

A Rely-Guarantee-Based Simulation for Verifying Concurrent Program Transformations

Hongjin Liang Xinyu Feng Ming Fu

School of Computer Science and Technology, University of Science and Technology of China
Hefei, Anhui 230026, China

lhj1018@mail.ustc.edu.cn xyfeng@ustc.edu.cn fuming@ustc.edu.cn

Abstract

Verifying program transformations usually requires proving that the resulting program (the target) refines or is equivalent to the original one (the source). However, the refinement relation between individual sequential threads cannot be preserved in general with the presence of parallel compositions, due to instruction reordering and the different granularities of atomic operations at the source and the target. On the other hand, the refinement relation defined based on fully abstract semantics of concurrent programs assumes arbitrary parallel environments, which is too strong and cannot be satisfied by many well-known transformations.

In this paper, we propose a Rely-Guarantee-based Simulation (RGSim) to verify concurrent program transformations. The relation is parametrized with constraints of the environments that the source and the target programs may compose with. It considers the interference between threads and their environments, thus is less permissive than relations over sequential programs. It is compositional *w.r.t.* parallel compositions as long as the constraints are satisfied. Also, RGSim does not require semantics preservation under all environments, and can incorporate the assumptions about environments made by specific program transformations in the form of rely/guarantee conditions. We use RGSim to reason about optimizations and prove atomicity of concurrent objects. We also propose a general garbage collector verification framework based on RGSim, and verify the Boehm *et al.* concurrent mark-sweep GC.

Categories and Subject Descriptors D.2.4 [Software Engineering]: Software/Program Verification – Correctness proofs, Formal methods; F.3.1 [Logics and Meanings of Programs]: Specifying and Verifying and Reasoning about Programs

General Terms Theory, Verification

Keywords Concurrency, Program Transformation, Rely-Guarantee Reasoning, Simulation

1. Introduction

Many verification problems can be reduced to verifying program transformations, *i.e.*, proving the target program of the transformation has no more observable behaviors than the source. Below we give some typical examples in concurrent settings:

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

POPL'12, January 25–27, 2012, Philadelphia, PA, USA.
Copyright © 2012 ACM 978-1-4503-1083-3/12/01...\$10.00

- *Correctness of compilation and optimizations of concurrent programs.* In this most natural program transformation verification problem, every compilation phase does a program transformation \mathbf{T} , which needs to preserve the semantics of the inputs.
- *Atomicity of concurrent objects.* A concurrent object or library provides a set of methods that allow clients to manipulate the shared data structure with abstract atomic behaviors [15]. Their correctness can be reduced to the correctness of the transformation from abstract atomic operations to concrete and executable programs in a concurrent context.
- *Verifying implementations of software transactional memory (STM).* Many languages supporting STM provide a high-level atomic block $\mathbf{atomic}\{\mathbb{C}\}$, so that programmers can assume the atomicity of the execution of \mathbb{C} . Atomic blocks are implemented using some STM protocol (*e.g.*, TL2 [11]) that allows very fine-grained interleavings. Verifying that the fine-grained program respects the semantics of atomic blocks gives us the correctness of the STM implementation.
- *Correctness of concurrent garbage collectors (GCs).* High-level garbage-collected languages (*e.g.*, Java) allow programmers to work at an abstract level without knowledge of the underlying GC algorithm. However, the concrete and executable low-level program involves interactions between the mutators and the collector. If we view the GC implementation as a transformation from high-level mutators to low-level ones with a concrete GC thread, the GC safety can be reduced naturally to the semantics preservation of the transformation.

To verify the correctness of a program transformation \mathbf{T} , we follow Leroy's approach [19] and define a refinement relation \sqsubseteq between the target and the source programs, which says the target has no more observable behaviors than the source. Then we can formalize the correctness of the transformation as follows:

$$\text{Correct}(\mathbf{T}) \triangleq \forall C, \mathbb{C}. C = \mathbf{T}(\mathbb{C}) \implies C \sqsubseteq \mathbb{C}. \quad (1.1)$$

That is, for any source program \mathbb{C} acceptable by \mathbf{T} , $\mathbf{T}(\mathbb{C})$ is a refinement of \mathbb{C} . When the source and the target are shared-state concurrent programs, the refinement \sqsubseteq needs to satisfy the following requirements to support effective proof of $\text{Correct}(\mathbf{T})$:

- Since the target $\mathbf{T}(\mathbb{C})$ may be in a different language from the source, the refinement should be general and independent of the language details.
- To verify fine-grained implementations of abstract operations, the refinement should support different views of program states and different granularities of state accesses at the source and the target levels.
- When \mathbf{T} is syntax-directed (and it is usually the case for parallel compositions, *i.e.*, $\mathbf{T}(\mathbb{C} \parallel \mathbb{C}') = \mathbf{T}(\mathbb{C}) \parallel \mathbf{T}(\mathbb{C}')$), a *com-*

positional refinement is of particular importance for modular verification of \mathbf{T} .

However, existing refinement (or equivalence) relations cannot satisfy all these requirements at the same time. Contextual equivalence, the canonical notion for comparing program behaviors, fails to handle different languages since the contexts of the source and the target will be different. Simulations and logical relations have been used to verify compilation [4, 16, 19, 21], but they are usually designed for sequential programs (except [21, 25], which we will discuss in Section 8). Since the refinement or equivalence relation between sequential threads cannot be preserved in general with parallel compositions, we cannot simply adapt existing work on sequential programs to verify transformations of concurrent programs. Refinement relations based on fully abstract semantics of concurrent programs are compositional, but they assume arbitrary program contexts, which is too strong for many practical transformations. We will explain the challenges in detail in Section 2.

In this paper, we propose a Rely-Guarantee-based Simulation (RGSim) for compositional verification of concurrent transformations. By addressing the above problems, we make the following contributions:

- RGSim parametrizes the simulation between concurrent programs with rely/guarantee conditions [17], which specify the interactions between the programs and their environments. This makes the corresponding refinement relation compositional *w.r.t.* parallel compositions, allowing us to decompose refinement proofs for multi-threaded programs into proofs for individual threads. On the other hand, the rely/guarantee conditions can incorporate the assumptions about environments made by specific program transformations, so RGSim can be applied to verify many practical transformations.
- Based on the simulation technique, RGSim focuses on comparing externally observable behaviors (*e.g.*, I/O events) only, which gives us considerable leeway in the implementations of related programs. The relation is mostly independent of the language details. It can be used to relate programs in different languages with different views of program states and different granularities of atomic state accesses.
- RGSim makes relational reasoning about optimizations possible in parallel contexts. We present a set of relational reasoning rules to characterize and justify common optimizations in a concurrent setting, including hoisting loop invariants, strength reduction and induction variable elimination, dead code elimination, redundancy introduction, *etc.*
- RGSim gives us a refinement-based proof method to verify fine-grained implementations of abstract algorithms and concurrent objects. We successfully apply RGSim to verify concurrent counters, the concurrent GCD algorithm, Treiber’s non-blocking stack and the lock-coupling list.
- We reduce the problem of verifying concurrent garbage collectors to verifying transformations, and present a general GC verification framework, which combines unary Rely-Guarantee-based verification [17] with relational proofs based on RGSim.
- We verify the Boehm *et al.* concurrent garbage collection algorithm [7] using our framework. As far as we know, it is the first time to formally prove the correctness of this algorithm.

In the rest of this paper, we first analyze the challenges for compositional verification of concurrent program transformations, and explain our approach informally in Section 2. Then we give the basic technical settings in Section 3 and present the formal definition of RGSim in Section 4. We show the use of RGSim to reason about

```

local r1;          local r2;
x := 1;            y := 1;
r1 := y;          || r2 := x;
if (r1 = 0) then  || if (r2 = 0) then
critical region   || critical region

```

(a) Dekker’s Mutual Exclusion Algorithm

```

x := x+1; || x := x+1;
vs.

```

```

local r1;          local r2;
r1 := x;           || r2 := x;
x := r1 + 1;       x := r2 + 1;

```

(b) Different Granularities of Atomic Operations

Figure 1. Equivalence Lost after Parallel Composition

optimizations in Section 5, verify atomicity of concurrent objects in Section 6, and prove the correctness of concurrent GCs in Section 7. Finally we discuss related work and conclude in Section 8.

2. Challenges and Our Approach

The major challenge we face is to have a compositional refinement relation \sqsubseteq between concurrent programs, *i.e.*, we should be able to know $\mathbf{T}(\mathbb{C}_1) \parallel \mathbf{T}(\mathbb{C}_2) \sqsubseteq \mathbb{C}_1 \parallel \mathbb{C}_2$ if we have $\mathbf{T}(\mathbb{C}_1) \sqsubseteq \mathbb{C}_1$ and $\mathbf{T}(\mathbb{C}_2) \sqsubseteq \mathbb{C}_2$.

2.1 Sequential Refinement Loses Parallel Compositionality

Observable behaviors of sequential imperative programs usually refer to their control effects (*e.g.*, termination and exceptions) and final program states. However, refinement relations defined correspondingly cannot be preserved after parallel compositions. It has been a well-known fact in the compiler community that sound optimizations for sequential programs may change the behaviors of multi-threaded programs [5]. The Dekker’s algorithm shown in Figure 1(a) has been widely used to demonstrate the problem. Re-ordering the first two statements of the thread on the left preserves its sequential behaviors, but the whole program can no longer ensure exclusive access to the critical region.

In addition to instruction reordering, the different granularities of atomic operations between the source and the target programs can also break the compositionality of program equivalence in a concurrent setting. In Figure 1(b), the target program at the bottom behaves differently from the source at the top (assuming each statement is executed atomically), although the individual threads at the target and the source have the same behaviors.

2.2 Assuming Arbitrary Environments is Too Strong

The problem with the refinement for sequential programs is that it does not consider the effects of threads’ intermediate state accesses on their parallel environments. People have given fully abstract semantics to concurrent programs (*e.g.*, [1, 8]). The semantics of a program is modeled as a set of execution traces. Each trace is an interleaving of state transitions made by the program itself and *arbitrary* transitions made by the environment. Then the refinement between programs can be defined as the subset relation between the corresponding trace sets. Since it considers all possible environments, the refinement relation has very nice compositionality, but unfortunately is too strong to formulate the correctness of many well-known transformations, including the four classes of transformations mentioned before:

- Many concurrent languages (*e.g.*, C++ [6]) do not give semantics to programs with data races (like the examples shown in Figure 1). Therefore the compilers only need to guarantee the semantics preservation of data-race-free programs.
- When we prove that a fine-grained implementation of a concurrent object is a refinement of an abstract atomic operation, we can assume that all accesses to the object in the context of the target program use the same set of primitives.
- Usually the implementation of STM (*e.g.*, TL2 [11]) ensures the atomicity of a transaction `atomic{C}` only when there are no data races. Therefore, the correctness of the transformation from high-level atomic blocks to fine-grained concurrent code assumes data-race-freedom in the source.
- Many garbage-collected languages are type-safe and prohibit operations such as pointer arithmetics. Therefore the garbage collector could make corresponding assumptions about the mutators that run in parallel.

In all these cases, the transformations of individual threads are allowed to make various assumptions about the environments. They do not have to ensure semantics preservation within all contexts.

2.3 Languages at Source and Target May Be Different

The use of different languages at the source and the target levels makes the formulation of the transformation correctness more difficult. If the source and the target languages have different views of program states and different atomic primitives, we cannot directly compare the state transitions made by the source and the target programs. This is another reason that makes the aforementioned subset relation between sets of program traces in fully abstract semantics infeasible. For the same reason, many existing techniques for proving refinement or equivalence of programs in the same language cannot be applied either.

2.4 Different Observers Make Different Observations

Concurrency introduces tensions between two kinds of observers: human beings (as external observers) and the parallel program contexts. External observers do not care about the implementation details of the source and the target programs. For them, intermediate state accesses (such as memory reads and writes) are silent steps (unobservable), and only external events (such as I/O operations) are observable. On the other hand, state accesses have effects on the parallel program contexts, and are not silent to them.

If the refinement relation relates externally observable event traces only, it cannot have parallel compositionality, as we explained in Section 2.1. On the other hand, relating all state accesses of programs is too strong. Any reordering of state accesses or change of atomicity would fail the refinement.

2.5 Our Approach

In this paper we propose a Rely-Guarantee-based Simulation (RGSim) \preceq between the target and the source programs. It establishes a weak simulation, ensuring that for every externally observable event made by the target program there is a corresponding one in the source. We choose to view intermediate state accesses as silent steps, thus we can relate programs with different implementation details. This also makes our simulation independent of language details.

To support parallel compositionality, our relation takes into account explicitly the expected interference between threads and their parallel environments. Inspired by the Rely-Guarantee (R-G) verification method [17], we specify the interference using rely/guarantee conditions. In Rely-Guarantee reasoning, the rely

condition R of a thread specifies the permitted state transitions that its environment may have, and its guarantee G specifies the possible transitions made by the thread itself. To ensure parallel threads can collaborate, we need to check the interference constraint, *i.e.*, the guarantee of each thread is permitted in the rely of every others. Then we can verify their parallel composition by separately verifying each thread, showing its behaviors under the rely condition indeed satisfy its guarantee. After parallel composition, the threads should be executed under their common environment (*i.e.*, the intersection of their relies) and guarantee all the possible transitions made by them (*i.e.*, the union of their guarantees).

Parametrized with rely/guarantee conditions for the two levels, our relation $(C, \mathcal{R}, \mathcal{G}) \preceq (\mathbb{C}, \mathbb{R}, \mathbb{G})$ talks about not only the target C and the source \mathbb{C} , but also the interference \mathcal{R} and \mathcal{G} between C and its target-level environment, and \mathbb{R} and \mathbb{G} between \mathbb{C} and its environment at the source level. Informally, $(C, \mathcal{R}, \mathcal{G}) \preceq (\mathbb{C}, \mathbb{R}, \mathbb{G})$ says the executions of C under the environment \mathcal{R} do not exhibit more observable behaviors than the executions of \mathbb{C} under the environment \mathbb{R} , and the state transitions of C and \mathbb{C} satisfy \mathcal{G} and \mathbb{G} respectively. RGSim is now compositional, as long as the threads are composed with well-behaved environments only. The parallel compositionality lemma is in the following form. If we know $(C_1, \mathcal{R}_1, \mathcal{G}_1) \preceq (\mathbb{C}_1, \mathbb{R}_1, \mathbb{G}_1)$ and $(C_2, \mathcal{R}_2, \mathcal{G}_2) \preceq (\mathbb{C}_2, \mathbb{R}_2, \mathbb{G}_2)$, and also the interference constraints are satisfied, *i.e.*, $\mathcal{G}_2 \subseteq \mathcal{R}_1$, $\mathbb{G}_1 \subseteq \mathbb{R}_2$, $\mathbb{G}_2 \subseteq \mathbb{R}_1$ and $\mathbb{G}_1 \subseteq \mathbb{R}_2$, we could get

$$(C_1 \parallel C_2, \mathcal{R}_1 \cap \mathcal{R}_2, \mathcal{G}_1 \cup \mathcal{G}_2) \preceq (\mathbb{C}_1 \parallel \mathbb{C}_2, \mathbb{R}_1 \cap \mathbb{R}_2, \mathbb{G}_1 \cup \mathbb{G}_2).$$

The compositionality of RGSim gives us a proof theory for concurrent program transformations.

Also different from fully abstract semantics for threads, which assumes arbitrary behaviors of environments, RGSim allows us to instantiate the interference \mathcal{R} , \mathcal{G} , \mathbb{R} and \mathbb{G} differently for different assumptions about environments, therefore it can be used to verify the aforementioned four classes of transformations. For instance, if we want to prove that a transformation preserves the behaviors of data-race-free programs, we can specify the data-race-freedom in \mathbb{R} and \mathbb{G} . Then we are no longer concerned with the examples in Figure 1, both of which have data races.

3. Basic Technical Settings

In this section, we present the source and the target programming languages. Then we define a basic refinement \sqsubseteq , which naturally says the target has no more externally observable event traces than the source. We use \sqsubseteq as an intuitive formulation of the correctness of transformations.

3.1 The Languages

Following standard simulation techniques, we model the semantics of target and source programs as labeled transition systems. Before showing the languages, we first define events and labels in Figure 2(a). We leave the set of events unspecified here. It can be instantiated by program verifiers, depending on their interest (*e.g.*, input/output events). A label that will be associated with a state transition is either an event or τ , which means the corresponding transition does not generate any event (*i.e.*, a silent step).

The target language, which we also call the low-level language, is shown in Figure 2(b). We abstract away the forms of states, expressions and primitive instructions in the language. An arithmetic expression E is modeled as a function from states to integers lifted with an undefined value \perp . Boolean expressions are modeled similarly. An instruction is a partial function from states to sets of label and state pairs, describing the state transitions and the events it generates. We use $\mathcal{P}(\cdot)$ to denote the power set. Unsafe executions

(Events) $e ::= \dots$ (Labels) $o ::= e \mid \tau$

(a) Events and Transition Labels

(LState) $\sigma ::= \dots$

(LExp) $E \in LState \rightarrow Int_{\perp}$

(LExp) $B \in LState \rightarrow \{\mathbf{true}, \mathbf{false}\}_{\perp}$

(LInstr) $c \in LState \rightarrow \mathcal{P}((Labels \times LState) \cup \{\mathbf{abort}\})$

(LStmt) $C ::= \mathbf{skip} \mid c \mid C_1; C_2 \mid \mathbf{if} (B) C_1 \mathbf{else} C_2$
 $\mid \mathbf{while} (B) C \mid C_1 \parallel C_2$

(LStep) $\rightarrow_L \in \mathcal{P}((LStmt/\{\mathbf{skip}\} \times LState) \times Labels$
 $\times ((LStmt \times LState) \cup \{\mathbf{abort}\}))$

(b) The Low-Level Language

(HState) $\Sigma ::= \dots$

(HExp) $\mathbb{E} \in HState \rightarrow Int_{\perp}$

(HExp) $\mathbb{B} \in HState \rightarrow \{\mathbf{true}, \mathbf{false}\}_{\perp}$

(HInstr) $c \in HState \rightarrow \mathcal{P}((Labels \times HState) \cup \{\mathbf{abort}\})$

(HStmt) $\mathbb{C} ::= \mathbf{skip} \mid c \mid \mathbb{C}_1; \mathbb{C}_2 \mid \mathbf{if} \mathbb{B} \mathbf{then} \mathbb{C}_1 \mathbf{else} \mathbb{C}_2$
 $\mid \mathbf{while} \mathbb{B} \mathbf{do} \mathbb{C} \mid \mathbb{C}_1 \parallel \mathbb{C}_2$

(HStep) $\rightarrow_L \in \mathcal{P}((HStmt/\{\mathbf{skip}\} \times HState) \times Labels$
 $\times ((HStmt \times HState) \cup \{\mathbf{abort}\}))$

(c) The High-Level Language

Figure 2. Generic Languages at Target and Source Levels

lead to **abort**. Note that the semantics of an instruction could be non-deterministic. Moreover, it might be undefined on some states, making it possible to model blocking operations such as acquiring a lock.

Statements are either primitive instructions or compositions of them. **skip** is a special statement used as a flag to show the end of executions. A single-step execution of statements is modeled as a labeled transition $_ \xrightarrow{L} _$, which is a triple of an initial program configuration (a pair of statement and state), a label and a resulting configuration. It is undefined when the initial statement is **skip**. The step aborts if an unsafe instruction is executed.

The high-level language (source language) is defined similarly in Figure 2(c), but it is important to note that its states and primitive instructions may be different from those in the low-level language. The compound statements are almost the same as their low-level counterparts. $\mathbb{C}_1; \mathbb{C}_2$ and $\mathbb{C}_1 \parallel \mathbb{C}_2$ are sequential and parallel compositions of \mathbb{C}_1 and \mathbb{C}_2 respectively. Note that we choose to use the same set of compound statements in the two languages for simplicity only. This is not required by our simulation relation, although the analogous program constructs of the two languages (e.g., parallel compositions $C_1 \parallel C_2$ and $C_1 \parallel \parallel C_2$) make it convenient for us to discuss the compositionality later.

Figure 3 shows part of the definition of $_ \xrightarrow{H} _$, which gives the high-level operational semantics of statements. We often omit the subscript H (or L) in $_ \xrightarrow{H} _$ (or $_ \xrightarrow{L} _$) and the label on top of the arrow when it is τ . The semantics is mostly standard. We only show the rules for primitive instructions and parallel compositions here. Note that when a primitive instruction c is blocked at state Σ (i.e., $\Sigma \notin \text{dom}(c)$), we let the program configuration reduce to itself. For example, the instruction `lock(1)` would be blocked when 1 is not 0, making it be repeated until 1 becomes 0; whereas `unlock(1)` simply sets 1 to 0 at any time and would never be blocked. Primitive instructions in the high-level and low-

level languages are *atomic* in the interleaving semantics. Below we use $_ \xrightarrow{*} _$ for zero or multiple-step transitions with no events generated, and $_ \xrightarrow{e} _$ for multiple-step transitions with *only one* event e generated.

3.2 The Event Trace Refinement

Now we can formally define the refinement relation \sqsubseteq that relates the set of externally observable event traces generated by the target and the source programs. A trace is a sequence of events e , and may end with a termination marker **done** or a fault marker **abort**.

(EvtTrace) $\mathcal{E} ::= \epsilon \mid \mathbf{done} \mid \mathbf{abort} \mid e::\mathcal{E}$

Definition 1 (Event Trace Set). $ETrSet_n(C, \sigma)$ represents a set of external event traces produced by C in n steps from the state σ :

- $ETrSet_0(C, \sigma) \triangleq \{\epsilon\}$;
- $ETrSet_{n+1}(C, \sigma) \triangleq$
 $\{\mathcal{E} \mid (C, \sigma) \xrightarrow{*} (C', \sigma') \wedge \mathcal{E} \in ETrSet_n(C', \sigma')$
 $\vee (C, \sigma) \xrightarrow{e} (C', \sigma') \wedge \mathcal{E}' \in ETrSet_n(C', \sigma') \wedge \mathcal{E} = e::\mathcal{E}'$
 $\vee (C, \sigma) \xrightarrow{*} \mathbf{abort} \wedge \mathcal{E} = \mathbf{abort}$
 $\vee C = \mathbf{skip} \wedge \mathcal{E} = \mathbf{done}\}$.

We define $ETrSet(C, \sigma)$ as $\bigcup_n ETrSet_n(C, \sigma)$.

We overload the notation and use $ETrSet(\mathbb{C}, \Sigma)$ for the high-level language. Then we define an event trace refinement as the subset relation between event trace sets, which is similar to Leroy's refinement property [19].

Definition 2 (Event Trace Refinement). We say (C, σ) is an *e-trace refinement* of (\mathbb{C}, Σ) , i.e., $(C, \sigma) \sqsubseteq (\mathbb{C}, \Sigma)$, if and only if

$$ETrSet(C, \sigma) \subseteq ETrSet(\mathbb{C}, \Sigma).$$

The refinement is defined for program configurations instead of for code only because the initial states may affect the behaviors of programs. In this case, the transformation \mathbf{T} should translate states as well as code. We overload the notation and use $\mathbf{T}(\Sigma)$ to represent the state transformation, and use $C \sqsubseteq_{\mathbf{T}} \mathbb{C}$ for

$$\forall \sigma, \Sigma. \sigma = \mathbf{T}(\Sigma) \implies (C, \sigma) \sqsubseteq (\mathbb{C}, \Sigma),$$

then $\text{Correct}(\mathbf{T})$ defined in formula (1.1) can be reformulated as

$$\text{Correct}(\mathbf{T}) \triangleq \forall C, \mathbb{C}. C = \mathbf{T}(\mathbb{C}) \implies C \sqsubseteq_{\mathbf{T}} \mathbb{C}. \quad (3.1)$$

4. The RGSim Relation

The e-trace refinement is defined directly over the externally observable behaviors of programs. It is intuitive, and also abstract in that it is independent of language details. However, as we explained before, it is *not* compositional *w.r.t.* parallel compositions. In this section we propose RGSim, which can be viewed as a compositional proof technique that allows us to derive the simple e-trace refinement and then verify the corresponding transformation \mathbf{T} .

4.1 The Definition

Our co-inductively defined RGSim relation is in the form of $(C, \sigma, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (C, \Sigma, \mathbb{R}, \mathbb{G})$, which is a simulation between program configurations (C, σ) and (\mathbb{C}, Σ) . It is parametrized with the rely and guarantee conditions at the low level and the high level, which are binary relations over states:

$$\mathcal{R}, \mathcal{G} \in \mathcal{P}(LState \times LState), \quad \mathbb{R}, \mathbb{G} \in \mathcal{P}(HState \times HState).$$

The simulation also takes two additional parameters: the *step invariant* α and the *postcondition* γ , which are both relations between the low-level and the high-level states.

$$\alpha, \gamma, \zeta \in \mathcal{P}(LState \times HState).$$

$$\begin{array}{c}
\frac{(\tau, \Sigma') \in c \Sigma}{(c, \Sigma) \longrightarrow (\mathit{skip}, \Sigma')} \quad \frac{(e, \Sigma') \in c \Sigma}{(c, \Sigma) \xrightarrow{e} (\mathit{skip}, \Sigma')} \quad \frac{\mathbf{abort} \in c \Sigma}{(c, \Sigma) \longrightarrow \mathbf{abort}} \quad \frac{\Sigma \notin \mathit{dom}(c)}{(c, \Sigma) \longrightarrow (c, \Sigma)} \\
\frac{}{(\mathit{skip} \parallel \mathit{skip}, \Sigma) \longrightarrow (\mathit{skip}, \Sigma)} \quad \frac{(\mathbb{C}_1, \Sigma) \longrightarrow (\mathbb{C}'_1, \Sigma')}{(\mathbb{C}_1 \parallel \mathbb{C}_2, \Sigma) \longrightarrow (\mathbb{C}'_1 \parallel \mathbb{C}_2, \Sigma')} \quad \frac{(\mathbb{C}_2, \Sigma) \longrightarrow (\mathbb{C}'_2, \Sigma')}{(\mathbb{C}_1 \parallel \mathbb{C}_2, \Sigma) \longrightarrow (\mathbb{C}_1 \parallel \mathbb{C}'_2, \Sigma')} \\
\frac{(\mathbb{C}_1, \Sigma) \xrightarrow{e} (\mathbb{C}'_1, \Sigma')}{(\mathbb{C}_1 \parallel \mathbb{C}_2, \Sigma) \xrightarrow{e} (\mathbb{C}'_1 \parallel \mathbb{C}_2, \Sigma')} \quad \frac{(\mathbb{C}_2, \Sigma) \xrightarrow{e} (\mathbb{C}'_2, \Sigma')}{(\mathbb{C}_1 \parallel \mathbb{C}_2, \Sigma) \xrightarrow{e} (\mathbb{C}_1 \parallel \mathbb{C}'_2, \Sigma')} \quad \frac{(\mathbb{C}_1, \Sigma) \longrightarrow \mathbf{abort} \quad \text{or} \quad (\mathbb{C}_2, \Sigma) \longrightarrow \mathbf{abort}}{(\mathbb{C}_1 \parallel \mathbb{C}_2, \Sigma) \longrightarrow \mathbf{abort}}
\end{array}$$

Figure 3. Operational Semantics of the High-Level Language

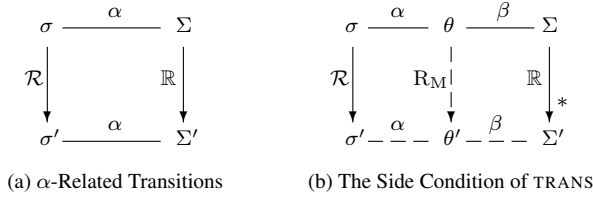


Figure 4. Related Transitions

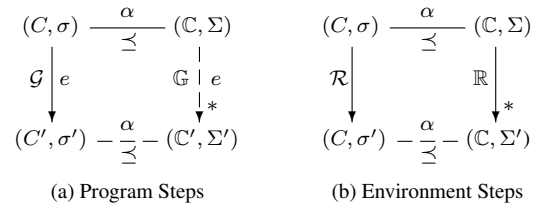


Figure 5. Simulation Diagrams of RGSim

Before we formally define RGSim in Definition 4, we first introduce the α -related transitions as follows.

Definition 3 (α -Related Transitions).

$$\langle \mathcal{R}, \mathbb{R} \rangle_\alpha \triangleq \{((\sigma, \sigma'), (\Sigma, \Sigma')) \mid (\sigma, \sigma') \in \mathcal{R} \wedge (\Sigma, \Sigma') \in \mathbb{R} \wedge (\sigma, \Sigma) \in \alpha \wedge (\sigma', \Sigma') \in \alpha\}.$$

$\langle \mathcal{R}, \mathbb{R} \rangle_\alpha$ represents a set of the α -related transitions in \mathcal{R} and \mathbb{R} , putting together the corresponding transitions in \mathcal{R} and \mathbb{R} that can be related by α , as illustrated in Figure 4(a). $\langle \mathcal{G}, \mathbb{G} \rangle_\alpha$ is defined in the same way.

Definition 4 (RGSim). Whenever $(C, \sigma, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (\mathbb{C}, \Sigma, \mathbb{R}, \mathbb{G})$, then $(\sigma, \Sigma) \in \alpha$ and the following are true:

1. if $(C, \sigma) \longrightarrow (C', \sigma')$, then there exist C' and Σ' such that $(\mathbb{C}, \Sigma) \longrightarrow^* (\mathbb{C}', \Sigma')$, $((\sigma, \sigma'), (\Sigma, \Sigma')) \in \langle \mathcal{G}, \mathbb{G}^* \rangle_\alpha$ and $(C', \sigma', \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (C', \Sigma', \mathbb{R}, \mathbb{G})$;
2. if $(C, \sigma) \xrightarrow{e} (C', \sigma')$, then there exist C' and Σ' such that $(\mathbb{C}, \Sigma) \xrightarrow{e}^* (\mathbb{C}', \Sigma')$, $((\sigma, \sigma'), (\Sigma, \Sigma')) \in \langle \mathcal{G}, \mathbb{G}^* \rangle_\alpha$ and $(C', \sigma', \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (C', \Sigma', \mathbb{R}, \mathbb{G})$;
3. if $C = \mathit{skip}$, then there exists Σ' such that $(\mathbb{C}, \Sigma) \longrightarrow^* (\mathit{skip}, \Sigma')$, $((\sigma, \sigma'), (\Sigma, \Sigma')) \in \langle \mathcal{G}, \mathbb{G}^* \rangle_\alpha$, $(\sigma, \Sigma') \in \gamma$ and $\gamma \subseteq \alpha$;
4. if $(C, \sigma) \longrightarrow \mathbf{abort}$, then $(\mathbb{C}, \Sigma) \longrightarrow^* \mathbf{abort}$;
5. if $((\sigma, \sigma'), (\Sigma, \Sigma')) \in \langle \mathcal{R}, \mathbb{R}^* \rangle_\alpha$, then $(C, \sigma', \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (\mathbb{C}, \Sigma', \mathbb{R}, \mathbb{G})$.

Then, $(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (\mathbb{C}, \mathbb{R}, \mathbb{G})$ iff for all σ and Σ , if $(\sigma, \Sigma) \in \zeta$, then $(C, \sigma, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (\mathbb{C}, \Sigma, \mathbb{R}, \mathbb{G})$. Here the *precondition* ζ is used to relate the initial states σ and Σ .

Informally, $(C, \sigma, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (\mathbb{C}, \Sigma, \mathbb{R}, \mathbb{G})$ says the low-level configuration (C, σ) is simulated by the high-level configuration (\mathbb{C}, Σ) with behaviors \mathcal{G} and \mathbb{G} respectively, no matter how their environments \mathcal{R} and \mathbb{R} interfere with them. It requires the following hold for every execution of C :

- Starting from α -related states, each step of C corresponds to zero or multiple steps of \mathbb{C} , and the resulting states are α -related too. If an external event is produced in the step of

C , the same event should be produced by \mathbb{C} . We show the simulation diagram with events generated by the program steps in Figure 5(a), where solid lines denote hypotheses and dashed lines denote conclusions, following Leroy's notations [19].

- The α relation reflects the abstractions from the low-level machine model to the high-level one, and is preserved by the related transitions at the two levels (so it is an *invariant*). For instance, when verifying a fine-grained implementation of sets, the α relation may relate a concrete representation in memory (e.g., a linked-list) at the low level to the corresponding abstract mathematical set at the high level.
- The corresponding transitions of C and \mathbb{C} need to be in $\langle \mathcal{G}, \mathbb{G}^* \rangle_\alpha$. That is, for each step of C , its state transition should satisfy the guarantee \mathcal{G} , and the corresponding transition made by the multiple steps of \mathbb{C} should be in the transitive closure of \mathbb{G} . The guarantees are abstractions of the programs' behaviors. As we will show later in the PAR rule in Figure 7, they will serve as the rely conditions of the sibling threads at the time of parallel compositions. Note that we do not need each step of \mathbb{C} to be in \mathbb{G} , although we could do so. This is because we only care about the coarse-grained behaviors (with mumbling) of the source that are used to simulate the target. We will explain more by the example (4.1) in Section 4.2.
- If C terminates, then \mathbb{C} terminates as well, and the final states should be related by the postcondition γ . We require $\gamma \subseteq \alpha$, i.e., the final state relation is not weaker than the step invariant.
- C is not safe only if \mathbb{C} is not safe either. This means the transformation should not make a safe high-level program unsafe at the low level.
- Whatever the low-level environment \mathcal{R} and the high-level one \mathbb{R} do, as long as the state transitions are α -related, they should not affect the simulation between C and \mathbb{C} , as shown in Figure 5(b). Here a step in \mathcal{R} may correspond to zero or multiple steps of \mathbb{R} . Note that different from the program steps, for the environment steps we do not require each step of \mathcal{R} to correspond to zero or multiple steps of \mathbb{R} . On the other hand, only requiring that \mathcal{R} be simulated by \mathbb{R} is not sufficient for parallel compositionality, which we will explain later in Section 4.2.

$$\begin{aligned}
\text{InitRel}_{\mathbf{T}}(\zeta) &\triangleq \forall \sigma, \Sigma. \sigma = \mathbf{T}(\Sigma) \implies (\sigma, \Sigma) \in \zeta \\
B \Leftrightarrow \mathbb{B} &\triangleq \{(\sigma, \Sigma) \mid B \sigma = \mathbb{B} \Sigma\} \quad B \mathbb{M} \mathbb{B} \triangleq \{(\sigma, \Sigma) \mid B \sigma \wedge \mathbb{B} \Sigma\} \\
\text{Intuit}(\alpha) &\triangleq \forall \sigma, \Sigma, \sigma', \Sigma'. (\sigma, \Sigma) \in \alpha \wedge \sigma \subseteq \sigma' \wedge \Sigma \subseteq \Sigma' \\
&\implies (\sigma', \Sigma') \in \alpha \\
\eta \# \alpha &\triangleq (\eta \cap \alpha) \subseteq (\eta \cup \alpha) \\
\beta \circ \alpha &\triangleq \{(\sigma, \Sigma) \mid \exists \theta. (\sigma, \theta) \in \alpha \wedge (\theta, \Sigma) \in \beta\} \\
\alpha \uplus \beta &\triangleq \{(\sigma_1 \uplus \sigma_2, \Sigma_1 \uplus \Sigma_2) \mid (\sigma_1, \Sigma_1) \in \alpha \wedge (\sigma_2, \Sigma_2) \in \beta\} \\
\text{Id} &\triangleq \{(\sigma, \sigma) \mid \sigma \in \text{LState}\} \quad \text{True} \triangleq \{(\sigma, \sigma') \mid \sigma, \sigma' \in \text{LState}\} \\
\mathbb{R}_{\mathbb{M}} \text{ isMidOf}(\alpha, \beta; \mathcal{R}, \mathbb{R}) &\triangleq \forall \sigma, \sigma', \Sigma, \Sigma'. \\
&((\sigma, \sigma'), (\Sigma, \Sigma')) \in \langle \mathcal{R}, \mathbb{R} \rangle_{\beta \circ \alpha} \\
&\implies \forall \theta. (\sigma, \theta) \in \alpha \wedge (\theta, \Sigma) \in \beta \\
&\implies \exists \theta'. ((\sigma, \sigma'), (\theta, \theta')) \in \langle \mathcal{R}, \mathbb{R}_{\mathbb{M}} \rangle_{\alpha} \wedge ((\theta, \theta'), (\Sigma, \Sigma')) \in \langle \mathbb{R}_{\mathbb{M}}, \mathbb{R} \rangle_{\beta}
\end{aligned}$$

Figure 6. Auxiliary Definitions for RGSim

Then based on the simulation, we hide the states by the precondition ζ and define the RGSim relation between programs only. By the definition we know $\zeta \subseteq \alpha$ if $(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G})$, *i.e.*, the precondition needs to be no weaker than the step invariant.

RGSim is sound *w.r.t.* the e-trace refinement (Definition 2). That is, $(C, \sigma, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (C, \Sigma, \mathbb{R}, \mathbb{G})$ ensures that (C, σ) does not have more observable behaviors than (C, Σ) .

Theorem 5 (Soundness). *If there exist $\mathcal{R}, \mathcal{G}, \mathbb{R}, \mathbb{G}, \alpha$ and γ such that $(C, \sigma, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma} (C, \Sigma, \mathbb{R}, \mathbb{G})$, then $(C, \sigma) \sqsubseteq (C, \Sigma)$.*

The soundness theorem can be proved by first strengthening the relies to the identity transitions and weakening the guarantees to the universal relations. Then we prove that the resulting simulation under identity environments implies the e-trace refinement.

For program transformations, since the initial state for the target program is transformed from the initial state for the source, we use $\text{InitRel}_{\mathbf{T}}(\zeta)$ (defined in Figure 6) to say the transformation \mathbf{T} over states ensures the binary precondition ζ .

Corollary 6. *If there exist $\mathcal{R}, \mathcal{G}, \mathbb{R}, \mathbb{G}, \alpha, \zeta$ and γ such that $\text{InitRel}_{\mathbf{T}}(\zeta)$ and $(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G})$, then $C \sqsubseteq_{\mathbf{T}} C$.*

4.2 Compositionality Rules

RGSim is compositional *w.r.t.* various program constructs, including parallel compositions. We present the compositionality rules in Figure 7, which gives us a relational proof method for concurrent program transformations.

As in the R-G logic [17], we require that the pre- and post-conditions be *stable* under the interference from the environments. Here we introduce the concept of stability of a relation ζ *w.r.t.* a set of transition pairs $\Lambda \in \mathcal{P}(\text{LState} \times \text{LState}) \times (\text{HState} \times \text{HState})$.

Definition 7 (Stability). $\text{Sta}(\zeta, \Lambda)$ holds iff for all σ, σ', Σ and Σ' , if $(\sigma, \Sigma) \in \zeta$ and $((\sigma, \sigma'), (\Sigma, \Sigma')) \in \Lambda$, then $(\sigma', \Sigma') \in \zeta$.

Usually we need $\text{Sta}(\zeta, \langle \mathcal{R}, \mathbb{R}^* \rangle_{\alpha})$, which says whenever ζ holds initially and \mathcal{R} and \mathbb{R}^* perform related actions, the resulting states still satisfy ζ . By unfolding $\langle \mathcal{R}, \mathbb{R}^* \rangle_{\alpha}$, we could see that α itself is stable *w.r.t.* any α -related transitions, *i.e.*, $\text{Sta}(\alpha, \langle \mathcal{R}, \mathbb{R}^* \rangle_{\alpha})$. Another simple example is given below, where both environments could increment x and the unary stable assertion $\{x \geq 0\}$ is lifted to the relation ζ :

$$\begin{aligned}
\zeta &\triangleq \{(\sigma, \Sigma) \mid \sigma(x) = \Sigma(x) \geq 0\} \quad \alpha \triangleq \{(\sigma, \Sigma) \mid \sigma(x) = \Sigma(x)\} \\
\mathcal{R} &\triangleq \{(\sigma, \sigma') \mid \sigma' = \sigma\{x \rightsquigarrow \sigma(x) + 1\}\} \\
\mathbb{R} &\triangleq \{(\Sigma, \Sigma') \mid \Sigma' = \Sigma\{x \rightsquigarrow \Sigma(x) + 1\}\}
\end{aligned}$$

We can prove $\text{Sta}(\zeta, \langle \mathcal{R}, \mathbb{R}^* \rangle_{\alpha})$. Stability of the pre- and post-conditions under the environments' interference is assumed as an implicit side condition at every proof rule in Figure 7, *e.g.*, we assume $\text{Sta}(\zeta, \langle \mathcal{R}, \mathbb{R}^* \rangle_{\alpha})$ in the SKIP rule. We also require implicitly that the relies and guarantees are closed over identity transitions, since stuttering steps will not affect observable event traces.

In Figure 7, the rules SKIP, SEQ, IF and WHILE reveal a high degree of similarity to the corresponding inference rules in Hoare logic. In the SEQ rule, γ serves as the postcondition of C_1 and \mathbb{C}_1 and the precondition of C_2 and \mathbb{C}_2 at the same time. The IF rule requires the boolean conditions of both sides to be evaluated to the same value under the precondition ζ . We give the definitions of the sets $B \Leftrightarrow \mathbb{B}$ and $B \mathbb{M} \mathbb{B}$ in Figure 6. The rule also requires the precondition ζ to imply the step invariant α . In the WHILE rule, the γ relation is viewed as a loop invariant preserved at the loop entry point, and needs to ensure $B \Leftrightarrow \mathbb{B}$.

Parallel compositionality. The PAR rule shows parallel compositionality of RGSim. The interference constraints say that two threads can be composed in parallel if one thread's guarantee implies the rely of the other. After parallel composition, they are expected to run in the common environment and their guaranteed behaviors contain each single thread's behaviors.

Note that, although RGSim does not require every step of the high-level program to be in its guarantee (see the first two conditions in Definition 4), this relaxation does not affect the parallel compositionality. This is because the target could have less behaviors than the source. To let $\mathbb{C}_1 \parallel \mathbb{C}_2$ simulate $C_1 \parallel C_2$, we only need a subset of the interleavings of \mathbb{C}_1 and \mathbb{C}_2 to simulate those of C_1 and C_2 . Thus the high-level relies and guarantees need to ensure the existence of those interleavings only. Below we give a simple example to explain this subtle issue. We can prove

$$(x := x + 2, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (x := x + 1; x := x + 1, \mathbb{R}, \mathbb{G}), \quad (4.1)$$

where the relies and the guarantees say x can be increased by 2 and α, ζ and γ relate x of the two sides:

$$\begin{aligned}
\mathcal{R} = \mathcal{G} &\triangleq \{(\sigma, \sigma') \mid \sigma' = \sigma \vee \sigma' = \sigma\{x \rightsquigarrow \sigma(x) + 2\}\}; \\
\mathbb{R} = \mathbb{G} &\triangleq \{(\Sigma, \Sigma') \mid \Sigma' = \Sigma \vee \Sigma' = \Sigma\{x \rightsquigarrow \Sigma(x) + 2\}\}; \\
\alpha = \zeta = \gamma &\triangleq \{(\sigma, \Sigma) \mid \sigma(x) = \Sigma(x)\}.
\end{aligned}$$

Note that the high-level program is actually finer-grained than its guarantee, but to prove (4.1) we only need the execution in which it goes two steps to the end without interference from its environment. Also we can prove $(\text{print}(x), \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (\text{print}(x), \mathbb{R}, \mathbb{G})$. Then by the PAR rule, we get

$$\begin{aligned}
& (x := x + 2 \parallel \text{print}(x), \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} \\
& (x := x + 1; x := x + 1) \parallel \text{print}(x), \mathbb{R}, \mathbb{G},
\end{aligned}$$

which does not violate the natural meaning of refinements. That is, all the possible external events produced by the low-level side can also be produced by the high-level side, although the latter could have more external behaviors due to its finer granularity.

Another subtlety in the RGSim definition is with the fifth condition over the environments, which is crucial for parallel compositionality. One may think a more natural alternative to this condition is to require that \mathcal{R} be simulated by \mathbb{R} :

$$\text{If } (\sigma, \sigma') \in \mathcal{R}, \text{ then there exists } \Sigma' \text{ such that } \\
(\Sigma, \Sigma') \in \mathbb{R}^* \text{ and } (C, \sigma', \mathcal{R}, \mathcal{G}) \preceq'_{\alpha; \zeta \times \gamma} (C, \Sigma', \mathbb{R}, \mathbb{G}). \quad (4.2)$$

We refer to this modified simulation definition as \preceq' . Unfortunately, \preceq' does not have parallel compositionality. As a counterexample, if the invariant α says the left-side x is not greater than the right-side x and the precondition ζ requires x of the two sides are equal, *i.e.*,

$$\alpha \triangleq \{(\sigma, \Sigma) \mid \sigma(x) \leq \Sigma(x)\} \quad \zeta \triangleq \{(\sigma, \Sigma) \mid \sigma(x) = \Sigma(x)\},$$

$$\begin{array}{c}
\frac{\zeta \subseteq \alpha}{(\text{skip}, \mathcal{R}, \text{ld}) \preceq_{\alpha; \zeta \times \zeta} (\text{skip}, \mathbb{R}, \text{ld})} \text{ (SKIP)} \quad \frac{(C_1, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (C_1, \mathbb{R}, \mathbb{G}) \quad (C_2, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma \times \eta} (C_2, \mathbb{R}, \mathbb{G})}{(C_1; C_2, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \eta} (C_1; C_2, \mathbb{R}, \mathbb{G})} \text{ (SEQ)} \\
\\
\frac{\zeta \subseteq (B \Leftrightarrow \mathbb{B}) \quad \zeta_1 = (\zeta \cap (B \wedge \mathbb{B})) \quad \zeta_2 = (\zeta \cap (\neg B \wedge \neg \mathbb{B})) \quad \zeta \subseteq \alpha}{(\text{if } (B) C_1 \text{ else } C_2, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (\text{if } \mathbb{B} \text{ then } C_1 \text{ else } C_2, \mathbb{R}, \mathbb{G})} \text{ (IF)} \\
\\
\frac{(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma_1 \times \gamma} (C, \mathbb{R}, \mathbb{G}) \quad \gamma \subseteq (B \Leftrightarrow \mathbb{B}) \quad \gamma_1 = (\gamma \cap (B \wedge \mathbb{B})) \quad \gamma_2 = (\gamma \cap (\neg B \wedge \neg \mathbb{B}))}{(\text{while } (B) C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \gamma \times \gamma_2} (\text{while } \mathbb{B} \text{ do } C, \mathbb{R}, \mathbb{G})} \text{ (WHILE)} \\
\\
\frac{\begin{array}{c} (C_1, \mathcal{R}_1, \mathcal{G}_1) \preceq_{\alpha; \zeta \times \gamma_1} (C_1, \mathbb{R}_1, \mathbb{G}_1) \quad (C_2, \mathcal{R}_2, \mathcal{G}_2) \preceq_{\alpha; \zeta \times \gamma_2} (C_2, \mathbb{R}_2, \mathbb{G}_2) \\ \mathcal{G}_1 \subseteq \mathcal{R}_2 \quad \mathcal{G}_2 \subseteq \mathcal{R}_1 \quad \mathbb{G}_1 \subseteq \mathbb{R}_2 \quad \mathbb{G}_2 \subseteq \mathbb{R}_1 \end{array}}{(C_1 \parallel C_2, \mathcal{R}_1 \cap \mathcal{R}_2, \mathcal{G}_1 \cup \mathcal{G}_2) \preceq_{\alpha; \zeta \times (\gamma_1 \cap \gamma_2)} (C_1 \parallel C_2, \mathbb{R}_1 \cap \mathbb{R}_2, \mathbb{G}_1 \cup \mathbb{G}_2)} \text{ (PAR)} \\
\\
\frac{\begin{array}{c} (C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G}) \\ (\zeta \cup \gamma) \subseteq \alpha' \subseteq \alpha \quad \text{Sta}(\alpha', \langle \mathcal{G}, \mathbb{G}^* \rangle_\alpha) \end{array}}{(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha'; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G})} \text{ (STREN-}\alpha\text{)} \quad \frac{\begin{array}{c} (C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G}) \\ \alpha \subseteq \alpha' \quad \text{Sta}(\alpha, \langle \mathcal{R}, \mathbb{R}^* \rangle_{\alpha'}) \end{array}}{(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha'; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G})} \text{ (WEAKEN-}\alpha\text{)} \\
\\
\frac{(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G}) \quad \zeta' \subseteq \zeta \quad \gamma \subseteq \gamma' \subseteq \alpha \quad \mathcal{R}' \subseteq \mathcal{R} \quad \mathbb{R}' \subseteq \mathbb{R} \quad \mathcal{G} \subseteq \mathcal{G}' \quad \mathbb{G} \subseteq \mathbb{G}'}{(C, \mathcal{R}', \mathcal{G}') \preceq_{\alpha; \zeta' \times \gamma'} (C, \mathbb{R}', \mathbb{G}')} \text{ (CONSEQ)} \\
\\
\frac{\begin{array}{c} (C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (C, \mathbb{R}, \mathbb{G}) \\ \eta \subseteq \beta \quad \text{Intuit}(\{\alpha, \zeta, \gamma, \beta, \eta, \mathcal{R}, \mathbb{R}, \mathcal{R}_1, \mathbb{R}_1\}) \\ \eta \# \{\zeta, \gamma, \alpha\} \quad \text{Sta}(\eta, \{\langle \mathcal{G}, \mathbb{G}^* \rangle_\alpha, \langle \mathcal{R}_1, \mathbb{R}_1^* \rangle_\beta\}) \end{array}}{(C, \mathcal{R} \uplus \mathcal{R}_1, \mathcal{G} \uplus \mathcal{G}_1) \preceq_{\alpha \uplus \beta; (\zeta \uplus \eta) \times (\gamma \uplus \eta)} (C, \mathbb{R} \uplus \mathbb{R}_1, \mathbb{G} \uplus \mathbb{G}_1)} \text{ (FRAME)} \quad \frac{\begin{array}{c} (C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (M, \mathbb{R}_M, \mathbb{G}_M) \\ (M, \mathbb{R}_M, \mathbb{G}_M) \preceq_{\beta; \delta \times \eta} (C, \mathbb{R}, \mathbb{G}) \\ \mathbb{R}_M \text{ isMidOf } (\alpha, \beta; \mathcal{R}, \mathbb{R}^*) \end{array}}{(C, \mathcal{R}, \mathcal{G}) \preceq_{\beta \circ \alpha; (\delta \circ \zeta) \times (\eta \circ \gamma)} (C, \mathbb{R}, \mathbb{G})} \text{ (TRANS)}
\end{array}$$

Figure 7. Compositionality Rules for RGSim

we could prove the following:

$$\begin{array}{l}
(x := x+1, \text{ld}, \text{True}) \preceq'_{\alpha; \zeta \times \alpha} (x := x+2, \text{ld}, \text{True}); \\
(\text{print}(x), \text{True}, \text{ld}) \preceq'_{\alpha; \zeta \times \alpha} (\text{print}(x), \text{True}, \text{ld}).
\end{array}$$

Here we use `ld` and `True` (defined in Figure 6) for the sets of identity transitions and arbitrary transitions respectively, and overload the notations at the low level to the high level. However, the following refinement does *not* hold after parallel composition:

$$\begin{array}{l}
(x := x+1 \parallel \text{print}(x), \text{ld}, \text{True}) \preceq'_{\alpha; \zeta \times \alpha} \\
(x := x+2 \parallel \text{print}(x), \text{ld}, \text{True}).
\end{array}$$

This is because the rely \mathcal{R} (or \mathbb{R}) is an abstraction of all the permitted behaviors in the environment of a thread. But a concrete sibling thread that runs in parallel may produce less transitions than \mathcal{R} (or \mathbb{R}). To obtain parallel compositionality, we need to ensure that the simulation holds for all concrete sibling threads. With our definition \preceq , the refinement $(\text{print}(x), \text{True}, \text{ld}) \preceq_{\alpha; \zeta \times \alpha} (\text{print}(x), \text{True}, \text{ld})$ is not provable because, after the environments' α -related transitions, the target may print a value smaller than the one printed by the source.

Other rules. We also develop some other useful rules about RGSim. For example, the `STREN- α` rule allows us to replace the invariant α by a stronger invariant α' . We need to check that α' is indeed an invariant preserved by the related program steps, *i.e.*, $\text{Sta}(\alpha', \langle \mathcal{G}, \mathbb{G}^* \rangle_\alpha)$ holds. Symmetrically, the `WEAKEN- α` rule requires α to be preserved by environment steps related by the weaker invariant α' . As usual, the pre/post conditions, the relies and the guarantees can be strengthened or weakened by the `CONSEQ` rule.

The `FRAME` rule allows us to use local specifications. When verifying the simulation between C and \mathbb{C} , we need to only talk about the locally-used resource in α , ζ and γ , and the local relies and guarantees \mathcal{R} , \mathcal{G} , \mathbb{R} and \mathbb{G} . Then the proof can be reused in contexts where some extra resource η is used, and the accesses of it

respect the invariant β and \mathcal{R}_1 , \mathbb{G}_1 and \mathbb{G}_1 . We give the auxiliary definitions in Figure 6. The disjoint union \uplus between states is lifted to state pairs. An intuitionistic state relation is monotone *w.r.t.* the extension of states. The disjointness $\eta \# \alpha$ says that any state pair satisfying both η and α can be split into two disjoint state pairs satisfying η and α respectively. For example, let $\eta \triangleq \{(\sigma, \Sigma) \mid \sigma(y) = \Sigma(y)\}$ and $\alpha \triangleq \{(\sigma, \Sigma) \mid \sigma(x) = \Sigma(x)\}$, then both η and α are intuitionistic and $\eta \# \alpha$ holds. We also require η to be stable under interference from the programs (*i.e.*, the programs do not change the extra resource) and the extra environments. We use $\eta \# \{\zeta, \gamma, \alpha\}$ as a shorthand for $(\eta \# \zeta) \wedge (\eta \# \gamma) \wedge (\eta \# \alpha)$. Similar representations are used in this rule.

Finally, the transitivity rule `TRANS` allows us to verify a transformation by using an intermediate level as a bridge. The intermediate environment \mathbb{R}_M should be chosen with caution so that the $(\beta \circ \alpha)$ -related transitions can be decomposed into β -related and α -related transitions, as illustrated in Figure 4(b). Here \circ defines the composition of two relations and `isMidOf` defines the side condition over the environments, as shown in Figure 6. We use θ for a middle-level state.

Soundness of all the rules in Figure 7 is proved in the technical report [20], showing that for each rule the premises imply the conclusion. The proofs are also mechanized in the Coq proof assistant [10].

Instantiations of relies and guarantees. We can derive the sequential refinement and the fully-abstract-semantics-based refinement by instantiating the rely conditions in RGSim. For example, the refinement (4.3) over closed programs assumes identity environments, making the interference constraints in the `PAR` rule unsatisfiable. This confirms the observation in Section 2.1 that the sequential refinement loses parallel compositionality.

$$(C, \text{ld}, \text{True}) \preceq_{\alpha; \zeta \times \gamma} (C, \text{ld}, \text{True}) \quad (4.3)$$

The refinement (4.4) assumes arbitrary environments, which makes the interference constraints in the PAR rule trivially true. But this assumption is too strong: usually (4.4) cannot be satisfied in practice.

$$(C, \text{True}, \text{True}) \preceq_{\alpha; \zeta \times \gamma} (\mathbb{C}, \text{True}, \text{True}) \quad (4.4)$$

4.3 A Simple Example

Below we give a simple example to illustrate the use of RGSim and its parallel compositionality in verifying concurrent program transformations. The high-level program $\mathbb{C}_1 \parallel \mathbb{C}_2$ is transformed to $C_1 \parallel C_2$, using a lock l to synchronize the accesses of the shared variable x . We aim to prove $C_1 \parallel C_2 \sqsubseteq_{\mathbf{T}} \mathbb{C}_1 \parallel \mathbb{C}_2$. That is, although $x := x + 2$ is implemented by two steps of incrementing x in C_2 , the parallel observer C_1 will not print unexpected values. Here we view output events as externally observable behaviors.

```

print(x);  |||  x := x + 2;
           ↓
lock(l);   |||  lock(l);
print(x);  |||  x := x+1; x := x+1;
unlock(l); |||  ⟨unlock(l); X := x;⟩

```

To facilitate the proof, we introduce an auxiliary shared variable X at the low level to record the value of x at the time when releasing the lock. It specifies the value of x outside every critical section, thus should match the value of the high-level x after every corresponding action. Here $\langle C \rangle$ means C is executed atomically.

By the soundness and compositionality of RGSim, we only need to prove simulations over individual threads, providing appropriate relies and guarantees. We first define the invariant α , which only cares about the value of x when the lock is free.

$$\alpha \triangleq \{(\sigma, \Sigma) \mid \sigma(x) = \Sigma(x) \wedge (\sigma(l) = 0 \implies \sigma(x) = \sigma(X))\}.$$

We let the pre- and post-conditions be α as well.

The high-level threads can be executed in arbitrary environments with arbitrary guarantees: $\mathbb{R} = \mathbb{G} \triangleq \text{True}$. The transformation uses the lock to protect every access of x , thus the low-level relies and guarantees are not arbitrary:

$$\begin{aligned} \mathcal{R} &\triangleq \{(\sigma, \sigma') \mid \sigma(l) = \text{cid} \implies \\ &\quad \sigma(x) = \sigma'(x) \wedge \sigma(X) = \sigma'(X) \wedge \sigma(l) = \sigma'(l)\}; \\ \mathcal{G} &\triangleq \{(\sigma, \sigma') \mid \sigma' = \sigma \vee \sigma(l) = 0 \wedge \sigma' = \sigma \{1 \rightsquigarrow \text{cid}\} \\ &\quad \vee \sigma(l) = \text{cid} \wedge \sigma' = \sigma \{x \rightsquigarrow \cdot\} \\ &\quad \vee \sigma(l) = \text{cid} \wedge \sigma' = \sigma \{1 \rightsquigarrow 0, X \rightsquigarrow \cdot\}\}. \end{aligned}$$

Every low-level thread guarantees that it updates x only when the lock is acquired. Its environment cannot update x or l if the current thread holds the lock. Here cid is the identifier of the current thread. When acquired, the lock holds the id of the owner thread.

Following the definition, we can prove $(C_1, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \alpha \times \alpha} (\mathbb{C}_1, \mathbb{R}, \mathbb{G})$ and $(C_2, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \alpha \times \alpha} (\mathbb{C}_2, \mathbb{R}, \mathbb{G})$. By applying the PAR rule and from the soundness of RGSim (Corollary 6), we know $C_1 \parallel C_2 \sqsubseteq_{\mathbf{T}} \mathbb{C}_1 \parallel \mathbb{C}_2$ holds for any \mathbf{T} that respects α .

Perhaps interestingly, if we omit the lock and unlock operations in C_1 , then $C_1 \parallel C_2$ would have more externally observable behaviors than $\mathbb{C}_1 \parallel \mathbb{C}_2$. This does *not* indicate the unsoundness of our PAR rule (which is sound!). The reason is that x might have different values on the two levels after the environments' α -related transitions, so that we cannot have $(\text{print}(x), \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \alpha \times \alpha} (\text{print}(x), \mathbb{R}, \mathbb{G})$ with the current definitions of α , \mathcal{R} and \mathcal{G} , even though the code of the two sides are syntactically identical.

More discussions. RGSim ensures that the target program preserves safety properties (including the partial correctness) of the source, but allows a terminating source program to be transformed

to a target having infinite silent steps. In the above example, this allows the low-level programs to be blocked forever (*e.g.*, at the time when the lock is held but never released by some other thread). Proving the preservation of the termination behavior would require liveness proofs in a concurrent setting (*e.g.*, proving the absence of deadlock), which we leave as future work.

In the next three sections, we show more serious examples to demonstrate the applicability of RGSim.

5. Relational Reasoning about Optimizations

As a general correctness notion of concurrent program transformations, RGSim establishes a relational approach to justify compiler optimizations on concurrent programs. Below we adapt Benton's work [3] on sequential optimizations to the concurrent setting.

5.1 Optimization Rules

Usually optimizations depend on particular contexts, *e.g.*, the assignment $x := E$ can be eliminated only in the context that the value of x is never used after the assignment. In a shared-state concurrent setting, we should also consider the parallel context for an optimization. RGSim enables us to specify various sophisticated requirements for the parallel contexts by rely/guarantee conditions. Based on RGSim, we provide a set of inference rules to characterize and justify common optimizations (*e.g.*, dead code elimination) with information of both the sequential and the parallel contexts. Due to the space limit, we only present some interesting rules here and leave other rules in the technical report [20]. Note in this section the target and the source programs are in the same language.

Sequential skip Law

$$\frac{(C_1, \mathcal{R}_1, \mathcal{G}_1) \preceq_{\alpha; \zeta \times \gamma} (C_2, \mathcal{R}_2, \mathcal{G}_2)}{(\text{skip}; C_1, \mathcal{R}_1, \mathcal{G}_1) \preceq_{\alpha; \zeta \times \gamma} (C_2, \mathcal{R}_2, \mathcal{G}_2)}$$

$$\frac{(C_1, \mathcal{R}_1, \mathcal{G}_1) \preceq_{\alpha; \zeta \times \gamma} (C_2, \mathcal{R}_2, \mathcal{G}_2)}{(C_1, \mathcal{R}_1, \mathcal{G}_1) \preceq_{\alpha; \zeta \times \gamma} (\text{skip}; C_2, \mathcal{R}_2, \mathcal{G}_2)}$$

Plus the variants with **skip** after the code C_1 or C_2 . That is, **skip**s could be arbitrarily introduced and eliminated.

Common Branch

$$\frac{\forall \sigma_1, \sigma_2. (\sigma_1, \sigma_2) \in \zeta \implies B \sigma_2 \neq \perp}{(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta_1 \times \gamma} (C_1, \mathcal{R}', \mathcal{G}') \quad \zeta_1 = (\zeta \cap (\text{true} \wedge B))}$$

$$\frac{(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta_2 \times \gamma} (C_2, \mathcal{R}', \mathcal{G}') \quad \zeta_2 = (\zeta \cap (\text{true} \wedge \neg B))}{(C, \mathcal{R}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (\text{if } (B) C_1 \text{ else } C_2, \mathcal{R}', \mathcal{G}')}$$

This rule says that, when the if-condition can be evaluated and both branches can be optimized to the same code C , we can transform the whole if-statement to C without introducing new behaviors.

Dead While

$$\frac{\zeta = (\zeta \cap (\text{true} \wedge \neg B)) \quad \zeta \subseteq \alpha \quad \text{Sta}(\zeta, \langle \mathcal{R}_1, \mathcal{R}_2^* \rangle \alpha)}{(\text{skip}, \mathcal{R}_1, \text{ld}) \preceq_{\alpha; \zeta \times \zeta} (\text{while } (B) \{C\}, \mathcal{R}_2, \text{ld})}$$

We can eliminate the loop, if the loop condition is **false** (no matter how the environments update the states) at the loop entry point.

Dead Code Elimination

$$\frac{(\text{skip}, \text{ld}, \text{ld}) \preceq_{\alpha; \zeta \times \gamma} (C, \text{ld}, \mathcal{G}) \quad \text{Sta}(\{\zeta, \gamma\}, \langle \mathcal{R}_1, \mathcal{R}_2^* \rangle \alpha)}{(\text{skip}, \mathcal{R}_1, \text{ld}) \preceq_{\alpha; \zeta \times \gamma} (C, \mathcal{R}_2, \mathcal{G})}$$

Intuitively $(\text{skip}, \text{ld}, \text{ld}) \preceq_{\alpha; \zeta \times \gamma} (C, \text{ld}, \mathcal{G})$ says that the code C can be eliminated in a sequential context where the initial and the final states satisfy ζ and γ respectively. If both ζ and γ are stable *w.r.t.* the interference from the environments \mathcal{R}_1 and \mathcal{R}_2 , then the code C can be eliminated in such a parallel context as well.

Redundancy Introduction

$$\frac{(c, \text{ld}, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (\text{skip}, \text{ld}, \text{ld}) \quad \text{Sta}(\{\zeta, \gamma\}, \langle \mathcal{R}_1, \mathcal{R}_2^* \rangle \alpha)}{(c, \mathcal{R}_1, \mathcal{G}) \preceq_{\alpha; \zeta \times \gamma} (\text{skip}, \mathcal{R}_2, \text{ld})}$$

As we lifted sequential dead code elimination, we can also lift sequential redundant code introduction to the concurrent setting, so long as the pre- and post-conditions are stable *w.r.t.* the environments. Note that here c is a single instruction, because we should consider the interference from the environments at every intermediate state when introducing a sequence of redundant instructions.

5.2 An Example of Invariant Hoisting

With these rules, we can prove the correctness of many traditional compiler optimizations performed on concurrent programs in appropriate contexts. Here we only give a small example of hoisting loop invariants. More examples (*e.g.*, strength reduction and induction variable elimination) can be found in the technical report [20].

| Target Code (C_1) | \Leftarrow | Source Code (C) |
|---|--------------|---|
| <pre>local t; t := x + 1; while(i < n) { i := i + t; }</pre> | \Leftarrow | <pre>local t; while(i < n) { t := x + 1; i := i + t; }</pre> |

When we do not care about the final value of t , it's not difficult to prove that the optimized code C_1 preserves the sequential behaviors of the source C [3]. But in a concurrent setting, safely hoisting the invariant code $t := x + 1$ also requires that the environment should not update x nor t .

$$\mathcal{R} \triangleq \{(\sigma, \sigma') \mid \sigma(x) = \sigma'(x) \wedge \sigma(t) = \sigma'(t)\}.$$

The guarantee of the program can be specified as arbitrary transitions. Since we only care about the values of i , n and x , the invariant relation α can be defined as:

$$\alpha \triangleq \{(\sigma_1, \sigma) \mid \sigma_1(i) = \sigma(i) \wedge \sigma_1(n) = \sigma(n) \wedge \sigma_1(x) = \sigma(x)\}.$$

We do not need special pre- and post-conditions, thus the correctness of the optimization is formalized as follows:

$$(C_1, \mathcal{R}, \text{True}) \preceq_{\alpha; \alpha \times \alpha} (C, \mathcal{R}, \text{True}). \quad (5.1)$$

We could prove (5.1) directly by the RGSim definition and the operational semantics of the code. But below we give a more convenient proof using the optimization rules and the compositionality rules instead. We first prove the following by the Dead-Code-Elimination and Redundancy-Introduction rules:

$$\begin{aligned} (t := x + 1, \mathcal{R}, \text{True}) &\preceq_{\alpha; \alpha \times \gamma} (\text{skip}, \mathcal{R}, \text{True}); \\ (\text{skip}, \mathcal{R}, \text{True}) &\preceq_{\alpha; \gamma \times \eta} (t := x + 1, \mathcal{R}, \text{True}), \end{aligned}$$

where γ and η specify the states at the specific program points:

$$\begin{aligned} \gamma &\triangleq \alpha \cap \{(\sigma_1, \sigma) \mid \sigma_1(t) = \sigma_1(x) + 1\}; \\ \eta &\triangleq \gamma \cap \{(\sigma_1, \sigma) \mid \sigma(t) = \sigma(x) + 1\}. \end{aligned}$$

After adding **skips** to C_1 and C to make them the same “shape”, we can prove the simulation by the compositionality rules SEQ and WHILE. Finally, we remove all the **skips** and conclude (5.1), *i.e.*, the correctness of the optimization in appropriate contexts. Since the relies only prohibit updates of x and t , we can execute C_1 and C concurrently with other threads which update i and n or read x , still ensuring semantics preservation.

6. Proving Atomicity of Concurrent Objects

A concurrent object provides a set of methods, which can be called in parallel by clients as the only way to access the object. RGSim gives us a refinement-based proof method to verify the atomicity

```

ADD(e) :                               RMV(e) :
0  atom {                               0  atom {
  S := S ∪ {e};                          S := S - {e};
}                                           }

```

(a) An Abstract Set

```

add(e) :                                rmv(e) :
  local x,y,z,u;                          local x,y,z,v;
0  <x := Head;>                            0  <x := Head;>
1  lock(x);                                1  lock(x);
2  <z := x.next;>                          2  <y := x.next;>
3  <u := z.data;>                          3  <v := y.data;>
4  while (u < e) {                        4  while (v < e) {
5    lock(z);                              5    lock(y);
6    unlock(x);                            6    unlock(x);
7    x := z;                               7    x := y;
8    <z := x.next;>                        8    <y := x.next;>
9    <u := z.data;>                        9    <v := y.data;>
}                                           }
10 if (u != e) {                          10 if (v = e) {
11   y := new();                            11   lock(y);
12   y.lock := 0;                          12   <z := y.next;>
13   y.data := e;                          13   <x.next := z;>
14   y.next := z;                          14   unlock(x);
15   <x.next := y;>                        15   free(y);
}                                           } else {
16   unlock(x);                            16   unlock(x);
}                                           }

```

(b) The Lock-Coupling List-Based Set

Figure 8. The Set Object

of implementations of the object: we can define abstract atomic operations in a high-level language as specifications, and prove the concrete fine-grained implementations refine the corresponding atomic operations when executed in appropriate environments. For instance, in Figure 8(a) we define two atomic set operations, ADD(e) and RMV(e). Figure 8(b) gives a concrete implementation of the set object using a lock-coupling list. Partial correctness and atomicity of the algorithm has been verified before [28, 29]. Here we show that its atomicity can also be verified using our RGSim by proving the low-level methods refine the corresponding abstract operations. We will discuss the key difference between the previous proofs and ours in Section 8.

We first take the generic languages in Figure 2, and instantiate the high-level program states below.

$$\begin{aligned} (HMem) \quad M_s, M_l &\in (Loc \cup PVar) \rightarrow HVal \\ (HThrds) \quad \Pi &\in ThrID \rightarrow HMem \\ (HState) \quad \Sigma &\in HThrds \times HMem \end{aligned}$$

The state consists of shared memory M_s (where the object resides) and a thread pool Π , which is a mapping from thread identifiers ($t \in ThrID$) to their memory M_l . The low-level state σ is defined similarly. We use m_s , m_l and π to represent the low-level shared memory, thread-local memory and the thread pool respectively.

To allow ownership transfer between the shared memory and thread-local memory, we use **atom**{ \mathbb{C} } \mathbb{A} (or $\langle C \rangle_{\mathcal{A}}$ at the low level) to convert the shared memory to local and then execute \mathbb{C} (or C) atomically. Following RGSep [29], an abstract transition $\mathbb{A} \in \mathcal{P}(HMem \times HMem)$ (or $\mathcal{A} \in \mathcal{P}(LMem \times LMem)$) is used to specify the effects of the atomic operation over the shared memory, which allows us to split the resulting state back into shared and local when we exit the atomic block. The atomic blocks are instantiations of the generic primitive operations c (or c) in Figure 2. Their semantics is shown in the technical report [20]. We omit the annotations \mathbb{A} and

| | |
|----------------------------------|--|
| $m_s \models \text{list}(x, A)$ | $\triangleq (m_s = \phi \wedge x = \text{null} \wedge A = \epsilon) \vee (\exists m'_s. \exists v. \exists y. \exists A'. m_s = m'_s \uplus \{x \rightsquigarrow (-, v, y)\} \wedge A = v :: A' \wedge m'_s \models \text{list}(y, A')$ |
| $\text{shared_map}(m_s, M_s)$ | $\triangleq \exists m'_s. \exists A. \exists x. m_s = m'_s \uplus \{\text{Head} \rightsquigarrow x\} \wedge (m'_s \models \text{list}(x, \text{MIN_VAL} :: A :: \text{MAX_VAL})) \wedge \text{sorted}(A) \wedge (\text{elems}(A) = M_s(\mathcal{S}))$ |
| $\text{local_map}(m_l, M_l)$ | $\triangleq m_l(\mathbf{e}) = M_l(\mathbf{e}) \wedge \exists m'_l. m_l = m'_l \uplus \{x \rightsquigarrow -, y \rightsquigarrow -, z \rightsquigarrow -, u \rightsquigarrow -, v \rightsquigarrow -\}$ |
| α | $\triangleq \{((\pi, m_s), (\Pi, M_s)) \mid \text{shared_map}(m_s, M_s) \wedge \forall t \in \text{dom}(\Pi). \text{local_map}(\pi(t), \Pi(t))\}$ |
| $\mathcal{G}_{\text{lock}}(t)$ | $\triangleq \{((\pi, m_s), (\pi, m'_s)) \mid \exists x, y, v. m_s(x) = (0, v, y) \wedge m'_s = m_s \{x \rightsquigarrow (t, v, y)\}\}$ |
| $\mathcal{G}_{\text{unlock}}(t)$ | $\triangleq \{((\pi, m_s), (\pi, m'_s)) \mid \exists x, y, v. m_s(x) = (t, v, y) \wedge m'_s = m_s \{x \rightsquigarrow (0, v, y)\}\}$ |
| $\mathcal{G}_{\text{add}}(t)$ | $\triangleq \{((\pi \uplus \{t \rightsquigarrow m_l\}, m_s), (\pi \uplus \{t \rightsquigarrow m'_l\}, m'_s)) \mid \exists x, y, z, u, v, w. m_s(x) = (t, u, z) \wedge m_s(z) = (-, w, -) \wedge m'_s = m_s \{x \rightsquigarrow (t, u, y)\} \uplus \{y \rightsquigarrow (0, v, z)\} \wedge (m'_l \uplus \{y \rightsquigarrow (0, v, z)\} = m_l) \wedge u < v < w\}$ |
| $\mathcal{G}_{\text{rmv}}(t)$ | $\triangleq \{((\pi \uplus \{t \rightsquigarrow m_l\}, m_s), (\pi \uplus \{t \rightsquigarrow m'_l\}, m'_s)) \mid \exists x, y, z, u, v. m_s(x) = (t, u, y) \wedge m_s(y) = (t, v, z) \wedge m'_s \uplus \{y \rightsquigarrow (t, v, z)\} = m_s \{x \rightsquigarrow (t, u, z)\} \wedge m'_l = m_l \uplus \{y \rightsquigarrow (t, v, z)\} \wedge v < \text{MAX_VAL}\}$ |
| $\mathcal{G}_{\text{local}}(t)$ | $\triangleq \{((\pi \uplus \{t \rightsquigarrow m_l\}, m_s), (\pi \uplus \{t \rightsquigarrow m'_l\}, m_s)) \mid \pi \in (\text{ThrdID} \rightarrow \text{LMem}) \wedge m_l, m'_l, m_s \in \text{LMem}\}$ |
| $\mathcal{G}(t)$ | $\triangleq \mathcal{G}_{\text{lock}}(t) \cup \mathcal{G}_{\text{unlock}}(t) \cup \mathcal{G}_{\text{add}}(t) \cup \mathcal{G}_{\text{rmv}}(t) \cup \mathcal{G}_{\text{local}}(t)$ |
| $\mathbb{R}(t)$ | $\triangleq \bigcup_{t' \neq t} \mathcal{G}(t')$ |
| $\mathbb{G}_{\text{add}}(t)$ | $\triangleq \{((\Pi \uplus \{t \rightsquigarrow M_l\}, M_s), (\Pi \uplus \{t \rightsquigarrow M'_l\}, M'_s)) \mid \exists e. M'_s = M_s \{\mathcal{S} \rightsquigarrow M_s(\mathcal{S}) \cup \{e\}\}\}$ |
| $\mathbb{G}_{\text{rmv}}(t)$ | $\triangleq \{((\Pi \uplus \{t \rightsquigarrow M_l\}, M_s), (\Pi \uplus \{t \rightsquigarrow M'_l\}, M'_s)) \mid \exists e. M'_s = M_s \{\mathcal{S} \rightsquigarrow M_s(\mathcal{S}) - \{e\}\}\}$ |
| $\mathbb{G}_{\text{local}}(t)$ | $\triangleq \{((\Pi \uplus \{t \rightsquigarrow M_l\}, M_s), (\Pi \uplus \{t \rightsquigarrow M'_l\}, M'_s)) \mid \Pi \in (\text{ThrdID} \rightarrow \text{HMem}) \wedge M_l, M'_l, M_s \in \text{HMem}\}$ |
| $\mathbb{G}(t)$ | $\triangleq \mathbb{G}_{\text{add}}(t) \cup \mathbb{G}_{\text{rmv}}(t) \cup \mathbb{G}_{\text{local}}(t)$ |
| $\mathbb{R}(t)$ | $\triangleq \bigcup_{t' \neq t} \mathbb{G}(t')$ |

Figure 9. Useful Definitions for the Lock-Coupling List

A in Figure 8, which are the same as the corresponding guarantees in Figure 9, as we will explain below.

In Figure 8, the abstract set is implemented by an ordered singly-linked list pointed to by a shared variable `Head`, with two sentinel nodes at the two ends of the list containing the values `MIN_VAL` and `MAX_VAL` respectively. Each list node is associated with a lock. Traversing the list uses “hand-over-hand” locking: the lock on one node is not released until its successor is locked. `add(e)` inserts a new node with value e in the appropriate position while holding the lock of its predecessor. `rmv(e)` redirects the predecessor’s pointer while both the node to be removed and its predecessor are locked.

We define the α relation, the guarantees and the relies in Figure 9. The predicate $m_s \models \text{list}(x, A)$ represents a singly-linked list in the shared memory m_s at the location x , whose values form the sequence A . Then the mapping `shared_map` between the low-level and the high-level shared memory is defined by only concerning about the value sequence on the list: the concrete list should be sorted and its elements constitute the abstract set. For a thread t ’s local memory of the two levels, we require that the values of e are the same and enough local space is provided for `add(e)` and `rmv(e)`, as defined in the mapping `local_map`. Then α relates the shared memory by `shared_map` and the local memory of each thread t by `local_map`.

The atomic actions of the algorithm are specified by $\mathcal{G}_{\text{lock}}$, $\mathcal{G}_{\text{unlock}}$, \mathcal{G}_{add} , \mathcal{G}_{rmv} and $\mathcal{G}_{\text{local}}$ respectively, which are all parameterized with a thread identifier t . For example, $\mathcal{G}_{\text{rmv}}(t)$ says that when holding the locks of the node y and its predecessor x , we can transfer the node y from the shared memory to the thread’s local memory. This corresponds to the action performed by the code of line 13 in `rmv(e)`. Every thread t is executed in the environment that any other thread t' can only perform those five actions, as defined in $\mathcal{R}(t)$. Similarly, the high-level $\mathbb{G}(t)$ and $\mathbb{R}(t)$ are defined according to the abstract `ADD(e)` and `RMV(e)`. The relies and guarantees are almost the same as those in the proofs in `RGSep` [28].

We can prove that for any thread t , the following hold:

$$\begin{aligned} (t.\text{add}(e), \mathcal{R}(t), \mathcal{G}(t)) &\preceq_{\alpha; \alpha \times \alpha} (t.\text{ADD}(e), \mathbb{R}(t), \mathbb{G}(t)); \\ (t.\text{rmv}(e), \mathcal{R}(t), \mathcal{G}(t)) &\preceq_{\alpha; \alpha \times \alpha} (t.\text{RMV}(e), \mathbb{R}(t), \mathbb{G}(t)). \end{aligned}$$

We give the detailed proofs in the technical report [20], which are done operationally based on the definition of `RGSim`.

By the compositionality and the soundness of `RGSim`, we know that the fine-grained operations (under the parallel environment \mathcal{R}) are simulated by the corresponding atomic operations (under the high-level environment \mathbb{R}), while \mathcal{R} and \mathbb{R} say all accesses to the set must be done through the add and remove operations. This gives us the atomicity of the concurrent implementation of the set object.

More examples. In the technical report [20], we also show the use of `RGSim` to prove the atomicity of other fine-grained algorithms, including the non-blocking concurrent counter [27], Treiber’s stack algorithm [26], and a concurrent GCD algorithm (calculating greatest common divisors).

7. Verifying Concurrent Garbage Collectors

In this section, we explain in detail how to reduce the problem of verifying concurrent garbage collectors to transformation verification, and use `RGSim` to develop a general GC verification framework. We apply the framework to prove the correctness of the Boehm *et al.* concurrent GC algorithm [7].

7.1 Correctness of Concurrent GCs

A concurrent GC is executed by a dedicate thread and performs the collection work in parallel with user threads (mutators), which access the shared heap via read, write and allocation operations. To ensure that the GC and the mutators share a coherent view of the heap, the heap operations from mutators may be instrumented with extra operations, which provide an interaction mechanism to allow arbitrary mutators to cooperate with the GC. These instrumented heap operations are called barriers (*e.g.*, read barriers, write barriers and allocation barriers).

The GC thread and the barriers constitute a concurrent garbage collecting system, which provides a higher-level user-friendly programming model for garbage-collected languages (*e.g.*, Java). In this high-level model, programmers feel they access the heap using regular memory operations, and are freed from manually disposing objects that are no longer in use. They do not need to consider the implementation details of the GC and the existence of barriers.

We could verify the GC system by using a Hoare-style logic to prove that the GC thread and the barriers satisfy their specifications. However, we say this is an indirect approach because it is unclear if the specified correct behaviors would indeed make the mutators happy and generate the abstract view for high-level programmers. Usually this part is examined by experts and then trusted.

Here we propose a more direct approach. We view a concurrent garbage collecting system as a transformation \mathbf{T} from a high-level garbage-collected language to a low-level language. A standard atomic memory operation at the source level is transformed into the corresponding barrier code at the target level. In the source level, we assume there is an *abstract GC thread* that magically turns unreachable objects into reusable memory. The abstract collector *AbsGC* is transformed into the concrete GC code C_{gc} running concurrently with the target mutators. That is,

$$\mathbf{T}(t_{gc}.AbsGC \parallel t_1.C_1 \parallel \dots \parallel t_n.C_n) \triangleq t_{gc}.C_{gc} \parallel t_1.\mathbf{T}(C_1) \parallel \dots \parallel t_n.\mathbf{T}(C_n),$$

where $\mathbf{T}(C)$ simply translates some memory access instructions in C into the corresponding barriers, and leaves the rest unchanged.

Then we reduce the correctness of the concurrent garbage collecting system to $\text{Correct}(\mathbf{T})$, saying that any mutator program will not have unexpected behaviors when executed using this system.

7.2 A General Framework

The compositionality of RGSim allows us to develop a general framework to prove $\text{Correct}(\mathbf{T})$, which cannot be done by monolithic proof methods. By the parallel compositionality of RGSim (the PAR rule in Figure 7), we can decompose the refinement proofs into proofs for the GC thread and each mutator thread.

Verifying the GC. The semantics of the abstract GC thread can be defined by a binary state predicate AbsGCStep :

$$\frac{(\Sigma, \Sigma') \in \text{AbsGCStep}}{(t_{gc}.AbsGC, \Sigma) \longrightarrow (t_{gc}.AbsGC, \Sigma')}$$

That is, the abstract GC thread always makes AbsGCStep to change the high-level state. We can choose different AbsGCStep for different GCs, but usually AbsGCStep guarantees not modifying reachable objects in the heap.

Thus for the GC thread, we need to show that C_{gc} is simulated by *AbsGC* when executed in their environments. This can be reduced to unary Rely-Guarantee reasoning about C_{gc} by proving $\mathcal{R}_{gc}; \mathcal{G}_{gc} \vdash \{p_{gc}\}C_{gc}\{q_{gc}\}$ in a standard Rely-Guarantee logic with proper \mathcal{R}_{gc} , \mathcal{G}_{gc} , p_{gc} and q_{gc} , as long as \mathcal{G}_{gc} is a concrete representation of AbsGCStep . The judgment says given an initial state satisfying the precondition p_{gc} , if the environment's behaviors satisfy \mathcal{R}_{gc} , then each step of C_{gc} satisfies \mathcal{G}_{gc} , and the postcondition q_{gc} holds at the end if C_{gc} terminates. In general, the collector never terminates, thus we can let q_{gc} be **false**. \mathcal{G}_{gc} and p_{gc} should be provided by the verifier, where p_{gc} needs to be general enough that can be satisfied by any possible low-level initial state. \mathcal{R}_{gc} encodes the possible behaviors of mutators, which can be derived, as we will show below.

Verifying mutators. For the mutator thread, since \mathbf{T} is syntax-directed on \mathbb{C} , we can reduce the refinement problem for arbitrary mutators to the refinement on each primitive instruction only, by the compositionality of RGSim. The proof needs proper rely/guarantee conditions. Let $\mathbb{G}(t.c)$ and $\mathcal{G}(t.\mathbf{T}(c))$ denote the guarantees of the source instruction c and the target code $\mathbf{T}(c)$ respectively. Then we can define the general guarantees for a mutator thread t :

$$\begin{aligned} \mathbb{G}(t) &\triangleq \bigcup_c \mathbb{G}(t.c); \\ \mathcal{G}(t) &\triangleq \bigcup_c \mathcal{G}(t.\mathbf{T}(c)). \end{aligned} \quad (7.1)$$

Its relies should include all the possible guarantees made by other threads, and the GC's abstract and concrete behaviors respectively:

$$\begin{aligned} \mathbb{R}(t) &\triangleq \text{AbsGCStep} \cup \bigcup_{t' \neq t} \mathbb{G}(t'); \\ \mathcal{R}(t) &\triangleq \mathcal{G}_{gc} \cup \bigcup_{t' \neq t} \mathcal{G}(t'). \end{aligned} \quad (7.2)$$

The \mathcal{R}_{gc} used to verify the GC code can now be defined below:

$$\mathcal{R}_{gc} \triangleq \bigcup_t \mathcal{G}(t). \quad (7.3)$$

The refinement proof also needs definitions of binary α , ζ and γ relations. The invariant α relates the low-level and the high-level states and needs to be preserved by each low-level step. In general, a high-level state Σ can be mapped to a low-level state σ by giving a concrete local store for the GC thread, adding additional structures in the heap (to record information for collection), renaming heap cells (for copying GCs), etc.. For each mutator thread t , the relations $\zeta(t)$ and $\gamma(t)$ need to hold at the beginning and the end of each basic transformation unit (every high-level primitive instruction in this case) respectively. We let $\gamma(t)$ be the same as $\zeta(t)$ to support sequential compositions. We require $\text{InitRel}_{\mathbf{T}}(\zeta(t))$ (see Figure 6), i.e., $\zeta(t)$ holds over the initial states. In addition, the target and the source boolean expressions should be evaluated to the same value under related states, as required in the IF and WHILE rules in Figure 7.

$$\text{Good}_{\mathbf{T}}(\zeta(t)) \triangleq \text{InitRel}_{\mathbf{T}}(\zeta(t)) \wedge \forall \mathbb{B}. \zeta(t) \subseteq (\mathbf{T}(\mathbb{B}) \Leftrightarrow \mathbb{B}) \quad (7.4)$$

Theorem 8 (Verifying Concurrent Garbage Collecting Systems).

If there exist \mathcal{R}_{gc} , \mathcal{G}_{gc} , p_{gc} , $\mathcal{R}(t)$, $\mathbb{R}(t)$, $\zeta(t)$ and α such that (7.1), (7.2), (7.3), (7.4) and the following hold:

1. (Verification of the GC code)
 $\mathcal{R}_{gc}; \mathcal{G}_{gc} \vdash \{p_{gc}\}C_{gc}\{\text{false}\};$
2. (Correctness of \mathbf{T} on mutator instructions)
 $\forall c. (t.\mathbf{T}(c), \mathcal{R}(t), \mathcal{G}(t)) \preceq_{\alpha; \zeta(t) \times \zeta(t)} (t.c, \mathbb{R}(t), \mathbb{G}(t));$
3. (Side conditions)
 $\mathcal{G}_{gc} \circ \alpha^{-1} \subseteq \alpha^{-1} \circ (\text{AbsGCStep})^*;$
 $\forall \sigma, \Sigma. \sigma = \mathbf{T}(\Sigma) \implies p_{gc} \sigma;$

then $\text{Correct}(\mathbf{T})$.

That is, to verify a concurrent garbage collecting system, we need to do the following:

- Define the α and $\zeta(t)$ relations, and prove the correctness of \mathbf{T} on high-level primitive instructions. Since \mathbf{T} preserves the syntax on most instructions, it's often immediate to prove the target instructions are simulated by their sources. But for instructions that are transformed to barriers, we need to verify the barriers that they implement both the source instructions (by RGSim) and the interaction mechanism (shown in their guarantees).
- Find some proper \mathcal{G}_{gc} and p_{gc} , and verify the GC code by R-G reasoning. We require the GC's guarantee \mathcal{G}_{gc} should not contain more behaviors than AbsGCStep (the first side condition), and C_{gc} can start its execution from any state σ transformed from a high-level one (the second side condition).

7.3 Application: Boehm *et al.* Concurrent GC Algorithm

We illustrate the applications of the framework (Theorem 8) by proving the correctness of a mostly-concurrent mark-sweep garbage collector proposed by Boehm *et al.* [7]. Variants of the algorithm have been used in practice (e.g., by IBM [2]). Due to the space limit, we only describe the proof sketch here. Details are presented in the technical report [20].

Overview of the GC algorithm. The top-level code of the GC thread is shown in Figure 10. In each collection cycle, after an

```

{wfstate}
0 Collection() {
1   local mstk: Seq(Int);
   Loop Invariant: {wfstate * (ownnp(mstk) ∧ mstk = ε)}
2   while (true) {
3     Initialize();
     {(wfstate ∧ reach_inv) * (ownnp(mstk) ∧ mstk = ε)}
4     Trace();
     {(wfstate ∧ reach_inv) * (ownnp(mstk) ∧ mstk = ε)}
5     CleanCard();
     {(wfstate ∧ reach_inv) * (ownnp(mstk) ∧ mstk = ε)}
     atomic{
6       ScanRoot();
       {∃X.(wfstate ∧ reach_rtnw_stk(X) ∧ stk_black(X))
        * (ownnp(mstk) ∧ mstk = X)}
7       CleanCard();
     }
     {(wfstate ∧ reach_black) * (ownnp(mstk) ∧ mstk = ε)}
8     Sweep();
   }
}
{false}

```

Figure 10. Outline of the GC Code and Proof Sketch

```

update(x.id, E) { // id ∈ {pt1, ..., ptm}
  atomic{ x.id := E; aux := x; }
  atomic{ x.dirty := 1; aux := 0; }
}

```

Figure 11. The Write Barrier for Boehm *et al.* GC

initialization process, the GC enters the concurrent mark-phase (line 4) and traces the objects reachable from the *roots* (i.e., the mutators' local pointer variables that may contain references to the heap objects). A *mark stack* (*mstk*) is used to do a depth-first tracing. During the tracing, the connectivity between objects might be changed by the mutators, thus a write barrier is required to notify the collector of those modified objects by dirtying the objects' tags (called cards). When the tracing is done, the GC suspends all the mutators and re-traces from the dirty objects that have been marked (called *card-cleaning*, line 6 and 7). The stop-the-world phase is implemented by **atomic**{*C*}. Finally, all the reachable objects are ensured marked and the GC performs the concurrent sweep-phase (line 8), in which unmarked objects are reclaimed. Usually in practice, there is also a concurrent card-cleaning phase (line 5) before the stop-the-world card-cleaning to reduce the pause time. The full GC code C_{gc} is given in the technical report [20].

The write barrier is shown in Figure 11, where the *dirty* field is set after modifying the object's pointer field. Here we use a write-only auxiliary variable *aux* for each mutator thread to record the current object that the mutator is updating. We add *aux* only for the purpose of verification, so that we can easily specify the fine-grained property of the write barrier in the guarantees that immediately after updating the pointer field, the thread would do nothing else except setting the corresponding *dirty* field. The GC does not use read barriers nor allocation barriers.

We first present the high-level and low-level program state models in Figure 12. The behaviors of the high-level abstract GC thread are defined as follows:

$$\text{AbsGCStep} \triangleq \{(\Pi, H), (\Pi, H') \mid \forall l. \text{reachable}(l)(\Pi, H) \implies H(l) = H'(l)\},$$

saying that, the mutator stores and the reachable objects in the heap are remained unmodified. Here $\text{reachable}(l)(\Pi, H)$ means the object at the location l is reachable in H from the roots in Π .

$$\begin{aligned} (HStore) S \in PVar \rightarrow HVal & \quad (HHeap) H \in Loc \rightarrow HObj \\ (HThrs) \Pi \in MutID \rightarrow HStore & \quad (HState) \Sigma \in HThrs \times HHeap \\ (LStore) s \in PVar \rightarrow LVal \times \{0, 1\} & \quad (LHeap) h \in [1..M] \rightarrow LObj \\ (LThrs) \pi \in ThrID \rightarrow LStore & \quad (LState) \sigma \in LThrs \times LHeap \end{aligned}$$

Figure 12. High-Level and Low-Level State Models

The transformation. The transformation \mathbf{T} is defined as follows. For *code*, the high-level abstract GC thread is transformed to the GC thread shown in Figure 10. Each instruction $x.id := \mathbb{E}$ in mutators is transformed to the write barrier, where *id* is a pointer field of x . Other instructions and the program structures of mutators are unchanged.

The following transformations are made over initial states.

- First we require the high-level initial state to be *well-formed*:

$$\text{wfstate}(\Pi, H) \triangleq \forall l. \text{reachable}(l)(\Pi, H) \implies l \in \text{dom}(H).$$
That is, reachable locations cannot be dangling pointers.
- High-level locations are transformed to integers by a bijective function $\text{Loc2Int} : Loc \leftrightarrow [0..M]$ satisfying $\text{Loc2Int}(\text{nil}) = 0$.
- Variables are transformed to the low level using an extra bit to preserve the high-level type information (0 for non-pointers and 1 for pointers).
- High-level objects are transformed to the low level by adding the *color* and *dirty* fields with initial values *WHITE* and 0 respectively. Other addresses in the low-level heap domain $[1..M]$ are filled out using unallocated objects whose *colors* are *BLUE* and all the other fields are initialized by 0. Here we use *BLACK* and *WHITE* for marked and unmarked objects respectively, and *BLUE* for unallocated memory.
- The concrete GC thread is given an initial store.

The formal definition of \mathbf{T} is included in the technical report [20].

To prove $\text{Correct}(\mathbf{T})$ in our framework, we apply Theorem 8, prove the refinement between low-level and high-level mutators, and verify the GC code using a unary Rely-Guarantee-based logic.

Refinement proofs for mutator instructions. We first define the α and $\zeta(t)$ relations.

$$\begin{aligned} \alpha \triangleq \{ & ((\pi \uplus \{t_{gc} \rightsquigarrow -\}, h), (\Pi, H)) \mid \\ & \forall t \in \text{dom}(\Pi). \text{store_map}(\pi(t), \Pi(t)) \\ & \wedge \text{heap_map}(h, H) \wedge \text{wfstate}(\Pi, H) \}. \end{aligned}$$

In α , the relation between low-level and high-level stores and heaps are enforced by *store_map* and *heap_map* respectively. Their definitions reflect the state transformations we describe above, ignoring the values of those high-level-invisible structures. It also requires the well-formedness of high-level states.

For each mutator thread t , the $\zeta(t)$ relation enforced at the beginning and the end of each transformation unit (each high-level instruction) is stronger than α . It requires that the value of the auxiliary variable *aux* (see Figure 11) be a null pointer (0^p):

$$\zeta(t) \triangleq \alpha \cap \{((\pi, h), (\Pi, H)) \mid \pi(t)(\text{aux}) = 0^p\}.$$

The refinement between the write barrier at the low level and the pointer update instruction at the high level is formulated as:

$$\begin{aligned} (t.\text{update}(x.id, E), \mathcal{R}(t), \mathcal{G}_{\text{write_barrier}}^t) & \preceq_{\alpha; \zeta(t) \times \zeta(t)} \\ (t.(x.id := \mathbb{E}), \mathbb{R}(t), \mathbb{G}_{\text{write_pt}}^t), & \end{aligned}$$

where $\mathcal{G}_{\text{write_barrier}}^t$ and $\mathbb{G}_{\text{write_pt}}^t$ are the guarantees of the two-step write barrier and the high-level atomic write operation respectively.

Since the transformation of other high-level instructions is identity, the corresponding refinement proofs are simple.

Rely-Guarantee reasoning about the GC code. The unary program logic we use to verify the GC thread is a standard Rely-Guarantee logic adapted to the target language. We describe states using separation logic assertions, as shown below:

$$p, q ::= B \mid \mathbf{t.own}_p(x) \mid \mathbf{t.own}_{np}(x) \mid E_1.id \mapsto E_2 \mid p * q \mid \dots$$

Following Parkinson *et al.* [23], we treat program variables as resource and use $\mathbf{t.own}_p(x)$ and $\mathbf{t.own}_{np}(x)$ for the thread t 's ownerships of pointers and non-pointers respectively. We omit the thread identifiers if these predicates hold for the current thread.

We first give the precondition and the guarantee of the GC. The GC starts its executions from a low-level *well-formed* state, *i.e.*, $p_{gc} \triangleq \mathbf{wfstate}$. Just corresponding to the high-level *wfstate* definition, the low-level *wfstate* predicate says none of the reachable objects are BLUE. We define \mathcal{G}_{gc} as follows:

$$\begin{aligned} \mathcal{G}_{gc} \triangleq & \{((\pi \uplus \{\mathbf{t}_{gc} \rightsquigarrow s\}, h), (\pi \uplus \{\mathbf{t}_{gc} \rightsquigarrow s'\}, h')) \\ & \mid \forall n. \mathbf{reachable}(n)(\pi, h) \\ & \implies [h(n)] = [h'(n)] \\ & \wedge h(n).\mathbf{color} \neq \mathbf{BLUE} \wedge h'(n).\mathbf{color} \neq \mathbf{BLUE}\}. \end{aligned}$$

The GC guarantees not modifying the mutator stores. For any mutator-reachable object, the GC does not update its fields coming from the high-level mutator, nor does it reclaim the object. Here $[-]$ lifts a low-level object to a new one that contains mutator data only.

The proof sketch is given in Figure 10. One of the key invariants used in the proof is $\mathbf{reach_inv}$, which says any WHITE reachable object can either be traced from a root object in a path on which every object is WHITE, or be reachable from a BLACK object whose pointer field was updated and \mathbf{dirty} bit was set to 1. Since the proof is done in the unary logic, the details here are orthogonal to our simulation-based proof (but it is RGSim that allows us to derive Theorem 8, which then links proofs in the unary logic with relational proofs). We give the complete proofs in the technical report [20].

8. Related Work and Conclusion

There is a large body of work on refinements and verification of program transformations. Here we only focus on the work most closely related to the typical applications discussed in this paper.

Verifying compilation and optimizations of concurrent programs. Compiler verification for concurrent programming languages can date back to work by Wand [31] and Gladstein *et al.* [14], which is about functional languages using message-passing mechanisms. Recently, Lochbihler [21] presents a verified compiler for Java threads and prove semantics preservation by a weak bisimulation. He views every heap update as an observable move, thus does not allow the target and the source to have different granularities of atomic updates. To achieve parallel compositionality, he requires the relation to be preserved by any transitions of shared states, *i.e.*, the environments are assumed arbitrary. As we explained in Section 2.2, this is a too strong requirement in general for many transformations, including the examples in this paper.

Burckhardt *et al.* [9] present a proof method for verifying concurrent program transformations on relaxed memory models. The method relies on a compositional trace-based denotational semantics, where the values of shared variables are always considered arbitrary at any program point. In other words, they also assume arbitrary environments.

Following Leroy's CompCert project [19], Ševčík *et al.* [25] verify compilation from a C-like concurrent language to x86 by

simulations. They focus on correctness of a particular compiler, and there are two phases in their compiler whose proofs are not compositional. Here we provide a general, compiler-independent, compositional proof technique to verify concurrent transformations.

We apply RGSim to justify concurrent optimizations, following Benton [3] who presents a declarative set of rules for sequential optimizations. Also the proof rules of RGSim for sequential compositions, conditional statements and loops coincide with those in relational Hoare logic [3] and relational separation logic [32].

Proving linearizability or atomicity of concurrent objects. Filipović *et al.* [13] show linearizability can be characterized in terms of an observational refinement, where the latter is defined similarly to our $\mathbf{Correct}(T)$. There is no proof method given to verify the linearizability of fine-grained object implementations.

Turon and Wand [27] propose a refinement-based proof method to verify concurrent objects. They first propose a simple refinement based on Brookes' fully abstract trace semantics [8], which is compositional but cannot handle complex algorithms (as discussed in Section 2.2). Their fenced refinement then uses rely conditions to filter out illegal environment transitions. The basic idea is similar to ours, and the refinement can also be used to verify Treiber's stack algorithm. However, it is "not a congruence for parallel composition". In their settings, both the concrete (fine-grained) and the abstract (atomic) versions of object operations need to be expressed in the same language. They also require that the fine-grained implementation should have only one update action over the shared state to correspond to the high-level atomic operation. These requirements and the lack of parallel compositionality limit the applicability of their method. It is unclear if the method can be used for general verification of transformations, such as concurrent GCs.

Elmas *et al.* [12] prove linearizability by incrementally rewriting the fine-grained implementation to the atomic abstract specification. Their behavioral simulation used to characterize linearizability is an event-trace subset relation with requirements on the orders of method invocations and returns. Their rules heavily rely on movers (*i.e.*, operations that can commute over any operation of other threads) and always rewrite programs to instructions, thus are designed specifically for atomicity verification.

In his thesis [28], Vafeiadis proves linearizability of concurrent objects in RGSep logic by introducing abstract objects and abstract atomic operations as auxiliary variables and code. The refinement between the concrete implementation and the abstract operation is implicitly embodied in the unary verification process, but is not spelled out formally in the meta-theory (*e.g.*, the soundness).

Verifying concurrent GCs. Vechev *et al.* [30] define transformations to generate concurrent GCs from an abstract collector. Afterwards, Pavlovic *et al.* [24] present refinements to derive concrete concurrent GCs from specifications. These methods focus on describing the behaviors of variants (or instantiations) of a correct abstract collector (or a specification) in a single framework, assuming all the mutator operations are atomic. By comparison, we provide a general correctness notion and a proof method for verifying concurrent GCs and the interactions with mutators (where the barriers could be fine-grained). Furthermore, the correctness of their transformations or refinements is expressed in a GC-oriented way (*e.g.*, the target GC should mark no less objects than the source), which cannot be used to justify other transformations.

Kapoor *et al.* [18] verify Dijkstra's GC using concurrent separation logic. To validate the GC specifications, they also verify a representative mutator in the same system. In contrast, we reduce the problem of verifying a concurrent GC to verifying a transformation, ensuring semantics preservation for *all* mutators. Our GC

verification framework is inspired by McCreight *et al.* [22], who propose a framework for separate verification of stop-the-world and incremental GCs and their mutators, but their framework does not handle concurrency.

Conclusion and future work. We propose RGSim to verify concurrent program transformations. By describing explicitly the interference with environments, RGSim is compositional, and can support many widely-used transformations. We have applied RGSim to reason about optimizations, prove atomicity of fine-grained concurrent algorithms and verify concurrent garbage collectors. In the future, we would like to further test its applicability with more applications, such as verifying STM implementations and compilers. It is also interesting to explore the possibility of building tools to automate the verification process.

Acknowledgments

We would like to thank Matthew Parkinson and anonymous referees for their suggestions and comments on earlier versions of this paper. This work is supported in part by grants from National Natural Science Foundation of China (NSFC) under Grant No. 60928004, 61073040 and 61103023. Xinyu Feng is also supported in part by NSFC under Grant No. 90818019, by Program for New Century Excellent Talents in Universities (NCET), and by the Fundamental Research Funds for the Central Universities.

References

- [1] M. Abadi and G. Plotkin. A model of cooperative threads. In *Proc. 36th ACM Symp. on Principles of Prog. Lang. (POPL'09)*, pages 29–40. ACM Press, January 2009.
- [2] K. Barabash, O. Ben-Yitzhak, I. Goft, E. K. Kolodner, V. Leikehman, Y. Ossia, A. Owshanko, and E. Petrank. A parallel, incremental, mostly concurrent garbage collector for servers. *ACM Trans. Program. Lang. Syst.*, 27(6):1097–1146, 2005.
- [3] N. Benton. Simple relational correctness proofs for static analyses and program transformations. In *Proc. 31th ACM Symp. on Principles of Prog. Lang. (POPL'04)*, pages 14–25. ACM Press, January 2004.
- [4] N. Benton and C.-K. Hur. Biorthogonality, step-indexing and compiler correctness. In *Proc. 14th ACM Int'l Conf. on Functional Prog. (ICFP'09)*, pages 97–108. ACM Press, September 2009.
- [5] H.-J. Boehm. Threads cannot be implemented as a library. In *Proc. 2005 ACM Conf. on Prog. Lang. Design and Impl. (PLDI'05)*, pages 261–268. ACM Press, June 2005.
- [6] H.-J. Boehm and S. V. Adve. Foundations of the C++ concurrency memory model. In *Proc. 2008 ACM Conf. on Prog. Lang. Design and Impl. (PLDI'08)*, pages 68–78. ACM Press, June 2008.
- [7] H.-J. Boehm, A. J. Demers, and S. Shenker. Mostly parallel garbage collection. In *Proc. 1991 ACM Conf. on Prog. Lang. Design and Impl. (PLDI'91)*, pages 157–164. ACM Press, June 1991.
- [8] S. D. Brookes. Full abstraction for a shared-variable parallel language. *Inf. Comput.*, 127(2):145–163, 1996.
- [9] S. Burckhardt, M. Musuvathi, and V. Singh. Verifying local transformations on relaxed memory models. In *Proc. 19th Int'l Conf. on Compiler Construction (CC'10)*, volume 6011 of *Lecture Notes in Computer Science*, pages 104–123. Springer, March 2010.
- [10] Coq Development Team. The Coq proof assistant reference manual. The Coq release v8.3, October 2010.
- [11] D. Dice, O. Shalev, and N. Shavit. Transactional locking II. In *Proc. 20th Int'l Symp. on Distributed Computing (DISC'06)*, volume 4167 of *Lecture Notes in Computer Science*, pages 194–208. Springer, September 2006.
- [12] T. Elmas, S. Qadeer, A. Sezgin, O. Subasi, and S. Tasiran. Simplifying linearizability proofs with reduction and abstraction. In *Proc. 16th Int'l Conf. on Tools and Algorithms for the Construction and Analysis of Systems (TACAS'10)*, volume 6015 of *Lecture Notes in Computer Science*, pages 296–311. Springer, March 2010.
- [13] I. Filipović, P. O'Hearn, N. Rinetzky, and H. Yang. Abstraction for concurrent objects. In *Proc. 18th European Symp. on Prog. (ESOP'09)*, volume 5502 of *Lecture Notes in Computer Science*, pages 252–266. Springer, March 2009.
- [14] D. S. Gladstein and M. Wand. Compiler correctness for concurrent languages. In *Proc. 1st Int'l Conf. on Coordination Languages and Models (COORDINATION'96)*, volume 1061 of *Lecture Notes in Computer Science*, pages 231–248. Springer, April 1996.
- [15] M. Herlihy and N. Shavit. *The Art of Multiprocessor Programming*. Morgan Kaufmann, April 2008.
- [16] C.-K. Hur and D. Dreyer. A Kripke logical relation between ML and assembly. In *Proc. 38th ACM Symp. on Principles of Prog. Lang. (POPL'11)*, pages 133–146. ACM Press, January 2011.
- [17] C. B. Jones. Tentative steps toward a development method for interfering programs. *ACM Trans. Program. Lang. Syst.*, 5(4):596–619, 1983.
- [18] K. Kapoor, K. Lodaya, and U. Reddy. Fine-grained concurrency with separation logic. *J. Philosophical Logic*, 40(5):583–632, 2011.
- [19] X. Leroy. A formally verified compiler back-end. *J. Autom. Reason.*, 43(4):363–446, December 2009.
- [20] H. Liang, X. Feng, and M. Fu. A rely-guarantee-based simulation for verifying concurrent program transformations. Technical report (with Coq implementation), University of Science and Technology of China, October 2011. <http://kyhcs.ustcsz.edu.cn/relconcur/rgsim>.
- [21] A. Lochbihler. Verifying a compiler for java threads. In *Proc. 20th European Symp. on Prog. (ESOP'10)*, volume 6012 of *Lecture Notes in Computer Science*, pages 427–447. Springer, March 2010.
- [22] A. McCreight, Z. Shao, C. Lin, and L. Li. A general framework for certifying garbage collectors and their mutators. In *Proc. 2007 ACM Conf. on Prog. Lang. Design and Impl. (PLDI'07)*, pages 468–479. ACM Press, June 2007.
- [23] M. Parkinson, R. Bornat, and C. Calcagno. Variables as resource in Hoare logics. In *Proc. 21th IEEE Symp. on Logic In Computer Science (LICS'06)*, pages 137–146. IEEE Computer Society, August 2006.
- [24] D. Pavlovic, P. Pepper, and D. R. Smith. Formal derivation of concurrent garbage collectors. In *Proc. 10th Int'l Conf. on Mathematics of Program Construction (MPC'10)*, volume 6120 of *Lecture Notes in Computer Science*, pages 353–376. Springer, June 2010.
- [25] J. Ševčík, V. Vafeiadis, F. Z. Nardelli, S. Jagannathan, and P. Sewell. Relaxed-memory concurrency and verified compilation. In *Proc. 38th ACM Symp. on Principles of Prog. Lang. (POPL'11)*, pages 43–54. ACM Press, January 2011.
- [26] R. K. Treiber. System programming: coping with parallelism. Technical Report RJ 5118, IBM Almaden Research Center, 1986.
- [27] A. Turon and M. Wand. A separation logic for refining concurrent objects. In *Proc. 38th ACM Symp. on Principles of Prog. Lang. (POPL'11)*, pages 247–258. ACM Press, January 2011.
- [28] V. Vafeiadis. Modular fine-grained concurrency verification. Technical Report UCAM-CL-TR-726, University of Cambridge, Computer Laboratory, July 2008.
- [29] V. Vafeiadis and M. J. Parkinson. A marriage of rely/guarantee and separation logic. In *Proc. 18th Int'l Conf. on Concurrency Theory (CONCUR'07)*, volume 4703, pages 256–271. Springer, 2007.
- [30] M. T. Vechev, E. Yahav, and D. F. Bacon. Correctness-preserving derivation of concurrent garbage collection algorithms. In *Proc. 2006 ACM Conf. on Prog. Lang. Design and Impl. (PLDI'06)*, pages 341–353. ACM Press, June 2006.
- [31] M. Wand. Compiler correctness for parallel languages. In *Proc. Conf. on Functional Prog. Lang. and Computer Architecture (FPCA'95)*, pages 120–134. ACM Press, June 1995.
- [32] H. Yang. Relational separation logic. *Theoretical Computer Science*, 375(1-3):308–334, 2007.