

Symbolic Reasoning for Automatic Signal Placement (Extended Version)

Kostas Ferles

The University of Texas at Austin
Austin, TX, USA
kferles@cs.utexas.edu

Isil Dillig

The University of Texas at Austin
Austin, TX, USA
isil@cs.utexas.edu

Jacob Van Geffen

The University of Texas at Austin
Austin, TX, USA
jsv@cs.utexas.edu

Yannis Smaragdakis

University of Athens
Athens, Greece
smaragd@di.uoa.gr

Abstract

Explicit signaling between threads is a perennial cause of bugs in concurrent programs. While there are several run-time techniques to *automatically* notify threads upon the availability of some shared resource, such techniques are not widely-adopted due to their run-time overhead. This paper proposes a new solution based on static analysis for automatically generating a performant *explicit-signal* program from its corresponding *implicit-signal* implementation. The key idea is to generate verification conditions that allow us to minimize the number of required signals and unnecessary context switches, while guaranteeing semantic equivalence between the source and target programs. We have implemented our method in a tool called EXPRESSO and evaluate it on challenging benchmarks from prior papers and open-source software. EXPRESSO-generated code significantly outperforms past automatic signaling mechanisms (avg. 1.56x speedup) and closely matches the performance of hand-optimized explicit-signal code.

Keywords Implicit signal monitors, concurrent programming, symbolic reasoning, verification conditions, abductive reasoning, monitor invariant

1 Introduction

A common challenge in concurrent programming is to coordinate access to shared resources and achieve correct synchronization between different threads. While there are many different language constructs that can be used to perform synchronization, a widely-established programming pattern is to encapsulate inter-thread coordination using *monitors* [23, 28, 39]. At a high level, a monitor encapsulates all shared state between threads and guarantees mutual exclusion. In addition, monitors perform synchronization between threads by blocking and unblocking them depending on the availability of some shared resource.

Broadly speaking, monitors can be classified into two categories, depending on the burden they impose on the system vs. the programmer [7]. In particular, *explicit-signal monitors*

typically employ *condition variables* to perform synchronization between threads and use an explicit “signal” construct to notify other threads when some shared resource becomes available. In contrast, implicit-signal (automatic) monitors provide a `waituntil(P)` construct such that any thread executing this statement blocks until predicate `P` becomes true. In implicit-signal monitors, there is no explicit signal construct, and it is the responsibility of the system to notify threads that are currently blocked on a predicate. To give the reader some intuition, Figure 1 shows the implementation of an implicit-signal monitor for the well-known readers-writers problem, and Figure 2 shows its corresponding implementation as an explicit-signal monitor.

As illustrated by the example from Figures 1 and 2, programming with implicit monitors is considerably easier because the programmer does not need to reason about when and which threads should be notified. In fact, it is well-known that many concurrency bugs are caused by erroneous signal placement in explicit-signal implementations [22, 31]. However, despite their easier programmability, implicit-signal monitors are not widely-used due to performance considerations. In particular, because the system needs to notify threads that are blocked on a predicate, run-time support for implicit-signal monitors may result in considerable overhead. For example, according to Buhr et al., automatic monitors can be 10-50 times slower than explicit signals [7]. Even though recent work by Hung and Garg proposes a more efficient implementation of automatic monitors [30], explicit-signal monitors still remain the de-facto synchronization mechanism in real-world concurrent programs.

In this paper, we propose a new solution —based on static analysis— to programming with implicit-signal monitors. Given the implementation of an implicit-signal monitor, our method automatically synthesizes an efficient and semantically equivalent explicit-signal implementation. We believe this approach has two advantages compared to prior run-time techniques: First, because our method does not require additional run-time book-keeping, it has the potential to be as efficient as a performant hand-written explicit-signal

implementation. Second, because the code generated by our system can be inspected and further refined by the programmer, it is more transparent compared to automatic-signaling systems that provide run-time instrumentation.

While it is straightforward to generate *any* semantically equivalent explicit-signal implementation of an automatic-signal monitor, a key consideration is the *efficiency* of the synthesized code. In particular, the synthesized code should *not* spuriously wake up threads that are blocked on a predicate that evaluates to false.¹ In practice, this means that the generated code should not notify threads blocked on a predicate P , if P is guaranteed to be false at the time of notification. Furthermore, whenever possible, the generated code should notify a single –rather than *all*– threads blocked on a predicate in order to avoid unnecessary context switches.

In addition to avoiding spurious wake-ups, another important efficiency consideration is to minimize the use of *conditional signals*, which notify other threads only if some condition evaluates to true. Because conditional signals require evaluating the truth value of (potentially complex) predicates at run-time, it is desirable to use unconditional signals whenever possible. In fact, while some run-time solutions, such as AutoSynch [30] avoid spurious wake-ups altogether, they may still incur significant overhead due to the frequent evaluation of predicates at run-time.

The solution that we adopt in this paper tries to minimize both the use of spurious wake-ups as well as conditional signals by performing precise static analysis of the monitor code. In particular, our method automatically generates Hoare triples, that, if valid, allow us to establish that a program fragment does *not* need to signal other threads waiting on a predicate. For program fragments where signaling may be necessary, our method generates additional Hoare triples whose validity allows us to minimize the use of conditional signals as well as *broadcast* operations that notify all threads.

In order to successfully discharge the generated Hoare triples, our method uses so-called *monitor invariants*, which are assertions that hold every time a thread enters or leaves the monitor. Our approach automatically infers these monitor invariants by combining abductive reasoning and predicate abstraction, allowing the synthesis of non-trivial invariants that involve disjunctions. Monitor invariants allow us to discharge verification conditions that could not be proven otherwise (e.g., by strengthening the precondition of the generated Hoare triples) and are therefore crucial for generating efficient explicit-signaling code.

We have implemented our proposed ideas in a tool called EXPRESSO and evaluate the efficiency of the code generated by EXPRESSO by comparing it against manually written explicit-signal monitors as well as the state-of-the-art

```

1  class RWLock {
2      unsigned int readers = 0;
3      boolean writerIn = false;
4
5      atomic void enterReader() {
6          waituntil(!writerIn);
7          readers++;
8      }
9      atomic void exitReader() {
10         if(readers > 0) readers--;
11     }
12     atomic void enterWriter() {
13         waituntil(readers == 0 && !writerIn);
14         writerIn = true;
15     }
16     atomic void exitWriter() {
17         writerIn = false;
18     } }

```

Figure 1. Implicit-signal monitor for readers-writers lock.

AutoSynch tool that provides run-time support for implicit-signal monitors. Our evaluation shows that the performance of the code synthesized by EXPRESSO is an average of 1.56x faster than AutoSynch and comparable to that of hand-written code.

In all, this paper makes the following key contributions:

- We propose a novel technique, based on static analysis, for generating efficient explicit-signal implementations of implicit-signal monitors.
- We show how the automatic signal placement problem can be reduced to proving the validity of certain kinds of Hoare triples in concurrent programs.
- We introduce the notion of *monitor invariants* and show how to automatically infer them using abductive reasoning and monomial predicate abstraction.
- We implement the proposed techniques in a tool called EXPRESSO and evaluate it by comparing against AutoSynch, a state-of-the-art runtime system for implicit-signal monitors, as well as hand-written code.

2 Overview of Technique

In this section, we give a high-level overview of our approach with the aid of the reference Readers-Writers example, shown in Figure 1. In particular, we explain the reasoning performed by EXPRESSO to automatically generate the code shown in Figure 2 by analyzing the implicit-signal monitor of Figure 1.

EXPRESSO starts its analysis by inferring a *monitor invariant*, which is an assertion that holds every time a thread enters or exits the monitor. For the code in Figure 1, EXPRESSO successfully infers the invariant $readers \geq 0$. Then, EXPRESSO uses this invariant to determine for each conditional critical section in Figure 1 (a) if signaling is necessary, (b) whether to signal or broadcast, and (c) whether to do so conditionally or unconditionally.

EnterReader. Consider a reader thread t_r executing the method `enterReader`. To generate explicit signaling code,

¹While a thread that is spuriously woken up will “go back to sleep”, this introduces significant overhead due to an unnecessary context-switch.

we need to determine whether t_r needs to notify any writer threads blocked on predicate $P_w = (readers = 0 \wedge \neg writerIn)$ at line 13. Towards this goal, we ask the following question: “Assuming that a writer thread t_w is blocked on P_w , is it possible that P_w becomes true after t_r executes the code in `enterReader`?”. If the answer to this question is “no”, we have established that t_r does not need to signal. Thus, to prove that no signals are necessary, EXPRESSO generates and checks the validity of the following Hoare triple:

$$\{readers \geq 0 \wedge \neg writerIn \wedge \neg P_w\} \text{ readers++ } \{\neg P_w\}$$

Here, the precondition states that (a) the monitor invariant holds when t_r enters the monitor, (b) $\neg writerIn$ must hold if t_r executes `readers++`, and (c) P_w is false, meaning that some writer thread may be blocked at line 13. The post-condition says that P_w continues to stay false after t_r exits the monitor. Since this Hoare triple is indeed valid, EXPRESSO establishes that no signaling is necessary. Observe that dropping the conjunct $readers \geq 0$ from the precondition would result in a Hoare triple that is *not* valid; thus, the monitor invariant is crucial for avoiding the signal operation in this example.

ExitReader. For the `exitReader` method, EXPRESSO needs to determine whether we should signal any reader threads blocked at line 6 or writer threads blocked at line 13. Using similar reasoning as in `enterReader`, it is easy to establish that we do not need to signal reader threads because `readers--` does not affect the truth value of the predicate `writerIn`. Now, to determine the necessity of signaling writer threads, EXPRESSO generates the following Hoare triple:

$$\{readers \geq 0 \wedge \neg P_w\} \text{ if}(readers > 0) \text{ readers-- } \{\neg P_w\}$$

Since this Hoare triple is not valid, signalling is necessary.

Next, EXPRESSO tries to determine whether it suffices to notify a *single* writer thread or we need to notify all writers blocked at line 13. To answer this question, we ask “Is it possible that P_w stays true after some writer thread t_w executes `enterWriter`?”. If not, we have proven that it is unnecessary (and wasteful) to wake up multiple threads, since P_w becomes false after the first writer thread executes. Thus, EXPRESSO generates and checks the following Hoare triple:

$$\{readers \geq 0 \wedge P_w\} \text{ writerIn} = \text{true } \{\neg P_w\}$$

Since this triple is valid, EXPRESSO has determined that broadcasting is *not* necessary.

Finally, EXPRESSO checks whether it can signal unconditionally, meaning that P_w is guaranteed to hold after the reader thread t_r exits the monitor. Towards this goal, we perform the following check:

$$\{readers \geq 0 \wedge \neg P_w\} \text{ if}(readers > 0) \text{ readers-- } \{P_w\}$$

This Hoare triple is not valid, so EXPRESSO signals conditionally in order to avoid a spurious wake-up.

```

1  class RWLock {
2      unsigned int readers = 0;
3      boolean writerIn = false;
4      Lock l = new ReentrantLock();
5      Condition readers = l.newCondition(),
6              writers = l.newCondition();
7      void enterReader() {
8          l.lock();
9          while(writerIn) readers.await();
10         readers++;
11         l.unlock();
12     }
13     void exitReader() {
14         l.lock();
15         if (readers > 0) readers--;
16         if (readers == 0) writers.signal();
17         l.unlock();
18     }
19     void enterWriter() {
20         l.lock();
21         while(readers != 0 || writerIn) writers.await();
22         writerIn = true;
23         l.unlock();
24     }
25     void exitWriter() {
26         l.lock();
27         writerIn = false;
28         if (readers == 0) writers.signal();
29         readers.signalAll();
30         l.unlock();
31     } }

```

Figure 2. Explicit-signal monitor for readers-writers lock.

EnterWriter. Using similar reasoning as in `enterReader`, EXPRESSO can establish that `enterWriter` does not need to signal any readers because the following Hoare triple is valid:

$$\{readers \geq 0 \wedge P_w \wedge writerIn\} \text{ writerIn} = \text{true } \{writerIn\}$$

ExitWriter. For `exitWriter`, EXPRESSO establishes that it is necessary to notify both reader and writer threads. Using similar reasoning as in `exitReader`, we can prove that broadcasting for all *writer* threads is not necessary, however, we need to perform conditional signaling to avoid spurious wake-ups. For *reader* threads, EXPRESSO determines that broadcasting is necessary since $\neg writerIn$ continues to hold after executing the statement `readers++`. Furthermore, since the Hoare triple

$$\{readers \geq 0 \wedge writerIn\} \text{ writerIn} = \text{false } \{\neg writerIn\}$$

is valid, EXPRESSO can establish that $\neg writerIn$ must be true after the writer thread exits. Thus, EXPRESSO instruments the code to signal reader threads unconditionally.

Summary. For this example, the code generated by EXPRESSO is precisely the same one as the human-written explicit-signal implementation shown in Figure 2. Observe that EXPRESSO can prove the gratuitousness of broadcasts, and it can also establish that `enterReader` and `enterWriter` do not need to signal. Finally, note that some of the Hoare triples generated by EXPRESSO could not be established without the useful monitor invariant $readers \geq 0$.

3 Source and Target Languages

In this section, we present some preliminary concepts related to concurrent programming and describe the source and target languages that can be used to implement implicit- and explicit-signal monitors respectively. The goal of the two languages presented here is to provide a unified theoretical framework suitable for automatic reasoning. In Section 6, we discuss how the target language can be instantiated in a concrete monitor implementation.

3.1 Preliminaries

In this paper, we consider a shared-memory concurrency model in which all accesses to shared resources occur inside a *monitor*. In other words, all variables accessed outside the monitor are assumed to be thread-local. We represent threads using integer identifiers drawn from the set $T \subseteq \mathbb{N}$. Because we do not impose any restrictions on the number of threads that can execute monitor code, our approach is applicable to parametrized concurrent programs.

We partition program variables used in the monitor into two disjoint sets, namely L and G , representing thread-local and shared (global) variables respectively. As stated by the definition below, the state σ of a monitor identifies the values of program variables for each thread.

Definition 3.1. (Monitor state) A monitor state, $\sigma : T \times (L \cup G) \rightarrow \mathbb{N}$, is a mapping from (thread identifier, monitor variable) pairs to a value. We require monitor states to agree on the values of shared variables for all threads; i.e.,

$$\forall t_1, t_2 \in T, v \in G. \sigma(t_1, v) = \sigma(t_2, v)$$

3.2 Source Language

Because our approach transforms an automatic-signal monitor to an explicit-signal one, we first present the source language in which automatic-signal monitors are implemented.

The syntax of our source language is presented in Figure 3. Since our implementation targets Java programs, we consider implicit-signal monitors written in a simple object-oriented language with Java-like syntax. In particular, an automatic-signal monitor consists of a set of field declarations and a set of *atomic* methods – i.e., the body of a method m executes without interruption unless the thread blocks on some *waituntil* statement whose corresponding predicate evaluates to false. To simplify presentation, we will assume that *local variables of different methods have unique names*.

The body of each monitor method is a sequence of statements of the form *waituntil*(p){ s }, where p is a predicate and s is a statement (assignment, store, sequence, loop etc.). Observe that a statement s is a special case of a *waituntil* statement whose corresponding predicate is *true*. We refer to predicate p as the *guard* of the *waituntil* construct and to statement s as its *body* and sometimes write $w = (p, s)$ to denote a *waituntil* statement with guard p and body s . Given a monitor M , we use the notation $CCRs(M)$ to represent the set of all *waituntil* statements used in any method in M .

Monitor M	$::=$ monitor M { $(fld \mid m)^*$ }
Field fld	$::=$ $f : \tau$
Method m	$::=$ atomic $m(\vec{v})$ { w^* ; return v ; }
WUntil w	$::=$ waituntil(p){ s }
Statement s	$::=$ $skip \mid s_1; s_2 \mid v := e \mid v.f := e$ \mid if (p) { s_1 } else { s_2 } \mid while (p) { s }

Figure 3. Implicit-signal monitor language. Here, e and p denote expressions and predicates respectively.

While *waituntil* statements can only appear as top-level statements in our source language, we note that this design decision does not sacrifice expressiveness. For example, consider the code snippet if (c) *waituntil*(p), which is not supported by our source language. Observe that the check if(c) can be moved outside of the monitor if c is not on shared data. On the other hand, if c does involve shared data, the condition is either checked in a logically racy way² or the programmer knows that c cannot change while the thread is blocked on p . In either case, the program’s logic is preserved or enhanced if condition c is moved inside the *waituntil* statement.

In the rest of this paper, we assume the standard semantics of *waituntil*(p){ s } statements where a thread t atomically performs the following actions: It first evaluates the boolean predicate p . If p evaluates to *true*, t also executes s immediately after the evaluation of p . Otherwise, t is blocked until p becomes true.

Since the semantics of statements s are standard, we do not present them in detail and use the notation $\langle s, t, \sigma \rangle \Downarrow \sigma'$ to indicate the resulting monitor state σ' when thread t executes statement s under initial state σ . Given a monitor state σ , thread t , and predicate p , we write $(\sigma, t) \models p$ if p evaluates to *true* and $(\sigma, t) \not\models p$ if p evaluates to false.

Monitor traces. To define the semantics of monitors, we first introduce the notion of a *monitor trace*. A monitor trace τ is a sequence of *monitor events* where each event e is a triple (t, w, b) where t is a thread identifier, w is a *waituntil* statement, and b is a boolean indicating whether the guard of w evaluates to *true* or *false*. In particular, the event $(t, w, false)$ indicates that thread t was blocked on the guard of w , whereas the event $(t, w, true)$ indicates that t was able to execute w in its entirety. Given an event $e = (t, w, b)$, we write \bar{e} to denote the pair (t, w) .

We say that a monitor trace is *syntactically well-formed* if it (a) respects the relative ordering of statements within a method, (b) obeys the requirement that a thread cannot execute method m' before finishing the execution of method m , and (c) satisfies the invariant that a thread exits the monitor either by blocking on a predicate or by finishing the

²i.e., condition c could have changed while the thread is waiting on p

$$\begin{array}{l}
(1a) \quad \frac{e = (t, w, false) \quad \bar{e} \notin \mathcal{B} \quad (\sigma, t) \not\models Guard(w)}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma, \epsilon, \mathcal{B} \cup \{\bar{e}\}, \mathcal{N})} \\
(1b) \quad \frac{e = (t, w, false) \quad \bar{e} \in \mathcal{N} \quad (\sigma, t) \not\models Guard(w)}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma, \epsilon, \mathcal{B}, \mathcal{N} \setminus \{\bar{e}\})} \\
(2a) \quad \frac{e = (t, w, true) \quad \bar{e} \notin \mathcal{B} \quad (\sigma, t) \models Guard(w) \quad \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \quad \mathcal{N}' = \{(t, w) \mid (t, w) \in \mathcal{B}, (\sigma', t) \models Guard(w)\}}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \epsilon, \mathcal{B}, \mathcal{N} \cup \mathcal{N}')} \\
(2b) \quad \frac{e = (t, w, true) \quad \bar{e} = \min(\mathcal{N}) \quad (\sigma, t) \models Guard(w) \quad \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \quad \mathcal{N}' = \{(t, w) \mid (t, w) \in \mathcal{B}, (\sigma', t) \models Guard(w)\}}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \epsilon, \mathcal{B} \setminus \{\bar{e}\}, (\mathcal{N} \cup \mathcal{N}') \setminus \{\bar{e}\})} \\
(3) \quad \frac{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')}{(\sigma, e :: \tau, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \tau, \mathcal{B}', \mathcal{N}')}
\end{array}$$

Figure 4. Transition relation for implicit-signal monitor traces. Given an event $e = (t, w, b)$, \bar{e} denotes (t, w) . We assume there is a total order relation $<$ between events and \min picks the minimum one with respect to $<$.

execution of a method. A more formal definition of syntactic well-formedness is presented in Appendix A.

Example 3.2. Consider the following monitor M , where we elide the “atomic” keywords for brevity:

```

monitor M {
  ...
  m1() {waituntil(x>0) {...}; waituntil{y>0}{...}}
  m2() {waituntil(z>0) {...}; waituntil{w>0}{...}}
}

```

Let us refer to the j 'th `waituntil` statement in method i as w_{ij} . The trace $[(1, w_{12}, true), (1, w_{11}, true)]$ is not syntactically well-formed since the same thread cannot execute w_{12} before w_{11} (i.e., it violates requirement (a)). Similarly, the trace $[(1, w_{11}, false), (1, w_{21}, true)]$ is also not syntactically well-formed since the same thread cannot execute method $m2$ before finishing the execution of $m1$ (violates (b)). Finally, the following trace is also not syntactically well-formed:

$[(1, w_{11}, false), (2, w_{21}, true), (1, w_{11}, true), (1, w_{12}, true)]$

In particular, it violates requirement (c) since thread 2 exists the monitor without getting blocked or finishing the execution of $m2$. On the other hand, the following trace is syntactically well-formed:

$[(1, w_{11}, false), (2, w_{21}, true), (2, w_{22}, false), (1, w_{11}, true), (1, w_{12}, true), (2, w_{22}, true)]$

In this trace, thread 1 attempts to execute the body of w_{11} but is blocked (i.e., $x > 0$ evaluates to false). Then, thread 2 executes the first `waituntil` statement in method $m2$, but

$$\begin{array}{l}
(1a) \quad \frac{e = (t, w, false) \quad \bar{e} \notin \mathcal{B} \quad (\sigma, t) \not\models Guard(w)}{(\sigma, e, \mathcal{B}, \mathcal{N}) \Longrightarrow (\sigma, \epsilon, \mathcal{B} \cup \{\bar{e}\}, \mathcal{N})} \\
(1b) \quad \frac{e = (t, w, false) \quad \bar{e} \in \mathcal{N} \quad (\sigma, t) \not\models Guard(w)}{(\sigma, e, \mathcal{B}, \mathcal{N}) \Longrightarrow (\sigma, \epsilon, \mathcal{B}, \mathcal{N} \setminus \{\bar{e}\})} \\
(2a) \quad \frac{e = (t, w, true) \quad \bar{e} \notin \mathcal{B} \quad (\sigma, t) \models Guard(w) \quad \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \quad \mathcal{N}_1 = GetSignals(w, \sigma', \mathcal{B}) \quad \mathcal{N}_2 = GetBroadcasts(w, \sigma', \mathcal{B})}{(\sigma, e, \mathcal{B}, \mathcal{N}) \Longrightarrow (\sigma', \epsilon, \mathcal{B}, \mathcal{N} \cup \mathcal{N}_1 \cup \mathcal{N}_2)} \\
(2b) \quad \frac{e = (t, w, true) \quad \bar{e} = \min(\mathcal{N}) \quad (\sigma, t) \models Guard(w) \quad \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \quad \mathcal{N}_1 = GetSignals(w, \sigma', \mathcal{B}) \quad \mathcal{N}_2 = GetBroadcasts(w, \sigma', \mathcal{B})}{(\sigma, e, \mathcal{B}, \mathcal{N}) \Longrightarrow (\sigma', \epsilon, \mathcal{B} \setminus \{\bar{e}\}, (\mathcal{N} \cup \mathcal{N}_1 \cup \mathcal{N}_2) \setminus \{\bar{e}\})} \\
(3) \quad \frac{(\sigma, e, \mathcal{B}, \mathcal{N}) \Longrightarrow (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')}{(\sigma, e :: \tau, \mathcal{B}, \mathcal{N}) \Longrightarrow (\sigma', \tau, \mathcal{B}', \mathcal{N}')}
\end{array}$$

Figure 5. Transition relation for explicit-signal monitor traces. Given an event $e = (t, w, b)$, \bar{e} denotes (t, w) .

gets blocked on the second one. After thread 2 executes w_{21} , $x > 0$ becomes true, and thread 1 is able to finish executing method $m1$. Finally thread 2 finishes executing method $m2$.

Semantics. We now define the semantics of implicit-signal monitors in terms of the feasibility of well-formed monitor traces. Given a monitor M and a monitor state σ , we say that a trace τ is *feasible* under σ iff (a) it is syntactically well-formed and (b) $(\sigma, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma', \epsilon, _, _)$ where \longrightarrow^* denotes the reflexive transitive closure of the transition relation \longrightarrow defined in Figure 4.

Transition relations for implicit-signal monitors are described in Figure 4 using judgments of the form

$$(\sigma, \tau, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \tau', \mathcal{B}', \mathcal{N}')$$

where \mathcal{B} and \mathcal{N} describe blocked and notified threads respectively. In particular, $(t, w) \in \mathcal{B}$ indicates that thread t is currently blocked on the predicate of w . In contrast, $(t, w) \in \mathcal{N}$ indicates that thread t should be woken up to recheck the predicate of w . The meaning of the judgment $(\sigma, \tau, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \tau', \mathcal{B}', \mathcal{N}')$ is that executing the first event e in τ under σ, \mathcal{B} , and \mathcal{N} yields a new state σ' as well as a new set of blocked and notified threads \mathcal{B}' and \mathcal{N}' respectively. We now explain the transition relations from Figure 4 in more detail.

According to rules (1a) and (1b), an event e of the form $(t, w, false)$ is only feasible when $(\sigma, t) \not\models Guard(w)$ (i.e., the predicate of w evaluates to false). If $\bar{e} = (t, w)$ was not previously in the blocked thread set \mathcal{B} , rule (1a) adds \bar{e} to \mathcal{B} .

$$\begin{aligned}
\text{Events}(\mathcal{B}, p) &= \{(t, w) \mid (t, w) \in \mathcal{B} \wedge \text{Guard}(w) = p\} \\
\text{GetSignals}(w, \sigma, \mathcal{B}) &= \{(t', w') \mid (p, c) \in \text{Signals}(w) \wedge (t', w') = \min(\text{Events}(\mathcal{B}, p)) \wedge (c = \checkmark \vee (\sigma, t') \models p)\} \\
\text{GetBroadcasts}(w, \sigma, \mathcal{B}) &= \{(t', w') \mid (p, c) \in \text{Broadcasts}(w) \wedge (t', w') \in \text{Events}(\mathcal{B}, p) \wedge (c = \checkmark \vee (\sigma, t') \models p)\}
\end{aligned}$$

Figure 6. Auxiliary functions used in Figure 5

If \bar{e} was already in \mathcal{B} , then e is only feasible if \bar{e} was “notified” by the system (i.e., $\bar{e} \in \mathcal{N}$ in rule (1b)).

The next two rules (2a) and (2b) state that an event $e = (t, w, \text{true})$ is only feasible when $(\sigma, t) \models \text{Guard}(w)$ (i.e., the predicate of w evaluates to true). Both rules execute the body of w to obtain a new monitor state σ' . Now, since the execution of w may cause the predicates of blocked threads to become true, \mathcal{N}' contains all (t, w) pairs that were previously in \mathcal{B} and whose predicates evaluate to true under σ' .

3.3 Target Language

Our target language is very similar to the source language from Figure 3, except that the body of `waituntil` statements contain explicit signals. In particular, a `waituntil` construct in the target language looks as follows:

$$\text{waituntil}(p)\{s; \text{signal}(S_1); \text{broadcast}(S_2)\}$$

Here, S_1 and S_2 are sets of pairs (p, c) where p is a predicate and $c \in \{?, \checkmark\}$. The informal semantics of `signal` and `broadcast` are as follows: If (p, \checkmark) is in S_1 , then the system will notify (i.e., wake up) a *single* thread blocked on predicate p . In contrast, if $(p, \checkmark) \in S_2$, then the system notifies *all* threads blocked on p . On the other hand, if $(p, ?)$ is in S_1 (resp. S_2), then p will be evaluated at run-time, and, if p evaluates to true, then one thread (resp. all threads) blocked on p will be notified. Given a `waituntil` statement in the target language, we write $\text{Signals}(w)$ to indicate S_1 and $\text{Broadcasts}(w)$ to represent S_2 .

We also describe the formal semantics of explicit-signal monitors in terms of monitor traces, where the definitions of *trace*, *event*, and *well-formedness* remain the same as in Section 3.2. However, the concept of *feasibility* is defined with respect to a different transition relation \Longrightarrow , shown in Figure 5. In particular, we say that an explicit-signal monitor trace τ is *feasible* if (a) it is syntactically well-formed, and (b) $(\sigma, \tau, \emptyset, \emptyset) \Longrightarrow^* (\sigma', \epsilon, _, _)$ where \Longrightarrow^* denotes the reflexive transitive closure of the transition relation \Longrightarrow from Figure 5.

The transition relation \Longrightarrow is defined similarly as \rightarrow except for events of the form (t, w, true) . In contrast to implicit signal monitors which wake up all threads whose predicates have become true, explicit-signal monitors decide which threads to notify based on $\text{Signals}(w)$ and $\text{Broadcasts}(w)$. In particular, rules (2a) and (2b) use auxiliary functions GetSignals and GetBroadcasts (defined in Figure 6) to decide which threads to add to the notification set \mathcal{N} . If $(p, c) \in \text{Signals}(w)$, then we notify a *single* event $(t', w') \in \mathcal{B}$ such that the predicate of w' is p . If $c = ?$, we additionally check that p evaluates to true under σ' before adding (t', w') to the notification set. The function GetBroadcasts is defined

similarly except that it notifies all threads blocked on the specified predicate rather than a single one.

3.4 Equivalence

We are now ready to define what it means for an implicit-signal monitor M from the source language and an explicit-signal monitor M' from the target language to be *equivalent*. Towards this goal, we first define a normal form for traces:

Definition 3.3. (Normalization) Let τ be an implicit-signal monitor trace. We say that τ is *normalized* with respect to monitor state σ if we can derive $(\sigma, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma, \epsilon, _, _)$ without using rule (1b) from Figure 4 in the derivation.

Since rule (1b) corresponds to a spurious notification, a trace is normalized if threads are woken up only when their predicates evaluate to true. Observe that we can always find a normalized feasible trace for any feasible trace by changing the order in which threads are woken up.³

Definition 3.4. (Equivalence) Let M, M' be implicit- and explicit-signal monitors respectively. We say that M and M' are semantically equivalent, written $M \sim M'$, iff for all monitor states σ and all well-formed traces τ , the following two conditions are satisfied:

1. If $(\sigma, \tau, \emptyset, \emptyset) \Longrightarrow^* (\sigma', \epsilon, _, _)$, then it is also the case that $(\sigma, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma', \epsilon, _, _)$.
2. If $(\sigma, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma', \epsilon, _, _)$ and τ is normalized with respect to σ , then $(\sigma, \tau, \emptyset, \emptyset) \Longrightarrow^* (\sigma', \epsilon, _, _)$.

Here, the first condition states that any feasible trace of the explicit-signal monitor M' must also be a feasible trace of its implicit version M . However, in general, we cannot expect the converse of this statement to hold: Since the explicit-signal monitor may be more efficient than its implicit-signal counterpart, we cannot require that all feasible traces of M to be also feasible in the explicit-monitor case. Thus, the second condition states that any *normalized* feasible trace of M should also be feasible in M' .

4 Signal Placement Algorithm

In this section, we describe our algorithm for automatically transforming an implicit-signal monitor M in the source language to an explicit-signal monitor M' in the target language. Our algorithm ensures that M and M' are equivalent in the sense of Definition 3.4 and also tries to minimize the number of spurious wake-ups and conditional signals in M' . We

³Recall that our notion of feasibility does not require the sets \mathcal{B}, \mathcal{N} to be empty. In particular, a trace τ is feasible under σ if $(\sigma, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma', \epsilon, \mathcal{B}, \mathcal{N})$ for any \mathcal{B} and \mathcal{N} . Therefore, a notification that would have been eliminated by rule (1b) can just be ignored indefinitely, i.e., remain in the \mathcal{N} set without affecting other transitions.

Algorithm 1 Signal Placement Algorithm

```

1: function PLACESIGNALS( $M, I$ )
2:   input:  $M$ , an implicit signal monitor
3:   input:  $I$ , a monitor invariant
4:   output:  $M'$ , an explicit signal monitor
5:    $\Sigma \leftarrow [w \mapsto \emptyset \mid w \in \text{CCRs}(M)]$ 
6:   for  $(w, p) \in \text{CCRs}(M) \times \text{Guards}(M)$  do
7:     if  $\vdash \{I \wedge \text{Guard}(w) \wedge \neg p\} \text{Body}(w) \{-p\} :$ 
8:       continue;
9:     if  $\vdash \{I \wedge \text{Guard}(w) \wedge \neg p\} \text{Body}(w)\{p\} :$ 
10:       $\text{cond} \leftarrow \checkmark$ 
11:     else
12:       $\text{cond} \leftarrow ?$ 
13:     if  $\forall (p, s') \in \text{CCRs}(M). \vdash \{I \wedge p\} s' \{-p\} :$ 
14:       $\text{bcast} \leftarrow \text{false}$ 
15:     else
16:       $\text{bcast} \leftarrow \text{true}$ 
17:       $\Sigma(w) \leftarrow \Sigma(w) \cup \{(p, \text{cond}, \text{bcast})\}$ 
18:   return Instrument( $M, \Sigma$ )
  
```

start with a basic version of the algorithm and then describe extensions and improvements later in this section.

4.1 Basic Algorithm

Our basic signal placement algorithm is shown in Algorithm 1. The PLACESIGNALS algorithm takes as input an implicit-signal monitor M as well as a *monitor invariant* I , which is an assertion that holds every time a thread enters and exits the monitor. Since automated inference of monitor invariants is described in the next section, we will assume that an oracle provides them for the time being. In this section, we further assume that guards used in `waituntil` statements do not contain thread-local variables. Given such an implicit-signal monitor M and its invariant I , PLACESIGNALS returns an explicit-signal monitor M' such that $M \sim M'$.

The algorithm maintains a mapping from each *conditional critical region* (CCR) (i.e., `waituntil` statement) w in M to a set of notifications that should be performed after executing w and before exiting the monitor. The algorithm represents these notification as triples of the form $(p, \text{cond}, \text{bcast})$, where p is a predicate, $\text{cond} \in \{?, \checkmark\}$ indicates whether the notification is conditional or unconditional, and bcast is a boolean indicating whether it is necessary to notify all threads blocked on p as opposed to a single one. Once the algorithm computes this mapping Σ , it instruments the original implicit-signal monitor M as shown in Figure 7 to obtain an explicit-signal monitor M' .

The key part of the PLACESIGNALS algorithm is the loop in lines 6–17. For each conditional critical region w and predicate p used in the monitor, the algorithm first decides whether w may need to notify threads blocked on predicate p . This decision is made based on the provability of the

$$\begin{array}{c}
 w = \text{waituntil}(p')\{s\} \\
 S_1 = \{(p, c) \mid (p, c, \text{false}) \in \Sigma(w)\} \\
 S_2 = \{(p, c) \mid (p, c, \text{true}) \in \Sigma(w)\} \\
 s' = \text{signal}(S_1); \text{broadcast}(S_2) \\
 \hline
 \Sigma \vdash w \rightsquigarrow \text{waituntil}(p')\{s; s'\} \\
 \\
 \Sigma \vdash w_1 \rightsquigarrow w'_1 \quad \dots \quad \Sigma \vdash w_n \rightsquigarrow w'_n \\
 \hline
 \Sigma \vdash w_1; \dots; w_n \rightsquigarrow w'_1; \dots; w'_n
 \end{array}$$

Figure 7. Performing instrumentation

following Hoare triple:

$$\{I \wedge \text{Guard}(w) \wedge \neg p\} \text{Body}(w) \{-p\}$$

Essentially, this triple says that executing the body of w in a state in which $\neg p$ holds ensures that predicate p continues to remain false. Hence, any thread t blocked on p will remain blocked after executing w , so there is no need to notify t . Observe that the precondition of the Hoare triples also assumes $I \wedge \text{Guard}(w)$ because (a) $\text{Guard}(w)$ is a prerequisite for executing the body of w and, (b) by definition of monitor invariant, I must hold before executing the body of any CCR.

Next, lines 9–12 determine whether the notification should be conditional or not. Recall that a conditional notification for predicate p checks whether p evaluates to true before waking up threads blocked on p . While conditional notifications prevent spurious wake-ups, it is desirable to avoid evaluating p at run-time if p is guaranteed to hold after executing w . Thus, line 9 checks the validity of the following Hoare triple:

$$\{I \wedge \text{Guard}(w) \wedge \neg p\} \text{Body}(w) \{p\}$$

In other words, assuming we execute w in a state where a thread is blocked on p , the execution of $\text{Body}(w)$ results in a state where p is true. Thus, there is no need to evaluate p at run time before signaling threads blocked on p .

The last part of Algorithm 1 (lines 13–16) determines whether we should notify *all* threads blocked on predicate p . Suppose there are n threads $T = \{t_1, \dots, t_n\}$ blocked on p , and suppose that an arbitrary thread t_i gets unblocked. If executing t_i is guaranteed to result in a state where predicate p is false, then it is not necessary to notify any of the remaining threads $T \setminus \{t_i\}$. Thus, the algorithm checks the following Hoare triple for *all* CCRs w' with guard p :

$$\{I \wedge p\} \text{Body}(w') \{-p\}$$

If this Hoare triple holds for all CCRs with guard p , then it is safe to signal rather than broadcast.

Theorem 4.1. ⁴ *Let $\text{PlaceSignals}(M, I) = M'$. If I is a correct monitor invariant and guards of CCRs in M do not contain thread-local variables, then $M \sim M'$.*

⁴The proofs of all theorems are in the appendix.

4.2 Handling Thread-Local Variables

To simplify presentation, our algorithm from Section 4.1 assumes that guards of CCRs in the input monitor do not contain thread-local variables. However, if the input monitor M does not satisfy this assumption, the explicit-signal monitor M' generated by Algorithm 1 may not be equivalent to M . We illustrate the problem using the following example:

Example 4.2. Consider the following monitor:

```
monitor M {
  int y=0;
  m1(int x) {waituntil(x < y) {x = y+1;} }
  m2() { y = y+2; }
```

Suppose we have threads t_1, t_2, t_3 , where t_1, t_2 are blocked in $m1$, and t_3 is executing $m2$, after which the value of y becomes 2. Further, suppose that the value of the thread-local variable x is 0 for t_1 and 1 for t_2 . Since the predicate $x < y$ has become true for both t_1, t_2 and executing t_1 does not change the value of the predicate in t_2 (and vice versa), t_3 should notify both threads. Thus, the explicit-signal monitor should use broadcast instead of signal.

However, recall that Algorithm 1 determines whether $m2$ should broadcast or signal by checking the validity of $\{x < y\} x = y + 1 \{x \geq y\}$. Since this Hoare triple is valid, we would erroneously conclude that it is safe for $m2$ to notify a single thread instead of all threads.

As illustrated by this example, Algorithm 1 is unsound when guards contain thread-local variables. To remedy this situation, we need to rename thread-local variables when checking validity. In particular, recall that PLACE SIGNALS checks the validity of Hoare triples of the form $\{P_1 \wedge P_2\} S \{Q\}$ where P_1 is an assumption about the currently running thread, whereas P_2 and Q are assumptions/assertions about some other thread. Since S and P_1 may refer to thread-local variables that are also used in P_2 and Q , we need to rename thread-local variables and check the validity of the following modified Hoare triple:

$$\{P_1 \wedge P_2[V'/V]\} S \{Q[V'/V]\}$$

where $V = \text{Locals}(P_2) \cup \text{Locals}(Q)$ and V' denotes a fresh set of variables not used elsewhere in P_1, P_2, S , and Q .⁵

4.3 Improvement over the Basic Algorithm

In this section, we consider an improvement over Algorithm 1 that aims to further reduce the number of broadcasts in the synthesized explicit-signal monitor. Recall that Algorithm 1 determines whether a CCR should notify one vs. all threads blocked on predicate p by checking the validity of the following Hoare triple for all CCRs w with guard p :

$$\{I \wedge p\} \text{Body}(w) \{\neg p\} \quad (1)$$

In some cases, it is possible to further strengthen the precondition of this Hoare triple. In particular, suppose that the signaling CCR is w' with body s' and guard p' and suppose

⁵Another subtlety with local variables is how to signal conditionally. We discuss evaluation of predicates with local variables in Section 6.

that ϕ is guaranteed to hold after executing s' . In general, we cannot assume ϕ in the pre-condition of Equation 1 because other threads may have invalidated ϕ before the notified thread has a chance to execute. However, if s commutes with the body of every other CCR in the monitor, then the monitor state after executing s for *any* interleaving is equivalent to one in which we execute s immediately after s' . In this case, we can safely assume that s executes immediately after s' since the resulting states are equivalent. This insight allows us to strengthen the precondition of Equation 1 by using the post-condition ϕ of the signaling thread.

To make this discussion more precise, let us define a predicate $\text{Comm}(w, M)$ as follows:

$$\begin{aligned} \text{Comm}(w, M) &\Leftrightarrow (\forall w' \in \text{CCRs}(M) \setminus \{w\}. \\ &\text{Body}(w'); \text{Body}(w) \equiv \text{Body}(w); \text{Body}(w')) \end{aligned}$$

Essentially, this predicate is true if the body of w commutes with every other CCR in the monitor. Now, using this definition, we can state a weaker sufficient condition for CCR w to notify one –rather than all– threads blocked on predicate p . In particular, we can change the condition at line 13 of Algorithm 1 to the following weaker one:

$$\begin{aligned} \forall w' = (p, s') \in \text{CCRs}(M). & \left(\vdash \{I \wedge p\} s' \{\neg p\} \vee \right. \\ & \left. (\text{Comm}(w', M) \wedge \vdash \{I \wedge \text{Guard}(w) \wedge \neg p\} \text{Body}(w); s' \{\neg p\}) \right) \quad (2) \end{aligned}$$

The first line of Equation 2 corresponds to the same check we perform at line 13 in Algorithm 1 to determine whether it is safe to signal rather than broadcast. However, if this condition does not hold, we may still be able to prove that broadcasting is unnecessary as long as s' commutes with every other CCR in the monitor and we can prove that p is falsified after executing $\text{Body}(w); s'$.

The correctness of Equation 2 follows from the following theorem (and the proof of Theorem 4.1):

Theorem 4.3. *Let $\tau = \tau_0 e$ be a monitor trace and let $\tau' = e \tau_0$ where $e = (t, w, b)$. If $(\sigma, \tau, \mathcal{B}, \mathcal{N}) \xrightarrow{*} (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')$ and $\text{Comm}(w, M)$, then we have $(\sigma, \tau', \mathcal{B}, \mathcal{N}) \xrightarrow{*} (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')$.*

Remark. Our discussion in this section assumes non-preemptive signal semantics [3] where a signaled thread is not guaranteed to consume the signal immediately. However, if we assume *preemptive* signal semantics, we can perform this optimization more liberally by only checking whether predicate p is invalidated by the sequential composition of the segment that produces the signal and $\text{Body}(w)$.

5 Inference of Monitor Invariants

Our signal placement algorithm from Section 4 relies on a monitor invariant I that holds at the entry and exit of every CCR. In this section, we describe our method for automatically inferring useful monitor invariants.

Our inference algorithm is property-directed in that it only infers invariants that are useful for proving the Hoare triples generated by the signal placement algorithm. Specifically,

Algorithm 2 Monitor Invariant Inference

```
1: function INFERMONITORINV( $M, \Theta$ )
2:   input:  $M$ , an implicit signal monitor
3:   input:  $\Theta$ , set of Hoare triples of the form  $\{P\} s \{Q\}$ 
4:   output:  $I$ , a monitor invariant
5:    $\Phi \leftarrow \emptyset$ 
6:   for  $\{P\} s \{Q\} \in \Theta$  do
7:      $\Phi \leftarrow \Phi \cup \text{abduce}(P, \text{wp}(s, Q))$ 
8:   do
9:      $\text{numPreds} \leftarrow |\Phi|$ 
10:    for  $\psi \in \Phi$  do
11:      if  $\not\models \{true\} \text{Ctr}(M) \{\psi\}$  :
12:         $\Phi \leftarrow \Phi \setminus \{\psi\}$ 
13:      continue;
14:       $I \leftarrow \bigwedge_{\psi_i \in \Phi} \psi_i$ 
15:      if  $\exists w \in \text{CCRs}(M). \not\models \{I \wedge \text{Guard}(w)\} \text{Body}(w) \{\psi\}$  :
16:         $\Phi \leftarrow \Phi \setminus \{\psi\}$ 
17:    while  $\text{numPreds} \neq |\Phi|$ 
18:  return  $I$ 
```

our inference engine uses *abductive reasoning* [16] to automatically infer predicates that are useful for proving a given set of Hoare triples. Given a universe of predicates Φ generated using abduction, it then infers the *strongest* conjunctive monitor invariant over predicates in Φ . Therefore, our invariant inference engine can be viewed as marrying the power of abductive reasoning with predicate abstraction [21, 37]. The advantage of this approach is two-fold: First, rather than relying on a hard-coded universe of predicate templates, our algorithm infers useful predicates automatically using abduction. Second, because the predicates inferred using abduction can involve disjunctions, the monitor invariants synthesized by our algorithm are not restricted to pure conjunctions.

With this intuition in mind, we now explain our INFERMONITORINV procedure from Algorithm 2 in more detail. This procedure takes two inputs, namely, an implicit-signal monitor M and a set Θ of Hoare triples of the form $\{P\} s \{Q\}$. Note that Θ simply corresponds to the set of Hoare triples generated by Algorithm 1, but with I set to *true*. The return value of INFERMONITORINV is a formula I representing a valid monitor invariant of M .

Conceptually, the INFERMONITORINV procedure operates in two phases. The first phase (lines 5–7) generates a universe Φ of candidate predicates, and the second phase (lines 8–17) performs fixed-point computation to infer the strongest conjunctive monitor invariant I over predicates in Φ .

In the first phase of the algorithm, we iterate over all Hoare triples $\{P\} s \{Q\}$ in Θ and look for a strengthening ψ of the precondition such that the Hoare triple $\{P \wedge \psi\} s \{Q\}$ becomes valid. Because the correctness of the Hoare triple $\{P\} s \{Q\}$ boils down to checking the validity of the formula $P \Rightarrow \text{wp}(s, Q)$, we can find a suitable strengthening of P by

solving the following abductive reasoning problem:

$$\begin{aligned} &\text{Find } \psi \text{ such that :} \\ &(1) P \wedge \psi \models \text{wp}(s, Q) \quad (2) P \wedge \psi \not\models \text{false} \end{aligned} \quad (3)$$

Here condition (1) states that $\{P \wedge \psi\} s \{Q\}$ is a valid Hoare triple, and (2) states that the speculated invariant ψ is consistent with precondition P . Since abductive reasoning is a well-studied problem, we use the *abduce* procedure described in prior work [16] to automatically infer candidate strengthenings ψ . Also, note that a call to *abduce* at line 7 may yield multiple predicates ψ_1, \dots, ψ_n , all of which constitute valid solutions for Equation 3.

Since the predicates Φ generated using abduction in lines 5–7 are merely *candidate* invariants, the next phase of the algorithm performs a fixed-point computation in which we drop every $\psi \in \Phi$ that is not a monitor invariant. Specifically, for each predicate ψ in Φ , we check whether (a) it holds initially (lines 11–13) and (b) whether it is preserved by each CCR in the monitor (lines 15–16). To determine whether ψ holds initially, we check the validity of the Hoare triple $\{true\} \text{Ctr}(M) \{\psi\}$, where $\text{Ctr}(M)$ denotes the constructor of M .⁶ If this Hoare triple is not valid, we simply drop ψ from set Φ . Next, to determine whether ψ is preserved by CCR w , we check the validity of the Hoare triple $\{I \wedge \text{Guard}(w)\} \text{Body}(w) \{\psi\}$, where I denotes the conjunction of all predicates in Φ . If this triple is invalid for any CCR in M , we again drop ψ from the set Φ . We then repeat this process until I satisfies both the initiation and consecution requirements. It is easy to see that formula I returned by INFERMONITORINV constitutes a valid monitor invariant.

6 Implementation

We have implemented our proposed method in a tool called EXPRESSO. Our implementation leverages the Soot program analysis infrastructure [46] and invokes the Z3 SMT solver [12] for checking logical validity. In what follows, we discuss some important design choices that are not addressed in previous sections.

Generating Java code. While the target language (IR) presented in Section 3.3 is convenient for describing our transformation, it does not yield valid Java code. Our implementation converts programs in this IR to valid Java code in the following manner. First, we associate a condition variable with the guard of every *waituntil* statement. Now, given a *waituntil* statement w with associated guard p , body s , and condition variable c , we then generate the following code⁷:

```
while(!p) {c.await();}; s
```

Furthermore, for each $(p_i, ?) \in \text{Signals}(w)$, we generate the code `if(pi) ci.signal()`, and for $(p_i, \checkmark) \in \text{Signals}(w)$, we emit `ci.signal()`. For each $(p_i, _) \in \text{Broadcasts}(w)$, we generate the same code where `signal` is replaced with `signalAll`.

⁶For simplicity, we assume a single constructor; if there are multiple ones, this triple needs to be checked for all constructors.

⁷Note that our implementation uses the `ReentrantLock` class.

Instrumentation for predicates with local variables. To support conditional signaling for predicates with local variables, EXPRESSO augments the monitor code with a data structure that tracks the values of local variables for any thread that is blocked on a predicate p . The code generated by EXPRESSO then uses this data structure to check whether p actually evaluates to true at program points that require conditional signaling for p .

Lazy broadcasts. EXPRESSO provides an option for performing broadcasts lazily. Consider a `waituntil` statement w such that $(p, _) \in \text{Broadcasts}(w)$. Rather than emitting the code `c.signalAll()` after the body of w , “lazy broadcast” notifies a single thread t blocked on p and ensures that t notifies all other threads by adding the instrumentation `if(p) c.signal()` after every `waituntil` statement with guard p . In our implementation, we enable this option by default to minimize context switches.

Discharging Hoare triples. EXPRESSO discharges any Hoare triple $\{P\} s \{Q\}$ by computing the weakest precondition of Q with respect to s and performing a validity check. Since s can contain pointers, EXPRESSO uses the points-to information provided by DOOP [6] to produce a whole-program model of the heap. In particular, given a store statement $v.f = e$, EXPRESSO generates additional statements of the form `if(v = xi) xi.f = e` where x_i is a potential alias of v .

7 Evaluation

We evaluate EXPRESSO by performing experiments that are designed to answer the following research questions:

- How does the code generated by EXPRESSO compare against hand-written explicit-signal code?
- How does our solution compare against run-time systems that provide support for implicit signals?
- How long does EXPRESSO take to generate code?

Benchmarks. The benchmarks used in our evaluation come from two different sources, namely all AutoSynch benchmarks from [30] and monitors collected from popular open-source projects from Github. We collected the Github benchmarks by writing a crawler that identifies potential monitors in Java programs using keywords such as `wait`, `signal`, `notify` etc. We then manually inspected these results in decreasing order of Github ranking (a mix of stars and forks) and identified self-contained modules (i.e., monitors) that encapsulate shared state. This process requires manual effort because we need to isolate the monitor code and insert it in a stress-testing harness.

Performance evaluation methodology. We evaluate performance using the same methodology used for evaluating AutoSynch in [30]. Specifically, we use *saturation tests* [8] in which threads only access the monitor and perform no extra work outside of the monitor. This set-up allows us to stress-test the monitor code and meaningfully compare

our solution with run-time solutions and near-optimal hand-written code (from the original AutoSynch benchmarks or from the GitHub project).

We perform measurements using the JMH framework [44], which is a benchmarking tool for rigorous measurement in JVM-based languages. All measurements are conducted on a 16-way (8 core x 2 SMT) Intel Xeon CPU E5-2640 v3 2.60GHz with 132 GB of memory using JDK 1.8.0_101-b13.

Performance results. The results of our performance evaluation are presented in Figures 8 and 9. Specifically, Figure 8 shows the results for the AutoSynch benchmarks, augmented with the readers-writers example of Section 2. Figure 9 presents results for monitors found in popular GitHub projects. Each graph plots the average time (in milliseconds) per monitor operation (e.g., `enterReader`) against the number of threads.

In virtually all cases, the performance of EXPRESSO-generated code is very close to hand-written explicit-signal code. The only significant differences are in the “H2O Barrier” benchmark under low concurrency and “Dining Philosophers” under high concurrency. In the latter, the explicit signalling code has knowledge of the problem structure itself, so it avoids all wakeups that do not lead to progress.

Comparing to AutoSynch, EXPRESSO outperforms it by 1.56x on average over all benchmarks. EXPRESSO significantly outperforms AutoSynch for about half the benchmarks of Figure 8, which are chosen or written by the AutoSynch authors themselves. On a few occasions, AutoSynch slightly outperforms EXPRESSO-generated code. As we discuss in Section 8, AutoSynch offers dynamic structures for quick inequality comparisons between shared variables and local values (which are captured as constants while the thread is waiting). This custom optimization can also be added to EXPRESSO but the emphasis of our work has been on statically eliminating unnecessary signalling, rather than minimizing the overhead of dynamic checks.

The monitors found in GitHub projects (Figure 9) are more representative of synchronization patterns in-the-wild. EXPRESSO performs very well on these benchmarks, matching hand-optimized code and significantly outperforming AutoSynch: by 1.62x on average, and up to 2.5x on a high-concurrency setting with 128 threads.

Upon closer inspection of these benchmarks, we observe that the symbolic reasoning needed to achieve the results from Figure 9 is far from trivial. As a simple example, “ConcurrencyThrottle” from the Spring framework has a waiting condition `threadCount < threadLimit` triggered by the statement `threadCount--` in the monitor exit operation. In order to avoid broadcasts, EXPRESSO needs to infer a monitor invariant that allows it to establish that whenever a thread enters the monitor, the waiting condition has to become true again due to a `threadCount++` operation. Because these increment and decrement operations are distant, symbolic reasoning has to model the semantics of all intervening program statements and establish that the operations commute. This kind

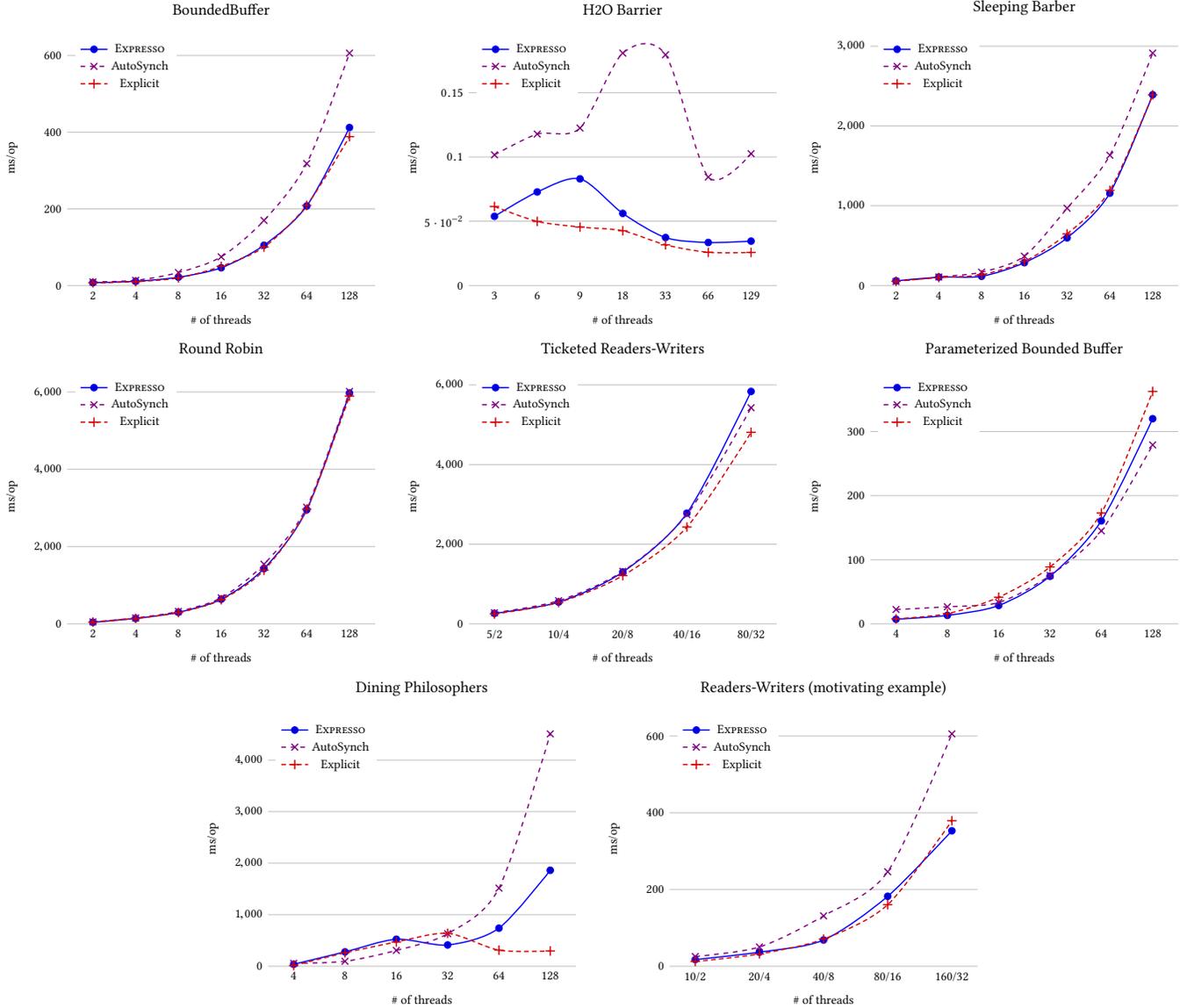


Figure 8. Performance over AutoSynch benchmarks and readers-writers example.

of reasoning is necessary for EXPRESSO to achieve the performance results from Figure 9 in all benchmarks. Furthermore, the inferred monitor invariants are often intricate—for instance, the “AsyncDispatch” invariant (shown in Appendix D) from the Gradle codebase has 22 sub-terms (12 equality/inequality comparisons and arithmetic operations and 10 logical connectives).

To summarize, these results demonstrate the plausibility of a practical and efficient implementation for implicit-signal monitors. In particular, the code generated by EXPRESSO is comparable to hand-written code even for saturation tests that stress-test the monitor. Furthermore, EXPRESSO’s

implicit-signal monitor implementation consistently outperforms AutoSynch on monitors extracted from real-world codebases such as Spring framework, Gradle, ExoPlayer, greenDAO, etc.

Analysis time. Table 1 shows the time that EXPRESSO takes to synthesize the explicit-signal code from its corresponding implicit-signal version for each benchmark. In most cases, the symbolic reasoning time is in the order of a few seconds. The only exception is the largest benchmark, AsyncDispatch, whose compilation takes 28.3 seconds. This example takes longer to analyze because some of the predicates depend on Java library operations that EXPRESSO also needs to analyze. Overall, these results demonstrate that EXPRESSO is practical

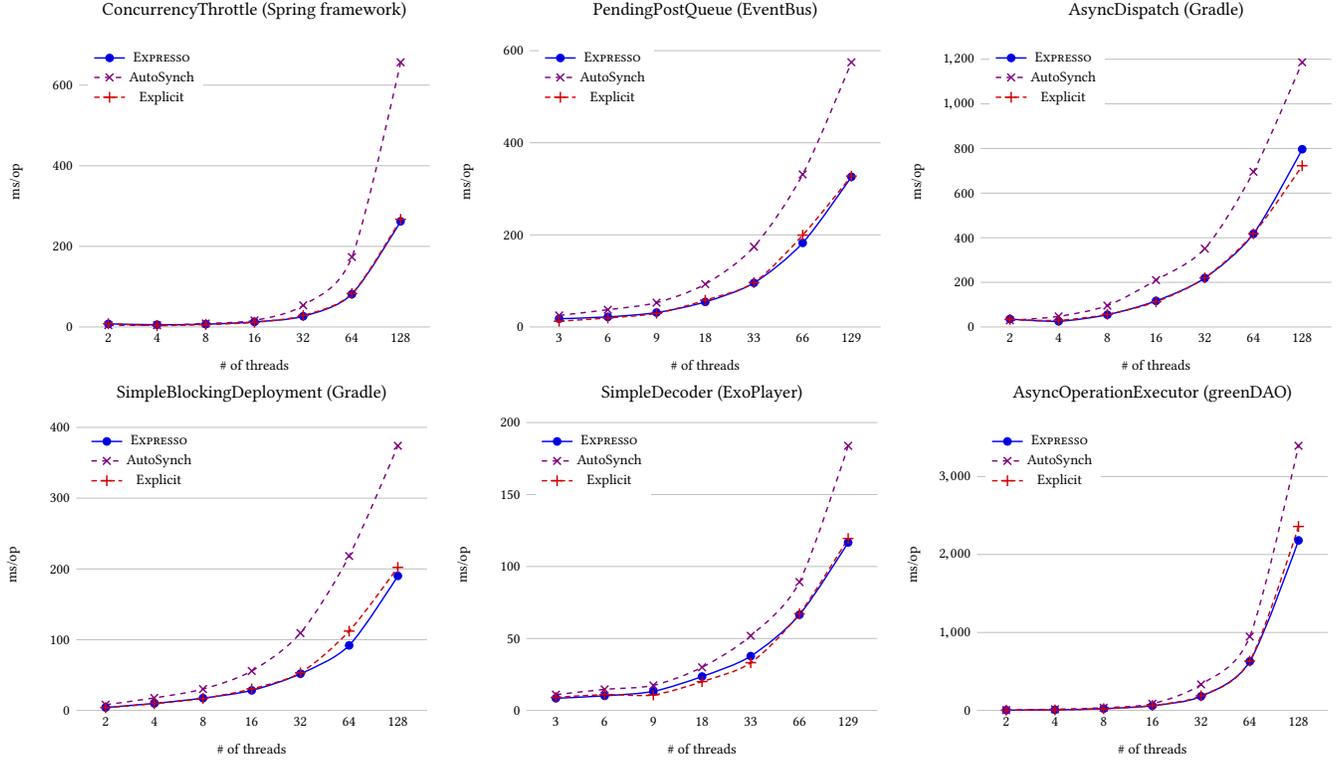


Figure 9. Performance over monitor code in popular GitHub projects.

and that it generates code whose performance is comparable to hand-written code in virtually all cases.

Generated code. We also assessed the quality of the code generated by EXPRESSO by manually inspecting the synthesized explicit-signal monitors. For most benchmarks obtained from Github projects, we found that the code generated by EXPRESSO is very similar to hand-written code. In the case of the AutoSynch benchmarks, however, we found some examples (e.g., Dining Philosophers) where the EXPRESSO-generated code differs significantly from hand-written code. For these benchmarks, manually-written code leverages dynamic data structures to achieve optimal signaling, whereas EXPRESSO uses the fixed strategy described in Section 6 for handling predicates with local variables.

8 Related Work

Our implicit-signal monitors have several relatives in the literature. We discuss representative past work, in loose thematic groupings of decreasing affinity with our work.

Language designs and run-time support for concurrency. There is a very rich literature on language support for concurrency, dating back to the early 1960s [3]. Dijkstra originally proposed the concept of *semaphores* to provide a friendlier and more efficient programming abstraction than the busy-wait design [15]. In later work, Hoare proposed the concept of *conditional critical regions (CCR)* [27], which overcome some of the difficulties associated with semaphores by

providing a more structured notation for specifying synchronization. In particular, every shared variable in a CCR must belong to a resource, and variables in a resource can only be accessed within so-called *region* statements of the form $\text{region } r \text{ when } B \rightarrow S$, where B is a guard and S is a statement. Concurrently to the introduction of CCRs, Dijkstra proposed the notion of *monitors*, which provide more structure than conditional critical regions and can be implemented as efficiently as semaphores [14]. To this day, monitors remain a

Benchmark	Time (sec.)
BoundedBuffer	2.5
H2OBarrier	2.3
Sleeping Barber	1.6
Round Robin	1.2
Ticketed Readers-Writers	3.8
Param. Bounded Buffer	2.5
Dining Philosophers	5.4
Readers-Writers	1.5
ConcurrencyThrottle	1.0
PendingPostQueue	0.5
AsyncDispatch	28.3
SimpleBlockingDeployment	0.4
SimpleDecoder	10.7
AsyncOperationExecutor	2.1

Table 1. Compilation time for benchmarks.

popular concurrent programming paradigm, and the “monitor pattern” is widely used in many programming languages, including Java and C++.

Even though early proposals for monitors advocate an implicit signaling mechanism [3, 4], most modern monitor implementations use explicit signaling due to performance considerations. More recent work in this area aim to popularize implicit signal monitors by providing a more efficient implementation [8, 30]. For example, the recent AutoSynch work [30] attempts to improve the efficiency of automatic signaling through a combination of efficient dynamic indexing and simple static analysis. AutoSynch offers sophisticated handling of local state in a thread. If threads wait on conditions based on standard equality/inequality patterns over local variables, the system dynamically snapshots the values of local variables and treats them as run-time constants—the variable values cannot have changed while the thread is waiting. As a result, the predicates have a standard structure of comparisons with constants. The efficient notification algorithm then works much like database indexing: it computes, given which shared values changed, what waiting predicates could possibly have been affected. Our approach is distinguished by its use of reasoning techniques for statically inferring when a condition must, may-not, or may have become true. In order to do so, our technique needs to take into consideration several issues that are not addressed by previous work, e.g., handling of memory aliases, interprocedural reasoning, etc. As shown in the evaluation, our method significantly outperforms the AutoSynch solution on many benchmarks and is comparable to hand-written code for most examples.

Transactional memory. CCRs have been recently reified in several language designs, such as that of Harris and Frasier [24]. Such designs can be viewed as special cases of *transactional memory (TM)* [25], which has attracted enormous attention in the literature, as an alternative of lock-based synchronization [38].

Our implicit-signal monitors are both *less* and *more* ambitious than transactional memory techniques, in different respects. Monitors do not attempt to automate-away lock-based synchronization: we explicitly require the programmer to declare different resources together with their synchronization policy. Concretely, each monitor can be thought of as a single lock, whereas transactional memory techniques do away with distinctions between locks in favor of a single atomic construct. The same essential feature is kept in Harris and Frasier’s conditional critical regions [24]: although atomic statements can have a condition associated with them, their code blocks are guaranteed to execute with atomic semantics, relative to any other atomic block, under whichever condition. Therefore, implicit-signal monitors are inherently lower-level than TM approaches, requiring more care on behalf of the programmer, but also imposing no overhead for maintaining atomic semantics. At the same time,

our implicit-signal monitors are more ambitious on the implementation side. Static reasoning over abstract conditions results in the removal of extraneous signaling overheads, enabling high-performance execution.

Automation of synchronization. Static analysis has been applied to the optimization of synchronization primitives in the past, mostly in the context of auto-locking or lock-inference techniques [11, 20, 26, 34, 41]. The language model these techniques implement is much closer to transactional memory than implicit-signal monitors: the analysis attempts to infer which locks protect which shared data, as well as which data are not thread-shared. Furthermore, these static analyses typically produce a whole-program model of shared data and do not reason over symbolic conditions, as in our approach. Techniques that infer synchronization given specifications and abstractions [10, 19, 47] often use similar reasoning techniques as our approach (e.g., SMT solvers), but start from much more abstract input and place an emphasis on correctness, not on approaching hand-written code performance. Similar comments apply to approaches that synthesize synchronization actions, given abstract models of program behavior, using control theory techniques [35, 48].

Program analysis for concurrency. A major application of static analysis in the domain of concurrency has been in guaranteeing safety or finding bugs. There are techniques for ensuring safe programming using advanced typing [5], analyses for static race detection [43], approaches to concurrency bug fixing [31], tools that flag suspicious concurrency patterns [29, 49], and much more. Additionally, dynamic techniques for concurrency bug detection [1, 13, 33, 40, 45] often benefit from symbolic reasoning, especially in approaches inspired by model checking or symbolic execution techniques [32, 36, 42]. Such past work is only superficially related to ours, since the aims, programming abstractions, and analysis techniques used are quite different.

Abductive reasoning in program analysis. Our proposed approach uses abductive reasoning for automatically inferring monitor invariants. The use of abduction in program analysis is quite common, especially in the context of modular analysis [2, 9, 18], loop invariant generation [17], and specification inference [2, 50]. However, our use of abductive reasoning differs from prior work in that we use abduction to generate candidate predicates over which we synthesize monitor invariants using predicate abstraction.

9 Limitations

While this paper takes a first step towards synthesizing efficient explicit-signal implementations of implicit-signal monitors, it makes several assumptions that allow the proposed approach to be practical. First, we assume that predicates appearing in `waituntil` statements are expressible in some first-order theory (e.g., linear arithmetic). Hence, if the predicate corresponds to the result of a complicated function

containing loops or recursion, our approach must be conservative and potentially broadcast at the end of every CCR. However, we believe that such cases are in practice, allowing EXPRESSO to generate efficient code for monitors taken from real-world applications.

Second, our approach assumes that monitors do not contain nested `waituntil` statements and that calls that occur inside the body of a monitor method are *closed calls* [3], meaning that they do not release access to the monitor. We believe that the latter assumption is quite realistic, since most modern implementations of monitor-like constructs also assume closed-call semantics. On the other hand, handling nested `waituntil` statements is an interesting direction for future work. In particular, we plan to investigate how to leverage more sophisticated forms of symbolic reasoning to analyze implicit-signal monitors with nested `waituntil` statements.

10 Conclusion

We have presented a new technique for automatic signalling in concurrent programs. Our method statically analyzes an implicit-signal monitor implementation to determine where signaling is necessary and automatically generates its corresponding explicit-signal version. Our method employs symbolic reasoning to avoid unnecessary context switches and run-time evaluation of predicates. We have implemented the proposed algorithm in a tool called EXPRESSO and evaluate it on 14 benchmarks. Our evaluation shows that the code generated by EXPRESSO is very similar to hand-written code and that it significantly outperforms a state-of-the-art run-time solution in most cases.

Acknowledgments

We would like to thank Doug Lea, the anonymous reviewers, and the members of the UTOPiA group for their insightful feedback. The first three authors of this paper are supported in part by NSF Award #1712060, NSF Award #1453386, AFRL Award #8750-15-2-0096, and the last author is supported by ERC grant 307334 (SPADE), a Facebook Research and Academic Relations Program award, and an Oracle Labs collaborative research grant.

References

- [1] Rahul Agarwal and Scott D. Stoller. 2006. Run-time Detection of Potential Deadlocks for Programs with Locks, Semaphores, and Condition Variables. In *Proceedings of the 2006 Workshop on Parallel and Distributed Systems: Testing and Debugging (PADTAD '06)*. ACM, New York, NY, USA, 51–60. <https://doi.org/10.1145/1147403.1147413>
- [2] Aws Albarghouthi, Isil Dillig, and Arie Gurfinkel. 2016. Maximal specification synthesis. In *ACM SIGPLAN Notices*, Vol. 51. ACM, 789–801.
- [3] Gregory R. Andrews. 1991. *Concurrent Programming: Principles and Practice*. Benjamin-Cummings Publishing Co., Inc., Redwood City, CA, USA.
- [4] Russell Atkinson and Carl Hewitt. [n. d.]. Synchronization in actor systems. In *POPL '77: Proc. 4th symposium on Principles of Programming Languages*. ACM Press, 267–280. <https://doi.org/10.1145/512950.512975>
- [5] Chandrasekhar Boyapati, Robert Lee, and Martin Rinard. 2002. Ownership types for safe programming: Preventing data races and deadlocks. In *ACM Sigplan Notices*, Vol. 37. ACM, 211–230.
- [6] Martin Bravenboer and Yannis Smaragdakis. 2009. Strictly Declarative Specification of Sophisticated Points-to Analyses. In *ACM Symposium on Object Oriented Programming: Systems, Languages, and Applications (OOPSLA)*. ACM Press, 243–262.
- [7] Peter A Buhr, Michel Fortier, and Michael H Coffin. 1995. Monitor classification. *ACM Computing Surveys (CSUR)* 27, 1 (1995), 63–107.
- [8] Peter A Buhr and Ashif S Harji. 2005. Implicit-signal monitors. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 27, 6 (2005), 1270–1343.
- [9] Cristiano Calcagno, Dino Distefano, Peter O’Hearn, and Hongseok Yang. 2009. Compositional shape analysis by means of bi-abduction. In *ACM SIGPLAN Notices*, Vol. 44. ACM, 289–300.
- [10] Pavol Černý, Thomas A Henzinger, Arjun Radhakrishna, Leonid Ryzhyk, and Thorsten Tarrach. 2013. Efficient synthesis for concurrency by semantics-preserving transformations. In *International Conference on Computer Aided Verification*. Springer, 951–967.
- [11] Sigmund Cheren, Trishul Chilimbi, and Sumit Gulwani. 2008. Inferring locks for atomic sections. In *SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. ACM, New York, NY, USA, 304–315. <https://doi.org/10.1145/1375581.1375619>
- [12] Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An efficient SMT solver. *Tools and Algorithms for the Construction and Analysis of Systems* (2008), 337–340.
- [13] Dongdong Deng, Wei Zhang, and Shan Lu. 2013. Efficient concurrency-bug detection across inputs. In *Acm Sigplan Notices*, Vol. 48. ACM, 785–802.
- [14] Edsger W Dijkstra. 1971. Hierarchical ordering of sequential processes. In *The origin of concurrent programming*. Springer, 198–227.
- [15] Edsger W. Dijkstra. 2002. *Cooperating Sequential Processes*. Springer New York, New York, NY, 65–138. https://doi.org/10.1007/978-1-4757-3472-0_2
- [16] Isil Dillig and Thomas Dillig. 2013. Explain: a tool for performing abductive inference. In *International Conference on Computer Aided Verification*. Springer, 684–689.
- [17] Isil Dillig, Thomas Dillig, Boyang Li, and Ken McMillan. 2013. Inductive invariant generation via abductive inference. In *Acm Sigplan Notices*, Vol. 48. ACM, 443–456.
- [18] Isil Dillig, Thomas Dillig, Boyang Li, Ken McMillan, and Mooly Sagiv. 2017. Synthesis of circular compositional program proofs via abduction. *International Journal on Software Tools for Technology Transfer* 19, 5 (2017), 535–547.
- [19] E Allen Emerson and Roopsha Samanta. 2011. An Algorithmic Framework for Synthesis of Concurrent Programs. In *ATVA*. Springer, 522–530.
- [20] Michael Emmi, Jeffrey S. Fischer, Ranjit Jhala, and Rupak Majumdar. 2007. Lock allocation. In *POPL '07: Proc. 34th symposium on Principles of Programming Languages*. ACM Press, 291–296.
- [21] Cormac Flanagan and K Rustan M Leino. 2001. Houdini, an annotation assistant for ESC/Java. In *International Symposium of Formal Methods Europe*. Springer, 500–517.
- [22] Rui Gu, Guoliang Jin, Linhai Song, Linjie Zhu, and Shan Lu. 2015. What change history tells us about thread synchronization. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. ACM, 426–438.
- [23] Per Brinch Hansen. 1973. *Operating system principles*. Prentice-Hall, Inc.
- [24] Tim Harris and Keir Fraser. 2003. Language support for lightweight transactions. In *OOPSLA '03: Proc. 18th conf. on Object-oriented Programming, Systems, Languages, and Applications*. ACM Press, Anaheim, CA, 388–402. <https://doi.org/10.1145/949305.949340>
- [25] Maurice Herlihy and J. Eliot B. Moss. 1993. Transactional Memory: Architectural Support for Lock-Free Data Structures. In *Proceedings of the 20th Annual International Symposium on Computer Architecture*. 289–300.

- [26] Michael Hicks, Jeffrey S. Foster, and Polyvios Prattikakis. 2006. Lock Inference for Atomic Sections. In *Proceedings of the First ACM SIGPLAN Workshop on Languages, Compilers, and Hardware Support for Transactional Computing*.
- [27] C. A. R. Hoare. 1971. Towards a theory of parallel programming. In *International Seminar on Operating System Techniques*.
- [28] Charles Antony Richard Hoare. 1974. Monitors: An operating system structuring concept. In *The origin of concurrent programming*. Springer, 272–294.
- [29] David Hovemeyer and William Pugh. 2004. Finding concurrency bugs in Java. In *Proc. of PODC*, Vol. 4.
- [30] Wei-Lun Hung and Vijay K. Garg. 2013. AutoSynch: An Automatic-signal Monitor Based on Predicate Tagging. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '13)*. ACM, New York, NY, USA, 253–262. <https://doi.org/10.1145/2491956.2462175>
- [31] Guoliang Jin, Wei Zhang, Dongdong Deng, Ben Liblit, and Shan Lu. 2012. Automated Concurrency-Bug Fixing.. In *OSDI*, Vol. 12. 221–236.
- [32] Pallavi Joshi, Mayur Naik, Chang-Seo Park, and Koushik Sen. 2009. CalFuzzer: An extensible active testing framework for concurrent programs. In *Computer Aided Verification*. Springer, 675–681.
- [33] Pallavi Joshi, Mayur Naik, Koushik Sen, and David Gay. 2010. An effective dynamic analysis for detecting generalized deadlocks. In *Proceedings of the eighteenth ACM SIGSOFT international symposium on Foundations of software engineering*. ACM, 327–336.
- [34] Vineet Kahlon. 2012. Automatic lock insertion in concurrent programs. In *Formal Methods in Computer-Aided Design (FMCAD), 2012*. IEEE, 16–23.
- [35] Terence Kelly, Yin Wang, Stéphane Lafortune, and Scott Mahlke. 2009. Eliminating concurrency bugs with control engineering. *Computer* 42, 12 (2009).
- [36] Sepideh Khoshnood, Markus Kusano, and Chao Wang. 2015. ConcBugAssist: constraint solving for diagnosis and repair of concurrency bugs. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis*. ACM, 165–176.
- [37] Shuvendu K Lahiri and Shaz Qadeer. 2009. Complexity and Algorithms for Monomial and Clausal Predicate Abstraction.. In *CADE*. Springer, 214–229.
- [38] James R. Larus and Ravi Rajwar. 2007. *Transactional Memory*. Morgan & Claypool.
- [39] Doug Lea. 1999. *Concurrent Programming in Java. Second Edition: Design Principles and Patterns* (2nd ed.). Addison-Wesley.
- [40] Brandon Lucia and Luis Ceze. 2009. Finding concurrency bugs with context-aware communication graphs. In *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*. ACM, 553–563.
- [41] Bill McCloskey, Feng Zhou, David Gay, and Eric Brewer. 2006. Autolocker: Synchronization inference for atomic sections. In *POPL '06: Proc. 33rd symposium on Principles of Programming Languages*. ACM Press, 346–358. <https://doi.org/10.1145/1111037.1111068>
- [42] Madanlal Musuvathi, Shaz Qadeer, Thomas Ball, Gerard Basler, Pira-manayagam Arumuga Nainar, and Iulian Neamtii. 2008. Finding and Reproducing Heisenbugs in Concurrent Programs.. In *OSDI*, Vol. 8. 267–280.
- [43] Mayur Naik, Alex Aiken, and John Whaley. 2006. *Effective static race detection for Java*. Vol. 41. ACM.
- [44] Aleksey Shipilev, Sergey Kuksenko, Anders Astrand, Staffan Freiberg, and Henrik Loeff. [n. d.]. OpenJDK: JMH. ([n. d.]). <http://openjdk.java.net/projects/code-tools/jmh/>
- [45] Yannis Smaragdakis, Jacob Evans, Caitlin Sadowski, Jaeheon Yi, and Cormac Flanagan. 2012. Sound Predictive Race Detection in Polynomial Time. In *ACM Symposium on Principles of Programming Languages (POPL)*. ACM Press, 387–400.
- [46] Raja Vallée-Rai, Phong Co, Etienne Gagnon, Laurie Hendren, Patrick Lam, and Vijay Sundareshan. 1999. Soot-a Java bytecode optimization framework. In *Proceedings of the 1999 conference of the Centre for Advanced Studies on Collaborative research*. IBM Press, 13.
- [47] Martin Vechev, Eran Yahav, and Greta Yorsh. 2010. Abstraction-guided Synthesis of Synchronization. In *Proceedings of the 37th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '10)*. ACM, New York, NY, USA, 327–338. <https://doi.org/10.1145/1706299.1706338>
- [48] Yin Wang, Stéphane Lafortune, Terence Kelly, Manjunath Kudlur, and Scott Mahlke. 2009. The theory of deadlock avoidance via discrete control. In *ACM SIGPLAN Notices*, Vol. 44. ACM, 252–263.
- [49] Weiwei Xiong, Soyeon Park, Jiaqi Zhang, Yuanyuan Zhou, and Zhiqiang Ma. 2010. Ad Hoc Synchronization Considered Harmful.. In *OSDI*, Vol. 10. 163–176.
- [50] Haiyan Zhu, Thomas Dillig, and Isil Dillig. 2013. Automated inference of library specifications for source-sink property verification. In *Asian Symposium on Programming Languages and Systems*. Springer, 290–306.

Appendix A: Well-formedness of Monitor Trace

In this section, we formalize what it means for a monitor trace τ to be syntactically well-formed. Towards this goal, we first define the projection of a trace τ onto a thread t , denoted as $\tau \downarrow t$.

Definition 10.1. The projection of trace τ onto thread t , denoted $\tau \downarrow t$, is defined as follows:

$$\begin{aligned} (t', w, true) \downarrow t &= [w] && \text{if } t' = t \\ (t', w, b) \downarrow t &= nil && \text{if } t' \neq t \text{ or } b = false \\ e :: \tau \downarrow t &= (e \downarrow t) :: (\tau \downarrow t) \end{aligned}$$

Definition 10.2. We say that $\tau \downarrow t$ is well-formed if $\tau \downarrow t = [S_1, \dots, S_n, S_{n+1}]$, each S_i where $i \in [1, n]$ corresponds to the body of some method, and S_{n+1} is a prefix of the body of some method.

Using these definitions, we can now define syntactic well-formedness of monitor traces as follows:

Definition 10.3. We say that a monitor trace τ is syntactically well-formed if the following conditions are satisfied:

1. For every thread t , $\tau \downarrow t = S_1; \dots; S_n; S_{n+1}$ is well-formed
2. If $\tau[i] = (t, w, true)$ and $\tau[i + 1] = (t', w', b)$, then either w is the last statement in its method or $t' = t$ and w' is the successor of w .

Appendix B: Soundness Proof

In this section, we prove the soundness of the PLACE SIGNAL algorithm. The proof of Theorem 4.1 follows immediately from Lemmas 10.4 and 10.5.

In the remainder of this section, we use the term *event* to also refer to pairs (t, w) without the corresponding boolean. We refer to t as *thread*(e) and, if $w = (p, s)$, we refer to p as *pred*(e). Given two sets $\mathcal{N}, \mathcal{N}'$, we will assume that the total order relation $<$ on events is defined in a way that respects the property $\mathcal{N}' \subseteq \mathcal{N} \Rightarrow \min(\mathcal{N}) = \min(\mathcal{N}')$.

Furthermore, for the rest of the section it is useful to keep in mind that the source language maintains the following invariant:

Invariant If $(\sigma_0, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma, _, \mathcal{B}, \mathcal{N})$, then:

- If $(t, w) \in \mathcal{N}$, then we have $(t, w) \in \mathcal{B}$.
- If $(t, w) \in \mathcal{B}$ and $(\sigma, t) \models \text{Guard}(w)$, then $(t, w) \in \mathcal{N}$.

Lemma 10.4. Let $M' = \text{PlaceSignals}(M, I)$. If I is a correct monitor invariant, then for all monitor states σ and all well-formed traces τ , if $(\sigma, \tau, \emptyset, \emptyset) \Longrightarrow^* (\sigma', \epsilon, _, _)$ then $(\sigma, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma', \epsilon, _, _)$.

Proof. Follows from Lemma 10.8. \square

Lemma 10.5. Let $M' = \text{PlaceSignals}(M, I)$. If I is a correct monitor invariant, then for all monitor states σ and all well-formed traces τ . If $(\sigma, \tau, \emptyset, \emptyset) \longrightarrow^* (\sigma', \epsilon, _, _)$ and τ is normalized with respect to σ , then $(\sigma, \tau, \emptyset, \emptyset) \Longrightarrow^* (\sigma', \epsilon, _, _)$.

Proof. Follows from Lemma 10.9. \square

Definition 10.6. (Invalidation). Let $e = (w, t)$ be a monitor event where $w = (p, s)$, t be a thread identifier, and p a predicate. We write $\text{Invalidate}(e, t', p')$, if, for any state σ such that $\langle s, t, \sigma \rangle \Downarrow \sigma'$, we have $(\sigma', t') \not\models p'$.

Definition 10.7. (Agreement). We say that $(\sigma, \mathcal{B}, \mathcal{N})$ agrees with $(\sigma, \mathcal{B}', \mathcal{N}')$, written $(\sigma, \mathcal{B}, \mathcal{N}) \sim (\sigma, \mathcal{B}', \mathcal{N}')$ if the following conditions are satisfied:

1. $\mathcal{B} = \mathcal{B}'$
2. $\mathcal{N} \supseteq \mathcal{N}'$
3. If $e \in (\mathcal{N} - \mathcal{N}')$, then either
 - a. $(\sigma, \text{thread}(e)) \not\models \text{pred}(e)$, or
 - b. $\exists e' \in \mathcal{N}'$. $e' < e \wedge \text{pred}(e) = \text{pred}(e') \wedge \text{Invalidate}(e', \text{thread}(e), \text{pred}(e))$

Lemma 10.8. If $(\sigma, \mathcal{B}, \mathcal{N}) \sim (\sigma, \mathcal{B}', \mathcal{N}')$ and

$$(\sigma, \tau, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma', \tau', \mathcal{B}_1, \mathcal{N}_1)$$

then we have

$$(\sigma, \tau, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \tau', \mathcal{B}_2, \mathcal{N}_2)$$

and $(\sigma', \mathcal{B}_1, \mathcal{N}_1) \sim (\sigma', \mathcal{B}_2, \mathcal{N}_2)$.

Proof. The proof is by induction on τ .

- **Base case 1a:** $\tau = e = (t, w, \text{false})$ and $(t, w) \notin \mathcal{B}'$. In this case, we have $(\sigma, t) \not\models \text{Guard}(w)$ and $(\sigma, e, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma, \epsilon, \mathcal{B}' \cup \{\bar{e}\}, \mathcal{N}')$. Since $\mathcal{B} = \mathcal{B}'$, we have $(t, w) \notin \mathcal{B}$. Thus, we also have $(\sigma, e, \mathcal{B}', \mathcal{N}) \longrightarrow (\sigma, \epsilon, \mathcal{B}' \cup \{\bar{e}\}, \mathcal{N})$. Since $\mathcal{N}, \mathcal{N}'$ satisfy conditions (2) and (3) of the agreement definition and the state σ has not changed, $\mathcal{N}, \mathcal{N}'$ continue to satisfy conditions (2) and (3). Thus, we have $(\sigma, \mathcal{B}' \cup \{\bar{e}\}, \mathcal{N}') \sim (\sigma, \mathcal{B} \cup \{\bar{e}\}, \mathcal{N})$.
- **Base case 1b:** $\tau = e = (t, w, \text{false})$ and $(t, w) \in \mathcal{N}'$. In this case, we have $(\sigma, e, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma, \epsilon, \mathcal{B}', \mathcal{N}' \setminus \{\bar{e}\})$. Since $\mathcal{N} \supseteq \mathcal{N}'$, $(t, w) \in \mathcal{N}$, we also have $(\sigma, e, \mathcal{B}', \mathcal{N}) \longrightarrow (\sigma, \epsilon, \mathcal{B}', \mathcal{N} \setminus \{\bar{e}\})$. Since $\mathcal{N}, \mathcal{N}'$ satisfy conditions (2) and (3) of the agreement definition and the state σ has not changed, $\mathcal{N}, \mathcal{N}'$ continue to satisfy conditions (2) and (3). Thus, $(\sigma, \mathcal{B}', \mathcal{N}' \setminus \{\bar{e}\}) \sim (\sigma, \mathcal{B}', \mathcal{N} \setminus \{\bar{e}\})$.

- **Base case 2a:** $\tau = e = (t, w, \text{true})$ and $\bar{e} \notin \mathcal{B}'$. In this case, we have:

$$\frac{\begin{array}{l} e = (t, w, \text{true}) \\ \bar{e} \notin \mathcal{B}' \quad (\sigma, t) \models \text{Guard}(w) \\ \langle \text{Body}(w), t, \sigma \rangle \Downarrow \sigma' \\ \mathcal{N}_2 = \text{GetSignals}(w, \sigma', \mathcal{B}) \\ \mathcal{N}_3 = \text{GetBroadcasts}(w, \sigma', \mathcal{B}) \end{array}}{(\sigma, e, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma', \epsilon, \mathcal{B}', \mathcal{N}' \cup \mathcal{N}_2 \cup \mathcal{N}_3)}$$

Since $\mathcal{B} = \mathcal{B}'$, this implies $\bar{e} \notin \mathcal{B}$, thus, we have:

$$\frac{\begin{array}{l} e = (t, w, \text{true}) \\ \bar{e} \notin \mathcal{B} \quad (\sigma, t) \models \text{Guard}(w) \\ \langle \text{Body}(w), t, \sigma \rangle \Downarrow \sigma' \\ \mathcal{N}_1 = \{(t, w) \mid (t, w) \in \mathcal{B}, (\sigma', t) \models \text{Guard}(w)\} \end{array}}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \epsilon, \mathcal{B}, \mathcal{N} \cup \mathcal{N}_1)}$$

Part (a) of the agreement definition trivially holds, since $\mathcal{B} = \mathcal{B}'$. For part (b), we need to show that $\mathcal{N}_1 \cup \mathcal{N} \supseteq \mathcal{N}_2 \cup \mathcal{N}_3 \cup \mathcal{N}'$. That is, if $e^* \in \mathcal{N}_2 \cup \mathcal{N}_3 \cup \mathcal{N}'$, we also have $e^* \in \mathcal{N}_1 \cup \mathcal{N}$. Case 1: Suppose $e^* \in \mathcal{N}'$. By agreement, we have $e^* \in \mathcal{N}$, thus part (b) of agreement holds in this case. Case 2: Suppose $e^* \in \mathcal{N}_2 \cup \mathcal{N}_3$. This means $(\text{pred}(e^*), c) \in P_1 \cup P_2$ where $P_1 = \text{Signals}(w)$ and $P_2 = \text{Broadcasts}(w)$. Case 2a: If $c = ?$, then e^* is only added to $\mathcal{N}_2 \cup \mathcal{N}_3$ if $(\sigma', \text{thread}(e^*)) \models \text{pred}(e^*)$. But then, we have $e^* \in \mathcal{N}_1$. Case 2b: If $c = \checkmark$, then the PLACE SIGNALS algorithm ensures that $\vdash \{I \wedge \neg \text{pred}(e^*)\} \text{Body}(w) \{ \text{pred}(e^*) \}$. By correctness of this Hoare triple, we have $(\sigma', \text{thread}(e^*)) \models \text{pred}(e^*)$. Thus, $e^* \in \mathcal{N}_1$.

For part (c) of the agreement definition, we need to show that if $e^* \in (\mathcal{N} \cup \mathcal{N}_1) \setminus (\mathcal{N}' \cup \mathcal{N}_2 \cup \mathcal{N}_3)$, then either (i) $(\sigma', \text{thread}(e^*)) \not\models \text{pred}(e^*)$, or (ii) $\exists e' \in \mathcal{N}_2 \cup \mathcal{N}_3 \cup \mathcal{N}$ such that $e' < e^* \wedge \text{pred}(e) = \text{pred}(e') \wedge \text{Invalidate}(e', \text{thread}(e^*), \text{pred}(e^*))$. We consider two cases:

- Case 1: $e^* \in \mathcal{N}$, but $e^* \notin \mathcal{N}' \cup \mathcal{N}_2 \cup \mathcal{N}_3$ (i.e., $e^* \notin \mathcal{N}'$, $e^* \notin \mathcal{N}_2$, $e^* \notin \mathcal{N}_3$). Since $e^* \in \mathcal{N} - \mathcal{N}'$, we have (from agreement) that either (a) $(\sigma, \text{thread}(e^*)) \not\models \text{pred}(e^*)$, or (b) $\exists e' \in \mathcal{N}'$ such that $e' < e^*$, $\text{pred}(e') = \text{pred}(e^*)$ and $\text{Invalidate}(e', \text{thread}(e^*), \text{pred}(e^*))$.
- * Case 1a: $(\sigma, \text{thread}(e^*)) \not\models \text{pred}(e^*)$. If we also have $(\sigma', \text{thread}(e^*)) \not\models \text{pred}(e^*)$, then part (c) of agreement definition obviously holds. If $(\sigma', \text{thread}(e^*)) \models \text{pred}(e^*)$, then $e^* \in \mathcal{N}_1$, and we argue correctness in Case 2.
- * Case 1b: $(\sigma, \text{thread}(e^*)) \not\models \text{pred}(e^*)$ and $\exists e' \in \mathcal{N}'$ such that $(e' < e^*, \text{pred}(e') = \text{pred}(e^*))$, and $\text{Invalidate}(e', \text{thread}(e^*), \text{pred}(e^*))$. Since such an e' is also in $\mathcal{N}' \cup \mathcal{N}_2 \cup \mathcal{N}_3$, part (c) of the agreement definition is preserved.
- Case 2: $e^* \in \mathcal{N}_1$, but $e^* \notin \mathcal{N}'$, $e^* \notin \mathcal{N}_2$, and $e^* \notin \mathcal{N}_3$. There are only two ways in which $e^* \in \mathcal{N}_1$ but not in \mathcal{N}_2 or \mathcal{N}_3 :
 - * Case 2a: $(\text{pred}(e^*), c) \in P_1 \cup P_2$ where $P_1 = \text{Signals}(w)$ and $P_2 \in \text{Broadcasts}(w)$. Using the invariant of the PLACE SIGNALS algorithm, this can only happen if $\vdash \{I \wedge \neg \text{pred}(e^*)\} \text{Body}(w) \{ \neg \text{pred}(e^*) \}$. If $\text{pred}(e^*)$ was

true before executing $Body(w)$ (i.e., $(\sigma, thread(e^*)) \models pred(e^*)$), then e^* would also be in \mathcal{N} , which implies it is also in \mathcal{N}' , so this case is not possible. Otherwise, if $(\sigma, thread(e^*)) \not\models pred(e^*)$, the correctness of the Hoare triple implies that $(\sigma', thread(e^*)) \not\models pred(e^*)$, which contradicts the fact that $e^* \in \mathcal{N}_1$.

- * Case 2b: $(pred(e^*), c) \in P_2$, but there exists another e' such that (i) $e' < e^*$, and (ii) $pred(e^*) = pred(e')$. The only way in which e^* is added to P_2 but not P_1 by the PLACESIGNALS algorithm is if the following Hoare triple holds for all waituntil statements w with guard $pred(e^*)$:

$$\{I \wedge Guard(w)\} Body(w) \{\neg Guard(w)\}$$

Since e' has the same guard as $pred(e^*)$, the validity of the Hoare triple implies $Invalidate(e', thread(e^*), pred(e^*))$. Thus part (c) of the agreement definition again holds.

- **Base case 2b:** $\tau = e = (t, w, true)$ and $\bar{e} = \min(\mathcal{N}')$. In this case, we have:

$$\frac{\begin{array}{l} e = (t, w, true) \\ \bar{e} = \min(\mathcal{N}') \quad (\sigma, t) \models Guard(w) \\ \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \\ \mathcal{N}_2 = GetSignals(w, \sigma', \mathcal{B}) \\ \mathcal{N}_3 = GetBroadcasts(w, \sigma', \mathcal{B}) \end{array}}{(\sigma, e, \mathcal{B}, \mathcal{N}') \Longrightarrow (\sigma', \epsilon, \mathcal{B} \setminus \{\bar{e}\}, (\mathcal{N}' \cup \mathcal{N}_2 \cup \mathcal{N}_3) \setminus \{\bar{e}\})}$$

Since $\mathcal{N}' \subseteq \mathcal{N}$, we have $\min(\mathcal{N}) = \min(\mathcal{N}')$; hence:

$$\frac{\begin{array}{l} e = (t, w, true) \\ \bar{e} = \min(\mathcal{N}) \quad (\sigma, t) \models Guard(w) \\ \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \\ \mathcal{N}_1 = \{(t, w) \mid (t, w) \in \mathcal{B}, (\sigma', t) \models Guard(w)\} \end{array}}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \epsilon, \mathcal{B} \setminus \{\bar{e}\}, (\mathcal{N} \cup \mathcal{N}_1) \setminus \{\bar{e}\})}$$

Because the sets $\mathcal{N}_1, \mathcal{N}_2, \mathcal{N}_3$ are constructed in exactly the same way as Base case 2a, the same argument also applies in this case. However, we additionally need to show that if $e^* \in \mathcal{N} \setminus \mathcal{N}'$, and $pred(e) = pred(e^*)$, then $(\sigma', thread(e^*)) \not\models pred(e^*)$. Since $e \in \mathcal{N}'$ and $pred(e^*) = pred(e)$ and $e < e^*$, we have $Invalidate(e, thread(e^*), pred(e))$. By definition of the $Invalidate$ predicate, this means that $(\sigma', thread(e^*)) \not\models pred(e^*)$.

- **Inductive step:** If trace τ contains multiple events, then we have:

$$\frac{(\sigma', e, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma_2, \epsilon, \mathcal{B}_2, \mathcal{N}_2)}{(\sigma', e :: \tau, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma_2, \tau, \mathcal{B}_2, \mathcal{N}_2)}$$

and

$$\frac{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma_1, \epsilon, \mathcal{B}_1, \mathcal{N}_1)}{(\sigma, e :: \tau, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma_1, \tau, \mathcal{B}_1, \mathcal{N}_1)}$$

Using the inductive hypothesis and the assumption that $(\sigma', \mathcal{B}', \mathcal{N}') \sim (\sigma, \mathcal{B}, \mathcal{N})$, we have $(\sigma_2, \mathcal{B}_2, \mathcal{N}_2) \sim (\sigma_1, \mathcal{B}_1, \mathcal{N}_1)$.

□

Lemma 10.9. *If $(\sigma, \mathcal{B}, \mathcal{N}) \sim (\sigma, \mathcal{B}', \mathcal{N}')$, τ is normalized, and*

$$(\sigma, \tau, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \tau', \mathcal{B}_1, \mathcal{N}_1)$$

then we have

$$(\sigma, \tau, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma', \tau', \mathcal{B}_2, \mathcal{N}_2)$$

and $(\sigma', \mathcal{B}_1, \mathcal{N}_1) \sim (\sigma', \mathcal{B}_2, \mathcal{N}_2)$.

Proof. The proof is by induction on τ .

- **Base case 1a:** $\tau = e = (t, w, false)$ and $(t, w) \notin \mathcal{B}$. In this case, we have $(\sigma, t) \not\models Guard(w)$ and $(\sigma, e, \mathcal{B}', \mathcal{N}') \longrightarrow (\sigma, \epsilon, \mathcal{B}' \cup \{\bar{e}\}, \mathcal{N}')$. Since $\mathcal{B} = \mathcal{B}'$, we have $(t, w) \notin \mathcal{B}'$. Thus, we also have $(\sigma, e, \mathcal{B}', \mathcal{N}) \Longrightarrow (\sigma, \epsilon, \mathcal{B}' \cup \{\bar{e}\}, \mathcal{N})$. Since $\mathcal{N}, \mathcal{N}'$ satisfy conditions (2) and (3) of the agreement definition and the state σ has not changed, $\mathcal{N}, \mathcal{N}'$ continue to satisfy conditions (2) and (3). Thus, we have $(\sigma, \mathcal{B}' \cup \{\bar{e}\}, \mathcal{N}') \sim (\sigma, \mathcal{B} \cup \{\bar{e}\}, \mathcal{N})$.
- **Base case 1b:** $\tau = e = (t, w, false)$ and $(t, w) \in \mathcal{N}$. We do not need to consider this case because it contradicts the assumption that the trace is normalized.

- **Base case 2a:** $\tau = e = (t, w, true)$ and $\bar{e} \notin \mathcal{B}$. In this case, we have:

$$\frac{\begin{array}{l} e = (t, w, true) \\ \bar{e} \notin \mathcal{B} \quad (\sigma, t) \models Guard(w) \\ \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \\ \mathcal{N}_1 = \{(t, w) \mid (t, w) \in \mathcal{B}, (\sigma', t) \models Guard(w)\} \end{array}}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \epsilon, \mathcal{B}, \mathcal{N} \cup \mathcal{N}_1)}$$

Since $\mathcal{B} = \mathcal{B}'$, this implies $\bar{e} \notin \mathcal{B}'$, thus, we also have:

$$\frac{\begin{array}{l} e = (t, w, true) \\ \bar{e} \notin \mathcal{B}' \quad (\sigma, t) \models Guard(w) \\ \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \\ \mathcal{N}_2 = GetSignals(w, \sigma', \mathcal{B}) \\ \mathcal{N}_3 = GetBroadcasts(w, \sigma', \mathcal{B}) \end{array}}{(\sigma, e, \mathcal{B}', \mathcal{N}') \Longrightarrow (\sigma', \epsilon, \mathcal{B}', \mathcal{N}' \cup \mathcal{N}_2 \cup \mathcal{N}_3)}$$

Thus, the proof is exactly the same as Base Case 2a of Lemma 10.8

- **Base case 2b:** $\tau = e = (t, w, true)$ and $\bar{e} = \min(\mathcal{N})$. In this case, we have:

$$\frac{\begin{array}{l} e = (t, w, true) \\ \bar{e} = \min(\mathcal{N}) \quad (\sigma, t) \models Guard(w) \\ \langle Body(w), t, \sigma \rangle \Downarrow \sigma' \\ \mathcal{N}_1 = \{(t, w) \mid (t, w) \in \mathcal{B}, (\sigma', t) \models Guard(w)\} \end{array}}{(\sigma, e, \mathcal{B}, \mathcal{N}) \longrightarrow (\sigma', \epsilon, \mathcal{B} \setminus \{\bar{e}\}, (\mathcal{N} \cup \mathcal{N}_1) \setminus \{\bar{e}\})}$$

The existence of e in \mathcal{N} does not guarantee the existence of e in \mathcal{N}' . We consider two cases, namely (1) $e \in \mathcal{N}'$ and (2) $e \notin \mathcal{N}'$. For case (1), the explicit-signal transitions also use rule (2b), thus the proof is exactly the same as Case 2b of the proof of Lemma 10.8.

We will now argue that case 2 (i.e., $e \notin \mathcal{N}'$) is not possible. Suppose $e \in \mathcal{N}$, but not in \mathcal{N}' . Then, by part (c) of the agreement definition, we have either (i) $(\sigma, t) \not\models \text{Guard}(w)$, or $\exists e' \in \mathcal{N}'$. $e' < e \wedge \text{pred}(e) = \text{pred}(e') \wedge \text{Invalidate}(e', t, \text{pred}(e))$. Observe that (i) contradicts the assumption $(\sigma, t) \models \text{Guard}(w)$. For (ii), we have $\exists e' \in \mathcal{N}'$. $e' < e \wedge \text{pred}(e) = \text{pred}(e') \wedge \text{Invalidate}(e', t, \text{pred}(e))$. But since $\mathcal{N} \supseteq \mathcal{N}'$, e' is also in \mathcal{N} and furthermore $e' < e$. But this contradicts the assumption that e' is the minimum element of \mathcal{N} .

- **Inductive step:** Same as the inductive step of the proof of Lemma 10.8.

□

Appendix C: Proof of Theorem 4.3

The proof is by induction on the length of τ_0 . For the base case, we have $\tau_0 = \epsilon$, thus $\tau' = \tau$ and the statement trivially holds. For the inductive step, we have $\tau_0 = e'\tau_1$. Suppose $(\sigma, e'\tau_1e, \mathcal{B}, \mathcal{N}) \xrightarrow{*} (\sigma, \epsilon, \mathcal{B}', \mathcal{N}')$ where $(\sigma, e', \mathcal{B}, \mathcal{N}) \xrightarrow{*} (\sigma_1, \epsilon, \mathcal{B}_1, \mathcal{N}_1)$ as well as $(\sigma_1, \tau_1e, \mathcal{B}_1, \mathcal{N}_1) \xrightarrow{*} (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')$. By the inductive hypothesis, we have $(\sigma_1, e\tau_1, \mathcal{B}_1, \mathcal{N}_1) \xrightarrow{*} (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')$. This implies $(\sigma, e'e\tau_1, \mathcal{B}, \mathcal{N}) \xrightarrow{*} (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')$. Finally, from the assumption $\text{Comm}(\text{Body}(e), M)$, we have $\text{Body}(e); \text{Body}(e') \equiv \text{Body}(e'); \text{Body}(e)$. Thus, we also have $(\sigma, ee'\tau_1, \mathcal{B}, \mathcal{N}) \xrightarrow{*} (\sigma', \epsilon, \mathcal{B}', \mathcal{N}')$

Appendix D: Sample Monitor Invariant

To give the reader some idea about the inferred monitor invariants, we show the invariant (in SMTLIB format) for the AsyncDispatch benchmark from Gradle:

```
(let ((a!1 (not (= (queue.size)
0))))
  (a!3 (not (>= (queue.size)
(+ 1
(maxQueueSize)))))
  (a!5 (not (>= (queue.size)
(maxQueueSize)))))
(let ((a!2 (not (or (= (state)
Stopped)
a!1))))
  (a!4 (or (= (queue.size)
0)
(= (queue.first)
0)
(= (state)
Stopped)
a!3))))
(let ((a!6 (and a!4
(or a!1
(= (state)
Stopped)
a!5))))
  (and (or a!2 a!6)
(not (<= (maxQueueSize)
0))
(>= (maxQueueSize)
0))
```

```
(>= (queue.size)
0))))))
```

As another example, we show the monitor invariant for the BoundedBuffer benchmark:

```
(let ((a!1 (not (>= buff.length
0)))
(a!2 (not (>= (count)
buff.length)))
(a!4 (>= (count)
(+ 1
(buff.length)))))
(let ((a!3 (or a!1
(not (or a!1 a!2))
(not (<= (count)
(- 1))))))
  (and a!3
(>= (count)
0)
(or a!1
(<= (count)
0)
(not a!4))))))
```