

Type Safe Nondeterminism – A Formal Semantics of Java Threads

Andreas Lochbihler *

Universität Passau
lochbihl@fim.uni-passau.de

Abstract

We present a generic framework to transform a single-threaded operational semantics into a semantics with interleaved execution of threads. Threads can be dynamically created and use locks for synchronisation. They can suspend themselves, be notified by other threads again, and interact via shared memory. We formalised this in the proof assistant Isabelle/HOL along with theorems to carry type safety proofs for the instantiating semantics (progress and preservation in the style of Wright and Felleisen [24]) over to the multithreaded case, thereby investigating the role of deadlocks and giving an explicit formalisation for them. We apply this framework to the Java thread model using an extension of the Jinja [12] source code semantics to have type safety for multithreaded Java machine-checked. The Java Memory Model is not included.

Categories and Subject Descriptors D.3.1 [Formal Definitions and Theory]: Semantics – type safety; F.3.2 [Semantics of Programming Languages]: Operational semantics; D.3.3 [Language Constructs and Features]: Concurrent programming structures

General Terms Languages, Theory

Keywords Threads, Deadlock, Type Safety, Java

1. Introduction

There are many formal approaches based on operational semantics and type systems to model certain safety properties of the Java programming language in the literature. A lot of work has been done in the context of type safety for sequential Java [5, 12, 19, 21, 22] on the source code level. In [19], Stärk et al. also give an abstract state machine semantics for threads, for which they show preservation of a number of invariants. Also, there is a large number of formal semantics for subsets of Java bytecode, see [3, Ch. 2] for an overview. All of them are single-threaded or lack other important features of concurrent Java, except for the formalisations by Belbidia and Debabbi [3] and the one by Liu and Moore [15], which contain a pretty comprehensive semantics of Java bytecode features. However, in both [3] and [15], only the semantics is given, but neither type system nor other safety features.

*This work was supported by DFG grant Sn11/10-1.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

FOOL '08 13 January, San Francisco, California, USA.
Copyright © 2008 ACM ... \$5.00

To the best of our knowledge, there is no formal type safety proof for multithreaded Java with dynamic thread creation and synchronisation, for neither source nor byte code, although Java's concurrency features are widely used in practise. Thus, the primary objective of this paper is a machine-checked type safety proof for a large and faithful subset of multithreaded Java, including classes with objects, fields and methods, inheritance with method overriding and dynamic dispatch, arrays, exception handling, dynamic thread creation, synchronisation via monitors and the wait/notify mechanism. To separate the semantics of these mechanisms from the low-level details of allowed compiler/runtime optimisations, we abstract from the complicated Java Memory Model (JMM) [9] by using a single shared memory like in [19]. Every properly synchronised program has the same semantics with or without the JMM.

Applying a standard approach by Wright and Felleisen [24] to show type safety, we experience the following in the multithreaded case: While the subject reduction theorem, i.e. preservation of well-typedness under reductions, is usually easily carried over from threads in isolation, the progress theorem, which shows that the evaluation of well-typed expressions does not get stuck, is particularly nontrivial in the presence of potential deadlocks and non-deterministic executions. In fact, most formalisations of type soundness for concurrent programming languages (e.g. [11, 13, 19, 20]) leave out the progress theorem or their notion of deadlock is given implicitly by the theorem's assumptions. This way, one cannot be sure that the theorem's notion coincides with the standard understanding of deadlock, especially because deadlocks can arise in many different ways (cf. Sec. 1.1 below). Thus, another major contribution of this paper is an explicit formalisation of deadlock in a theorem prover. Moreover, we then prove type safety with respect to this notion.

The basis of our work is the Jinja project [12], which already contains most features of the Java language subset that are not related to threads. On the source code level, Jinja consists of both a big-step and a small-step semantics, which are shown to be equivalent, a type system and a type safety proof through progress and preservation in the style of [24]. We add support for both threads and arrays, but since the latter is pretty orthogonal to the development here, we do not show this feature in detail. On the way, we develop a generic formal framework for lifting a sequential operational semantics to the concurrent case, which gives modularity in proving progress and preservation with respect to the sequential and concurrent aspects. We then instantiate this framework with the modified small-step semantics of Jinja to carry both progress and preservation over to the multithreaded case, and type safety then easily follows. In modelling threads, we closely follow Ch. 17 in the Java Language Specification (JLS) [9] for Java 5. Our semantics faithfully covers arbitrary dynamic thread creation, synchronisation on monitors and the wait/notify mechanism. Interaction takes place via shared memory. We have not included the deprecated methods `stop`, `suspend` and `resume`, neither `yield` and `sleep`, which are just hints to schedulers from which we abstract. Neither have we in-

Thread (I)	Thread (II)	Thread (III)
synchronized (f) {	synchronized (e) {	synchronized (g) {
synchronized (g) {	synchronized (f) {	synchronized (e) {
...
g.wait();	.	g.notify();
...
}	}	}
}	}	}

Figure 1. Three Java threads with different deadlock possibilities.

cluded `interrupt` and `join`, but we believe that they can be added easily. We have formalised every notion and proved every lemma and theorem in this paper in the theorem prover Isabelle/HOL [17], i.e. every single proof is machine-checked. Even though the presentation is quite technical at times, we show only the most important definitions, lemmata and theorems. The framework formalisation consists of about 5000 lines of definitions and proof scripts with approximately 250 lemmata. Multithreading for Jinja added about another 5000 lines to the Jinja source code formalisation. The current development of the formalisation is available online [16].

This contribution is organised as follows: The rest of this section gives an example on deadlocks and introduces some notation. The framework semantics is formalised in detail in Sec. 2. Sec. 3 contains the formal treatment of deadlocks. The machinery for lifting well-formedness constraints from the instantiating semantics to the multithreaded level is outlined in Sec. 4. In Sec. 5, we extend the Jinja semantics, apply the framework to it and show type safety. Sec. 6 discusses related work and Sec. 7 concludes this paper, glancing at future directions.

1.1 Deadlock example

For an example on how different forms of deadlock can arise in Java, consider Fig. 1. For a start, suppose there are only threads (I) and (II), which share the object referenced by `f`. Thread (I) first locks the monitor for the object `f`, then acquires the lock on object `g`'s monitor, and later suspends itself to the wait set of `g`'s monitor, thereby releasing the lock on `g` again. Once it is woken up (by some notify invocation on `g` by some other thread),¹ it then competes for the lock on `g` again until it reacquires it and then finishes execution after releasing the locks on both `g` and `f`. Thread (II) has to acquire locks on objects `e` and `f` to execute its critical section. Now, if we start (I) first and then interleave threads (I) and (II), both end up in a deadlock: (I) is first to acquire the lock on `f`, so later it suspends itself to `g`'s wait set and waits for being woken up by some other thread, but there is only thread (II) which waits on the lock on `f`, which is still held by (I). Conversely, if (II) acquires the lock on `f` first, (I) will end up deadlocked in the wait set because there will be no other thread to wake up (I) again. Note that in this case, there is only a single thread in deadlock. This is not possible if deadlock is due to locks only.

Now, suppose thread (III) is run in parallel with (I) and (II) and shares objects `e` and `g`. (III) acquires locks on `g` and `e` and possibly wakes up thread (I). However, e.g., if thread (I) acquires its first lock and next, thread (II) acquires its first lock and then (III) does so, all threads end up in deadlock, because they are waiting cyclically on each other: (I) is waiting for the lock on `g`, which is held by (III), but (III) is waiting on the lock on `e`, which is held by (II), and (II) itself is waiting for the lock on `f` held by (I). There are many more deadlock possibilities in this example, but not all schedules lead to deadlock: If e.g. (II) acquires both locks first and (I) acquires the lock on `g` before (III) does, there is no possibility for a deadlock.

¹ We abstract from “spurious wake-ups” which are allowed by the JLS [9].

1.2 Notation

Both our framework and the Jinja semantics are formalised in the proof assistant Isabelle/HOL, i.e. all formulae and propositions are written in HOL, which is close to standard mathematical notation. This section introduces further non-standard notation, a few basic data types and their primitive operations.

Types is the set of all types which contains, in particular, the type of truth values `bool`, natural numbers `nat`, and integers `int`. The space of total functions is denoted by \Rightarrow . Type variables are written `'a`, `'b` etc. The notation $t :: \tau$ means that the HOL term t has type τ .

Pairs come with two projection function $fst :: 'a \times 'b \Rightarrow 'a$ and $snd :: 'a \times 'b \Rightarrow 'b$. We identify tuples with pairs nested to the right: (a, b, c) is identical to $(a, (b, c))$ and $'a \times 'b \times 'c$ to $'a \times ('b \times 'c)$.

Sets (type `'a set`) follow the usual mathematical convention.

Lists (type `'a list`) come with the empty list `[]`, the infix constructor `·` for consing and the infix `@` that concatenates two lists. Variable names ending in “s” usually stand for lists or maps (see below) and $|xs|$ is the length of xs . If $n < |xs|$, then $xs[n]$ denotes the n -th element of xs . The standard function $map f xs$ applies the function f to every element in the list xs . $flatten xss$ flattens a list xss of lists by concatenating all its elements. $list-all2 P xs ys$ is shorthand for $|xs| = |ys| \wedge (\forall i < |xs|. P xs[i] ys[i])$.

Function update is defined as follows: Let $f :: 'a \Rightarrow 'b$, $a :: 'a$ and $b :: 'b$. Then $f(a := b) \equiv \lambda x. if x = a then b else f x$.

The **option** data type `'a option = None | Some 'a` adjoins a new element `None` to a type `'a`. All existing elements in type `'a` are also in `'a option`, but are prefixed by `Some`. For succinctness, we write $[a]$ for `Some a`. Hence, for example, `bool option` has the values `None`, `[True]` and `[False]`.

Case distinctions on data types are written with guards. For x of option type, e.g., $case x of None \Rightarrow f | [y] \Rightarrow g y$ means that if x is `None` then the expression is f , and if x is some y then the expression is $g y$ where g may refer to the value y of the `Some` constructor.

Partial functions are modelled as functions of type `'a \Rightarrow 'b option` where `None` represents undefined and $f x = Some y$ means x is mapped to y . Instead of `'a \Rightarrow 'b option` we write `'a \rightarrow 'b` and call such functions **maps**. We write $f(x \mapsto y)$ as shorthand for $f(x := [y])$. The map $\lambda x. None$ is written `empty` and `empty(...)`, where `...` are updates, abbreviates to `[...]`. For example, `empty(x \mapsto y)` becomes `[x \mapsto y]`.

The **indefinite description** $\epsilon x. Q x$ is known as Hilbert’s ϵ -operator. It denotes some x such that $Q x$ is true, provided one exists.

2. Formalisation of the framework semantics

Our framework consists of two major parts. On the one hand, there is the framework semantics `redT` which handles the threads for an arbitrary (single-threaded) semantics. It is instantiated with an operational small-step semantics r such that every single thread can be reduced according to r . `redT` takes care of all bookkeeping, e.g. managing locks and thread reductions. On the other hand, we give a number of well-formedness conditions, both for the initial program configuration and for the instantiating semantics, which enable us to easily carry type safety proofs for the instantiating semantics (when we view it as a single-threaded semantics) to the multithreaded setting of the framework. In this section, we only treat the framework semantics, well-formedness conditions are presented in Sec. 3.2 and 4.

Our framework operates on a small-step operational semantics which is modelled by a set of single-step reductions. When applying the framework, we always deal with two different small-step semantics. We do not use a big-step semantics because any such semantics would have to appeal to small-step semantics features to allow for interleaved evaluations of different threads.

2.1 The instantiating semantics

The semantics with which we instantiate the framework contains all single-step reductions for a single thread in isolation. We call such a semantics an **instantiating semantics** and usually denote it by r . Formally, a single-step reduction in rP , where P is a program environment parameter to the semantics, is a tuple $((e, c, x), tas, (e', c', x'))$, which we write as $P \vdash \langle e, (c, x) \rangle \text{---}tas\text{---}r \rightarrow \langle e', (c', x') \rangle$. Intuitively, this denotes that in r with program environment P , the expression e with common/shared memory c and local state x can reduce to the expression e' with new shared memory c' and local state x' . tas is a list of thread actions which tell the framework semantics to create a new thread, to acquire a lock, to wake another thread, etc. r has type $(\prime p, \prime e, \prime l, \prime t, \prime x, \prime c, \prime w)$ semantics, which is shorthand for $\prime p \Rightarrow ((\prime e \times \prime c \times \prime x) \times (\prime e, \prime l, \prime t, \prime x, \prime c, \prime w)$ thread-action list $\times \prime e \times \prime c \times \prime x$) set, where we assign the following meaning to type parameters:

- $\prime p$ Type of the program environment, e.g. Java class definitions
- $\prime e$ Type of a thread expression
- $\prime l$ Type of the locks threads acquire and release, e.g. Java monitors
- $\prime t$ Type of the thread IDs, by which threads are identified
- $\prime x$ Type of the state information that is local to a single thread, e.g. local variables
- $\prime c$ Type of the common/shared memory, e.g. type of the heap
- $\prime w$ Type of the identifiers for wait sets, e.g. monitors in Java

In what follows, to denote concrete values of these types we use the same letters (without the leading \prime , but with various decorations), e.g. t, t', t_1 , etc. for thread IDs.

2.2 The multithreaded semantics

The framework semantics takes an instantiating semantics r as a parameter and forms the set of single-step reductions for the multithreaded case, in which a **state** consists of four components:

1. The **lock status** in the first component, usually denoted by ls , stores in a map of type $\prime l \rightarrow \prime t \times nat$ (denoted by $(\prime l, \prime t)$ locks) for every lock l how many times it is held by a thread, if any. A thread (identified by its thread ID t) holding the lock l ($n + 1$) times is represented by $ls \ l = \lfloor (t, n) \rfloor$. If l is not held by any thread, then $ls \ l = None$. Using a map ensures that a lock is always held by at most one thread at a time.
2. The second component, usually denoted by es , stores the **thread information** in a map of type $\prime t \rightarrow \prime e \times \prime x$, denoted by $(\prime e, \prime t, \prime x)$ thread-info. For every ID t for which a thread is present, es maps t to the current expression of the thread and the thread-local state.
3. The third component of type $\prime c$ is the **shared memory**, which we usually denote with c .
4. The last component keeps track of the **wait sets**. Every thread can be in at most one wait set at a time. We therefore model them as a map ws of type $\prime t \rightarrow \prime w$ (denoted by $(\prime w, \prime t)$ waitsets) where a thread t is waiting in the wait set w iff $ws \ t = \lfloor w \rfloor$. t is ready for execution iff $ws \ t = None$.

Suppose, e.g., in Fig. 1 with threads (I) and (II) only, we are in the deadlock state as described in the first scenario in Sec. 1.1. Suppose e references an object at address e , and similarly for f and g . Then, this state is represented by the tuple (ls, es, c, ws) where $ls = \lfloor f \mapsto (I, 0), e \mapsto (II, 0) \rfloor$, i.e. the locks on f and e are held once by (I) and (II) resp., $es = \lfloor I \mapsto (sync(locked(f)) \{ sync(g) \{ \dots \} \}, x_1), II \mapsto (sync(locked(e)) \{ sync(f) \{ \dots \} \}, x_2) \rfloor$ stores the thread expressions and the local data for (I) and (II). The shared memory $c = \lfloor e \mapsto Obj \dots, f \mapsto Obj \dots, g \mapsto Obj \dots \rfloor$ contains the objects referenced by e, f , and g and (I) is in the wait set of g : $ws = \lfloor I \mapsto g \rfloor$

When a single thread expression is reduced in the instantiating semantics r , it can ask the framework to perform finitely many thread actions of type $(\prime e, \prime l, \prime t, \prime x, \prime c, \prime w)$ thread-action, which can alter the state of the locks, threads and wait sets. Note that these actions are the only means of “communication” between the two semantics. Since this is unidirectional, the framework semantics can only transfer information to the instantiating semantics by picking one reduction offered by the instantiating semantics. Hence, the instantiating semantics must anticipate in its reductions all possible answers it is willing to accept from the framework semantics. At the moment, our framework supports eight different **thread actions**, which can be split into three groups:

Locking *Lock* l acquires a lock on l for the current thread t . If l is held by another thread, this reduction is not possible. For *Unlock* l , which releases one lock on l , t must hold at least one lock on l . *UnlockFail* l is only possible if t does not hold a lock on l , i.e. releasing a lock on l would fail.

Thread creation *NewThread* $t \ e \ c \ x$ creates a new thread with ID t , initial expression e and local state x . c must be the equal to the common memory after the reduction step of the executing thread. The new thread is ready for execution and does not hold any locks. *NewThreadFail* tests whether all IDs of type $\prime t$ are assigned to a thread, i.e., normal thread creation would fail.

Wait sets *Suspend* w inserts the current thread t in the wait set w , any previous assignment of the current thread to a wait set is lost. *Notify* w non-deterministically wakes up one of the threads in the wait set w . If w is empty, no thread is woken up. *NotifyAll* w wakes up all threads in the wait set w .

Note that most actions with preconditions are paired with another one whose precondition is the negation of the former’s. A thread should always know, e.g., if it is able to release a lock that it ought to have acquired before – if not, i.e., the lock status is inconsistent with the thread’s state or the thread simply does not care whether it owns the locks it tries to release, the framework semantics can tell the instantiating semantics. The only exception is the *Lock* action, since a thread cannot decide on its own if it will be able to acquire a lock. We deliberately chose this asymmetry in the mutually-exclusive locks because they are the source of deadlocks in practice.

Similarly, *NewThread* $t \ e \ c \ x$ contains more information than immediately necessary, e.g. the shared memory c , which is saved for later use. The thread ID for the new thread is actually assigned by the framework semantics, which picks non-deterministically a fresh thread ID if available, but since some semantics might be interested to know the IDs of those threads they have created, we include this extra information in the action itself already. Hence, a sensible semantics should - when it creates a new thread - offer a reduction for every possible thread ID. Similarly, when the maximum number of threads is reached, the *NewThreadFail* action tells the semantics so.

The list tas of thread actions issued by a reduction step of the instantiating semantics tells the framework what changes to do to the multithreaded state. The framework then performs these changes using three update functions for single thread actions:

- $ls \overset{\ell}{\rightsquigarrow}_t ta$ updates the lock state ls according to ta for thread t ,
- $es \rightsquigarrow ta$ adds the new thread in ta to the thread information es ,
- $ws \bullet \rightsquigarrow_t ta$ changes the wait sets ws according to ta requested by thread t .

Their definitions are shown in Fig. 2, the cases in which the map is not updated have been omitted. The function *lock-lock* $ls \ t \ l$ (*unlock-lock* $ls \ t \ l$) increases (decreases) the lock count of t on l , *new-thread-id* es gives an unassigned thread ID if there is one left. There are also functions which fold the update functions over

$ls \overset{\ell}{\rightsquigarrow}_t \text{Lock } l$	$\equiv \text{lock-lock } ls \ t \ l$
$ls \overset{\ell}{\rightsquigarrow}_t \text{Unlock } l$	$\equiv \text{unlock-lock } ls \ t \ l$
$ls \overset{\ell}{\rightsquigarrow}_t \square$	$\equiv ls$
$ls \overset{\ell}{\rightsquigarrow}_t ta \text{-}tas$	$\equiv (ls \overset{\ell}{\rightsquigarrow}_t ta) \overset{\ell}{\rightsquigarrow}_t tas$
$es \rightsquigarrow \text{NewThread } t' \ e \ c \ x$	$\equiv \text{if new-thread-id } es = \lfloor t' \rfloor$ $\quad \text{then } es(t' \mapsto (e, x)) \text{ else arbitrary}$
$es \rightsquigarrow \text{NewThreadFail}$	$\equiv \text{if new-thread-id } es = \text{None}$ $\quad \text{then } es \text{ else arbitrary}$
$es \overset{[\rightsquigarrow]}{\square}$	$\equiv es$
$es \overset{[\rightsquigarrow]}{ta \text{-}tas}$	$\equiv (es \rightsquigarrow ta) \overset{[\rightsquigarrow]}{tas}$
$ws \bullet \rightsquigarrow_t \text{Notify } w$	$\equiv \text{if } \exists t. ws \ t = \lfloor w \rfloor$ $\quad \text{then let } t = \epsilon t. ws \ t = \lfloor w \rfloor$ $\quad \quad \text{in } ws(t := \text{None})$ $\quad \text{else } ws$
$ws \bullet \rightsquigarrow_t \text{NotifyAll } w$	$\equiv \lambda t. \text{if } ws \ t = \lfloor w \rfloor \text{ then None else } ws \ t$
$ws \bullet \rightsquigarrow_t \text{Suspend } w$	$\equiv ws(t \mapsto w)$
$ws \overset{[\bullet \rightsquigarrow]}{\square}$	$\equiv ws$
$ws \overset{[\bullet \rightsquigarrow]}{ta \text{-}tas}$	$\equiv (ws \bullet \rightsquigarrow_t ta) \overset{[\bullet \rightsquigarrow]}{tas}$
$(ls, es, ws) \overset{[\rightsquigarrow]}{ta}$	$\equiv (ls \overset{\ell}{\rightsquigarrow}_t ta, es \overset{[\rightsquigarrow]}{ta}, ws \overset{[\bullet \rightsquigarrow]}{ta})$

Figure 2. Update functions for locks, thread data and wait sets.

$ls \vdash_t \text{Lock } l \ \checkmark_\ell$	$\equiv \text{may-lock } ls \ t \ l$
$ls \vdash_t \text{Unlock } l \ \checkmark_\ell$	$\equiv \text{has-lock } ls \ t \ l$
$ls \vdash_t \text{UnlockFail } l \ \checkmark_\ell$	$\equiv \neg \text{has-lock } ls \ t \ l$
$ls \vdash_t \square \ \checkmark_\ell$	$\equiv \text{True}$
$ls \vdash_t ta \text{-}tas \ \checkmark_\ell$	$\equiv ls \vdash_t ta \ \checkmark_\ell \wedge (ls \overset{\ell}{\rightsquigarrow}_t ta) \vdash_t tas \ \checkmark_\ell$
$es, c \vdash \text{NewThread } t \ e \ c' \ x \ \checkmark_\gamma$	$\equiv \text{new-thread-id } es = \lfloor t \rfloor \wedge c = c'$
$es, c \vdash \text{NewThreadFail} \ \checkmark_\gamma$	$\equiv \text{new-thread-id } es = \text{None}$
$es, c \vdash \square \ \checkmark_\gamma$	$\equiv \text{True}$
$es, c \vdash ta \text{-}tas \ \checkmark_\gamma$	$\equiv es, c \vdash ta \ \checkmark_\gamma \wedge (es \rightsquigarrow ta), c \vdash tas \ \checkmark_\gamma$

Figure 3. Predicates on the state satisfying the preconditions of a thread action list

lists of thread actions tas , denoted by $ls \overset{\ell}{\rightsquigarrow}_t tas$, $es \overset{[\rightsquigarrow]}{tas}$, and $ws \overset{[\bullet \rightsquigarrow]}{tas}$ resp. They are combined in one state update function $(ls, es, ws) \overset{[\rightsquigarrow]}{tas}$ for thread actions tas . Let us look briefly at wake-ups: The equation for $\text{Notify } w$ says that if there is a thread waiting in w then pick any such thread, say t , and remove it from w , otherwise do nothing. Similarly, $\text{NotifyAll } w$ removes all threads from w and leaves all other threads unchanged.

Before updating the multithreaded state (ls, es, c, ws) , we first have to check whether all changes in tas requested by the thread t can actually be granted to it. For this purpose, we introduce the predicates $ls \vdash_t ta \ \checkmark_\ell$ and $es, c \vdash ta \ \checkmark_\gamma$ for an action ta . Again, they are paired with list versions $ls \vdash_t tas \ \checkmark_\ell$ and $es, c \vdash tas \ \checkmark_\gamma$, respectively, which work through the list of thread actions, temporarily updating the state after each successful check. We show their definitions in Fig. 3, leaving out the cases in which the predicates are constantly true. $\text{may-lock } ls \ t \ l$ checks that no thread other than t holds a lock on l in ls , $\text{has-lock } ls \ t \ l$ that t does so. Note that we do not have to check for conditions on the wait set because wait set actions have no preconditions.

If a list of thread actions contains at least one action whose preconditions are not met by the current state at that position, then the reduction step cannot be executed that time. This is a powerful means to an instantiating semantics for checking that

some preconditions on the locks hold, even without changes to the lock status. It can check e.g. whether it holds (does not hold) the lock l by using the list $[\text{Unlock } l, \text{Lock } l]$ ($[\text{UnlockFail } l]$), without altering the locks. Note that the order of thread actions in the list can be important: $[\text{Lock } l, \text{Unlock } l]$, e.g., consists of the same actions, but checks that no other thread is holding a lock on l . In the same way, $[\text{Unlock } l, \text{Unlock } l, \text{UnlockFail } l, \text{Lock } l]$ releases one lock on l , but also tests that it has had exactly two before the reduction.

In case the list contains multiple *Suspend* actions, the thread will be listed only in the wait set specified by the last such action after the reduction. It is not forbidden to have a *Suspend* w action followed by some *Notify* w or *NotifyAll* w action later in the same list. Note that this is the only case in which a thread can possibly wake up itself, because whenever a thread is in a wait set, this thread is not considered for execution by the framework.

Now, we put together everything we have so far to obtain the set of reductions in the framework semantics redT for the instantiating semantics r in the program environment P . A single reduction step in $\text{redT } r \ P$ is denoted by $P, r \vdash \langle ls|es, c|ws \rangle -t, tas \rightarrow \langle ls'|es', c'|ws' \rangle$, which means that in the multithreaded state (ls, es, c, ws) , thread t can reduce with thread actions tas which yields the multithreaded state (ls', es', c', ws') . The only reduction rule for the **framework semantics** then reads:

$$\frac{P \vdash \langle e, (c, x) \rangle -tas \rightarrow \langle e', (c', x') \rangle \quad \begin{array}{l} es \ t = \lfloor (e, x) \rfloor \quad ws \ t = \text{None} \\ ls \vdash_t tas \ \checkmark_\ell \end{array}}{P, r \vdash \langle ls|es, c|ws \rangle -t, tas \rightarrow \langle ls'|es', c'|ws' \rangle} \quad \begin{array}{l} (ls', es', ws') = (ls, es(t \mapsto (e', x')), ws) \overset{[\rightsquigarrow]}{tas} \\ P, r \vdash \langle ls|es, c|ws \rangle -t, tas \rightarrow \langle ls'|es', c'|ws' \rangle \end{array}$$

Intuitively, $P, r \vdash \langle ls|es, c|ws \rangle -t, tas \rightarrow \langle ls'|es', c'|ws' \rangle$ is a reduction in $\text{redT } r \ P$ iff

- there is a thread t in es , say $es \ t = \lfloor (e, x) \rfloor$,
- which is not in a wait set in ws and
- which can do a reduction $P \vdash \langle e, (c, x) \rangle -tas \rightarrow \langle e', (c', x') \rangle$ in r with thread actions tas such that lock and thread creation conditions are met by the state (ls, es, c, ws) and
- the old state – where the thread information es gets updated with the single-thread reduction result (e', x') – is updated according to tas and combined with the new shared memory c' to yield (ls', es', c', ws')

We write reductions in the transitive and reflexive closure of $P, r \vdash \langle \cdot | \cdot | \cdot \rangle -\cdot, \cdot \rightarrow \langle \cdot | \cdot | \cdot \rangle$ as $P, r \vdash \langle ls|es, c|ws \rangle -tas \rightarrow^* \langle ls'|es', c'|ws' \rangle$. tas now is a list of pairs of thread IDs and thread action lists and keeps track of the thread ID and thread actions for every reduction step, i.e., in step n , thread $fst \ tas_{[n]}$ reduces with the actions $snd \ tas_{[n]}$. Note that we do not model a specific scheduler here, i.e. a reschedule is possible between any two single-step reductions. However, restricting to a specific scheduling scheme can be done easily by selecting the appropriate reduction tuples from the framework semantics, e.g. based on the data given in tas .

Similarly, every reduction step of the instantiating semantics is considered to be atomic. Conversely, there is no direct means to force multiple reduction steps of a thread being atomic. However, an instantiating semantics can easily use the locking mechanism to ensure atomicity, if necessary: Replace, e.g., the lock type l by $l \ \text{option}$ to introduce one extra lock. Then, every atomic reduction sequence should start (end) with a *Lock None* (*Unlock None*) action. If all single-step reductions, which are supposed to be atomic themselves, also prefix and postfix their thread actions by *Lock None* and *Unlock None* respectively, every thread can only be reduced if it can lock *None*. Since atomic sequences hold *None* in their intermediate steps (and locks are mutually exclusive), rescheduling is not possible inside atomic actions.

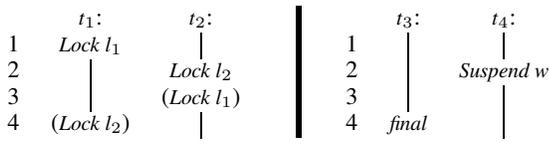


Figure 4. Two thread schedules that lead to deadlocked threads.

3. Type safety with deadlocks

In the sequential case, type safety is usually shown in the syntactic way [24] by showing progress and preservation. Progress means that every well-formed and well-typed expression that is not final can be reduced, preservation requires that well-typedness is preserved under reductions and the expression’s type may become only more specific. Usually, preservation can easily be carried over from the instantiating semantics. Deadlocks, however, can break the progress property.

3.1 Formalising deadlocks

A thread is said to be in deadlock if it is waiting for an event that will never occur. In operating systems, deadlock of processes has four preconditions on resource usage: mutual exclusion, hold and wait, circular waiting, and no preemption (cf. e.g. [18]). In our setting, we have two different possibilities for deadlocks:

- A thread that is in a wait set is deadlocked if all threads are either in a wait set or have completed its execution or are already deadlocked themselves.
- A thread is deadlocked if it is waiting to obtain a lock which is held by another thread which itself is deadlocked, i.e. threads are waiting circularly on each other.

We assume that the function *final* tells whether a given thread’s expression has completed its execution and is not supposed to be reduced any further. An instantiating semantics must define *final* appropriately. The predicate *wf-final* then ensures that final expressions indeed cannot be reduced.

Consider, e.g., the schedules shown in Fig. 4. On the left-hand side, thread t_1 acquires the lock l_1 , then thread t_2 acquires the lock l_2 . To proceed, t_2 needs the lock l_1 , too, so the *Lock* l_1 is postponed. However, t_1 requests the lock l_2 which is held by t_2 , hence both threads are deadlocked. On the right-hand side, we see an example with wait sets. Suppose there are only two threads, t_3 and t_4 . Some time while being reduced, t_4 suspends itself to the wait set w . Once t_3 has reduced to a final value, t_4 is deadlocked because there are no more threads to wake it up again.

In our framework, things are a little more tricky because threads can atomically request and release any number of different locks at one reduction step, and the instantiating semantics need not be deterministic, i.e. a thread usually can reduce in many different ways requesting many different locks. To get a hold on this, we introduce two abstractions in our formalisation:

- $P, r \vdash \langle e, (c, x) \rangle L \lambda$ denotes that in the semantics r with program P , the expression e can reduce in the state (c, x) with a thread action list which contains at least one *Lock* l action for every lock l in L and in which all *Lock* actions are on locks in L .
- $P, r \vdash \langle e, (c, x) \rangle \lambda$ denotes that in the semantics r with program P , the thread action list of every possible reduction of e in state (c, x) contains at least one *Lock* action.

Note that we do not care about unlock actions because only a thread itself would be able to remedy the missing lock, not others. With these two abstractions, we can now define the set of **threads in deadlock** $deadlocked\ r\ P\ l_s\ e_s\ w_s\ c$ as a co-inductive set. Fig. 5 shows the introduction rules, where the predicate $final(es \setminus M)$ checks

$$\begin{array}{c}
 es\ t = \llbracket (e, x) \rrbracket \quad P, r \vdash \langle e, (c, x) \rangle \lambda \quad P, r \vdash \langle e, (c, x) \rangle L \lambda \\
 \forall L. P, r \vdash \langle e, (c, x) \rangle L \lambda \longrightarrow \\
 \exists t' \in deadlocked\ r\ P\ l_s\ e_s\ w_s\ c. t' \neq t \wedge (\exists l \in L. has-lock\ l_s\ t'\ l) \\
 \hline
 t \in deadlocked\ r\ P\ l_s\ e_s\ w_s\ c \\
 \\
 es\ t = \llbracket (e, x) \rrbracket \quad final(es \setminus deadlocked\ r\ P\ l_s\ e_s\ w_s\ c) \quad ws\ t = \llbracket w \rrbracket \\
 \hline
 t \in deadlocked\ r\ P\ l_s\ e_s\ w_s\ c
 \end{array}$$

Figure 5. Introduction rules for the coinductive set $deadlocked\ r\ P\ l_s\ e_s\ w_s\ c$

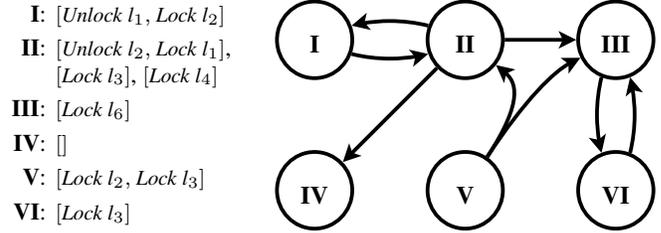


Figure 6. Deadlock example.

whether each thread of es is either final or its thread ID is in M . The first rule ensures that those threads t are in deadlock which must obtain a lock and can be reduced and for every reduction of which there must be another thread t' in deadlock which holds one of the locks requested by that reduction. The second rule says that if a thread is in a wait set with all other non-final threads being either also in a wait set or deadlocked, then it is deadlocked. This is the case since there is no rule that eliminates that thread from the coinductive set, which is in fact a greatest fixpoint. We consider the set of threads to be closed w.r.t. the outside, i.e. we cannot add a spinning thread to the system, which would “undeadlock” again waiting threads. As a minor technical detail, we consider a final thread which has suspended itself in its last reduction to a wait set to be deadlocked, too.

Consider Fig. 6 for an example of different deadlock situations. Suppose there are six threads which at the moment can reduce with the thread action lists shown on the left-hand side. If there are multiple lists for one thread, then there is one reduction for each list. Suppose that no thread is in a wait set and that the i -th thread holds the lock l_i . The graph on the right-hand side shows which thread is waiting to obtain a lock held by another thread. Then, thread III and VI are waiting on each other without other reduction options. Clearly, both of them are deadlocked. Although I and II are also waiting on each other, they are not deadlocked at the moment: II has two more reduction options. Waiting on lock l_3 will be in vain, because III is deadlocked. However, IV is not waiting on anyone, hence II may still hope to obtain the lock l_4 some time. Since thread action lists must be executed atomically, we may not appropriately interleave the action lists of I and II. Note that V is waiting *simultaneously* on II and III. Since III is already in deadlock, so is V. Clearly, IV is not deadlocked, since the empty list is always possible. Now, suppose thread IV is in a wait set. Then, all threads are deadlocked, since every thread except IV is waiting on some other thread to release a lock, and the only thread that could be reduced (i.e. IV) is waiting for some other thread waking it up.

For the framework semantics working properly, we need two basic well-formedness conditions. First, we require (via the predicate $es \vdash_f l_s \checkmark$) that all locks are held by actual, non-final threads:

$$\begin{array}{l}
 es \vdash_f l_s \checkmark \equiv \forall l. \text{case } l_s\ l \text{ of } None \Rightarrow True \\
 \quad \mid \llbracket (t, n) \rrbracket \Rightarrow \exists e\ x. es\ t = \llbracket (e, x) \rrbracket \wedge \neg final\ e
 \end{array}$$

$$\frac{t' \notin \text{deadlocked } r P l s e s w s c \quad e s t' = [(e, x)] \quad \neg \text{final } e \quad e s \vdash_f l s \checkmark}{\exists t e x. e s t = [(e, x)] \wedge \neg \text{final } e \wedge w s t = \text{None} \wedge (P, r \vdash \langle e, (c, x) \rangle \wr) \wedge (\exists L. P, r \vdash \langle e, (c, x) \rangle L \wr) \longrightarrow (\exists L. P, r \vdash \langle e, (c, x) \rangle L \wr \wedge (\forall l \in L. \text{may-lock } l s t l))}$$

Figure 7. Formal restatement of Lem. 3

Second, let $wf\text{-lock } r P c c'$ denote that both $P, r \vdash \langle e, (c, x) \rangle L \wr$ and $P, r \vdash \langle e, (c, x) \rangle \wr$ are invariant if c is changed (via reduction by another thread) to c' . Then, we show the following properties for the set of deadlocked threads. The first lemma is shown by case analysis on t being deadlocked. Lem. 2 and 3 are shown by coinduction on the set of deadlocked threads.

LEMMA 1. *Deadlocked threads are irreducible in the framework semantics. Formally:*

$$\frac{t \in \text{deadlocked } r P l s e s w s c}{\neg P, r \vdash \langle l s | e s, c | w s \rangle -t, t a s \rightarrow \langle l s' | e s', c' | w s' \rangle}$$

LEMMA 2. *Deadlocked threads remain deadlocked after reductions of other threads.*

$$\frac{\frac{P, r \vdash \langle l s | e s, c | w s \rangle -t, t a s \rightarrow \langle l s' | e s', c' | w s' \rangle \quad wf\text{-final } r \quad wf\text{-lock } r P c c'}{\text{deadlocked } r P l s e s w s c \subseteq \text{deadlocked } r P l s' e s' w s' c'}}{}$$

LEMMA 3 (cf. Fig. 7). *If there is a non-final thread which is not deadlocked – and all locks are held by non-final threads – then there is a thread which is not final, not in a wait set and if it can be reduced and must obtain a lock when being reduced then there is a reduction whose Lock actions can all be performed.*

3.2 Well-formedness conditions for progress

To show progress in the multithreaded case, we naturally need some progress lemma for the instantiating semantics. We write $wf\text{-progress } r P e s c$ to denote that all non-final threads in $e s$ can be reduced in r if the shared memory is c . However, this is not sufficient for progress in the multithreaded case. Suppose, e.g., there is only a single thread t which tries to release a lock l by the $Unlock\ l$ action although it is not holding l . If this was the only reduction possible for t , then this would break the progress property of the framework semantics, because the expression in the instantiating semantics is not stuck, but its embedding as a thread in the multithreaded semantics is. Therefore, we introduce the predicate $ex\text{-red } r P l s e s c$ which requires that for every thread t in $e s$ which is reducible in r with thread actions $t a s$ there is also a reduction in r with thread actions $t a s'$ where

- all actions in $t a s'$ that are related to thread creation are possible,
- for every lock l , all (un)lock actions on l in $t a s'$ are possible or the first one not possible is a $Lock\ l$ action, and
- for every $Lock\ l$ action in $t a s'$ that is not preceded by a $Unlock\ l$ action in $t a s'$ there must already be a $Lock\ l$ action in $t a s$.

$ex\text{-red } r P l s e s c$ ensures that if a thread t in $e s$ can be reduced in r , then either t is deadlocked in $(l s, e s, c, w s)$ or there exists a thread in $e s$ and a reduction of it in r that is possible in $(l s, e s, c, w s)$. However, it is frequently more convenient to show for every thread ID t a different kind of well-formedness denoted by $wf\text{-r-progress } r P l s e s c t$, which implies $ex\text{-red}$ (cf. Lem. 4). Intuitively, $wf\text{-r-progress } r P l s e s c t$ imposes six constraints provided that there is a thread t in $e s$, say $e s t = [(e, x)]$:

- The length of the thread action list for reductions of e in state (c, x) must be bounded.
- For all reductions of t with actions $t a s$, all $NewThread$ actions must have the same shared memory as the resulting state.

- For every reduction of t which creates a new thread with ID t' , there must also be a reduction creating a thread with any other thread ID and, if all thread IDs are used up, another one which has a $NewThreadFail$ action instead.
- Conversely, if not all thread IDs are used, for every $NewThreadFail$ action, there must be a reduction actually creating a new thread with an arbitrary ID.
- If there is a reduction of t with actions $t a s$ which contain an $Unlock\ l$ action, there must be a reduction with an $UnlockFail\ l$ action instead, if the $Unlock\ l$ action is not possible.
- Conversely, if $UnlockFail\ l$ is requested by some reduction of t , but t holds sufficiently many locks on l , then there is also a reduction with an $Unlock\ l$ action instead.

In fact, both $ex\text{-red } r P l s e s c$ and $wf\text{-r-progress } r P l s e s c t$ are more complex because we allow for partially rearranging of thread actions in the substitute reductions, but this technical detail is not important for the further development.

The next lemma shows that $wf\text{-r-progress } r P l s e s c t$ for every thread t is indeed sufficient for $ex\text{-red } r P l s e s c$. It is shown by induction on prefixes of thread action lists and a large case distinction.

LEMMA 4. If $\forall t. wf\text{-r-progress } r P l s e s c t$ then $ex\text{-red } r P l s e s c$.

With Lem. 1, 2, and 3, we finally show the following key theorem in proving progress for the instantiating semantics:

THEOREM 1 (Progress). *Let t be a non-final thread which is not deadlocked. Suppose r satisfies the progress conditions $wf\text{-progress}$ and $ex\text{-red } r P l s e s c$. Then, if all locks are held by non-final threads, the multithreaded semantics can make progress. Formally:*

$$\frac{e s t = [(e, x)] \quad \neg \text{final } e \quad t \notin \text{deadlocked } r P l s e s w s c \quad wf\text{-progress } r P e s c \quad e s \vdash_f l s \checkmark \quad ex\text{-red } r P l s e s c}{\exists t' t a s' e s' l s' w s' c'. P, r \vdash \langle l s | e s, c | w s \rangle -t', t a s' \rightarrow \langle l s' | e s', c' | w s' \rangle}$$

Note that our deadlock formalisation is independent of any type system at all. Although we do not require progress for deadlocked threads, this therefore does not open up possible holes in the type system. Moreover, for a thread in isolation, we still require the progress property, independent of being deadlocked.

4. Lifting thread-local well-formedness conditions

Often one wants to consider only program expressions which satisfy some well-formedness condition. This section introduces the machinery provided by the framework to lift such conditions to the multithreaded case at no cost. An instantiating semantics can define a predicate Q of type $'e \Rightarrow 'c \Rightarrow 'x \Rightarrow bool$ and the operator $\uparrow\uparrow$ automatically lifts Q to a predicate of type $(e, 't, 'x) \text{ thread-info} \Rightarrow 'c \Rightarrow bool$. Formally:

$$\uparrow\uparrow Q \uparrow e s c \equiv \forall t. \text{case } e s t \text{ of None} \Rightarrow \text{True} \mid [(e, x)] \Rightarrow Q e c x$$

Suppose, e.g., that on the single-thread level, we have a definite assignment requirement, i.e. every variable must be assigned before being used. Suppose further, the predicate \mathcal{D} , which takes an expression e and a store for local variables x , guarantees that evaluating e satisfies the definite assignment condition with variables in the local store x having already been initialised. Then, $\uparrow\uparrow \mathcal{D} \uparrow e s c$ says that every thread t in $e s$, say $e s t = [(e, x)]$, satisfies $\mathcal{D} e x$.

In the context of type safety proofs, such well-formedness conditions are usually preserved under reductions in r . The predicate $P, r \vdash Q \checkmark \rightarrow$ gives three conditions that are sufficient for Q being also preserved under reduction in the framework semantics:

1. Q must be preserved under reductions of the instantiating semantics.
2. Q must also hold for new threads at the time of creation.
3. Q is preserved even if another thread, which also satisfies Q , changes the shared memory in a single-step reduction in r .

The next lemma shows that a predicate Q satisfying $P, r \vdash Q \checkmark \rightarrow$ is in fact preserved under multithreaded reductions, both single-step and multi-step. The first part is shown by a case analysis and induction on tas , the second by induction on the number of reduction steps.

LEMMA 5 (Soundness of the predicate $P, r \vdash \cdot \checkmark \rightarrow$).

$$\frac{P, r \vdash \langle ls | es, c | ws \rangle - t, tas \rightarrow \langle ls' | es', c' | ws' \rangle \quad \uparrow Q \uparrow es \ c \quad P, r \vdash Q \checkmark \rightarrow}{\uparrow Q \uparrow es' \ c'} \\ \frac{P, r \vdash \langle ls | es, c | ws \rangle - tas \rightarrow * \langle ls' | es', c' | ws' \rangle \quad \uparrow Q \uparrow es \ c \quad P, r \vdash Q \checkmark \rightarrow}{\uparrow Q \uparrow es' \ c'}$$

Returning to the example above, if we can show for r that $\mathcal{D}ex$ is preserved under reductions and all new threads also satisfy the predicate, then $P, r \vdash \mathcal{D} \checkmark \rightarrow$ holds and $\uparrow \mathcal{D} \uparrow es \ c$ is preserved under multithreaded reductions, too.

Similarly, such predicates on the thread level sometimes also needs some extra data, which is thread-specific, but invariant under reductions, e.g. a typing environment for local variables. We model such extra invariant data as maps from thread IDs to some type $'i$. As before, suppose Q is a predicate of type $'i \Rightarrow 'e \Rightarrow 'c \Rightarrow 'x \Rightarrow bool$ on the invariant data and the thread state. We automatically lift Q to the framework semantics, denoted by $\uparrow Q \uparrow$ of type $(t \rightarrow 'i) \Rightarrow ('e, 't, 'x) \text{ thread-info} \Rightarrow 'c \Rightarrow bool$.

We say a **map** I (type $t \rightarrow 'i$) **to invariant data** is well-formed w.r.t. a thread map es , denoted by $I \vdash_i es \checkmark$, iff I and es are defined on the same set of thread IDs. The partial order **extends** $I \trianglelefteq I'$ on maps to invariant data denotes that I' is an extension of I , i.e. I' is defined whenever I is defined, and in that case, I and I' coincide.

Let $I [\mathcal{I} \rightsquigarrow]_Q \text{ tas}$ denote the extension of I with invariant data for all threads created in tas . These updates then preserve well-formedness of maps to invariant data, which is shown by case analysis, induction on tas , and by induction on the number of reduction steps.

LEMMA 6.

$$\frac{P, r \vdash \langle ls | es, c | ws \rangle - t, tas \rightarrow \langle ls' | es', c' | ws' \rangle \quad es \vdash_i I \checkmark}{es' \vdash_i (I [\mathcal{I} \rightsquigarrow]_Q \text{ tas}) \checkmark} \\ \frac{P, r \vdash \langle ls | es, c | ws \rangle - tas \rightarrow * \langle ls' | es', c' | ws' \rangle \quad es \vdash_i I \checkmark}{es' \vdash_i (I [\mathcal{I} \rightsquigarrow]_Q \text{ flatten}(\text{map snd } tas)) \checkmark}$$

Here we see the reason for storing the shared memory in *NewThread* actions. Only with this trick are we able to update I correctly in the case of multiple reductions because the data chosen for a new thread may depend on the shared memory at thread creation time.

Now, let the predicate $P, r, R \vdash Q \checkmark \rightarrow$ impose the following conditions on an predicate Q sufficient for Q being preserved under a multithreaded reduction in $redT r P$, provided that the well-formedness predicate R holds in the initial state. (This extra predicate R is, strictly speaking, not necessary for the proofs, but allows for some modularity in applications.)

1. Q is preserved under reductions in $r P$ for threads satisfying R .
2. For all actions *NewThread* $t \ e \ c \ x$ issued by threads satisfying R , there is some invariant data i for (e, c, x) with $Q i \ e \ c \ x$.
3. Q is unaffected by changes to the shared memory by another thread that also satisfies both Q and R .

The next lemma shows that these conditions are sufficient for Q being preserved under multithreaded reductions. Again, the first part is shown by a case analysis and induction on tas , the second part by induction on the number of reduction steps.

LEMMA 7 (Soundness of the predicate $P, r, \cdot \vdash \cdot \checkmark \rightarrow$).

$$\frac{P, r \vdash \langle ls | es, c | ws \rangle - t, tas \rightarrow \langle ls' | es', c' | ws' \rangle \quad \uparrow Q \uparrow I \ e \ s \ c \quad \uparrow R \uparrow es \ c \quad P, r, R \vdash Q \checkmark \rightarrow}{\uparrow Q \uparrow (I [\mathcal{I} \rightsquigarrow]_Q \text{ tas}) \ e \ s' \ c'} \\ \frac{P, r \vdash \langle ls | es, c | ws \rangle - tas \rightarrow * \langle ls' | es', c' | ws' \rangle \quad \uparrow Q \uparrow I \ e \ s \ c \quad \uparrow R \uparrow es \ c \quad P, r, R \vdash Q \checkmark \rightarrow \quad P, r, R \vdash Q \checkmark \rightarrow}{\uparrow Q \uparrow (I [\mathcal{I} \rightsquigarrow]_Q \text{ flatten}(\text{map snd } tas)) \ e \ s' \ c'}$$

5. Multithreaded Jinja

In this section, we instantiate the framework with an extension of the Jinja source code semantics [12] for Java threads (without the JMM) and present the formalisation in detail. In modelling Java threads, we closely follow Ch. 17 in the Java Language Specification [9] for Java 5, making minor abstractions where special cases would have unnecessarily complicated the formalisation. First, we explain how Java threads are modelled in the framework. Then, we present the well-formedness conditions we need to impose on Jinja programs and expressions in order to show type safety (progress and preservation) in the last part.

In Jinja, source code expressions have the type *expr*, which is a data type with a constructor for every kind of operation, i.e. for creating objects and arrays, for casts, literals, access and assignment to local variables, arrays and fields, binary operations, method call, nested blocks, sequential composition, if and while statements, exception throwing and catching, and synchronized blocks. An expression is considered to be final iff it is a value or a thrown exception object on the heap. Although the syntax is quite different from Java, a compiler can translate any program which only uses Jinja features to the Jinja syntax.

5.1 Modelling the Java thread concept in Jinja

In our formalisation, the “program” P (type *J-prog*) contains the class declarations. Every thread stores its expression (*expr*) and its local variables (*locals*); since method calls are dynamically inlined, we do not need an explicit call stack. The shared memory (*heap*) is a map from addresses (*addr*) to heap objects (*heapobj*), thread IDs are natural numbers. In Java, only monitors – of which every object has one – can be locked, so locks are addresses, too. So are the wait set identifiers, because every monitor manages its own wait set. Thus, the Jinja semantics *red* is of type (*J-prog, expr, addr, nat, locals, heap, addr*) *semantics* and, accordingly, a reduction in it is denoted by $P \vdash (e, (h, x)) -tas-red \rightarrow (e', (h', x'))$.

Having explained the different instantiations of type parameters to the framework, we now present the introduction rules for the set of reductions *red P* which are relevant for our modelling Java threads. Fig. 8 and Fig. 9 show the reduction rules for method calls and the synchronized statement, respectively. We use Jinja syntax for expressions (e.g. $e.M(es)$ invokes the method M on the object designated by the expression e with parameter list es), see the Jinja semantics [12] for syntax and all other reduction rules, none of which issues thread actions on its own.

$$\begin{array}{c}
\frac{h a = \lfloor \text{Obj } C \text{ fs} \rfloor \quad P \vdash C \text{ sees } M: Ts \rightarrow T = (pns, \text{body}) \text{ in } D \quad |vs| = |pns| \quad |Ts| = |pns| \quad \neg P \vdash C \preceq^* \text{Thread} \vee M \neq \text{start}}{P \vdash \langle \text{addr } a \cdot M(\text{map } \text{Val } vs), (h, x) \rangle - [] - \text{red} \rightarrow \langle \text{blocks } (\text{this} \cdot pns, \text{Class } D \cdot Ts, \text{Addr } a \cdot vs, \text{body}), (h, x) \rangle} \text{CALL} \\
\frac{h a = \lfloor \text{Obj } C \text{ fs} \rfloor \quad P \vdash C \preceq^* \text{Thread}}{P \vdash \langle \text{addr } a \cdot \text{start}(\square), (h, x) \rangle - [\text{NewThread } t (\text{Var } \text{this} \cdot \text{run}(\square)) h [\text{this} \mapsto \text{Addr } a]] - \text{red} \rightarrow \langle \text{unit}, (h, x) \rangle} \text{NT1} \\
\frac{h a = \lfloor \text{Obj } C \text{ fs} \rfloor \quad P \vdash C \preceq^* \text{Thread}}{P \vdash \langle \text{addr } a \cdot \text{start}(\square), (h, x) \rangle - [\text{NewThreadFail}] - \text{red} \rightarrow \langle \text{THROW } \text{OutOfMemory}, (h, x) \rangle} \text{NT2} \\
\frac{h a = \lfloor q \rfloor}{P \vdash \langle \text{addr } a \cdot \text{wait}(\square), (h, x) \rangle - [\text{Suspend } a, \text{Unlock } a, \text{Lock } a] - \text{red} \rightarrow \langle \text{unit}, (h, x) \rangle} \text{W1} \\
\frac{h a = \lfloor q \rfloor}{P \vdash \langle \text{addr } a \cdot \text{wait}(\square), (h, x) \rangle - [\text{UnlockFail } a] - \text{red} \rightarrow \langle \text{THROW } \text{IllegalMonitorState}, (h, x) \rangle} \text{W2} \\
\frac{h a = \lfloor q \rfloor}{P \vdash \langle \text{addr } a \cdot \text{notify}(\square), (h, x) \rangle - [\text{Notify } a, \text{Unlock } a, \text{Lock } a] - \text{red} \rightarrow \langle \text{unit}, (h, x) \rangle} \text{N1} \\
\frac{h a = \lfloor q \rfloor}{P \vdash \langle \text{addr } a \cdot \text{notify}(\square), (h, x) \rangle - [\text{UnlockFail } a] - \text{red} \rightarrow \langle \text{THROW } \text{IllegalMonitorState}, (h, x) \rangle} \text{N2} \\
\frac{h a = \lfloor q \rfloor}{P \vdash \langle \text{addr } a \cdot \text{notifyAll}(\square), (h, x) \rangle - [\text{NotifyAll } a, \text{Unlock } a, \text{Lock } a] - \text{red} \rightarrow \langle \text{unit}, (h, x) \rangle} \text{NA1} \\
\frac{h a = \lfloor q \rfloor}{P \vdash \langle \text{addr } a \cdot \text{notifyAll}(\square), (h, x) \rangle - [\text{UnlockFail } a] - \text{red} \rightarrow \langle \text{THROW } \text{IllegalMonitorState}, (h, x) \rangle} \text{NA2}
\end{array}$$

Figure 8. Jinja reduction rules for method call.

5.1.1 Thread creation

In Java, threads are associated with objects of the class `Thread` or its subclasses. A new thread is spawned by invoking the `Thread` object’s `start` method. Hence, in our model, we treat objects of class `Thread` (or subclasses thereof) like standard objects, except that the well-formedness condition requires that no class comparable to `Thread` in the type hierarchy declares a method called `start`. Only when invoking the `start` method on such an object, rules NT1 and NT2 apply instead of CALL, which is for invoking nonsynthesized methods and dynamically inlines the called method’s code. In contrast, NT1 generates a new thread with an arbitrary thread ID, initial expression `Var this-run(\square)` (which represents the call to the thread object’s `run` method), the heap as the shared memory component, and, as thread local information, `this` initialised to the correct thread object, given by its address. Similarly, NT2 models the reduction in case there are no more free thread IDs, which results in an *OutOfMemory* exception being thrown in the spawning thread.

5.1.2 Wait, notify and notifyAll

We cannot implement the *wait*, *notify*, and *notifyAll* methods of Object inside the Jinja language, because they are synthesized just like spawning a thread. Hence, we include a pair of extra rules for each of these (rules W1, W2, N1, N2, NA1 and NA2), emulating the behaviour of the natives methods in the JVM:

1. One (W1, N1, NA1) models normal execution with three thread actions: The first one tells the framework to manipulate the wait sets according to the method’s meaning, the *Unlock* and subsequent *Lock* actions are to check if the current thread holds a lock on monitor associated with the wait set. If not, this reduction will not be chosen by the multithreaded semantics.
2. In that case, the other rule (W2, N2, NA2) instead raises an *IllegalMonitorState* exception if the thread does not hold a lock on the monitor associated with the object. The *UnlockFail* action ensures that this rule can only be chosen if the thread does not hold a lock on the monitor.

$$\begin{array}{c}
\frac{P \vdash \langle o', s \rangle - \text{tas} - \text{red} \rightarrow \langle o'', s' \rangle \quad \neg \text{lock-granted } o'}{P \vdash \langle \text{sync}(o') e, s \rangle - \text{tas} - \text{red} \rightarrow \langle \text{sync}(o'') e, s' \rangle} \text{S1} \\
P \vdash \langle \text{sync}(\text{null}) e, s \rangle - [] - \text{red} \rightarrow \langle \text{THROW } \text{NullPointer}, s \rangle \text{S2} \\
P \vdash \langle \text{sync}(\text{throw } a) e, s \rangle - [] - \text{red} \rightarrow \langle \text{throw } a, s \rangle \text{S3} \\
P \vdash \langle \text{sync}(\text{addr } a) e, s \rangle - [\text{Lock } a] - \text{red} \rightarrow \langle \text{sync}(\text{locked}(a)) e, s \rangle \text{S4} \\
\frac{P \vdash \langle e, s \rangle - \text{tas} - \text{red} \rightarrow \langle e', s' \rangle \quad \forall \text{tas}'. \text{tas} \neq \text{Suspend } a \cdot \text{tas}' \wedge \text{tas} \neq \text{UnlockFail } a \cdot \text{tas}'}{P \vdash \langle \text{sync}(\text{locked}(a)) e, s \rangle - \text{tas} - \text{red} \rightarrow \langle \text{sync}(\text{locked}(a)) e', s' \rangle} \text{S5} \\
\frac{P \vdash \langle e, s \rangle - \text{tas} - \text{red} \rightarrow \langle e', s' \rangle \quad \text{tas} = \text{Suspend } a \cdot \text{tas}'}{P \vdash \langle \text{sync}(\text{locked}(a)) e, s \rangle - \text{tas} @ [\text{Unlock } a] - \text{red} \rightarrow \langle \text{sync}(\text{addr } a) e', s' \rangle} \text{S6} \\
P \vdash \langle \text{sync}(\text{locked}(a)) \text{Val } v, s \rangle - [\text{Unlock } a] - \text{red} \rightarrow \langle \text{Val } v, s \rangle \text{S7} \\
P \vdash \langle \text{sync}(\text{locked}(a)) \text{throw } a', s \rangle - [\text{Unlock } a] - \text{red} \rightarrow \langle \text{throw } a', s \rangle \text{S8}
\end{array}$$

Figure 9. Jinja reduction rules for the synchronized statement.

Note that the framework semantics does not allow for “spurious wake-ups”, which are permitted by the JLS [9].

5.1.3 Synchronisation

Synchronisation in Java is done via synchronized statements. Sec. 14.19 in the JLS [9] determines its behaviour: “A synchronized statement acquires a mutual-exclusion lock on behalf of the executing threads, executes a block and then releases the lock.” Fig. 9 shows the rules for the synchronized statement. A synchronized statement is written `sync(o') e` where `o'` is the expression for the object on whose monitor a lock is acquired and `e` is the block’s expression. To remember syntactically that a synchronized statement has already acquired the lock on `a`, we use the expression `locked(a)` at `o'`. The predicate *lock-granted o'* checks whether `o'` is of the form `locked(a)` for some `a`. We now discuss the rules in Fig. 9 in detail.

S1 reduces the monitor subexpression, provided it is not the expression denoting the granted lock on the monitor. If the monitor subexpression reduces to the *null* value, a *NullPointerException* exception is thrown (S2). If an exception is raised while reducing the monitor subexpression, the same exception is propagated by S3. If it reduces to some address a , the thread can only reduce further (S4) by acquiring the lock on a .² Once the lock has been granted, the body e of the synchronized statement is reduced: If e 's reduction does not request that the thread is suspended to the wait set a , nor that the thread does not hold a lock on a , $\text{sync}(\text{locked}(a)) e$ is reduced accordingly by S5. If e 's reduction corresponds to a call to the *wait* method of a , which is characterised by a *Suspend a* action at the action list's head, additionally, the lock on a is released (by appending an *Unlock a* action to the end) (S6).

If the body reduces to a value in a normal way, the lock is released and the synchronized statement returns the value by S7. If an exception is thrown while the body is reduced, the lock on a is released and the exception propagated (S8).

Since synchronized methods are just syntactic sugar for ordinary methods with their whole body inside a synchronized statement on *this*, we omit this option in our formalisation for simplicity.

5.2 Well-formedness constraints

In Sec. 4, we have shown how local well-formedness constraints can be lifted to the multithreaded case. In this section, we introduce the well-formedness conditions for Jinja and show how they are lifted in the above way.

There are already a number of constraints for well-formed Jinja programs (see [12] for the details). For the threads, we need to impose some more on the class declarations. We say, a class with name C is well-formed if it is a well-formed Jinja class and satisfies additionally:

- if C is *Object*, it must not have fields or methods which are not native (and thus would be hardwired into the semantics)
- if C is *Thread*, it must have a *run* method with no parameters and return type *Void*,
- if C is comparable to *Thread* in the subtype relation, it must not have a *start* method (which is hardwired in the semantics), and
- it must not have methods called *wait*, *notify*, nor *notifyAll*.

A well-formed program P , denoted by *wf-J-prog* P , consists of well-formed class declarations only.

Since we use the special expression $\text{locked}(a)$ in monitor expressions of synchronized statements to remember that the lock on the monitor a has already been granted, we must also ensure that $\text{locked}(a)$ cannot occur in a monitor expression via subexpression reduction (rule S1 in Fig. 9). To this end, we define the predicate $\vdash e \sqrt{\&}$ for an expression e . Intuitively, $\vdash e \sqrt{\&}$ ensures that in an

²In the formalisation by Cenciarelli et al. [4], in rule **SYN2**, the monitor subexpression may evaluate to an address only if the lock can be acquired at the same time. This simplification, which we avoid, may have some strange effects. For instance, suppose we have a thread θ_1 with expression

$$\text{sync}(o.f) \{ o.f = \text{new Object}(); \text{sync}(o) \{ \dots \} \}$$

and another one θ_2 with expression

$$\text{sync}(o) \{ \text{sync}(o.f) \{ \dots \} \}$$

where o is a shared object, say O , with a field f initially referencing another object, say F . If θ_1 first obtains the lock on F , which is referenced by $o.f$, and then θ_2 is evaluated as far as possible, i.e. it is then waiting on the lock on F , this schedule leads to a deadlock in our semantics. In [4], however, θ_2 cannot evaluate that far, but stops evaluating $o.f$ at $O.f$. Once θ_1 changes $O.f$ to another object F' , θ_2 can continue with a lock on F' , i.e., there is no possibility for a deadlock here. Although the JMM introduces much complexity about which values a thread sees when, an example with a volatile field f reveals the same issue in the presence of the memory model.

expression e with multiple subexpressions, if an explicit address value occurs in such a subexpression e' , then all other subexpressions of e which are evaluated before e' according to the evaluation order imposed by the semantics must have already been completely evaluated. For example, in an array assignment $a[i] := e$, if a has evaluated to some address a and the index expression i has not yet been completely reduced to an integer, e must not contain an explicit address. We write $\vdash \cdot \uparrow \sqrt{\&}$ for $\vdash \cdot \sqrt{\&}$ lifted to the multithreaded case with the \uparrow operator.

LEMMA 8. If $P \vdash \langle o', s \rangle - \text{tas} - \text{red} \rightarrow \langle o'', s' \rangle$ and $\vdash \text{sync}(o') e \sqrt{\&}$ then $\neg \text{lock-granted } o''$.

By this lemma, shown by case analysis on the reduction steps, $\vdash \cdot \sqrt{\&}$ in fact does the job. The next lemma, which proves that $\vdash \cdot \sqrt{\&}$ is also preserved under reductions for well-formed programs, is shown by induction on the set of reductions.

LEMMA 9. If *wf-J-prog* P then $P, \text{red} \vdash (\vdash \cdot \sqrt{\&}) \checkmark \rightarrow$.

Locks are stored twice in the multithreaded Jinja semantics: On the one hand, there is the map ls for storing locks in the framework. On the other hand, we remember locks in the monitor subexpression of synchronized statements, which are stored in the map es . Naturally, we want them to be consistent, which is expressed by the predicate $es \vdash_e ls \checkmark$:

$$es \vdash_e ls \checkmark \equiv \forall t. \text{case } es \ t \text{ of } None \Rightarrow \forall l. \neg \text{has-lock } ls \ t \ l \mid [(e, x)] \Rightarrow \forall l. \text{if } 0 < \ell \ e \ l \text{ then } ls \ l = [(t, \ell \ e \ l - 1)] \text{ else } \neg \text{has-lock } ls \ t \ l$$

The function $\ell \ e \ l$ counts the number of $\text{sync}(\text{locked}(l))$ subexpressions in e . Note that $es \vdash_e ls \checkmark$ implies $es \vdash_f ls \checkmark$:

LEMMA 10. If $es \vdash_e ls \checkmark$ then $es \vdash_f ls \checkmark$.

LEMMA 11. For well-formed programs and well-formed thread expressions, $\cdot \vdash_e \cdot \checkmark$ is preserved under reductions:

$$\frac{\text{wf-J-prog } P \quad P, \text{red} \vdash \langle ls | es, c | ws \rangle - t, \text{tas} \rightarrow \langle ls' | es', c' | ws' \rangle \quad \vdash es \uparrow \sqrt{\&} \uparrow \quad es \vdash_e ls \checkmark}{es' \vdash_e ls' \checkmark}}{\text{wf-J-prog } P \quad P, \text{red} \vdash \langle ls | es, c | ws \rangle - \text{tas} \rightarrow * \langle ls' | es', c' | ws' \rangle \quad \vdash es \uparrow \sqrt{\&} \uparrow \quad es \vdash_e ls \checkmark}{es' \vdash_e ls' \checkmark}}$$

This lemma is shown by an extensive case analysis and induction on the thread action list, and the number of reduction steps.

Another feature of both Jinja and Java is definite assignment (cf. example in Sec. 4): It ensures that every variable must be assigned before being used. As before, we write $\mathcal{D} \ e \ x$, which takes an expression e and a store for local variables x , for the definite assignment test. The details for \mathcal{D} can be found in [12], the definitions for the synchronized statements are straight forward. The next lemma, which is shown by rule induction on the reductions, shows that definite assignment is preserved under reductions for well-formed programs.

LEMMA 12. If *wf-J-prog* P then $P, \text{red} \vdash \mathcal{D} \checkmark \rightarrow$.

5.3 The type safety proof

In this section, we introduce the Jinja type system and show the progress and subject reduction theorem, which we combine in the end to get type safety.

5.3.1 The type system

Jinja has two primitive types *Boolean* and *Integer*, for *Unit* the type *Void*, and reference types *NT*, *Class C*, and $T[\]$, where $\cdot[\]$ is a recursive type constructor for arrays. This way, array types of arbitrarily many dimensions are allowed in Jinja whereas the Java VM supports at most 255 dimensions [14].

$$\begin{array}{c}
P \vdash T \leq T \quad P \vdash NT \leq \text{Class } C \quad \frac{P \vdash A \leq B}{P \vdash A[] \leq B[]} \\
\hline
\frac{P \vdash C \preceq^* D}{P \vdash \text{Class } C \leq \text{Class } D} \quad P \vdash NT \leq A[] \quad P \vdash A[] \leq \text{Class Object}
\end{array}$$

Figure 10. Jinja subtyping rules

$$\begin{array}{c}
\frac{P, E \vdash e :: \text{Class } C \quad P \vdash C \text{ sees } M: Ts \rightarrow T = (pns, \text{body}) \text{ in } D \quad \text{list-all2 } (\lambda e T. P, E \vdash e :: T) \text{ es } Ts' \quad P \vdash Ts' [\leq] Ts}{P, E \vdash e \cdot M(\text{es}) :: T} \text{WTC1} \\
\frac{P, E \vdash e :: \text{Class } C \quad P \vdash C \preceq^* \text{Thread}}{P, E \vdash e \cdot \text{start}(\emptyset) :: \text{Void}} \text{WTC2} \\
\frac{P, E \vdash e :: T \quad \text{is-ref } T \quad T \neq NT}{P, E \vdash e \cdot \text{wait}(\emptyset) :: \text{Void}} \text{WTC3} \\
\frac{P, E \vdash e :: T \quad \text{is-ref } T \quad T \neq NT}{P, E \vdash e \cdot \text{notify}(\emptyset) :: \text{Void}} \text{WTC4} \\
\frac{P, E \vdash e :: T \quad \text{is-ref } T \quad T \neq NT}{P, E \vdash e \cdot \text{notifyAll}(\emptyset) :: \text{Void}} \text{WTC5} \\
\frac{P, E \vdash o' :: T \quad \text{is-ref } T \quad T \neq NT \quad P, E \vdash e :: T'}{P, E \vdash \text{sync}(o') e :: T'} \text{WTS}
\end{array}$$

Figure 11. Jinja typing rules for method call and the synchronized statement.

$$\begin{array}{c}
\frac{P, E, h \vdash e :: NT \quad \text{list-all2 } (\lambda e T. P, E, h \vdash e :: T) \text{ es } Ts}{P, E, h \vdash e \cdot M(\text{es}) :: T} \text{WTCN} \\
\frac{P, E, h \vdash o' :: NT \quad P, E, h \vdash e :: T}{P, E, h \vdash \text{sync}(o') e :: T'} \text{WTSN}
\end{array}$$

Figure 12. Jinja runtime typing rules for the *null* type in method call and the synchronized statement.

On the types of a program P , we define a **widening relation** whose rules are shown in Fig. 10: $P \vdash T \leq T'$ denotes that T is a subtype of T' . Subtyping behaves as in Java: The subclass relation $P \vdash \cdot \preceq^*$ is injected into the subtype relation, NT is a subtype of reference types, arrays are subtypes of *Class Object* and $\cdot []$ is covariant. We extend \leq to lists pointwise, denoted by $[\leq]$.

Jinja has a static type system with typing judgements of the form $P, E \vdash e :: T$ where P is the Jinja program, E the typing environment, i.e. a map from variable names to types, e is the expression to be typed and T is e 's Jinja type. Fig. 11 shows the typing rules for method calls and the synchronized statement. Due to space limitations, we cannot show all of the typing rules. For the remaining rules, see [12].

Rule WTC1 is the standard Jinja rule for method call where $P \vdash C \text{ sees } M: Ts \rightarrow T = (pns, \text{body}) \text{ in } D$ denotes that in program P , class C sees the method M with parameter types Ts and return type T in class D which has parameter names pns and body body , taking method overriding into account. $\text{list-all2 } (\lambda e T. P, E \vdash e :: T) \text{ es } Ts'$ denotes that $|es| = |Ts'|$ and the parameters es have types Ts' . As before, this rule cannot subsume the hard-wired methods for thread creation, *wait*, *notify*, and *notifyAll*, because they do not have an implementation in the Jinja language. Thus, we introduce one extra rule for each of them where *is-ref* T is a predicate for reference types (rules WTC2, WTC3, WTC4 and WTC5). The rule for the synchronized statement WTS is straightforward. Note that we include

the condition $T \neq NT$ in these rules to disallow expressions like *null-wait*(\emptyset), since *null* cannot directly be dereferenced in Java [9].

However, the static type system is not preserved under reductions. In particular, during reduction, explicit addresses may occur inside expressions which are not typable without knowing the heap. Thus, we also have another (runtime) type system [5], which takes the heap into account and which is denoted by $P, E, h \vdash e :: T$. It relaxes some constraints imposed by the static type system, for details see [12]. Most importantly, in addition to each rule in Fig. 11, which also exist in the runtime type system with the h added to them, there is another rule for typing the case when an expression typed with a non-*null* reference type reduces to *null*. Fig. 12 shows the typing rules for method call WTCN, which also subsumes the new native methods, and the synchronized statement WTSN. Note that, e.g., the expression *null-wait*(\emptyset) now has arbitrary type, not even a unique least type, however, since this is only technical device for the type safety proof, this does not matter to us when typing an expression statically in the type system above.

LEMMA 13. *If $P, E \vdash e :: T$ then $P, E, h \vdash e :: T$.*

5.3.2 Progress

For both progress and preservation, the initial configuration must conform to the type constraints given by the program in various aspects. First, all field contents of all objects on the heap h must have a type which conforms to the type given in the field declaration, which we denote by $P \vdash h \checkmark$. Similarly, every local variable in x must hold a value which is compatible with the typing environment E , denoted by $P, h \vdash x (\leq) E$. If both heap h and local store x are conform, we write $P, E \vdash (h, x) \checkmark$ as shorthand. Details for conformance can be found in [12], the extension to the Jinja thread semantics is straightforward. For progress, we lift typability to the multithreaded setting. $P, h \vdash \cdot :: es \checkmark$ holds iff for every thread t , there exists a typing environment in which t 's expression is (runtime) typable. With the runtime type system, we can carry over the induction proof on *red P* for progress to the thread extension under the same conditions:

LEMMA 14 (Progress for threads in isolation).

$$\frac{\text{wf-J-prog } P \quad P \vdash h \checkmark \quad P, E, h \vdash e :: T \quad \mathcal{D} e x \quad \neg \text{final } e}{\exists e' s' \text{tas}. P \vdash \langle e, (h, x) \rangle \text{-tas-red} \rightarrow \langle e', s' \rangle}$$

Hence, we have:

LEMMA 15 (*wf-progress* for *red P*).

$$\frac{\text{wf-J-prog } P \quad P \vdash h \checkmark \quad P, h \vdash \cdot :: es \checkmark \quad \uparrow \mathcal{D} \uparrow es h}{\text{wf-progress red } P \text{ es } h}$$

For progress of the framework semantics, we also need to show that the instantiating Jinja semantics is well-behaved:

LEMMA 16.

If $\vdash es \uparrow \checkmark \& \uparrow$ and $es \vdash e \text{ ls } \checkmark$ then *wf-r-progress red P ls es h t*.

Together with Thm. 1 and Lem. 4, we get progress for the framework semantics:

THEOREM 2 (Progress).

$$\frac{\uparrow \mathcal{D} \uparrow es h \quad \text{wf-J-prog } P \quad P, h \vdash \cdot :: es \checkmark \quad \uparrow \mathcal{D} \uparrow es h \quad \vdash es \uparrow \checkmark \& \uparrow \quad es \vdash e \text{ ls } \checkmark \quad P \vdash h \checkmark}{\exists t \text{tas } es' \text{ls}' \text{ws}' h'. P, \text{red} \vdash \langle \text{ls} | es, h | \text{ws} \rangle \text{-t, tas} \rightarrow \langle \text{ls}' | es', h' | \text{ws}' \rangle}$$

5.3.3 Preservation

For preservation, we combine conformance conditions and typability in a single predicate. First, let $P, E, h \vdash e \leq T$ be a shorthand for $\exists T'. (P, E, h \vdash e :: T' \wedge P \vdash T' \leq T)$, i.e. e is typable in the environment E and heap h with a type T' which is a subtype of T . However,

this is not the only such condition where the typing environment E is used: $P, E \vdash (h, x) \checkmark$ also depends on E . All in all, we define $P, (E, T) \vdash e, x, h \checkmark \equiv P, E \vdash (h, x) \checkmark \wedge P, E, h \vdash e \leq T$ and lift it to the multithreaded setting:

$$P, Es \vdash es, h \uparrow \checkmark \uparrow \checkmark \equiv \uparrow (\lambda(E, T) e h x. P, (E, T) \vdash e, x, h \checkmark) \uparrow Es es h$$

Note that this time, since the typing environment and the initial type must remain the same during reduction, we use a map to invariant data Es which is a pair of a typing environments E and a type T . Hence, for every thread t , $Es t$ stores both the typing environment for t 's local variables and the type of t 's initial expression.

The next lemma, which follows from a number of different other preservation lemmata not presented here, shows that $P, Es \vdash es, h \uparrow \checkmark \uparrow \checkmark$ is invariant for well-formed programs if threads respect definite assignment.

LEMMA 17. If $wf\text{-}J\text{-}prog P$ then $P, red, \mathcal{D} \vdash (P, \cdot \vdash \cdot, \cdot \checkmark) \checkmark \rightarrow$.

Subject reduction, the key lemma for preservation, immediately follows from this using Lem. 7:

THEOREM 3 (Subject reduction).

$$\frac{wf\text{-}J\text{-}prog P \quad \frac{P, red \vdash \langle ls | es, h | ws \rangle -t, tas \rightarrow \langle ls' | es', h' | ws' \rangle}{\uparrow \mathcal{D} \uparrow es h \quad P, Es \vdash es, h \uparrow \checkmark \uparrow \checkmark}}{P, (Es [T \rightsquigarrow] P, \cdot \vdash \cdot, \cdot \checkmark) \vdash es', h' \uparrow \checkmark \uparrow \checkmark}$$

5.3.4 Type safety

We use the preservation lemmata for the various conditions to lift all single-step reductions to the transitive reflexive closure by induction. With Thm. 2 and 3, this gives the final type safety theorem:

THEOREM 4 (Type safety).

$$\frac{wf\text{-}J\text{-}prog P \quad \begin{array}{l} es \vdash_i Es \checkmark \quad P, Es \vdash es, h \uparrow \checkmark \uparrow \checkmark \quad \uparrow \mathcal{D} \uparrow es h \quad es \vdash_e ls \checkmark \\ \vdash es \uparrow \checkmark \& \uparrow \quad P, red \vdash \langle ls | es, h | ws \rangle -tas \rightarrow * \langle ls' | es', h' | ws' \rangle \\ \nexists t \ ta \ es'' \ ls'' \ ws'' \ h'', P, red \vdash \langle ls' | es', h' | ws' \rangle -t, ta \rightarrow \langle ls'' | es'', h'' | ws'' \rangle \\ Es' = Es [T \rightsquigarrow] P, \cdot \vdash \cdot, \cdot \checkmark \quad \checkmark \text{flatten} (map \ snd \ tas) \end{array}}{Es \leq Es' \wedge \begin{array}{l} (\forall t \ e'. \exists x'. es' t = \lfloor (e', x') \rfloor) \rightarrow \\ (\exists v. e' = Val v \wedge (\exists E T. Es' t = \lfloor (E, T) \rfloor \wedge P, h' \vdash v \leq T)) \vee \\ (\exists a. e' = Throw a \wedge a \in dom h') \vee \\ (t \in \text{deadlocked red } P \ ls' \ es' \ ws' \ h' \wedge (\exists E T. Es' t = \lfloor (E, T) \rfloor \wedge \\ P, E, h' \vdash e' \leq T)) \end{array}}$$

Let us examine what Thm. 4 states: Suppose we have:

- a well-formed Jinja program P , and
- there is a typing environment and an initial type in Es for thread ID t iff there is a thread with ID t in es , and
- for each thread t in es , say $es t = \lfloor (e, x) \rfloor$ and $Es t = \lfloor (E, T) \rfloor$, its expression e and local variable store x is typable in the type environment E and the heap h with a subtype of T , and
- every thread passes the definite assignment check, and
- locks are held only by non-final threads, and
- no thread can be reduced such that the monitor subexpression of a synchronized statement reduces to $locked(a)$ for some address a , and
- the state (ls, es, c, ws) reduces to some state (ls', es', c', ws') with thread action lists tas such that
- (ls', es', h', ws') is in normal form.

Let Es' denote the map to typing environments and initial types Es which is extended with the environments and types for threads that are newly created in tas . Note that for threads t that have already

been present in Es , $Es' t = Es t$, by $Es \leq Es'$, i.e. type environments and initial types have remained the same for them. Then, for every thread t in es' , say $es' t = \lfloor (e', x') \rfloor$, one of the following cases holds:

1. e' is a final value v , whose type is a subtype of t 's initial type, or
2. e' throws an object that exists on the heap h , or
3. t is deadlocked and e' is typable with a subtype of t 's initial type.

6. Related work

There are a number of formal semantics for Java on the source code level [21, 5, 22, 12], all of which model different subsets of sequential Java. Our basis for the sequential part is the Jinja semantics [12], which is a successor to the Bali project [22]. In [4], Cenciarelli et al. give a formal semantics of multithreaded Java on the source code level, which includes most Java thread features such as dynamic thread running and stopping, synchronisation via monitors and the wait/notify mechanism. In particular, using event spaces, they carefully model the memory model for Java 2 [8], which is now out-dated. However, they neither give a type system nor do they prove any meta-theoretic results on their semantics. In contrast, they say:

Event spaces are not necessarily “complete” [...] In fact, there are well-formed event spaces which are not completable, and this complicates the meta-theory of the semantics.

Since we are aiming for type safety via progress and preservation, we cannot resort to their semantics because non-completable event spaces would break the progress proof.

Stärk et al. [19] also present a multithreaded semantics of Java (without the JMM) based on abstract state machines together with a proof for preservation. However, they do not consider deadlocks, neither do they give a proof for progress. Moreover, their proofs are not checked by a theorem prover.

On the bytecode level, Belblidia and Debabbi present a formal small-step semantics for Java bytecode [3] which also features threads, but not the JMM. Like our approach, they have a semantics for threads in isolation and a second layer which manages the threads and receives thread actions, which they call labels, from them. In contrast to our framework semantics, at most one action can be issued at a time, but their single-thread semantics already takes care of the locks, which are stored in the shared memory, i.e. they only have actions for creating, killing, blocking and notifying threads. Yet, they do not model the wait/notify mechanism, which – strictly speaking – is not an integral part of Java bytecode, but of the `java.lang` package [14]; their block and notify actions are used by the second layer to keep track of which threads are ready for execution. Like Cenciarelli et al. in [4], they only give the semantics, but no type system and no discussion about deadlocks. Our framework semantics also handles the locks, i.e. an instantiating semantics need not care about other threads, and is generic and far more versatile, in particular, because arbitrary lists of thread actions can be passed in a single reduction. Their second layer, however, is tailored to their single-thread semantics for Java bytecode.

In [15], Liu and Moore present a monolithic formal semantics for multithreaded Java bytecode in form of an interpreter which also models class loading and initialisation. They aim for verifying JVM implementations w.r.t. the JVM specification and small Java programs in ACL2. However, they do not give a type system, nor do they comment on the deadlock issues.

Apart from showing type safety for the Java type system, type systems have proven useful for other safety features. Flanagan and Abadi [6] came up with an object calculus and a type system with dependent types to ensure that data races in accessing object

members cannot occur. A data race occurs if two threads can access a location simultaneously without synchronisation, which can result in a corrupted state. Object members are annotated with locks' names, the type system ensures that accessing a member is only possible if the specified lock is held by the thread. An appropriate subject reduction theorem shows soundness. However, they do not provide any progress results. Flanagan and Freund [7] translated this calculus to full Java bytecode and implemented it in the `rccjava` tool. In [10], Grossman extends their approaches [6, 7] to multithreaded Cyclone, which is a type safe variant of C. He also shows the progress property that no well-typed thread can get “badly stuck”. A thread t is “badly stuck” iff t cannot reduce any further even if t could acquire an arbitrary additional lock and, in case being final, t still holds some lock. Together with the subject reduction theorem, type safety, i.e. all threads reachable from a well-typed thread via reductions are not badly stuck, follows. In general, being in deadlock is stronger than being badly stuck because the latter does not involve the aspect of circular waiting.

Programs without data races are a desired kind of Java programs: The new JMM ensures that such programs have sequentially consistent behaviour. Aspinall and Ševčík [2] have formalised data race freedom and the memory model in Isabelle/HOL and proved this guarantee.

There are also approaches to eliminate potential deadlocks via type systems: For example, Suenaga and Kobayashi [20] propose a process calculus with thread creation, interrupts (which can be temporarily disabled), and synchronisation via structured locking. They assign to each syntactic occurrence of a lock a unique level tag. Their type system remembers bounds on the level of acquired locks in effect labels and ensures deadlock freedom by requiring that locks must be acquired in ascending order. Their deadlock formalisation is on the syntactic level only: A set of threads is in deadlock iff every “reducible” subexpression of the threads' expressions must be a synchronisation statement which has to acquire a lock which is already held. Consequently, they cannot express that some threads are in deadlock while others are still active. Moreover, in their model, a thread cannot acquire a lock multiple times as it is the case in Java. In contrast to that, our approach is defined in terms of the semantics, handles the extra cases introduced by wait sets and computes whether a given thread is in deadlock – even if some threads are not yet in deadlock.

Another interesting issue is to ensure that every lock acquired is eventually released and that a lock must have been acquired before being released, i.e. to avoid objects being locked forever and `IllegalMonitorStateException` being raised at `monitorexit` instructions. For Java on the source level, this is no issue because locks are reliably acquired and released in a structured fashion by `synchronized` blocks and methods, but when we want to extend our formalisation to bytecode, `monitorenter` need not always be matched by `monitorexit`, which violates the $es \vdash_{f} ls \checkmark$ well-formedness condition of our framework. In particular, if an exception is raised while the lock is held and the exception handler is outside the synchronized section, things get tricky. Iwama and Kobayashi [11] propose to tag every object with a usage label which specifies a policy on how this object may be locked. A type system, for which they also give a type inference algorithm, guarantees that method implementations respect the usage tags, and the subject reduction theorem ensures soundness. However, their subset of Java bytecode is rather restricted in not allowing object fields, method invocations and inheritance. To tackle this problem of proper lock acquisition and release, Laneve [13] presents an operational semantics and a type system for a slightly larger subset of Java bytecode which includes both synchronisation and the wait/notify mechanism. It enforces the structured

locking principle that is known from Java source code. Hence, the soundness proof gives that well-typed programs are free of `IllegalMonitorStateExceptions`, even for calls to `wait` and `notify`. He also discusses the intricacies introduced by exception handling in detail. However, his type system is not Java's and he does not mention progress.

7. Conclusion and future work

We have presented a general framework for adding concurrency to a single-threaded, operational semantics, in the proof assistant Isabelle/HOL. The framework semantics unhinges thread management and synchronisation from the instantiating semantics, thus allowing for modular proofs of meta-theoretic properties such as type safety. Further, we have given a formal definition of deadlock in terms of the semantics. In extending the Jinja source code semantics and applying the framework to it, we have type safety for multithreaded Java machine-checked, via progress and preservation.

For the future, we plan to extend the Jinja bytecode semantics to threads and to apply the framework there, too. The formalisation of byte code instructions related to concurrency by Belblidia and Debbabi [3] will be a good start there. Since our framework requires that every thread releases all locks it has acquired before terminating, we need to formalise how to ensure proper locking and unlocking on the bytecode level. To this end, we may recourse on the ideas by Iwama et al. [11] or Laneve [13]. When we then also add the JMM, Aspinall's and Ševčík's work [2] will be a good basis. Beyond Java, we also plan to add posix-like threads to a formal semantics of C++ [23] in Isabelle/HOL.

Acknowledgments

We would like to thank Daniel Wasserrab and Dennis Giffhorn for their valuable comments and the inspiring discussions we had with them. Also, we would like to thank the anonymous referees for their useful remarks.

References

- [1] J. Alves-Foss, editor. *Formal Syntax and Semantics of Java*. LNCS 1523. Springer, 1999.
- [2] D. Aspinall and J. Ševčík. Formalising Java's Data Race Free Guarantee. In *TPHOLS*, pages 22–37, 2007.
- [3] N. Belblidia and M. Debbabi. A Dynamic Operational Semantics for JVM. *Journal of Object Technology*, 6(3):71–100, 2007.
- [4] P. Cenciarelli, A. Knapp, B. Reus, and M. Wirsing. An Event-Based Structural Operational Semantics of Multi-Threaded Java. In Alves-Foss [1], pages 157–200.
- [5] S. Drossopoulou and S. Eisenbach. Describing the Semantics of Java and Proving Type Soundness. In Alves-Foss [1], pages 41–82.
- [6] C. Flanagan and M. Abadi. Object Types against Races. In *CONCUR*, pages 288–303, 1999.
- [7] C. Flanagan and S. N. Freund. Type-Based Race Detection for Java. In *PLDI*, pages 219–232, 2000.
- [8] J. Gosling, B. Joy, G. Steele, and G. Bracha. *The Java Language Specification, Second Edition*. Addison-Wesley, 2000.
- [9] J. Gosling, B. Joy, G. Steele, and G. Bracha. *The Java Language Specification, Third Edition*. Addison-Wesley, 2005.
- [10] D. Grossman. Type-Safe Multithreading in Cyclone. In *TLDI*, pages 13–25, 2003.
- [11] F. Iwama and N. Kobayashi. A New Type System for JVM Lock Primitives. In *ASIA-PEPM*, pages 71–82, 2002.

- [12] G. Klein and T. Nipkow. A Machine-Checked Model for a Java-Like Language, Virtual Machine and Compiler. *TOPLAS*, 28(4):619–695, 2006.
- [13] C. Laneve. A Type System for JVM Threads. Technical Report TCS 290, University of Bologna, 2003.
- [14] T. Lindholm and F. Yellin. *The Java Virtual Machine Specification, Second Edition*. Addison-Wesley, 1999.
- [15] H. Liu and J S. Moore. Executable JVM Model for Analytical Reasoning: A Study. In *IVME*, pages 15–23, 2003.
- [16] A. Lochbihler. Jinja With Threads. In Gerwin Klein, Tobias Nipkow, and Lawrence Paulson, editors, *The Archive of Formal Proofs*. <http://afp.sf.net/entries/JinjaThreads.shtml>, 2007. Formal proof development.
- [17] T. Nipkow, L. C. Paulson, and M. Wenzel. *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. LNCS 2283. Springer, 2002.
- [18] G. Nutt. *Operating Systems*. Addison-Wesley, 2nd edition, 2000.
- [19] R. F. Stärk, J. Schmid, and E. Börger. *Java and the Java Virtual Machine*. Springer, 2001.
- [20] K. Suenaga and N. Kobayashi. Type-Based Analysis of Deadlock for a Concurrent Calculus with Interrupts. In *ESOP*, pages 490–504, 2007.
- [21] D. Syme. Proving Java Type Soundness. In Alves-Foss [1], pages 83–118.
- [22] D. von Oheimb and T. Nipkow. Machine-Checking the Java Specification: Proving Type-Safety. In Alves-Foss [1], pages 119–156.
- [23] D. Wasserrab, T. Nipkow, G. Snelting, and F. Tip. An Operational Semantics and Type Safety Proof for Multiple Inheritance in C++. In *OOPSLA*, pages 345–362, 2006.
- [24] A. K. Wright and M. Felleisen. A Syntactic Approach to Type Soundness. *Information and Computation*, 115(1):38–94, 1994.