

# Towards Certified Compositional Compilation for Concurrent Programs

ANONYMOUS AUTHOR(S)

Certified compositional compilation is important for establishing end-to-end guarantees for certified systems consisting of separately compiled modules. In this paper, we propose a language-independent framework consisting of the key semantics components and verification steps that bridge the gap between the compilers for sequential programs and for (race-free) concurrent ones, so that the efforts on certified sequential compilation can be reused for concurrent programs. One of the key contributions of the framework is a novel footprint-preserving compositional simulation as the compilation correctness criterion. With our framework, we have verified the correctness of the CompCert x86 backend (including all the translation phases from Cminor to x86, and two optimization phases) for compositional compilation of race-free concurrent programs.

## 1 INTRODUCTION

Compositional compilation is important for real-world systems, which usually consist of multiple program modules that need to be compiled independently. Correct compilation then needs to guarantee that the target modules can work together and preserve the semantics of the source program as a whole. It requires not only that individual modules be compiled correctly, but also that the expected interaction between modules be preserved at the target.

CompCert [Leroy 2009a], the most well-known certified realistic compiler, establishes the semantics preservation property for compilation of *sequential* Clight programs, but with no explicit support of separate compilation. To support general separate compilation, Stewart et al. [2015] develop Compositional CompCert, which allows the modules to call each other's external functions. Like CompCert, Compositional CompCert only supports sequential programs too. There each module interacts with others only at certain program points specified in the module code (i.e., at external calls only). Also there needs to be one-to-one correspondence between the interaction points in the target program and those in the source. However, in interleaving concurrency, a thread (viewed as a module) can be preempted by others at any non-deterministically chosen program points, and the target program is usually more fine-grained and have more interaction points than the source. It is unclear if their approach can be applied to compilation of concurrent programs.

Stewart et al. [2015] do argue that Compositional CompCert may be extended for certified compositional compilation of data-race-free (DRF) concurrent programs. They argue that, for DRF concurrent programs, the behaviors of the threads under the standard interleaving semantics should be equivalent to those in some non-preemptive semantics where the control of the CPU switches between threads at certain designated program points. Since a thread cannot be interrupted between two consecutive switch points in the non-preemptive semantics, the code segment between the two switch points can be compiled as sequential code. The sequential compilation is sound as long as the switch points are viewed as external function calls so that optimizations do not go beyond them. Although the argument is plausible, there are still significant challenges in building fully certified compositional compilers for DRF concurrent programs:

- We need a proper formulation of the non-preemptive semantics and the notion of DRF. On the one hand, the formulation relies on the synchronization constructs in the language and the notion of footprints (i.e. the memory locations accessed in each step), On the other hand, like the interaction semantics in Compositional CompCert [Stewart et al. 2015], the formulation

50 should be *language-independent* to support general compositional compilation, where the  
51 modules can be implemented in different languages.

- 52 • We need to prove that DRF programs behave the same in the standard preemptive semantics  
53 and in our non-preemptive semantics, preferably in a language-independent setting. Also  
54 the equivalence should be strong enough to preserve the termination of programs, a natural  
55 correctness requirement for certified compilation. Although such semantics equivalence  
56 has been known as a folklore theorem, we have not seen any mechanized proofs of such  
57 *termination-preserving* semantics equivalence in a language-independent setting. As we will  
58 explain in Sec. 2.2, the proofs can actually be quite challenging, and the complexity is affected  
59 by many factors of language semantics design, such as semantics for memory allocation.
- 60 • We need to prove the compilation preserves DRF, which ensures that the target compiled  
61 from a DRF source is DRF too. Then the behaviors of the target program in the preemptive  
62 semantics are the same as those in the non-preemptive semantics. However, since a data  
63 race involves the behaviors of at least two threads, DRF is not a thread-local property. Then,  
64 how do we prove DRF-preservation of a compositional compilation, which compiles one  
65 module/thread at a time without knowledge of other modules/threads? That is, we need a  
66 compositional proof of the non-local DRF-preservation property.
- 67 • Concurrency introduces non-determinism, which makes it difficult to directly reuse the  
68 downward simulation proofs in CompCert and Compositional CompCert. Instead of proving  
69 the source program simulates the target (i.e., upward simulation), CompCert proves the  
70 reverse direction (downward simulation), and then derive the semantics equivalence based on  
71 the determinism of the language semantics. Since the source is usually more coarse-grained  
72 than the target, the downward simulation is simpler to prove than the upward one. However,  
73 it is unclear if such a proof strategy can still be applied for concurrent languages, whose  
74 semantics is non-deterministic.

75 We will explain the challenges in detail in Sec. 2, and discuss more related work on certified  
76 compilation in Sec. 9.

77 In this paper, we propose a language-independent framework consisting of the key semantics  
78 components and verification steps that bridge the gap between compilation for sequential programs  
79 and for DRF concurrent programs. We also apply our framework to verify the correctness of  
80 CompCert x86 backend for compositional compilation of DRF programs. Our work is based on  
81 previous work on certified compilation, but makes the following new contributions:

- 82 • We design a compositional *footprint-preserving simulation* as the correctness formulation of  
83 separate compilation for sequential modules. As an extension of the simulation in Composi-  
84 tional CompCert, our simulation considers module interactions at both external function calls  
85 and synchronization points, thus is compositional with respect to both module linking and  
86 non-preemptive concurrency. It also requires that the footprints of the steps made by the  
87 source and target modules be related, where footprints refer to the set of memory locations  
88 accessed at each step of execution. This way we reduce the proof of DRF-preservation for  
89 whole programs into proofs of *local* footprint preservation.
- 90 • We work with an *abstract* programming language, which is not tied to any specific synchroni-  
91 zation constructs such as locks but uses abstract labels to model how such constructs interact  
92 with other modules. It also abstracts away the concrete primitives that accesses memory.  
93 We introduce the notion of *well-defined languages* to enforce a set of constraints over the  
94 state transitions and the related footprints, which actually give an extensional interpretation  
95 of footprints. These constraints are satisfied by various real languages such as Cminor and  
96 x86 assembly. With the abstract language, we study the equivalence between preemptive  
97

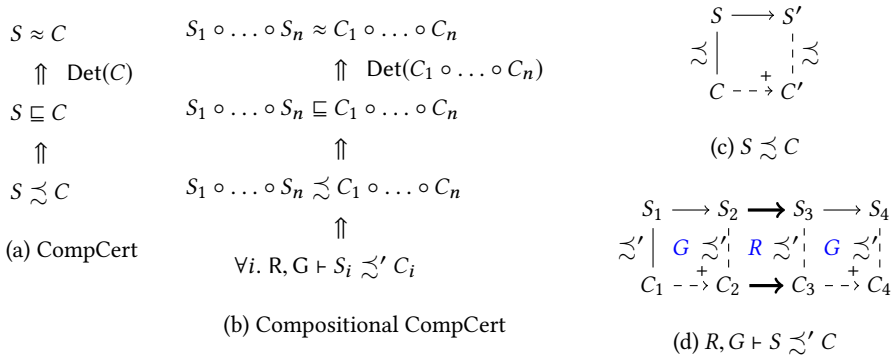


Fig. 1. Proof structures of certified compilation.

and non-preemptive semantics, the equivalence between DRF and NPDRF (the notion of race-freedom defined in the non-preemptive setting [Xiao et al. 2018]), and the properties of our new simulation. As a result, the lemmas in our proof framework are re-usable when instantiating to real languages.

- Putting all these together, our framework (see Fig. 2) is the first to build certified compositional compilation for concurrent programs *from sequential compilation*. It highlights the importance of DRF preservation for correct compilation. It also shows a possible way to adapt the existing work of CompCert and Compositional CompCert to interleaving concurrency.
- As an instantiation of our language-independent compilation verification framework, we instantiate the source and target languages as Cminor and x86 assembly with lock prefix. We have successfully proved that the CompCert-3.0.1 [2017] x86 backend (including all the translation phases and one optimization phase) satisfy our compilation correctness criterion. In the verification of the CompCert compilation phases, we reuse a considerable amount of the original CompCert proofs, with minor adjustment for footprint-preservation. The proofs for each phase take less than one person week on average.

In the rest of this paper, we first analyze the challenge and give an overview of our approach in Sec. 2. We give the basic technical setting in Sec. 3, including the preemptive semantics and the refinement definition. Sec. 4 presents the non-preemptive semantics, which is the basis for both our new simulation (Sec. 5) and the NPDRF definition (Sec. 6). We show the final theorem in Sec. 7, and discuss the implementation details in Sec. 8. We conclude and discuss related work in Sec. 9.

## 2 INFORMAL DEVELOPMENT

Below we first give an overview of the main ideas in CompCert [Leroy 2009a] and Compositional CompCert [Stewart et al. 2015] as starting points for our work. Then we explain the challenges and our ideas in reusing them to build certified compositional compilation for concurrent programs.

### 2.1 Background

**2.1.1 CompCert.** The pioneer work on CompCert [Leroy 2009a,b] builds certified compilation for sequential programs. Fig. 1(a) shows its key proof structure.

The compilation  $Comp$  is correct, if for every source program  $S$ , the compiled code  $C$  preserves the semantics of  $S$ . That is,  $Correct(Comp) \stackrel{\text{def}}{=} \forall S, C. Comp(S) = C \implies S \approx C$ . Here the semantics preservation  $S \approx C$  requires  $S$  and  $C$  have the same sets of observable event traces:

$$S \approx C \text{ iff } \forall \mathcal{B}. Etr(S, \mathcal{B}) \iff Etr(C, \mathcal{B}).$$

Here we write  $Etr(S, \mathcal{B})$  to mean that an execution of  $S$  produces the observable event trace  $\mathcal{B}$ , and likewise for  $C$ . The observable events include control effects (e.g., termination and exceptions) and input-output events, which in CompCert correspond to invocations of external functions.

To verify  $S \approx C$ , CompCert relies on the determinism of the target language (written as  $\text{Det}(C)$  in Fig. 1(a)) and proves only the downward direction  $S \sqsubseteq C$ , i.e.,  $S$  is a refinement of  $C$ .

$$S \sqsubseteq C \text{ iff } \forall \mathcal{B}. Etr(S, \mathcal{B}) \implies Etr(C, \mathcal{B}).$$

The determinism  $\text{Det}(C)$  ensures that  $C$  admits only one observable behavior, so one can derive the upward refinement  $S \sqsupseteq C$  from  $S \sqsubseteq C$ . The latter is then proved by constructing a (downward) simulation relation  $S \lesssim C$ . Depicted in Fig. 1(c), the simulation  $S \lesssim C$  establishes a *consistency relation* between  $S$  and  $C$ , which is always preserved under some correspondence between the executions of  $S$  and  $C$ . Every step of  $S$  must correspond to zero-or-more steps of  $C$ . Figure 1(c) shows the case when the source step corresponds to multiple target steps.

The simulation  $\lesssim$  serves as a proof technique for verifying whole-program compilation, but it is not compositional and *cannot* be used for verifying separate compilation. This is because  $\lesssim$  does not take into account the interactions with other modules which may update the shared resource.

**2.1.2 Compositional CompCert.** Compositional CompCert [Stewart et al. 2015] supports separate compilation by re-defining the simulation relation for modules. Figure 1(b) shows the proof structure of Compositional CompCert. Suppose we have separate compilation  $Comp_1, \dots, Comp_n$ . Each  $Comp_i$  transforms a source module  $S_i$  to a target module  $C_i$ . The overall compilation is correct if, when linked together, the target modules  $C_1 \circ \dots \circ C_n$  preserve the semantics of the source modules  $S_1 \circ \dots \circ S_n$  (here we write  $\circ$  as the module-linking operator). For example, the following program consists of two modules. The function  $f$  in module  $S1$  calls the external function  $g$ . The external module  $S2$  may access the variable  $b$  in  $S1$ .

```

173 // Module S1
174 extern void g(int *x);
175 int f(){
176     int a = 0, b = 0;
177     g(&b);
178     return a + b; }
179
180 // Module S2
181 int g(int *x){
182     *x = 3;
183 }

```

(2.1)

Suppose the two modules  $S1$  and  $S2$  are independently compiled to the target modules  $C1$  and  $C2$ . Aligned with the external call in  $S1$ ,  $C1$  also calls the external function  $g$  in  $C2$ . The correctness of the overall compilation requires  $(S1 \circ S2) \approx (C1 \circ C2)$ .

With the determinism of the target modules, we only need to prove the downward refinement  $(S1 \circ S2) \sqsubseteq (C1 \circ C2)$ , which is reduced to proving  $(S1 \circ S2) \lesssim (C1 \circ C2)$ , just as in CompCert. Ideally we hope to know  $(S1 \circ S2) \lesssim (C1 \circ C2)$  from  $S1 \lesssim C1$  and  $S2 \lesssim C2$ , and ensure the latter two by correctness of  $Comp_1, \dots, Comp_n$ . However, the CompCert simulation  $\lesssim$  is not compositional.

To achieve compositionality, Compositional CompCert defines the simulation relation  $\lesssim'$  shown in Fig 1(d). It is parameterized with the interactions between modules, formulated as the rely/guarantee conditions  $R$  and  $G$  [Jones 1983]. The rely condition  $R$  of a module specifies the permitted state transitions of its environment (i.e., other modules that may be linked together) at both the source and target levels, and the guarantee  $G$  specifies the possible transitions made by the module itself. The simulation diagram in Fig 1(d) requires that the steps of the current module (the thin arrows) be allowed in  $G$ , and the simulation  $\lesssim'$  be preserved by the environment steps  $R$  (the thick arrows). Note that the  $R$  steps happen only at the *external function calls* of the current module. It models the general callee behaviors. Such a simulation is compositional as long as the rely/guarantee conditions of linked modules are compatible.

Compositional CompCert proves that the CompCert compilation phases satisfy the new simulation  $\approx'$  (which is stronger than the CompCert simulation  $\approx$ ). The intuition is that the compiler optimizations do not go beyond external calls (unless only local variables get involved). That is, for the above example (2.1), the compiler cannot do optimizations based on the assumption that  $b$  is 0 at the last line of the function  $f$ .

Since the  $R$  steps happen only at the *external function calls* in Compositional CompCert, it cannot be applied to concurrent programs, where module/thread interactions may occur at any program point. However, if we consider race-free concurrent programs only, where threads are properly synchronized, we may consider the interleaving at certain synchronization points only. It is a well-known folklore theorem that DRF programs in interleaving semantics behave the same as in non-preemptive semantics. For instance, the following program (2.2) uses a lock to synchronize the accesses of the shared variables, and it is race-free. Its behaviors are the same as those when the threads yield controls at `lock()` and `unlock()` only. That is, we can view lines 1–2 and lines 4–5 in either thread as sequential code that will not be interfered by the other. The interactions occur only at the boundaries of the critical regions.

1	<code>r1 = 1;</code>		<code>r2 = 2;</code>	(2.2)
2	<code>r1 = r1 + 1;</code>		<code>r2 = r2 + 1;</code>	
3	<code>lock();</code>		<code>lock();</code>	
4	<code>x = 1;</code>		<code>x = 2;</code>	
5	<code>y = x + 1;</code>		<code>y = x + 1;</code>	
6	<code>unlock();</code>		<code>unlock();</code>	

Intuitively, we can use Compositional CompCert to compile the program, where the code segment between two consecutive switch points is compiled as sequential code. By viewing the switch points as special external function calls, the simulation  $\approx'$  can be applied to relate the non-preemptive executions of the source and target modules.

Although the idea is plausible, we have to address several key challenges to really apply it to build certified compositional compilers, as we explain below.

## 2.2 Challenges and Our Approaches in Verifying Concurrent Program Compilation

**2.2.1 How to give language-independent formulation of DRF and non-preemptive semantics?** The interaction semantics in Compositional CompCert describes module interaction without referring to the concrete languages used to implement the modules. This allows composition of modules implemented in different languages. We would like to follow the semantics framework, but how do we define DRF and non-preemptive semantics if we do not even know the concrete synchronization constructs and the commands that access memory?

*Our solution.* We extend the module-local semantics in Compositional CompCert so that each local step of a module reports its footprints, i.e. the memory locations it accesses. Instead of relying on the concrete memory-access commands to define what valid footprints are, we introduce the notion of *well-defined languages* (in Sec. 3) to specify the requirements over the state transitions and the related footprints. For instance, we require the behavior of each step is affected by the read set only, and each step does *not* touch the memory content outside of the write set. When we instantiate the framework with real languages, we prove they satisfy these requirements.

Besides, we also allow module-local steps to generate messages `EntAtom` and `ExtAtom` to indicate the boundary of the atomic operations inside the module. The concrete commands that generate these messages are unspecified, which can be different in different modules. In Sec. 8.1 we show how the lock-prefixed x86 instructions as synchronization constructs generate these messages.

246 2.2.2 *What memory model to use in the proofs?* The choice of memory models could greatly  
 247 affect the complexity of proofs. For instance, using the same memory model as CompCert allows  
 248 us to reuse CompCert proofs, but it also causes many problems. The CompCert memory model  
 249 records the next available block number `nextblock` for memory allocation. Using the model under  
 250 a concurrent setting may allow all threads to share one `nextblock`. Then allocation in one thread  
 251 would affect the subsequent allocations in other threads. This breaks the property that re-ordering  
 252 non-conflicting operations from different threads would *not* affect the final states, which is a key  
 253 lemma we rely on to prove the equivalence between preemptive and non-preemptive semantics for  
 254 DRF programs. In addition, sharing the `nextblock` by all threads also means we have to keep track  
 255 of the ownership of each allocated block when we reason about footprints.

256 *Our solution.* We decide to use a different memory model from the one of CompCert. We reserve  
 257 separate address spaces  $F$  for memory allocation in different threads (see Sec. 3). Therefore allocation  
 258 of one thread would not affect behaviors of others. This greatly simplifies the semantics and the  
 259 proofs, but also makes it almost impossible to reuse CompCert proofs, as we explain in Sec. 8.2. We  
 260 address this problem by establishing some semantics equivalence between our memory model and  
 261 the CompCert memory model (shown in Sec. 8.2.1).

262 2.2.3 *How to compositionally prove DRF-preservation?* The simulation  $\approx'$  in Compositional  
 263 CompCert is not sufficient to ensure DRF-preservation. As we have explained, DRF is a whole-  
 264 program property, and so is DRF-preservation. To support separate compilation, we need to reduce  
 265 the requirement of DRF-preservation on whole programs to some requirements on single threads.  
 266 In particular, we hope to encode the requirements in the thread-local and module-local simulation.

267 *Our solution.* We propose a new compositional simulation  $\approx$  which extends  $\approx'$  with the require-  
 268 ments of footprint-preservation on single threads. In detail, based on the simulation diagram in  
 269 Fig. 1(d), we additionally require *footprint consistency* saying that the target  $C$  should have the same  
 270 or smaller footprints than the source  $S$  during related transitions. For instance, when compiling  
 271 lines 4-5 of the left thread in (2.2), the target is only allowed to read  $x$  and write to  $x$  and  $y$ .

272 Note that we check footprint consistency at switch points only. This way we allow compiler  
 273 optimizations as long as they do not go beyond the switch points. For the example in (2.2), the  
 274 target of lines 4-5 of the left thread could be  $y=2; x=1$  where the writes to  $x$  and  $y$  are swapped.

275 2.2.4 *How to flip refinement/simulation with non-deterministic behaviors?* As we explained, the  
 276 last steps of CompCert and Compositional CompCert in Fig. 1 derive semantics equivalence  $\approx$  (or  
 277 the upward refinement  $\sqsupseteq$ ) from the downward refinement  $\sqsubseteq$  using determinism of target programs.  
 278 Actually the simulations  $\approx$  and  $\approx'$  can also be flipped if the target programs are deterministic.  
 279 It is unclear if the refinement or simulation can still be flipped in the concurrent settings where  
 280 programs have non-deterministic behaviors. The problem is that the target program can be more  
 281 fine-grained and have more non-deterministic interleavings than the source.

282 *Our solution.* Data-race-freedom and the non-preemptive semantics come to the rescue. For  
 283 DRF programs, the switch points in the target are aligned with those in the source under the  
 284 non-preemptive semantics. Also the target and source programs can always make the same non-  
 285 deterministic choices of thread switching. Thus we are able to flip the downward simulation to  
 286 derive upward simulation, assuming determinism of each target module.

## 291 2.3 The framework and key semantics components

292 Figure 2 shows the semantics components and proof steps in our framework. The ultimate goal of  
 293 our compilation correctness proof is to show the semantics preservation of the source and target  
 294

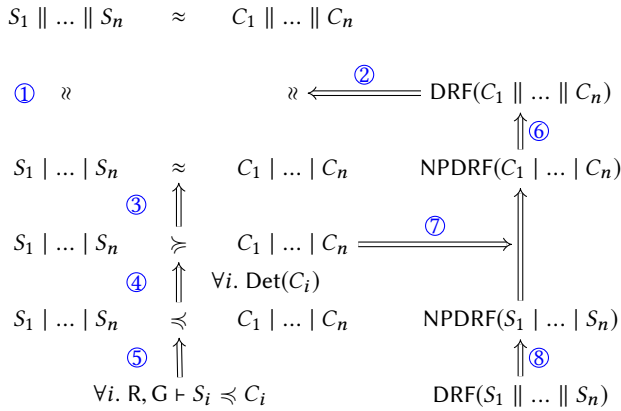


Fig. 2. Our framework

programs (i.e.,  $S_1 \parallel \dots \parallel S_n \approx C_1 \parallel \dots \parallel C_n$  at the top of the figure). This follows the correctness of the separate compilation of each module, formulated as  $R, G \vdash S_i \preceq C_i$  (the bottom left), which is our new footprint-preserving module-local simulation explained above. We do the proofs in the following steps. Note that the double arrows in the figure represent logical implication.

**First**, we restrict the compilation to source preemptive programs that are data-race-free (i.e.,  $DRF(S_1 \parallel \dots \parallel S_n)$  at the right bottom of the figure). Then from the equivalence between the preemptive and non-preemptive semantics, we derive ①, the equivalence between  $S_1 \parallel \dots \parallel S_n$  and  $S_1 \mid \dots \mid S_n$ , the latter representing non-preemptive execution of the threads. Similarly, if we have  $DRF(C_1 \parallel \dots \parallel C_n)$  (at the top right), we can derive ②. With ① and ②, we can derive the semantics preservation  $\approx$  between preemptive programs from the semantics preservation between their non-preemptive counterparts.

**Second**, DRF of the target programs is obtained through the right path ⑥, ⑦ and ⑧. We define a notion of DRF for non-preemptive programs (called NPDRF in the figure), making it equivalent to DRF, from which we can derive ⑥ and ⑧. To know  $NPDRF(C_1 \mid \dots \mid C_n)$  from  $NPDRF(S_1 \mid \dots \mid S_n)$ , we need a DRF-preserving simulation  $\preceq$  between non-preemptive programs (see ⑦).

**Third**, the DRF-preserving simulation  $\preceq$  for non-preemptive whole programs can be derived by composing our footprint-preservation local simulation (step ⑤). Given the the downward whole-program simulation, we flip it to get an *upward* one (step ④), with the extra premise that the local execution in each target module is deterministic. Using the simulation in both directions we derive the equivalence (step ③).

Note that the notations used here are simplified ones to give a semi-formal overview of the key ideas. We may use different notations in the formal development in the following sections.

### 3 THE LANGUAGE AND THE PREEMPTIVE SEMANTICS

#### 3.1 The Abstract Language

Figure 3 shows the syntax and the state model of an abstract language for preemptive concurrent programming. A program  $P$  consists of  $n$  threads running in parallel. Each thread starts execution from an entry  $f$ , which points to the code segment defined in the code  $\pi$  of a module declared in  $\Pi$ . Each module declaration is a triple consisting of the language declaration  $tl$ , the global environment  $ge$  containing the addresses of static global variables declared in the module, and the code  $\pi$  of the module. Here we use  $\mathcal{P}(Addr)$  to represent the powerset of memory addresses  $Addr$ .

344	(Prog)	$P, \mathbb{P} ::= \mathbf{let} \Pi \mathbf{in} f_1 \parallel \dots \parallel f_n$	(Entry)	$f \in \mathit{String}$
345	(MdSet)	$\Pi, \Gamma ::= \{(t_1, ge_1, \pi_1), \dots, (t_m, ge_m, \pi_m)\}$	(Module)	$\pi, \gamma ::= \dots$
346	(Lang)	$tl, sl ::= (\mathit{Module}, \mathit{Core}, \mathit{InitCore}, \mapsto)$	(Core)	$\kappa, \mathbb{k} ::= \dots$
347		$ge \in \mathcal{P}(\mathit{Addr})$		$\mathit{InitCore} \in \mathit{Module} \rightarrow \mathit{Entry} \rightarrow \mathit{Core}$
348		$\mapsto \in \mathit{FList} \times (\mathit{Core} \times \mathit{State}) \rightarrow \{\mathcal{P}((\mathit{Msg} \times \mathit{FtPrt}) \times (\mathit{Core} \times \mathit{State})), \mathbf{abort}\}$		
349				
350	(ThrdID)	$t \in \mathbb{N}$	(Addr)	$l ::= \dots$
351	(State)	$\sigma, \Sigma \in \mathit{Addr} \rightarrow_{\mathit{fin}} \mathit{Val}$	(Val)	$v ::= l \mid \dots$
352	(FtPrt)	$\delta, \Delta ::= (rs, ws) \quad \text{where } rs, ws \in \mathcal{P}(\mathit{Addr})$	(FList)	$F, \mathbb{F} \in \mathcal{P}^\omega(\mathit{Addr})$
353	(Msg)	$\iota ::= \tau \mid e \mid \mathbf{ret} \mid \mathit{EntAtom} \mid \mathit{ExtAtom}$	(Event)	$e ::= \dots$
354	(Config)	$\phi, \Phi ::= (\kappa, \sigma) \mid \mathbf{abort}$		

Fig. 3. The Abstract Concurrent Language

Since different modules may be written in different languages, we define the abstract module language  $tl$  as a tuple  $(\mathit{Module}, \mathit{Core}, \mathit{InitCore}, \mapsto)$ .  $\mathit{Module}$  describes the syntax of programs. Following Compositional CompCert [Stewart et al. 2015],  $\mathit{Core}$  is the set of internal “core” states, which can be instantiated to control continuations, instruction streams, register files, etc. The function  $\mathit{InitCore}$  is called whenever a thread is created or a cross-module external function call is made. Given a module  $\pi$  and an entry name  $f$ ,  $\mathit{InitCore}$  returns the initial “core” state  $\kappa$  (which is undefined if the entry is not contained in the module). In this paper we mainly focus on compositional compilation of concurrent programs, where threads can be from different modules. Since cross-module external calls are mostly orthogonal, we omit them to simplify the presentation. *We do support external calls in our Coq implementation in the same way as in Compositional CompCert.* The labelled transition “ $\mapsto$ ” models the local execution of a module, which we explain below. To instantiate a language  $tl$ , one needs to provide concrete definitions for each component described above.

*Module-local semantics.* The local execution steps inside a module is modeled as a labeled transition in the form of  $F \vdash_{tl} (\kappa, \sigma) \xrightarrow[\delta]{\tau} (\kappa', \sigma')$  (or  $F \vdash_{tl} (\kappa, \sigma) \xrightarrow[\delta]{\iota} \mathbf{abort}$  if the step goes wrong). In addition to its core state  $\kappa$ , each module can also access the memory state  $\sigma$ , which is a *finite* partial mapping from memory addresses to values. We leave the meta-type  $\mathit{Addr}$  undefined, which can be instantiated for specific languages. For instance, each address in the Cminor memory model can be instantiated into a pair of a block number and an offset. Each step may change the core state and the memory state into  $\kappa'$  and  $\sigma'$  respectively. It is also labeled with a message  $\iota$  and a footprint  $\delta$ . Note that the step relation can be *non-deterministic*. That is, given a pair  $(\kappa, \sigma)$ , there can be more than one resulting states  $(\kappa', \sigma')$ , and the corresponding messages  $\iota$  and footprints  $\delta$  can be different. *To avoid clutter, we usually omit the parameter  $tl$  in the judgment.*

Each module also has a *free list*  $F$ . It is the pool of memory addresses from which fresh memory cells are allocated. It can be viewed as the preserved space for allocating local stack frames, where “ $F\text{-dom}(\sigma)$ ” is the set of free addresses, shown in Fig. 4 as the part outside of the boundary of  $\sigma$ . Initially we require  $F \cap \text{dom}(\sigma) = \emptyset$ , and the only memory accessible by the module is the static global variables declared in the  $ge$  of all modules, represented as the *shared* part  $S$  in Fig. 4. The local execution of a module may allocate memory from its  $F$ , which enlarges the state  $\sigma$ . The memory allocated from

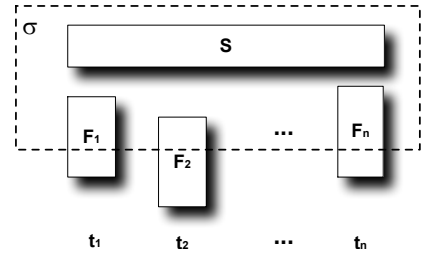


Fig. 4. The state model



$F$  is exclusively owned by this module. We allow the set  $F$  to be *infinite*.  $F$  for different modules must be *disjoint*.

The messages  $\iota$  contain information about the module-local steps. To simplify the presentation, we only consider externally observable events  $e$  (such as outputs), termination of threads (`ret`), and the beginning and the end of *atomic blocks* (`EntAtom` and `ExtAtom`). Any other steps are silent, labeled with  $\tau$ . The label  $\tau$  is often omitted for cleaner presentation. Atomic blocks are the language constructs to ensure sequential execution of code blocks inside them. They can be implemented differently in different module languages. In the example given in Sec. 8.1 we instantiate atomic blocks with `lock` prefixed instructions in x86. The messages define the protocols of communications with the global whole-program semantics (presented below). All the module languages use the same message formats. They allow us to abstract away details of the module languages, and focus on the interactions with other modules and the external observer (the latter observes  $e$  only).

The footprint  $\delta$  is defined as a pair  $(rs, ws)$ , which records the memory locations *read* and *written* in this step. Recording the footprint allows us to discuss races between threads in the following sections. We write `emp` for the special footprint where both the read and write sets are empty. Below we may directly use  $\delta$  as a set, which is a shorthand for  $\delta.rs \cup \delta.ws$ .

*Conventions.* We use two sets of symbols to distinguish the source and the target level notations. The blackboard bold or capital letters (e.g.,  $\mathbb{P}$ ,  $\mathbb{F}$ ,  $\mathbb{k}$  and  $\Sigma$ ) are used for the source, while their counterparts (e.g.,  $P$ ,  $F$ ,  $\kappa$  and  $\sigma$ ) are for the target. The set of modules at source is written as  $\Gamma$ , to distinguish from the target  $\Pi$ . Similarly,  $\gamma$  is a source module while  $\pi$  is a target one.

*Well-defined languages.* Although the abstract module language  $tl$  can be instantiated with different real languages, the instantiation needs to satisfy certain basic requirements. We define these requirements as *well-defined languages* below in Def. 1. It gives us an extensional interpretation of footprints. It is also used to prove properties of DRF programs in the sections below.

**Definition 1** (Well-Defined Languages).  $\text{wd}(tl)$  iff, for any execution step  $F \vdash (\kappa, \sigma) \xrightarrow[\delta]{\iota} (\kappa', \sigma')$  in this language, all of the following hold (some auxiliary definitions are in Fig. 5):

- (1)  $\text{forward}(\sigma, \sigma')$ ;
- (2)  $\text{LEffect}(\sigma, \sigma', \delta, F)$
- (3) for any  $\sigma_1$ , if  $\text{LEqPre}(\sigma, \sigma_1, \delta, F)$ , then there exists  $\sigma'_1$  such that  $F \vdash (\kappa, \sigma_1) \xrightarrow[\delta]{\iota} (\kappa', \sigma'_1)$  and  $\text{LEqPost}(\sigma', \sigma'_1, \delta, F)$ .
- (4) for any  $\delta_0$  and  $\sigma_1$ , if  $\delta_0 = \bigcup \{ \delta'' \mid \exists \kappa'', \sigma''. F \vdash (\kappa, \sigma) \xrightarrow[\delta'']{\tau} (\kappa'', \sigma'') \}$  and  $\text{LEqPre}(\sigma, \sigma_1, \delta_0, F)$ , then  $\forall \kappa'_1, \sigma'_1, \iota'_1, \delta''_1. F \vdash (\kappa, \sigma_1) \xrightarrow[\delta''_1]{\iota'_1} (\kappa'_1, \sigma'_1) \implies \exists \sigma''. F \vdash (\kappa, \sigma) \xrightarrow[\delta''_1]{\iota'_1} (\kappa'_1, \sigma'')$ .

It requires that a step may enlarge the memory domain but cannot reduce it (see Item (1), where  $\text{forward}(\sigma, \sigma')$  is defined in Fig. 5), and the additional memory should be allocated from  $F$  and included in the write set (as required in Item (2)). This requirement follows the CompCert memory model where memory disposal does not really remove the locations from the memory (they just become invalid). Item (2) also requires that the memory out of the write set should keep unchanged, as described by  $\sigma \xrightarrow{\text{dom}(\sigma) - \delta.ws} \sigma'$  (which is defined at the top of Fig. 5). Item (3) says the memory updates and allocation only depend on the memory content in the read set, and the set of memory locations already allocated from  $F$ . Item (4) requires that the non-determinism of each step is not affected by memory contents outside of all the possible read sets. Here  $\delta_0$  is the union of footprints in all possible steps. Then for any state  $\sigma_1$  with the same contents in the read set of  $\delta_0$  and the same set of allocated addresses as in  $\sigma$ , it does not generate any new behavior that cannot be generated

$$\begin{aligned}
442 \quad \sigma &\stackrel{rs}{=} \sigma' \quad \text{iff} \quad \forall l \in rs. l \notin (\text{dom}(\sigma) \cup \text{dom}(\sigma')) \vee l \in (\text{dom}(\sigma) \cap \text{dom}(\sigma')) \wedge \sigma(l) = \sigma'(l) \\
443 \quad \delta &\subseteq \delta' \quad \text{iff} \quad (\delta.rs \subseteq \delta'.rs) \wedge (\delta.ws \subseteq \delta'.ws) \quad \delta \cup \delta' \stackrel{\text{def}}{=} (\delta.rs \cup \delta'.rs, \delta.ws \cup \delta'.ws) \\
444 \\
445 \quad \text{forward}(\sigma, \sigma') &\text{ iff } (\text{dom}(\sigma) \subseteq \text{dom}(\sigma')) \\
446 \quad \text{LEqPre}(\sigma_1, \sigma_2, \delta, F) &\text{ iff } \sigma_1 \stackrel{\delta.rs}{=} \sigma_2 \wedge \\
447 &\quad (\text{dom}(\sigma_1) \cap \delta.ws) = (\text{dom}(\sigma_2) \cap \delta.ws) \wedge (\text{dom}(\sigma_1) \cap F) = (\text{dom}(\sigma_2) \cap F) \\
448 \quad \text{LEqPost}(\sigma'_1, \sigma'_2, \delta, F) &\text{ iff } \sigma'_1 \stackrel{\delta.ws}{=} \sigma'_2 \wedge (\text{dom}(\sigma'_1) \cap F) = (\text{dom}(\sigma'_2) \cap F) \\
449 \quad \text{LEffect}(\sigma_1, \sigma_2, \delta, F) &\text{ iff } \sigma_1 \stackrel{\text{dom}(\sigma_1) - \delta.ws}{=} \sigma_2 \wedge (\text{dom}(\sigma_2) - \text{dom}(\sigma_1)) \subseteq (\delta.ws \cap F) \\
450
\end{aligned}$$

Fig. 5. Notations and auxiliary definitions about states and footprints

by  $\sigma$ . This property ensures that, for data-race-free programs, the possible execution steps of a thread cannot be affected by the interleaving with other threads.

### 3.2 Global preemptive semantics

$$\begin{aligned}
456 \quad (\text{World}) \quad W, \mathbb{W} &::= (T, t, d, \sigma) & (\text{AtomBit}) \quad d &::= 0 \mid 1 \\
457 \\
458 \quad (\text{ThrdPool}) \quad T, \mathbb{T} &::= \{t_1 \rightsquigarrow (tl_1, F_1, \kappa_1), \dots, t_n \rightsquigarrow (tl_n, F_n, \kappa_n)\} \\
459 \quad (\text{GMsg}) \quad o &::= \tau \mid e \mid \text{sw} & (\text{ASet}) \quad S, \mathbb{S} &\in \mathcal{P}(\text{Addr})
\end{aligned}$$

Fig. 6 defines a set of global semantics rules to manipulate the preemption among threads. As shown above, the global world  $W$  consists of the thread pool  $T$ , the ID  $t$  of the thread currently being executed, a bit  $d$  indicating whether the current thread is in an atomic block or not, and the memory state  $\sigma$ . The thread pool  $T$  maps a thread ID to a triple recording the module language  $tl$ , the free list  $F$ , and the current core state  $\kappa$ .

The Load rule in Fig. 6 shows the initialization of the world from the program and an initial program state  $\sigma$ . We assume the code  $\pi$  of modules in  $\Pi$  has disjoint entries (i.e. code labels). Therefore, for each entry  $f_i$  of the thread  $i$ , we can find at most one module declaration in  $\Pi$  containing it. The core state  $\kappa_i$  is created through the `InitCore` function of the corresponding language. For each thread, we also assign a local address space  $F$  for allocation of stack frames. The local address spaces for the threads must be disjoint, and initially they are disjoint with the domain of  $\sigma$ . We non-deterministically pick a thread  $t$  as the current thread. The bit  $d$  is set to 0, indicating that the current thread is *not* in the atomic block.

The transition rules of the whole world is given in Fig. 6. Like local transitions, global transitions are also labeled with footprints and messages  $o$ . Here  $o$  marks a silent  $\tau$  step, a step with external event  $e$ , or a switch (sw) step, as defined above. Different from the messages  $\iota$  that record both external events and inter-module communications, the global messages  $o$  only record the externally observable events, i.e., events that are observable to the human being sitting in front of the computer or printer. The switch message sw is an exception, which is for verification purpose only.

Each global step executes the module locally and processes the message of the local transition. The  $\tau$ -step rule and the Print rule show the local  $\tau$ -step and the step generating an external event, respectively. Note that the Print rule requires that the flag  $d$  must be 0, i.e., external events can be generated only outside of atomic blocks. Also the step generating an external event does *not* access memory, so its footprint is empty (emp).

The EntAt and ExtAt rules correspond to the entry and exit of atomic blocks, respectively. The flag  $d$  is flipped in the two steps. Since context-switch can be done only when  $d$  is 0, as required by the Switch rule below, we know a thread in its atomic block cannot be preempted.

When the current thread terminates, we either remove it from the thread pool and then switch to another thread if there is one (see the Term rule), or terminate the whole program if not, as shown in the Done rule. The Term rule generates the sw message to record the context switch.

491	for all $i$ and $j$ in $\{1, \dots, n\}$ , and $i \neq j$ :		
492	$F_i \cap F_j = \emptyset \quad \text{dom}(\sigma) \cap F_i = \emptyset \quad tl_i.\text{InitCore}(\pi_i, f_i) = \kappa_i$ , where $(tl_i, ge_i, \pi_i) \in \Pi$		
493	$T = \{1 \rightsquigarrow (tl_1, F_1, \kappa_1), \dots, n \rightsquigarrow (tl_n, F_n, \kappa_n)\} \quad t \in \{1, \dots, n\}$		
494	$(\text{let } \Pi \text{ in } f_1 \parallel \dots \parallel f_n, \sigma) \xrightarrow{\text{load}} (T, t, 0, \sigma)$		Load
495			
496	$T(t) = (tl, F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\delta]{\tau} (\kappa', \sigma')$	$T(t) = (tl, F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{e} (\kappa', \sigma)$	
497	$T' = T\{t \rightsquigarrow (tl, F, \kappa')\}$	$T' = T\{t \rightsquigarrow (tl, F, \kappa')\}$	
498	$(T, t, d, \sigma) \xrightarrow[\delta]{\tau} (T', t, d, \sigma')$		$\tau$ -step
499			
500	$T(t) = (tl, F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\text{EntAtom}} (\kappa', \sigma)$		
501	$T' = T\{t \rightsquigarrow (tl, F, \kappa')\}$	$T(t) = (tl, F, \kappa) \quad t' \in \text{dom}(T) \quad t$	
502	$(T, t, 0, \sigma) \xrightarrow[\text{emp}]{\tau} (T', t, 1, \sigma)$		EntAt
503			
504	$T(t) = (tl, F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\text{ExtAtom}} (\kappa', \sigma)$		
505	$T' = T\{t \rightsquigarrow (tl, F, \kappa')\}$	$F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\text{ret}} (\kappa', \sigma)$	Term
506	$(T, t, 0, \sigma) \xrightarrow[\text{emp}]{\tau} (T', t, 0, \sigma)$		Done
507			
508	$T(t) = (tl, F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\tau} \mathbf{done}$		
509			
510	$t' \in \text{dom}(T)$	$T(t) = (tl, F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\delta]{\tau} \mathbf{abort}$	Abort
511	$(T, t, 0, \sigma) \xrightarrow[\text{emp}]{\text{sw}} (T, t', 0, \sigma)$		Switch
512			
513			

Fig. 6. The Preemptive Global Semantics

The Switch rule shows that context switch can occur at any program point outside of atomic blocks ( $d = 0$ ). This also indicates that the semantics is preemptive. The step is also marked with the sw message. The Abort rule says the whole program aborts if a local module aborts.

Below we write  $F \vdash \phi \xrightarrow[\delta]{\tau} \phi'$  for multiple silent-step transitions, where  $\delta$  is the accumulation of the footprints generated.  $F \vdash \phi \xrightarrow[\delta]{\tau} \phi'$  is for zero or multiple silent-step transitions, where  $\delta$  is emp for the case of zero step. Similarly, for global steps, we write  $W \xrightarrow[\delta]{\tau} W'$  for multiple silent-step transitions. Besides, we also write  $W \Rightarrow^+ W'$  for multiple steps that either are silent or produce sw events. It must contain at least one silent step. The meanings of  $W \Rightarrow^+ \mathbf{abort}$  and  $W \Rightarrow^+ \mathbf{done}$  are similar.  $W \xrightarrow[e]{\tau} W'$  represents multiple steps with exactly one  $e$  event produced (while other steps either are silent or produce sw events).

### 3.3 Event-Trace Refinement and Equivalence

The correctness of compilation for concurrent programs is defined as the event-trace refinement (or equivalence) between source and target programs. An externally observable event trace  $\mathcal{B}$  is a finite or infinite sequence of external events  $e$ , and may end with a termination marker **done** or an abortion marker **abort**. It is co-inductively defined as follows.

$$(\text{EvtTrace}) \quad \mathcal{B} ::= \mathbf{done} \mid \mathbf{abort} \mid \epsilon \mid e :: \mathcal{B} \quad (\text{co-inductive})$$

$$\text{ProgEtr}(P, \sigma, \mathcal{B}) \text{ iff } \exists W. ((P, \sigma) \xrightarrow{\text{load}} W) \wedge \text{Etr}(W, \mathcal{B})$$

$$\frac{W \Rightarrow^+ \mathbf{abort}}{\text{Etr}(W, \mathbf{abort})} \quad \frac{W \Rightarrow^+ \mathbf{done}}{\text{Etr}(W, \mathbf{done})} \quad \frac{W \xrightarrow[e]{\tau} W' \quad \text{Etr}(W', \mathcal{B})}{\text{Etr}(W, e :: \mathcal{B})} \quad \frac{W \Rightarrow^+ W' \quad \text{Etr}(W', \epsilon)}{\text{Etr}(W, \epsilon)}$$

We use  $\text{ProgEtr}(P, \sigma, \mathcal{B})$  to say that the trace  $\mathcal{B}$  can be produced by executing  $P$  with the initial state  $\sigma$ . The co-inductive definition of  $\text{Etr}(W, \mathcal{B})$  says that  $\mathcal{B}$  can be produced by executing  $W$ . Note we distinguish the traces of non-terminating (diverging) executions from those of terminating ones. If the execution of  $W$  diverges, its observable event trace  $\mathcal{B}$  is either of infinitely length, or finite but does not end with **done** or **abort** (called *silent divergence*, see the right-most rule above).

Then we define the refinement  $(\mathbb{P}, \Sigma) \sqsupseteq (P, \sigma)$  and the equivalence  $(\mathbb{P}, \Sigma) \approx (P, \sigma)$  below. They ensure that if  $(P, \sigma)$  has a diverging execution, so does  $(\mathbb{P}, \Sigma)$ . Thus the refinement and the equivalence relations preserve the termination of  $(\mathbb{P}, \Sigma)$ .

**Definition 2** (Event-Trace Refinement and Equivalence).

$(\mathbb{P}, \Sigma) \sqsupseteq (P, \sigma)$  iff  $\forall \mathcal{B}. \text{ProgEtr}(P, \sigma, \mathcal{B}) \implies \text{ProgEtr}(\mathbb{P}, \Sigma, \mathcal{B})$ .

$(\mathbb{P}, \Sigma) \approx (P, \sigma)$  iff  $\forall \mathcal{B}. \text{ProgEtr}(P, \sigma, \mathcal{B}) \iff \text{ProgEtr}(\mathbb{P}, \Sigma, \mathcal{B})$ .

Following CompCert, the compilation correctness assumes safety of the source programs. Below we use the event traces to define  $\text{Safe}(\mathbb{W})$  and  $\text{Safe}(\mathbb{P}, \Sigma)$ .

**Definition 3** (Safety).  $\text{Safe}(\mathbb{W})$  iff  $\neg \exists tr. \text{Etr}(\mathbb{W}, tr :: \text{abort})$ .

$\text{Safe}(\mathbb{P}, \Sigma)$  iff  $(\exists \mathbb{W}. (\mathbb{P}, \Sigma) \xrightarrow{\text{load}} \mathbb{W})$  and  $\forall \mathbb{W}. ((\mathbb{P}, \Sigma) \xrightarrow{\text{load}} \mathbb{W}) \implies \text{Safe}(\mathbb{W})$ .

## 4 THE NON-PREEMPTIVE SEMANTICS

A key step in our framework is to reduce the semantics preservation under the preemptive semantics to the semantics preservation in non-preemptive semantics. In this section we define the global non-preemptive semantics, where a thread interacts with other threads at only synchronization points (i.e., when it enters and exits atomic blocks, and outputs). The non-preemptive semantics is the basis for both our new simulation (see Sec. 5) and our NPDRF definition (see Sec. 6).

(NPProg)  $\hat{P} ::= \text{let } \Pi \text{ in } f_1 \mid \dots \mid f_n$       (NPWorld)  $\widehat{W}, \widehat{W}' ::= (T, t, \mathfrak{d}, \sigma)$

(AtomBitMap)  $\mathfrak{d} ::= \{t_1 \rightsquigarrow d_1, \dots, t_n \rightsquigarrow d_n\}$

To distinguish from the preemptive concurrency, we write  $\text{let } \Pi \text{ in } f_1 \mid \dots \mid f_n$  for the program with non-preemptive semantics, denoted by  $\hat{P}$ . As shown above, the non-preemptive global world  $\widehat{W}$  is defined similarly as the preemptive world  $W$ , except that  $\widehat{W}$  keeps an atomic bit map  $\mathfrak{d}$  recording whether each thread's next step is inside an atomic block. We need to record the atomic bits of all threads because the context switch may occur when a thread just enters an atomic block.

Fig. 7 defines the non-preemptive global steps  $\widehat{W} \xrightarrow[\delta]{\circ} \widehat{W}'$ . There is no rule like Switch of the preemptive semantics, since context switch occurs only at synchronization points in the non-preemptive setting. The rules  $\text{Print}_{\text{np}}$ ,  $\text{EntAt}_{\text{np}}$ ,  $\text{ExtAt}_{\text{np}}$  and  $\text{Term}_{\text{np}}$  execute one step of the current thread  $t$ , and then non-deterministically switch to a thread  $t'$  (which could just be  $t$ ). The corresponding global steps produce the sw events (or the external event  $e$  in the  $(\text{Print}_{\text{np}})$  rule). Note that in the  $\text{EntAt}_{\text{np}}$  rule, a thread may switch before it executes the body of the atomic block. Thus the global step needs to set the corresponding atomic bit in  $\mathfrak{d}$ , indicating that the thread must be inside an atomic block when it regains the control later. Other rules are very similar to their counterparts in the preemptive semantics in Fig. 6, and are not presented here.

Similar to Def. 2, we define  $(\hat{\mathbb{P}}, \Sigma) \sqsupseteq (\hat{P}, \sigma)$  and  $(\hat{\mathbb{P}}, \Sigma) \approx (\hat{P}, \sigma)$  for the refinement and equivalence between programs under non-preemptive semantics.

$$\begin{array}{c}
589 \quad \text{for all } i \text{ and } j \text{ in } \{1, \dots, n\}, \text{ and } i \neq j: \\
590 \quad F_i \cap F_j = \emptyset \quad \text{dom}(\sigma) \cap F_i = \emptyset \quad tl_i.\text{InitCore}(\pi_i, f_i) = \kappa_i, \quad \text{where } (tl_i, ge_i, \pi_i) \in \Pi \\
591 \quad T = \{1 \rightsquigarrow (tl_1, F_1, \kappa_1), \dots, n \rightsquigarrow (tl_n, F_n, \kappa_n)\} \quad t \in \{1, \dots, n\} \quad d = \{t_1 \rightsquigarrow 0, \dots, t_n \rightsquigarrow 0\} \\
592 \quad \hline \\
593 \quad \text{(let } \Pi \text{ in } f_1 \mid \dots \mid f_n, \sigma) \xrightarrow{\text{load}} (T, t, d, \sigma) \quad \text{Load}_{\text{np}} \\
594 \quad \hline \\
595 \quad \begin{array}{c} T(t) = (tl, F, \kappa) \quad d(t) = 0 \\ F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{e} (\kappa', \sigma) \quad t' \in \text{dom}(T) \end{array} \\
596 \quad \hline \\
597 \quad \begin{array}{c} T' = T\{t \rightsquigarrow (tl, F, \kappa')\} \\ (T, t, d, \sigma) \xrightarrow[\text{emp}]{e} (T', t', d, \sigma) \end{array} \quad \text{Print}_{\text{np}} \\
598 \quad \hline \\
599 \quad \begin{array}{c} T(t) = (tl, F, \kappa) \quad d(t) = 0 \\ F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\text{ret}} (\kappa', \sigma) \quad t' \in \text{dom}(T \setminus t) \end{array} \\
600 \quad \hline \\
601 \quad \begin{array}{c} T' = T \setminus t \\ (T, t, d, \sigma) \xrightarrow[\text{emp}]{\text{sw}} (T \setminus t, t', d \setminus t, \sigma) \end{array} \quad \text{Term}_{\text{np}} \\
602 \quad \hline \\
603 \quad \begin{array}{c} T(t) = (tl, F, \kappa) \quad d(t) = 0 \\ F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\text{EntAtom}} (\kappa', \sigma) \end{array} \\
604 \quad \hline \\
605 \quad \begin{array}{c} T' = T\{t \rightsquigarrow (tl, F, \kappa')\} \quad t' \in \text{dom}(T) \\ (T, t, d, \sigma) \xrightarrow[\text{emp}]{\text{sw}} (T', t', d\{t \rightsquigarrow 1\}, \sigma) \end{array} \quad \text{EntAt}_{\text{np}} \\
606 \quad \hline \\
607 \quad \begin{array}{c} T(t) = (tl, F, \kappa) \quad d(t) = 1 \\ F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\text{ExtAtom}} (\kappa', \sigma) \end{array} \\
608 \quad \hline \\
609 \quad \begin{array}{c} T' = T\{t \rightsquigarrow (tl, F, \kappa')\} \quad t' \in \text{dom}(T) \\ (T, t, d, \sigma) \xrightarrow[\text{emp}]{\text{sw}} (T', t', d\{t \rightsquigarrow 0\}, \sigma) \end{array} \quad \text{ExtAt}_{\text{np}}
\end{array}$$

Fig. 7. The Non-Preemptive Global Semantics

## 5 THE FOOTPRINT-PRESERVING COMPOSITIONAL SIMULATION

In this section, we define a module-local simulation as the correctness obligation of each module's compilation, which is compositional and preserves footprints, allowing us to derive a whole-program simulation that preserves data-race-freedom. We will discuss compositionality in Sec. 5.2 and postpone the discussions of DRF and NPDRF preservation to Sec. 6.

### 5.1 The Module-Local Simulation

As informally explained in Sec. 2, the simulation establishes a *consistency relation* between executions of the source module  $\gamma$  and the target one  $\pi$ . To achieve compositionality, our simulation uses rely/guarantee conditions to specify the interactions between the current module and its environment at switch points. The consistency relation should be preserved under the environment transitions allowed in the rely condition. So the key to define the local simulation is to define a proper consistency relation and proper rely/guarantee conditions.

*Footprint consistency.* As in CompCert, the consistency relation require that the source and the target generate the same external events. In addition, we also require that the target has the same or smaller footprints than the source, which is important to ensure DRF-preservation. Recall that the memory accessible by a thread  $t_i$  consists of two parts, the shared memory  $\mathbb{S}$  and the local memory allocated from  $\mathbb{F}_i$ , as shown in Fig. 4. DRF informally requires that the threads never access the shared memory in  $\mathbb{S}$  at the same time where at least one such access is a write.

We introduce the triple  $\mu$  below to record the key information about the shared memory at the source and the target.

$$\mu \stackrel{\text{def}}{=} (\mathbb{S}, S, f) \quad \text{where } \mathbb{S}, S \in \mathcal{P}(\text{Addr}) \text{ and } f \in \text{Addr} \rightarrow \text{Addr}$$

Here  $\mathbb{S}$  and  $S$  specify the shared memory locations at the source and the target respectively. The partial mapping  $f$  maps locations at the source level to those at the target. We require  $\mu$  to be well-formed, defined as  $\text{wf}(\mu)$  in Fig. 8. It says that the domain of  $f$  is  $\mathbb{S}$ ,  $f$  is injective, and maps shared locations (in  $\mathbb{S}$ ) to shared locations (in  $S$ ). Here  $f\{\mathbb{S}\}$  returns the set of target locations that are mapped from locations in  $\mathbb{S}$ .

$f\{\mathbb{S}\} \stackrel{\text{def}}{=} \{l' \mid \exists l. (l \in \mathbb{S}) \wedge f(l) = l'\} \quad f|_{\mathbb{S}} \stackrel{\text{def}}{=} \{(l, f(l)) \mid l \in (\mathbb{S} \cap \text{dom}(f))\}$   
 $\text{wf}(\mu) \text{ iff } \text{injective}(f) \wedge \text{dom}(f) = \mathbb{S} \wedge f\{\mathbb{S}\} \subseteq \mathbb{S} \quad \text{where } \mu = (\mathbb{S}, S, f)$   
 $\text{FPmatch}(\mu, \Delta, \delta) \text{ iff } (\delta.rs \cap S \subseteq f\{\Delta.rs\}) \wedge (\delta.ws \cap S \subseteq f\{\Delta.ws\}) \quad \text{where } \mu = (\mathbb{S}, S, f)$   
 $\text{closed}(\mathbb{S}, \Sigma) \text{ iff } \text{cl}(\mathbb{S}, \Sigma) \subseteq \mathbb{S}$   
 $\text{cl}(\mathbb{S}, \Sigma) \stackrel{\text{def}}{=} \bigcup_k \text{cl}_k(\mathbb{S}, \Sigma)$ , where  $\text{cl}_k(\mathbb{S}, \Sigma)$  is inductively defined:  
 $\text{cl}_0(\mathbb{S}, \Sigma) \stackrel{\text{def}}{=} \mathbb{S} \quad \text{cl}_{k+1}(\mathbb{S}, \Sigma) \stackrel{\text{def}}{=} \{l' \mid \exists l. (l \in \text{cl}_k(\mathbb{S}, \Sigma)) \wedge \Sigma(l) = l'\}$   
 $\text{initM}(\varphi, ge, \Sigma, \sigma) \text{ iff } ge \subseteq \text{dom}(\Sigma) \wedge \text{closed}(\text{dom}(\Sigma), \Sigma) \wedge \text{dom}(\sigma) = \varphi\{\text{dom}(\Sigma)\} \wedge \text{Inv}(\varphi, \Sigma, \sigma)$   
 $\text{HG}(\Delta, \Sigma', \mathbb{F}, \mathbb{S}) \text{ iff } \Delta \subseteq (\mathbb{F} \cup \mathbb{S}) \wedge \text{closed}(\mathbb{S}, \Sigma')$   
 $\text{LG}(\mu, (\delta, \sigma', F), (\Delta, \Sigma')) \text{ iff } \delta \subseteq (F \cup \mu.S) \wedge \text{closed}(\mu.S, \sigma') \wedge \text{FPmatch}(\mu, \Delta, \delta) \wedge \text{Inv}(\mu.f, \Sigma', \sigma')$   
 $\text{R}(\Sigma, \Sigma', \mathbb{F}, \mathbb{S}) \text{ iff } (\Sigma \stackrel{\mathbb{F}}{=} \Sigma') \wedge \text{closed}(\mathbb{S}, \Sigma') \wedge \text{forward}(\Sigma, \Sigma')$   
 $\text{Rely}(\mu, (\Sigma, \Sigma', \mathbb{F}), (\sigma, \sigma', F)) \text{ iff } \text{R}(\Sigma, \Sigma', \mathbb{F}, \mu.S) \wedge \text{R}(\sigma, \sigma', F, \mu.S) \wedge \text{Inv}(\mu.f, \Sigma', \sigma')$   
 $\text{Inv}(f, \Sigma, \sigma) \text{ iff } \forall l, l'. (l \in \text{dom}(\Sigma) \wedge f(l) = l') \implies (l' \in \text{dom}(\sigma) \wedge \Sigma(l) \stackrel{f}{\longmapsto} \sigma(l'))$   
 $v_1 \stackrel{f}{\longmapsto} v_2 \text{ iff } (v_1 \notin \text{Addr}) \wedge (v_1 = v_2) \vee (v_1, v_2 \in \text{Addr} \wedge f(v_1) = v_2)$

Fig. 8. Footprint Matching and Rely/Guarantee Conditions in Our Simulation

Then, given footprints  $\Delta$  and  $\delta$ , we define their consistency with respect to  $\mu$  as  $\text{FPmatch}(\mu, \Delta, \delta)$  in Fig. 8. It says the *shared* locations in  $\delta$  must be contained in  $\Delta$ , modulo the mapping  $\mu.f$ . We only consider the shared locations in  $\mu.S$  because accesses of local memory would *not* cause races.

*Rely/guarantee conditions.* We use rely and guarantee conditions to specify the module interaction protocols. One of the most important protocol is to enforce the view of accessibility of shared and local memory in Fig. 4. More specifically, when the execution of a module switches to external modules, it expects them to keep its local memory (in  $\mathbb{F}$ ) intact. In addition, although the external modules may update the shared memory  $\mathbb{S}$ , they must preserve certain properties of  $\mathbb{S}$ . One important property is that  $\mathbb{S}$  cannot contain memory pointers pointing to local memory cells in any  $\mathbb{F}_i$ <sup>1</sup>. Otherwise a thread  $t_j$  can update the local memory in  $\mathbb{F}_i$  by tracing the pointers in  $S$ . This requirement is formalized as  $\text{closed}(\mathbb{S}, \Sigma)$  in Fig. 8, which says the closure of addresses reachable from  $\mathbb{S}$  must be no bigger than  $\mathbb{S}$ . We encode these requirements in the rely condition, and guarantee conditions are defined correspondingly to ensure the rely is satisfied.

We define the module-local downward simulation  $(sl, ge, \gamma) \preceq_{\varphi} (tl, ge', \pi)$  below, which relates the executions of the source module  $(sl, ge, \gamma)$  and the target module  $(tl, ge', \pi)$ . The injective function  $\varphi$  maps source addresses to the target ones.

**Definition 4** (Module-Local Downward Simulation).

$(sl, ge, \gamma) \preceq_{\varphi} (tl, ge', \pi)$  iff for all  $f, \mathbb{k}, \Sigma, \sigma, \mathbb{F}, F$ , and  $\mu = (\text{dom}(\Sigma), \text{dom}(\sigma), \varphi|_{\text{dom}(\Sigma)})$ , if  $sl.\text{InitCore}(\gamma, f) = \mathbb{k}, \varphi\{ge\} = ge'$ ,  $\text{initM}(\varphi, ge, \Sigma, \sigma)$ , and  $\mathbb{F} \cap \text{dom}(\Sigma) = F \cap \text{dom}(\sigma) = \emptyset$ , then there exist  $\kappa$  and  $i \in \text{index}$  such that  $tl.\text{InitCore}(\pi, f) = \kappa$ , and

$$(\mathbb{F}, (\mathbb{k}, \Sigma), \text{emp}) \preceq_{\mu}^i (F, (\kappa, \sigma), \text{emp}),$$

where  $(\mathbb{F}, (\mathbb{k}, \Sigma), \Delta) \preceq_{\mu}^i (F, (\kappa, \sigma), \delta)$ , is defined in Def. 5 below.

It says that if we take any function entry  $f$  and the corresponding initial core states  $\mathbb{k}$  and  $\kappa$  at the source and the target respectively, then with any states  $(\Sigma$  and  $\sigma)$  and free lists  $(\mathbb{F}$  and  $F)$  satisfying

<sup>1</sup> This means we do not allow escape of pointers pointing to stack variables.

some initial constraints, we can establish the simulation  $(\mathbb{F}, (\mathbb{k}, \Sigma), \text{emp}) \preceq_{\mu}^i (F, (\kappa, \sigma), \text{emp})$ , relating the local module configurations, which is defined in Def. 5 and is explained below. Here we require the initial states  $\Sigma$  and  $\sigma$  be related through  $\text{initM}$  defined in Fig. 8. The last condition in  $\text{initM}$  is the invariant  $\text{Inv}$  defined at the bottom of Fig. 8. It says that the location  $f(l)$  must be in  $\text{dom}(\sigma)$  if  $l$  is in  $\text{dom}(\Sigma)$ , and the memory contents  $\Sigma(l)$  and  $\sigma(f(l))$  must be equal modulo  $\mu.f$ .

**Definition 5.** We define  $(\mathbb{F}, (\mathbb{k}, \Sigma), \Delta_0) \preceq_{\mu}^i (F, (\kappa, \sigma), \delta_0)$  as the largest relation such that, whenever  $(\mathbb{F}, (\mathbb{k}, \Sigma), \Delta_0) \preceq_{\mu}^i (F, (\kappa, \sigma), \delta_0)$ , then the following are true:

- (1) for all  $\mathbb{k}'$ ,  $\Sigma'$  and  $\Delta$ , if  $\mathbb{F} \vdash (\mathbb{k}, \Sigma) \xrightarrow{\tau}_{\Delta} (\mathbb{k}', \Sigma')$  and  $(\Delta_0 \cup \Delta) \subseteq (\mathbb{F} \cup \mu.\mathbb{S})$ , then one of the following holds:
  - (a)  $\exists j < i$ .  $(\mathbb{F}, (\mathbb{k}', \Sigma'), \Delta_0 \cup \Delta) \preceq_{\mu}^j (F, (\kappa, \sigma), \delta_0)$ , or
  - (b) there exist  $\kappa'$ ,  $\sigma'$ ,  $\delta$  and  $j$  such that the following are true:
    - (i)  $F \vdash (\kappa, \sigma) \xrightarrow{\tau}_{\delta}^+ (\kappa', \sigma')$ ;
    - (ii)  $(\delta_0 \cup \delta) \subseteq (F \cup \mu.S)$  and  $\text{FPmatch}(\mu, \Delta_0 \cup \Delta, \delta_0 \cup \delta)$ ; and
    - (iii)  $(\mathbb{F}, (\mathbb{k}', \Sigma'), \Delta_0 \cup \Delta) \preceq_{\mu}^j (F, (\kappa', \sigma'), \delta_0 \cup \delta)$ .
- (2) for all  $\mathbb{k}'$  and  $\iota$ , if  $\mathbb{F} \vdash (\mathbb{k}, \Sigma) \xrightarrow[\text{emp}]{\iota} (\mathbb{k}', \Sigma)$ ,  $\iota \neq \tau$ , and  $\text{HG}(\Delta_0, \Sigma, \mathbb{F}, \mu.\mathbb{S})$ , then there exist  $\kappa'$ ,  $\delta$ ,  $\sigma'$  and  $\kappa''$  such that the following are true:
  - (a)  $F \vdash (\kappa, \sigma) \xrightarrow{\tau}_{\delta}^* (\kappa', \sigma')$ , and  $F \vdash (\kappa', \sigma') \xrightarrow[\text{emp}]{\iota} (\kappa'', \sigma')$ , and
  - (b)  $\text{LG}(\mu, (\delta_0 \cup \delta, \sigma', F), (\Delta_0, \Sigma))$ , and
  - (c) for all  $\sigma''$  and  $\Sigma'$ , if  $\text{Rely}(\mu, (\Sigma, \Sigma', \mathbb{F}), (\sigma', \sigma'', F))$ , then there exists  $j$  such that  $(\mathbb{F}, (\mathbb{k}', \Sigma'), \text{emp}) \preceq_{\mu}^j (F, (\kappa'', \sigma''), \text{emp})$ .

The simulation  $(\mathbb{F}, (\mathbb{k}, \Sigma), \Delta_0) \preceq_{\mu}^i (F, (\kappa, \sigma), \delta_0)$  carries  $\Delta_0$  and  $\delta_0$ , the footprints accumulated until now at the source and the target, respectively. The definition follows the diagram in Fig. 1(d). For every  $\tau$ -step in the source (case 1), if the newly generated footprints and the accumulated  $\Delta_0$  are *in scope* (i.e. every location must either be from the freelist space  $\mathbb{F}$  of current thread, or from the shared memory  $\mu.\mathbb{S}$ ), then the step should correspond to zero or multiple  $\tau$ -steps in the target, and the simulation holds over the resulting states with the accumulated footprints and a new index  $j$ . Here we use  $\Delta$  as a shorthand for  $(\Delta.rs \cup \Delta.ws)$ . If the source step corresponds to zero target step (case 1-a), the index  $j$  must be strictly smaller than  $i$ . Here the indices  $i$  and  $j$  belong to a well-founded set  $\text{index}$  that has no infinite decreasing sequences. The use of a smaller index  $j$  in this case ensures non-terminating source module can only be simulated by non-terminating target.

If the source step corresponds to at least one target steps (case 1-b), the index  $j$  can be arbitrary. In this case we require the footprints at the target are also in scope, and they must be consistent with the source level footprints (see our explanation of  $\text{FPmatch}$  before). The accumulation of footprints allows us to establish  $\text{FPmatch}$  for compiler optimizations that reorder the instructions.

At the switch points when the source generates a non-silent message  $\iota$  (case 2), if the footprints and states satisfy  $\text{HG}$ , the target must be able to generate the same  $\iota$  (after zero or multiple silent steps), and the accumulated footprints and the state satisfy  $\text{LG}$ . As defined in Fig. 8, both  $\text{HG}$  and  $\text{LG}$  require the footprints are in scope and the shared memory is closed.  $\text{LG}$  additionally requires the footprints at the target and the source satisfy  $\text{FPmatch}$ , and the states satisfy  $\text{Inv}$ .

One may wonder if it is possible to check  $\text{FPmatch}$  at the switch points only. However, executions of non-terminating modules would never reach a switch point. That is why we have to also require  $\text{FPmatch}$  during internal  $\tau$ -steps (case 1).

At the switch points we also need to consider the interaction with the environment (i.e. other modules or threads). For any environment steps at the source and the target, if they satisfy the  $\text{rely}$

condition, then the simulation holds over the new states, with some index  $j$  and *empty* footprints — Since the effects of the current thread have been made visible to the environments at the switch point, we can clear the accumulated footprints. The rely condition is defined in Fig. 8. As explained before, it requires the local memory is untouched, the shared memory is closed, and the domain of states is not reduced (see the definition of forward in Fig. 5). Also it requires that the invariant  $\text{Inv}$  be preserved over the new states at the source and the target.

Note that each case in Def. 5 has prerequisite about the source level footprints (e.g. the footprints are in scope or satisfy HG). On the one hand, this makes the simulation weak and easy to prove. We do not need to prove these requirements indeed hold for each compilation phase since they are premises. On the other hand, after we apply transitivity of the simulation and prove that the target generated after the multi-phase compilation simulates the source, we need to additionally prove that these requirements indeed hold at the source level, to make the simulation meaningful instead of being vacuously true. Following the approach by Stewart et al. [2015], we formalize these requirements separately as *ReachClose* (Def. 7) in the next subsection.

Our simulation is transitive. One can decompose the whole compiler correctness proofs into proofs for individual compilation phases.

**Lemma 6** (Transitivity).  $\forall sl, sl', tl, \gamma, \gamma', \pi.$

if  $(sl, ge, \gamma) \preceq_{\varphi} (sl', ge', \gamma')$  and  $(sl', ge', \gamma') \preceq_{\varphi'} (tl, ge'', \pi)$ , then  $(sl, ge, \gamma) \preceq_{\varphi' \circ \varphi} (tl, ge'', \pi)$ .

## 5.2 Compositionality and the Non-Preemptive Global Simulation

The whole program *downward* simulation  $\hat{\mathbb{P}} \preceq_{\varphi} \hat{P}$  relates the execution of the whole source program  $\mathbb{P}$  and target program  $P$ . It serves as an intermediate result for proving event trace refinement and data-race-freedom preservation, as shown in Fig. 2.

Due to space limit, we omit the definition of the whole program simulation, which is similar to the module local simulation, except the case for environment interference (rely steps), which is unnecessary for whole program simulation. Every source step should correspond to multiple target steps. As in module local simulation, we always require the target footprints matches those of the source, defined as *FPmatch*. Also the footprint should always be in scope. As a special requirement for the whole program simulation, we always require the source and the target do lock-step context switch and they always switch to the same thread.

*Compositionality.* Following Liang et al. [2012], the module-local simulations could be composed to derive the whole-program simulation by proving the Rely condition of a module is guaranteed by other modules' guarantee conditions HG and LG. However, as explained at the end of Sec. 5.1, in our module-local simulation, LG is established if the source module satisfies HG. We require that HG always hold during the execution of the source, and formulate this requirement as *ReachClose* in Def. 7. It is a simplified version of the *reach-close* concept by Stewart et al. [2015] (simplified because we do not allow the leak of local stack pointers into the shared memory).

**Definition 7** (Reach Closed Module).  $\text{ReachClose}(sl, ge, \gamma)$  iff , for all  $f, \mathbb{k}, \Sigma, \mathbb{F}$  and  $\mathbb{S}$ , if  $sl.\text{InitCore}(\gamma, f) = \mathbb{k}, ge \subseteq \mathbb{S} = \text{dom}(\Sigma), \mathbb{F} \cap \mathbb{S} = \emptyset$ , and  $\text{closed}(\mathbb{S}, \Sigma)$ , then  $\text{RC}(\mathbb{F}, \mathbb{S}, (\mathbb{k}, \Sigma))$ .

Here RC is defined as the largest relation such that, whenever  $\text{RC}(\mathbb{F}, \mathbb{S}, (\mathbb{k}, \Sigma))$ , then for all  $\Sigma'$  such that  $\text{R}(\Sigma, \Sigma', \mathbb{F}, \mathbb{S})$ , and for all  $\mathbb{k}', \Sigma', \Sigma'', \iota$  and  $\Delta$  such that  $\mathbb{F} \vdash (\mathbb{k}, \Sigma') \xrightarrow[\Delta]{\iota} (\mathbb{k}', \Sigma'')$ , we have  $\text{HG}(\Delta, \Sigma'', \mathbb{F}, \mathbb{S})$ , and  $\text{RC}(\mathbb{F}, \mathbb{S}, (\mathbb{k}', \Sigma''))$ .

The relation  $\text{RC}(\mathbb{F}, \mathbb{S}, (\mathbb{k}, \Sigma))$  essentially says during every step of the execution of  $(\mathbb{k}, \Sigma)$ , HG always holds over the resulting footprints  $\Delta$  and states, even with possible interference from the environment, as long as the environment steps satisfy the rely condition R defined in Fig. 8.

Assuming all source modules are *ReachClose*, we can prove the Compositionality Lemma.



**Lemma 8** (Compositionality, ⑤ in Fig. 2). For any  $f_1, \dots, f_n, \varphi, \Gamma = \{(sl_1, ge_1 \gamma_1), \dots, (sl_m, ge_m \gamma_m)\}$ ,  
 $\Pi = \{(tl_1, ge'_1, \pi_1), \dots, (tl_m, ge'_m, \pi_m)\}$  if

$$\forall i \in \{1, \dots, m\}. \text{wd}(sl_i) \wedge \text{wd}(tl_i) \wedge \text{ReachClose}(sl_i, ge_i, \gamma_i) \wedge (sl_i, ge_i, \gamma_i) \preceq_\varphi (tl_i, ge'_i, \pi_i),$$

then  $\text{let } \Gamma \text{ in } f_1 \mid \dots \mid f_n \preceq_\varphi \text{let } \Pi \text{ in } f_1 \mid \dots \mid f_n$ .

### 5.3 Flip of the Non-Preemptive Global Simulation

As explained in Sec. 2.2, with determinism of the target modules, one is able to flip the downward simulation to derive upward simulation for safe source programs. Then the upward event trace refinement follows the upward simulation.

We define the whole program upward simulation as  $\hat{P} \leq_\varphi \hat{\mathbb{P}}$ . Due to space limit, we omit its concrete definition here. It is similar to the downward simulation  $\hat{\mathbb{P}} \preceq_\varphi \hat{P}$ , with the positions of source and target swapped. Note that we do *not* flip the FPmatch condition, since we always require footprint of target program being a refinement of the footprint of the source program, in order to prove DRF of the target program. Correspondingly, the address mapping  $\varphi$  is *not* flipped either.

**Definition 9** (Deterministic Languages).  $\text{det}(tl)$  iff, for all configuration  $\phi$  in  $tl$  (see the definition of  $\phi$  in Fig. 3), and for all  $F$ ,  $F \vdash \phi \xrightarrow{\delta_1} \phi_1 \wedge F \vdash \phi \xrightarrow{\delta_2} \phi_2 \implies \phi_1 = \phi_2 \wedge \iota_1 = \iota_2 \wedge \delta_1 = \delta_2$ .

With determinism of the target module languages, we are able to prove Lemma 10.

**Lemma 10** (Flip, ④ in Fig. 2). For any  $f_1, \dots, f_n, ge, \varphi, \Gamma = \{(sl_1, ge_1, \gamma_1), \dots, (sl_m, ge_m, \gamma_m)\}$ ,  
 $\Pi = \{(tl_1, ge'_1, \pi_1), \dots, (tl_m, ge'_m, \pi_m)\}$ , if  $\forall i. \text{det}(tl_i)$ , and

$$\begin{aligned} & \text{let } \Gamma \text{ in } f_1 \mid \dots \mid f_m \preceq_\varphi \text{let } \Pi \text{ in } f_1 \mid \dots \mid f_m, \\ \text{then} & \quad \text{let } \Pi \text{ in } f_1 \mid \dots \mid f_m \leq_\varphi \text{let } \Gamma \text{ in } f_1 \mid \dots \mid f_m. \end{aligned}$$

*Soundness.* The non-preemptive global simulation ensures the refinement. Before presenting the soundness lemma, we first lift the refinement  $(P, \Sigma) \sqsubseteq (\mathbb{P}, \sigma)$  (see Def. 2) to  $P \sqsubseteq_\varphi \mathbb{P}$  as follows. Similarly,  $(\hat{P}, \Sigma) \sqsubseteq (\hat{\mathbb{P}}, \sigma)$  is lifted to  $\hat{P} \sqsubseteq_\varphi \hat{\mathbb{P}}$ .

$$\begin{aligned} P \sqsubseteq_\varphi \mathbb{P} \text{ iff } \forall \Sigma, \sigma. \text{initM}(\varphi, \text{GE}(\mathbb{P}, \Gamma), \Sigma, \sigma) \implies (P, \sigma) \sqsubseteq (\mathbb{P}, \Sigma) \\ \text{where } \text{GE}(\{(tl_1, ge_1, \pi_1), \dots, (tl_m, ge_m, \pi_m)\}) \stackrel{\text{def}}{=} \bigcup_{i=1}^m ge_i \end{aligned}$$

**Lemma 11** (Soundness, ③ in Fig. 2). If  $\hat{P} \leq_\varphi \hat{\mathbb{P}}$ , then  $\hat{P} \sqsubseteq_\varphi \hat{\mathbb{P}}$ .

## 6 DATA-RACE-FREEDOM

Informally, a data race occurs when two threads concurrently access the same memory location and at least one of the accesses is a write. A program is DRF if it never generates data races in all possible executions. Below we first define the conflict of footprints.

$$\begin{aligned} \delta_1 \frown \delta_2 \quad \text{iff} \quad (\delta_1.\text{ws} \cap \delta_2 \neq \emptyset) \vee (\delta_2.\text{ws} \cap \delta_1 \neq \emptyset) \\ (\delta_1, d_1) \frown (\delta_2, d_2) \quad \text{iff} \quad (\delta_1 \frown \delta_2) \wedge (d_1 = 0 \vee d_2 = 0) \end{aligned}$$

Two footprints  $\delta_1$  and  $\delta_2$  are conflicting, denoted as  $\delta_1 \frown \delta_2$ , if the write set  $\text{ws}$  of one of them overlaps with the read set or the write set of the other. Recall that, when used as a set,  $\delta$  represents  $\delta.\text{rs} \cup \delta.\text{ws}$ . Since we do *not* treat accesses of the same memory location inside atomic blocks as a race, we instrument a footprint  $\delta$  with the atomic bit  $d$  to record whether the footprint is generated inside an atomic block ( $d = 1$ ) or not ( $d = 0$ ). Two instrumented footprints  $(\delta_1, d_1)$  and  $(\delta_2, d_2)$  are conflicting if  $\delta_1$  and  $\delta_2$  are conflicting and at least one of  $d_1$  and  $d_2$  is 0.

We formulate data races in Fig. 9(a) for preemptive semantics. A program  $P$  with initial memory state  $\sigma$  is racy  $((P, \sigma) \implies \text{Race})$  if its execution reaches a program configuration  $W'$  that steps to

$$\begin{array}{c}
834 \\
835 \\
836 \\
837 \\
838 \\
839 \\
840 \\
841 \\
842 \\
843 \\
844 \\
845 \\
846 \\
847 \\
848 \\
849 \\
850 \\
851 \\
852 \\
853 \\
854 \\
855 \\
856 \\
857 \\
858 \\
859 \\
860 \\
861 \\
862 \\
863 \\
864 \\
865 \\
866 \\
867 \\
868 \\
869 \\
870 \\
871 \\
872 \\
873 \\
874 \\
875 \\
876 \\
877 \\
878 \\
879 \\
880 \\
881 \\
882
\end{array}$$

$$\begin{array}{c}
\frac{(P, \sigma) \xrightarrow{\text{load}} W \quad W \Rightarrow^* W' \quad W' \Longrightarrow \text{Race}}{(P, \sigma) \Longrightarrow \text{Race}} \\
\frac{t_1 \neq t_2 \quad (\delta_1, d_1) \frown (\delta_2, d_2) \quad \text{predict}(W, t_1, (\delta_1, d_1)) \quad \text{predict}(W, t_2, (\delta_2, d_2))}{W \Longrightarrow \text{Race}} \text{ Race} \\
\frac{W = (T, \_, 0, \sigma) \quad T(t) = (F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\delta]{\tau} (\kappa', \sigma')}{\text{predict}(W, t, (\delta, 0))} \text{ Predict-0} \\
\frac{W = (T, \_, 0, \sigma) \quad T(t) = (F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\text{emp}]{\text{EntAtom}} (\kappa', \sigma) \quad F \vdash (\kappa', \sigma) \xrightarrow[\delta]{\tau} (\kappa'', \sigma'')}{\text{predict}(W, t, (\delta, 1))} \text{ Predict-1} \\
\text{(a) Data Race in Preemptive Semantics} \\
\frac{(\hat{P}, \sigma) \xrightarrow{\text{load}} \hat{W} \quad \hat{W} \Rightarrow^* \hat{W}' \quad \hat{W}' \Longrightarrow \text{Race}}{(\hat{P}, \sigma) \Longrightarrow \text{Race}} \\
\frac{(\hat{P}, \sigma) \xrightarrow{\text{load}} \hat{W} \quad t_1 \neq t_2 \quad (\delta_1, d_1) \frown (\delta_2, d_2) \quad \text{NPpredict}(\hat{W}, t_1, (\delta_1, d_1)) \quad \text{NPpredict}(\hat{W}, t_2, (\delta_2, d_2))}{(\hat{P}, \sigma) \Longrightarrow \text{Race}} \\
\frac{\hat{W} \xrightarrow[\text{emp}]{o} \hat{W}' \quad o \neq \tau \quad t_1 \neq t_2 \quad (\delta_1, d_1) \frown (\delta_2, d_2) \quad \text{NPpredict}(\hat{W}', t_1, (\delta_1, d_1)) \quad \text{NPpredict}(\hat{W}', t_2, (\delta_2, d_2))}{\hat{W} \Longrightarrow \text{Race}} \text{ Race}_{\text{np}} \\
\frac{\hat{W} = (T, \_, d, \sigma) \quad T(t) = (F, \kappa) \quad F \vdash (\kappa, \sigma) \xrightarrow[\delta]{\tau} (\kappa', \sigma') \quad d(t) = d}{\text{NPpredict}(\hat{W}, t, (\delta, d))} \text{ Predict}_{\text{np}} \\
\text{(b) Data Race in Non-Preemptive Semantics}
\end{array}$$

Fig. 9. Predictive Semantics for Defining Race

Race ( $W' \Longrightarrow \text{Race}$ ). In the Race rule,  $W$  steps to Race if there are conflicting footprints of two threads predicted from the current configuration ( $\text{predict}(W, t, (\delta, d))$ ). The predictions are only performed at the switch points, i.e., the program points outside atomic blocks where  $d = 0$ , since the intermediate states inside atomic blocks are not visible to other threads. Footprints of steps before reaching the next switch point are able to be predicted using the rules Predict-0 or Predict-1. Predict-0 is used for predicting footprint of the next  $\tau$ -step of thread  $t$ , where there is a switch point (in preemptive semantics) right after the step. Predict-1 applies if thread  $t$  enters an atomic block, the footprint generated by any number of steps inside the atomic block can be predicted. Note that we do not insist on predicting the footprint generated by the execution of the *whole* atomic block because the code inside the atomic block may be nonterminating.

A program  $P$  with an initial state  $\sigma$  is DRF if it never steps to Race. A program  $P$  is DRF if for any proper initial state  $\sigma$ ,  $\text{DRF}(P, \sigma)$ :

$$\begin{array}{c}
877 \\
878 \\
879
\end{array}$$

$$\begin{array}{l}
\text{DRF}(P, \sigma) \text{ iff } \neg((P, \sigma) \Longrightarrow \text{Race}) \\
\text{DRF}(P) \text{ iff } \forall \sigma. (\text{GE}(P, \text{II}) \subseteq \text{locs}(\sigma)) \wedge (\text{cl}(\text{dom}(\sigma), \sigma) = \text{dom}(\sigma)) \implies \text{DRF}(P, \sigma)
\end{array}$$

Data races in the non-preemptive semantics ( $(\hat{P}, \sigma) \Longrightarrow \text{Race}$ ) are defined in Fig. 9(b). Similar to the definition for the preemptive semantics, a non-preemptive program configuration  $\hat{W}$  steps to

883 Race ( $\widehat{W} \stackrel{\circ}{\Longrightarrow} \text{Race}$ ) if two threads are predicted having conflicting footprints, as shown in rule  
 884  $\text{Race}_{\text{np}}$ . Predictions are made only at the switch points of non-preemptive semantics ( $\widehat{W} \stackrel{\circ}{\underset{\text{emp}}{\Longrightarrow}} \widehat{W}'$   
 885 where  $\circ \neq \tau$ ), i.e., program points at atomic block boundaries, after an observable event, or after  
 886 thread termination. The prediction rule  $\text{Predict}_{\text{np}}$  is a unified version of its counterpart  $\text{Predict-0}$   
 887 and  $\text{Predict-1}$ , where  $\text{cl}$  indicates whether the predicted steps are inside an atomic block. Similar to  
 888 the  $\text{Predict-1}$  rule, we do *not* insist on predicting the footprint generated by the execution reaching  
 889 the next switch point, because the code segments between switch points could be nonterminating.  
 890 In addition to at the switch points, we also need to be able to perform a prediction at the initial  
 891 state as well (the second rule in Fig. 9), because the first executing thread is non-deterministically  
 892 picked at the beginning, which has similar effect as thread switching.

893 A program  $\hat{P}$  with the initial state  $\sigma$  is NPDRF if it never steps to  $\text{Race}$ , and a program  $\hat{P}$  is  
 894 NPDRF if, for any proper initial state  $\sigma$ ,  $\text{NPDRF}(\hat{P}, \sigma)$ :

$$895 \quad \text{NPDRF}(\hat{P}, \sigma) \text{ iff } \neg((\hat{P}, \sigma) \stackrel{\circ}{\Longrightarrow} \text{Race})$$

$$896 \quad \text{NPDRF}(P) \text{ iff } \forall \sigma. (\text{GE}(P, \Pi) \subseteq \text{locs}(\sigma)) \wedge (\text{cl}(\text{dom}(\sigma), \sigma) = \text{dom}(\sigma)) \implies \text{NPDRF}(P, \sigma)$$

897 Note that defining NPDRF is not for studying the absence of data races in the non-preemptive  
 898 semantics (which is probably not very interesting since the execution is non-preemptive anyway).  
 899 Rather, it is intended to serve as an equivalent notion of DRF but formulated in the non-preemptive  
 900 semantics. The following lemma shows the equivalence.

901 **Lemma 12** (Equivalence between DRF and NPDRF, ⑥ and ⑧ in Fig. 2).

902 For any  $f_1, \dots, f_m, \sigma, \Pi = \{(tl_1, \pi_1), \dots, (tl_m, \pi_m)\}$  such that  $\forall i. \text{wd}(tl_i)$ ,

$$903 \quad \text{DRF}(\mathbf{let} \ \Pi \ \mathbf{in} \ f_1 \parallel \dots \parallel f_m, \sigma) \iff \text{NPDRF}(\mathbf{let} \ \Pi \ \mathbf{in} \ f_1 \mid \dots \mid f_m, \sigma).$$

904 As mentioned in Sec. 2, we need the compilation to preserve DRF of the source, which should be  
 905 ensured by our upward whole program simulation  $\hat{P} \leq_{\varphi} \hat{\mathbb{P}}$ . The following Lemma 13 shows the  
 906 simulation preserves NPDRF of the source. Together with Lemma 12, we know it preserves DRF.

907 **Lemma 13** (NPDRF Preservation, ⑦ in Fig. 2).

908 For any  $\hat{\mathbb{P}}, \hat{P}, \varphi, \Sigma$  and  $\sigma$ , if  $\hat{P} \leq_{\varphi} \hat{\mathbb{P}}$ ,  $\text{initM}(\varphi, \text{GE}(\hat{\mathbb{P}}, \Gamma), \Sigma, \sigma)$ , and  $\text{NPDRF}(\hat{\mathbb{P}}, \Sigma)$ , then  $\text{NPDRF}(\hat{P}, \sigma)$ .

909 The lemma below shows the semantics equivalence for DRF programs.

910 **Lemma 14** (Equivalence between Preemptive and Non-Preemptive semantics, ① and ② in Fig. 2).

911 For any  $\Pi, f_1, \dots, f_m, \sigma$ , if  $\text{DRF}(\mathbf{let} \ \Pi \ \mathbf{in} \ f_1 \parallel \dots \parallel f_m, \sigma)$ , then

$$912 \quad (\mathbf{let} \ \Pi \ \mathbf{in} \ f_1 \mid \dots \mid f_m, \sigma) \approx (\mathbf{let} \ \Pi \ \mathbf{in} \ f_1 \parallel \dots \parallel f_m, \sigma).$$

## 913 7 THE FINAL THEOREM

914 Putting all the previous results together, we are able to prove our final theorem, i.e., certified  
 915 sequential compositional compilers could correctly compile data-race-free concurrent programs by  
 916 compiling each module separately.

917 Before presenting the theorem, we first model a sequential compiler  $\text{SeqComp}$  as a code trans-  
 918 formation function  $\text{CodeT}$  with a data transformation function  $\varphi$ . Here  $\varphi$  maps the addresses  
 919 in the global environments  $ge$ . It may not be an identity function when the compiler performs  
 920 optimizations on global environments, such as eliminating unused global variables.

$$921 \quad \text{SeqComp} ::= (\text{CodeT}, \varphi) \quad \text{CodeT} \in \text{Module} \rightarrow \text{Module} \quad \varphi \in \text{Addr} \rightarrow \text{Addr}$$

922 As the key proof obligation, we need to verify that each  $\text{SeqComp}$  satisfies  $\text{Correct}$ . We define  
 923  $\text{Correct}$  as follows using our footprint-preserving module-local simulation.

924

**Definition 15** (Sequential Compiler Correctness).  $\text{Correct}(\text{SeqComp}, sl, tl)$  iff

$$\forall \gamma, \pi, ge, ge'. \text{SeqComp.CodeT}(\gamma) = \pi \wedge \text{SeqComp.}\varphi\{ge\} = ge' \implies (sl, ge, \gamma) \preceq_{\text{SeqComp.}\varphi} (tl, ge', \pi).$$

The desired correctness property  $\text{GCorrect}$  of concurrent program compilation is the semantics preservation of whole programs, i.e., every target concurrent program is an event-trace refinement of the source. We formulate  $\text{GCorrect}$  in Def. 16. Here we require all the sequential compilers to agree on the transformation  $\varphi$  of global environments (see Def. 16(1)).

**Definition 16** (Concurrent Compiler Correctness).

$\text{GCorrect}((\text{SeqComp}_1, sl_1, tl_1), \dots, (\text{SeqComp}_m, sl_m, tl_m))$  iff

for any  $f_1, \dots, f_n$ ,  $\Gamma = \{(sl_1, ge_1, \gamma_1), \dots, (sl_m, ge_m, \gamma_m)\}$ ,  $\Pi = \{(tl_1, ge'_1, \pi_1), \dots, (tl_m, ge'_m, \pi_m)\}$ ,  $\varphi$ , if

- (1)  $\forall i \in \{1, \dots, m\}. (\text{SeqComp}_i.\text{CodeT}(\gamma_i) = \pi_i) \wedge (\text{SeqComp}_i.\varphi = \varphi) \wedge \text{injective}(\varphi)$ ,
- (2)  $\text{DRF}(\mathbf{let} \Gamma \mathbf{in} f_1 \parallel \dots \parallel f_n)$ , and  $\text{Safe}(\mathbf{let} \Gamma \mathbf{in} f_1 \parallel \dots \parallel f_n)$ ,
- (3)  $\forall i \in \{1, \dots, m\}. \text{ReachClose}(sl_i, ge_i, \gamma_i)$ ,

then  $\mathbf{let} \Gamma \mathbf{in} f_1 \parallel \dots \parallel f_n \sqsupseteq_{\varphi} \mathbf{let} \Pi \mathbf{in} f_1 \parallel \dots \parallel f_n$ .

Our final theorem is then formulated as Thm. 17. It says if a set of sequential compilers are certified to satisfy our correctness obligation  $\text{Correct}$ , the source and target languages  $sl_i$  and  $tl_i$  are well-defined, and the target languages are deterministic, then the sequential compilers as a whole is  $\text{GCorrect}$  for compiling concurrent programs. The proof simply applies the lemmas that correspond to ①-⑧ in Fig. 2.

**Theorem 17** (Final Theorem).

For any  $\text{SeqComp}_1, \dots, \text{SeqComp}_m, sl_1, \dots, sl_m, tl_1, \dots, tl_m$  such that for any  $i \in \{1, \dots, m\}$  we have  $\text{wd}(sl_i)$ ,  $\text{wd}(tl_i)$ ,  $\text{det}(tl_i)$ , and  $\text{Correct}(\text{SeqComp}_i, sl_i, tl_i)$ , then

$$\text{GCorrect}((\text{SeqComp}_1, sl_1, tl_1), \dots, (\text{SeqComp}_m, sl_m, tl_m)).$$

## 8 FRAMEWORK INSTANTIATION AND COMPCERT BACKEND VERIFICATION

We apply our compiler verification framework to prove the correctness of  $\text{CompCert-3.0.1}$  x86 backend [Leroy 2009a; CompCert Developers 2017] for compositional compilation of DRF programs. To demonstrate the support of cross-language inter-module interaction, we provide synchronization modules implemented in x86 assembly and link them with the modules compiled by  $\text{CompCert}$ .

### 8.1 Language Instantiations

Below we instantiate our abstract languages as  $\text{Cminor}$  and x86 assembly.  $\text{Cminor}$  is a low-level imperative language, and serves as the source language of several compiler backend of the  $\text{CompCert}$  compiler. Although it is a sequential language, we can write concurrent  $\text{Cminor}$  programs by parallel compositions of sequential  $\text{Cminor}$  threads. Inter-thread synchronization can be achieved through the external call mechanism to call functions outside  $\text{Cminor}$  modules.

We give a tiny example in Fig. 10, where the synchronization functions  $\text{sc\_cas}$  and  $\text{sc\_store}$  are implemented in a separate x86 module  $\pi_{x86}$  (see Fig. 10(b)). Their signatures have been exported to  $\text{Cminor}$  in Fig. 10(a). These functions are similar to the C-11 SC-atomic primitives [Batty et al. 2011]. The function  $\text{sc\_cas}$  tests if the value in the destination memory location (1st argument) is the same with the expected value (2nd argument). If true it stores the new value (3rd argument) into the destination address. The  $\text{sc\_store}$  function simply stores the second argument to the destination memory location (1st argument). Their implementation in Fig. 10(b) utilizes the lock-prefixed instructions. With these external x86 functions, the  $\text{Cminor}$  program implements a simple test-and-set lock and a counter  $\text{inc}$  that increments the shared variable  $x$  in the lock-protected

```

981         void    sc_store (int32_t *, int32_t)
982         int     sc_cas   (int32_t *, int32_t, int32_t)
983                 (a) sc_atomic signatures
984
985         sc_store:                                int32_t x = 0, l = 0;
986             movl 4(%esp), %eax                    int lock(){ return sc_cas(&l, 0, 1); }
987             movl 8(%esp), %ecx
988             lock xchgl %ecx, (%eax)                void unlock(){ sc_store(&l, 0); }
989             retl
990
991         sc_cas:                                    void inc(){
992             movl 4(%esp), %edx                      int32_t tmp;
993             movl 8(%esp), %eax                      while(!lock());
994             movl 12(%esp), %ecx                     tmp = x; x ++;
995             lock cmpxchgl %ecx, (%edx)              unlock();
996             sete %al                                print(tmp);
997             retl                                    }
998
999         (b) sc_atomic x86 implementation, module  $\pi_{x86}$                 (c) example Cminor module  $\gamma_C$ 

```

Fig. 10. An example program with implementation of C11-like atomics

critical region and prints the old value of  $x$ . Here we present the Cminor program in C syntax. An example whole program  $\mathbb{P}$  is **let**  $\{\pi_{x86}, \gamma_C\}$  **in**  $\text{inc}() \parallel \text{inc}()$ .

*The Cminor Language.* The selected parts of our Cminor instantiation are represented in Fig. 11(a). To instantiate our abstract language of Fig. 3, we instantiate the *Module* with the same syntax as the CompCert Cminor language. A core state  $\kappa$  is a pair of a local state  $c$  and an index  $N$  of type  $\mathbb{N}$  indicating the position of the next block in the freelist  $F$  to be allocated. Here  $F$  is defined as a sequence of block numbers. The local state  $c$  is instantiated the same way as the Cminor interaction semantics in Compositional CompCert [Stewart et al. 2015]. The  $\text{InitCore}$  function is also adapted from Compositional CompCert, with the index  $N$  initialized to 0. We omit their definitions here.

The instantiation of the local transition semantics instruments Compositional CompCert’s interaction semantics with footprints, which is determined by the memory locations (base-offset pairs, following the CompCert block-based memory model) accessed in each step. For instance, the allocation (e.g., for allocating stack frames) in our semantics takes the next available block number ( $b = F(N)$ ) as the base address and then increments  $N$  so that it points to the next available block number in  $F$ . The memory  $\sigma$  is extended with the newly allocated memory cells, resulting in  $\sigma\{(b, 0) \rightsquigarrow \text{undef}, \dots, (b, n - 1) \rightsquigarrow \text{undef}\}$ , where  $n$  is the size of allocated block. The write set of the footprint includes all the addresses with base address  $b$ , and the read set is empty.

As we mentioned at the beginning of this subsection, Cminor does not support synchronization operations within the language, therefore the operational semantics of Cminor does not generate steps with labels  $\text{EntAtom}$  or  $\text{ExtAtom}$ .

*x86 Assembly with lock-prefixed Instructions.* Figure 11(b) presents the selected syntax and key semantics components of x86 assembly with lock-prefixed instructions. In addition to the instructions that are already implemented in CompCert, we introduce the lock-prefixed instructions ( $\text{lockxchg}$ ,  $\text{lockxadd}$ ,  $\text{lockcmpxchg}$ ) for synchronization.

The lock prefix asserts a  $\text{lock\#}$  signal (or something alike in later IA-32 implementations), which guarantees exclusive memory access while executing the accompanying instruction. The behavior of lock-prefixed instructions could be modeled by our abstract language with  $\text{EntAtom}$  and  $\text{ExtAtom}$  events. To implement the semantics of lock-prefixed instructions, we introduce a

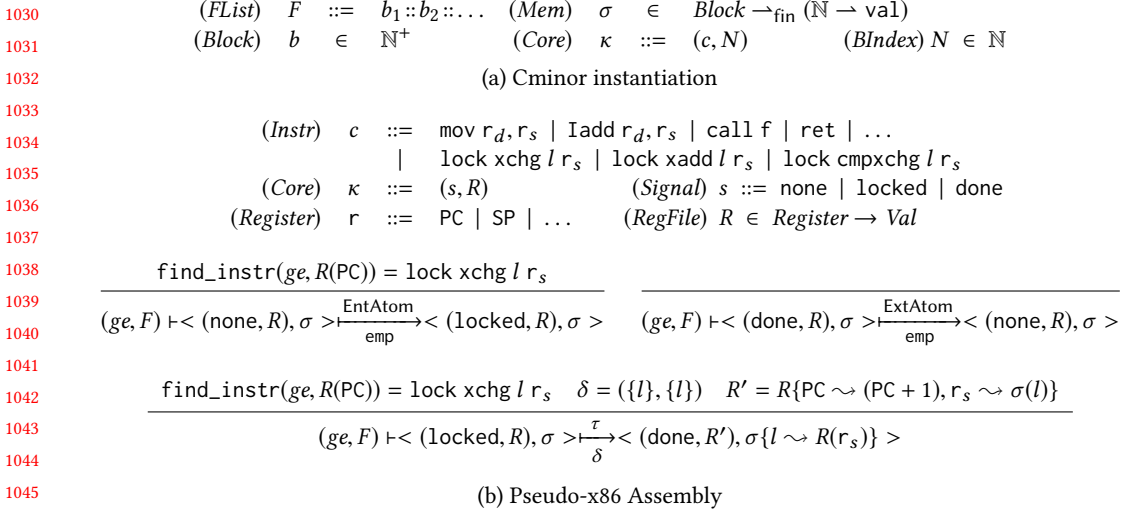


Fig. 11. Language instantiations

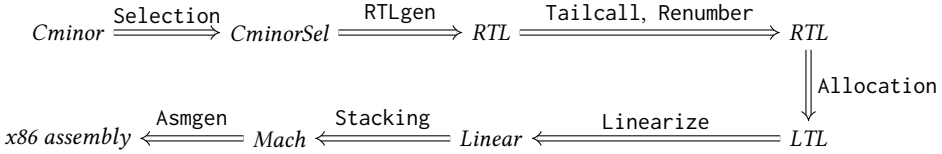


Fig. 12. Proved CompCert Compilation Passes

1056 *Signal* flag in the core state  $\kappa$ . The flag has three possible values: none means a normal state with  
 1057 no lock prefix asserted; locked means a lock# signal is asserted; and done indicates that we have  
 1058 finished executing the accompanying instruction and is about to unset the lock# signal. We give  
 1059 the operational semantics rules for lockxchg as an example at the bottom of Fig 11(b).  
 1060

## 1061 8.2 Adapting CompCert Compiler

1062 As mentioned at the beginning of this section, we adapt the original CompCert compiler passes for  
 1063 compiling Cminor modules to x86 assembly. The compilation passes (shown in Fig. 12) include  
 1064 all translation passes and two optimization passes (Tailcall and Renumber), which are proved to  
 1065 be correct with respect to our correctness judgement. Other optimization passes have not been  
 1066 proved yet and are left as future work.  
 1067

1068 Given the program **let**  $\{y_1, \dots, y_l, \pi_{l+1}, \dots, \pi_m\}$  **in**  $f_1 \parallel \dots \parallel f_n$  consisting of Cminor modules  
 1069  $y_i$  and x86 modules  $\pi_j$  (we omit the corresponding language definitions *sl* and *ge* in the modules to  
 1070 simplify the presentation), the compilation Comp is formulated as

1071  $Comp(\mathbf{let} \{y_1, \dots, y_l, \pi_{l+1}, \dots, \pi_m\} \mathbf{in} f_1 \parallel \dots \parallel f_n) \stackrel{\text{def}}{=} \mathbf{let} \{CminorTrans(y_1), \dots, CminorTrans(y_l), IdTrans(\pi_{l+1}), \dots, IdTrans(\pi_m)\} \mathbf{in} f_1 \parallel \dots \parallel f_n$

1074 where CminorTrans is the adapted compiler consisting of the original CompCert backend passes,  
 1075 and IdTrans is the identity translation which returns the x86 assembly module unchanged. The  
 1076 state transition function  $\varphi$  of both compilers are instantiated as identity function.

1077 We have proved that CminorTrans satisfies our Correct condition:

**Lemma 18** (CminorTrans Correctness).  $\text{Correct}(\text{CminorTrans}, \text{Cminor}, \text{x86asm})$ .

We also prove the correctness of IdTrans, the well-definedness of Cminor and the x86 assembly language, and the determinism of x86 assembly. Together with our framework's final theorem (Theorem 17), we proved the following result:

**Theorem 19** (Compilation Correctness).

$\text{GCorrect}((\text{CminorTrans}, \text{Cminor}, \text{x86asm}), (\text{IdTrans}, \text{x86asm}, \text{x86asm}))$ .

To prove Lemma 18, we try to reuse as much the original CompCert correctness proofs as possible, but have encountered two major challenges: (1) Many CompCert lemmas rely on the specific definition of the CompCert memory model, which is different from ours; and (2) Footprints are new in our setting and it seems we need to re-prove most existing lemmas to support footprints. Below we explain these challenges and show our efforts to reuse CompCert proofs.

*8.2.1 Reusing CompCert Proofs by Converting Memory Layout.* The memory layout in our semantics is different from the CompCert memory model. CompCert memory maintains a `nextblock` field indicating the next block to be allocated. Starting from 0, `nextblock` is incremented after each allocation. Therefore *the sequence of memory allocations in an execution get consecutive natural numbers as block numbers*. As a consequence,  $b$  is a valid block number if  $b < \text{nextblock}$ . Also a block with a smaller block number must be allocated earlier than those with bigger block numbers.

But these assumptions do not hold in our model. As we explain in Sec. 8.1, each allocation takes the block number  $F(N)$  and then increments the index  $N$ , but the block numbers on  $F$  are *not* consecutive. Actually we do not even assume an increasing order of the block numbers. Recall that in our model each thread has its own freelist  $F$ , which can be an infinite sequence of block numbers. Also the freelists of different threads must be disjoint. This means we *cannot* make a freelist  $F$  a infinite sequence of consecutive natural numbers to directly simulate CompCert.

Unfortunately, CompCert correctness proofs heavily rely on its allocation strategy, making it difficult for us to reuse the proofs. For instance, the check of block validity is used extensively in CompCert's fundamental libraries (e.g., `Memory.v` for memroy operations and their properties, and `Separation.v` for separation logic style predicates) and the compilation correctness proofs. Modifying the validity check would affect most of these proofs.

*Lifting simulations to CompCert memory model.* Although our memory model is different from that of CompCert, we can define a bijection between memories under the two models. As a result, the behaviors of a thread under our model are equivalent to its behaviors under CompCert model, and our module-local simulation can be derived from a simulation based on the CompCert model.

In detail, we first define auxiliary semantics for the source and target languages based on CompCert memory model. We denote the CompCert memory by the hat-notation  $\hat{\Sigma}$  and  $\hat{\sigma}$ , and the corresponding core states under CompCert memory model by  $\hat{\mathbb{k}}$  and  $\hat{\kappa}$ . The lemmas below show the correspondence between our original semantics and the auxiliary semantics.

**Lemma 20** (Lifting Cminor). For any Cminor module  $(ge, \gamma)$  and function entry  $f$ , for any  $\mathbb{k}, \Sigma, \hat{\Sigma}$ , if  $\text{InitCore}(\gamma, f) = \mathbb{k}$  and  $ge \subseteq \text{dom}(\Sigma)$ , then there exist  $\hat{\mathbb{k}}, \hat{\Sigma}, \mu_s$  and  $f$  such that  $\mu_s = (\text{dom}(\Sigma), \text{dom}(\hat{\Sigma}), f)$  and  $(\mathbb{F}, (\mathbb{k}, \Sigma)) \preceq_{\mu_s} (\hat{\mathbb{k}}, \hat{\Sigma})$ . Here  $f$  is a bijection from  $\text{dom}(\Sigma)$  to  $\text{dom}(\hat{\Sigma})$ .

It says that, starting from some well-formed initial state, the execution in our semantics is simulated by the execution in the auxiliary semantics. Here  $(\mathbb{F}, (\mathbb{k}, \Sigma)) \preceq_{\mu_s} (\hat{\mathbb{k}}, \hat{\Sigma})$  is defined as a simple lock-step simulation relation, where the corresponding steps in the two semantics generate footprints  $\Delta$  and  $\hat{\Delta}$  such that  $\text{FPmatch}(\mu_s, \Delta, \hat{\Delta})$  holds.

The next lemma shows simulation in the reverse direction in the target language (x86 assembly).

```

1128 Lemma sel_expr_correct:
1129   forall sp e m a v fp,
1130   Cminor.eval_expr sge sp e m a v ->
1131   Cminor.eval_expr_fp sge sp e m a fp ->
1132   forall e' le m', env_lessdef e e' -> Mem.extends m m' ->
1133   exists v', exists fp',
1134   eval_expr tge sp e' m' le (sel_expr a) v' /\ Val.lessdef v v'
1135   /\ eval_expr_fp tge sp e' m' le (sel_expr a) fp' /\ FP.subset fp' fp.

```

Fig. 13. Coq code example

Library Code	Spec		Proof	
	CompCert	Ours	CompCert	Ours
Selectionproof.v	336	500	647	780
RTLgenproof.v	428	543	821	857
Tailcallproof.v	173	328	275	403
ReNUMBERproof.v	86	245	117	356
Allocproof.v	704	785	1410	1696
Linearizeproof.v	236	371	349	732
Stackingproof.v	730	1154	1108	2015
AsmgENproof.v	208	881	571	583
Compositionality (Lemma 8)		580		2249
DRF preservation (Lemma 13)		358		1142
Semantics equivalence (Lemma 14)		1529		4742
Lifting (Theorem 22)		828		1795

Fig. 14. Lines of code (using coqwc) for selected parts of the Coq implementation

**Lemma 21** (Lifting x86asm). For any x86asm module  $(ge, \pi)$  and function entry  $f$ , for any  $\kappa, \sigma, \hat{\sigma}$ , if  $\text{InitCore}(\pi, f) = \kappa$  and  $ge \subseteq \text{dom}(\sigma)$ , then there exist  $\hat{\kappa}, \hat{\sigma}, \mu_t$  and  $f$  such that  $\mu_t = (\text{dom}(\hat{\sigma}), \text{dom}(\sigma), f)$  and  $(\hat{\kappa}, \hat{\sigma}) \preceq_{\mu_t} (F, (\kappa, \sigma))$ . Here  $f$  is a bijection from  $\text{dom}(\hat{\sigma})$  to  $\text{dom}(\sigma)$ .

Theorem 22 below says we can derive the local simulation  $(\mathbb{F}, (\mathbb{k}, \Sigma), \Delta) \preceq_{\mu} (F, (\kappa, \sigma), \delta)$  by proving the simulation  $(\hat{\mathbb{k}}, \hat{\Sigma}, \hat{\Delta}) \preceq_{\hat{\mu}} (\hat{\kappa}, \hat{\sigma}, \hat{\delta})$  instead. Here  $(\hat{\mathbb{k}}, \hat{\Sigma}, \hat{\Delta}) \preceq_{\hat{\mu}} (\hat{\kappa}, \hat{\sigma}, \hat{\delta})$  is the same as our local simulation, except using CompCert memory and the corresponding core states instead of our memory with freelists. As a result, most CompCert libraries and compilation proofs could be reused without modification.

**Theorem 22** (Lifting). If  $(\mathbb{F}, (\mathbb{k}, \Sigma)) \preceq_{\mu_s} (\hat{\mathbb{k}}, \hat{\Sigma}), \text{FPmatch}(\mu_s, \Delta, \hat{\Delta}), (\hat{\kappa}, \hat{\sigma}) \preceq_{\mu_t} (F, (\kappa, \sigma)), \text{FPmatch}(\mu_t, \hat{\delta}, \delta)$ , and  $\mu = \mu_t \circ \hat{\mu} \circ \mu_s$ , then  $(\hat{\mathbb{k}}, \hat{\Sigma}, \hat{\Delta}) \preceq_{\hat{\mu}} (\hat{\kappa}, \hat{\sigma}, \hat{\delta}) \implies (\mathbb{F}, (\mathbb{k}, \Sigma), \Delta) \preceq_{\mu} (F, (\kappa, \sigma), \delta)$ . Here  $\mu_2 \circ \mu_1 = (S_1, S_2, f_2 \circ f_1)$ , if  $\mu_1 = (S_1, S, f_1)$  and  $\mu_2 = (S, S_2, f_2)$ .

**8.2.2 Footprint Preservation.** Although CompCert does not model footprints, many of the definitions and lemmas can be slightly modified to support footprint preservation. For instance, the Selection phase selects appropriate machine operations for operations in Cminor, and generates CminorSel code. One of the key lemmas, `sel_expr_correct`, is shown in Fig. 13, with our newly-added code highlighted as blue texts. It says the selected expression would evaluate to a value refined by the Cminor expression. We simply extends the lemma by requiring *the selected expression has smaller footprint while evaluating on related memory*.

Adapting the simulation invariants (the `match_state` relations) in CompCert proofs is also straightforward. We just need to instrument `match_state` with footprint relations.

### 8.3 Proof Efforts in Coq

Statistics of our Coq implementation and proofs are depicted in Fig 14. We can see that adapting compilation correctness proofs from CompCert is relatively lightweight. For most phases our proofs



1177 are within 300 lines of code more than the original CompCert proofs. The Stacking phase introduces  
 1178 more additional proofs, mostly caused by arguments marshalling for supporting cross-language  
 1179 linking, following Compositional CompCert. In our experience, adapting CompCert’s original  
 1180 compilation proofs to our settings takes less than one person week per translation phase (except  
 1181 for Stacking). For simpler phases such as Tailcallproof.v, Linearizeproof.v, Allocproof.v,  
 1182 and RTLgenproof.v, it takes less than one person day per phase.

1183 In contrast, implementing the framework and proving its correctness are much more chal-  
 1184 lenging, which took us about 1 person year. In particular, proving the equivalence between the  
 1185 non-preemptive and the preemptive semantics under the DRF assumption took us much more time  
 1186 than expected, although it seems to be a well-known folklore theorem. The co-inductive proofs  
 1187 there involve a large number of non-trivial cases of reordering threads’ execution.

1188

## 1189 9 RELATED WORK

1190 *Compiler verification.* There has been various work extending CompCert [Leroy 2009a] to support  
 1191 separate compilation or concurrency. SepCompCert [Kang et al. 2016] extends CompCert with the  
 1192 support of syntactical linking. Their approach requires all the compilation units be compiled by  
 1193 CompCert. They do not support cross-language linking or concurrency as we do.

1194 CompCertX [Gu et al. 2015] extends CompCert semantics with abstract layers for verified  
 1195 separate compilation. It does not support general DRF concurrent program compilation as we do.

1196 Compositional CompCert [Beringer et al. 2014; Stewart et al. 2015] introduces interaction se-  
 1197 mantics to support cross-language interactions. Its module-linking theorem allows escape of stack  
 1198 pointers, which we do not support. They also conjecture that their approach can be extended to  
 1199 verify compositional compilation of well-synchronized concurrent programs. This work addresses  
 1200 the key challenges to do so and proposes a verification framework to achieve this goal.

1201 CompCertTSO [Ševčík et al. 2013] extends CompCert to compile concurrent C-like programs in a  
 1202 relaxed memory model. It focuses on the correctness of a particular compiler and does not support  
 1203 cross-language linking. Also there are two compilation phases whose proofs are not compositional.

1204 Vellvm [Zhao et al. 2012, 2013] proves correctness of several optimization passes for *sequential*  
 1205 LLVM programs. Perconti and Ahmed [2014] verify separate compilation by embedding languages  
 1206 in a combined language. They do not support concurrency either. Ševčík [2011] studies safety of  
 1207 a class of optimizations in concurrent settings using an abstract trace semantics. It is unclear if  
 1208 his approach can be applied to verify general compilation. Lochbihler [2010] verifies a compiler  
 1209 for concurrent Java programs. He requires the target and the source always be lock-step on heap  
 1210 updates, which makes his simulation compositional but restricted. It is unclear how to apply his  
 1211 approach to verify more complex compilers like CompCert.

1212 *Non-preemptive semantics and data-race-freedom.* Non-preemptive (or cooperative) semantics has  
 1213 been developed in various settings for various purposes (e.g., [Abadi and Plotkin 2009; Boudol 2007;  
 1214 Li and Zdancewic 2007; Vouillon 2008; Yi et al. 2011]). Both Ferreira et al. [2010] and Xiao et al. [2018]  
 1215 study the relationships between non-preemptive semantics and DRF, but they do not give any  
 1216 mechanized proofs of termination-preserving semantics equivalence as in our work. DRFx [Marino  
 1217 et al. 2010] proposes a concept called Region-Conflict-Freedom, which looks similar to our NPDRF,  
 1218 but there is no formal operational formulation as we do.

1219

1220 *Compilation validation.* Validation is another technique to ensure correctness of compilation.  
 1221 CompCert employs a verified validator [Rideau and Leroy 2010] for register allocation. More recently  
 1222 Crellvm [Kang et al. 2018] is developed for verified credible compilation of LLVM IR. The work is  
 1223 based on the Vellvm semantics with no support for concurrency. It is not clear how to apply the  
 1224 validation techniques to concurrent settings.

1225

## REFERENCES

- 1226  
1227 Martin Abadi and Gordon Plotkin. 2009. A Model of Cooperative Threads. In *Proceedings of the 36th Annual ACM SIGPLAN-*  
1228 *SIGACT Symposium on Principles of Programming Languages (POPL '09)*. ACM, New York, NY, USA, 29–40.
- 1229 Mark Batty, Scott Owens, Susmit Sarkar, Peter Sewell, and Tjark Weber. 2011. Mathematizing C++ Concurrency. In *POPL*.  
1230 55–66.
- 1231 Lennart Beringer, Gordon Stewart, Robert Dockins, and Andrew W. Appel. 2014. Verified Compilation for Shared-Memory  
1232 C. In *ESOP*. 107–127.
- 1233 Gérard Boudol. 2007. Fair Cooperative Multithreading. In *CONCUR 2007 – Concurrency Theory*, Luís Caires and Vasco T.  
1234 Vasconcelos (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 272–286.
- 1235 Rodrigo Ferreira, Xinyu Feng, and Zhong Shao. 2010. Parameterized Memory Models and Concurrent Separation Logic. In  
1236 *ESOP*, Andrew D. Gordon (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 267–286.
- 1237 Ronghui Gu, Jérémie Koenig, Tahina Ramananandro, Zhong Shao, Xiongnan (Newman) Wu, Shu-Chun Weng, Haozhong  
1238 Zhang, and Yu Guo. 2015. Deep Specifications and Certified Abstraction Layers. In *POPL*. 595–608.
- 1239 Cliff B. Jones. 1983. Tentative Steps Toward a Development Method for Interfering Programs. *ACM Trans. Program. Lang.*  
1240 *Syst.* 5, 4 (1983), 596–619.
- 1241 Jeehoon Kang, Yoonseung Kim, Chung-Kil Hur, Derek Dreyer, and Viktor Vafeiadis. 2016. Lightweight Verification of  
1242 Separate Compilation. In *POPL*. 178–190.
- 1243 Jeehoon Kang, Yoonseung Kim, Youngju Song, Juneyoung Lee, Sanghoon Park, Mark Dongyeon Shin, Yonghyun Kim,  
1244 Sungkeun Cho, Joonwon Choi, Chung-Kil Hur, and Kwangkeun Yi. 2018. Crellvm: Verified Credible Compilation for  
1245 LLVM. In *PLDI 2018*. 631–645.
- 1246 Xavier Leroy. 2009a. Formal verification of a realistic compiler. *Commun. ACM* 52, 7 (2009), 107–115.
- 1247 Xavier Leroy. 2009b. A Formally Verified Compiler Back-end. *J. Autom. Reason.* 43 (December 2009), 363–446. Issue 4.
- 1248 Peng Li and Steve Zdancewic. 2007. Combining Events and Threads for Scalable Network Services Implementation and  
1249 Evaluation of Monadic, Application-level Concurrency Primitives. In *PLDI '07*. 189–199.
- 1250 Hongjin Liang, Xinyu Feng, and Ming Fu. 2012. A Rely-guarantee-based Simulation for Verifying Concurrent Program  
1251 Transformations. In *POPL*. 455–468.
- 1252 Andreas Lochbihler. 2010. Verifying a Compiler for Java Threads. In *ESOP*. 427–447.
- 1253 Daniel Marino, Abhayendra Singh, Todd Millstein, Madanlal Musuvathi, and Satish Narayanasamy. 2010. DRFX: A Simple  
1254 and Efficient Memory Model for Concurrent Programming Languages. In *PLDI '10*. 351–362.
- 1255 James T. Perconti and Amal Ahmed. 2014. Verifying an Open Compiler Using Multi-language Semantics. In *Programming*  
1256 *Languages and Systems*, Zhong Shao (Ed.). 128–148.
- 1257 Silvano Rideau and Xavier Leroy. 2010. Validating Register Allocation and Spilling. In *Compiler Construction*, Rajiv Gupta  
1258 (Ed.). 224–243.
- 1259 Jaroslav Ševčík. 2011. Safe optimisations for shared-memory concurrent programs. In *PLDI*. 306–316.
- 1260 Jaroslav Ševčík, Viktor Vafeiadis, Francesco Zappa Nardelli, Suresh Jagannathan, and Peter Sewell. 2013. CompCertTSO: A  
1261 Verified Compiler for Relaxed-Memory Concurrency. *J. ACM* 60, 3 (2013), 22.
- 1262 Gordon Stewart, Lennart Beringer, Santiago Cuellar, and Andrew W. Appel. 2015. Compositional CompCert. In *POPL*.  
1263 275–287.
- 1264 CompCert Developers. 2017. CompCert-3.0.1. <http://compcert.inria.fr/release/compcert-3.0.1.tgz>
- 1265 Jérôme Vouillon. 2008. Lwt: A Cooperative Thread Library. In *ML*. 3–12.
- 1266 Siyang Xiao, Hanru Jiang, Hongjin Liang, and Xinyu Feng. 2018. Non-Preemptive Semantics for Data-Race-Free Programs.  
1267 In *ICTAC*. to appear.
- 1268 Jaeheon Yi, Caitlin Sadowski, and Cormac Flanagan. 2011. Cooperative Reasoning for Preemptive Execution. In *PPoPP*.  
1269 147–156.
- 1270 Jianzhou Zhao, Santosh Nagarakatte, Milo M.K. Martin, and Steve Zdancewic. 2012. Formalizing the LLVM Intermediate  
1271 Representation for Verified Program Transformations. In *POPL*. 427–440.
- 1272 Jianzhou Zhao, Santosh Nagarakatte, Milo M.K. Martin, and Steve Zdancewic. 2013. Formal Verification of SSA-based  
1273 Optimizations for LLVM. In *PLDI*. 175–186.
- 1274