

# Unifying Refinement and Hoare-Style Reasoning in a Logic for Higher-Order Concurrency

Aaron Turon

MPI-SWS  
turon@mpi-sws.org

Derek Dreyer

MPI-SWS  
dreyer@mpi-sws.org

Lars Birkedal

Aarhus University  
birkedal@cs.au.dk

## Abstract

Modular programming and modular verification go hand in hand, but most existing logics for concurrency ignore two crucial forms of modularity: *higher-order functions*, which are essential for building reusable components, and *granularity abstraction*, a key technique for hiding the intricacies of fine-grained concurrent data structures from the clients of those data structures. In this paper, we present CaReSL, the first logic to support the use of granularity abstraction for modular verification of higher-order concurrent programs. After motivating the features of CaReSL through a variety of illustrative examples, we demonstrate its effectiveness by using it to tackle a significant case study: the first formal proof of (partial) correctness for Hendler *et al.*'s “flat combining” algorithm.

**Categories and Subject Descriptors** D.3.1 [Programming Languages]: Formal Definitions and Theory; F.3.1 [Logics and Meanings of Programs]: Specifying and Verifying and Reasoning about Programs

**Keywords** Contextual refinement; higher-order functions; fine-grained concurrency; separation logic; Kripke logical relations.

## 1. Introduction

Over the past decade, a number of Hoare logics have been developed to cope with the complexities of concurrent programming [20, 32, 7, 3, 18, 4, 27]. Unsurprisingly, the name of the game in these logics is improving support for *modular* reasoning along a variety of dimensions. Concurrent Abstract Predicates (CAP) [3], for example, utilizes a deft combination of *separation logic* [26] (for spatially-modular reasoning about resources and ownership), *rely-guarantee* [15] (for thread-modular reasoning in the presence of interference), and *abstract predicates* [21] (for hiding invariants about a module's private data structures from its clients).

At the same time, of course, the importance of modularity is not restricted to verification. Programmers use modularity in the *design* of their concurrent programs, precisely to enable reasoning about individual program components in relative isolation. And indeed, certain aspects of advanced concurrency logics, such as their aforementioned use of abstract predicates, are geared toward building proofs that reflect the data abstraction inherent in well-designed programs.

We contend, however, that existing concurrency logics are not exploiting the modular design of sophisticated concurrent programs to full effect. In particular, we observe that there are two crucial dimensions of modular concurrent programming that existing logics provide no way to leverage in the construction of proofs: namely, **higher-order functions** and **granularity abstraction**.

**Higher-order concurrency** Higher-order functional abstraction is of course one of the most basic hammers in the modern programmer's toolkit for writing reusable and modular code. Moreover, a number of concurrent programming patterns rely on it: work stealing [2], Concurrent ML-style events [25], concurrent iterators [16], parallel evaluation strategies [29], and monadic approaches to concurrent programming [10], just to name a few.

Yet, verification of higher-order concurrent programs remains a largely unexplored topic. To our knowledge, only a few existing logics can handle higher-order concurrent programs [27, 17, 14], but these logics are limited in other ways—in particular, they do not presently support verification of “fine-grained” concurrent ADTs. This leads us directly to the second limitation we observe of the state of the art, concerning granularity abstraction.

**Granularity abstraction via contextual refinement** An easy way to adapt a sequential mutable data structure for concurrent access is to employ *coarse-grained* synchronization: use a single global lock, and instrument each of the operations on the data structure so that they acquire the lock before they begin and release it after they complete. On the other hand, more sophisticated implementations of concurrent data structures employ *fine-grained* synchronization: they protect different parts of a data structure with different locks, or avoid locking altogether, so that threads can access disjoint pieces of the data structure in parallel.

There may seem at first glance to be a fundamental trade-off here. Fine-grained synchronization enables parallelism, but makes the data structures that use it very tricky to reason about directly, due to their complex internal coordination between threads. Coarse-grained synchronization sequentializes access to the data structure, which is bad for parallelism but perfect for client-side reasoning, since it enables clients to reason about concurrent accesses as if each operation takes effect atomically.

Fortunately, modular programming comes to the rescue. In particular, so long as tricky uses of fine-grained synchronization are confined to the hidden state of a carefully crafted ADT, it is possible to prove that the fine-grained implementation of the ADT is a *contextual refinement* of some coarse-grained implementation. Contextual refinement means that, assuming clients only access the ADT through its abstract interface (so that the state really is hidden), every behavior that clients can observe of the fine-grained implementation is also observable of the coarse-grained one. Thus, clients can pretend, for the purpose of simplifying their own verification, that they are using the coarse-grained version, yet be sure that their code will still be correct when linked with the more efficient fine-grained version. This is what we call *granularity abstraction*. (Note: gran-

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. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ICFP '13, September 25–27, 2013, Boston, MA, USA.

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-2326-0/13/09...\$15.00.

<http://dx.doi.org/10.1145/2500365.2500600>

ularity abstraction is similar to *atomicity abstraction* [19], but is more general in that, as we will see in the iterator example later in this section, it applies even if the target of the abstraction is only *somewhat* coarse-grained.)

To illustrate this point more concretely, let us consider a simple motivating example of reasoning about Treiber’s stack [28]. (We will in fact use this very example as part of a larger case study later in the paper.) Treiber’s stack is a fine-grained implementation of a concurrent stack ADT. Instead of requiring concurrently executing push and pop operations to contend for a global lock on the whole stack (as a coarse-grained implementation would), Treiber’s implementation allows them to race to access the head of the stack using compare-and-set (CAS). (The implementation of Treiber’s stack is shown in Figure 9, and discussed in detail in §3.3.)

Now, the reader may expect that stacks should admit a canonical, principal specification (spec, for short), perhaps something like the following “precise” spec which tracks the exact contents  $s$  of the stack using the abstract predicate  $\text{Con}(s)$ :

$$\begin{aligned} & \{\text{Con}(s)\} \text{push}(x) \{\text{Con}(x :: s)\} \\ & \{\text{Con}(s)\} \text{pop}() \{\text{ret. } (\text{ret} = \text{none} \wedge s = \text{nil} \wedge \text{Con}(s)) \\ & \quad \vee (\exists x, s'. \text{ret} = \text{some}(x) \wedge s = x :: s' \wedge \text{Con}(s'))\} \end{aligned}$$

The trouble with using this spec in a concurrent setting is that knowledge about the exact contents of the stack is not stable under interference from other threads. As a result, some concurrency logics prohibit this spec altogether. Others permit the spec, but force the  $\text{Con}(s)$  predicate to be treated as a resource that only one thread can own at a time, thus effectively preventing any concurrent access to the stack and defeating the point of using Treiber’s stack in the first place!

We want instead to capture the idea that the client threads of a data structure interact with it according to some (application-specific) protocol. Take, for example, the following “per-item” spec, which abstracts away from the LIFO nature of the stack and instead imposes an item-level protocol:

$$\begin{aligned} & \forall x. \{p(x)\} \text{push}(x) \{\text{true}\} \\ & \wedge \{\text{true}\} \text{pop}() \{\text{ret. } \text{ret} = \text{none} \vee (\exists x. \text{ret} = \text{some}(x) \wedge p(x))\} \end{aligned}$$

Given an arbitrary predicate  $p$  of the client’s choice, this per-item spec asserts that the stack contains only elements that satisfy the predicate  $p$ . It is pleasantly simple, and sufficient for the purposes of the case study we present later in the paper. It should be clear, however, that this “per-item” spec is far from a canonical or principal specification of stacks: the same spec would also be satisfied, for instance, by queues.

This brings us to our key point: different clients have different needs, and so we may want to verify the stack ADT against a range of different Hoare-style specs, but we do not want to have to re-verify *Treiber’s implementation* each time. Ideally, we would like to modularly decompose the proof effort into two parts:

1. Prove that Treiber’s implementation is a contextual refinement of a coarse-grained, lock-based implementation, which serves as a simple *reference implementation* of the stack ADT.
2. Use Hoare-style reasoning to verify that this reference implementation satisfies the various specs of interest to clients.

This decomposition engenders a clean separation of concerns, confining the difficulty of reasoning about Treiber’s particular implementation<sup>1</sup> to the proof of refinement, and simplifying the verification of the stack ADT against different client specs.

The story we have just told is, *in principle*, nothing new. The ability to prove granularity abstraction via contextual refinement has been accepted in the concurrency literature as a useful correct-

<sup>1</sup> Treiber’s stack is actually one of the easiest fine-grained data structures to reason about, but our general line of argument applies equally well to more sophisticated implementations, such as the HSY elimination stack [12].

ness criterion for tricky concurrent data structures, precisely because of the modular decomposition of proof effort that it ought to facilitate [13, 8, 9]. Unfortunately, despite the utility of such a modular decomposition in theory, *no existing concurrency logic* actually supports it in practice. In particular, very few systems support proofs of contextual refinement at all, and those few that do support refinement proofs—such as the recent work of Liang and Feng [18]—do not provide a means of composing refinement with client-side Hoare-style verification in a unified logic.

**Granularity abstraction for higher-order functions** Although supporting granularity abstraction is already a challenge for first-order concurrent programs, it becomes even more interesting for higher-order concurrent programs.

Suppose, for instance, we wished to add a higher-order iterator, *iter*, to the concurrent stack ADT. (Such concurrent iterators are already commonplace [16].) Adding *iter* to Treiber’s implementation is trivial:  $\text{iter}(f)$  will (without acquiring any lock) just read the current head pointer of the stack and apply  $f$  to each element of the stack accessible from it. But how do we specify the behavior of this iterator? What reference implementation should we use in proving granularity abstraction? Unlike for push and pop, it does not make sense for the reference implementation of  $\text{iter}(f)$  to be “maximally” coarse-grained (*i.e.*, to execute entirely within a critical section) because that will not correspond to the reality of the implementation we just described. In particular, the Treiber implementation does *not* freeze modifications to the stack while  $f$  is being applied to each element, so it does not contextually refine the maximally coarse-grained implementation of iteration.

Rather, what the Treiber iterator guarantees is that  $f$  will be applied to all the elements that were accessible from some node that was the head of the stack *at some point in time*. A clean way to specify this behavior is via a reference implementation that (1) acquires the lock, (2) takes a “snapshot” (*i.e.*, makes a copy) of the stack, (3) releases the lock, and (4) iterates  $f$  over the snapshot. (For the code of this reference implementation, see Figure 9.) What makes this reference implementation so intriguing is that it is only *somewhat* coarse-grained, yet as we will show later in the paper, it is nonetheless quite useful as a target for granularity abstraction.

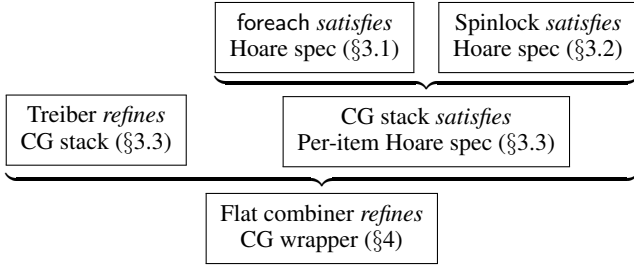
This example demonstrates the flexibility of refinement, a flexibility which has not heretofore been tested, since no one has previously applied granularity abstraction to higher-order ADTs.

## CaReSL: A logic for higher-order concurrency

In this paper, we present CaReSL (pronounced “carousel”), the first logic to support the use of granularity abstraction for modular verification of higher-order concurrent programs.

In providing both refinement and Hoare-style reasoning, CaReSL builds on ideas from two distinct lines of research:

- *Kripke logical relations* [22, 1, 5]. Logical relations are a fundamental technique for proving refinement in languages with *e.g.*, higher-order functions, polymorphism, and recursive types. Logical relations explain observable behavior in terms of the logical interpretation of each type. For example, two functions of type  $\tau \rightarrow \tau'$  are logically related if giving them related inputs of type  $\tau$  *implies* that they will produce related results of type  $\tau'$ . *Kripke* logical relations are defined relative to a “possible world” that describes the relationship between the internal, hidden state of a program and that of its spec.
- *Concurrent separation logics* [20, 32, 3]. Separation logic is a reimagining of Hoare logic in which assertions are understood relative to some *portion* of the heap, with the separating conjunction  $P * Q$  dividing up the heap between assertions  $P$  and  $Q$ . The motivation for separation is *local* reasoning: pre- and post-conditions need only mention the relevant portion of the



**Figure 1.** The structure of the case study—and paper

heap, and the rest of the heap can be presumed invariant. Most concurrent separation logics also add some form of protocols (*e.g.*, rely-guarantee), which facilitate compositional reasoning by constraining the potential interference between threads operating concurrently on some *shared* portion of the heap.

Rather than a combination of these approaches, CaReSL is an attempt to unify them. So, for example, the logic does *not* include refinement as a primitive notion. Instead, refinement is derived from Hoare-style reasoning: to prove that  $e$  refines  $e'$ , one merely proves a particular Hoare triple about  $e$ , allowing the tools of concurrent separation logic to be exploited in the proof of logical relatedness. In the other direction, CaReSL is a modal logic with the possible worlds of Kripke logical relations, which can in turn be used to express the shared-state protocols of concurrent separation logics. The unification yields surprising new techniques, such as the use of modal *necessity* ( $\Box$ ) to hide resources that an ADT’s client would otherwise have to thread through their reasoning (§3.2).<sup>2</sup>

The semantics of CaReSL is derived directly from the model of Turon *et al.* [31], who presented the first formal proofs of refinement for fine-grained data structures in an ML-like setting (although they did not use it to reason about higher-order examples). Compared with Turon *et al.*’s work, which was purely semantic, CaReSL provides a *syntactic* theory for carrying out refinement proofs at a much higher level of abstraction. This proof theory, in turn, is inspired by Dreyer *et al.*’s LADR [5], a modal logic for reasoning about contextual equivalence of (sequential) stateful ML programs. CaReSL exemplifies how the key ideas of LADR are just as relevant to—and arguably even more compelling when adapted to—the concurrent setting.

To demonstrate the effectiveness of CaReSL, we use it to tackle a significant case study: namely, the first formal proof of (partial) correctness for Hendler *et al.*’s “flat combining” algorithm [11]. Flat combining provides a generic way to turn a sequential ADT into a relatively efficient concurrent one, by having certain threads perform—in one fell swoop—the combined actions requested by a whole bunch of other threads. The flat combining algorithm is interesting not only because it itself is a rather sophisticated higher-order function, but also because it is modularly assembled from other higher-order components, including a fine-grained stack with concurrent iterator. We therefore take the opportunity to verify flat combining in a modular way that mirrors the modular structure of its implementation.

Figure 1 illustrates the high-level structure of our proof, which involves an intertwined application of refinement and Hoare-style verification. For example, we prove that Treiber’s stack refines the coarse-grained (CG) reference implementation of stacks; we then use Hoare-style reasoning to prove that the reference implementation satisfies the per-item spec; and finally we rely on the per-item spec in the proof that the flat combining algorithm refines an even

<sup>2</sup>Previously, such hiding required the subtle *anti-frame* rule [24].

## Syntax

Val  $v ::= () \mid \mathbf{true} \mid \mathbf{false} \mid (v, v) \mid \mathbf{inj}_i v \mid \mathbf{rec} f(x).e \mid \Lambda.e \mid \ell$   
 Exp  $e ::= v \mid \mathbf{if} e \mathbf{then} e \mathbf{else} e \mid e e \mid e \_ \mid (e, e) \mid \mathbf{prj}_i e \mid \mathbf{inj}_i e$   
 $\quad \mid \mathbf{case}(e, \mathbf{inj}_1 x \Rightarrow e, \mathbf{inj}_2 y \Rightarrow e) \mid \mathbf{new} e \mid \mathbf{get} e \mid e := e$   
 $\quad \mid \mathbf{CAS}(e, e, e) \mid \mathbf{newLcl} \mid \mathbf{getLcl}(e) \mid \mathbf{setLcl}(e, e) \mid \mathbf{fork} e$

CType  $\sigma ::= \mathbf{B} \mid \mathbf{ref} \tau \mid \mathbf{refLcl} \tau \mid \mu\alpha.\sigma$

Type  $\tau ::= \sigma \mid \mathbf{1} \mid \alpha \mid \tau \times \tau \mid \tau + \tau \mid \mu\alpha.\tau \mid \forall\alpha.\tau \mid \tau \rightarrow \tau$

ECTx  $K ::= [] \mid \mathbf{if} K \mathbf{then} e \mathbf{else} e \mid K e \mid v K \mid \dots$

## Typing contexts

$\Delta ::= \cdot \mid \Delta, \alpha \quad \Gamma ::= \cdot \mid \Gamma, x : \tau \quad \Omega ::= \Delta; \Gamma$

## Well-typed expressions

$\Delta; \Gamma \vdash e : \tau$

$\frac{\Omega \vdash e : \mathbf{ref} \sigma \quad \Omega \vdash e_o : \sigma \quad \Omega \vdash e_n : \sigma}{\Omega \vdash \mathbf{CAS}(e, e_o, e_n) : \mathbf{B}} \quad \frac{\Omega \vdash e : \forall\alpha.\tau}{\Omega \vdash e \_ : \tau[\tau'/\alpha]}$

## Machine configurations

TLV  $L \in \mathbb{N}^{\text{fin}}$  Value TPool  $\mathcal{E} \in \mathbb{N}^{\text{fin}}$  Expression

Heaps  $h \in \mathbb{N}^{\text{fin}}$  (Val  $\cup$  TLV) Config  $\varsigma ::= h; \mathcal{E}$

**Thread-local lookup**  $[L](i) \triangleq \begin{cases} \text{none} & i \notin \text{dom}(L) \\ \text{some}(L(i)) & \text{otherwise} \end{cases}$

## Pure reduction

$e \xrightarrow{\text{pure}} e'$

$(\mathbf{rec} f(x).e) v \xrightarrow{\text{pure}} e[\mathbf{rec} f(x).e/f, v/x] \quad (\Lambda.e) \_ \xrightarrow{\text{pure}} e$

## Per-thread reduction

$h; e \xrightarrow{i} h'; e'$

$h; e \xrightarrow{i} h; e' \quad \text{when } e \xrightarrow{\text{pure}} e'$   
 $h; \mathbf{newLcl} \xrightarrow{i} h \uplus [\ell \mapsto \emptyset]; \ell$   
 $h \uplus [\ell \mapsto L]; \mathbf{setLcl}(\ell, v) \xrightarrow{i} h \uplus [\ell \mapsto L[i := v]]; ()$   
 $h; \mathbf{getLcl}(\ell) \xrightarrow{i} h; v \quad \text{when } h(\ell) = L, [L](i) = v$   
 $h; \mathbf{CAS}(\ell, v_o, v_n) \xrightarrow{i} h; \mathbf{false} \quad \text{when } h(\ell) \neq v_o$   
 $h \uplus [\ell \mapsto v_o]; \mathbf{CAS}(\ell, v_o, v_n) \xrightarrow{i} h \uplus [\ell \mapsto v_n]; \mathbf{true}$

## General reduction

$h; \mathcal{E} \rightarrow h'; \mathcal{E}'$

$\frac{h; e \xrightarrow{i} h'; e'}{h; \mathcal{E} \uplus [i \mapsto K[e]] \rightarrow h'; \mathcal{E} \uplus [i \mapsto K[e']]}$

$h; \mathcal{E} \uplus [i \mapsto K[\mathbf{fork} e]] \rightarrow h; \mathcal{E} \uplus [i \mapsto K[()]] \uplus [j \mapsto e]$

**Figure 2.** The calculus:  $F^{\text{ML}}$  with **fork**, **CAS**, and thread-local refs

higher-level abstraction. Every component of the case study involves higher-order code and therefore higher-order specifications.

## 2. The programming model

The programming language for our program logic is sketched in Figure 2; the full details are in the appendix [30]. It has a standard core: the polymorphic lambda calculus with sums, products, references, and equi-recursive types. We elide all type annotations in terms, but polymorphism is nevertheless introduced and eliminated by explicit type abstraction ( $\Lambda.e$ ) and application ( $e \_$ ). To this standard core, we add concurrency (through a **fork** construct), atomic compare-and-set (**CAS**) and thread-local references (type **refLcl**  $\tau$ ).

While normal references provide shared state through which threads can communicate, thread-local references allow multiple threads to access disjoint values in memory (a different one for each thread) using a common location. Formally, each thread-local reference has an associated *thread ID-indexed* table  $L$ , which is initially empty. Thus for **refLcl**  $\tau$ , the **getLcl** operation returns an

“option” of type  $1 + \tau$ , reflecting whether the thread-local reference has been initialized for the thread that invoked it.<sup>3</sup> Thread-local references are commonly used when a single data structure needs mutable, thread-local storage for an unbounded number of threads; the flat combining algorithm presented in §4 embodies this technique.

We define a small-step, call-by-value operational semantics, using evaluation contexts  $K$  to specify a left-to-right evaluation order within each thread. The  $\xrightarrow{\text{pure}}$  relation gives the reductions that do not interact with the heap, while  $\xrightarrow{s}$  gives general, single-thread reductions for the thread with ID  $i$ . These reductions are then lifted to general reduction  $\rightarrow$  of *machine configurations*  $\zeta$ , which consist of a heap  $h$  together with an ID-indexed pool  $\mathcal{E}$  of threads.

The type system imposes two important restrictions. First, recursive types  $\mu\alpha.\tau$  must be *guarded*, meaning that all free occurrences of  $\alpha$  in  $\tau$  must appear under a non- $\mu$  type constructor. Second, because **CAS** exposes physical equality comparison of word-sized values at the hardware level, it can only be applied to references of *comparable type*  $\sigma$ , which we take to be base types and locations.

If  $\Omega \vdash e_1 : \tau$  and  $\Omega \vdash e_2 : \tau$ , we say  $e_1$  *contextually refines*  $e_2$ , written  $\Omega \models e_1 \preceq e_2 : \tau$ , if, for every pair of thread IDs  $i$  and  $j$ ,

$$\forall C : (\Omega, \tau) \rightsquigarrow (\emptyset, \mathbf{B}). \forall n. \forall \mathcal{E}_1. \emptyset; [i \mapsto C[e_1]] \xrightarrow{*} h_1; [i \mapsto v] \uplus \mathcal{E}_1 \implies \exists \mathcal{E}_2. \emptyset; [j \mapsto C[e_2]] \xrightarrow{*} h_2; [j \mapsto v] \uplus \mathcal{E}_2$$

where  $C : (\Omega, \tau) \rightsquigarrow (\Omega', \tau')$  is the standard typing judgment for contexts (guaranteeing that  $\Omega \vdash e : \tau$  implies  $\Omega' \vdash C[e] : \tau'$ ). This definition resembles the standard one for sequential languages, but take note that only the termination behavior of the main thread is observed; as with most languages in practice, we assume that once the main thread has finished, any forked threads are aborted.

### 3. CaReSL

While CaReSL ultimately provides mixed refinement and Hoare-style reasoning about higher-order concurrent programs, it is easiest to understand by first considering less powerful “sublogics”, and then gradually incorporating more powerful features:

- We begin with a core logic that provides separation-logic reasoning for sequential (but higher-order) programs (§3.1). In this core logic, both Hoare triples and heap assertions live in a single syntactic class: they are propositions. Thus propositions include *e.g.*, implications between Hoare triples, which are ubiquitous when reasoning about higher-order code. Propositions also include first- and second-order quantification, guarded recursion, separating conjunction, and modal necessitation—all standard concepts whose interplay we explain in §3.1.
- We then add concurrency, leading to a concurrent separation logic for higher-order programs (§3.2). The new concept needed to deal with concurrency is that of a *protocol* governing some shared region of the heap. We express protocols through transition systems whose states provide an abstraction of the heap region’s actual state. Threads can gain or lose *roles* in the protocol, which determine the transitions those (and other) threads are permitted to take. Traditional rely/guarantee reasoning can be recovered as a particular, restricted way of using protocols.
- Finally, we show how to reason about refinement, yielding a **Concurrent and Refined Separation Logic (CaReSL)** for higher-order programs (§3.3). As explained in the introduction, CaReSL does not simply add refinement as a primitive notion on top of a concurrent separation logic. Instead, CaReSL treats refinement as a derived notion, expressed as a particular Hoare-style specification. In order to do so, however, CaReSL does require one further extension: the notion of “spec resources”,

first introduced by Turon *et al.* [31], which allow pieces of a program configuration (heap and threads) to be manipulated as logical resources within Hoare-style proofs.

While CaReSL is the main contribution of the paper, its sublogics are of independent interest. We explain them by way of idiomatic higher-order examples, each of which serves as a component of the case study in §4.

#### 3.1 Core logic: Sequential higher-order programs

To motivate the basic ingredients of core CaReSL, we begin with the quintessential higher-order stateful combinator, `foreach`:

$$\begin{aligned} \text{foreach} & : \forall\alpha. (\alpha \rightarrow \mathbf{1}) \rightarrow \text{list}(\alpha) \rightarrow \mathbf{1} & \text{list}(\alpha) \triangleq \mu\beta. \mathbf{1} + (\alpha \times \beta) \\ \text{foreach} & \triangleq \Lambda. \lambda f. \mathbf{rec} \text{ loop}(l). \mathbf{case} \ l \ \text{of} \ \text{none} & \Rightarrow () \\ & | \text{some}(x, n) \Rightarrow f(x); \text{loop}(n) \end{aligned}$$

A specification of `foreach` should explain that `foreach f l` “lifts” the imperative behavior of  $f$  (which works on elements) to  $l$  (a list of elements). But to do so, it needs to *quantify* over the unknown behavior of an unknown function in a way that can be lifted to lists.

One possibility is to characterize  $f$  by means of a Hoare triple, while ultimately quantifying over the pre- and post-conditions of that triple. Suppose that  $p(x)$  and  $q(x)$  are predicates on values  $x$ , and that  $r$  is an arbitrary proposition, such that

$$\forall x. \{p(x) * r\} f(x) \{q(x) * r\}$$

The idea is that when  $f$  is applied to an element  $x$  of the list, it may assume both that  $p(x)$  holds and that an invariant  $r$  holds of a disjoint portion of the heap. If  $f(x)$  transforms  $p(x)$  to  $q(x)$  while maintaining the invariant, then `foreach` lifts that behavior to an entire list—regardless of what the predicates actually are. To capture the assumption that  $p(x)$  holds of each element of a list  $l$ , we need something like the following recursive predicate:<sup>4</sup>

$$\text{Map}_p(l) \triangleq l = \text{none} \vee (\exists x, n. l = \text{some}(x, n) * p(x) * \text{Map}_p(n))$$

Note that, thanks to the use of the separating conjunction  $*$ , for each  $x$  in  $l$  there is a *disjoint region* of the heap satisfying  $p(x)$ .

These ideas lead to the following spec for `foreach`:

$$\begin{aligned} \forall p, q, r. \forall f. (\forall x. \{p(x) * r\} f(x) \{q(x) * r\}) \\ \Rightarrow \forall l. \{ \text{Map}_p(l) * r \} \text{foreach} \ f \ l \ \{ \text{Map}_q(l) * r \} \end{aligned}$$

Formalizing this sketch will require a logic with a number of basic features: we need to be able to mix Hoare triples with other kinds of connectives (*e.g.*, implication), to quantify over both values (*e.g.*,  $f$  and  $x$ ) and predicates (*e.g.*,  $p$  and  $q$ ), and to define recursive predicates (*e.g.*, `Map`). With these goals in mind, we now explore CaReSL’s design more systematically.

**Syntax** Figure 3 presents core CaReSL—the fragment of the logic appropriate for reasoning about sequential, higher-order code. CaReSL is a multi-sorted second-order logic, meaning that its syntax is stratified into the following three layers.

First, there are *terms*  $M, N$ , which come in a variety of *sorts*  $\Sigma$ , including the values and expressions of the language, as well as thread-local storage (with operations  $[-]$ ,  $\uplus$  and  $:=$ ). The judgment  $\mathcal{X} \vdash M : \Sigma$  gives the sort of a term (where  $\mathcal{X}$  is a variable context giving the sorts of term and predicate variables).

Second, there are *propositions*  $P, Q, R$ , which include the connectives of multi-sorted first-order logic (*e.g.*,  $\forall X \in \Sigma. P$ ) and second-order logic (*e.g.*,  $\exists p \in \mathbb{P}(\Sigma). P$ ) with their standard meanings. The judgment  $\mathcal{X} \vdash P : \mathbb{B}$  asserts that the proposition  $P$  is well-sorted.

Third, there are *predicates*  $\varphi, \psi$ , which are just propositions parameterized by a term: a predicate of sort  $\mathbb{P}(\Sigma)$  can be introduced by  $(X \in \Sigma). P$ , which binds an unknown term  $X$  of sort  $\Sigma$  in  $P$ ,

<sup>3</sup>Throughout, we let  $\text{none} \triangleq \text{inj}_1 ()$  and  $\text{some}(e) \triangleq \text{inj}_2 e$ .

<sup>4</sup>Throughout, names of functions and types are written in lower case sans-serif, while predicate names are Capitalized Serif.

### Syntax

Terms  $M ::= X \mid (M, M) \mid n \mid e \mid M \uplus M \mid [M](M) \mid M[M := M]$

Sorts  $\Sigma ::= \mathbf{1} \mid \Sigma \times \Sigma \mid \text{Nat} \mid \text{Val} \mid \text{Exp} \mid \text{AtomicExp} \mid \text{LclStorage}$

Preds  $\varphi ::= p \mid (X \in \Sigma).P \mid \mu p \in \mathbb{P}(\Sigma).\varphi$

Propositions  $P ::=$

True  $\mid P \Rightarrow P \mid P \wedge P \mid P \vee P \mid \exists X \in \Sigma.P \mid \forall X \in \Sigma.P \mid \triangleright P \mid \varphi(M) \mid \exists p \in \mathbb{P}(\Sigma).P \mid \forall p \in \mathbb{P}(\Sigma).P \mid P * P \mid M \hookrightarrow_1 M \mid A$

Necessary (“always”) propositions  $A ::=$

$\Box P \mid M \xrightarrow{\text{pure}} M \mid M = M \mid \langle P \rangle M \Rightarrow_{\text{is}} M \mid \langle \varphi \rangle \mid \{P\} M \Rightarrow M \mid \{\varphi\}$

### Contexts

$\mathcal{X} ::= \cdot \mid \mathcal{X}, X : \Sigma \mid \mathcal{X}, p : \mathbb{P}(\Sigma) \quad \mathcal{P} ::= \cdot \mid \mathcal{P}, P \quad \mathcal{C} ::= \mathcal{X}; \mathcal{P}$

### Well-sorted terms

$\mathcal{X} \vdash M : \Sigma$   
 $\mathcal{X}, X : \Sigma \vdash X : \Sigma$   
 $\frac{\mathcal{X} \vdash M : \text{LclStorage} \quad \mathcal{X} \vdash N : \text{Nat} \quad \dots}{\mathcal{X} \vdash [M](N) : \text{Val}}$

**Well-sorted props.**  $\mathcal{X} \vdash P : \mathbb{B}$     **Well-sorted preds.**  $\mathcal{X} \vdash \varphi : \mathbb{P}(\Sigma)$

Figure 3. Core CaReSL: Syntax

and eliminated by  $\varphi(M)$  (also written  $M \in \varphi$ ) which substitutes  $M$  for the parameter of  $\varphi$ . Because sorts include unit and products, predicates can express *relations* of arbitrary arity. The judgment  $\mathcal{X} \vdash \varphi : \mathbb{P}(\Sigma)$  asserts that  $\varphi$  is a well-sorted predicate over terms of sort  $\Sigma$ .

In general, term variables are written  $X$ . But to avoid cluttering our rules and proofs with sort annotations, we use variables  $x, y, z$  for sort Val and  $i, j, k, l$  for sort Nat. We abuse notation, writing *e.g.*,  $v$ ,  $e$  or  $L$  to stand for a term of sort Val, Exp or LclStorage, respectively.

In addition to this standard logical setup, CaReSL adopts key connectives from separation logic. In general, these connectives will refer to a variety of “resources” that will be introduced as we go along. In *core* CaReSL, however, the only resource is the heap. The *points-to* assertion  $\ell \hookrightarrow_1 v$  holds of any heap containing a location  $\ell$  that points to the value  $v$ . (Ignore the 1 subscript for now; we will return to it in §3.3). The *separating conjunction*  $P * Q$  holds if the currently-owned resources (here, a portion of the heap) can be split into two disjoint parts satisfying  $P$  and  $Q$  respectively.

While the truth of some propositions (*e.g.*,  $\ell \hookrightarrow_1 v$ ) is *contingent* on the presence of certain resources, others (*e.g.*,  $M = N$ ) are *necessary*: if they hold, they do so regardless of the currently-owned resources, and therefore will continue to hold in any future state. (Such propositions are called “pure” in separation logic parlance.) The syntactic subcategory  $A$  of necessary propositions includes claims about term equality, about the operational semantics ( $M \xrightarrow{\text{pure}} N$ ), and Hoare triples. Arbitrary propositions can be made necessary via the  $\Box$  (“always”) modality, where  $\Box P$  holds if  $P$  holds for all possible resources. As we will see shortly, necessary propositions play by special rules: they can move freely through Hoare triples and separating conjunctions.

Ultimately, CaReSL distinguishes between triples about general expressions  $e$  and those about *atomic* expressions  $a$  (which execute in a single step). Since this distinction is motivated by concurrency, we postpone its explanation to §3.2. We include the distinction syntactically in core CaReSL, but it can be safely ignored for now.

Atomic expressions  $a$  have the following grammar:

**new**  $v \mid$  **get**  $v \mid v := v \mid$  **CAS**  $(v, v, v) \mid$  **newLcl**  $\mid$  **getLcl**  $(v) \mid$  **setLcl**  $(v, v)$

The propositions for atomic triples<sup>5</sup>  $\langle P \rangle i \Rightarrow_{\text{is}} a \langle \varphi \rangle$  and general triples  $\{P\} i \Rightarrow e \{\varphi\}$  are both parameterized by a *thread ID*  $i$ ; the expression may access thread-local storage, in which case its behavior is ID-dependent. (When the thread ID doesn’t matter, we

<sup>5</sup> Again, ignore the IS subscript until §3.3.

### Logical axioms and rules

$$\frac{\mathcal{X}; A \vdash P}{\mathcal{X}; \mathcal{P}, A \vdash \Box P} \Box\text{I} \quad \frac{\mathcal{C} \vdash P}{\mathcal{C} \vdash \triangleright P} \text{MONO} \quad \frac{\mathcal{C}, \triangleright P \vdash P}{\mathcal{C} \vdash P} \text{LÖB}$$

$$\Box(P \Rightarrow Q) \Rightarrow \Box P \Rightarrow \Box Q \quad A * P \Leftrightarrow A \wedge P$$

$$A \Rightarrow \Box A \quad \text{True} * P \Leftrightarrow P$$

$$\Box P \Rightarrow P \quad \triangleright(P * Q) \Leftrightarrow \triangleright P * \triangleright Q$$

$$\frac{\mathcal{C}, p : \mathbb{P}(\Sigma) \vdash P}{\mathcal{C} \vdash \forall p \in \mathbb{P}(\Sigma). P} \forall_2\text{I} \quad M \in (X \in \Sigma).P \Leftrightarrow P[M/X]$$

$$M \in (\mu p.\varphi) \Leftrightarrow M \in \varphi[\mu p.\varphi/p]$$

Figure 4. Core CaReSL: Selected logical axioms and rules

write  $\{P\} e \{\varphi\}$  as short for  $\forall i. \{P\} i \Rightarrow e \{\varphi\}$ .) In addition, since we are working with an expression language (as opposed to a command language), postconditions are *predicates* over the return value of the expression, rather than simply propositions about the final state of the heap. But when an expression returns unit, we often abuse notation and write a proposition instead.

In order to support recursive assertions, the logic includes *guarded recursion* ( $\mu p \in \mathbb{P}(\Sigma). P$ ), which entails the following tradeoff. On the one hand, guarded recursion allows occurrences of the recursive predicate  $p$  to appear negatively, which is crucial for modelling recursive types (§3.3) but is usually prohibited for lack of monotonicity. On the other hand, the recursion is “guarded” in that references to  $p$  in  $P$  must appear under the  $\triangleright$  modality. A proposition  $\triangleright Q$  represents the present knowledge that  $Q$  holds *later*, *i.e.*, after at least one step of computation. Guarded recursion supports a coinductive style of reasoning: to prove  $P$  one can assume  $\triangleright P$ , but this assumption is only useful after at least one step of computation. As we explain in §5, our use of guarded recursion descends from a line of work on *step-indexed* logical relations, but the interaction with Hoare triples is a novelty of CaReSL.

**Proof rules** The main judgment of CaReSL is written  $\mathcal{C} \vdash P$ , where  $\mathcal{C} = \mathcal{X}; \mathcal{P}$  is a *combined* context of annotated term/predicate variables  $\mathcal{X}$  and propositional assumptions  $\mathcal{P}$ . The meaning is the usual one: for any way of instantiating the variables in  $\mathcal{X}$ , if the hypotheses  $\mathcal{P}$  are true for a given resource (*i.e.*, for the moment, a given heap), then  $P$  is true of the same resource. We implicitly assume  $\mathcal{X} \vdash Q : \mathbb{B}$  for every proposition  $Q$  in  $\mathcal{P}$  and likewise that  $\mathcal{X} \vdash P : \mathbb{B}$  holds.

A few basic logical axioms and proof rules for core CaReSL are shown in Figure 4. Axioms hold in an arbitrary well-sorted context. The axioms include all the standard ones for a multi-sorted second-order logic (we show only  $\forall_2\text{I}$ ), as well as several characterizing separating conjunction and the  $\Box$  and  $\triangleright$  modalities. We’ll just mention the highlights:

- Because True is a unit for separating conjunction (and every proposition implies True), propositions are *affine*: we can “throw away” resources, because  $(P * Q) \Rightarrow (\text{True} * Q) \Rightarrow Q$ .
- The two conjunctions  $*$  and  $\wedge$  are identical if at least one of their operands is a necessary proposition. Consequently, necessary propositions can be “freely copied”:  $A \Rightarrow A * A$ .
- The LÖB rule provides a coinductive reasoning principle: to prove  $P$ , you may assume  $P$ —but only under the  $\triangleright$  modality, which guards use of the assumption until at least one step of computation has taken place. On the other hand, MONO says that any proposition can be *weakened* to one that is guarded. (Both MONO and LÖB are inherited from LADR [5].) We will momentarily see how  $\triangleright$  is eliminated in Hoare-style reasoning.
- The  $\triangleright$  modality distributes over  $*$ , as it does with most other connectives, except implication and recursion.

$$\begin{array}{c}
\text{Atomic Hoare logic} \quad (\text{where } \text{IS} ::= \text{I} \mid \text{S} \text{ as explained in } \S 3.3) \\
\begin{array}{l}
\langle \text{True} \rangle i \Vdash_{\text{IS}} \text{new } v \quad \langle \text{ret. ret} \hookrightarrow_{\text{IS}} v \rangle \\
\langle v \hookrightarrow_{\text{IS}} v' \rangle i \Vdash_{\text{IS}} \text{get } v \quad \langle \text{ret. ret} = v' * v \hookrightarrow_{\text{IS}} v' \rangle \\
\langle v \hookrightarrow_{\text{IS}} - \rangle i \Vdash_{\text{IS}} v := v' \quad \langle \text{ret. ret} = () * v \hookrightarrow_{\text{IS}} v' \rangle \\
\langle \text{True} \rangle i \Vdash_{\text{IS}} \text{newLcl} \quad \langle \text{ret. ret} \hookrightarrow_{\text{IS}} \emptyset \rangle \\
\langle v \hookrightarrow_{\text{IS}} L \rangle i \Vdash_{\text{IS}} \text{getLcl}(v) \quad \langle \text{ret. ret} = [L](i) * v \hookrightarrow_{\text{IS}} L \rangle \\
\langle v \hookrightarrow_{\text{IS}} L \rangle i \Vdash_{\text{IS}} \text{setLcl}(v, v') \quad \langle \text{ret. ret} = () * v \hookrightarrow_{\text{IS}} L[i := v'] \rangle \\
\langle v \hookrightarrow_{\text{IS}} v' \rangle i \Vdash_{\text{IS}} \quad \langle \text{ret. } (v' = v_o * \text{ret} = \text{true} * v \hookrightarrow_{\text{IS}} v_n) \rangle \\
\text{CAS}(v, v_o, v_n) \quad \vee \quad \langle \text{ret. } (v' \neq v_o * \text{ret} = \text{false} * v \hookrightarrow_{\text{IS}} v') \rangle
\end{array} \\
\frac{\mathcal{X}; P \vdash P' \quad \mathcal{X}; \mathcal{P} \vdash (P') i \Vdash_{\text{IS}} a \langle x. Q' \rangle \quad \mathcal{X}, x; Q' \vdash Q}{\mathcal{X}; \mathcal{P} \vdash (P) i \Vdash_{\text{IS}} a \langle x. Q \rangle} \text{ACSQ} \\
\frac{C \vdash (P) i \Vdash_{\text{IS}} a \langle x. Q \rangle}{C \vdash (P * R) i \Vdash_{\text{IS}} a \langle x. Q * R \rangle} \text{AFRAME} \quad (\text{ADISJ, AEX elided}) \\
\text{General Hoare logic} \\
\frac{C \vdash (P) i \Vdash_{\text{IS}} a \langle Q \rangle}{C \vdash \{\triangleright P\} i \Vdash a \{Q\}} \text{PRIVATE} \quad \frac{C \vdash e \xrightarrow{\text{pure}} e'}{C \vdash \{P\} i \Vdash e' \{ \varphi \}} \text{PURE} \\
\frac{C \vdash \{P\} i \Vdash e \{x. Q\} \quad C, x \vdash \{Q\} i \Vdash K[x] \{R\}}{C \vdash \{P\} i \Vdash K[e] \{R\}} \text{BIND} \\
C \vdash \{\text{True}\} i \Vdash v \{x. x = v\} \text{RET} \quad (\text{CSQ, FRAME, DISJ, EX elided}) \\
\frac{C \vdash A}{C \vdash \{P * A\} i \Vdash e \{ \varphi \}} \text{AIN} \quad \frac{C, A \vdash \{P\} i \Vdash e \{ \varphi \}}{C \vdash \{P * A\} i \Vdash e \{ \varphi \}} \text{AOUT} \\
\text{Derivable rules} \\
\frac{C, f, \forall x. \{\triangleright P\} i \Vdash f x \{ \varphi \} \vdash \forall x. \{P\} i \Vdash e \{ \varphi \}}{C \vdash \forall x. \{\triangleright P\} i \Vdash (\text{rec } f(x).e) x \{ \varphi \}} \text{REC} \\
\frac{C \vdash \{P\} i \Vdash e \{x. \exists y. (x = \text{inj}_1 y * \triangleright Q_1) \vee (x = \text{inj}_2 y * \triangleright Q_2)\} \quad \forall k \in \{1, 2\}: C, x, y \vdash \{x = \text{inj}_k y * Q_k\} i \Vdash e_k \{ \varphi \}}{C \vdash \{P\} i \Vdash \text{case}(e, \text{inj}_1 y \Rightarrow e_1, \text{inj}_2 y \Rightarrow e_2) \{ \varphi \}}
\end{array}$$

Figure 5. Core CaReSL: Hoare logic

The rules for atomic triples (Figure 5) are formulated in the standard style of separation logic. They transcribe the operational semantics of atomic expressions, mentioning only the part of the heap relevant to the expression. Atomic Hoare logic supports the usual rules—consequence, framing, disjunction,  $\exists$ -elimination—with one important exception: it does not support a sequencing rule, since a sequence of atomic expressions is not atomic.

All atomic expressions take exactly one step to execute, and in so doing allow us to peel off a layer of the  $\triangleright$  modality. To cut down on clutter, the precondition in an atomic triple is *implicitly* understood as being under one  $\triangleright$  modality. The rule PRIVATE, which lifts an atomic triple to a general one, makes this implicit assumption explicit. (In core CaReSL, *all* resources are private; §3.2 adds shared resources.) To handle pure reduction steps (like  $\beta$ -reduction), the PURE rule appeals directly to the operational semantics using the necessary proposition  $e \xrightarrow{\text{pure}} e'$ . The rest of the rules for general Hoare triples are mostly standard; we show only the nonstandard rules in Figure 5.

In an expression language, the monadic nature of Hoare logic becomes visible: the BIND rule replaces the usual rule of sequencing, while RET is used to inject a value into a Hoare triple.

We also have a rule allowing necessary propositions to move freely from the proof context into preconditions (AIN), and vice versa (AOUT). In general, any contingent assumptions like  $x \hookrightarrow_1 \exists$  need to be given explicitly in the precondition of a Hoare triple, because the truth of such statements can change over time; the triple

says that it is only usable at times when its precondition holds. But in the specific case of necessary propositions, we can do better: we know that if such a proposition happens to be true now, it will be true forever, and so it does not need to be given explicitly in the precondition. As we will see in §3.2, these rules will allow us to completely hide pieces of state that are known to always obey a certain protocol.

Finally, using LÖB and PURE, together with standard Hoare rules, we can derive specialized rules for constructs like recursive functions and pattern matching—the two derived rules in Figure 5.

**Verifying foreach** Having seen core CaReSL in detail, we can now return to the foreach example. First, we need to rewrite our sketch of the specification more formally. The Map predicate needs to employ guarded recursion:

$$\begin{aligned}
\text{Map}_\varphi &\triangleq \mu m \in \mathbb{P}(\text{Val}). (l \in \text{Val}). \\
& l = \text{none} \vee (\exists x, n. l = \text{some}(x, n) * \varphi(x) * \triangleright m(n))
\end{aligned}$$

while the foreach spec should be annotated with sorts:

$$\begin{aligned}
&\forall p, q \in \mathbb{P}(\text{Val}). \forall r \in \mathbb{P}(\mathbf{1}). \forall f. (\forall x. \{p(x) * r\} f(x) \{q(x) * r\}) \\
&\Rightarrow \forall l. \{\text{Map}_p(l) * r\} \text{foreach } \_ f l \{\text{Map}_q(l) * r\}
\end{aligned}$$

To prove foreach correct, we use a kind of Hoare “proof outline”, annotating each program point with a proposition:

$\text{foreach} \triangleq \Lambda. \lambda f. \text{rec loop}(l).$

<b>Prop context:</b> $\forall x. \{p(x) * r\} i \Vdash f(x) \{q(x) * r\}$ $\forall n. \{\triangleright (\text{Map}_p(n) * r)\} i \Vdash \text{loop } n \{\text{Map}_q(n) * r\}$	<b>Variables:</b> $f, i, p, q, r, l, \text{loop}$
--	---

$$\begin{aligned}
&\{\text{Map}_p(l) * r\} \\
&\text{case } l \text{ of none} \Rightarrow \{l = \text{none} * \text{Map}_p(l) * r\} () \{\text{Map}_q(l) * r\} \\
& \quad | \text{some}(x, n) \Rightarrow \{l = \text{some}(x, n) * \text{Map}_p(l) * r\} \\
& \quad \quad \{l = \text{some}(x, n) * p(x) * \triangleright \text{Map}_p(n) * r\} \\
& \quad \quad f(x); \\
& \quad \quad \{l = \text{some}(x, n) * q(x) * \triangleright \text{Map}_p(n) * r\} \\
& \quad \quad \text{loop}(n) \\
& \quad \quad \{l = \text{some}(x, n) * q(x) * \text{Map}_q(n) * r\} \\
& \quad \quad \{\text{Map}_q(l) * r\}
\end{aligned}$$

Proof outlines implicitly apply the Hoare logic rule corresponding to each language construct: for recursive functions and pattern matching we use the derived rules shown earlier; for an atomic expression we use the corresponding axiom; when going under a  $\Lambda$  or  $\lambda x$  we use the PURE rule and prove a triple about the function as applied to  $\_$  or  $x$  respectively. The frame rule is applied implicitly as needed, while uses of consequence are shown by writing a sequence of assertions (each one implying the next). We do not explicitly write the thread ID in a proof outline, but it is always clear from context (and, for our examples, always the variable  $i$ ).

We supplement traditional proof outlines with boxed *context assertions* spelling out an extension to both the variable context  $\mathcal{X}$  and proposition context  $\mathcal{P}$ . Extending the context with new variables introduces a universal quantification (using *e.g.*, the rule  $\forall_2\text{I}$  in Figure 4), while adding propositions introduces an implication. We use the AIN rule implicitly to bring necessary consequences of the proof context into the Hoare-style outline.

The proof for foreach is quite straightforward: the initial case analysis on the input list allows us to expand the definition of the Map predicate, which for the nonempty case gives us the necessary knowledge to execute  $f$  on an element of the list; note the heavy use of the frame rule to add invariant propositions. The use of  $\triangleright$  in Map is neatly dispatched by the use of the REC rule for foreach.

### 3.2 Adding concurrency: Protocols for shared state

To reason about concurrency, we need to reason about the protocols governing shared (and often hidden) state. Take, for example, the following higher-order spinlock:

$\text{tryAcq} \triangleq \lambda x. \text{CAS}(x, \text{false}, \text{true})$   
 $\text{acq} \triangleq \text{rec loop}(x). \text{if tryAcq}(x) \text{ then } () \text{ else loop}(x)$   
 $\text{rel} \triangleq \lambda x. x := \text{false}$   
 $\text{mkSync} : \mathbf{1} \rightarrow \forall \alpha. \forall \beta. (\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \beta)$   
 $\text{mkSync} \triangleq \lambda(). \text{let lock} = \text{new false in}$   
 $\quad \Lambda. \Lambda. \lambda f. \lambda x. \text{acq}(\text{lock}); \text{let } z = f(x) \text{ in rel}(\text{lock}); z$

Each invocation of  $\text{mkSync}$  creates a new *wrapper* that closes over a fresh, hidden lock. The wrapper can then be used to add simple synchronization to an arbitrary function. There are, of course, a variety of ways to use synchronization, but a particularly common one is to protect access to some shared resource characterized by an invariant  $p$ —an idea that leads to the following specification:

$\forall p \in \mathbb{P}(\mathbf{1}). \{p\} \text{mkSync } () \{s. \text{Syncer}_p(s)\}$   
 $\text{where } \text{Syncer}_p \triangleq (s \in \text{Val}). \square \forall f, x, q, r. \{p * q\} f x \{z. p * r(z)\}$   
 $\quad \Rightarrow \{q\} s \_ \_ f x \{z. r(z)\}$

The spec says that when  $\text{mkSync}()$  is executed, the client can choose some invariant resource  $p$ , giving up control over the resource in exchange for a synchronization wrapper  $s$ . When  $s$  is later applied to a function  $f$ , it provides  $f$  with exclusive access to the resource  $p$  (seemingly out of thin air), which  $f$  can use however it pleases, so long as the invariant is reestablished on completion.

Intuitively, the reason that  $\text{mkSync}$  satisfies its specification is that the lock itself is *hidden*: all access to it is mediated through the wrapper  $s$ , which the client can only apply to invariant-preserving functions. Hiding enables  $\text{mkSync}$  to maintain an internal *protocol* in which, whenever the lock is free, the invariant  $p$  holds. To express this protocol, as well as the more sophisticated protocols needed for fine-grained concurrency (§3.3, §4), CaReSL provides a syntactic account of the semantic protocols of Turon *et al.* [31].

**Protocols** A protocol  $\pi$  governs a shared resource abstractly, by means of a set of protocol states ( $S$ ) equipped with a transition relation ( $\rightsquigarrow$ ). Each state  $s \in S$  has an associated proposition  $\varphi(s)$  giving its concrete interpretation in terms of *e.g.*, heap resources. The idea is then that any thread is allowed to update the shared resource, so long as at each atomic step those concrete updates correspond to a permitted abstract transition.

In general, although all threads must abide by a given protocol, not all of them play the same role within it. For example, a protocol governing a lock might have two states, Locked and Unlocked, with transitions in both directions, but we usually want to allow only the thread that acquired the lock to be allowed to release it (as is the case in the  $\text{mkSync}$  example). To allow threads to take on or relinquish roles dynamically, protocols employ *tokens* that individual threads may own. By taking a transition, a thread may “earn” a token from the protocol; conversely, certain transitions require a thread to “pay” by giving up a previously-earned token. For the locking protocol, we use a single token *Lock*:



In the Unlocked state, the protocol itself owns the *Lock*, which we show by annotating the state with the token. A thread that transitions from Unlocked to Locked then earns (takes ownership of) the *Lock*. Conversely, to transition back to the Unlocked state, a thread must transfer the *Lock* back to the protocol—a move only possible for the thread that owns the *Lock*.

Formally, a protocol  $\pi = (S, \rightsquigarrow, \mathcal{T}, \varphi)$  is given by a transition system ( $S$  and  $\rightsquigarrow$ ), a function  $\mathcal{T}$  giving the set of tokens owned by the protocol at each state, and a predicate  $\varphi$  on states giving their interpretations. To be able to talk about states and tokens in the logic, we add sorts State and TokSet, for which we will use the metavariables  $s$  and  $T$  respectively; see Figure 6. We leave the grammar of terms for these sorts open-ended, implicitly extending

## Protocols

Protocol  $\pi ::= (S, \rightsquigarrow, \mathcal{T}, \varphi)$  where  $S \subseteq \text{State}$ ,  $\rightsquigarrow \subseteq S \times S$ ,  
 $\mathcal{T} \in S \rightarrow \text{TokSet}$ ,  $\varphi$  token-pure

$(s; T) \rightsquigarrow_{\pi} (s'; T') \triangleq s (\pi. \rightsquigarrow) s' \wedge \pi. \mathcal{T}(s) \uplus T = \pi. \mathcal{T}(s') \uplus T'$   
 $\text{frame}_{\pi}(s; T) \triangleq (s; \text{AllTokens} - (\pi. \mathcal{T}(s) \uplus T))$   
 $(s; T) \sqsubseteq_{\pi}^{\text{guar}} (s'; T') \triangleq (s; T) \rightsquigarrow_{\pi}^* (s'; T')$   
 $(s; T) \sqsubseteq_{\pi}^{\text{rely}} (s'; T') \triangleq \text{frame}_{\pi}(s; T) \rightsquigarrow_{\pi}^* \text{frame}_{\pi}(s'; T')$

## Syntax

Sort  $\Sigma ::= \dots \mid \text{State} \mid \text{TokSet}$   
Prop  $P ::= \dots \mid \boxed{M}_{\pi}^M \mid \text{tid}(M)$   
AbProp  $A ::= \dots \mid M \sqsubseteq_{\pi}^{\text{rely}} M \mid M \sqsubseteq_{\pi}^{\text{guar}} M \mid \text{TokPure}(P)$

## Hoare logic

(where  $\pi[b] \triangleq \exists s. b = (s, -) \wedge \pi. \varphi(s)$ )

$\frac{C \vdash \{P\} i \Rightarrow e \{x. Q * \pi[b]\}}{C \vdash \{P\} i \Rightarrow e \{x. Q * \exists n. \boxed{b}_{\pi}^n\}}$  NEWISL

$\frac{C \vdash \forall b \sqsubseteq_{\pi}^{\text{rely}} b_0. (\pi[b] * P) i \Rightarrow_1 a \{x. \exists b' \sqsubseteq_{\pi}^{\text{guar}} b. \pi[b'] * Q\}}{C \vdash \{\boxed{b_0}_{\pi}^n * \triangleright P\} i \Rightarrow a \{x. \exists b'. \boxed{b'}_{\pi}^n * Q\}}$  UPDISL

$\frac{C \vdash \{P * \text{tid}(j)\} j \Rightarrow e \{\text{ret. ret} = ()\}}{C \vdash \{P\} i \Rightarrow \text{fork } e \{\text{ret. ret} = ()\}}$  FORK

## Logical axioms and rules

$\boxed{s_1; T_1}_{\pi}^n * \boxed{s_2; T_2}_{\pi}^n \Leftrightarrow \exists s. \boxed{s; T_1 \uplus T_2}_{\pi}^n \wedge (s; T_1) \sqsubseteq_{\pi}^{\text{rely}} (s_1; T_1) \wedge (s; T_2) \sqsubseteq_{\pi}^{\text{rely}} (s_2; T_2)$  SPLITISL

Figure 6. CaReSL: Incorporating concurrency and protocols

it as needed for particular transition systems.<sup>6</sup> Thus, we can give an interpretation  $\text{LockInterp}_p$  for the lock protocol that is appropriate for an instance of  $\text{mkSync}$  protecting an invariant  $p$ :

$\text{LockInterp}_p \triangleq (s \in \text{State}). s = \text{Locked} \vee (s = \text{Unlocked} * p)$

The combination of transition systems and tokens gives rise to *token-sensitive transitions*. A transition from state/tokens  $(s; T)$  to state/tokens  $(s'; T')$  is permitted by  $\pi$ , written  $(s; T) \rightsquigarrow_{\pi} (s'; T')$ , so long as the *law of token conservation* holds: the disjoint union of the thread’s tokens and the protocols tokens  $T \uplus \pi. \mathcal{T}(s)$  before the transition must be the same as the disjoint union  $T' \uplus \pi. \mathcal{T}(s')$  afterwards.<sup>7</sup> Token-sensitive transitions constrain both what a thread can do (the *guarantee* moves  $\sqsubseteq_{\pi}^{\text{guar}}$  enabled by its tokens) and what its environment can do (the *rely* moves  $\sqsubseteq_{\pi}^{\text{rely}}$  enabled by the tokens owned by other threads). See Figure 6.

**Island assertions** In CaReSL, all resources are either privately owned by a thread, or else governed by a shared protocol. When a heap assertion like  $x \hookrightarrow_1 3$  appears in the pre- or postcondition of a triple, it is understood as asserting private ownership over the corresponding portion of the heap; no other thread is allowed to access it. Thus the rules of core CaReSL are immediately sound in a concurrent setting: there is no interference to account for.

To talk about shared resources, CaReSL includes *island assertions*  $\boxed{b}_{\pi}^n$  (similar to region assertions in RGSep/CAP). As the name suggests, each island is an independent region of the heap governed by its own laws (the protocol  $\pi$ ). The number  $n$  is the “name” of the island, which is used to distinguish between multiple islands with the same protocol; we often leave off the name when

<sup>6</sup> To be completely formal, we could allow each transition system to come equipped with its own explicit grammar of states and tokens.

<sup>7</sup> We use notation like  $\pi. \mathcal{T}$  to extract named components from a tuple.

it is existentially quantified. The term  $b$  (of sort  $\text{State} \times \text{TokSet}$ ) asserts private ownership of a set of tokens,<sup>8</sup> and acts as a “rely-lower bound” on the state of the protocol: the current state of the protocol is some  $b' \sqsubseteq_{\pi}^{\text{rely}} b$ . Thus, in CaReSL every assertion about shared resources is automatically stable under interference.

While a thread’s island assertions cannot be invalidated by a change its environment makes, they *can* be invalidated by a change the thread itself makes. For example, if a thread owns *Lock* in Locked state of the mkSync protocol, the environment cannot change the state at all—but the thread itself can move to the Unlocked state. There is, however, a special class of island assertions which are “completely” stable, *i.e.*, that act as necessary propositions: island assertions that do not claim any tokens. To understand why, consider that when no tokens are owned, the  $\sqsubseteq_{\pi}^{\text{rely}}$  relation degenerates to the underlying transition system of the protocol, which means it contains all possible moves that any thread can make. Formally, we have  $\boxed{(M; \emptyset)}_{\pi}^n \Rightarrow \square \boxed{(M; \emptyset)}_{\pi}^n$ .

When island assertions *do* claim ownership over tokens (a form of resource), they can be meaningfully combined by separating conjunction; see the SPLITISL rule in Figure 6, which takes into account the fact that the assertions give only rely-lower bounds.<sup>9</sup>

There are two Hoare logic rules for working with islands (Figure 6). The rule NEWISL creates a new island starting with an initial state and token set bound  $b$ ; the resources necessary to satisfy the protocol’s initial state must be present as private resources, and afterwards will be shared. (The shorthand  $\pi \llbracket b \rrbracket$  just gives the interpretation at the state in  $b$ .) Once an island is established, the only way to access the shared resources it governs is through the UPDISL rule, which can only be applied to an *atomic* expression. Starting from an initial bound  $b_0$ , the atomic expression  $a$  might be (instantaneously) executed in any rely-reachable  $b$ ; for each such state, the expression is granted exclusive access to the shared resource, but in its single atomic step, it must make a guarantee move in the protocol. This rule reveals the important semantic difference between atomic and general triples: through the UPDISL rule, atomic triples gain access to the concrete resources owned by shared islands, while general triples only have access to island assertions.

**Thread creation** In a protocol, threads gain and lose roles (tokens) by making deliberate moves within the protocol. But there is also a role that every thread plays *automatically*: the role of being a thread with a certain ID. To support protocols that use thread IDs (such as the one in §4), CaReSL builds in a proposition  $Tid(j)$  that asserts the existence of a thread  $j$ , and acts as an *uncopyable* resource:  $Tid(j) * Tid(j) \Leftrightarrow \text{False}$ . The resource is introduced by the FORK rule, which also allows the parent thread to transfer an arbitrary  $P$  to the newly-forked thread. (The parent can keep resources for itself via the FRAME rule.) The trivial postconditions in FORK may seem alarming, but they reflect the fact that the language has no **join** mechanism: the *only* form of communication between threads is shared state, mediated by a shared protocol.

**Verifying mkSync** The interpretations  $\pi.\varphi(s)$  of protocol states in CaReSL inherit a limitation from their semantic treatment by Turon *et al.* [31]: they must be token-pure, *i.e.*, they cannot assert ownership of island tokens. The necessary proposition  $\text{TokPure}(P)$  can be used to assert that, for example, an unknown proposition is token-pure and thus safe to use in an interpretation, and we need such a restriction on  $p$  in mkSync:

$$\forall p \in \mathbb{P}(\mathbf{1}). \text{TokPure}(p) \Rightarrow \{p\} \ i \Rightarrow \text{mkSync}() \ \{s. \text{Syncer}_p(s)\} \\ \text{where } \text{Syncer}_p \triangleq (s \in \text{Val}). \ \square \forall f, x, q, r. \ \{p * q\} \ f \ x \ \{z. p * r(z)\} \\ \Rightarrow \{q\} \ s \ \_ \_ \ f \ x \ \{z. r(z)\}$$

<sup>8</sup> An island assertion is only satisfied if the asserted token set is disjoint from the tokens owned by the protocol at the asserted state.

<sup>9</sup> We assume that terms  $M$  include disjoint union on token sets.

The hidden protocol  $\text{LockProt}_p$  for mkSync just puts together the pieces we have already seen:

$$\text{LockProt}_p \triangleq (\{\text{Locked}, \text{Unlocked}\}, \rightsquigarrow, \mathcal{T}, \text{LockInterp}_p)$$

where  $\rightsquigarrow$  and  $\mathcal{T}$  are given by the transition system for locking shown above.

The high-level proof outline for mkSync is straightforward:

$$\lambda(). \ \boxed{\text{Prop context: TokPure}(p) \quad \text{Variables: } p} \\ \{p\} \ \text{let lock} = \text{new false} \ \{p * \text{lock} \hookrightarrow, \text{false}\} \\ \{\exists n. \ \boxed{\text{Unlocked}; \emptyset}_{\text{LockProt}_p}^n\} \\ \Lambda. \Lambda. \lambda f. \lambda x. \\ \boxed{\text{Prop context:} \quad \text{Variables: } f, x, q, r} \\ \square \ \boxed{\text{Unlocked}; \emptyset}_{\text{LockProt}_p} \ \{p * q\} \ f \ x \ \{z. p * r(z)\} \\ \{q\} \ \{\boxed{\text{Unlocked}; \emptyset}_{\text{LockProt}_p} * q\} \\ \text{acq}(\text{lock}); \ \{\boxed{\text{Locked}; \text{Lock}}_{\text{LockProt}_p} * p * q\} \\ \text{let } z = f(x) \ \{\boxed{\text{Locked}; \text{Lock}}_{\text{LockProt}_p} * p * r(z)\} \\ \text{rel}(\text{lock}); \ \{\boxed{\text{Unlocked}; \emptyset}_{\text{LockProt}_p} * r(z)\} \\ z \ \{z. r(z)\}$$

After allocating the hidden lock, we move it into a fresh island using NEWISL, and then move that island assertion (wrapped with  $\square$ ) into the proof context (using AOUT) before reasoning about the returned wrapper function. In this way the protocol is completely hidden from the client, yet available to the closure we return to the client—mirroring the fact that lock is hidden from the client but available in the closure. Since the spec  $\text{Syncer}_p(s)$  is a necessary proposition, the client may move the spec into its proof context, and thus freely use the synchronization wrapper  $s$  without threading *any* assertion about it through its proof. (Previous logics, like CAP, require at least that the client thread through an abstract predicate standing for the lock.)

When verifying the closure, we begin with a precondition  $q$  of the client’s choice, but then apply AIN to load the hidden island assertion into the precondition. The subsequent lines use the locking protocol to acquire the resource  $p$ , execute the client’s function  $f$ , and then return  $p$  to the protocol. The tokens in the island assertions, which say what the *thread* owns at each point, complement those (in  $\pi.\mathcal{T}$ ) labelling the corresponding states in the protocol.

The triples for acq and rel must ultimately be proved by appeal to the UPDISL rule. For acq, the *bound* on the protocol is just  $(\text{Unlocked}; \emptyset)$ , which means that the *actual* state might be either Locked or Unlocked. The **CAS** within acq will return:

- **true** if successful; the state must have been Unlocked. We make a guarantee move to  $(\text{Locked}; \text{Lock})$ , gaining the token.
- **false** if it fails; the state must have been Locked. We do not change the state, and the acq function retries by calling loop.

On the other hand, for rel, the bound  $(\text{Locked}; \text{Lock})$  means that the protocol *must* be in state Locked: the environment cannot move to Unlocked because it does not own *Lock*. But since rel *does* own *Lock*, it can simply update the lock to **false**, corresponding to a guarantee move to the Unlocked state in the protocol.

### 3.3 Adding refinement: Reasoning about specification code

At this point, we have seen the fragment of CaReSL providing Hoare-style specs and proofs for higher-order concurrent programs, as exemplified (in a simple way) by the mkSync example. This section explains the other major component of the logic: higher-order granularity abstraction, exemplified (in a simple way) by the Treiber stack with iterator.



## Syntax

$$\begin{array}{l} \Sigma ::= \dots \mid \text{EvalCtx} \quad P ::= \dots \mid M \hookrightarrow_s M \mid M \mapsto_s M \\ M ::= \dots \mid K \mid M[M] \quad A ::= \dots \mid P \rightarrow_s P \mid M \bowtie M \end{array}$$

## Spec rewriting

$$\frac{\mathcal{C} \vdash e \xrightarrow{\text{pure}} e'}{\mathcal{C} \vdash (P * j \mapsto_s K[e]) \rightarrow_s (P * j \mapsto_s K[e'])} \text{SPURE}$$

$$\frac{\mathcal{C} \vdash (P) j \mapsto_s a \langle x. Q \rangle}{\mathcal{C} \vdash (P * j \mapsto_s K[a]) \rightarrow_s (\exists x. Q * j \mapsto_s K[x])} \text{SPRIM}$$

$$\mathcal{C} \vdash (j \mapsto_s K[\text{fork } e]) \rightarrow_s (j \mapsto_s K[()] * k \mapsto_s e) \text{SFORK}$$

## Hoare logic

(AEXECSPEC elided)

$$\frac{\mathcal{C} \vdash \{P\} i \mapsto e \{x. Q\} \quad \mathcal{C} \vdash Q \rightarrow_s R}{\mathcal{C} \vdash \{P\} i \mapsto e \{x. R\}} \text{EXECSPEC}$$

$$\frac{\mathcal{C}, j \bowtie j' \vdash \{P * \text{Tid}(j) * j' \mapsto_s e_s\} j \mapsto_{e_1} \{x. x = ()\}}{\mathcal{C} \vdash \{P * i' \mapsto_s K[\text{fork } e_s]\} i \mapsto_{\text{fork } e_1} \{i' \mapsto_s K[()]\}} \text{FORKS}$$

Figure 7. CaReSL: Incorporating spec resources

**Spec resources** While it is possible to formulate a relational version of Hoare logic (with Hoare *quadruples* [33]), or to develop special-purpose logics for refinement [5], our goal with CaReSL is to support *both* standard Hoare-style reasoning and refinement reasoning in a single unified logic. In particular, we want a treatment of refinement that re-uses as much Hoare-style reasoning as possible. To this end, we adapt the idea of *specification resources*, first proposed by Turon *et al.* [31] as a way of proving refinement when threads engage in “cooperation” (a point we return to in §4). We will show how spec resources make it possible to state and prove refinements as an *entirely* derived concept on top of the Hoare logic we have already built up, in particular allowing protocols and triples to serve double-duty when reasoning about spec code.

To prove that  $e_1$  refines  $e_s$ , one must (intuitively) show that every behavior observable of  $e_1$  is observable of  $e_s$  as well. Our strategy is to encode these proof obligations into certain Hoare triples about the execution of  $e_1$ , but with pre- and postconditions instrumented with pieces of the corresponding spec state—both heap and code—which we treat as resources in CaReSL. These spec resources are entirely a fiction of the logic: they do not reflect anything about the *physical* state of  $e_1$ , but they must be used through logical rules that enforce the operational semantics for  $e_s$ .

There are two basic spec resource assertions: the (*spec*) *points-to* assertion  $\ell \hookrightarrow_s v$ , which claims ownership of a fragment of the spec heap containing a location  $\ell$  that points to the value  $v$ , and the *spec thread* assertion  $i \mapsto_s e$ , which claims ownership of a spec thread with ID  $i$  executing expression  $e$ . The separating conjunction  $P * Q$  works in the usual way with these resources, dividing up the spec heap and thread pool between the propositions  $P$  and  $Q$ . We also add a final sort, EvalCtx, and terms  $K$  and  $M[N]$  for expressing and combining them, which makes it possible to *abstract* over the evaluation context for some spec thread. (By convention, the variable  $\kappa$  is always of sort EvalCtx.)

In addition, CaReSL provides a necessary proposition  $P \rightarrow_s Q$ , which says that the spec resources owned by  $P$  can, according to the operational semantics, take a step to those owned by  $Q$ . All other resources must be left invariant. The rule SPURE lifts the pure stepping relation from the operational semantics directly. The rule SPRIM, on the other hand, *re-uses* atomic triples to support reasoning about atomic spec expressions. (The subscript IS in the laws of atomic triples allows them to be used in *either* implementation mode I or spec mode S.) The EXECSPEC Hoare rule (and identical AEXECSPEC rule for atomic triples) allows the specification to be

## Relating expressions

$$(e_1, e_s) \downarrow \varphi \triangleq \forall (i \bowtie j), \kappa.$$

$$\{ \text{Tid}(i) * j \mapsto_s \kappa[e_s] \} i \mapsto_{e_1} \{ x_1. \exists x_s. \varphi(x_1, x_s) * \text{Tid}(i) * j \mapsto_s \kappa[x_s] \}$$

## Relating values

$$\llbracket \mathbf{1} \rrbracket \triangleq (x_1, x_s). x_1 = x_s = ()$$

$$\llbracket \mu\alpha.\tau \rrbracket \triangleq \mu\alpha.\llbracket \tau \rrbracket$$

$$\llbracket \mathbf{B} \rrbracket \triangleq (x_1, x_s). (x_1 = x_s = \text{true}) \vee (x_1 = x_s = \text{false})$$

$$\llbracket \alpha \rrbracket \triangleq \alpha$$

$$\llbracket \tau_1 \times \tau_2 \rrbracket \triangleq (x_1, x_s). (\text{prj}_1 x_1, \text{prj}_1 x_s) \downarrow \llbracket \tau_1 \rrbracket \wedge (\text{prj}_2 x_1, \text{prj}_2 x_s) \downarrow \llbracket \tau_2 \rrbracket$$

$$\llbracket \tau_1 + \tau_2 \rrbracket \triangleq (x_1, x_s). (\triangleright \exists (y_1, y_s) \in \llbracket \tau_1 \rrbracket. x_1 = \text{inj}_1 y_1 \wedge x_s = \text{inj}_1 y_s) \vee (\triangleright \exists (y_1, y_s) \in \llbracket \tau_2 \rrbracket. x_1 = \text{inj}_2 y_1 \wedge x_s = \text{inj}_2 y_s)$$

$$\llbracket \tau \rightarrow \tau' \rrbracket \triangleq (x_1, x_s). \square \forall y_1, y_s. (\triangleright (y_1, y_s) \in \llbracket \tau \rrbracket) \Rightarrow (x_1 y_1, x_s y_s) \downarrow \llbracket \tau' \rrbracket$$

$$\llbracket \forall \alpha.\tau \rrbracket \triangleq (x_1, x_s). \square \forall \alpha. \text{Type}(\alpha) \Rightarrow (x_1 -, x_s -) \downarrow \llbracket \tau \rrbracket$$

$$\llbracket \text{ref } \tau \rrbracket \triangleq (x_1, x_s). \text{Inv}(\exists (y_1, y_s) \in \llbracket \tau \rrbracket. x_1 \hookrightarrow_1 y_1 * x_s \hookrightarrow_2 y_s)$$

$$\llbracket \text{refLcl } \tau \rrbracket \triangleq (x_1, x_s). \text{Inv} \left( \begin{array}{l} \exists L_1, L_s. x_1 \hookrightarrow_1 L_1 * x_s \hookrightarrow_2 L_s * \\ \forall (i \bowtie j). (\llbracket L_1 \rrbracket(i), \llbracket L_s \rrbracket(j)) \in \llbracket \mathbf{1} + \tau \rrbracket \end{array} \right)$$

## Invariant protocols

$$\text{Inv}(P) \triangleq \exists n. \boxed{0; \emptyset}^n_{\{\emptyset, \emptyset, [0 \rightarrow \emptyset], (s). P}}$$

## Type interpretations

$$\text{Type}(\varphi) \triangleq \square \forall (x_1, x_s) \in \varphi. \square (x_1, x_s) \in \varphi$$

## Logical relatedness (i.e., refinement)

$$\bar{\alpha}; \bar{x} : \bar{\tau} \vdash e_1 \preceq e_s : \tau \triangleq$$

$$\bar{\alpha}, \bar{x}, x_1, x_s; \text{Type}(\alpha), (x_1, x_s) \in \llbracket \tau \rrbracket \vdash (e_1[x_1/x], e_s[x_s/x]) \downarrow \llbracket \tau \rrbracket$$

Figure 8. Encoding refinement in CaReSL

“executed” within any postcondition, *i.e.*, at any point in a proof. Since EXECSPEC can be applied repeatedly, a single step of some implementation code—*e.g.*, an atomic expression—can correspond to an arbitrary number of spec steps.

Putting these pieces together, a triple like

$$\{ j \mapsto_s e_s \} i \mapsto_{e_1} \{ x_1. \exists x_s. j \mapsto_s x_s * \varphi(x_1, x_s) \}$$

says that if  $e_1$  produces some value  $x_1$ , then  $e_s$  can be executed to produce some value  $x_s$  such that  $\varphi(x_1, x_s)$  holds—exactly the kind of observational claim we set out to make.

## Encoding refinement

To give a full account of refinement, we also need to ensure that updates to public reference cells by the implementation are matched in lock-step by updates by the spec, which we will do by imposing a protocol governing public reference cells. And to account for thread-local references, we must also track a *correspondence* between implementation and specification thread IDs, which allows us to state invariants connecting the storage on either side. The (necessary) assertion  $i \bowtie j$  asserts that the implementation thread  $i$  and spec thread  $j$  belong to the bijective correspondence, *i.e.*,  $(i \bowtie j) \wedge (i' \bowtie j') \Rightarrow (i = i' \Leftrightarrow j = j')$ . Correspondences are introduced using the rule FORKS, a variant of FORK that jointly creates fresh implementation and spec threads.

And that’s all: using these ingredients, we can *encode* refinement by simply writing down a particular predicate in CaReSL.

The encoded predicate expresses a logical relation—a variant of the relation given semantically by Turon *et al.* [31]—following the approach first laid out by Plotkin and Abadi [23]. We define the relation in stages (see Figure 8).

First, we have the proposition  $(e_1, e_s) \downarrow \varphi$ , which is a more general version of the triple we suggested above: it includes assumptions about thread IDs, and permits compositional reasoning about specs by quantifying over an unknown evaluation context  $\kappa$ . The expressions do not begin with private ownership of any heap resources because, when proving refinement, we must assume that any pre-existing state is shared with the context. As we will see in a moment, this shared state is governed by an extremely liberal protocol; all we can assume about the context is that it is well-typed.

Second, we have the binary predicate  $\llbracket \tau \rrbracket$ , which is satisfied by  $(v_1, v_s)$  when the observations a context can make of  $v_1$  can

```

stacks : ∀α. (α → 1) × (1 → (1 + α)) × ((α → 1) → 1)
stacks ≜ Λ. let sync = mkSync(), hds := new (none)
  let push = λx. hds := some(x, get hds)
  let pop = λ(). case get hds of none ⇒ none
    | some(x, n) ⇒ hds := n; some(x)
  let snap = sync - (λ(). get hds)
  let iter = λf. foreach - f (snap())
  in (sync - push, sync - pop, iter)

stacki ≜ Λ. let hdi = new nil
  where nil ≜ new none,
  cons e ≜ new some(e)
  let push = rec try(x).
  let c = get hdi, n = cons(x, c)
  in if CAS(hdi, c, n) then () else try(x),
  let pop = rec try().
  let c = get hdi in case get c of
    none ⇒ none
  | some(d, n) ⇒ if CAS(hdi, c, n) then some(d) else try()
  let iter = λf. let rec loop(c) = case get c of
    none ⇒ () | some(d, n) ⇒ f(d); loop(n)
  in loop(get hdi)
  in (push, pop, iter)

```

**Figure 9.** Coarse- and fine-grained stacks, with iterators

also be made of  $v_s$ —the heart of the logical relation. It is defined by induction on the structure of  $\tau$ , since the ways a value can be observed depend on its type:

For **ground types**, the context can observe the exact value, so only equal values are in the interpretation.

For **product types**, the context can project both sides of the product, so two values are related iff each of their projections are related; similarly for **sums**. The context can observe **functions** only by applying them, so one function is related to another iff, when applied to related values, they produce related results.

**Recursive and polymorphic types** are interpreted via guarded recursion and second-order quantification, respectively; we pun type variables  $\alpha, \beta$  as predicate variables, which are implicitly assumed to be of sort  $\mathbb{P}(\text{Val} \times \text{Val})$ , and for quantification explicitly required to satisfy  $\text{Type}(\alpha)$ . These two constraints guarantee that the relation  $\llbracket \tau \rrbracket$  is a resource-insensitive, binary predicate on values. The resource-insensitivity reflects the fact that refinement between values must hold in an arbitrary context of observation/usage—including contexts that freely copy the values in question—which means that we can assume nothing about privately-owned resources.

The context can interact with values of **reference type** by either reading from or writing to them at any time. Thus, references assert the existence of a simple *invariant* protocol with a single state, whose interpretation says that related locations must continuously point to related values. For **thread-local references**, the statement is qualified by indexing the storage table by related thread IDs.

The uses of  $\triangleright$  throughout are necessary to ensure that recursive predicate variables  $\alpha$  only occur in guarded positions. Note that island assertions are implicitly guarded, as are uses of  $(e_1, e_s) \downarrow \varphi$  where  $e_1$  is not a value, because any  $\alpha$  mentioned therein will not be interpreted until after a step of computation.

Finally,  $e_1$  is *logically related* to  $e_s$  at some context  $\Omega$  and type  $\tau$  if  $\Omega \vdash e_1 \preceq e_s : \tau$ , which is shorthand for a use of the  $\mathcal{C} \vdash P$  judgment of CaReSL. The proof context closes all type variables  $\alpha$  with arbitrary type interpretations and all term variables  $x$  with *pairs* of term variables related at the appropriate types. Thus, a term with a free variable of type **ref**  $\tau$  will gain access in its proof context to a shared invariant protocol governing the corresponding location. The protocol in turn forces updates to the reference in the implementation to occur in lock-step with those in the spec. *Private* references, *i.e.*, those allocated *within* the implementation or spec, face no such requirement—unless or until they are exposed.

### Treiber’s stack, with an iterator

We now sketch a simple refinement proof. Figure 9 gives two stack implementations. The first,  $\text{stack}_s$ , is a coarse-grained reference implementation (which we think of as a specification). Its representation includes a mutable reference  $\text{hd}_s$  to an immutable list, as well as a synchronization wrapper  $\text{sync}$  provided by  $\text{mkSync}$ . The exported functions to push and pop from the stack simply wrap the corresponding updates to  $\text{hd}_s$  with synchronization. But the iterator is *not* fully synchronized: it takes an atomic snapshot of the stack, then calls its argument  $f$  on each element of the snapshot without holding any locks. So by the time  $f$  is called, the contents of the stack may have changed.

The second implementation,  $\text{stack}_i$ , is Treiber’s stack supplemented with an iterator. Treiber’s stack [28] is one of the simplest lock-free data structures—a kind of “Hello World” for concurrent program logics. Like  $\text{stack}_s$ , Treiber’s stack maintains a reference  $\text{hd}_i$  to a list that it uses to represent the stack. Instead of using a lock, however, it updates  $\text{hd}_i$  directly (and atomically) using **CAS**, which requires the contents of  $\text{hd}_i$  to have a *comparable* type, *i.e.*, to be testable for pointer equality. Thus  $\text{hd}_i$  is of type **ref**  $\text{clist}(\alpha)$ , where  $\text{clist}(\alpha) \triangleq \mu\beta. \text{ref } (1 + (\alpha \times \beta))$ . The push and pop functions employ *optimistic concurrency*: instead of acquiring a lock to inform other threads that they are going to make a change, they just take a snapshot of the current  $\text{hd}_i$ , compute a new value, and use **CAS** to install the new value if the identity of the  $\text{hd}_i$  has not changed. Intuitively, this strategy works because if the  $\text{hd}_i$ ’s identity has not changed, then *nothing* about the stack has changed: *all* mutation is performed by identity-changing updates to  $\text{hd}_i$ . The iterator simply walks over the stack’s contents starting from a (possibly stale!) snapshot of  $\text{hd}_i$ .

To prove that  $\text{stack}_i$  refines  $\text{stack}_s$ , we begin as follows:

```

stacki ≜ Λ. Prop context: Type(α), i ⊲ j   Variables: α, i, j, κ
{ j ⇒s κ[stacks -] } let hdi = new nil
{ j ⇒s κ[stacks -] * ∃x. hdi ↦1 x * x ↦1 none } (EXECSPEC)
{ ∃ hds, lock. hds ↦s none * lock ↦s false * j ⇒s κ[stacks'] * }
{ ∃x. hdi ↦1 x * x ↦1 none }

```

We have executed the “preambles” of both the implementation and spec—the code that allocates their respective hidden state. (Let  $\text{stack}'_s$  denote the rest of the spec code after the **let**-binding of  $\text{sync}$  and  $\text{hd}_s$ .) At this point, we will use NEWISL to move *all* of this state into a hidden island, with a protocol that we can use when proving refinement for each exported function.

The states  $s$  of the protocol are maps  $\text{Loc} \stackrel{\text{fin}}{\mapsto} \text{Val}$  which simply record the identity and contents of every node added to  $\text{stack}_i$ . The transition relation  $s \rightsquigarrow s' \triangleq s \subseteq s'$  captures the idea that, once a node has appeared in the stack, its contents never change. There are no tokens. We interpret a state  $s$  as follows:

$$\varphi \triangleq (s). \text{lock} \hookrightarrow_s \text{false} * \bigstar_{\ell \in \text{dom}(s)} \ell \hookrightarrow_1 s(\ell) * \exists x_1, x_s. \text{hd}_i \hookrightarrow_1 x_1 * \text{hd}_s \hookrightarrow_s x_s * \text{Link}(s, x_1, x_s)$$

Because the lock is protected by the protocol, we can only execute the exported spec functions with AEXECSPEC, as part of an application of UPDISL; by stipulating  $\text{lock} \hookrightarrow_s \text{false}$ , the interpretation furthermore requires that we run the spec in “big steps”, starting and ending in unlocked states. It also says that the implementation and spec stack contents are *linked*, in the following sense:

$$\text{Link} \triangleq \mu p. (s, x_1, x_s). \exists x'_1, s'. s = [x_1 \mapsto x'_1] \uplus s' * (x'_1 = \text{none} * x_s = \text{none}) \vee \left( \begin{array}{l} \exists y_1, z_1. x'_1 = \text{some}(y_1, z_1) * \\ \exists y_s, z_s. x_s = \text{some}(y_s, z_s) * \\ (y_1, y_s) \in \alpha * \triangleright p(s', z_1, z_s) \end{array} \right)$$

Linking requires that the stack state  $s$  contain an entire linked list starting from  $\text{hd}_i$ , and that each element of the list be related (at type  $\alpha$ ) to the corresponding element of the list in  $\text{hd}_s$ . *But it does so without mentioning the heap at all!* This is the key: it means that  $\text{Link}(s, x_1, x_s) \Rightarrow \square \text{Link}(s, x_1, x_s)$ , so the **Link** assertion can

be freely copied in the proof for `iter` when the traversal begins. Since the abstract state  $s$  of the stack can only grow, all of the nodes mentioned in this “snapshot” of `Link` are guaranteed to still be available in the region of the heap governed by the protocol, with their original values, as traversal proceeds—even if, in the meantime, the nodes are no longer reachable from `hd`.

The full proof outlines are available in the appendix [30].

**Combining refinement and Hoare-style reasoning** Suppose a client of Treiber’s stack wishes to use it only in a weak way, via a “per-item” spec similar to the one suggested in the introduction:

$$\begin{aligned} \text{PerItem} &\triangleq (e \in \text{Expr}). \forall p, q \in \mathbb{P}(\text{Val}). \\ &(\forall x. \text{TokPure}(p(x)) \wedge (p(x) \Rightarrow \Box q(x))) \Rightarrow \\ &\{\text{True}\} e \{\text{obj}. \text{obj} = (\text{add}, \text{rem}, \text{iter}) \wedge \Box P\}, \quad \text{where} \\ P &\triangleq \forall x. \{p(x)\} \text{add}(x) \{\text{ret}. \text{ret} = ()\} \\ &\wedge \{\text{True}\} \text{rem}() \{\text{ret}. \text{ret} = \text{none} \vee (\exists x. \text{ret} = \text{some}(x) * p(\text{ret}))\} \\ &\wedge \forall f, r. (\forall x. \{q(x) * r\} f(x) \{r\}) \Rightarrow \{r\} \text{iter}(f) \{r\} \end{aligned}$$

This “per-item” spec associates a resource  $p$  with each element of the stack, which is transferred to and from the data structure when adding or removing elements. If, in addition, each per-item resource entails some permanent knowledge  $\Box q$ , then that knowledge can be safely assumed by a function concurrently iterating over the stack, even if the resources originally supporting it have been consumed. As we will see in the case study (§4), this spec for iteration is useful for associating permanent, token-free knowledge about some other protocol with each item that appears in the stack.

While the per-item spec could in principle be applied directly to Treiber’s stack (*i.e.*, we could prove  $\text{PerItem}(\text{stack}_i \_)$ ), doing so would repeat much of the verification we just performed. On the other hand, proving  $\text{PerItem}(\text{stack}_s \_)$  is nearly trivial: we just instantiate the spinlock spec given in §3.2 with the representation invariant  $\text{Rep}_p \triangleq \exists l. \text{hd}_s \mapsto_1 l * \text{Map}_p(l)$ . To verify `iter`, we need only observe that  $\text{Map}_p(l) \Rightarrow \Box \text{Map}_q(l)$ , *i.e.*, the snapshot of the stack provides a list of necessarily true associated facts. (The details are in the appendix [30].) By mixing refinement and Hoare logic, we can thus significantly modularize our verification effort.

### 3.4 Soundness

The model theory of CaReSL is based directly on the semantic model of Turon *et al.* [31], with minor adjustments to accommodate the assertions  $\text{Tid}(i)$  and  $i \bowtie j$  necessary for reasoning about thread-local storage. (Since this paper is focused on the complementary aim of giving a syntactic proof theory, we do not delve into the model here; it can be found in full in the appendix [30].)

Soundness for CaReSL encompasses two theorems. **First**, that  $\mathcal{C} \vdash P$  implies  $\mathcal{C} \models P$ , where the latter is the semantic entailment relation of the model. Proving this theorem requires validating, in particular, the key Hoare logic rules, which we do in the appendix; these proofs resemble the proofs of key lemmas supporting Turon *et al.*’s model. **Second**, that  $\cdot \vdash (\Omega \vdash e_1 \leq e_s : \tau)$  implies  $\Omega \models e_1 \leq e_s : \tau$ , *i.e.*, that the logical relation is sound for contextual refinement. Again, this proof follows the soundness proof of Turon *et al.*, except that we can carry it out at a much higher level of abstraction using CaReSL’s proof rules.

## 4. Case study: Flat combining

Using CaReSL’s combination of Hoare-style and refinement reasoning, we can verify higher-order concurrent algorithms in layers of abstraction—and this section shows how to do it.

We study the recent *flat combining* construction of Hendler *et al.* [11], which takes an arbitrary sequential data structure and transforms it into a concurrent one in which all operations appear to take effect atomically. One can do so by merely wrapping each operation with synchronization on a global lock—indeed, this is

$$\begin{aligned} \text{flat}_s &\triangleq \text{mkSync}() : [\forall \alpha. \forall \beta. (\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \beta)] \\ \text{flat}_t &\triangleq \Lambda[\alpha]. \Lambda[\beta]. \lambda f : [\alpha \rightarrow \beta]. \quad (\text{Annotated with types for clarity}) \\ &\text{let lock} : [\text{ref } \mathbf{B}] = \text{new } (\text{false}) \\ &\text{let lclOps} : [\text{refLcl } op] = \text{newLcl} \quad (\text{where } op \triangleq \text{ref } (\alpha + \beta)) \\ &\text{let } (\text{add}, \_, \text{iter}) = \text{stack}_i [op] \\ &\text{let install} : [\alpha \rightarrow op] = \lambda \text{req}. \text{case getLcl}(lclOps) \\ &\quad \text{of some}(o) \Rightarrow o := \text{inj}_1 \text{ req}; o \\ &\quad \mid \text{none} \Rightarrow \text{let } o = \text{new } (\text{inj}_1 \text{ req}) \text{ in setLcl}(lclOps, o); \text{add}(o); o \\ &\text{let doOp} : [op \rightarrow \mathbf{1}] = \lambda o. \text{case get}(o) \\ &\quad \text{of inj}_1 \text{ req} \Rightarrow o := \text{inj}_2 f(\text{req}) \\ &\quad \mid \text{inj}_2 \text{ res} \Rightarrow () \\ &\text{in } \lambda x : [\alpha]. \text{let } o = \text{install}(x) \\ &\quad \text{let rec loop}() = \text{case get}(o) \\ &\quad \quad \text{of inj}_1 \_ \Rightarrow \text{if not}(\text{get lock}) \text{ and tryAcq}(\text{lock}) \text{ then} \\ &\quad \quad \quad \text{iter doOp}; \text{rel}(\text{lock}); \text{loop}() \\ &\quad \quad \quad \text{else loop}() \\ &\quad \mid \text{inj}_2 \text{ res} \Rightarrow \text{res} \\ &\text{in loop}() \end{aligned}$$

**Figure 10.** Flat combining: Spec and implementation

what our *spec* does—but flat combining takes a cache-friendly approach intended for hard-to-parallelize data structures.

The basic idea is to allow concurrent threads to advertise, through a lock-free side channel, their intent to perform an operation on the data structure. One of the threads then becomes the *combiner*: it locks the data structure and services the requests of all the other threads, one at a time (though new requests may be published concurrently with the combiner’s execution). The algorithm exhibits relatively good cache behavior for two reasons: (1) most of the time, operations do not need to execute *any* **CASES** to complete and (2) the combining thread enjoys good cache locality when executing the operations of the sequential data structure. In practice, flat combining yields remarkably strong performance, even when compared against completely lock-free data structures.

**The algorithm and its spec** The flat combining algorithm was originally given as informal prose, so our first task is to formalize its implementation and find an appropriate spec, while making its higher-order structure explicit. We do so in Figure 10.

We model an “arbitrary sequential data structure” as a closure  $f$  of some type  $\alpha \rightarrow \beta$ , where  $\alpha$  represents the input to an operation to be performed,  $\beta$  represents the result, and calling the function performs the operation by imperatively updating some state hidden within the closure. In other words,  $\alpha \rightarrow \beta$  represents the type of an *object*. The goal of flat combining, then, is to wrap this object with (cache-efficient) synchronization. Its spec,  $\text{flat}_s$ , simply uses  $\text{mkSync}$  to provide generic synchronization via global locking.

Our flat combining implementation  $\text{flat}_t$  uses three data structures to control synchronization.<sup>10</sup> First, it uses a global lock to ensure that only one thread at a time acts as the combiner. Second, it uses a Treiber’s stack to enable threads to *enroll* in the data structure. To enroll, a thread inserts an *operation record* of type  $\text{ref } (\alpha + \beta)$ , by which it can both advertise requests (of type  $\alpha$ ) and then await the results (of type  $\beta$ ). The combiner performs these requests by iterating over the stack: for each record containing an  $\text{inj}_1 \text{ req}$ , it applies  $f$  to  $\text{req}$ , obtaining result  $\text{res}$  and updating the record to  $\text{inj}_2 \text{ res}$ . Finally, the thread-local reference `lclOps` allows threads to re-use their operation record once they have enrolled. Note that operation records can be added to the stack or reset from  $\beta$  to  $\alpha$  at *any* time—even when the combiner holds the lock.

<sup>10</sup> This implementation simplifies Hendler *et al.*’s description in three respects: it uses a stack rather than a queue, operation records are never de-enrolled, and the combiner does not coalesce multiple operations into a single step. The first simplification makes little difference, because ultimately we give a modular proof using our per-item spec, which could be applied to a queue as well. The other two need only minor changes to the protocol.

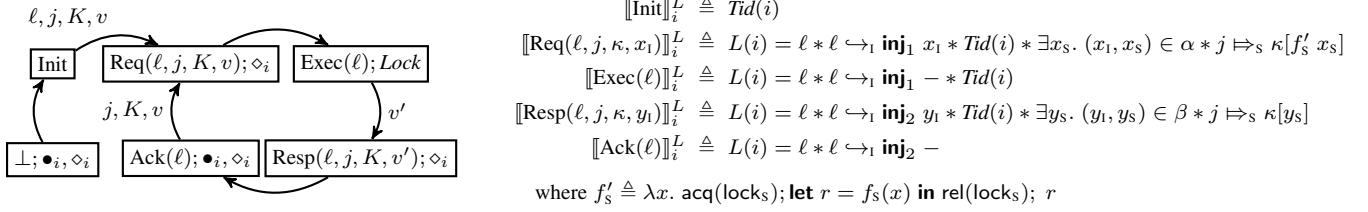


Figure 11. The protocol for the operation record of thread  $i$

Unfortunately, if we set out to prove the refinement

$$\cdot \vdash \text{flat}_1 \preceq \text{flat}_s : \forall \alpha. \forall \beta. (\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \beta)$$

we will run headlong into a problem: it does not hold! To wit: let

$$C \triangleq \text{let } r = \text{newLcl}, f = [] \_ \_ (\lambda(). \text{getLcl}(r)) \\ \text{in fork } (\text{setLcl}(r, \text{true}); f()); \text{ setLcl}(r, \text{false}); f()$$

When this context is applied to  $\text{flat}_s$ , it always returns some (**false**), because the final  $f()$  is always executed on the thread whose local  $r$  is **false**. But when applied to  $\text{flat}_1$ , it can also return some (**true**): the forked thread might act as the combiner, performing *both* applications of  $f$  and thus using its own thread-local value for  $r$ .

The crux of the problem is that thread-local storage allows functions to observe the identity of the thread executing them. We need to rule out this kind of side effect. We do so by unrolling the definition of refinement for functions, and simply removing access to ID-related knowledge, leading to a *qualified refinement*:

$$\forall \alpha, \beta, f_1, f_s. \\ \text{Type}(\alpha) \wedge \text{Type}(\beta) \wedge \square(\forall(x_1, x_s) \in \alpha. (f_1 x_1, f_s x_s) \downarrow^{\text{pure}} \beta) \\ \Rightarrow \cdot \vdash \text{flat}_1 \_ \_ f_1 \preceq \text{flat}_s \_ \_ f_s : \alpha \rightarrow \beta \\ \text{where } (e_1, e_s) \downarrow^{\text{pure}} \varphi \triangleq \forall i, j, \kappa. \\ \{j \mapsto_s \kappa[e_s]\} i \mapsto e_1 \{x_1. \exists x_s. \varphi(x_1, x_s) * j \mapsto_s \kappa[x_s]\}$$

This “pure” notion of related expressions embodies a semantic version of purity in a type-and-effect system where the only effect is “uses thread-local storage”. It coincides *exactly* with the expression relation in Turon *et al.* [31], where thread-local storage was not present. It is easy to prove, using CaReSL, that for any well-typed  $f : \tau \rightarrow \tau'$  that does not use thread-local storage (*i.e.*, an “effect-free”  $f$ ) we have  $\square \forall(x_1, x_s) \in \llbracket \tau \rrbracket. (f x_1, f x_s) \downarrow^{\text{pure}} \llbracket \tau' \rrbracket$ .

**The protocol** To prove the qualified refinement, we need to formulate a protocol characterizing the algorithm. We do so by giving a local protocol governing the operation record of each thread, and then tying these local protocols together into a single global one. Figure 11 gives the local protocol for a thread with ID  $i$ , where the states have the following informal meanings:

$\perp$	thread $i$ has not yet created its operation record
Init	thread $i$ is creating its operation record
Req	thread $i$ has requested an operation be performed
Exec	the combiner is executing thread $i$ 's operation
Resp	thread $i$ 's operation has been completed
Ack	thread $i$ has acknowledged the completion

The cycle in the protocol reflects that, once thread  $i$  has enrolled an operation record, it can reuse that record indefinitely. The labels on transitions signify *branching* on the values of the listed variables. Note that, while  $\ell$  (the location of the operation record) is chosen arbitrarily during initialization, it must remain fixed thereafter.

Movement through the local protocol encompasses two roles: that of thread  $i$ , and that of the combiner. The protocol guarantees: (1) when thread  $i$  begins execution of the flat combining algorithm, the current state is either  $\perp$  or Ack, meaning that either there is no operation record associated with the thread, or that the associated record is ready to be re-used; (2) only thread  $i$  can move to or from Init and Ack; and (3) only the combiner can move to or from Exec.

To fully understand how these constraints are enforced, we must take into account both the tokens and the state interpretations of the protocol. The tokens include  $\bullet_i$  and  $\diamond_i$  (one for each thread  $i$ , used only in  $i$ 's local protocol) and a single global token  $Lock$  representing ownership of the combiner's lock.

The interpretations  $\llbracket - \rrbracket_i^L$  of the non- $\perp$  states are shown in the figure; the parameter  $L$  represents the contents of  $\text{lclOps}$ . With the exception of  $\perp$  and Ack, all states assert ownership of  $Tid(i)$ . On the other hand, when thread  $i$  begins executing the algorithm, it will have ownership of  $Tid(i)$  (recall the definition of refinement, §3.3). Thus, thread  $i$  can assume that the protocol is in state  $\perp$  or Ack—and only thread  $i$  can take a step away from these states. For thread  $i$  to take such a step, it must “prove” that it *is* thread  $i$  by giving up  $Tid(i)$ , but in return it gains the token  $\bullet_i$  as a “receipt” that can later be traded back for  $Tid(i)$  by moving (back) to Ack. All told, the ownership of  $Tid(i)$  and the corresponding token  $\bullet_i$  account for the first two of the guarantees listed above.

The third guarantee is achieved using a similar strategy: to move to the Exec state, the combiner must “prove” that it holds the lock by temporarily giving up the (global)  $Lock$  token, but it receives the (local) token  $\diamond_i$  as a receipt. Subsequently, only the combiner can move to Resp, making the opposite trade.

After initialization, each state also asserts that the location  $\ell$  of the operation record for thread  $i$  both corresponds to  $L(i)$  and points to an appropriate sum injection (depending on the state).

A final aspect of the local protocol is capturing the *cooperation* inherent in flat combining: the combiner executes code on behalf of other threads. Cooperation is difficult for compositional refinement techniques to handle, because such techniques usually require proving that for *each piece* of implementation code the corresponding piece of spec code behaves the same way—and thus they typically provide no way to account for the *mismatch* between implementation and spec that cooperation entails. While our definition of refinement (§3.3) likewise imposes a one-to-one relationship between implementation and spec code in its pre- and post-conditions, our treatment of spec code as a resource allows ownership to be transferred to other threads *during* the execution of the implementation.

Thus, the Req and Resp states of our protocol assert *shared* ownership of  $i$ 's spec code (running in spec thread  $j$ ). In moving from Req to Exec the combiner thread must take ownership of the spec code  $\kappa[f'_s x_s]$ ; and to subsequently move to Resp, the combiner must actually evaluate this code on thread  $i$ 's behalf until  $f'_s x_s$  reaches a value  $y_s$ . Subsequently, when moving to the Ack state, thread  $i$  regains ownership of both  $Tid(i)$  and its spec code, by giving up its “receipt” token,  $\bullet_i$ .

The global transition system is then a *product construction*:

$$\text{State space } \left\{ (S, s) \mid \begin{array}{l} S \in \mathbb{N} \overset{\text{fin}}{\text{OpState}}, \quad s \in \text{LockState}, \\ \text{at most one } S(i) \text{ or } s \text{ owns } Lock \end{array} \right\} \\ \text{Transitions } (S, s) \rightsquigarrow (S', s') \triangleq s \rightsquigarrow_{\text{lock}}^? s' \wedge \forall i. S(i) \rightsquigarrow_i^? S'(i) \\ \text{Owned tokens } \mathcal{T}(S, s) \triangleq \mathcal{T}_{\text{lock}}(s) \cup \bigcup_i \mathcal{T}_i(S(i))$$

A global state includes a collection  $S$  of local states from the operation record protocol (Figure 11), one for each thread ID. In giving the state space, we pun the  $\perp$  state with an “undefined”

input to a partial function—and thus, we require that only a finite number of threads  $i$  have a non- $\perp$  state in  $S$ . Global states also have a single state  $s$  drawn from the standard locking protocol (§3.2), reflecting the state of the global lock. The global *Lock* token is *shared* amongst all of the combined protocols: in a valid global state, at most one of the local states  $S(i)$  or  $s$  may claim the *Lock* token. A transition in the global protocol allows each local protocol to make at most one transition. Finally, the tokens owned in a global state are just the union of all the tokens claimed by the local states.

Recall that in the refinement proof for Treiber’s stack (§3.3), the spec lock is treated continuously as a protocol-owned resource with value **false**, capturing the atomicity of the implementation code: if an implementation step coincides with *any* spec steps, it must coincide with an entire critical section’s worth, going from unlocked to unlocked spec states in big steps. The combiner implementation, by contrast, is executing a function  $f_1$  that is not necessarily atomic; all we know is that  $f_1$  refines  $f_s$ . We must therefore allow the spec lock to be held for multiple implementation steps, which we do by interpreting lock states as follows:

$$\begin{aligned} \llbracket \text{Unlocked} \rrbracket &\triangleq \text{lock} \hookrightarrow_1 \mathbf{false} * \text{lock}_s \hookrightarrow_s \mathbf{false} \\ \llbracket \text{Locked} \rrbracket &\triangleq \text{lock} \hookrightarrow_1 \mathbf{true} \end{aligned}$$

In short, the combiner gains private ownership of both the *Lock* token and the spec lock itself (*i.e.*, the internal lock used by `mkSync`; see  $f'_s$  in Figure 11) when it acquires the implementation lock.

Finally, we lift these local interpretations to interpret global states via the following predicate:

$$(S, s). \llbracket s \rrbracket * \exists L. \text{!} \text{!} \text{Ops} \hookrightarrow_1 L * \bigstar_{i \in \text{dom}(S)} \llbracket S(i) \rrbracket_i^L$$

**The proof** We close with a high-level view of the proof itself; details, as usual, are in the appendix. Unrolling the statement of qualified refinement, we need to prove

$$\{ \text{Tid}(i) * j \mapsto_s \kappa[\text{flat}_s \dots f_s] \} i \mapsto \text{flat}_1 \dots f_1 \\ \{ x_1. \exists x_s. (x_1, x_s) \in \llbracket \alpha \rightarrow \beta \rrbracket * \text{Tid}(i) * j \mapsto_s \kappa[x_s] \}$$

under the assumptions

$$i \bowtie j, \text{Type}(\alpha), \text{Type}(\beta), \square(\forall(x_1, x_s) \in \alpha. (f_1 x_1, f_s x_s) \downarrow^{\text{pure}} \beta)$$

The proof begins in much the same way as the refinement proof for Treiber’s stack (§3.3): we execute the **let**-bound expressions that allocate the hidden state in both the implementation and spec, *i.e.*,

$$\text{lock} \hookrightarrow_1 \mathbf{false} * \text{!} \text{!} \text{Ops} \hookrightarrow_1 \emptyset * \exists \text{lock}_s. j \mapsto_s \kappa[f'_s] * \text{lock}_s \hookrightarrow_s \mathbf{false}$$

where  $\text{lock}_s$  is the lock allocated by `mkSync` and  $f'_s$  is as in Figure 11. These resources are precisely what we need to apply NEWISL for our global protocol, moving them from private to shared ownership. Letting  $\pi$  be the global protocol defined above, we can claim  $\square(\emptyset, \text{Unlocked}; \emptyset) \stackrel{\pi}{\pi}$ , *i.e.*, permanent knowledge that the global protocol exists. We must then, in the context of this island and our previous assumptions, show that the closure returned by  $\text{flat}_1$  refines the one returned by  $\text{flat}_s$ , *i.e.*,  $f'_s$ , at type  $\alpha \rightarrow \beta$ .

We first verify a version  $\text{flat}'_1$  of the flat combining algorithm that is identical to the one in Figure 10, except that it uses the coarse-grained stack, allowing us to use the per-item spec of §3.3. We instantiate the per-item spec using the same predicate  $\text{Op}$  for per-item resources and iteration knowledge ( $p$  and  $q$ , respectively, in the per-item spec):  $\text{Op} \triangleq (\ell). \exists k. \boxed{[k \mapsto \text{Req}(\ell, -, -, -)]}_{\pi}^{11}$ . This is a *local* assertion about the global protocol, in that the states of threads other than  $k$  can be in any rely-future state of  $\perp$ , *i.e.*, any state whatsoever. The  $\text{Op}$  predicate just claims that location  $\ell$  is an initialized operation record for *some* thread. By the per-item spec, all locations inserted into the stack must satisfy this property—and since we have  $\text{Op}(\ell) \Rightarrow \square \text{Op}(\ell)$ , the property can be assumed even when iterating over stale elements of the stack. Operation records are created via `install`, whose specification consumes  $\text{Tid}(i)$  and associated spec resources in exchange for the “receipt”  $\bullet_i$ :

$$\{ \text{Tid}(i) * j \mapsto_s \kappa[f'_s x_s] \} i \mapsto \text{install } x \left\{ o. \boxed{[i \mapsto \text{Req}(o, j, \kappa, x)]}_{\pi} \bullet_i \right\}$$

<sup>11</sup> We leave off the locking state when it is `Unlocked`.

for any  $(x, x_s) \in \alpha$  (*i.e.*, for any related arguments). The combiner actually performs operations via `doOp`,

$$\{ \text{Op}(o) * P \} \text{doOp } o \{ P \} \text{ where } P \triangleq \boxed{\emptyset; \text{Lock}}_{\pi}^n * \text{lock}_s \hookrightarrow_s \mathbf{false}$$

which assumes that the given location  $o$  is a valid operation record, and that the invoking thread owns the *Lock* token as well as the spec lock. The shape of the spec for `doOp` exactly matches that required for iteration in the per-item spec (§3.3), with  $P$  serving as the loop invariant, thus allowing us to deduce  $\{ P \} \text{iter doOp } \{ P \}$ . These auxiliary specs make it straightforward to prove refinement for the closures returned by  $\text{flat}'_1$  (with coarse-grained stack) and  $\text{flat}_s$ .

Finally, suppose we plug the combiner into a client context  $C$  that provides a suitable argument  $f$ . The above proof (together with soundness, §3.4) allows us to deduce  $\models C[\text{flat}'_1] \preceq C[\text{flat}_s] : \tau$ . Likewise, refinement for Treiber’s stack allows us to deduce  $\models C[\text{flat}_1] \preceq C[\text{flat}'_1] : \tau$ . Since refinement is *transitive*, we can compose these together to conclude  $\models C[\text{flat}_1] \preceq C[\text{flat}_s] : \tau$ .<sup>12</sup>

## 5. Related work

In many respects, CaReSL builds directly on the semantic foundation that we and colleagues laid in our prior work [31]. There, we developed a relational Kripke model of a language very similar to the one considered here, and showed how granularity abstraction for several sophisticated (but structurally simple) fine-grained data structures could be established by direct appeal to the model. The present work is a natural continuation of that work. First, we provide a *logic* that greatly simplifies reasoning compared to working with the model; such a proof theory is essential for scaling the verification method to large examples. Second, we use our logic not just to prove granularity abstraction results in isolation (as we did before), but also to facilitate the modular verification of a higher-order, structurally complex program. Along the way, we also extend our prior model to handle thread-local storage.

The logic itself draws inspiration from LADR [5], which in turn provided a proof theory for reasoning about a *sequential* relational Kripke model (ADR [1]). Aside from incorporating concurrency, CaReSL offers a deeper unification of refinement and Hoare logic by treating refinement as a mode of use of Hoare logic.

**Higher-order functions and concurrency** While there has been stunning progress in logics for concurrency over the last decade, very few of these logics meaningfully support *higher-order* programming, and among those that do, none supports reasoning about fine-grained synchronization or granularity abstraction.

The logic that comes closest is *higher-order concurrent abstract predicates* (HOCAP), first proposed for reasoning about first-order code [4] (the “higher-order” refers to the logic of predicates) and very recently applied to higher-order code as well [27]. HOCAP, like its predecessor CAP [3], accounts for concurrency by way of “shared region” assertions that constrain the possible updates to a shared resource. Our island assertions resemble shared region assertions—indeed, we have adopted notation suggesting as much—but work at a higher level of abstraction (*i.e.*, protocol states), separating *knowledge* bounding the state of the protocol (treated as a copyable assertion) from the *rights* to change the state (treated as a linear resource: tokens); see [31] for a more detailed comparison. Because CAP lacks granularity abstraction, it is difficult to give a single “principal” specification for a data structure. Instead, one gives specialized specs (like the per-item spec for stacks) reflecting particular usages envisioned for a client—which means new client scenarios necessitate new proofs. HOCAP attempts to

<sup>12</sup> We are appealing to *semantic* refinement here to take advantage of transitivity. This reasoning can be internalized in CaReSL by adding a proposition for semantic refinement and axioms for refinement soundness and transitivity, but there is little to be gained from doing so.

address this shortcoming by explicitly quantifying over the way a client uses a data structure, but (1) this introduces the potential for a problematic circularity, which clients must explicitly prove does not arise, and (2) it is not clear how to scale the approach to handle *cooperation* between threads (as in the flat combiner).

Other concurrency logics that support higher-order functions—such as Hobor *et al.*'s extension of concurrent separation logic [14], or the very recent *Subjective Concurrent Separation Logic* [17]—support only reasoning about simple lock-based synchronization, and do not enable the kinds of refinement proofs we have presented.

**Granularity abstraction** Herlihy and Wing's seminal notion of *linearizability* [13] has long been held as the gold standard of correctness for concurrent data structures, but as Filipović *et al.* argue [8], what clients of these data structures *really* want is a contextual refinement property. Filipović *et al.* go on to show that, under certain (strong) assumptions about a programming language, linearizability implies contextual refinement for that language. More recently, Gotsman and Yang generalized both linearizability and Filipović *et al.*'s result (the so-called *abstraction theorem*) to include potential ownership transfer of memory between concurrent data structures and their clients in a first-order setting [9]. While in principle proofs of linearizability can be composed with these results to support granularity abstraction, no existing logic has provided support for doing so. We found it simpler to work with refinement directly (in particular, to *encode* it directly into a Hoare logic), rather than reasoning about linearizability. CaReSL enables clients to layer ownership-transferring protocols *on top* of a coarse-grained reference implementation, *after* applying granularity abstraction, as we showed with the per-item spec (§3.3).

The only Hoare logic we are aware of that can (formally) prove refinement results is Liang and Feng's new logic [18] (extending LRG [7]), which is inspired by the use of ghost state in Vafeiadis's thesis [32]. The logic is powerful enough to deal with a range of sophisticated, fine-grained concurrent algorithms, but it is limited to *first-order* code. In addition, the specifications used in refinement are not reference implementations, but are instead essentially atomic Hoare triples. While such specifications are appealingly abstract, they present two limitations: (1) they do not model the more general notion of granularity abstraction (supporting only *atomicity* abstraction, as we explained in footnote 1) and (2) they do not support the kind of transitive composition of proofs that we used in our case study. As a result, it is unclear how to use the logic to build modular proofs layering Hoare logic and refinement.

A radically different approach to atomicity abstraction is Lipton's method of *reduction* [19], which is based on showing that an atomic step *commutes* with all concurrent activity, and can therefore be coalesced into a larger atomic block with *e.g.*, code that is sequenced after it. Elmas *et al.* developed a logic for proving linearizability via a combination of reduction and abstraction [6], which in some ways resembles our interleaved application of refinement and Hoare-style reasoning, but with a rather different way of *proving* refinement. It is limited, however, to first-order code and atomicity refinement, and like HOCAP it is not powerful enough to handle fine-grained algorithms that employ cooperation. Moreover, it does not allow linearizability proofs to be composed transitively.

**Acknowledgments** We would like to thank David Swasey for his careful reading of both the paper and its appendix. This work was partially funded by the EC FET project ADVENT.

## References

- [1] A. Ahmed, D. Dreyer, and A. Rossberg. State-dependent representation independence. In *POPL*, 2009.
- [2] R. D. Blumofe, C. F. Joerg, B. C. Kuszmaul, C. E. Leiserson, K. H. Randall, and Y. Zhou. Cilk: An efficient multithreaded runtime system. *JPDC*, 37(1):55–69, Aug. 1996.
- [3] T. Dinsdale-Young, M. Dodds, P. Gardner, M. Parkinson, and V. Vafeiadis. Concurrent abstract predicates. In *ECOOP*, 2010.
- [4] M. Dodds, S. Jagannathan, and M. Parkinson. Modular reasoning for deterministic parallelism. In *POPL*, 2011.
- [5] D. Dreyer, G. Neis, A. Rossberg, and L. Birkedal. A relational modal logic for higher-order stateful ADTs. In *POPL*, 2010.
- [6] T. Elmas, S. Qadeer, A. Sezgin, O. Subasi, and S. Tasiran. Simplifying linearizability proofs with reduction and abstraction. In *TACAS*, 2010.
- [7] X. Feng. Local rely-guarantee reasoning. In *POPL*, 2009.
- [8] I. Filipović, P. O'Hearn, N. Rinetzkzy, and H. Yang. Abstraction for concurrent objects. *Theoretical Computer Science*, 411, 2010.
- [9] A. Gotsman and H. Yang. Linearizability with ownership transfer. In *CONCUR*, 2012.
- [10] T. Harris, S. Marlow, S. Peyton-Jones, and M. Herlihy. Composable memory transactions. In *PPOPP*, 2005.
- [11] D. Hendler, I. Ince, N. Shavit, and M. Tzafrir. Flat combining and the synchronization-parallelism tradeoff. In *SPAA*, 2010.
- [12] D. Hendler, N. Shavit, and L. Yerushalmi. A scalable lock-free stack algorithm. In *SPAA*, 2004.
- [13] M. P. Herlihy and J. M. Wing. Linearizability: a correctness condition for concurrent objects. *TOPLAS*, 12(3):463–492, 1990.
- [14] A. Hobor, A. W. Appel, and F. Z. Nardelli. Oracle semantics for concurrent separation logic. In *ESOP*, 2008.
- [15] C. B. Jones. Tentative steps toward a development method for interfering programs. *TOPLAS*, 5(4):596–619, 1983.
- [16] D. Lea. The java.util.concurrent ConcurrentHashMap.
- [17] R. Ley-Wild and A. Nanevski. Subjective auxiliary state for coarse-grained concurrency. In *POPL*, 2013.
- [18] H. Liang and X. Feng. Modular verification of linearizability with non-fixed linearization points. In *PLDI*, 2013.
- [19] R. J. Lipton. Reduction: a method of proving properties of parallel programs. *Commun. ACM*, 18(12):717–721, 1975.
- [20] P. W. O'Hearn. Resources, concurrency, and local reasoning. *Theor. Comput. Sci.*, 375(1-3):271–307, 2007.
- [21] M. Parkinson and G. Bierman. Separation logic and abstraction. In *POPL*, 2005.
- [22] A. M. Pitts and I. Stark. Operational reasoning for functions with local state. In *HOOTS*, 1998.
- [23] G. Plotkin and M. Abadi. A logic for parametric polymorphism. In *TLCA*, 1993.
- [24] F. Pottier. Hiding local state in direct style: a higher-order anti-frame rule. In *LICS*, 2008.
- [25] J. H. Reppy. *Higher-order concurrency*. PhD thesis, Cornell University, 1992.
- [26] J. C. Reynolds. Separation logic: A logic for shared mutable data structures. In *LICS*, 2002.
- [27] K. Svendsen, L. Birkedal, and M. Parkinson. Modular reasoning about separation of concurrent data structures. In *ESOP*, 2013.
- [28] R. Treiber. Systems programming: coping with parallelism. Technical report, Almaden Research Center, 1986.
- [29] P. W. Trinder, K. Hammond, H.-W. Loidl, and S. L. Peyton Jones. Algorithm + strategy = parallelism. *JFP*, 8(1):23–60, Jan. 1998.
- [30] A. Turon, D. Dreyer, and L. Birkedal. Unifying refinement and Hoare-style reasoning in a logic for higher-order concurrency: Appendix. <http://www.mpi-sws.org/~turon/caresl/appendix.pdf>.
- [31] A. Turon, J. Thamsborg, A. Ahmed, L. Birkedal, and D. Dreyer. Logical relations for fine-grained concurrency. In *POPL*, 2013.
- [32] V. Vafeiadis. *Modular fine-grained concurrency verification*. PhD thesis, University of Cambridge, 2008.
- [33] H. Yang. Relational separation logic. *TCS*, 375(1-3):308–334, 2007.