



CompCertOC: Verified Compositional Compilation of Multi-threaded Programs with Shared Stacks

LING ZHANG, Shanghai Jiao Tong University, China

YUTING WANG*, Shanghai Jiao Tong University, China

YALUN LIANG, Shanghai Jiao Tong University, China

ZHONG SHAO, Yale University, USA

It is a long-standing open problem to support verified compilation of multi-threaded programs *compositionally* when sharing of stack data between threads is allowed. Although certain solutions exist on paper, none of them is completely formalized because of the difficulty in simultaneously enabling sharing and forbidding modification of stack memory in presence of arbitrary memory operations (e.g., pointer arithmetic). We present a compiler verification framework that solves this open problem in the setting of *cooperative multi-threading*. To address the challenges of sharing stack data, we introduce *threaded Kripke memory relations* (TKMR) to support both protection and sharing of stacks in a multi-stack memory model. We further introduce *threaded forward simulations* parameterized by TKMR to capture semantics preservation for compiling program modules in multi-threaded contexts. We show that threaded forward simulations are both *horizontally composable*—thereby enabling the compositional verification of open threads and heterogeneous modules—and *vertically composable*—thereby enabling composition of compiler correctness for multiple compiler passes. Furthermore, threaded forward simulations can be converted into backward simulations. We apply this framework to 18 passes of CompCert to get CompCertOC, the first optimizing verified compiler that supports compositional verification of cooperative multi-threaded programs with shared stacks.

CCS Concepts: • **Theory of computation** → **Program verification; Concurrency**; • **Software and its engineering** → **Correctness; Formal software verification; Compilers**.

Additional Key Words and Phrases: Verified Compositional Compilation, CompCert, Cooperative Multithreading, Stack Sharing

ACM Reference Format:

Ling Zhang, Yuting Wang, Yalun Liang, and Zhong Shao. 2025. CompCertOC: Verified Compositional Compilation of Multi-threaded Programs with Shared Stacks. *Proc. ACM Program. Lang.* 9, PLDI, Article 173 (June 2025), 24 pages. <https://doi.org/10.1145/3729276>

1 Introduction

Software often consists of modules written in different languages (i.e., heterogeneous modules). To ensure the verified properties hold for their compiled binary code, it is essential to support *verified compositional compilation* (VCC), i.e., separate verification of the compilers for heterogeneous modules and composition of compiler correctness theorems. The past decade has witnessed great

*Yuting Wang is the corresponding author.

Authors' Contact Information: [Ling Zhang](#), John Hopcroft Center for Computer Science, School of Computer Science, Shanghai Jiao Tong University, Shanghai, China, ling.zhang@sjtu.edu.cn; [Yuting Wang](#), John Hopcroft Center for Computer Science, School of Computer Science, Shanghai Jiao Tong University, Shanghai, China, yuting.wang@sjtu.edu.cn; [Yalun Liang](#), Zhiyuan College, Shanghai Jiao Tong University, Shanghai, China, liangyalun@sjtu.edu.cn; [Zhong Shao](#), Department of Computer Science, Yale University, New Haven, USA, zhong.shao@yale.edu.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2025 Copyright held by the owner/author(s).

ACM 2475-1421/2025/6-ART173

<https://doi.org/10.1145/3729276>

```

1 void *thread(void *p) { (*(int *)p)++; }
2 int main() { int i = 0; pthread_t tid;
3             pthread_create(&tid, NULL, thread, (void *)&i);
4             pthread_join(tid, NULL); printf("%d\n", i); }

```

Fig. 1. A Concurrent Program not Supported by VCC

progress towards achieving this goal [Gu et al. 2015; Koenig and Shao 2021; Song et al. 2020; Stewart et al. 2015; Wang et al. 2019; Zhang et al. 2024]. These projects target realistic optimizing compilers and are built upon CompCert, the state-of-the-art verified compiler [Leroy 2023].

It has been long believed that VCC can support concurrent programs. For *cooperative multitasking*, when a single concurrent thread (or process) is under focus, it behaves like an open module interacting with other threads through external function calls because context switches behave like function calls in this case. Therefore, compiler correctness for individual threads and processes can be proved separately and composed into that for the complete program with cooperative concurrent semantics. To apply such compiler correctness to concurrent program verification with preemptive semantics, the equivalence between cooperative and preemptive semantics is then proved by exploiting the well-behavedness of source and target programs. For example, in CASCompCert [Jiang et al. 2019a], this is done by exploiting a DRF-SC (i.e., data-race free programs are sequentially consistent) theorem to shuffle context switches to call sites to thread primitives. Since this proof is more about program verification rather than compiler verification, we shall focus on cooperative concurrency in the remaining discussion.

Although the existing works have realized the above vision to a certain extent [Gu et al. 2018; Jiang et al. 2019a; Zha et al. 2022], they are unable to support a common practice in multi-threaded programming, i.e., sharing of stack-allocated data among threads.¹ A simple example for illustrating this problem is depicted in Fig. 1. It is representative of concurrent programs that delegate tasks to child threads by sharing stack data. The main thread creates a child thread for increasing its stack variable `i` by passing its address, which effectively results in sharing of `i`. After the task is completed, the main thread prints the result. In essence, the two threads behave like two modules sharing the same code but with different entry points (`main` and `thread`). More importantly, their executions are interleaved because of concurrency. A context switch starting from one thread and ending with switching back behaves like an external call of that thread. Therefore, compiler correctness for the two threads should simply be those for modules focusing on different entry functions, which are then composed to form the correctness for compiling the whole program.

Unfortunately, none of the existing approaches to VCC supports this seemingly trivial example, because there are fundamental conflicts between supporting shared stack-data and the critical compositionality properties for compiler correctness, i.e., *horizontal and vertical compositionality*:

Vertical Compositionality: $L_1 \leq L_2 \Rightarrow L_2 \leq L_3 \Rightarrow L_1 \leq L_3$

Horizontal Compositionality: $L_1 \leq L'_1 \Rightarrow L_2 \leq L'_2 \Rightarrow L_1 \oplus L_2 \leq L'_1 \oplus L'_2$

Here, a *refinement* relation $L_1 \leq L_2$ (usually a *forward simulation*) encodes the compiler correctness between cooperative semantics of open modules before (L_1) and after (L_2) compilation. Vertical compositionality ensures that the refinements for two adjacent compiler passes can be combined into a single refinement. It is critical for supporting multi-pass compilers. Horizontal compositionality ensures that the parallel refinements for different open modules can be combined together where $L_1 \oplus L_2$ denotes the semantic composition of L_1 and L_2 . It is critical for compiling heterogeneous

¹A solution exists on paper [Jiang et al. 2019b] which is not formalized due to excessive complexity. See §7 for details.

modules where L_1 and L_2 may denote semantics of modules written in different languages with complicated interactions (e.g., mutual calls). Both compositionalitys require appropriate description of the *rely conditions* the open modules depend upon and the *guarantee conditions* they ensure. A great amount of work has been devoted to ensuring both compositionalitys with or without concurrency [Gu et al. 2015; Hur et al. 2012; Koenig and Shao 2021; Neis et al. 2015; Patterson and Ahmed 2019; Song et al. 2020; Stewart et al. 2015; Zhang et al. 2024].

Continue with the example in Fig. 1. The child thread should be able to modify the publicly accessible stack variable `i` of its parent. On the other hand, it should not modify any private stack variable of its parent (e.g., `tid`). The requirement for sharing public stack data while protecting private stack data imposes additional rely-guarantee conditions. Because it is unclear how to formalize these rely-guarantee conditions so that they are preserved under horizontal and vertical composition, existing works on VCC for concurrency (among them the most well-known are Thread-safe CompCertX [Gu et al. 2018] and CASCompCert [Jiang et al. 2019a]) simply forbid stack memory of one thread (`i` for the main thread in Fig. 1) to be leaked to and modified by another thread. This solution deviates from the common assumption that stack memory should be treated no differently from other memory such as heap or global memory (e.g., POSIX thread APIs), and cannot support common programming idioms like in Fig. 1 or *scoped threads* in Rust [Library 2024].

In this paper, we present a formalized solution to the open problem of sharing stack-data in VCC for *cooperative multi-threaded* programs. Our contributions are summarized as follows:

- We introduce *threaded Kripke memory relations* (TKMR), a formalization of rely-guarantee conditions for multi-threaded programs that simultaneously supports sharing of public stack data and protection of private data. On top of that, we introduce *threaded forward simulations*, a notion of compiler correctness as forward simulations for *cooperative concurrent semantics* that supports shared stacks, open threads written in heterogeneous modules, and both horizontal and vertical compositionality. Moreover, they can be flipped into backward simulations in the end by following the same idea of CompCert.
- We verify that 18 of CompCert’s total 20 passes satisfy threaded forward simulations². By vertical compositionality we get CompCertOC, the first verified optimizing compiler that supports multi-threaded open programs with stack sharing. The formal development is based on two existing extensions of CompCert: Nominal CompCert [Wang et al. 2022] for supporting thread-local stacks and CompCertO [Koenig and Shao 2021; Zhang et al. 2024] for supporting VCC. It is fully formalized in Coq based on CompCert v3.13.
- To demonstrate the effectiveness of the above framework, we apply it to verify non-trivial multi-threaded programs that combine heterogeneity (modules written in different languages such as C and assembly) and stack sharing (sharing stack data between parent and child threads using POSIX thread APIs). For this, we develop semantics for describing programs using POSIX thread APIs and the techniques for composing compiler correctness to derive both forward and backward simulations between complete multi-threaded programs.

Note that our compiler correctness theorems include both forward and backward simulation of cooperative semantics. As mentioned above, following the existing approaches [Gu et al. 2018; Jiang et al. 2019a], the connection with preemptive semantics should be handled in program verification frameworks and left for future work (see the discussion in §7 for more details). In the rest of the paper, we first introduce the key ideas for enabling stack sharing in VCC in §2. We then discuss the technical background and challenges in §3. We elaborate on our ideas and contributions in §4, §5 and §6. Finally, we discuss the related work in §7.

²The two missing passes are `SimplExpr` and `UnusedGlob`. We explain why they are omitted in §5.

```

1 /* client.c */
2 #define N 5
3 typedef struct {
4     int *input, *result, size; } Arg;
5 void* server(void *a);
6
7 int main() {
8     pthread_t a;
9     int input[N]={1,2,3,4,5}, result[N];
10    int mask = 0, i;
11    Arg arg = {input,result,N};
12
13    pthread_create(&a,0,server,&arg);
14    for (i = 0; i < N; i++)
15        { mask += input[i]; yield(); }
16    pthread_join(a, NULL);
17    for (i = 0; i < N; i++) {
18        result[i] = result[i] & mask;
19        printf("%d; ", result[i]); }
20 }

```

```

1 /* server.c */
2 void encrypt (int i, int *r);
3 void* server(void *a) {
4     int *i = ((Arg *)a)->input;
5     int *r = ((Arg *)a)->result;
6     int size = ((Arg *)a)->size;
7     for (int j = 0; j < size; j++) {
8         encrypt(i[j], r+j);
9         yield(); }
10    return NULL;
11 }
12 /* encrypt.s */
13 key: .long 42
14 encrypt:
15 Pallocframe 16 8 0 // allocate frame
16 Pmov key RAX
17 Pxor RAX RDI // result = i XOR key
18 Pmov RDI (RSI) // stores to (RSI)
19 Pfreeframe 16 8 0 // free frame
20 Pret

```

(a) Client running on the main thread

(b) Server running on the child thread

Fig. 2. An Example of Multi-Threaded Program with Stack Sharing

2 Key Ideas

2.1 A Running Example

To illustrate the key ideas, we introduce a running example of multi-threaded program composed of open modules in Fig. 2. It consists of a client written in C (Fig. 2a) which runs on the main thread and a server written in C and assembly (Fig. 2b) which runs on the child thread. We use POSIX primitives `pthread_create` and `pthread_join` to implement thread creation and waiting. As we work with cooperative concurrency, we use a `yield` primitive to encode explicit context switches.

This example encrypts N input integers, where the encryption happens both at the server side (XORing inputs with a fixed encryption key 42) and at the client side (encoding the results of the server with a mask calculated from the inputs). The main function uses variables `input` and `result` as its buffers. The addresses and sizes of these buffers are passed to server through `pthread_create`. Concurrent to the execution of server, the client calculates a mask from the input values which is used to encode the results after the `pthread_join`. The function `encrypt` is written in assembly. It receives input value i and the pointer to result r in RDI and RSI respectively, does the encryption using XOR and stores the result. The function `server` simply calls the encryption function written in assembly for each input value. We assume that context switches happen at the end of each iteration by calling `yield`. As we can see, although the input and result arrays reside in the client stack, they are shared and modified by the server thread.

Despite being small, our running example contains all the features we care about in this work, including heterogeneous modules, dynamic thread creation, stack sharing, and context switches between threads. Therefore, it is non-trivial and none of the existing approaches can support its VCC. We shall use this example to illustrate the *feasibility* of our approach in VCC of concurrency with stack-sharing. We leave the application to larger-scale program verification for future work.

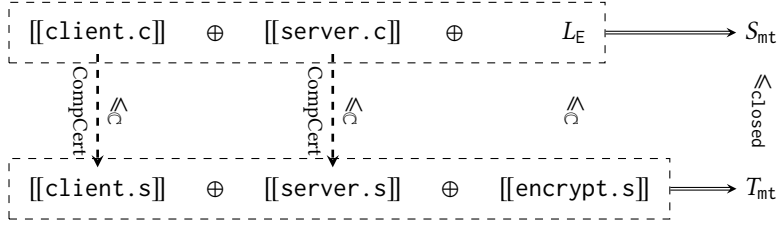


Fig. 3. Verification of the Running Example

2.2 Motivations

Following the standard practice of VCC, the schema for verifying the running example is depicted in Fig. 3. We first verify the correctness for compiling individual modules which is denoted by a refinement relation \leq_C between source and target semantics ($\llbracket M \rrbracket$ denotes the semantics of module M). We then compose individual refinements into \leq_{closed} for the whole program. In the first step, we treat all thread primitives such as `yield`, `thread_create` and `thread_join` as external calls, which are turned into internal executions like other external calls in the second step.

Unfortunately, existing techniques for VCC of sequential programs do not work for thread primitives. To see that, we first review the most recent approach to VCC developed by Zhang et al. [2024]. In this approach, semantics of open modules are implemented as open *labeled transition systems*. Refinements between them are defined as forward simulations (which may be flipped into backward simulations) using rely-guarantee conditions which have been widely studied [Koenig and Shao 2021; Liang et al. 2012; Song et al. 2020; Stewart et al. 2015; Zhang et al. 2024].

Fig. 4 depicts a simulation between two open semantics L_1 and L_2 . The source semantics L_1 is invoked by a query (i.e., function call) q_1 and may issue external call q'_1 after internal execution. When the external function returns with a reply r'_1 , L_1 continues to execute. After several external calls, it finally returns a reply r_1 to the initial query. The target semantics L_2 shares the same structure. Its queries and replies are related to the source via simulation invariants denoted by the vertical double arrows. Certain *rely conditions* for external calls are assumed so that the invariant between r'_1 and r'_2 still holds when they return. In turn, *guarantee conditions* from q_1 and q_2 to r_1 and r_2 need to be proved so that L_1 and L_2 satisfy the rely conditions for their callers. It is exactly those rely-guarantee conditions that enable the composition of refinements in Fig. 3.

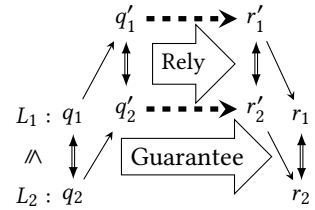


Fig. 4. Rely-Guarantee Simulation

The most important job of the rely-guarantee conditions is to distinguish between *private* memory local to modules or threads and *public* memory shared between them, and to protect the former from modification by the environment. For example, the rely condition for calling `encrypt` in $\llbracket \text{server.c} \rrbracket$ requires that the local variable j is unchanged during the call because j could be optimized into a callee-saved register. The simulation no longer holds if j is changed by `encrypt`.

For this, *Kripke memory relations*, or KMRs are proposed by CompCertO [Koenig and Shao 2021]. In a KMR, a world type W is given s.t. each world $w = (j, m_1, m_2) \in W$ contains a pair of source and target memory states m_1 and m_2 related by a memory injection function j (a structure-preserving embedding of source memory states into target [Leroy et al. 2012]). The evolution of memory across external calls is then defined as the *accessibility* relation between worlds (written as $w \rightsquigarrow w'$) which enforces necessary memory protection. The most general KMR is called `injp`. Its key idea is to

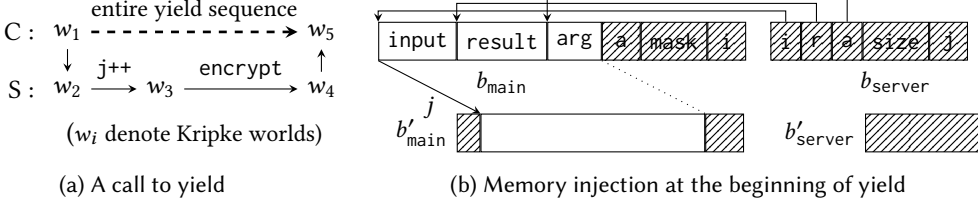


Fig. 6. An Example Showing Inadequacy of injp for Multithreading

define the *private* and *public* regions for the source memory m_1 and the target m_2 w.r.t. the domain and range of injection j for any given world (j, m_1, m_2) . For example, in Fig. 5, m_1 and m_2 contain memory regions to the left of the dashed vertical line. They are related by j before an external call. The regions in the domain and range of j (white areas) are *public* because they represent memory reachable from source and target function arguments. The regions in m_1 *unmapped* by j (shaded areas in m_1) and those *out-of-reach* from j in m_2 (shaded areas in m_2) are *private* memory of the caller because they are either removed from the source memory or inserted into the target memory to store private data (e.g., spilled registers) by the compiler. During the external call, new memory blocks may be allocated, resulting in updated memories m'_1 and m'_2 and injection j' . The accessibility $(j, m_1, m_2) \rightsquigarrow (j', m'_1, m'_2)$ ensures that the private regions in m_1 and m_2 are unchanged in m'_1 and m'_2 , i.e., they cannot be modified by the external call (callee).

Although injp works well for sequential programs, it does not correctly describe the rely-guarantee conditions for context switches when they are treated as external calls. For example, assume the execution of our running example has reached the following state: the server thread has previously switched to the client by calling yield at line 9 in Fig. 2b, and the client thread is about to switch back to the server at line 15 in Fig. 2a. The subsequent call to yield executes $j++$, calls encrypt in the server thread, and returns to the client by another yield, as shown in Fig. 6a where the vertical arrows stand for context switches.

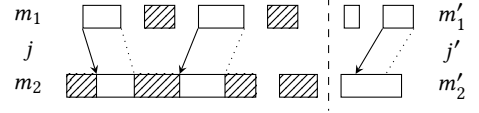


Fig. 5. Kripke Worlds Related by injp

A snapshot of the memory state right before yield is given in Fig. 6b where b_{main} and b_{server} depict the client and server stack frames at the source, respectively; b'_{main} and b'_{server} are target frames after compilation. Note that for simplicity we omit the stack frame for encrypt. Because input, result and arg are leaked public memory, they are preserved by compilation and reside in the domain of injection j . By definition, injp already enables sharing of those public stack regions across threads. Therefore, encrypt is allowed to modify their contents. However, except for input, result and arg, all the remaining memory regions are private: source private variables including j are turned into register values and target ones are newly added private stack regions. Therefore, $j++$ is illegal by injp. Since yield as an external call violates the rely condition that private memory cannot be modified, we fail to prove the compiler correctness.

2.3 Key Idea 1: Threaded Kripke Memory Relations

It is easy to observe that calling yield is actually not the same as calling a regular function. The yield sequence in Fig. 6a begins and ends in the middle of a function call. Such internal execution of a function call should be allowed to modify the private variables on its own stack frame (e.g.,

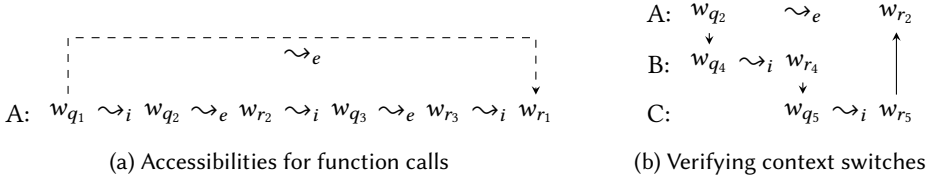


Fig. 7. Rely-Guarantee Reasoning for Function Calls and Context Switches

j). The problem is that `injp` only captures the rely-guarantee of an entire function call, not its internal steps, which is *too strong* for context switches.

Our solution is to revise KMR and the memory model so that they capture the desired rely-guarantee memory protection for both regular function calls and context switches (caused by `yield`, `thread_create`, `thread_join`, etc.). The enhanced KMR is called *threaded Kripke memory relations* (TKMR). In TKMR, a new accessibility relation called *internal accessibility* (denoted by \sim_i) is introduced whose only purpose is to describe the guarantees provided by internal executions. That is, internal steps in the *current thread* can modify *its own private stack*, but never modify the private stack of *the other threads*. For example, the server can perform `j++` internally, but it can never modify `mask`. Now, as a *rely* condition for `yield`, the original accessibility \sim is in conflict with \sim_i as it does not allow modification to *any* private memory. To solve this problem, TKMR relaxes \sim to only protect the private memory of the *current* thread, resulting in *external accessibility* (denoted by \sim_e). This relaxation is possible because simulation invariants for the current thread can be preserved by context switches or external function calls if its own private state is not modified by them; modification of the private stacks of the other threads is harmless (e.g., `j++` in Fig. 6a). In summary, the new accessibility relations are defined as follows:

- (1) $w \sim_e w'$ holds if the *private* memory of the *current* thread in w is unchanged in w' ;
- (2) $w \sim_i w'$ holds if the *private* memory of the *other* threads in w is unchanged in w' .

To formalize the above definitions, we need to distinguish stack memory blocks allocated for different threads. This is difficult in the original CompCert where memory blocks are assigned positive numbers sequentially as their identifiers. We adopt a Nominal Memory Model [Wang et al. 2022] to divide the memory space into global memory and individual stack memory for each thread. Each stack block is named additionally using the thread id of its owner. The memory state also includes a current thread id to distinguish thread local blocks from the other stack blocks.

With the devices of TKMR in place, we first show what the rely-guarantee conditions for regular function calls look like. An example of memory evolution is depicted in Fig. 7a, which describes a function call by thread A. w_{q_1} represents the world from initial calls (e.g., q_1 and q_2 in Fig. 4). There are two external calls q_2 and q_3 during the function execution. Assume q_2 is a `yield` and q_3 is a regular function call (e.g., `encrypt`). We use \sim_e to uniformly describe their *rely* conditions, fulfilling our purpose to treat context switches as regular calls. \sim_i captures the *internal guarantee* which defines legal memory evolution of internal execution starting from either the initial call (q_1) to an external call (q_2), or a function reply (r_2) to another external call (q_3), or from a function reply (r_3) to the final return (r_1). Finally, the *external guarantee* $w_{q_1} \sim_e w_{r_1}$ for the entire function call is proved by reasoning about its behavior like in CompCertO.

We then illustrate how rely conditions for context switches are verified using internal accessibility. Fig. 7b shows a possible execution sequence for the `yield` from q_2 to r_2 in Fig. 7a. Context switches between different threads are denoted by vertical arrows. For example, w_{q_2} and w_{q_4} contain the same memory states except that the current thread is changed from A to B. Between the initial call

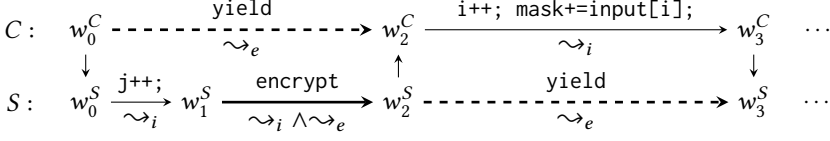


Fig. 9. TKMR Accessibility Relations for the Running Example

to yield by thread A and the return to A, the rely condition \sim_e requires the private and local stack memory of thread A unchanged. This condition is satisfied by the internal guarantee of thread B and C as follows. By definition, $w_{q_4} \sim_i w_{r_4}$ ensures the private stacks of threads A and C are not modified by B. Similarly, $w_{q_5} \sim_i w_{r_5}$ ensures the private stacks of threads A and B are not modified by C. Therefore, the private stack of A is not modified throughout the entire yield sequence. Note that we assume context switches will not further trigger external calls. This is ensured by first performing linking of modules and then linking of threads, as we shall discuss below.

2.4 Key Idea 2: Threaded Forward Simulations

To formalize the correctness for compiling individual threads where thread primitives are treated as external functions, we introduce an enhancement to the rely-guarantee simulation in Fig. 4 called *threaded forward simulations*. As depicted in Fig. 8, it is a forward simulation of *cooperative concurrent semantics* that relates the source semantics L_1 and target semantics L_2 through a TKMR. Its rely condition embeds external accessibility \sim_e to uniformly describe the requirement for sharing and protecting stack memory for both function calls and context switches. The internal guarantee (I.G.) embeds internal accessibility \sim_i for memory protection by internal steps. Finally, the external guarantee embeds \sim_e to capture the guarantee for regular function calls.

Threaded forward simulations satisfy the vertical compositionality as described in the introduction. It means that the correctness of individual compiler passes can be composed into a single threaded forward simulation (e.g., \leq_C in Fig. 3). Moreover, threaded forward simulations are also horizontally composable, thereby enabling the compiler correctness of individual threads or modules to be composed into correctness for compiling whole programs. Technically speaking, this composition has two steps: *module linking* and *thread linking*. The first step links regular function calls, resulting in refinements for source and target semantics whose only external interaction is calling thread primitives. The second step links thread primitives: by exploiting the internal guarantee and external rely conditions of TKMR, a forward simulation for *closed multi-threaded semantics* which concretely defines program loading and thread primitives is generated. Threaded forward simulations can also be flipped into backward simulations and composed into closed backward simulation by exploiting the *determinacy* of target semantics and *receptiveness* of source semantics like in CompCert.

2.5 Verification of the Running Example

We illustrate how to verify our running example with the above key ideas by following the pattern in Fig. 3. First, we need a verified compiler for compiling `client.c` and `server.c` into assembly code. For this, we have developed CompCertOC which consists of 18 passes of CompCert's 20

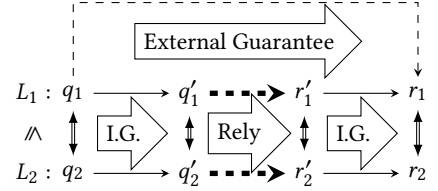


Fig. 8. Threaded Forward Simulation

verified passes. We prove each pass correct w.r.t a TKMR. By vertical compositionality, we get CompCertOC's compiler correctness as the threaded forward simulation $\leq_{\mathbb{C}}$. Second, we verify that the assembly module `encrypt.s` is related to a specification L_E through the same $\leq_{\mathbb{C}}$. Finally, we perform module linking to link the three modules into a complete program, and thread linking to give closed semantics to `pthread_create`, `pthread_join` and `yield` and to derive the closed forward simulation \leq_{closed} . A closed backward simulation \geq_{closed} can also be obtained by first flipping $\leq_{\mathbb{C}}$ and then performing horizontal composition.

The most critical part of the verification is to prove that each open module running on a single thread satisfies the rely-guarantee conditions of TKMR (i.e., internal and external accessibilities). To see how it works, let us revisit our failed example of proving compilation of `yield` at the end of §2.2. A complete loop of context switches between the two threads is depicted in Fig. 9. The execution follows the solid arrows where the vertical ones stand for context switches and horizontal ones stand for internal execution. The private and public memory are determined by injection as depicted in Fig. 6b. Now, consider proving the correctness for compiling `server.c` which treats `encrypt` as an external function. We first need to prove $w_0^S \leadsto_i w_1^S$ holds for internal execution, i.e., `j++` does not modify any private variables `a`, `mask` and `i` in b_{main} of the main thread. This is obviously true. Similarly, we can prove that `encrypt` also satisfies the internal accessibility $w_1^S \leadsto_i w_2^S$. By linking with the module `encrypt.s`, we expose this internal accessibility. Then, we derive the external accessibility $w_0^C \leadsto_e w_2^C$ for the entire `yield` sequence from $w_0^S \leadsto_i w_1^S$ and $w_1^S \leadsto_i w_2^S$ because they guarantee the internal steps of the server do not modify any private stack variable of the client. Note that $w_0^C \leadsto_e w_2^C$ holds because, unlike the original accessibility, \leadsto_e does not require `j` to be unchanged. Symmetrically, $w_2^S \leadsto_e w_3^S$ is satisfied by the internal guarantees of the main thread.

The changes from KMR to TKMR may look simple at first. However, their formal development is quite challenging and requires substantial efforts, which we elaborate in the following sections.

3 Background and Challenges

3.1 Background

3.1.1 Memory Models. The block-based memory model of CompCert [Leroy et al. 2012] defines a memory state m (of type `mem`) as a disjoint set of *memory blocks*. A pointer (b, o) points to the o -th byte from block b where b has type `block` and $o \in \mathbb{Z}$ is an integer. Values of type `val` are either 32- or 64-bit integers or floats, pointers ($\text{Vptr}(b, o)$) or undefined values (`Vundef`). During the compilation, source and target memories are related via *injection functions* $j : \text{block} \rightarrow [\text{block} \times \mathbb{Z}]$ where $[\]$ is overloaded for the `Option` type and its `Some` constructor. $j(b) = \emptyset$ stands for b being removed by compilation and $j(b) = [(b', o)]$ if the source block b is embedded at (b', o) in the target memory. The type of j is given the name `meminj`. The values v_1 and v_2 are related by j (denoted as $v_1 \hookrightarrow_v^j v_2$) if either $v_1 = \text{Vundef}$, $v_1 = v_2$ for scalar values, or $v_1 = \text{Vptr}(b_1, o_1)$, $v_2 = \text{Vptr}(b_2, o_2)$ and $j(b_1) = [(b_2, o_2 - o_1)]$, i.e., pointers are related by injection. Two memory states m_1 and m_2 are related by j (denoted by $m_1 \hookrightarrow_m^j m_2$) if j is a homomorphism from values in m_1 to those in m_2 .

Nominal Memory Model [Wang et al. 2022] is an extension of the block-based memory model using *nominal techniques* [Pitts 2016]. It introduces a nominal interface for the memory model such that the type of memory blocks `block` can be further instantiated by other types. This enables separation of different memory regions (e.g., global memory and stacks for individual threads) by instantiating `block` which we shall exploit later.

3.1.2 Open Simulation of Sequential Semantics. Our work is based on the framework for open simulation in CompCertO [Koenig and Shao 2021]. A *language interface* $A = \langle A^q, A^r \rangle$ is a pair of sets where A^q and A^r denote the type of queries and replies for an open module. The language interface for C semantics is defined as $C = \langle \text{val} \times \text{sig} \times \text{val}^* \times \text{mem}, \text{val} \times \text{mem} \rangle$. A query $v_f[\text{sg}](\vec{v})@m$ consists

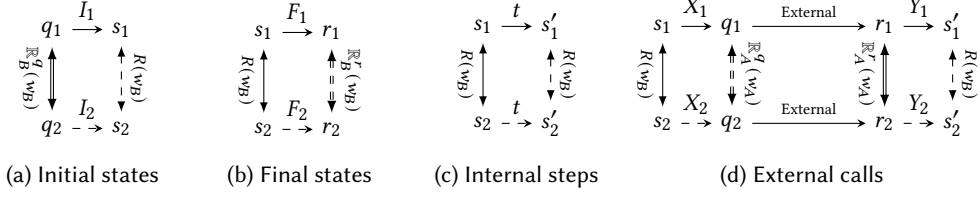


Fig. 10. Simulation Diagrams for Open Forward Simulation

of function pointer, signature, arguments and memory. A reply $v'@m'$ contains a return value together with the updated memory. The interface for assembly is $\mathcal{A} = \langle \text{regset} \times \text{mem}, \text{regset} \times \text{mem} \rangle$ where its queries and replies have the same form $rs@m$ where rs maps registers to values.

Open labeled transition systems (LTS) represent semantics of modules that may accept queries and provide replies at the *incoming side* and provide queries and accept replies at the *outgoing side* (i.e., calling external functions). An open LTS $L : A \Rightarrow B$ is a tuple $\langle D, S, I, \rightarrow, F, X, Y \rangle$ where A and B are the language interfaces for outgoing and incoming sides, $D \subseteq B^q$ a set of initial queries, S a set of internal states, $I \subseteq D \times S$ and $F \subseteq S \times B^r$ transition relations for incoming queries and replies, $X \subseteq S \times A^q$ and $Y \subseteq S \times A^r \times S$ transitions for outgoing queries and replies, and $\rightarrow \subseteq S \times \mathcal{E} \times S$ internal transitions emitting events of type \mathcal{E} . Note that $(s, q^O) \in X$ iff an outgoing query q^O happens at s ; $(s, r^O, s') \in Y$ iff s receives the reply r^O and resumes with an updated state s' .

Kripke relations are used to describe evolution of program states in open simulations. A Kripke relation $R : W \rightarrow \{S \mid S \subseteq A \times B\}$ is a family of relations indexed by a Kripke world W ; for simplicity, we define $\mathcal{K}_W(A, B) = W \rightarrow \{S \mid S \subseteq A \times B\}$. The *simulation convention* relating two language interfaces A_1 and A_2 is defined by $\mathbb{R} : A_1 \Leftrightarrow A_2 = \langle W, \mathbb{R}^q : \mathcal{K}_W(A_1^q, A_2^q), \mathbb{R}^r : \mathcal{K}_W(A_1^r, A_2^r) \rangle$

A open forward simulation (denoted by $L_1 \leq_{\mathbb{R}_A \rightarrow \mathbb{R}_B} L_2$) is defined as [Koenig and Shao 2021]:

Definition 3.1. Given $L_1 : A_1 \Rightarrow B_1$, $L_2 : A_2 \Rightarrow B_2$, $\mathbb{R}_A : A_1 \Leftrightarrow A_2$ and $\mathbb{R}_B : B_1 \Leftrightarrow B_2$, $L_1 \leq_{\mathbb{R}_A \rightarrow \mathbb{R}_B} L_2$ holds if there is some Kripke relation $R \in \mathcal{K}_{W_B}(S_1, S_2)$ that satisfies:

- (1) $\forall w_B \ q_1 \ q_2, (q_1, q_2) \in \mathbb{R}_B^q(w_B) \Rightarrow (q_1 \in D_1 \Leftrightarrow q_2 \in D_2)$
- (2) $\forall w_B \ q_1 \ q_2 \ s_1, (q_1, q_2) \in \mathbb{R}_B^q(w_B) \Rightarrow (q_1, s_1) \in I_1 \Rightarrow \exists s_2, (s_1, s_2) \in R(w_B) \wedge (q_2, s_2) \in I_2.$
- (3) $\forall w_B \ s_1 \ s_2 \ r_1, (s_1, s_2) \in R(w_B) \Rightarrow (s_1, r_1) \in F_1 \Rightarrow \exists r_2, (r_1, r_2) \in \mathbb{R}_B^r(w_B) \wedge (s_2, r_2) \in F_2.$
- (4) $\forall w_B \ s_1 \ s_2 \ t \ s', (s_1, s_2) \in R(w_B) \Rightarrow s_1 \xrightarrow{t} s' \Rightarrow \exists s'_2, (s'_1, s'_2) \in R(w_B) \wedge s_2 \xrightarrow{t} s'_2.$
- (5) $\forall w_B \ s_1 \ s_2 \ q_1, (s_1, s_2) \in R(w_B) \Rightarrow (s_1, q_1) \in X_1 \Rightarrow$
 $\exists w_A \ q_2, (q_1, q_2) \in \mathbb{R}_A^q(w_A) \wedge (s_2, q_2) \in X_2 \wedge$
 $\forall r_1 \ r_2 \ s'_1, (r_1, r_2) \in \mathbb{R}_A^r(w_A) \Rightarrow (s_1, r_1, s'_1) \in Y_1 \Rightarrow \exists s'_2, (s'_1, s'_2) \in R(w_B) \wedge (s_2, r_2, s'_2) \in Y_2.$

Property (1) requires that L_1 and L_2 accept related queries. The remaining properties are illustrated in Fig. 10 where dashed arrows represent existentially quantified relations. Property (2) initializes the simulation invariant $R(w_B)$ from related incoming queries; (3) requires that the final replies are related by the initial world w_B ; (4) requires internal execution to preserve the invariant $R(w_B)$; (5) states that if L_1 issues external call q_1 , L_2 should issue q_2 related to q_1 by some world w_A derived from the current program states. When the external calls return, the invariant $R(w_B)$ should be established again. Note that the *rely-guarantee* conditions in Fig. 4 is implicitly defined in the Kripke relation of replies $\mathbb{R}^r(w)$ by recording initial memory states in w [Koenig and Shao 2021].

3.1.3 Kripke Memory Relations and Direct Refinements. For the subsequent formal development, we describe the most general definition of KMR in jp [Zhang et al. 2024], which is for memory protection in verification of the composition of sequential heterogeneous modules.

Definition 3.2 (Kripke Memory Relation). A Kripke Memory Relation is a tuple $\langle W, f, \sim, R \rangle$ where W is a set of worlds, $f : W \rightarrow \text{meminj}$ a function for extracting injections from worlds, $\sim \subseteq W \times W$ an accessibility relation between worlds, and $R : \mathcal{K}_W(\text{mem}, \text{mem})$ a Kripke relation over memory states that is compatible with the memory operations. We write $w \sim w'$ for $(w, w') \in \sim$.

Definition 3.3 (Kripke Relation with Memory Protection). $\text{injp} = \langle W_{\text{injp}}, f_{\text{injp}}, \sim_{\text{injp}}, R_{\text{injp}} \rangle$ where $W_{\text{injp}} = (\text{meminj} \times \text{mem} \times \text{mem})$, $f_{\text{injp}}(j, _, _) = j$, $(m_1, m_2) \in R_{\text{injp}}(j, m_1, m_2) \Leftrightarrow m_1 \xrightarrow{j}_m m_2$ and $(j, m_1, m_2) \sim_{\text{injp}} (j', m'_1, m'_2) \Leftrightarrow j \subseteq j' \wedge \text{unmapped}(j) \subseteq \text{unchanged-on}(m_1, m'_1) \wedge \text{out-of-reach}(j, m_1) \subseteq \text{unchanged-on}(m_2, m'_2)$.

$\text{unchanged-on}(m, m')$ denotes the footprint of memory cells whose permissions and values are not changed from m to m' . $(b_1, o_1) \in \text{unmapped}(j)$ iff $j(b_1) = \emptyset$ (i.e., b_1 is removed from memory by turning into temporary or register values). $(b_2, o_2) \in \text{out-of-reach}(j, m_1)$ means if $j(b_1) = \lfloor (b_2, o_2 - o_1) \rfloor$ then (b_1, o_1) is not a valid memory location in m_1 (i.e., (b_2, o_2) cannot be reached from any value in m_1 via injection). unmapped and out-of-reach exactly define the *private* regions in source and target memories as depicted in Fig. 5. Therefore, the accessibility \sim_{injp} ensures that private regions are unchanged between function calls and returns (i.e., they are protected).

Definition 3.4 (Direct Refinement). A *direct refinement* is a forward simulation $L_c \leq_{\text{CA}_{\text{injp}} \rightarrow \text{CA}_{\text{injp}}} L_{\text{asm}}$ which directly relates C and assembly semantics. The simulation convention CA_{injp} is implicitly defined using injp as: $\text{CA}_{\text{injp}} = \langle W_{\text{CA}}, \mathbb{R}_{\text{CA}}^q, \mathbb{R}_{\text{CA}}^r \rangle$ where $W_{\text{CA}} = (W_{\text{injp}}, \text{sig}, \text{regset})$. A world $w = ((j, m_1, m_2), \text{sg}, \text{rs})$ consists of related memories, the signature of incoming function and the initial register values. $(v_f[\text{sg}](\vec{v})@m_1, \text{rs}@m_2) \in \mathbb{R}_{\text{CA}}^q(w)$ holds if 1) w corresponds to the initial queries, i.e. $w = ((j, m_1, m_2), \text{sg}, \text{rs})$, 2) the memories and arguments are related by j , and 3) the private stack data before the call is indeed in the private memory regions by injp . $(\text{res}@m'_1, \text{rs}'@m'_2) \in \mathbb{R}_{\text{CA}}^r(w)$ holds if the registers are protected according to the calling convention of CompCert and private memory are protected by injp , i.e., $(j, m_1, m_2) \sim_{\text{injp}} (j', m'_1, m'_2)$ and $(m'_1, m'_2) \in R_{\text{injp}}(j', m'_1, m'_2)$.

3.2 Challenges

Challenges in our formal development of VCC with shared stacks can be summarized at three levels. First and foremost, we need to extend KMR and forward simulation to the threaded versions as described in §2.3 and §2.4. As TKMR requires distinguishing thread local stack blocks from other memory blocks, we need an enriched memory model for its formal definition. Second, with TKMR as the new rely-guarantee conditions, it is unclear whether horizontal and vertical compositionality of forward simulations still hold. In particular, context switches between threads are very different from function calls: the beginning of context switch for one thread is *not* the beginning of context switch for the target thread, but its *end*. Therefore, conflating yield with function calls requires *symmetric* simulation conventions which were not presented in CompCertO. This change combined with the complexity of TKMR incurs major technical challenges in proving the horizontal and vertical compositionality. Finally, realizing our ideas in a realistic optimizing compiler such as CompCert is highly non-trivial. We need to adapt CompCertO with symmetric rely-guarantee interfaces and to adapt the proofs for CompCert's compiler passes and optimizations to fit into threaded open simulations. Below we discuss our formal development to address these challenges.

4 Threaded Kripke Memory Relations and Simulations

To address the first two challenges, we first introduce a multi-threaded memory model. We then extend injp into threaded Kripke memory relation tinjp with internal and external accessibilities. Next, we define threaded simulation conventions and forward simulations using TKMR. Their designs are non-trivial because the internal accessibilities do not fit into the pattern of function calls. Finally, we discuss how to prove the compositionality of threaded forward simulations.

4.1 A Multi-Threaded Memory Model

We implement multi-threaded memory model as an instance of the Nominal Memory Model. The block type `block` is defined as $\text{block} = \text{nat} \times \text{positive}$. A block name $b = (t, p)$ consists of its thread id t and a block identifier p . We write $\text{tid}(b)$ for the thread id of b . A memory state m contains a *support* for recording the history of memory blocks allocated in different regions of memory. Its type is $\text{sup} = \text{nat} \times \text{list}(\text{list positive})$.³ A support (t, stacks) consists of the current running thread t and lists of blocks stacks . A thread id t is always non-zero. $\text{stacks}[0]$ represents the list of global memory blocks allocated so far, and $\text{stacks}[t] (t > 0)$ contains stack blocks allocated for thread t . A block $(t, p) \in (t', \text{stacks})$ is valid iff $p \in \text{stacks}[t]$. When a new block is allocated in a memory state with support (t, stacks) , it is given the fresh name $\text{fresh_block}((t, \text{stacks})) = (t, \max(\text{stacks}(t)) + 1)$. We shall write $\text{tid}(m)$ for the current thread id of the support in memory state m .

With the above memory model, it is possible to decide whether b is allocated by the current thread by comparing $\text{tid}(b)$ with $\text{tid}(m)$, which is essential for defining accessibilities in TKMR. The semantics of `thread_create` and `yield` affect only the support (t, stacks) : `yield` changes the current thread id t to some t' where $1 \leq t' < \text{length}(\text{stacks})$. `thread_create` appends an empty list after stacks as the name space for the new thread. Notice that, because we have explicitly maintained stacks for individual threads, our memory model is more faithful to how OS views the memory state than previous work. See §7 for more details.

4.2 Threaded Kripke Memory Relations

We define threaded Kripke memory relation by extending Definition 3.2 with two accessibilities as a tuple $\langle W, f, \rightsquigarrow_i, \rightsquigarrow_e, R \rangle$. TKMR with memory protection is defined as follows.

Definition 4.1 (Threaded Kripke Memory Relation with Memory Protection). The threaded version of `injp` is defined as $\text{tinjp} = \langle W_{\text{tinjp}}, f_{\text{tinjp}}, \rightsquigarrow_i^{\text{tinjp}}, \rightsquigarrow_e^{\text{tinjp}}, R_{\text{tinjp}} \rangle$ where the related memories have the same thread id, i.e., $R_{\text{tinjp}}(j, m_1, m_2) \Leftrightarrow m_1 \hookrightarrow_m^j m_2$ and $\text{tid}(m_1) = \text{tid}(m_2)$. W_{tinjp} and f_{tinjp} are the same as in `injp` (Definition 3.3). The internal and external accessibilities are defined as follows:

$$\begin{aligned}
 (j, m_1, m_2) \rightsquigarrow_i^{\text{tinjp}} (j', m'_1, m'_2) &\Leftrightarrow j \subseteq j' \wedge \text{tid}(m_1) = \text{tid}(m'_1) \wedge \\
 &\quad \text{unmapped}(j) \cap \text{thr-ext}(m_1) \subseteq \text{unchanged-on}(m_1, m'_1) \wedge \\
 &\quad \text{out-of-reach}(j, m_1) \cap \text{thr-ext}(m_2) \subseteq \text{unchanged-on}(m_2, m'_2). \\
 (j, m_1, m_2) \rightsquigarrow_e^{\text{tinjp}} (j', m'_1, m'_2) &\Leftrightarrow j \subseteq j' \wedge \text{tid}(m_1) = \text{tid}(m'_1) \wedge \\
 &\quad \text{unmapped}(j) \cap \text{thr-int}(m_1) \subseteq \text{unchanged-on}(m_1, m'_1) \wedge \\
 &\quad \text{out-of-reach}(j, m_1) \cap \text{thr-int}(m_2) \subseteq \text{unchanged-on}(m_2, m'_2).
 \end{aligned}$$

This is the formal definition of \rightsquigarrow_i and \rightsquigarrow_e we presented in §2.3. The *private* regions in (j, m_1, m_2) are already defined by $\text{unmapped}(j) \subseteq m_1$ and $\text{out-of-reach}(j, m_1) \subseteq m_2$ via the injection function j in Definition 3.3. $\text{thr-ext}(m) \subseteq m$ stands for the *thread-external* blocks, i.e., $b \in \text{thr-ext}(m) \Leftrightarrow \text{tid}(b) \neq \text{tid}(m)$. Similarly, *thread-internal* blocks are defined as $\text{thr-int}(m) \subseteq m$ where $b \in \text{thr-int}(m) \Leftrightarrow \text{tid}(b) = \text{tid}(m)$. Note that both \rightsquigarrow_i and \rightsquigarrow_e require the thread id unchanged. The current thread id tid is changed by context switches. As a result, although we have a static depiction of private memory (e.g., in Fig. 6b, j is always private as it is in $\text{unmapped}(j)$), the accessibilities ensure that private stack memory is allowed to be modified by its owner and protected from the other threads.

³The terminology *support* comes from nominal techniques [Pitts 2016].

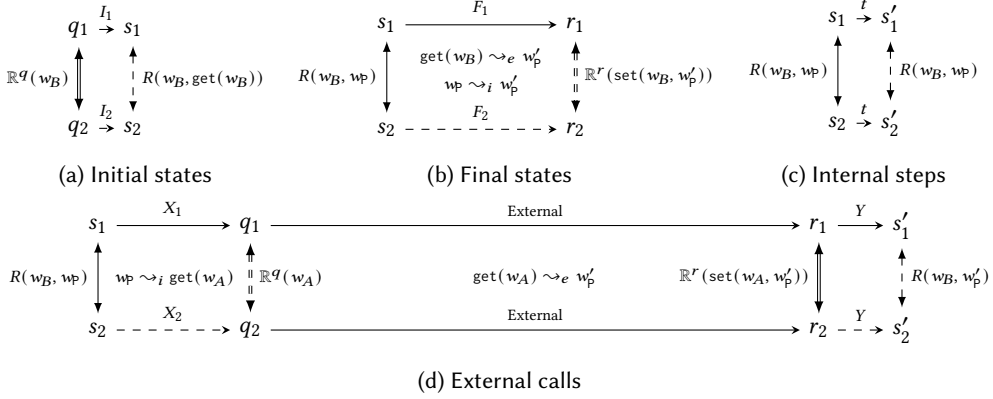


Fig. 11. Simulation Diagrams for Threaded Forward Simulation

4.3 Threaded Forward Simulations

Definition 4.2 (Threaded Simulation Conventions). A threaded simulation convention $\mathbb{R}_P : A_1 \Leftrightarrow A_2 = \langle W, \mathbb{R}^q : \mathcal{K}_W(A_1^q, A_2^q), \mathbb{R}^r : \mathcal{K}_W(A_1^r, A_2^r), P \rangle$ is defined using a TKMR $P = \langle W_P, f_P, \sim_i^P, \sim_e^P, R_P \rangle$, where W_P is a sub-world of W equipped with operations $\text{get} : W \rightarrow W_P$ and $\text{set} : W \rightarrow W_P \rightarrow W$.

Compared with the original simulation conventions, there are two differences in the threaded version. First, the accessibility \sim is implicitly encoded in \mathbb{R}^r with a \Diamond modality in the former, while the threaded simulation conventions explicitly contain a TKMR P with internal and external accessibilities. This exposure is necessary for defining the threaded forward simulations as we shall discuss shortly. For simplicity, we often drop the superscript P in \sim_i^P and \sim_e^P in our discussions. Second, threaded simulation conventions provide explicit get and set operators for extracting Kripke worlds W_P for memory states from Kripke worlds W for program states which may contain additional information such as registers and function signatures. This is necessary because internal execution between context switches may relate functions with entirely different arguments and signatures (e.g., encrypt and yield). To describe these internal guarantees, we need to extract pure memory states from program states through get and set .

Definition 4.3 (Threaded Forward Simulations). A threaded forward simulation $L_1 \leq_{\mathbb{R}_P} L_2$ is defined between open semantics with the same language interfaces on both sides, i.e., for $L_1 : A_1 \Rightarrow A_1$ and $L_2 : A_2 \Rightarrow A_2$. The internal invariant is defined using a Kripke Relation with two worlds $\mathcal{K}_{W_1, W_2}(A, B) : W_1 \rightarrow W_2 \rightarrow \{S \mid S \subseteq A \times B\}$. Given the threaded simulation convention $\mathbb{R}_P : A_1 \Leftrightarrow A_2$, $L_1 \leq_{\mathbb{R}_P} L_2$ holds if there exists some Kripke relation $R \in \mathcal{K}_{W, W_P}(S_1, S_2)$ that satisfies:

- (1) $\forall w_B \ q_1 \ q_2, (q_1, q_2) \in \mathbb{R}_B^q(w_B) \Rightarrow (q_1 \in D_1 \Leftrightarrow q_2 \in D_2)$
- (2) $\forall w_B \ q_1 \ q_2 \ s_1, (q_1, q_2) \in \mathbb{R}_B^q(w_B) \Rightarrow (q_1, s_1) \in I_1 \Rightarrow \exists s_2, (s_1, s_2) \in R(w_B, \text{get}(w_B)) \wedge (q_2, s_2) \in I_2.$
- (3) $\forall w_B \ wp \ s_1 \ s_2 \ r_1, (s_1, s_2) \in R(w_B, wp) \Rightarrow (s_1, r_1) \in F_1 \Rightarrow \exists wp' \ r_2, (r_1, r_2) \in \mathbb{R}_B^r(\text{set}(w_B, wp')) \wedge (s_2, r_2) \in F_2 \wedge \text{get}(w_B) \sim_e wp' \wedge wp \sim_i wp'.$
- (4) $\forall w_B \ wp \ s_1 \ s_2 \ t \ s'_1, (s_1, s_2) \in R(w_B, wp) \Rightarrow s_1 \xrightarrow{t} s'_1 \Rightarrow \exists s'_2, (s'_1, s'_2) \in R(w_B, wp) \wedge s_2 \xrightarrow{*} s'_2.$
- (5) $\forall w_B \ wp \ s_1 \ s_2 \ q_1, (s_1, s_2) \in R(w_B, wp) \Rightarrow (s_1, q_1) \in X_1 \Rightarrow \exists w_A \ q_2, (q_1, q_2) \in \mathbb{R}_A^q(w_A) \wedge (s_2, q_2) \in X_2 \wedge wp \sim_i \text{get}(w_A) \wedge \forall wp' \ r_1 \ r_2 \ s'_1, \text{get}(w_A) \sim_e wp' \Rightarrow (r_1, r_2) \in \mathbb{R}_A^r(\text{set}(w_A, wp')) \Rightarrow (s_1, r_1, s'_1) \in Y_1 \Rightarrow \exists s'_2, (s'_1, s'_2) \in R(w_B, wp') \wedge (s_2, r_2, s'_2) \in Y_2.$

This definition is an enhancement to Definition 3.1. To deal with symmetric behaviors of context switches as indicated in §3.2, the same simulation convention \mathbb{R}_P is used for both incoming and outgoing sides (Symmetric interfaces are not a limitation as any LTS $L : A \rightarrow B$ can be embedded into the interface $A \cup B \rightarrow A \cup B$). The differences from Definition 3.1 are highlighted in red. The key to proving threaded forward simulation is to verify that internal and external accessibilities hold. For this, the simulation invariant $R(w_B, w_P)$ should remember the relation between source and target initial queries (when entering an open module for the first time) in w_B and the relation between source and target memories *when the last time the module is entered* in w_P (either by an initial function call or a return from an external function call). When $(s_1, s_2) \in R(w_B, w_P)$ holds, it implies that the *current* Kripke world $w_s = (j, m_1, m_2)$ (where m_1 and m_2 are the memory states in s_1 and s_2 , respectively, and j is the current memory injection) satisfies both the external guarantee from w_B and the internal guarantee from w_P , i.e., $\text{get}(w_B) \sim_e w_s$ and $w_P \sim_i w_s$. By maintaining both guarantees as the execution goes, the threaded forward simulation is proved.

Based on the above observations, we elaborate on the ideas behind properties (2-5) which are illustrated in Fig. 11 (property (1) is the same as before). Property (2) initializes the simulation invariant R using the memories $\text{get}(w_B)$ in the initial queries (Fig. 11a). Dually, property (3) describes returning replies r_1 and r_2 related by the updated world $\text{set}(w_B, w'_P)$ from the final invariant R (Fig. 11b). The final memory world w'_P should satisfy both the *external guarantee* $\text{get}(w_B) \sim_e w'_P$ for the whole function call and the *internal guarantee* $w_P \sim_i w'_P$ for the internal execution since last time the module was entered. Using the definition of R mentioned above, they can be satisfied by picking the current memory world w_s for w'_P .

Property (4) states that, although the current memories are changed by internal steps, the accessibilities $\text{get}(w_B) \sim_e w_s$ and $w_P \sim_i w_s$ still hold for the same w_B and w_P (Fig. 11d). For this, one needs to prove that an internal step cannot change the thread local private memory in w_B , and it cannot change the thread external private memory in w_P . In other words, a function call can only change the private stack regions both owned by the current thread and allocated after the initial query. In our running example, the main function can modify the loop invariant i allocated by itself but can not change j because it is owned by another thread. The encrypt function also cannot modify j because it is already allocated by the server on its private stack before calling encrypt.

Finally, for property (5) of external calls (Fig. 11d), when the semantics issue queries q_1 and q_2 related by w_A , the *internal guarantee* of recent internal execution is enforced by $w_P \sim_i \text{get}(w_A)$. After the external call, we require that the replies are related by $\text{set}(w_A, w'_P)$ where w'_P stands for the updated memory states. Given the rely condition $\text{get}(w_A) \sim_e w'_P$, we reestablish the internal invariant $(s'_1, s'_2) \in R(w_B, w'_P)$ using the memory world w'_P from replies. Note that one needs to ensure thread local private memory is not changed by the environment (i.e., $w_B \sim_e w'_P$ holds). For the function server in running example, the rely conditions on encrypt and yield ensure that *any* thread local blocks are unchanged. Therefore, the private memory from its caller is properly protected during the whole execution of server. Note also the necessity of get and set here: we need to relate the memory state w_P initially extracted from w_B to a different world w_A .

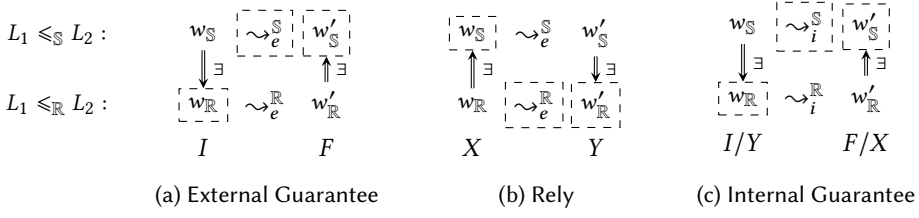
4.4 Compositionality of Threaded Forward Simulations

4.4.1 *Vertical Composition.* A naive composition theorem for threaded forward simulations holds:

THEOREM 4.4. *If $L_1 : A_1 \rightarrow A_1, L_2 : A_2 \rightarrow A_2, L_3 : A_3 \rightarrow A_3$ and $\mathbb{R} : A_1 \Leftrightarrow A_2, \mathbb{S} : A_2 \Leftrightarrow A_3$, then*

$$L_1 \leq_{\mathbb{R}} L_2 \Rightarrow L_2 \leq_{\mathbb{S}} L_3 \Rightarrow L_1 \leq_{\mathbb{R} \cdot \mathbb{S}} L_3.$$

Here, $(_ \cdot _)$ is the concatenation of threaded simulation convention s.t. $\mathbb{R} \cdot \mathbb{S} = \langle W_{\mathbb{R}} \times W_{\mathbb{S}}, \mathbb{R}^q \cdot \mathbb{S}^q, \mathbb{R}' \cdot \mathbb{S}', P_{\mathbb{R}} \cdot P_{\mathbb{S}} \rangle$ where $(q_1, q_3) \in \mathbb{R}^q \cdot \mathbb{S}^q \Leftrightarrow \exists q_2, (q_1, q_2) \in \mathbb{R}^q(q_1, q_2)$ and $(q_2, q_3) \in \mathbb{S}^q(q_2, q_3)$.

Fig. 12. A Refinement of Threaded Simulation Convention $R \sqsubseteq S$

(same for $\mathbb{R}' \cdot \mathbb{S}'$). The TKMRs are also composed where $(w_1, w_2) \sim_i (w'_1, w'_2)$ is defined as $w_1 \sim_i w_2 \wedge w'_1 \sim_i w'_2$ (same for \sim_e). Since the concatenated simulation convention is defined as the composition of two relations with an existential quantifier in the middle, one can simply instantiate the existential quantifier using the states of L_2 to complete the proof.

However, this composed interface exposes the internal compilation and weakens the compiler correctness. To get a direct refinement between C and assembly semantics, we need to define the refinement between simulation conventions. Given $\mathbb{R}, \mathbb{S} : A_1 \Leftrightarrow A_2$, we write $\mathbb{R} \sqsubseteq \mathbb{S}$ for that \mathbb{R} is refined by \mathbb{S} . The following refinement of threaded simulations can be proved:

THEOREM 4.5. *Given $L_1 : A_1 \Rightarrow A_1, L_2 : A_2 \Rightarrow A_2$ and $\mathbb{R}, \mathbb{S} : A_1 \Leftrightarrow A_2$,*

$$\mathbb{R} \sqsubseteq \mathbb{S} \Rightarrow L_1 \leq_{\mathbb{R}} L_2 \Rightarrow L_1 \leq_{\mathbb{S}} L_2$$

The key ideas behind its proof are depicted in Fig. 12. To refine the interface \mathbb{R} of the threaded simulation to \mathbb{S} , we need to construct multiple worlds with proper accessibilities (denoted using dashed boxes in Fig. 12) as rely-guarantee conditions. In Fig. 12a, given the world w_S of initial queries (I), we first construct w_R to relate the queries such that the property (2) of $L_1 \leq_{\mathbb{R}} L_2$ is satisfied. When the semantics return with replies related by w'_R (F), we need to construct an updated w'_S for these replies such that the external guarantee \sim_e^S is satisfied. Symmetrically, in Fig. 12b, when the semantics issue external calls related by w_R (X), we construct w_S for external calls using interface \mathbb{S} . After external calls (Y), we construct w'_R and prove \sim_e^R as the rely condition for $L_1 \leq_{\mathbb{R}} L_2$. For threaded simulations, we also need to prove the refinement of internal guarantees, as depicted in Fig. 12c. Note that the internal guarantee has four variants: an open module can begin its partial execution by I or Y and end by X or F . Notice that since internal and external guarantees may share the same starting world, we need to prove Fig. 12c together with Fig. 12a and Fig. 12b in one pass.

With the above refinement, vertical composition of threaded forward simulations becomes possible. For example, a common convention is $\text{ctinjp} : C \Leftrightarrow C$ for simulations between semantics with C interfaces using `tinjp`. We can prove its transitivity as the following refinement theorem:

THEOREM 4.6. $\text{ctinjp} \cdot \text{ctinjp} \sqsubseteq \text{ctinjp}$

Then, given two threaded forward simulations $L_1 \leq_{\text{ctinjp}} L_2$ and $L_2 \leq_{\text{ctinjp}} L_3$, we first compose them into $L_1 \leq_{\text{ctinjp} \cdot \text{ctinjp}} L_3$ by Theorem 4.4, and then derive $L_1 \leq_{\text{ctinjp}} L_3$ by applying Theorem 4.5 to Theorem 4.6. As we shall see in §5, other simulation conventions can support vertical composition in a way similar to `ctinjp`.

4.4.2 Horizontal Composition. The horizontal composition of threaded forward simulations can be divided into two steps. We first compose open modules assuming that they run on the same thread. Then, we compose open semantics into closed semantics running on a multi-threaded abstract machine. Their closed simulation can be derived from the threaded forward simulations. We call the first step *module linking* and the second step *thread linking*. We discuss their correctness below.

THEOREM 4.7 (CORRECTNESS OF MODULE LINKING). *Let $L_1 \oplus L_2$ denote the semantic linking of L_1 and L_2 . Given $L_1, L'_1 : A_1 \rightarrow A_1$, $L_2, L'_2 : A_2 \rightarrow A_2$ and threaded simulation convention $\mathbb{R} : A_1 \Leftrightarrow A_2$, the simulations can be horizontally composed as:*

$$L_1 \leq_{\mathbb{R}} L_2 \Rightarrow L'_1 \leq_{\mathbb{R}} L'_2 \Rightarrow L_1 \oplus L'_1 \leq_{\mathbb{R}} L_2 \oplus L'_2$$

We omit the concrete definition of the semantic linking $L_1 \oplus L_2$. Its proof is similar to the horizontal composition in CompCertO [Koenig and Shao 2021]. The main difference is that we also need to compose *internal guarantee* conditions, which is achieved with the transitivity of \sim_i .

The multi-threaded semantics is defined as a closed LTS with a non-deterministic thread-scheduling strategy. A closed LTS is a tuple $\langle S, I, F, \rightarrow \rangle$ containing program state S , predicates for initial and final state I, F and the internal transition \rightarrow .

We define *thread linking function* MT_C and $\text{MT}_{\mathcal{A}}$ for transferring C and assembly open semantics into multi-threaded semantics. We omit the concrete definition for simplicity and present the key ideas using C semantics; the construction for assembly is similar. Given $L : C \rightarrow C$, $\text{MT}_C(L) = \langle S_M, I_M, F_M, \rightarrow_M \rangle$. The program state S_M includes the current thread id t and a list of *thread states* which are defined using the program state of L . The current global memory can be found in the current thread state. The initial state I_M and final state F_M are defined using the *query* and *reply* of main function on the main thread. The main thread applies I (from L) to take a query to main as initialization and finally issues a reply using F to produce the return value. The transition step \rightarrow_M has three cases. The internal step is defined as a step of L on the current running thread. The thread creation step happens if the current thread calls `pthread_create` (as an external call). Its arguments and the current memory state are used to create a new thread state. Finally, for context switches, the current thread can switch out by calling `yield` or `thread_join`, or ending its execution. Similarly, the target thread is waiting for a reply of a previous `yield`, a `thread_join`, or thread initialization. The execution simply changes the current thread id, pass it (by constructing a query or reply) to wake up the target thread, and put the current thread into a waiting state.

The correctness of thread linking is stated as follows.

THEOREM 4.8 (CORRECTNESS OF THREAD LINKING).

$$\forall L_1 : C \rightarrow C, L_2 : \mathcal{A} \rightarrow \mathcal{A}, L_1 \leq_C L_2 \Rightarrow \text{MT}_C(L_1) \leq_{\text{closed}} \text{MT}_{\mathcal{A}}(L_2)$$

Here, \mathbb{C} is the simulation convention derived from CompCertOC which we define in §5.3. The closed simulation \leq_{closed} is exactly the final forward simulation in the original CompCert. This theorem is proved by exploiting the rely-guarantee conditions of threaded forward simulations, especially the internal guarantees, as depicted in Fig. 7b.

5 CompCert for Open and Concurrent Modules

We now discuss how to verify the compiler passes of CompCert and compose them into the threaded forward simulation \leq_C as the correctness of CompCertOC. CompCert compiles C programs into Asm programs through 20 passes, including optimization passes on the RTL intermediate language. Most of the extensions to CompCert start from the Clight language which is compiled from C source through the first pass named `SimplExpr` by extracting side effects from expressions. Our development is based on the latest CompCertO [Zhang et al. 2024] which also starts from Clight and hence contains 19 passes of CompCert. However, it is recently discovered that the optimization pass `Unusedglob` in CompCertO for removing unused static definitions relies on an incorrect axiom about the transformation of global symbols made by CompCert. As a result, the axioms for CompCertO lead to inconsistency. Fortunately, because this axiom is only used for proving `Unusedglob`, consistency can be recovered by simply removing `Unusedglob` from CompCertO. Therefore, CompCertOC contains 18 passes of CompCert. We believe this omission is orthogonal

Table 1. Significant Passes of CompCertOC

Languages/Passes	Interfaces/Conventions	Languages/Passes	Interfaces/Conventions
Clight	$C \rightarrow C$	Inlining	C_{tinjp}
Self-Sim	ro	Constprop	$ro \cdot C_{\text{tinjp}}$
SimplLocals	C_{tinjp}	CSE	$ro \cdot C_{\text{tinjp}}$
Csharpminor	$C \rightarrow C$	Deadcode	$ro \cdot C_{\text{tinjp}}$
Cminorgen	C_{tinjp}	Allocation	$wt \cdot C_{\text{ext}} \cdot CL$
Cminor	$C \rightarrow C$	LTL	$\mathcal{L} \rightarrow \mathcal{L}$
Selection	$wt \cdot C_{\text{ext}}$	Tunneling	ltl_{ext}
CminorSel	$C \rightarrow C$	Linear	$\mathcal{L} \rightarrow \mathcal{L}$
RTLgen	C_{ext}	Stacking	LM_{tinjp}
RTL	$C \rightarrow C$	Mach	$M \rightarrow M$
Self-Sim	C_{tinjp}	Asmggen	$mach_{\text{ext}} \cdot MA$
Tailcall	C_{ext}	Asm	$\mathcal{A} \rightarrow \mathcal{A}$

to our work as Unusedglob only deals with *global memory* while CompCertOC deals with shared *stack memory*. A solution to add back Unusedglob is left for future work.

In the subsections below, we first discuss the verification of individual passes of CompCertOC in §5.1. We then discuss properties for refining simulation conventions used in CompCertOC in §5.2. Finally, we discuss the composition of the correctness of individual passes into a threaded forward simulation with the simulation convention \mathbb{C} that directly relates C and assembly semantics in §5.3.

5.1 Threaded Forward Simulations for Individual Passes

The important compiler passes and their simulation types are listed in Table 1 together with the intermediate languages and their interfaces (which are in bold fonts). The passes on the right follow those on the left in the compilation chain. The passes marked in red are optimizations for RTL. We insert two self-simulations (in blue) for refining the composed simulation convention into a direct refinement. Passes using the identity simulation do not affect the composition and hence are omitted in Table 1 (including Cshmggen, Renumber, Linearize, CleanupLabels and Debugvar).

5.1.1 Simulation Conventions and Semantic Invariants. We first discuss the simulation conventions presented in Table 1. The simulation conventions in the form of $I_P : \mathcal{X} \Leftrightarrow \mathcal{X}$ use TKMR P to relate the source and target semantics with the language interface \mathcal{X} for language I . \mathcal{L} is the interface for the LTL intermediate language. ext is a simplified version of tinjp which assumes source and target memories have the same structure. Therefore no private memory needs to be protected and the only rely-guarantee condition is the invariance of thread id.

$CL : C \Leftrightarrow \mathcal{L}$ and $MA : M \Leftrightarrow \mathcal{A}$ are *structural conventions* for relating queries and replies in different intermediate languages. They assume that the source and target memory states are the same. Another structural convention $LM_{\text{tinjp}} : \mathcal{L} \Leftrightarrow M$ for the Stacking pass is parameterized by tinjp . This is because the source and target memories are related by an injection as the stack pointer, return address and outgoing arguments are allocated on the stack after this pass. Therefore, the private stack in the target memory needs to be protected by tinjp .

ro and wt are *semantic invariants* relating queries and replies in the same language. They describe invariants of an open semantics. wt ensures that the arguments and return values are well-typed. ro ensures that the values of read-only *constants* are never modified by execution.

5.1.2 Verification of Individual Passes. We revise the correctness proofs for the compiler passes in CompCertO so that they fit into the framework of threaded forward simulations. The difficulties

come from two limitations in the original proofs. First, many passes in CompCertO are proved using *different* simulation conventions for the incoming and outgoing sides. Second, in the original proofs, memory protection is only provided at the outgoing side, not at the incoming side. Both limitations need to be lifted because the symmetry of context switches demands the *same* simulation convention and memory protection for *both* the incoming and outgoing sides. Moreover, according to the discussion in §4.3 (after Definition 4.3), to establish a threaded forward simulation for a compiler pass, we need to find a simulation invariant R s.t. $(s_1, s_2) \in R(w_B, w_P)$ iff the external and internal accessibilities $\text{get}(w_B) \rightsquigarrow_e w_s$ and $w_P \rightsquigarrow_i w_s$ hold for the current Kripke world w_s extracted from s_1 and s_2 . The most difficult passes to prove are optimizations involving static analysis (ConstProp, CSE and DeadCode) and Stacking for concretely laying out the stack frames. For the former, we need to prove that the blocks whose addresses are never taken (hence may be optimized away) are always labeled as private in the memory world w_P in $R(w_B, w_P)$. For the latter, because the original Stacking pass uses different simulation conventions for the incoming and outgoing sides, we have to do a substantial rewrite of the proof with the new symmetric structural convention LM_{tinjp} and prove that the newly added private stack is protected by an invariant R .

5.2 Properties for Refining Thread Simulation Conventions

We adopt the techniques discussed in §4.4.1 to compose and refine the correctness of individual passes into that of CompCertOC. For this, we discuss the necessary refinement lemmas for threaded simulation conventions. First, the following lemmas are for refining C-level simulation conventions:

LEMMA 5.1. (1) $c_{\text{tinjp}} \cdot c_{\text{tinjp}} \sqsubseteq c_{\text{tinjp}}$ (2) $c_{\text{ext}} \cdot c_{\text{ext}} \sqsubseteq c_{\text{ext}}$
 (3) $c_{\text{tinjp}} \cdot c_{\text{ext}} \cdot c_{\text{tinjp}} \sqsubseteq c_{\text{tinjp}}$ (4) $ro \cdot c_{\text{tinjp}} \cdot ro \cdot c_{\text{tinjp}} \sqsubseteq ro \cdot c_{\text{tinjp}}$.

Property (1) is exactly Theorem 4.6. (2) and (3) are trivial because ext does not protect any private memory. (4) is an extension of (1) for transitively composing the ro invariant with tinjp .

LEMMA 5.2. For $K \in \{\text{ext}, \text{tinjp}\}$, (1) $c_K \cdot \text{CL} \equiv \text{CL} \cdot \text{ltl}_K$ (2) $\text{mach}_K \cdot \text{MA} \equiv \text{MA} \cdot \text{asm}_K$.

Here, $\mathbb{R} \equiv \mathbb{S} \Leftrightarrow \mathbb{R} \sqsubseteq \mathbb{S} \wedge \mathbb{S} \sqsubseteq \mathbb{R}$. Since the structural conventions CL and MA assume the same memory for source and target, we can freely move conventions with TKMR through them. The refinement lemmas for LM_{tinjp} are as follows:

LEMMA 5.3. (1) $\text{ltl}_{\text{tinjp}} \cdot \text{ltl}_{\text{ext}} \cdot \text{LM}_{\text{tinjp}} \sqsubseteq \text{LM}_{\text{tinjp}}$ (2) $\text{CL} \cdot \text{LM}_{\text{tinjp}} \cdot \text{MA} \sqsubseteq \text{CA}_{\text{tinjp}}$

Property (1) is based on the same compositionality of TKMR as property (3) in Lemma 5.1. Property (2) composes the structural relations to hide the intermediate languages.

Finally, the permutation and elimination lemmas for wt can be proved as follows:

LEMMA 5.4. (1) $ro \cdot wt \equiv wt \cdot ro$ (2) For $K \in \{\text{ext}, \text{tinjp}\}$, $c_K \cdot wt \sqsubseteq wt \cdot c_K$
 (3) $wt \cdot c_{\text{tinjp}} \cdot wt \sqsubseteq wt \cdot c_{\text{tinjp}}$

5.3 Composing the Threaded Forward Simulations

For composing individual correctness proofs, we first insert the self-simulations as presented in Table 1. The self-simulation for each intermediate language can be easily proved. The trivial vertical composition of simulation conventions by Theorem 4.4 for all passes results in:

$$\mathbb{R} = ro \cdot c_{\text{tinjp}} \cdot c_{\text{tinjp}} \cdot wt \cdot c_{\text{ext}} \cdot c_{\text{ext}} \cdot c_{\text{tinjp}} \cdot c_{\text{ext}} \cdot c_{\text{tinjp}} \cdot ro \cdot c_{\text{tinjp}} \cdot ro \cdot c_{\text{tinjp}} \cdot ro \cdot c_{\text{tinjp}} \cdot wt \cdot c_{\text{ext}} \cdot \text{CL} \cdot \text{ltl}_{\text{ext}} \cdot \text{LM}_{\text{tinjp}} \cdot \text{mach}_{\text{ext}} \cdot \text{MA}$$

Our goal is to refine \mathbb{R} into the simulation convention \mathbb{C} for the whole compiler:

$$\mathbb{C} = ro \cdot wt \cdot \text{CA}_{\text{tinjp}} \cdot \text{asm}_{\text{ext}}$$

Note that ro and wt apply to source C programs and asm_{ext} applies to the target assembly program. They are all irrelevant because ro and wt are always satisfied by any well-behaved C semantics, and any well-behaved assembly semantics is self-simulating. Therefore, the only important component is CA_{tinjp} which *directly* relates C and assembly semantics through $tinjp$. In particular, $CA_{tinjp} : C \Leftrightarrow \mathcal{A} = \langle W_{CA}, \mathbb{R}_{CA}^q, \mathbb{R}_{CA}^r, tinjp \rangle$ where $W_{CA} = (W_{tinjp}, sig, regset)$. The get and set operations for the sub world W_{tinjp} are simply defined as getting and setting the first component of $w \in W_{CA}$. \mathbb{R}_{CA}^q and \mathbb{R}_{CA}^r are almost the same as in Definition 3.4. The only difference is that the accessibilities are explicitly exposed in $tinjp$.

We prove $\mathbb{R} \sqsubseteq \mathbb{C}$ by a sequence of refined simulation conventions $\mathbb{R} \sqsubseteq \mathbb{R}_1 \sqsubseteq \dots \sqsubseteq \mathbb{R}_n \sqsubseteq \mathbb{C}$:

- (1) $ro \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{C}_{tinjp} \cdot wt \cdot \textcolor{red}{C}_{ext} \cdot \textcolor{red}{C}_{ext} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{C}_{ext} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot wt \cdot \textcolor{red}{C}_{ext} \cdot CL \cdot ltl_{ext} \cdot LM_{tinjp} \cdot mach_{ext} \cdot MA$
- (2) $ro \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{wt} \cdot \textcolor{red}{C}_{ext} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{wt} \cdot \textcolor{red}{C}_{ext} \cdot CL \cdot ltl_{ext} \cdot LM_{tinjp} \cdot mach_{ext} \cdot MA$
- (3) $\textcolor{red}{ro} \cdot \textcolor{red}{wt} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{wt} \cdot \textcolor{red}{C}_{ext} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{C}_{ext} \cdot CL \cdot ltl_{ext} \cdot LM_{tinjp} \cdot mach_{ext} \cdot MA$
- (4) $wt \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{C}_{ext} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{C}_{ext} \cdot CL \cdot ltl_{ext} \cdot LM_{tinjp} \cdot mach_{ext} \cdot MA$
- (5) $\textcolor{red}{wt} \cdot \textcolor{red}{ro} \cdot \textcolor{red}{C}_{tinjp} \cdot \textcolor{red}{C}_{ext} \cdot CL \cdot ltl_{ext} \cdot LM_{tinjp} \cdot \textcolor{red}{mach}_{ext} \cdot MA$
- (6) $ro \cdot wt \cdot CL \cdot \textcolor{red}{ltl}_{tinjp} \cdot \textcolor{red}{ltl}_{ext} \cdot \textcolor{red}{LM}_{tinjp} \cdot \textcolor{red}{MA} \cdot asm_{ext}$
- (7) $ro \cdot wt \cdot CA_{tinjp} \cdot asm_{ext}$

The red letters mark the simulation conventions refined in each step. In step (1), we use all properties in Lemma 5.1 to compose conventions involving C_{tinjp} . In step (2-3), we use the properties in Lemma 5.4 to move wt upwards, compose them using C_{tinjp} , and commute the first ro with wt . In step (4), we do the composition of TKMRs again using Lemma 5.1. In step (5), we transfer the c simulation conventions to ltl and the $mach_{ext}$ to asm using Lemma 5.2. The invariants for the source program are swapped again using property (1) in Lemma 5.4. Finally, we compose the conventions on ltl and structural conventions CL and MA with LM_{tinjp} using Lemma 5.3 to get \mathbb{C} .

The final correctness of CompCertOC is presented below:

THEOREM 5.5 (COMPILER CORRECTNESS OF COMPCERTOC).

$$\forall (M : \text{Clight})(M' : \text{Asm}), \text{CompCertOC}(M) = M' \Rightarrow \llbracket M \rrbracket \leq_{\mathbb{C}} \llbracket M' \rrbracket \wedge \llbracket M \rrbracket \geq_{\mathbb{C}} \llbracket M' \rrbracket.$$

The threaded forward simulation $\llbracket M \rrbracket \leq_{\mathbb{C}} \llbracket M' \rrbracket$ is easily proved by applying Theorem 4.5 to $\mathbb{R} \sqsubseteq \mathbb{C}$. We also get a threaded *backward* simulation $\geq_{\mathbb{C}}$ between the source and target modules. It is defined as a generalization of backward simulations in CompCert. We are able to flip threaded forward simulation into backward simulation by showing the source semantics (i.e., Clight) is *receptive* and the target (i.e., Asm) is *determinate*. The proof is almost the same as in CompCert.

6 Application and Evaluation

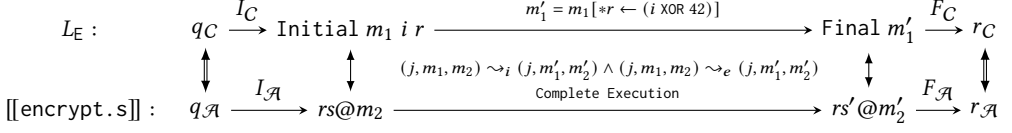
In this section, we formally verify the running example as depicted in Fig. 3. We also evaluate our Coq development of CompCertOC comparing to CompCertO [Zhang et al. 2024].

6.1 Verification of Heterogeneous Multi-threaded Modules

We first define the C-level specification L_E for `encrypt.s`. Here X_E and Y_E are empty sets and omitted because `encrypt` does not contain external calls.

Definition 6.1. LTS of L_E is defined as

$$\begin{aligned} S_E &:= \{\text{Initial } m \text{ i } r, \text{Final } m\}; \\ I_E &:= \{(\text{Vptr}(b_e, 0)[\text{int} \rightarrow \text{ptr} \rightarrow \text{void}]([i, r])@m, \text{Initial } m \text{ i } r)\}; \\ \rightarrow_E &:= \{\text{Initial } m \text{ i } r, \text{Final } m'\} \mid m' = m[*r \leftarrow (i \text{ XOR } 42)]\} \\ F_E &:= \{(\text{Final } m, \text{Vundef}@m)\}. \end{aligned}$$

Fig. 13. Simulation Diagram for $L_E \leq_{CA_{tijnjp}} [[\text{encrypt.s}]]$

We then prove the threaded forward simulation of encrypt as:

LEMMA 6.2. $L_E \leq_{\mathbb{C}} [[\text{encrypt.s}]]$.

PROOF SKETCH. Observe $\mathbb{C} = \text{ro} \cdot \text{wt} \cdot CA_{tijnjp} \cdot \text{asm}_{\text{ext}}$. The source invariants on L_E hold according to its definition. The self-simulation using asm_{ext} holds for any CompCert assembly. The proof of $L_E \leq_{CA_{tijnjp}} [[\text{encrypt}]]$ is depicted in Fig. 13. We need to verify the internal and external guarantee conditions of the whole execution which is straightforward because encrypt only changes the public memory pointed by r in both source and target semantics. Note that even if encrypt can yield back to other threads or allocate stack memory, this proof will not be more complicated by much. This is because yield is treated as a regular call. With public stack memory, we only need to modify L_E to allow it be shared and the remaining stack be protected. \square

To propagate the simulation to linked assembly code, we also prove the equivalence of syntactic and semantics linking of assembly modules:

THEOREM 6.3. $\forall (M \ M' : \text{Asm}), [[M]] \oplus [[M']] \leq_{\text{id}} [[M + M']]$.

To get the final correctness result, we first perform module linking by applying Theorem 4.7 to compose Lemma 6.2 and the compiler correctness of `client.c` and `server.c` obtained by applying Theorem 5.5. These results are then vertically extended to the linked target assembly by using Theorem 6.3 (id is absorbed using $\mathbb{C} \cdot \text{id} \sqsubseteq \mathbb{C}$). It gets us the following result:

LEMMA 6.4 (THREADED FORWARD SIMULATION FOR COMPOSED MODULES).

$$[[\text{client.c}]] \oplus [[\text{server.c}]] \oplus L_E \leq_{\mathbb{C}} [[\text{client.s} + \text{server.s} + \text{encrypt.s}]]$$

As mentioned in §5.3, we are able to flip threaded forward simulations into backward simulations by showing the composed C-level semantics is receptive. Moreover, our correctness of thread linking (Theorem 4.8) remains valid for backward simulation with additional determinacy and receptiveness requirements on the source semantics, which can be satisfied by the composed C-level semantics. Therefore, we are able to prove the closed backward simulation for our running example, which has the same definition as in the original CompCert.

THEOREM 6.5 (BACKWARD SIMULATION FOR THE COMPLETE PROGRAM).

$$\text{MT}_{\mathbb{C}}([[\text{client.c}]] \oplus [[\text{server.c}]] \oplus L_E) \geq_{\text{closed}} \text{MT}_{\mathcal{A}}([[\text{client.s} + \text{server.s} + \text{encrypt.s}]])$$

6.2 Evaluation

Our Coq development took about 10 person-months and contains 21.4k lines of code (LOC) on top of the latest CompCertO [Zhang et al. 2024]. We added 1.0k LOC to implement the multi-stack memory model based on Nominal Memory Model, 2.6k LOC to define the framework of threaded simulations with module linking. We modified the proof of 13 compilation passes (listed in Table 1 using TKMR) by adding 2.8k LOC on top of 19.3k LOC in CompCertO, which is about 14% increase. We added 7.6k LOC to verify the refinements of TKMRs and thread simulation conventions. We defined the multi-threaded semantics based on open semantics and verified the thread linking by 6.3k LOC. Finally, we verified the running example with 1.1k LOC.

Table 2. Comparison of Work on Verified Compositional Compilation

	CCC	CCM	CCO	CCTSO	CCX	CASCC	CCOC
Cooperative Concurrency	No	No	No	No	Yes	Yes	Yes
Preemptive Concurrency	No	No	No	Yes	No	Yes	No
Stack Sharing	No	No	No	Yes	No	No	Yes
Direct Refinement	No	No	Yes	Closed	No	No	Yes
Vertical Composition	Yes	RUSC	Yes	No	CCAL	Yes	Yes
Horizontal Composition	Yes	RUSC	Yes	No	CCAL	Yes	Yes
Heterogeneous Modules	Yes	Yes	Yes	No	Yes	Yes	Yes

7 Related Work

We are concerned with VCC of imperative programs with pointers. Compared to languages without pointers, this significantly complicates memory protection on shared stacks. We compare CompCertOC (CCOC) with other work on VCC not supporting concurrency, including Compositional CompCert (CCC) [Stewart et al. 2015], CompCertM (CCM) [Song et al. 2020], CompCertO (CCO) [Koenig and Shao 2021; Zhang et al. 2024], and those with concurrency, including Thread-safe CompCertX (CCX) [Gu et al. 2018], CASCompCert (CASCC) [Jiang et al. 2019a] and CompCertTSO (CCTSO) [Sevcik et al. 2013]. Table 2 summarizes this comparison where each row displays one feature the verified compilers support or not (texts in brown indicate that a weaker variant of the feature is supported). *Direct refinement* means that the source and target semantics are related at their *native* language interfaces and without exposing any intermediate semantics; it is important for compositionality and usability by third-parties. We elaborate on the comparison below.

Verified Compositional Compilation without Concurrency. Compositional CompCert (CCC) is the pioneering work on VCC of imperative programs [Stewart et al. 2015]. Its *interaction semantics* provide the simulation-based foundation for VCC which is adopted and enhanced by many projects (including this one). Its main limitation is that every language semantics and simulation must conform to the C interface, with which direct refinements cannot be supported. CompCertM (CCM) fixes this problem by introducing *Refinement Under Self-related Contexts* or RUSC [Song et al. 2020] which is the collection of all intermediate semantics and simulation relations in compilation. However, a direct refinement is still not possible because the fixed collection exposes intermediate semantics. Instead, CompCertO (CCO) defines the transitive composition of all intermediate simulations as compiler correctness [Koenig and Shao 2021]. It is later shown that this transitive composition is equivalent to a direct refinement by exploiting the memory protection of Kripke Memory Relation [Zhang et al. 2024]. Therefore, CompCertO for the first time supports all the important properties for verified compilation of heterogeneous modules (i.e., the last four rows in Table 2) and provides the basis for developing CompCertOC.

Verified Compositional Compilation with Concurrency. An early work on verified compilation of concurrent programs is CompCertTSO (CCTSO) [Sevcik et al. 2013] which is a whole-program compiler for multi-threaded programs. Although CCTSO supports stack sharing, it does not support open modules or separate compilation due to the closed nature of its compiler correctness.

Thread-safe CompCertX (CCX) [Gu et al. 2018] extends CompCertX [Gu et al. 2015]—an extension of CompCert supporting mixed C and assembly programs—with the support of open threads. Its compositionality is related to Concurrent Certified Abstraction Layers (CCAL) [Gu et al. 2016] which disallows mutual calls between modules (i.e., no callback function). Moreover, because CCX puts all memory blocks into a single space like CompCert, it cannot support thread-local stacks. To avoid the complexity brought by this limitation, CCX forbids sharing of stack data.

CASCompCert (CASCC) is the first extension of CompCert that supports general recursion and concurrency [Jiang et al. 2019a]. Being built on Compositional CompCert, it inherits its limitations including forcing languages into a C-like interface, resulting in loss of direct refinements. CASCC supports thread-local stacks by dividing CompCert’s memory state into stack and global memory space. However, it is unclear how to enable stack sharing while maintaining compositionality in CASCC. An appendix of 22 pages in the technical report of CASCC describes a solution on paper [Jiang et al. 2019b]. However, due to its complexity, this solution has not been formalized.

Therefore, CompCertOC is the first extension of CompCert that supports concurrency with stack sharing. Its correctness is described as *forward simulation* between *cooperative semantics*. The original CompCert has shown that forward simulations can be flipped into backward simulations if the source semantics is *receptive* and the target one is *determinate*. The source C semantics and the target assembly semantics in CompCertOC, CASCC and CCX all satisfy those properties. Therefore, their forward simulations can all be flipped into backward ones.

In both CASCC and CCX, a compiler correctness theorem is first stated against cooperative semantics where context switches are treated as external function calls. Then, a separate layer of verification framework on top of this compiler correctness proves that cooperative semantics is equivalent to preemptive ones. This separation makes the complexity of verified compilers much more manageable. We follow the same approach in this work. To extend our results to support preemptive semantics, we could follow CASCC to prove the equivalence between cooperative and preemptive semantics by exploiting a DRF-SC theorem to shuffle context switch points to thread primitives. We could also follow CCX to prove this equivalence in a program verification framework like CCAL. Those extra proofs are about program verification and left for future work.

Like CompCert, all stack blocks in CCX are in the same memory space [Gu et al. 2018]. Therefore, the stack blocks for different threads interleave with each other, making it difficult to link with a realistic machine model. CASCC alleviates this problem by creating separate memory space for stacks in an ad-hoc manner. In CompCertOC, we directly extend CompCert’s memory model with a multi-stack name space. Because this extension is based on a nominal interface, it does not change the memory operations or break any existing properties of CompCert’s memory model. Moreover, this treatment of stacks has already enabled elimination of pseudo instructions in CompCert’s assembly code and further compilation to real machine code as demonstrated by Wang et al. [2022]. Therefore, our multi-stack memory model is closer to how an OS views the stack (each thread has its own stack) [Gu et al. 2018] and may enable further compilation to realistic machine models.

Relationship with other Concurrency Models and Separation Logics. Although we only consider sequential consistency in this work, we believe our ideas also apply to stronger concurrency models such as linearizability. Linearizable objects can be described as transition systems where methods take effect non-deterministically and atomically [Herlihy and Wing 1990]. Therefore, our module and thread linking theorems for closed programs may be extended to support well-encapsulated objects. A more difficult problem is to support relaxed memory models. Based on existing work [Cho et al. 2022; Kang et al. 2017; Lee et al. 2020; Zha et al. 2022], we need to develop new rely-guarantee simulations that allow both stack sharing and out-of-order execution, which is left for future work.

One may wonder if the rely-guarantee reasoning can be avoided by using separation logics. From our experience, the rely-guarantee reasoning in compiler verification needs to be ignorant of program structures and contexts, so that the compiler correctness can remain fully composable. By contrast, the algebras in separation logic are more geared towards reasoning about *specific programs* (i.e., program verification). It is not clear how to directly apply them to substitute rely-guarantee reasoning. Recently, Conditional Contextual Refinement [Song et al. 2023] has combined separation logics with refinements. Its application to compiler verification needs further investigation.

Data-Availability Statement

The Coq artifact containing the formal developments in this paper is available on Zenodo [Zhang et al. 2025].

Acknowledgments

We would like to thank the anonymous referees for their helpful feedback which improved this paper significantly. This work is supported by the National Natural Science Foundation of China (NSFC) under Grant No. 62372290 and 62002217. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funding agency.

References

- Minki Cho, Sung-Hwan Lee, Dongjae Lee, Chung-Kil Hur, and Ori Lahav. 2022. Sequential reasoning for optimizing compilers under weak memory concurrency. In *Proc. 2022 ACM Conference on Programming Language Design and Implementation (PLDI'22)*. ACM, New York, NY, USA, 213–228. doi:10.1145/3519939.3523718
- Ronghui Gu, Jérémie Koenig, Tahina Ramananandro, Zhong Shao, Xiongnan(Newman) Wu, Shu-Chun Weng, Haozhong Zhang, and Yu Guo. 2015. Deep Specifications and Certified Abstraction Layers. In *Proc. 42nd ACM Symposium on Principles of Programming Languages (POPL'15)*, Sriram K. Rajamani and David Walker (Eds.). ACM, New York, NY, USA, 595–608. doi:10.1145/2775051.2676975
- Ronghui Gu, Zhong Shao, Hao Chen, Xiongnan (Newman) Wu, Jieung Kim, Vilhelm Sjöberg, and David Costanzo. 2016. CertiKOS: An Extensible Architecture for Building Certified Concurrent OS Kernels. In *Proc. 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI'16)*. USENIX Association, GA, 653–669.
- Ronghui Gu, Zhong Shao, Jieung Kim, Xiongnan (Newman) Wu, Jérémie Koenig, Vilhelm Sjöberg, Hao Chen, David Costanzo, and Tahnia Ramananandro. 2018. Certified Concurrent Abstraction Layers. In *Proc. 2018 ACM Conference on Programming Language Design and Implementation (PLDI'18)*, Jeffrey S. Foster and Dan Grossman (Eds.). ACM, New York, NY, USA, 646–661. doi:10.1145/3192366.3192381
- Maurice Herlihy and Jeannette M. Wing. 1990. Linearizability: A Correctness Condition for Concurrent Objects. *ACM Trans. Program. Lang. Syst.* 12, 3 (1990), 463–492.
- Chung-Kil Hur, Derek Dreyer, Georg Neis, and Viktor Vafeiadis. 2012. The Marriage of Bisimulations and Kripke Logical Relations. In *Proc. 39th ACM Symposium on Principles of Programming Languages (POPL'12)*, John Field and Michael Hicks (Eds.). ACM, New York, NY, USA, 59–72. doi:10.1145/2103656.2103666
- Hanru Jiang, Hongjin Liang, Siyang Xiao, Junpeng Zha, and Xinyu Feng. 2019a. Towards Certified Separate Compilation for Concurrent Programs. In *Proc. 2019 ACM Conference on Programming Language Design and Implementation (PLDI'19)*, Kathryn S. McKinley and Kathleen Fisher (Eds.). ACM, New York, NY, USA, 111–125. doi:10.1145/3314221.3314595
- Hanru Jiang, Hongjin Liang, Siyang Xiao, Junpeng Zha, and Xinyu Feng. 2019b. Towards Certified Separate Compilation for Concurrent Programs. <https://plax-lab.github.io/publications/ccc/ccc-tr.pdf>
- Jeehoon Kang, Chung-Kil Hur, Ori Lahav, Viktor Vafeiadis, and Derek Dreyer. 2017. A Promising Semantics for Relaxed-memory Concurrency. In *Proc. 44th ACM Symposium on Principles of Programming Languages (POPL'17)*, Giuseppe Castagna and Andrew D. Gordon (Eds.). ACM, New York, NY, USA, 175–189. doi:10.1145/3009837.3009850
- Jérémie Koenig and Zhong Shao. 2021. CompCertO: Compiling Certified Open C Components. In *Proc. 2021 ACM Conference on Programming Language Design and Implementation (PLDI'21)*. ACM, New York, NY, USA, 1095–1109. doi:10.1145/3453483.3454097
- Sung-Hwan Lee, Minki Cho, Anton Podkopaev, Soham Chakraborty, Chung-Kil Hur, Ori Lahav, and Viktor Vafeiadis. 2020. Promising 2.0: Global Optimizations in Relaxed Memory Concurrency. In *Proc. 2020 ACM Conference on Programming Language Design and Implementation (PLDI'20)*. ACM, New York, NY, USA, 362–376. doi:10.1145/3385412.3386010
- Xavier Leroy. 2005–2023. The CompCert Verified Compiler. <https://compcert.org/>.
- Xavier Leroy, Andrew W. Appel, Sandrine Blazy, and Gordon Stewart. 2012. *The CompCert Memory Model, Version 2*. Research Report RR-7987. INRIA. 26 pages. <https://hal.inria.fr/hal-00703441>
- Hongjin Liang, Xinyu Feng, and Ming Fu. 2012. A rely-guarantee-based simulation for verifying concurrent program transformations. In *Proc. 39th ACM Symposium on Principles of Programming Languages (POPL'12)*. ACM, New York, NY, USA, 455–468. doi:10.1145/2103621.2103711
- Rust Standard Library. 2024. std::thread::scope. <https://doc.rust-lang.org/std/thread/fn.scope.html>
- Georg Neis, Chung-Kil Hur, Jan-Oliver Kaiser, Craig McLaughlin, Derek Dreyer, and Viktor Vafeiadis. 2015. Pilsner: a Compositionally Verified Compiler for a Higher-Order Imperative Language. In *Proc. 2015 ACM SIGPLAN International*

- Conference on Functional Programming (ICFP'15)*, Kathleen Fisher and John H. Reppy (Eds.). ACM, New York, NY, USA, 166–178. doi:10.1145/2784731.2784764
- Daniel Patterson and Amal Ahmed. 2019. The Next 700 Compiler Correctness Theorems (Functional Pearl). *Proc. ACM Program. Lang.* 3, ICFP, Article 85 (August 2019), 29 pages. doi:10.1145/3341689
- Andrew M. Pitts. 2016. Nominal Techniques. *ACM SIGLOG News* 3, 1 (2016), 57–72. doi:10.1145/2893582.2893594
- Jaroslav Sevcik, Viktor Vafeiadis, Francesco Zappa Nardelli, Suresh Jagannathan, and Peter Sewell. 2013. CompCertTSO: A Verified Compiler for Relaxed-Memory Concurrency. *J. ACM* 60, 3, Article 22 (June 2013), 50 pages. doi:10.1145/2487241.2487248
- Youngju Song, Minki Cho, Dongjoo Kim, Yonghyun Kim, Jeehoon Kang, and Chung-Kil Hur. 2020. CompCertM: CompCert with C-Assembly Linking and Lightweight Modular Verification. *Proc. ACM Program. Lang.* 4, POPL, Article 23 (January 2020), 31 pages. doi:10.1145/3371091
- Youngju Song, Minki Cho, Dongjae Lee, Chung-Kil Hur, Michael Sammler, and Derek Dreyer. 2023. Conditional Contextual Refinement. *Proc. ACM Program. Lang.* 7, POPL, Article 39 (January 2023), 31 pages. doi:10.1145/3571232
- Gordon Stewart, Lennart Beringer, Santiago Cuellar, and Andrew W. Appel. 2015. Compositional CompCert. In *Proc. 42nd ACM Symposium on Principles of Programming Languages (POPL'15)*. ACM, New York, NY, USA, 275–287. doi:10.1145/2676726.2676985
- Yuting Wang, Pierre Wilke, and Zhong Shao. 2019. An Abstract Stack Based Approach to Verified Compositional Compilation to Machine Code. *Proc. ACM Program. Lang.* 3, POPL, Article 62 (January 2019), 30 pages. doi:10.1145/3290375
- Yuting Wang, Ling Zhang, Zhong Shao, and Jérémie Koenig. 2022. Verified Compilation of C Programs with a Nominal Memory Model. *Proc. ACM Program. Lang.* 6, POPL, Article 25 (January 2022), 31 pages. doi:10.1145/3498686
- Junpeng Zha, Hongjin Liang, and Xinyu Feng. 2022. Verifying Optimizations of Concurrent Programs in the Promising Semantics. In *Proc. 2022 ACM Conference on Programming Language Design and Implementation (PLDI'22)*. ACM, New York, NY, USA, 903–917. doi:10.1145/3519939.3523734
- Ling Zhang, Yuting Wang, Yalun Liang, and Zhong Shao. 2025. CompCertOC: Verified Compositional Compilation of Multi-Threaded Programs with Shared Stacks (Artifact). <https://doi.org/10.5281/zenodo.15201756>
- Ling Zhang, Yuting Wang, Jinhua Wu, Jérémie Koenig, and Zhong Shao. 2024. Fully Composable and Adequate Verified Compilation with Direct Refinements between Open Modules. *Proc. ACM Program. Lang.* 8, POPL, Article 72 (January 2024), 31 pages. doi:10.1145/3632914

Received 2024-11-15; accepted 2025-03-06