enumerating sound asynchronizations, which follows essentially the same schema, except that the problem of data race detection is delegated to a static analysis.

### 3 ASYNCHRONOUS PROGRAMS

Fig. 5 lists the syntax of a simple programming language used to formalize our approach. A *program* is defined by set of methods, including a distinguished main, which are classified as *synchronous* or
 *asynchronous*. Synchronous methods execute immediately as they are invoked and run continuously
 until completion. Asynchronous methods, marked using the keyword async, can run only partially

294

285

286

287 288

299

300

301

 295
  $\langle prog \rangle$  ::= program  $\langle md \rangle$  

 296
  $\langle md \rangle$  ::= method  $\langle m \rangle \langle inst \rangle | async method <math>\langle m \rangle \langle inst \rangle | \langle md \rangle$ ;  $\langle md \rangle$  

 297
  $\langle inst \rangle$  ::=  $\langle x \rangle := \langle le \rangle | \langle r \rangle := \langle x \rangle | \langle r \rangle := call \langle m \rangle | return | await \langle r \rangle | await * | if <math>\langle le \rangle \{\langle inst \rangle\}$  else { $\langle inst \rangle$ }

 298
 | while  $\langle le \rangle \{\langle inst \rangle\} | \langle inst \rangle; \langle inst \rangle$ 

Fig. 5. Syntax.  $\langle m \rangle$ ,  $\langle x \rangle$ , and  $\langle r \rangle$  represent method names, program and local variables, resp.  $\langle le \rangle$  is an expression over local variables, or \* which is non-deterministic choice.

and be interrupted when executing an await. Only asynchronous methods are allowed to use
 await, and all methods using await must be defined as asynchronous. We assume that methods
 are not (mutually) recursive. A program is called *synchronous* if it is a set of synchronous methods.

A method consists of a method name from a set  $\mathbb{M}$  and a method body, i.e., a list of statements. 305 These statements use a set  $\mathbb{PV}$  of *program variables*, which can be accessed from different methods 306 (ranged over using x, y, z,...), and a set  $\mathbb{LV}$  of method *local variables* (ranged over using r,  $r_1$ , 307  $r_2,\ldots$ ). We assume that input/return parameters are modeled using dedicated program variables. 308 We assume that each method call returns a *unique task identifier* from a set  $\mathbb{T}$ , which is used to 309 record control dependencies imposed by awaits (for uniformity, synchronous methods return a 310 task identifier as well). Our language includes assignments to local/program variables, awaits, 311 return statements, while loops, and conditionals. We assume that variables take values from a 312 data domain  $\mathbb{D}$ , which includes  $\mathbb{T}$  to account for variables storing task identifiers. The assignment to 313 a local variable  $\langle r \rangle := \langle x \rangle$ , where x is a program variable, is called a *read* of  $\langle x \rangle$  and an assignment 314 to a program variable  $\langle x \rangle := \langle le \rangle$  is called a *write* to  $\langle x \rangle$ . A *base* method is a method whose body 315 does not contain method calls. 316

Asynchronous methods. Asynchronous methods can use awaits to wait for the completion of a task (invocation) while *the control is passed to their caller*. The parameter *r* of the await specifies the id of the awaited task. As a sound abstraction of awaiting the completion of an IO operation (reading or writing a file, an http request, etc.), which we do not model explicitly, we use a variation await \*. This has a non-deterministic effect of either continuing to the next statement in the same method (as if the IO operation already completed), or passing the control to the caller (as if the IO operation is still pending).

For example, Fig. 6 lists the modeling in our language of IO 324 methods ReadToEndAsync and GetStringAsync used in Fig. 1. 325 We use program variables to represent system resources such as 326 the network or the file system. The await for the completion of 327 accesses to such resources is modeled by await \*. This enables 328 capturing racing accesses to system resources in asynchronous 329 executions. GetStringAsync contains a read of the resource WWW 330 (for world wide web) at some input url. Parameters or return 331 values are modeled using program variables. ReadToEndAsync 332 is modeled using reads/writes of the index/content of the input 333 stream, and await \* models the await for their completion. 334

```
async method GetStringAsync() {
   await *;
   retVal = WWW[url_Input];
   return
}
async method ReadToEndAsync() {
   await *;
   ind = Stream.index;
   len = Stream.content.Length;
   if (ind >= len)
      retVal = ""; return
   Stream.index = len;
   retVal = Stream.content(ind,len);
   retUarn }
```

We assume that the body of every asynchronous method m Fig. 6. Modeling IO operations.

satisfies several well-formedness syntactic constraints, defined on its control-flow graph (CFG). We
 recall that each node of the CFG represents a basic block of code (a maximal-length sequence of
 branch-free code), and nodes are connected by directed edges which represent a possible transfer
 of control between blocks. Thus,

- (1) every call r := call m' uses a distinct variable r (to store task identifiers),
  - (2) every CFG block containing an await r is dominated by the CFG block containing the call
    - $r := call \dots$  (i.e., every CFG path from the entry to the await has to pass through the call),
- 342 343

341

(3) every CFG path starting from a block containing a call r := call ... to the exit has to pass through an await r statement.

The first condition simplifies the technical exposition, while the last two ensure that r stores a valid task identifier when executing an await r, and that every asynchronous invocation is awaited before the caller finishes. Languages like C# or Javascript do not enforce the latter constraint, but it is considered bad practice due to possible exceptions that may arise in the invoked task and which are not caught. In this work, we forbid passing task identifiers as method parameters or return values (which is possible in C#). An await r statement is said to *match* an r := call m' statement.

For example, the program on the left of Fig. 7 does not satisfy the second condition above since await r can be reached without entering the loop. The program in the center of Fig. 7 does not satisfy the third condition since we can reach the end of the method without entering the if branch and



thus, without executing await r. The program on the right of Fig. 7 satisfies both conditions.

Semantics. A program configuration is a tuple (g, stack, pending, completed, c-by, w-for) where 360 g is composed of the valuation of the program variables, stack is the call stack, pending is the 361 set of asynchronous tasks, e.g., continuations predicated on the completion of some method call, 362 completed is the set of completed tasks, c-by represents the relation between a method call and its 363 caller, and w-for represents the control dependencies imposed by await statements. The activation 364 frames in the call stack and the asynchronous tasks are represented using triples  $(i, m, \ell)$  where 365  $i \in \mathbb{T}$  is a task identifier,  $m \in \mathbb{M}$  is a method name, and  $\ell$  is a valuation of local variables, including 366 as usual a dedicated program counter. The set of completed tasks is represented as a function 367 completed :  $\mathbb{T} \to \{\top, \bot\}$  such that completed $(i) = \top$  when *i* is completed and completed $(i) = \bot$ , 368 otherwise. We define c-by and w-for as partial functions  $\mathbb{T} \to \mathbb{T}$  with the meaning that c-by(i) = j, 369 resp., w-for(i) = j, iff i is called by j, resp., i is waiting for j. We set w-for(i) = \* if the task i was 370 interrupted because of an await \* statement. 371

The semantics of a program *P* is defined as a labeled transition system (LTS)  $[P] = (\mathbb{C}, \operatorname{Act}, \operatorname{ps}_0, \rightarrow$ ) where  $\mathbb{C}$  is the set of program configurations, Act is a set of transition labels called *actions*,  $\operatorname{ps}_0$  is the initial configuration, and  $\rightarrow \subseteq \mathbb{C} \times \operatorname{Act} \times \mathbb{C}$  is the transition relation. Each program statement is interpreted as a transition in [P]. The set of actions is defined by:

Act ={ $(i, ev) : i \in \mathbb{T}, ev \in \{ rd(x), wr(x), call(j), await(k), return, cont : j \in \mathbb{T}, k \in \mathbb{T} \cup \{ * \}, x \in \mathbb{PV} \}$ 

The transition relation  $\rightarrow$  is defined in Fig. 8. Transition labels are written on top of  $\rightarrow$ .

Transitions labeled by (i, rd(x)) and (i, wr(x)) represent a read and a write accesses to the program variable x, respectively, executed by the task (method call) with identifier *i*. A transition labeled by (i, call(j)) corresponds to the fact that task *i* executes a method call that results in creating a task *j*. Task *j* is added on the top of the stack of currently executing tasks, declared pending (setting completed(*j*) to  $\perp$ ), and c-by is updated to track its caller (c-by(*j*) = *i*). A transition (i, return) represents the return from task *i*. Task *i* is removed from the stack of currently executing tasks, and completed(*i*) is set to  $\top$  to record the fact that task *i* is finished.

A transition (i, await(j)) corresponds to task *i* waiting asynchronously for task *j*. Its effect depends on whether task *j* is already completed. If this is the case (i.e., completed $[j] = \top$ ), task *i* continues and executes the next statement. Otherwise, task *i* executing the await is removed from the stack and added to the set of pending tasks, and w-for is updated to track the waitingfor relationship (w-for(*i*) = *j*). Similarly, a transition (*i*, await(\*)) corresponds to task *i* waiting asynchronously for the completion of an unspecified task. Non-deterministically, task *i* continues

378

344

$$\frac{m \in \mathbb{M} \quad r := x \in inst(\ell(pc)) \quad \ell' = \ell[r \mapsto g(x), pc \mapsto next(\ell(pc))]}{(g, (i, m, \ell) \circ stack, *, *, *)} \frac{(i, rd(x))}{(g, (i, m, \ell') \circ stack, *, *, *)} \frac{(i, rd(x))}{(g, (i, m, \ell') \circ stack, *, *, *)} \frac{(i, rd(x))}{(g, (i, m, \ell') \circ stack, *, *, *)} \frac{(i, rd(x))}{(g, (i, m, \ell') \circ stack, *, *, *)} \frac{(i, rd(x))}{(g, (i, m, \ell') \circ stack, *, *, *)}$$

$$\frac{m \in \mathbb{M} \quad x := le \in inst(\ell(pc)) \quad \ell' = \ell[pc \mapsto next(\ell(pc))] \quad g' = g[x \mapsto \ell(le)]}{(g, (i, m, \ell) \circ stack, *, *, *)} \frac{(i, rd(x))}{(g, (i, m, \ell') \circ stack, *, *, *)}$$

$$r := call m \in inst(\ell(pc)) \quad \ell_0 = int(g, m) \quad j \in T fresh \quad \ell' = \ell[r \mapsto j, pc \mapsto next(\ell(pc))]}{(g, (j, m, \ell) \circ stack, *, completed, c - by, *) \quad \frac{(i, call(j))}{(g, (i, m, \ell) \circ (j, m', \ell') \circ stack, *, completed', r - by'; j \mapsto i]}$$

$$\frac{m \in \mathbb{M} \wedge return \in inst(\ell(pc)) \quad completed' = completed[j \mapsto T]}{(g, (i, m, \ell) \circ stack, *, completed, c - by, *) \quad \frac{(i, call(\ell))}{(g, (i, m, \ell) \circ stack, *, completed, *, *)} \frac{(g, (i, m, \ell) \circ stack, *, completed, *, *)}{(g, (i, m, \ell) \circ stack, *, completed, *, *) \quad \frac{(i, return)}{(g, (i, m, \ell) \circ stack, *, completed, *, *)} \frac{m \in \mathbb{M} \quad await r \in inst(\ell(pc)) \quad completed(\ell(r)) = T \quad \ell' = \ell[pc \mapsto next(\ell(pc))]}{(g, (i, m, \ell) \circ stack, *, completed, *, *) \quad \frac{(i, await(\ell(r)))}{(g, (i, m, \ell') \circ stack, *, completed, *, *)} \frac{m \in \mathbb{M} \quad await r \in inst(\ell(pc)) \quad completed(\ell(r)) = T \quad \ell' = \ell[pc \mapsto next(\ell(pc))]}{(g, (i, m, \ell) \circ stack, pending, completed, *, w - for) \quad (i, await(\ell(r)))} (g, (i, m, \ell') \circ stack, *, completed, *, w - for)} \frac{m \in \mathbb{M} \quad await r \in inst(\ell(pc)) \quad u^{-1} = \ell[pc \mapsto next(\ell(pc))]}{(g, (i, m, \ell) \circ stack, pending, completed, *, w - for) \quad (i, await(e))} (g, (i, m, \ell') \circ stack, *, *, *)} \frac{m \in \mathbb{M} \quad await * \in inst(\ell(pc)) \quad w - for' = w - for[i \mapsto *] \quad \ell' = \ell[pc \mapsto next(\ell(pc))]}{(g, (i, m, \ell) \circ stack, pending, *, *, w - for) \quad (i, await(e))} (g, (i, m, \ell) \circ stack, *, *, *)} \frac{m \in \mathbb{M} \quad w - for(i) = j \quad completed(j) = T \quad m \in \mathbb{M} \quad w - for(i) = j \quad m \in \mathbb{M} \quad w - for(i) = s \quad (i, cont)} (g, (i, m, \ell) \circ stack, pending, *, *, w - for) \quad (i, cont)} (g, (i, m, \ell) \circ stack, pending,$$

Fig. 8. Program semantics. For a function f, we use  $f[a \mapsto b]$  to denote a function g such that g(c) = f(c)for all  $c \neq a$  and g(a) = b. The function inst returns the instruction at some given control location while next gives the next instruction to execute. We use  $\circ$  to denote sequence concatenation. We use init to represent the initial state of a method call.

to the next statement, or task *i* is interrupted and transferred to the set of pending tasks (w-for(i) is set to \*).

A transition (i, cont) represents the scheduling of the continuation of task *i*. There are two 426 cases depending on whether *i* waited for the completion of another task *j* modeled explicitly in 427 the language (i.e., w-for(i) = j), or an unspecified task (i.e., w-for(i) = \*). In the first case, the 428 transition is enabled only when the call stack is empty and the task *j* is completed. In the second 429 case, the transition is enabled without any additional requirements. The latter models the fact that 430 methods implementing IO operations (waiting for unspecified tasks in our language) are executed 431 in background threads and can interleave with the main thread (that executes the Main method). 432 Although this may seem restricted because we do not allow arbitrary interleavings between such 433 methods, it is however complete when focusing on the existence of data races as in our approach. 434

An execution of *P* is a sequence  $\rho = ps_0 \xrightarrow{a_1} ps_1 \xrightarrow{a_2} \dots$  of transitions starting in the initial configuration  $ps_0$  and leading to a configuration ps where the call stack and the set of pending tasks are empty.  $\mathbb{C}[P]$  denotes the set of all program variable valuations included in configurations that are reached in executions of *P*. Since we are only interested in reasoning about the sequence of actions  $a_1 \cdot a_2 \cdot \dots$  labeling the transitions of an execution, we will call the latter an execution as well. The set of executions of a program *P* is denoted by  $\mathbb{E}x(P)$ .

423

424

443	$a_1 \leq_{\rho} a_2$	$a_1$ occurs before $a_2$ in $\rho$
444	$a_1 \sim a_2$	$a_1 = (i, ev) \text{ and } a_2 = (i, ev')$
445	$(a_1, a_2) \in MO$	$a_1 \sim a_2 \wedge a_1 \leq_{\rho} a_2$
446	$(a_1, a_2) \in \mathrm{CO}$	$(a_1, a_2) \in MO \lor (a_1 = (i, call(j)) \land a_2 = (j, \_)) \lor (\exists a_3. (a_1, a_3) \in CO \land (a_3, a_2) \in CO)$
447	$(a_1, a_2) \in SO$	$(a_1, a_2) \in \mathrm{CO} \lor (\exists a_3. (a_1, a_3) \in \mathrm{SO} \land (a_3, a_2) \in \mathrm{SO})$
448		$\lor (a_1 = (j, \_) \land a_2 = (i, \_) \land \exists a_3 = (i, call(j)). a_3 \leq_{\rho} a_2)$
449	$(a_1, a_2) \in HB$	$(a_1, a_2) \in CO \lor (\exists a_3. (a_1, a_3) \in HB \land (a_3, a_2) \in HB)$
450		$\vee (a_1 = (j, \_) \land a_2 = (i, \_) \land \exists a_3 = (i, \operatorname{await}(j)). a_3 \neq a_2 \land a_3 \leq_{\rho} a_2)$
451		$\lor$ $(a_1 = (j, \text{await}(i'))$ is the first await in $j \land a_2 = (i, \_) \land \exists a_3 = (i, \text{call}(j))$ . $a_3 \leq_{\rho} a_2)$

Table 1. Strict partial orders included in a trace.

452 **Traces.** The *trace* of an execution  $\rho \in \mathbb{E}x(P)$  is a tuple  $tr(\rho) = (\rho, MO, CO, SO, HB)$  of strict partial 453 orders between the actions in  $\rho$  defined in Table 1. The *method invocation order* MO records the 454 order between actions in the same invocation, and the call order CO is an extension of MO that 455 additionally orders actions before an invocation with respect to those inside that invocation. The synchronous happens-before order SO orders the actions in an execution as if all the invocations 456 457 were synchronous (even if the execution may contain asynchronous ones). It is an extension of CO where additionally, every action inside a callee is ordered before the actions following its invocation 458 459 in the caller. The (asynchronous) happens-before order HB contains typical control-flow constraints: it is an extension of CO where every action a inside an asynchronous invocation is ordered before 460 the corresponding await in the caller, and before the actions following its invocation in the caller 461 if *a* precedes an await in MO (an invocation can be interrupted only when executing an await). 462 463  $\mathbb{T}r(P)$  is the set of traces of a program *P*.



tions in the trace.

Fig. 9 shows a trace where two statements are linked by a dotted arrow if the corresponding actions are related by MO, a dashed arrow if the corresponding actions are related by the CO but not by MO, and a solid arrow if the corresponding actions are related by the HB but not by CO.

### SYNTHESIZING ASYNCHRONOUS PROGRAMS 4

Fig. 9. A trace of an asynchronous pro-We define the synthesis problem we investigate in this gram. Arrows between statements denote work. Given a synchronous program P and a subset of base relations between the corresponding acmethods  $L \subseteq P$ , the goal is to synthesize *all* asynchronous programs  $P_a$  that are equivalent to P and that are obtained

by substituting every method in L with an equivalent asynchronous version. The base methods are intended to be models of standard library calls (e.g., IO operations) in a practical context, and asynchronous versions are defined by inserting await \* statements (in the original synchronous code). We use P[L] to emphasize a subset of base methods L of a program P. Also, we will call L a library. A library is called (a)synchronous when all methods are (a)synchronous.

### 4.1 Asynchronizations of a Synchronous Program

Let P[L] be a synchronous program, and 482  $L_a$  a set of asynchronous methods obtained 483 from those in L by inserting at least one 484 await \* statement in their body (and adding 485 the keyword async). We assume that each 486 method in  $L_a$  corresponds to a method in L 487 with the same name, and vice-versa. A pro-488 gram  $P_a[L_a]$  is called an *asynchronization* of 489 490

<pre>method Main {</pre>	<pre>async method Main {</pre>	<pre>async method Main {</pre>
r1 = call m;	r1 = call m; await r1;	r1 = call m;
r2 = x;	r2 = x;	r2 = x; await r1;
}	}	}
<pre>method m() {</pre>	async method m {     await *	async method m {     await *
retVal = x;	retVal = x;	retVal = x;
<pre>x = input;</pre>	<pre>x = input;</pre>	<pre>x = input;</pre>
return;	return;	return;
}	3	3

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

Fig. 10. A program and its asynchronizations.

471

472

473

474

475

476

477

478

479

480

491 P[L] with respect to  $L_a$  if it is a syntactically

<sup>492</sup> correct program obtained by replacing the

<sup>493</sup> methods in *L* with those in  $L_a$  and adding

await statements as necessary. More precisely, let  $L^* \subseteq P$  be the set of all methods of P that transitively call methods of L. Formally,  $L^*$  is the smallest set of methods that includes L and that satisfies the following: if a method m calls a method  $m \in L^*$ , then  $m \in L^*$ . Then,  $P_a[L_a]$  is an *asynchronization* of P[L] with respect to  $L_a$  if it is obtained from P as follows:

- All methods in L\* \ L are declared as asynchronous (we assume that every call to an asynchronous method is followed by an await statement, and any method that uses await must be declared as asynchronous).
  - For each invocation r := call m of a method  $m \in L^*$ , add await statements await r satisfying the well-formedness syntactic constraints described in Section 3.

For instance, Fig. 10 lists a synchronous program and its two asynchronizations, where  $L = L^* = \{m\}$ . Asynchronizations differ only in the await placement.

Async  $[P, L, L_a]$  is the set of all asynchronizations of P[L] w.r.t.  $L_a$ . The *strong* asynchronization, denoted by strongAsync  $[P, L, L_a]$ , is an asynchronization where every added await *immediately* follows the matching method call. The strong asynchronization reaches exactly the same set of program variable valuations as the original program.

### 4.2 **Problem Definition**

498

499

500

501

502

509

We investigate the problem of enumerating *all* asynchronizations of a given program w.r.t. a given asynchronous library, which are *sound*, in the sense that they do not admit data races. Two actions  $a_1$  and  $a_2$  in a trace  $\tau = (\rho, MO, CO, SO, HB)$  are *concurrent* if  $(a_1, a_2) \notin HB$  and  $(a_2, a_1) \notin HB$ .

Definition 4.1 (Data Race). An ansynchronous program  $P_a$  admits a data race  $(a_1, a_2)$ , where  $(a_1, a_2) \in SO$ , if  $a_1$  and  $a_2$  are two concurrent actions of a trace  $\tau \in Tr(P_a)$ , and  $a_1$  and  $a_2$  are read or write accesses to the same program variable x, and at least one of them is a write.

For example, the program on the right of Fig. 10 admits a data race between the actions that correspond to x = input and  $r^2 = x$ , respectively, in a trace where the call to *m* is suspended when it reaches await \* and the control is transferred to Main which executes  $r^2 = x$ . Traces of *synchronous* programs can *not* contain concurrent actions, and therefore they do not admit data races. Note that also strongAsync[*P*, *L*, *L*<sub>a</sub>] does not admit data races.

<sup>522</sup> An asynchronization  $P_a[L_a]$  is called *sound* when  $P_a[L_a]$  does not admit data races. Absence of <sup>523</sup> data races implies equivalence to the original program, in the sense of reaching the same set of <sup>524</sup> configurations (program variable valuations).

LEMMA 4.2.  $P[L] \equiv P_a[L_a]$  implies  $\mathbb{C}[P[L]] = \mathbb{C}[P_a[L_a]]$ , for every  $P_a[L_a] \in \text{Async}[P, L, L_a]$ 

PROOF. We have to show that for any asynchronization  $P_a$  of a program P, if  $P_a$  does not admit data races then  $\mathbb{C}[P_a] = \mathbb{C}[P]$ . Let  $\rho$  be an execution of  $P_a$  that reaches a configuration  $ps \in \mathbb{C}[P_a]$ . We show that actions in  $\rho$  can be reordered such that any action that occurs in  $\rho$  between (i, call(j))and (j, return) is not of the form  $(i, \_)$  (i.e., the task j is executed synchronously). If an action  $(i, \_)$ occurs in  $\rho$  between (i, call(j)) and (j, return), then it must be concurrent with (j, return). Since  $P_a$  does not admit data races, an execution  $\rho'$  resulting from  $\rho$  by reordering any two concurrent actions reaches the same configuration ps as  $\rho$ . Therefore, there exists an execution  $\rho''$  where the actions that occur between any (i, call(j)) and (j, return) are not of the form  $(i, \_)$ . This is also an execution of P (modulo removing the awaits which have no effect), which implies  $ps \in \mathbb{C}[P]$ .  $\Box$ 

Definition 4.3. Given a synchronous program P[L], and an asynchronous library  $L_a$ , the asychronization synthesis problem asks to enumerate all sound asynchronizations in Async[ $P, L, L_a$ ].

537 538 539

525

526

527

528

529

530

531

532

533

534

There are other approaches for checking behavioral equivalence, e.g., reachable states equivalence, 540 and input-output equivalence, however, they are either undecidable or decidable but highly complex. 541 Our approach offers a stronger notion of equivalence that is simpler to check. 542

### ENUMERATING SOUND ASYNCHRONIZATIONS 544

We describe an algorithm for solving the asynchronization synthesis problem. This algorithm relies on a partial order between asynchronizations that guides the enumeration of possible solutions. It computes optimal solutions (according to this order) repeatedly under different bounds, until exploring the whole space.

### **Optimal Asynchronization** 5.1

We define a partial order on the space of asynchronizations which takes into account the distance between call statements and corresponding await statements. 552

An await statement  $s_w$  in a method m of an asynchronization  $P_a[L_a] \in Async[P, L, L_a]$  covers a read/write statement s in P if there exists a path in the CFG of m from the call statement matching  $s_w$  to  $s_w$  that contains s. The set of statements covered by an await  $s_w$  is denoted by Cover $(s_w)$ .

We compare asynchronizations in terms of sets of statements covered by awaits that match the same call from the original synchronous program P[L]. Since asynchronizations are obtained by adding awaits, every call statement in an asynchronization  $P_a[L_a] \in Async[P, L, L_a]$  corresponds to a *fixed* call in P[L].

Definition 5.1. For two asynchronizations  $P_a, P'_a \in Async[P, L, L_a], P_a$  is less asynchronous than  $P'_a$ , denoted by  $P_a \leq P'_a$ , iff for every await statement  $s_w$  in  $P_a$ , there exists an await statement  $s_w$ in  $P'_a$  that matches the same call as  $s_w$ , such that  $Cover(s_w) \subseteq Cover(s'_w)$ .

For example, the two asynchronous programs in Fig. 10 are ordered by  $\leq$  since Cover(await r1) = {} in the first and Cover(await r1) = {r2 = x} in the second.

Note that the strong asynchronization is less asynchronous than any other asynchronization. Also, note that  $\leq$  has a unique maximal element that is called the weakest asynchronization and denoted by weakAsync[ $P, L, L_a$ ]. For instance, the program on the right of Fig. 10 is the weakest asynchronization of the synchronous program on the left of the figure.

**Relative Optimality.** A crucial property of this partial order is that for every asynchronization  $P_a$ , there exists a *unique* maximal asynchronization that is smaller than  $P_a$  (w.r.t.  $\leq$ ) and that is sound. Formally, given  $P_a \in \text{Async}[P, L, L_a]$ , an asynchronization  $P'_a$  is called an *optimal asynchronization of P* relative to  $P_a$  if (1)  $P'_a \leq P_a$ ,  $P'_a$  is sound, and (2)  $P'_a$  is maximal among other sound asynchronizations smaller than  $P_a$ , i.e.,  $\forall P''_a \in \text{Async}[P, L, L_a]$ .  $P''_a$  is sound and  $P''_a \leq P_1 \Rightarrow P''_a \leq P'_a$ .

The following lemma shows that for a given  $P_a$  there exists a unique  $P'_a$  that is an optimal asynchronization of P relative to  $P_a$ . The existence is implied by the fact that strongAsync[P, L,  $L_a$ ] is the bottom element of  $\leq$ . To prove uniqueness, we assume by contradiction that there exist two incomparable optimal asynchronizations  $P_a^1$  and  $P_a^2$  and select the first await statement  $s_w^1$ , according to the control-flow of the sequential program, that is placed in different positions in the two programs. Assume that  $s_w^1$  is closer to its matching call in  $P_a^1$ . Then, we move  $s_w^1$  in  $P_a^1$  further away from its matching call to the same position as in  $P_a^2$ . This modification does not introduce data races since  $P_a^2$  is data race free. Thus, the resulting program is data race free, bigger than  $P_a^1$ , and smaller than  $P_a$  w.r.t.  $\leq$  contradicting the fact that  $P_a^1$  is an optimal asynchronization.

LEMMA 5.2. Given an asynchronization  $P_a \in \text{Async}[P, L, L_a]$ , there exists a unique program  $P'_a$  that is an optimal asynchronization of P relative to  $P_a$ .

**PROOF.** Since strongAsync[ $P, L, L_a$ ] is the bottom element of  $\leq$ , then there always exists a sound asynchronization smaller than  $P_a$ . Assume by contradiction that there exist two distinct programs

1:12

543

545

546

547

548 549

550

551

553

554

555

556

557

558

559 560

561 562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

Algorithm 1 An algorithm for enumerating all sound asynchronizations (these asynchronizations 589 are obtained as a result of the **output** instruction). OptRelative returns the optimal asynchroniza-590 tion of *P* relative to  $P_a$ 591

592	1: p	<b>rocedure</b> AsyncSynthesis( $P_a$ , $s_w$ )
593	2:	$P'_a \leftarrow \text{OptRelative}(P_a);$
594	3:	output $P'_a$ ;
595	4:	$\mathcal{P} \leftarrow ImmPred(P'_a, s_w);$
500	5:	for each $(P''_a, s''_w) \in \mathcal{P}$
596	6:	AsyncSynthesis( $P''_a, s''_w$ );
597		

598

599

600

601

602

603

604

605

606

607

608

609

610

611 612 613  $P_a^1$  and  $P_a^2$  that are both optimal asynchronizations of P relative to  $P_a$ . Let  $\rho^1$  (resp.,  $\rho^2$ ) be an execution of  $P_a^1$  (resp.,  $P_a^2$ ) where every await \* does not suspend the execution of the current task, i.e.,  $\rho^1$  and  $\rho^2$  simulate the synchronous execution of *P*. Let  $s^1_w$  be the statement corresponding to the first await action in  $\rho^1$  such that (1) there exists an await action in  $\rho^2$  with the corresponding await statement  $s_w^2$ , such that  $s_w^1$  and  $s_w^2$  match the same call in *P*, and  $Cover(s_w^1) \subset Cover(s_w^2)$ (this holds because  $P_a^1$  and  $P_a^2$  are distinct asynchronizations of the same synchronous program, thus  $Cover(s_w^1)$  and  $Cover(s_w^2)$  must be comparable), and (2) for every other await statement  $s_w^3$  in  $P_a^1$  that generates an await action which occurs before the await action of  $s_w^1$  in  $\rho^1$ , there exists an await statement  $s_w^4$  in  $P_a^2$  matching the same call in *P*, such that  $Cover(s_w^3) = Cover(s_w^4)$ .

Let  $P_a^3$  be the program obtained from  $P_a^1$  by moving the await  $s_w^1$  down (further away from the matching call) such that  $Cover(s_w^1) = Cover(s_w^2)$ . Moving an await down can only create data races between actions that occur after the execution of the matching call. Then,  $P_a^3$  contains a data race iff there exists an execution  $\rho$  of  $P_a^3$  and two concurrent actions  $a_1$  and  $a_2$  that occur between the action (i, await(j)) generated by  $s_w^1$  and the action (i, call(j)) of the call matching  $s_w^1$ , such that:

$$((i, call(j)), a_1) \in CO, (a_1, a_w) \notin HB, ((i, call(j)), a_2) \in CO \text{ and } (a_2, (i, await(j))) \in HB$$

614 where the action  $a_w$  corresponds to the first await action in the task *j*. Let  $s_w$  be the statement 615 corresponding to the action  $a_w$ . Since the only difference between  $P_a^3$  and  $P_a^2$  is the placement of 616 awaits then  $((i, call(j)), a_1) \in CO$  and  $((i, call(j)), a_2) \in CO$  hold in any execution  $\rho'$  of  $P_a^2$  that 617 contains the actions  $a_1$  and  $a_2$ . Also, note that since  $a_w$  occurs in the task j that the action of  $s_w^1$  is 618 waiting for. This implies that in  $\rho^1$  the action of  $s_w$  occurs before the action of  $s_w^1$  in  $\rho^1$ . Therefore, 619 by the definition of  $s_w^1$  we have that  $s_w$  in  $P_a^1$  covers the same set of statements as the corresponding 620  $s'_w$  in  $P_a^2$  that matches the same call as  $s_w$ . Consequently,  $(a_1, a'_w) \notin HB$  and  $(a_2, (i, await(j))) \in HB$ 621 hold in any execution  $\rho'$  of  $P_a^2$  that contains the actions  $a_1$  and  $a_2$  ( $a'_w$  is the action of  $s'_w$ ). Thus, 622 there exists an execution  $\rho'$  of  $P_a^2$  such that the actions  $a_1$  and  $a_2$  are concurrent. This implies that if  $P_a^3$  admits a data race, then  $P_a^2$  admits a data race between actions generated by the same statements. 623 624 As  $P_a^2$  is data race free, we get that  $P_a^3$  is data race free as well. Since  $P_a^1 < P_a^3$ , we get that  $P_a^1$  is not 625 optimal, which contradicts the hypothesis. 626

### 5.2 Enumeration Algorithm

Our algorithm for enumerating all sound asynchronizations is given in Algorithm 1 as a recursive 628 procedure AsyncSynthesis that we describe in two phases. 629

First, we ignore the second argument of AsyncSynthesis (written in blue), which represents an 630 await instruction. For an asynchronization  $P_a$ , AsyncSynthesis outputs all sound asynchronizations 631 that are smaller than  $P_a$  w.r.t.  $\leq$ . It uses OptRelative to compute the optimal asynchronization 632  $P'_a$  of P relative to  $P_a$ , and then, calls itself recursively for all immediate predecessors of  $P'_a$  w.r.t. 633  $\leq$ . AsyncSynthesis outputs all sound asynchronizations of P when given as input the weakest 634 asynchronization of P. The delay complexity of this algorithm remains exponential in general, 635 because it may output a sound asynchronization multiple times. Indeed, because asynchronizations 636

are only partially ordered by  $\leq$ , different chains of recursive calls starting in different immediate predecessors may end up outputting the same asynchronization. For instance, looking at the asynchronizations of our motivating example on the right of Fig. 4, the asynchronization (0, 0, 0, 0) will be outputted twice because it is an immediate predecessor of both (0, 1, 0, 0) and (1, 0, 0, 0).

To avoid outputting the same solution twice, we use a refinement of the above that restricts the set of immediate predecessors available for a (recursive) call of AsyncSynthesis. This is based on a strict total order  $\prec_w$  between awaits in a program  $P_a$  that follows a topological ordering of its inter-procedural CFG, i.e., if  $s_w$  occurs before  $s'_w$  in the body of a method *m*, then  $s_w \prec_w s'_w$ , and if  $s_w$  occurs in a method m and  $s'_w$  occurs in a method m' s.t. m (indirectly) calls m', then  $s_w \prec_w s'_w$ . Therefore, AsyncSynthesis takes an await statement  $s_w$  as a second parameter, which is initially the maximal element w.r.t.  $\prec_w$ , and it calls itself only on immediate predecessors of an optimal solution obtained by moving up an await  $s''_w$  smaller than or equal to  $s_w$  w.r.t.  $\prec_w$ . The recursive call on that predecessor will receive as input  $s'_{w}$ . Formally, this relies on a function ImmPred that returns pairs of immediate predecessors and await statements defined as follows: 

$$\operatorname{ImmPred}(P'_a, s_w) = \{ (P''_a, s''_w) : P''_a < P'_a \text{ and } \forall P'''_a \in \operatorname{Async}[P, L, L_a]. P'''_a < P'_a \Longrightarrow P'''_a \le P''_a \\ \operatorname{And} s''_w \le w s_w \text{ and } P''_a \in P'_a \uparrow s''_w \}$$

 $(P'_a \uparrow s''_w \text{ is the set of asynchronizations obtained from } P'_a \text{ by changing only the position of } s''_w, moving it up w.r.t. the position in <math>P'_a$ ). For instance, looking at immediate predecessors of (1, 1, 0, 0) on the right of Fig. 4, (0, 1, 0, 0) is obtained by moving the *first* await in  $\prec_w$  and therefore, after computing the optimal solution relative to it, which is itself, it will explore no more immediate predecessors (ImmPred returns  $\emptyset$  because the input  $s_w$  is the minimal element of  $\prec_w$ , and it is already immediately after the matching call). Its immediate predecessor will be explored when recursing on (1, 0, 0, 0). The complexity analysis also relies on a property of the optimal asynchronization relative to an immediate predecessor: if the predecessor is defined by moving an await  $s''_w$ , then the optimal asynchronization is obtained by moving only awaits smaller than  $s''_w$  w.r.t.  $\prec_w$ .

LEMMA 5.3. If  $P''_a$  is an immediate predecessor of a sound asynchronization  $P'_a$ , which is defined by moving an await  $s''_w$  in  $P'_a$  up, then the optimal sound asynchronization relative to  $P''_a$  is obtained by moving only awaits smaller than  $s''_w$  w.r.t.  $\prec_w$ .

PROOF. Moving an await up in  $P'_a$  can only create data races between actions that occur after the execution of this await (because the invocation is suspended earlier). The only possible repairs of these data races consists in either moving  $s''_w$  down which results in  $P'_a$  or moving up some other awaits that occur in methods that (indirectly) call the method in which  $s''_w$  occurs. The first case is not applicable because it gives a program that is not smaller than  $P''_a$ . In the second case, every await  $s'_w$  that is moved up occurs in a method that (indirectly) calls the method in which  $s''_w$  occurs, and therefore,  $s'_w$  is smaller than  $s''_w$  w.r.t.  $\prec_w$ .

We show that Algorithm 1 returns all sound asynchronizations when called with the weakest asynchronization and the maximum await in  $\prec_w$ . Lemma 5.3 shows that after having computed an optimal sound asynchronization  $P'_a$  in a recursive call with parameter  $s_w$  any smaller sound asynchronization is also smaller than some predecessor in ImmPred( $P'_a$ ,  $s_w$ ). Thus, the restriction to a subset of predecessors is without loss of completeness (see also the supplementary material).

THEOREM 5.4. ASYNCSYNTHESIS(weakAsync[ $P, L, L_a$ ],  $s_w$ ), where  $s_w$  is maximal in weakAsync[ $P, L, L_a$ ] w.r.t.  $\prec_w$ , outputs all sound asynchronizations of P[L] w.r.t.  $L_a$ .

Viewing asychronization synthesis as an enumeration problem, the following theorem states its *delay* complexity in terms of an oracle  $O_{opt}$  that returns an optimal asynchronization relative to a given one. This follows from the fact that Algorithm 1 cannot return the same asynchonization

twice and computing immediate predecessors is polynomial time. The former is a consequence of Lemma 5.3. Thus, let  $P_a^1$  and  $P_a^2$  be two predecessors in ImmPred( $P'_a$ ,  $s_w$ ) obtained by moving up the awaits  $s_w^1$  and  $s_w^2$ , respectively, and assume that  $s_w^1 < w s_w^2$ . By Lemma 5.3, all solutions computed in the recursive call on  $P_a^1$  will have  $s_w^2$  placed as in  $P'_a$  while all the solutions computed in the recursive call on  $P_a^2$  will have  $s_w^2$  closer to the matching call. Therefore, the sets of solutions computed in these two branches of the recursion are distinct and the same solution cannot be outputted twice.

THEOREM 5.5. The delay complexity of the asychronization synthesis problem is polynomial time modulo  $O_{opt}$ .

### 696 6 COMPUTING OPTIMAL ASYNCHRONIZATIONS

We describe an approach for computing the optimal asynchronization relative to a given synchronization  $P_a$ , which can be seen as a way of repairing  $P_a$  so that it becomes data-race free. Intuitively, we repeatedly eliminate data races in  $P_a$  by moving certain await statements closer to the matching calls. The data races in  $P_a$  (if any) are enumerated in a certain order that prioritizes data races between actions that occur first in executions of the original synchronous program.

### 6.1 Data Race Ordering

693

694

695

702

703

We define an order between data races of asynchronizations based on the order between actions in executions of the synchronous program P. This order relates data races in possibly different executions or asynchronizations (of the same program), which is possible because each action in a data race corresponds to a statement in P (a read or a write to a program variable).

For two read/write statements *s* and *s'*, s < s' denotes the fact that there is an execution of *P* in which the *first* time *s* is executed occurs before the *first* time *s'* is executed. For two actions *a* and *a'* in an execution/trace of an asynchronization, generated<sup>2</sup> by two read/write statements *s* and *s'*, resp., we use  $a <_{SO} a'$  to denote the fact that s < s' and either  $s' \neq s$  or s' is reachable from *s* in the interprocedural<sup>3</sup> control-flow graph of *P* without taking any back edge<sup>4</sup>.

For a *deterministic* synchronous program (admitting a single execution),  $a \prec_{SO} a'$  iff  $s \prec s'$ . For non-deterministic programs, when *s* and *s'* are contained in a loop body, it is possible that  $s \prec s'$  and  $s' \prec s$ . For instance, the statements r1 = x and r2 = y of the program in Fig. 11 can be executed in different orders depending on the number of loop iterations and whether the if branch is entered during the first loop iteration. In this case, we use the control-flow order to break the tie between *a* and *a'*.

The order between data races corresponds to the colexicographic order induced by  $\prec_{SO}$ . This is a partial order since actions may originate from different control-flow paths and are incomparable w.r.t.  $\prec_{SO}$ .

*Definition 6.1 (Data Race Order).* Given two races  $(a_1, a_2)$  and  $(a_3, a_4)$  admitted by (possibly different) asynchronizations of a synchronous program *P*, we have that  $(a_1, a_2) \prec_{SO} (a_3, a_4)$  iff  $a_2 \prec_{SO} a_4$ , or  $a_2 = a_4$  and  $a_1 \prec_{SO} a_3$ .

async method Main { r1 = call m;if \* r2 = x;x = r2 + 1;else r3 = x:await r1; } async method m { await \* retVal = x; x = input;return; } Fig. 12

method Main {
 while \*

if \*

Fig. 11

r1 = x;

r2 = y; }

*Example 6.2.* For the program in Fig. 12, we have the following order between data races: (x = input, r2 = x)  $<_{SO}$  (retVal = x, x = r2 + 1) because r2 = x is executed before the write <sup>2</sup>Each action labels a transition in the operational semantics (Section 3), and each transition corresponds to executing a

735

724

725

726

statement. This statement is said to generate the action.

<sup>&</sup>lt;sup>3</sup>The interprocedural graph is the union of the control-flow graphs of each method along with edges from call sites to entry nodes, and from exit nodes to return sites.

 <sup>&</sup>lt;sup>4</sup>A back edge points to a block that has already been met during a depth-first traversal of the control-flow graph, and
 corresponds to loops.

 $x = r^2 + 1$  in the original synchronous program (for simplicity we use statements instead of 736 actions). However, the data races  $(x = input, r^2 = x)$  and  $(x = input, r^3 = x)$  are incomparable. 737

#### **Repairing Data Races** 739 6.2

740 Repairing a data race  $(a_1, a_2)$  reduces to modifying the position of a certain await. In general, we can either move an await down (further away from the matching call), for instance in the method executing  $a_1$ , or move an await up (closer to the matching call), for instance in the method executing  $a_2$ . For example, the data race between x = input and  $r^2 = x$  on the right of Fig. 10 can 744 be repaired by either moving await \* in *m* after the write x = input, so the call to *m* is suspended 745 later, or await r1 in Main before r2 = x, to restrict the set of statements that can execute before m746 finishes. In the following, we consider only repairs where awaits are moved up. The "completeness" 747 of this set of repairs follows from the particular order in which we enumerate data races. Intuitively, 748 moving the other await down would introduce a data race we have already repaired.

In general,  $a_1$  may not occur in a method m' that is called directly by m, as in Fig. 10, but in another method called by m' or even further down the call tree. Similarly,  $a_2$  may not be part of m, but it may be included in another method called by m after calling m' (but before await r), and so on. Next, we describe precisely the transformation that suffices to repair a given data race.

Any two racing actions have a common ancestor in the call order CO which is a call action. This is at least the call action of main. The least common ancestor of  $a_1$  and  $a_2$  in CO among call actions is denoted by LCA<sub>CO</sub>( $a_1, a_2$ ). Formally, LCA<sub>CO</sub>( $a_1, a_2$ ) is a call action  $a_c = (i, \text{call}(j))$  s.t.  $(a_c, a_1) \in CO$ ,  $(a_c, a_2) \in CO$ , and for each other call action  $a'_c$ , if  $(a_c, a'_c) \in CO$  then  $(a'_c, a_1) \notin CO$ . For instance, the call action corresponding to r1 = call m on the right of Fig. 10 is the *least* common ancestor of the racing actions discussed above. The following lemma (see the supplementary material for a proof) shows that this is the asynchronous call for which the matching await must be moved in order to repair a given data race. It also identifies the position where the await matching  $LCA_{CO}(a_1, a_2)$ should be moved in order to repair the data race. Intuitively, this is just before  $a_2$  if  $a_2$  is in the same method as  $LCA_{CO}(a_1, a_2)$ , or more generally, just before the last statement in the same method which precedes  $a_2$  in the call order. On the right of Fig. 10, await r1 has to be moved before the statement  $r^2 = x$ , which plays the role of  $a_2$ .

LEMMA 6.3. Let  $(a_1, a_2)$  be a data race in a trace  $\tau$  of an asynchronization  $P_a$ , and  $a_c = (i, call(j)) =$  $LCA_{CO}(a_1, a_2)$ . Then,  $\tau$  contains a unique action  $a_w = (i, await(j))$  and a unique action a such that:

•  $(a, a_w) \in MO$ , and a is the latest action in the method order MO such that  $(a_c, a) \in MO$  and  $(a, a_2) \in CO^*$  (CO<sup>\*</sup> denotes the reflexive closure of CO).

Lemma 6.3 identifies a sufficient transformation for repairing a data race  $(a_1, a_2)$ : moving the await  $s_w$  generating the action  $a_w$  just before the statement s generating a. This is sufficient because it ensures that every statement that follows  $LCA_{CO}(a_1, a_2)^5$  in call order will be executed before a and before any statement which succeeds *a* in call order, including  $a_2$ . Note that moving the await  $a_w$  anywhere after a will not affect the concurrency between  $a_1$  and  $a_2$ .

The pair  $(s_c, s)$ , where  $s_c$  is the call statement generat-779 ing  $LCA_{CO}(a_1, a_2)$ , is called the *root cause* of the data race 780  $(a_1, a_2)$ . Let RepDRace $(P_a, s_c, s)$  be the maximal asynchronization  $P'_a$  smaller than  $P_a$  w.r.t.  $\leq$ , s.t. no await statement 782

741 742 743

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

781

783 784

async method Main { async method Main { r1 = call m;r1 = call m;if \* if \* r2 = x;await r1; else r2 = x;r3 = y;else await r1; r3 = y;await r1; } async method m { } await \* async method m { retVal = x; await \* x = input;retVal = x; x = input;return: } return: }

Fig. 13. Examples of asynchronizations.

<sup>&</sup>lt;sup>5</sup>We abuse the terminology and make no distinction between statements and actions.

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

785	Algorithm 2 The procedure OPTRELATIVE for computing the optimal asynchronization of P relative
786	to $P_a$ . ROOTCAUSEMINDRACE( $P'_a$ ) returns the root cause of a minimal data race of $P'_a$ w.r.t. $\prec_{SO}$ , or
787	$\perp$ if $P'_a$ is data race free.

788	1: p	<b>rocedure</b> OptRelative( $P_a$ )
789	2:	$P'_a \leftarrow P_a$
790	3:	$root \leftarrow \text{RootCauseMinDRace}(P'_a)$
791	4:	while $root \neq \bot$
702	5:	$P'_a \leftarrow RepDRace(P'_a, root)$
172	6:	$root \leftarrow \text{RootCauseMinDRace}(P'_a)$
793	7:	return $P'_a$
704		

795 matching  $s_c$  occurs after s on a CFG path. When the control-flow graph of the method contains 796 branches, the construction of RepDRace( $P_a, s_c, s$ ) consists of (1) replacing all await statements 797 matching  $s_c$  that are reachable in the CFG from s with a single await statement placed just before s, 798 and (2) adding additional await statements in branches that "conflict" with the branch containing 799 s. This is to ensure the syntactic constraints described in §3. These additional await statements are 800 at maximal distance from the corresponding call statement because of the maximality requirement. 801 For instance, to repair the data race between  $r^2 = x$  and x = input in the program on the left of 802 Fig. 13, the statement await r1 must be moved before r2 = x in the if branch, which implies that 803 another await must be added on the else branch. The result is given on the right of Fig. 13.

The following shows that repairing a minimal data race cannot introduce smaller data races (w.r.t.  $\prec_{SO}$ ), which ensures some form of monotonicity when repairing minimal data races iteratively.

LEMMA 6.4. Let  $P_a$  be an asynchronization,  $(a_1, a_2)$  a data race in  $P_a$  that is minimal w.r.t.  $\prec_{SO}$ , and  $(s_c, s)$  the root cause of  $(a_1, a_2)$ . Then, RepDRace $(P_a, s_c, s)$  does not admit a data race that is smaller than  $(a_1, a_2)$  w.r.t.  $\prec_{SO}$ .

**PROOF.** The only modification in the program  $P'_a$  = 810 RepDRace( $P_a$ ,  $s_c$ , s) compared to  $P_a$  is the movement 811 of the await  $s_w$  matching the call  $s_c$  to be before 812 the statement s in a method m. The concurrency 813 added in  $P'_a$  that was not possible in  $P_a$  is between 814 actions (a', a'') generated by statements s' and s'', 815 respectively, as shown in Fig. 14. W.l.o.g., we assume 816 that  $(a', a'') \in$  SO. The statements  $s_1$  and  $s_2$  are those 817 generating  $a_1$  and  $a_2$ , respectively. The statement 818

s' is related by  $CO^*$  to some statement in *m* that



Fig. 14. An excerpt of an asynchronous program.

follows *s*, and *s*'' is related by CO<sup>\*</sup> to some statement that follows the call to *m* in the caller of *m*. Note that *s*' is ordered by < after  $s_2$ . Since  $(a_1, a_2) \in$  SO and  $(a', a'') \in$  SO then  $s_2 < s''$  and  $s_1 < s'$ . Thus, any new data race (a', a'') in  $P'_a$  that was not reachable in  $P_a$  is bigger than  $(a_1, a_2)$ .

### 6.3 A Procedure for Computing Optimal Asynchronizations

Given an asynchronization  $P_a$ , the procedure OPTRELATIVE in Algorithm 2 computes the optimal asynchronization relative to  $P_a$  by repairing data races iteratively until the program becomes data race free. The following theorem states that correctness of this procedure.

THEOREM 6.5. Given an asynchronization  $P_a \in \text{Async}[P, L, L_a]$ , OptRelative( $P_a$ ) returns the optimal asynchronization of P relative to  $P_a$ .

PROOF. We need to show that any immediate successor  $P_a^1$  of the output  $P'_a = \text{OPTRELATIVE}(P_a)$  that is also smaller than  $P_a$  (w.r.t.  $\leq$ ) admits data races. By the definition of  $\leq$ ,  $P_a^1$  is obtained by moving exactly one await statement  $s_w$  in a method m of  $P'_a$  further away from the matching call  $s_c$ .

804

805

806

807

808

809

819

820

821

822 823

824

825

826

827

828

829

830

Since  $P_a^1 \leq P_a$ , the position of  $s_w$  in the output  $P_a'$  is due to repairing a data race between two actions 834  $a_1$  and  $a_2$  with a root-cause  $(s_c, s)$ , for some s, on some program  $P'_a \leq P''_a \leq P_a$ . We show that these 835 actions form a data race in  $P_a^1$ . These actions are reachable in an execution of  $P_a^1$  because every 836 method m' that is called by m between  $s_c$  and  $s_w$  ( $s_c$  included), or that follows m' in the call-graph 837 of  $P_a^1$  (or  $P_a''$ ) has exactly the same code as in  $P_a''$ , i.e., the placement of the awaits in those methods 838 is the same as in  $P''_a$  (call graphs remain identical between different asynchronizations). This is due 839 to the fact that any data race that would lead to moving an await in one of those methods is before 840  $(a_1, a_2)$  in the order  $\prec_{SO}$ . Since  $s_w$  in  $P_a^1$  is placed after s, we get that  $a_1$  and  $a_2$  are also concurrent 841 in that execution of  $P_a^1$ , which concludes the proof. 842 

OPTRELATIVE( $P_a$ ) iterates the process of repairing a data race a number of times which is linear in the size of the input. Indeed, each iteration of the loop results in moving an await closer to the matching call and before at least one more statement from the original synchronous program P.

The fact that data races are enumerated in the order defined by  $\prec_{SO}$  guarantees a bound on the number of times an await matching the same call is moved during the execution of OPTRELATIVE( $P_a$ ). In general, this bound is the number of statements covered by all the awaits matching the call in the input program  $P_a$ . Actually, this is a rather coarse bound. A more refined analysis has to take into account the number of branches in the CFGs. For programs without conditionals or loops, every await is moved at most once during the execution of OPTRELATIVE( $P_a$ ). In the presence of branches, a call to an asynchronous method may match multiple await statements (one for each CFG path starting from the call), and the data races that these await statements may create may be incomparable w.r.t.  $\prec_{SO}$ . Therefore, for a call statement  $s_c$ , let  $|s_c|$  be the sum of  $|Cover(s_w)|$  for every await  $s_w$  matching  $s_c$  in  $P_a$ .

LEMMA 6.6. For any asynchronization  $P_a \in \text{Async}[P, L, L_a]$  and call statement  $s_c$  in  $P_a$ , the while loop in OptRelative( $P_a$ ) does at most  $|s_c|$  iterations that result in moving an await matching  $s_c$ .

PROOF. We consider first the case without conditionals or loops, and we show by contradiction that every await statement  $s_w$  is moved at most once during the execution of OPTRELATIVE( $P_a$ ), i.e., there exists at most one iteration of the while loop which changes the position of  $s_w$ . Suppose that the contrary holds for an await  $s_w$ . Let  $(a_1, a_2)$ , and  $(a_3, a_4)$  be the data races repaired by the first and second moves of  $s_w$ , respectively. By Lemma 6.3, there exist two actions a and a' such that

$$(a_c, a) \in MO, (a, a_2) \in CO^*, (a, a_w) \in MO \text{ and } (a_c, a') \in MO, (a', a_4) \in CO^*, (a', a_w) \in MO$$

where  $a_w = (i, \operatorname{await}(j))$  and  $a_c = (i, \operatorname{call}(j))$  are the asynchronous call action and the matching await action. Let  $s_2$  and  $s_4$  be the statements generating the two actions  $a_2$  and  $a_4$ , respectively. Then, we have either  $s_2 < s_4$  or  $s_2 = s_4$ , and both cases imply that  $(a, a') \in MO^*$ . Thus, moving the await statement generating  $a_w$  before the statement generating a implies that it is also placed before the statement generating a' (that occurs after a in the same method). Thus, the first move of the await  $s_w$  repaired both data races, which is contradiction.

In the presence of conditionals or loops, moving an await up in one branch may correspond to adding multiple awaits in the other conflicting branches. Also, one call in the program may correspond to multiple awaits on different branches. However, every repair of a data race consists in moving one await closer to the matching call  $s_c$  and before one more statement covered by some await matching  $s_c$  in the input  $P_a$ .

### 6.4 Computing Root Causes of Minimal Data Races

We present a reduction from the problem of computing root causes of minimal data races to
reachability (assertion checking) in sequential programs. This reduction builds on a program
instrumentation for checking if there exists a minimal data race that involves two given statements

882

877

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858 859

860

861

862

863

```
Replace every statement ``await r'' with:
883
             Add before s_1:
         1
                                                                           13
         2
             if ( lastTaskDelayed == ⊥ && * )
                                                                                  if( r == lastTaskDelayed ) then
                                                                           14
884
               lastTaskDelayed := myTaskId();
         3
                                                                           15
                                                                                     if ( !DescendantDidAwait )
885
         4
               DescendantDidAwait := thisHasDoneAwait;
                                                                           16
                                                                                       DescendantDidAwait := thisHasDoneAwait;
886
         5
                                                                                     lastTaskDelayed := myTaskId();
               return
                                                                           17
                                                                           18
                                                                                     return
887
         7
             Add before s2:
                                                                           19
                                                                                  else
888
         8
               if ( task_s<sub>c</sub> == myTaskId() )
                                                                           20
                                                                                     thisHasDoneAwait := true
889
         9
                 s := s_2;
        10
               assert (lastTaskDelayed == ⊥ || !DescendantDidAwait); 22
                                                                                Add before every statement ``r := call m'':
890
                                                                                  if ( task_{s_c} == myTaskId() ) then
                                                                           23
891
                                                                           24
                                                                                     s := this statement;
892
                                                                           26
                                                                                Add after every statement ``r := call m'':
893
                                                                                  if ( r == lastTaskDelayed )
                                                                           27
894
                                                                           28
                                                                                     sc := this statement;
895
                                                                           29
                                                                                     task_s<sub>c</sub> := myTaskId();
```

Fig. 15. A program instrumentation for computing the root cause of a minimal data race between the statements  $s_1$  and  $s_2$  (if any). All variables except for thisHasDoneAwait are program (global) variables. thisHasDoneAwait is a local variable. The value  $\perp$  represents an initial value of a variable. The variables  $s_c$ and s store the (program counters of the) statements representing the root cause. The method myTaskId returns the id of the current task.

 $(s_1, s_2)$ , whose correctness relies on the assumption that another pair of statements cannot produce a smaller data race. This instrumentation is used in an iterative process where pairs of statements are enumerated according to the colexicographic order induced by  $\prec$ . This specific enumeration ensures that the assumption made for the correctness of the instrumentation is satisfied.

Given an asynchronization  $P_a$ , the instrumentation described in Fig. 15 represents a synchronous 905 program where all await statements are replaced with synchronous code (lines 14-20). This 906 instrumentation simulates asynchronous executions of  $P_a$  where methods may be only partially 907 executed, modeling await interruptions. It reaches an error state (see the assert at line 10) when 908 an action generated by  $s_1$  is concurrent with an action generated by  $s_2$ , which represents a data 909 race, provided that  $s_1$  and  $s_2$  access a common program variable (these statements are assumed to 910 be given as input). Also, the values of  $s_c$  and s when reaching the assertion violation represent the 911 root-cause of this data race. 912

The instrumentation simulates an execution of  $P_a$  to search for a data race as follows (we discuss the identification of the root-cause afterwards):

- It executes under the synchronous semantics until an instance of  $s_1$  is non-deterministically chosen as a candidate for the first action in the data race ( $s_1$  can execute multiple times if it is included in a loop for instance). The current invocation is interrupted when it is about to execute this instance of  $s_1$  and its task id  $t_0$  is stored into lastTaskDelayed (see lines 2–5).
  - Every invocation that transitively called  $t_0$  is interrupted when an await for an invocation in this call chain (whose task id is stored into lastTaskDelayed) would have been executed in the asynchronization  $P_a$  (see line 18).
    - Every other method invocation is executed until completion as in the synchronous semantics.
- When reaching  $s_2$ , if  $s_1$  has already been executed (lastTaskDelayed is not  $\perp$ ) and at least one invocation has only partially been executed, which is recorded in the boolean flag DescendantDidAwait and which means that  $s_1$  is concurrent with  $s_2$ , then the instrumentation stops with an assertion violation.

A subtle point is that the instrumentation may execute code that follows an await r even if the task r has been executed only partially, which would not happen in an execution of the original  $P_a$ . Here, we rely on the assumption that there exist no data race between that code and the rest of the task r. Such data races would necessarily involve two statements which are before  $s_2$  w.r.t.  $\prec$ .

919

920

921

Therefore, the instrumentation is correct only if it is applied by enumerating pairs of statements  $(s_1, s_2)$  w.r.t. the colexicographic order induced by  $\prec$ .

Next, we describe the computation of the root-cause, i.e., the updates on the variables  $s_c$  and s. By definition, the statement  $s_c$  in the root-cause should be a call that makes an invocation that is in the call stack when  $s_1$  is reached. This can be checked using the variable lastTaskDelayed that stores the id of the last such invocation popped from the call stack (see the test at line 27). The statement s in the root-cause can be any call statement that has been executed in the same task as  $s_c$  (see the test at line 23), or  $s_2$  itself (see line 9).

Let  $[[P_a, s_1, s_2]]$  denote the instrumentation in Fig. 15. We say that the values of  $s_c$  and s when reaching the assertion violation are the root cause computed by this instrumentation. The following theorem states its correctness.

THEOREM 6.7. If  $[[P_a, s_1, s_2]]$  reaches an assertion violation, then it computes the root cause of a minimal data race, or there exists  $(s_3, s_4)$  such that  $[[P_a, s_3, s_4]]$  reaches an assertion violation and  $(s_3, s_4)$  is before  $(s_1, s_2)$  in colexicographic order w.r.t. <.

Based on Theorem 6.7, we define an implementation of the procedure ROOTCAUSEMINDRACE ( $P_a$ ) used in computing optimal asynchronizations (Algorithm 2) as follows:

- For all pairs of read or write statements  $(s_1, s_2)$  in colexicographic order w.r.t.  $\prec$ .
  - If  $[[P_a, s_1, s_2]]$  reaches an assertion violation, then
    - \* return the root cause computed by  $[[P_a, s_1, s_2]]$
- return  $\perp$

The order  $\prec$  between read or write statements can be computed using a quadratic number of reachability queries in the synchronous program *P*. Therefore,  $s \prec s'$  iff an instrumentation of *P* that sets a flag when executing *s* and asserts that this flag is not set when executing *s'* reaches an assertion violation. The following theorem states the correctness of the procedure above.

THEOREM 6.8. ROOTCAUSEMINDRACE( $P_a$ ) returns the root cause of a minimal data race of  $P_a$  w.r.t.  $\prec_{SO}$ , or  $\perp$  if  $P'_a$  is data race free.

This procedure performs a quadratic number of reachability queries in sequential programs.

THEOREM 6.9. The complexity of ROOTCAUSEMINDRACE is polynomial time modulo an oracle for the reachability problem in sequential programs.

## 6.5 Asymptotic Complexity of Asynchronization Synthesis

We state the complexity of the asynchronization synthesis problem. Theorem 5.5 shows that its delay complexity is polynomial modulo the complexity of OPTRELATIVE in Algorithm 2, which by the results in this section, reduces to a polynomial number of reachability queries in sequential programs. The reachability problem is PSPACE-complete for finite-state sequential programs [Godefroid and Yannakakis 2013].

THEOREM 6.10. The output complexity<sup>6</sup> and delay complexity of the asynchronization synthesis problem is polynomial time modulo an oracle for reachability in sequential programs, and PSPACE for finite-state programs.

This result is optimal, i.e., checking whether there exists a sound asynchronization which is different from the trivial strong synchronization is PSPACE-hard (follows from a reduction from the reachability problem).

THEOREM 6.11. Checking whether there exists a sound asynchronization different from the strong <u>asynchronization is PSPACE-complete</u>.

<sup>6</sup>Note that all asynchronizations can be enumerated with polynomial space.

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

943

944

945

946

947

948

949

950

951 952

953

954

955

956

957

958

959

960

961 962

963

964

965

966

967

968

969 970

971

972

973

974

975

976

977

**PROOF.** For hardness, checking if a sequential program *P* reaches a particular control location 981  $\ell$  can be reduced to the non-existence of a non-trivial sound asynchronization of a program P'982 983 defined as follows: (1) define a new method *m* that writes to a new program variable *x*, and insert a call to *m* followed by a write to x at location  $\ell$ , and (2) insert a write to x after every call statement 984 that calls a method in  $\{m'\}^*$ , where m' is the method containing  $\ell$ . Let  $m_a$  be an asynchronous 985 version of *m* obtained by inserting an await \* at the beginning. Then,  $\ell$  is reachable in *P* iff the 986 only sound asynchronization of P' w.r.t.  $\{m_a\}$  is the strong asynchronization. 987 П

## **OPTIMAL ASYNCHRONIZATIONS USING DATA-FLOW ANALYSIS**

990 We present a procedure for computing sound asynchronizations, based on a bottom-up interprocedural data-flow analysis. It computes optimal asynchronizations for abstractions of programs 991 where every Boolean condition in if-then-else statements or while loops is replaced with the 992 993 non-deterministic choice \*.

For a program P, we define an abstraction  $P^{\#}$  where every conditional if  $\langle le \rangle \{S_1\}$  else  $\{S_2\}$  is 994 rewritten to if  $\{S_1\}$  else  $\{S_2\}$ , and every while  $\langle le \rangle \{S\}$  is rewritten to if  $\{S\}$ . Besides adding 995 the non-deterministic choice \*, loops are unrolled exactly once. Every asynchronization  $P_a$  of P 996 corresponds to an abstraction  $P_a^{\#}$  obtained by applying exactly the same rewriting. 997

 $P^{\#}$  is a sound abstraction of P in terms of sound asynchronizations it admits. Unrolling loops once is sound because every asynchronous call in a loop iteration should be waited for in the same 1000 iteration (see the syntactic constraints in §3).

THEOREM 7.1. If  $P_a^{\#}$  is a sound asynchronization of  $P^{\#}$  w.r.t.  $L_a$ , then  $P_a$  is a sound asynchronization of P w.r.t.  $L_a$ .

We present a procedure for computing optimal asynchronizations of  $P^{\#}$ , relative to a given 1004 asynchronization  $P_a^{\#}$ . This procedure traverses methods of  $P_a^{\#}$  in a bottom-up fashion, detects data 1005 races using summaries of read/write accesses computed using a straightforward data-flow analysis, 1006 and repairs data races using the schema presented in Section 6.2. Applying this procedure to a real 1007 programming language requires an alias analysis to detect statements that may access the same 1008 memory location (this is trivial in our language whose purpose is to simplify the exposition). 1009

We consider an enumeration of methods called *bottom-up order*, which is the reverse of a 1010 topological ordering of the call graph<sup>7</sup>. For each method m, let  $\mathcal{R}(m)$  be the set of program variables 1011 that *m* can read, which is defined as the union of  $\mathcal{R}(m')$  for every method *m'* called by *m* and the set 1012 of program variables read in statements in the body of m. The set of variables  $\mathcal{W}(m)$  that m can write 1013 is defined in a similar manner. We define RW-var $(m) = (\mathcal{R}(m), \mathcal{W}(m))$ . We extend the notation 1014 RW-var to statements as follows: RW-var( $\langle r \rangle := \langle x \rangle$ ) = ( $\{x\}, \emptyset$ ), RW-var( $\langle x \rangle := \langle le \rangle$ ) = ( $\emptyset, \{x\}$ ), 1015 RW-var(r := call m) = RW-var(m), and RW-var $(s) = (\emptyset, \emptyset)$ , for any other type of statement s. 1016 Also, let CRW-var(*m*) be the set of read or write accesses that *m* can do and that can be concurrent 1017 with accesses that a caller of *m* can do after calling *m*. These correspond to read/write statements 1018 that follow an await in m, or to accesses in CRW-var(m') for a method m' called by m. These sets 1019 of accesses can be computed using the following data-flow analysis: for all methods  $m \in P_a^{\#}$  in 1020 bottom-up order, and for each statement *s* in the body of *m* from begin to end, 1021

1022 1023

1025

1026

988

989

998

999

1001

1002

1003

• CRW- $var(m) \leftarrow CRW$ - $var(m) \cup CRW$ -var(m')

- 1024 • If s is reachable from an await statement in the CFG of m
  - CRW-var $(m) \leftarrow$  CRW-var $(m) \cup$  RW-var(s)

• If s is a call to m' and s is not reachable from an await in the CFG of m

<sup>&</sup>lt;sup>7</sup>The nodes of the call graph are methods and there is an edge from a method  $m_1$  to a method  $m_2$  if  $m_1$  contains a call 1027 statement that calls  $m_2$ . 1028

We use  $(\mathcal{R}_1, \mathcal{W}_1) \bowtie (\mathcal{R}_2, \mathcal{W}_2)$  to denote the fact that  $\mathcal{W}_1 \cap (\mathcal{R}_2 \cup \mathcal{W}_2) \neq \emptyset$  or  $\mathcal{W}_2 \cap (\mathcal{R}_1 \cup \mathcal{W}_1) \neq \emptyset$ (i.e., a conflict between read/write accesses). We define the procedure OPTRELATIVE<sup>#</sup> that given an asynchronization  $P_a^{\#}$  works as follows:

• For all methods  $m \in P_a^{\#}$  in bottom-up order, and for each statement *s* in the body of *m* from begin to end,

- If s occurs between r := call m' and await r (for some m'), and RW-var(s)  $\bowtie$ 

1049 1050

1051

1052

1053

1054 1055

1033

CRW-var(m'), then  $P_a^{\#} \leftarrow \text{RepDRace}(P_a^{\#}, r := \text{call } m', s)$ 

• Return  $P_a^{\#}$ 

1038 The following theorem states the correctness of OptRelative<sup>#</sup>. This procedure repairs data races 1039 in an order which is  $\prec_{SO}$  with some exceptions that do not affect optimality, i.e., the number of 1040 times an await matching the same call can be moved. For instance, if a method *m* calls two other 1041 methods  $m_1$  and  $m_2$  in this order, the procedure above may handle  $m_2$  before  $m_1$ , i.e., repair data 1042 races between actions that originate from  $m_2$  before data races that originate from  $m_1$ , although 1043 the former are bigger than the latter in  $\leq_{SO}$ . This does not affect optimality because those repairs 1044 are "independent", i.e., any repair in  $m_2$  cannot influence a repair in  $m_1$ , and vice-versa. The crucial 1045 point is that this procedure repairs data races between actions that originate from a method *m* 1046 before data races that involve actions in methods preceding *m* in the call graph, which are bigger 1047 in  $\prec_{SO}$  than the former. 1048

THEOREM 7.2. OPTRELATIVE<sup>#</sup> ( $P_a^{\#}$ ) returns an optimal asynchronization relative to  $P_a^{\#}$ .

Since OptRelative<sup>#</sup> is based on a single bottom-up traversal of the call graph of the input asynchronization  $P_a^{\#}$ . Theorem 5.5 implies the following result.

THEOREM 7.3. The delay complexity of the asynchronization synthesis problem restricted to abstracted programs  $P^{\#}$  is polynomial time.

### 1056 8 MULTI-THREADED REFACTORINGS

We discuss an extension of our framework to *multi-threaded refactorings* that rewrite a sequential program into a multi-threaded program where every method invocation is executed on a different thread. A caller can wait for a callee to complete using a join primitive. A start primitive for spawning a new thread is the counterpart of an asynchronous call while join is the counterpart of await. For instance, Fig. 16 lists a sequential program, a possible asynchonization, and a multi-thread refactoring (both refactorings place the awaits/joins as far away as possible from the calls).

An important difference between start/join and async/await is the happens-before order relation. 1063 For instance, the asynchronization on the center of Fig. 16 assigns 1 to x (line 11) before it assigns 1064 2 to x (line 4), as in the original sequential program. However, the multi-thread program on the 1065 right of Fig. 16 may execute these two assignments in any order, and admits a behavior that is not 1066 possible in the sequential program (assigning 2 before assigning 1). Repairing this data-race consists 1067 in moving the join at line 5 to occur before assigning 2 to x at line 4. In general, the happens-before 1068 order is weaker compared to an analogous asynchronization, where awaits are placed as the joins, 1069 which implies that any multi-threaded refactoring can be rewritten to an asynchronization. The 1070 vice-versa may not be possible as shown in this example. 1071

Despite this difference, it can still be proved that there exists a unique multi-threaded refactoring that is sound, i.e., does not admit data races, and optimal, i.e., maximizes the distance between start and join, a result similar to Lemma 5.2. Assuming by contradiction the existence of two incomparable optimal and sound refactorings, one can show that moving a join in one refactoring further away from the matching call as in the other refactoring does not introduce data races (contradicting optimality). To compute optimal and sound multi-threaded refactorings, one can

<sup>1078</sup> 

```
1079
                                                                                          1 void Main() {
         1 void Main() {
                                            1 async Task MainAsync() {
         2 F();
                                            2 Task t1 = F();
                                                                                          2 Thread thr1 = new Thread(F);
1080
                                                                                          3 thr1.Start();
1081
         4 x = 2;
                                            4 \times = 2:
                                                                                          4 \times = 2:
1082
                                                                                          5 thr1.Join();
                                            5 await t1;
         6 }
                                            6 }
                                                                                          6 }
1083
1084
         8 void F() {
                                            8 async Task F() {
                                                                                          8 void F() {
1085
         9 IO();
                                            9 Task t2 = IOAsync();
                                                                                          9 Thread thr2 = new Thread(IO);
                                                                                         10 thr2.Start();
1086
        11 \times = 1;
                                           11 \times = 1;
                                                                                         11 \times = 1;
1087
                                           12 await t2;
                                                                                         12 thr2.Join();
1088
        13 }
                                           13 }
                                                                                         13 }
```

Fig. 16. A synchronous C# program, an asynchronization, and a multi-threaded refactoring.

apply the same iterative process of repairing data-races (the happens-before reflects multi-threading
 instead of async/await), prioritizing data races involving statements that would execute first in the
 sequential program. The repairing of a data-race is similar and consists in moving a join up.

In contrast to async/await, moving a join up does not introduce new data races (since no new parallelism is introduced). This implies that all the predecessors of a sound multi-threaded refactoring are also sound, i.e., the set of sound multi-threaded refactorings is downward closed.

## 1098 9 EXPERIMENTAL EVALUATION

We present an empirical evaluation of our asynchronization enumeration approach, where optimal asynchronizations are computed using the data-flow analysis described in Section 7. We consider a benchmark consisting mostly of asynchronous C# programs extracted from open-source GitHub projects. We evaluate the effectiveness of our technique in reproducing the original program as an asynchronization of a program where asynchronous calls are reverted to synchronous calls, along with other sound asynchronizations.

Implementation. We developped a prototype tool that relies on the Roslyn .NET compiler plat form [Roslyn 2021] to construct CFGs for methods in a given C# program. This prototype sup ports C# programs written in SSA form that include basic conditional or looping constructs and
 async/await as concurrency primitives. It assumes that any alias information is provided apriori;
 these constraints can be removed in the future with more engineering effort. Object fields are
 interpreted as program variables in the terminology of the program syntax in Section 3 (data races
 concern accesses to object fields).

The tool takes as input a possibly asynchronous program, and a mapping between synchronous and asynchronous variations of base methods in this program. It reverts every asynchronous call to a synchronous call, and it enumerates sound asynchronizations of the obtained program (using Algorithm 1).

Benchmark. Our evaluation uses a benchmark outlined in Table 2. This contains 5 synthetic 1117 examples (variations of the program in Fig. 1), 9 programs extracted from open-source C# GitHub 1118 projects (their name is a prefix of the repository name), and 2 programs inspired by questions on 1119 stackoverflow.com about async/await in C# (their name ends in Stackoverflow). Overall, there 1120 are 13 base methods involved in computing asynchronizations of these programs (that have both 1121 synchronous and asynchronous versions), which come from 5 C# libraries (System.IO, System.Net, 1122 Windows.Storage, Microsoft.WindowsAzure.Storage, and Microsoft.Azure.Devices). They are modeled 1123 as described in Section 3. 1124

**Evaluation.** The last five columns of Table 2 list data concerning the application of our tool. The column async lists the number of outputted sound asynchronizations. In general, the number of

1127

1089

1090

1128	Table 2. Empirical results. Syntactic characteristics of input programs: lines of code (loc), number of meth-
1129	ods (m), number of method calls (c), number of asynchronous calls (ac), number of awaits that could be
1130	placed at least one statement away from the matching call (await#). Data concerning the enumeration of
1121	asynchronizations: number of awaits that were placed at least one statement away from the matching call
1151	(await), number of races discovered and repaired (races), number of statements that the awaits in the optimal
1132	asynchronization are covering more than in the input program (cover), number of computed asynchronizations
1133	(async), and running time (t).
1124	

1134											
1135	Program	loc	m	с	ac	await#	await	races	cover	async	t(s)
1136	SyntheticBenchmark-1	77	3	6	5	4	4	5	0	9	5
1150	SyntheticBenchmark-2	115	4	12	10	6	3	3	0	8	5
1137	SyntheticBenchmark-3	168	6	16	13	9	7	4	0	128	9
1138	SyntheticBenchmark-4	171	6	17	14	10	8	5	0	256	55
1150	SyntheticBenchmark-5	170	6	17	14	10	8	9	0	272	138
1139	Azure-Remote	520	10	14	5	0	0	0	0	1	5
1140	Azure-Webjobs	190	6	14	6	1	1	0	1	3	4
1140	FritzDectCore	141	7	11	8	1	1	0	1	2	5
1141	MultiPlatform	53	2	6	4	2	2	0	2	4	5
1149	NetRpc	887	13	18	11	4	1	3	0	3	5
1174	TestAZureBoards	43	3	3	3	0	0	0	0	1	4
1143	VBForums-Viewer	275	7	10	7	3	2	1	1	6	5
1144	Voat	178	3	6	5	2	1	1	1	4	10
1144	WordpressRESTClient	133	3	10	8	4	2	1	0	4	5
1145	ReadFile-Stackoverflow	47	2	3	3	1	0	1	0	1	6
1146	UI-Stackoverflow	50	3	4	4	3	3	3	0	12	5

1147 asynchronizations depends on the number of invocations (column ac in Table 1) and the size of the 1148 code blocks between an invocation and the instruction using its return value (column await# gives 1149 the number of non-empty blocks). The number of sound asynchronizations depends roughly, on 1150 how many of these code blocks are racing with the method body. These asynchronizations contain 1151 awaits that are at a non-zero distance from the matching call (non-zero values in column await) 1152 and for many Github programs, this distance is bigger than in the original program (non-zero 1153 values in column cover)<sup>8</sup>. This shows that we are able to increase the distances between awaits 1154 and their matching calls for those programs. On average the distance between awaits and their 1155 matching calls in optimal asynchronizations for non synthetic benchmarks is 1.27 statements. 1156

With few exceptions, each program admits multiple sound asynchronizations (values in column async bigger than one), which makes the focus on the delay complexity relevant. Also, this leaves the possibility of making a choice based on other criteria, e.g., performance metrics. While asynchronizations are computed statically, their performance can be derived only dynamically (executing them). In general, we are not aware of any syntactic criteria that can guide towards computing a best solution w.r.t. performance in practice. These results show that our techniques have the potential of becoming the basis of a refactoring tool allowing programmers to improve their usage of the async/await primitives. The artifacts are available in an anonymous GitHub repository [Experiments 2021].

### 10 RELATED WORK

There are many works on synthesizing or repairing concurrent programs in the standard multithreading model, e.g., automatic parallelization in compilers [Bacon et al. 1994; Blume et al. 1996;
Han and Tseng 2001], or synchronization synthesis [Bloem et al. 2014; Cerný et al. 2015, 2013, 2014;
Clarke and Emerson 2008; Gupta et al. 2015; Manna and Wolper 1984; Vechev et al. 2009, 2010].
Our paper focuses on the use of the async/await primitives which poses specific challenges that
are not covered in these works. For instance, synthesizing lock placements does not admit unique
optimal solutions w.r.t. a syntactic order as for async/await.

1157

1158

1159

1160

1161

1162

1163

1164

1165

1:24

<sup>&</sup>lt;sup>1175</sup> <sup>8</sup>The synthetic examples are weakest asynchronizations to start with.

<sup>1176</sup> 

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

Program Refactoring. A number of program refactoring tools have been proposed for convert-1177 ing C# programs using explicit callbacks into async/await programs [Okur et al. 2014], Android 1178 1179 programs using AsyncTask into programs that use IntentService [Lin et al. 2015], or sequential applications into parallel applications using concurrent libraries for Java [Dig et al. 2009]. The 1180 C# related tool [Okur et al. 2014], which is the closest to our work, makes it possible to repair 1181 misusage of async/await that might result in deadlocks. Their repairing mechanism is based on 1182 forcing the continuations after the first await to run on background threads. This tool cannot 1183 modify procedure calls to be asynchronous as in our work. Compared to all these works, we give 1184 an algorithmic framework with precise specifications and complexity analysis. 1185

Data Race Detection. There are many works that study dynamic data race detection using 1186 happens-before and lock-set analysis, or timing-based detection, e.g., [Flanagan and Freund 2009; 1187 Kini et al. 2017; Li et al. 2019; Raman et al. 2010; Smaragdakis et al. 2012]. [Raman et al. 2010] 1188 proposes a dynamic data race detector for async-finish task-parallel programs by adapting the 1189 algorithm proposed in [Feng and Leiserson 1997] that computes abstract summaries of parallel 1190 tasks. [Li et al. 2019] presents a testing technique for finding data races in C# and F# programs, 1191 based on inserting timing delays in unsafe methods (e.g., methods that access memory without 1192 locking), and a monitor for finding data races. These methods could be used to approximate our 1193 reduction from data race checking to reachability in sequential programs. 1194

A number of works [Blackshear et al. 2018; Engler and Ashcraft 2003; Liu and Huang 2018] 1195 propose static analyses for finding data races. [Blackshear et al. 2018] designs a compositional data 1196 race detector for multi-threaded Java programs, based on an inter-procedural analysis assuming 1197 that any two public methods can execute in parallel. Similar to [Santhiar and Kanade 2017], they 1198 precompute method summaries in order to extract potential racy accesses. These approaches are 1199 similar to the analysis we present in Section 7, but they concern a different programming model. 1200 Analyzing Asynchronous Programs. There exist several works that propose program analyses 1201 for various classes of asynchronous programs. [Bouajjani and Emmi 2012; Ganty and Majumdar 1202 2012] give complexity results for the reachability problem, and [Santhiar and Kanade 2017] proposes 1203 a static analysis for deadlock detection in C# programs that use both asynchronous and synchronous 1204 wait primitives. This work relies on the static analysis introduced in [Madhavan et al. 2012] for 1205 computing method summaries in terms of points-to relations. [Bouajjani et al. 2017] investigates 1206 the problem of checking whether Java UI asynchronous programs have the same set of behaviors 1207 as sequential programs where roughly, asynchronous tasks are executed synchronously. 1208

### 11 CONCLUSION

1211 We have proposed a framework for refactoring sequential programs to equivalent asynchronous 1212 programs that rely on the async/await primitives. We have determined precise complexity bounds 1213 for the problem of computing a sound asynchronization that maximizes the distance between 1214 asyncs and awaits, which in theory, increases the level of parallelism, and the problem of com-1215 puting all sound asynchronizations. The latter problem is useful in a context where performance 1216 measures cannot be derived statically, which is usually the case, and makes it possible to compute 1217 a sound asynchronization that maximizes performance by separating concerns (enumerate sound 1218 asynchronizations and evaluate performance separately). We have also investigated the related 1219 problem of synthesizing sound multi-threaded refactorings where every method call is executed by 1220 a different thread, showing that our techniques extend quite easily, which witnesses the "robustness" 1221 of our framework. On the practical side, we have introduced an approximated synthesis procedure 1222 based on data-flow analysis that we implemented and evaluated on a benchmark of non-trivial C# 1223 programs extracted from open-source repositories. 1224

1225

Sidi Mohamed Beillahi, Ahmed Bouajjani, Constantin Enea, and Shuvendu Lahiri

The asynchronous programs rely exclusively on async/await and are deadlock-free by definition.
 Deadlocks can occur in a mix of async/await with "explicit" multi-threading that includes blocking
 wait primitives. Our paper deals with these two paradigms separately, but extending our approach
 for such programs is an interesting direction for future work.

### 1231 REFERENCES

- David F. Bacon, Susan L. Graham, and Oliver J. Sharp. 1994. Compiler Transformations for High-Performance Computing.
   ACM Comput. Surv. 26, 4 (1994), 345–420. https://doi.org/10.1145/197405.197406
- Gavin M. Bierman, Claudio V. Russo, Geoffrey Mainland, Erik Meijer, and Mads Torgersen. 2012. Pause 'n' Play: Formalizing
   Asynchronous C#. In ECOOP 2012 Object-Oriented Programming 26th European Conference, Beijing, China, June
   11-16, 2012. Proceedings (Lecture Notes in Computer Science), James Noble (Ed.), Vol. 7313. Springer, 233–257. https:
   //doi.org/10.1007/978-3-642-31057-7\_12
- Sam Blackshear, Nikos Gorogiannis, Peter W. O'Hearn, and Ilya Sergey. 2018. RacerD: compositional static race detection.
   *Proc. ACM Program. Lang.* 2, OOPSLA (2018), 144:1–144:28. https://doi.org/10.1145/3276514
- Roderick Bloem, Georg Hofferek, Bettina Könighofer, Robert Könighofer, Simon Ausserlechner, and Raphael Spork. 2014.
   Synthesis of synchronization using uninterpreted functions. In *Formal Methods in Computer-Aided Design, FMCAD 2014, Lausanne, Switzerland, October 21-24, 2014*. IEEE, 35–42. https://doi.org/10.1109/FMCAD.2014.6987593
- William Blume, Ramon Doallo, Rudolf Eigenmann, John Grout, Jay Hoeflinger, Thomas Lawrence, Jaejin Lee, David A.
   Padua, Yunheung Paek, William M. Pottenger, Lawrence Rauchwerger, and Peng Tu. 1996. Parallel Programming with
   Polaris. *IEEE Computer* 29, 12 (1996), 87–81. https://doi.org/10.1109/2.546612
- Ahmed Bouajjani and Michael Emmi. 2012. Analysis of recursively parallel programs. In *Proceedings of the 39th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2012, Philadelphia, Pennsylvania, USA, January 22-28, 2012*, John Field and Michael Hicks (Eds.). ACM, 203–214. https://doi.org/10.1145/2103656.2103681
- Ahmed Bouajjani, Michael Emmi, Constantin Enea, Burcu Kulahcioglu Ozkan, and Serdar Tasiran. 2017. Verifying Robustness
   of Event-Driven Asynchronous Programs Against Concurrency. In Programming Languages and Systems 26th European
   Symposium on Programming, ESOP 2017, Held as Part of the European Joint Conferences on Theory and Practice of Software,
   ETAPS 2017, Uppsala, Sweden, April 22-29, 2017, Proceedings (Lecture Notes in Computer Science), Hongseok Yang (Ed.),
   Vol. 10201. Springer, 170–200. https://doi.org/10.1007/978-3-662-54434-1\_7
- Pavol Cerný, Edmund M. Clarke, Thomas A. Henzinger, Arjun Radhakrishna, Leonid Ryzhyk, Roopsha Samanta, and Thorsten Tarrach. 2015. From Non-preemptive to Preemptive Scheduling Using Synchronization Synthesis. In Computer Aided Verification - 27th International Conference, CAV 2015, San Francisco, CA, USA, July 18-24, 2015, Proceedings, Part II (Lecture Notes in Computer Science), Daniel Kroening and Corina S. Pasareanu (Eds.), Vol. 9207. Springer, 180–197. https://doi.org/10.1007/978-3-319-21668-3\_11
- Pavol Cerný, Thomas A. Henzinger, Arjun Radhakrishna, Leonid Ryzhyk, and Thorsten Tarrach. 2013. Efficient Synthesis for Concurrency by Semantics-Preserving Transformations. In *Computer Aided Verification - 25th International Conference, CAV 2013, Saint Petersburg, Russia, July 13-19, 2013. Proceedings (Lecture Notes in Computer Science)*, Natasha Sharygina and Helmut Veith (Eds.), Vol. 8044. Springer, 951–967. https://doi.org/10.1007/978-3-642-39799-8\_68
- Pavol Cerný, Thomas A. Henzinger, Arjun Radhakrishna, Leonid Ryzhyk, and Thorsten Tarrach. 2014. Regression-Free
   Synthesis for Concurrency. In Computer Aided Verification 26th International Conference, CAV 2014, Held as Part of the
   Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 18-22, 2014. Proceedings (Lecture Notes in Computer Science),
   Armin Biere and Roderick Bloem (Eds.), Vol. 8559. Springer, 568–584. https://doi.org/10.1007/978-3-319-08867-9\_38
- Edmund M. Clarke and E. Allen Emerson. 2008. Design and Synthesis of Synchronization Skeletons Using Branching Time
   Temporal Logic. In 25 Years of Model Checking History, Achievements, Perspectives (Lecture Notes in Computer Science),
   Orna Grumberg and Helmut Veith (Eds.), Vol. 5000. Springer, 196–215. https://doi.org/10.1007/978-3-540-69850-0\_12
- Danny Dig, John Marrero, and Michael D. Ernst. 2009. Refactoring sequential Java code for concurrency via concurrent libraries. In 31st International Conference on Software Engineering, ICSE 2009, May 16-24, 2009, Vancouver, Canada, Proceedings. IEEE, 397–407. https://doi.org/10.1109/ICSE.2009.5070539
- Dawson R. Engler and Ken Ashcraft. 2003. RacerX: effective, static detection of race conditions and deadlocks. In *Proceedings* of the 19th ACM Symposium on Operating Systems Principles 2003, SOSP 2003, Bolton Landing, NY, USA, October 19-22,
   2003, Michael L. Scott and Larry L. Peterson (Eds.). ACM, 237–252. https://doi.org/10.1145/945445.945468
- 1269 Experiments. 2021. https://github.com/asynchronizations/artifact
- Mingdong Feng and Charles E. Leiserson. 1997. Efficient Detection of Determinacy Races in Cilk Programs. In Proceedings of the 9th Annual ACM Symposium on Parallel Algorithms and Architectures, SPAA '97, Newport, RI, USA, June 23-25, 1997, Charles E. Leiserson and David E. Culler (Eds.). ACM, 1–11. https://doi.org/10.1145/258492.258493
- Cormac Flanagan and Stephen N. Freund. 2009. FastTrack: efficient and precise dynamic race detection. In Proceedings of
   the 2009 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2009, Dublin, Ireland,

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

1:26

<sup>1274</sup> 

1275 June 15-21, 2009, Michael Hind and Amer Diwan (Eds.). ACM, 121–133. https://doi.org/10.1145/1542476.1542490

- 1276Pierre Ganty and Rupak Majumdar. 2012. Algorithmic verification of asynchronous programs. ACM Trans. Program. Lang.1277Syst. 34, 1 (2012), 6:1–6:48. https://doi.org/10.1145/2160910.2160915
- Patrice Godefroid and Mihalis Yannakakis. 2013. Analysis of Boolean Programs. In *Tools and Algorithms for the Construction* and Analysis of Systems - 19th International Conference, TACAS 2013, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2013, Rome, Italy, March 16-24, 2013. Proceedings (Lecture Notes in Computer Science), Nir Piterman and Scott A. Smolka (Eds.), Vol. 7795. Springer, 214–229. https://doi.org/10.1007/978-3-642-36742-7\_16
- Ashutosh Gupta, Thomas A. Henzinger, Arjun Radhakrishna, Roopsha Samanta, and Thorsten Tarrach. 2015. Succinct
   Representation of Concurrent Trace Sets. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2015, Mumbai, India, January 15-17, 2015, Sriram K. Rajamani and David* Walker (Eds.). ACM, 433–444. https://doi.org/10.1145/2676726.2677008
- Hwansoo Han and Chau-Wen Tseng. 2001. A Comparison of Parallelization Techniques for Irregular Reductions. In
   *Proceedings of the 15th International Parallel & Distributed Processing Symposium (IPDPS-01), San Francisco, CA, USA, April* 23-27, 2001. IEEE Computer Society, 27. https://doi.org/10.1109/IPDPS.2001.924963
- Dileep Kini, Umang Mathur, and Mahesh Viswanathan. 2017. Dynamic race prediction in linear time. In *Proceedings of the* 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017, Barcelona, Spain, June 18-23, 2017, Albert Cohen and Martin T. Vechev (Eds.). ACM, 157–170. https://doi.org/10.1145/3062341.3062374
- Guangpu Li, Shan Lu, Madanlal Musuvathi, Suman Nath, and Rohan Padhye. 2019. Efficient scalable thread-safety-violation
   detection: finding thousands of concurrency bugs during testing. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP 2019, Huntsville, ON, Canada, October 27-30, 2019,* Tim Brecht and Carey Williamson (Eds.). ACM,
   162–180. https://doi.org/10.1145/3341301.3359638
- Yu Lin, Semih Okur, and Danny Dig. 2015. Study and Refactoring of Android Asynchronous Programming (T). In 30th IEEE/ACM International Conference on Automated Software Engineering, ASE 2015, Lincoln, NE, USA, November 9-13, 2015, Myra B. Cohen, Lars Grunske, and Michael Whalen (Eds.). IEEE Computer Society, 224–235. https://doi.org/10.1109/ ASE.2015.50
- Bozhen Liu and Jeff Huang. 2018. D4: fast concurrency debugging with parallel differential analysis. In *Proceedings of the* 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018, Philadelphia, PA, USA, June 18-22, 2018, Jeffrey S. Foster and Dan Grossman (Eds.). ACM, 359–373. https://doi.org/10.1145/3192366.3192390
- Ravichandhran Madhavan, G. Ramalingam, and Kapil Vaswani. 2012. Modular Heap Analysis for Higher-Order Programs. In
   Static Analysis 19th International Symposium, SAS 2012, Deauville, France, September 11-13, 2012. Proceedings (Lecture Notes
   in Computer Science), Antoine Miné and David Schmidt (Eds.), Vol. 7460. Springer, 370–387. https://doi.org/10.1007/978 3-642-33125-1\_25
- Zohar Manna and Pierre Wolper. 1984. Synthesis of Communicating Processes from Temporal Logic Specifications. ACM Trans. Program. Lang. Syst. 6, 1 (1984), 68–93. https://doi.org/10.1145/357233.357237
- Semih Okur, David L. Hartveld, Danny Dig, and Arie van Deursen. 2014. A study and toolkit for asynchronous programming
   in c#. In 36th International Conference on Software Engineering, ICSE '14, Hyderabad, India May 31 June 07, 2014, Pankaj
   Jalote, Lionel C. Briand, and André van der Hoek (Eds.). ACM, 1117–1127. https://doi.org/10.1145/2568225.2568309
- Raghavan Raman, Jisheng Zhao, Vivek Sarkar, Martin T. Vechev, and Eran Yahav. 2010. Efficient Data Race Detection for
   Async-Finish Parallelism. In *Runtime Verification First International Conference, RV 2010, St. Julians, Malta, November 1-4, 2010. Proceedings (Lecture Notes in Computer Science),* Howard Barringer, Yliès Falcone, Bernd Finkbeiner, Klaus
   Havelund, Insup Lee, Gordon J. Pace, Grigore Rosu, Oleg Sokolsky, and Nikolai Tillmann (Eds.), Vol. 6418. Springer,

1309 368-383. https://doi.org/10.1007/978-3-642-16612-9\_28

- 1310 Roslyn. 2021. https://github.com/dotnet/roslyn
- Anirudh Santhiar and Aditya Kanade. 2017. Static deadlock detection for asynchronous C# programs. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017, Barcelona, Spain, June 18-23, 2017, Albert Cohen and Martin T. Vechev (Eds.). ACM, 292–305.* https://doi.org/10.1145/3062341.3062361
- Yannis Smaragdakis, Jacob Evans, Caitlin Sadowski, Jaeheon Yi, and Cormac Flanagan. 2012. Sound predictive race detection
   in polynomial time. In *Proceedings of the 39th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2012, Philadelphia, Pennsylvania, USA, January 22-28, 2012,* John Field and Michael Hicks (Eds.). ACM, 387–400.
   https://doi.org/10.1145/2103656.2103702
- Martin T. Vechev, Eran Yahav, and Greta Yorsh. 2009. Inferring Synchronization under Limited Observability. In *Tools and Algorithms for the Construction and Analysis of Systems, 15th International Conference, TACAS 2009, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2009, York, UK, March 22-29, 2009. Proceedings (Lecture Notes in Computer Science)*, Stefan Kowalewski and Anna Philippou (Eds.), Vol. 5505. Springer, 139–154. https://doi.org/10.1007/978-3-642-00768-2\_13
- 1321Martin T. Vechev, Eran Yahav, and Greta Yorsh. 2010. Abstraction-guided synthesis of synchronization. In Proceedings of the<br/>37th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2010, Madrid, Spain, January
- 1323

	1:28	Sidi Mohamed Beillahi, Ahmed Bouajjani, Constantin Enea, and Shuvendu Lahiri
1204	17 02 00	10 Manual V. Hamman arilda and Lura Dalahang (Eda.). ACM 207-228 https://doi.org/10.1145/170/200.170/228
1324	17-23, 20	<i>To</i> , Manuel V. Hermeneginuo and Jens Faisberg (Eds.). ACM, 527–558. https://doi.org/10.1143/1700299.1700558
1325		
1320		
1327		
1320		
1329		
1330		
1331		
1332		
1334		
1335		
1336		
1337		
1338		
1339		
1340		
1341		
1342		
1343		
1344		
1345		
1346		
1347		
1348		
1349		
1350		
1351		
1352		
1353		
1354		
1355		
1356		
1357		
1358		
1359		
1360		
1361		
1362		
1363		
1364		
1365		
1366		
1367		
1368		
1369		
1370		
13/1		
1372	Drog ACM	Program Long Vol 1 No CONE Article 1 Dublication Jates Language 2019
	TIUC. ACM	1 rogram. Lang., vol. 1, no. Cont., Article 1. 1 ubilication date: January 2010.