

Modular Verification of Concurrent Assembly Code with Dynamic Thread Creation and Termination

Xinyu Feng and Zhong Shao

Department of Computer Science, Yale University
New Haven, CT 06520-8285, U.S.A.
{feng, shao}@cs.yale.edu

Abstract

Proof-carrying code (PCC) is a general framework that can, in principle, verify safety properties of arbitrary machine-language programs. Existing PCC systems and typed assembly languages, however, can only handle *sequential* programs. This severely limits their applicability since many real-world systems use some form of concurrency in their core software. Recently Yu and Shao proposed a logic-based “type” system for verifying concurrent assembly programs. Their thread model, however, is rather restrictive in that no threads can be created or terminated dynamically and no sharing of code is allowed between threads. In this paper, we present a new formal framework for verifying general multi-threaded assembly code with unbounded dynamic thread creation and termination as well as sharing of code between threads. We adapt and generalize the rely-guarantee methodology to the assembly level and show how to specify the semantics of thread “fork” with argument passing. In particular, we allow threads to have different assumptions and guarantees at different stages of their lifetime so they can co-exist with the dynamically changing thread environment. Our work provides a foundation for certifying realistic multi-threaded programs and makes an important advance toward generating proof-carrying concurrent code.

Categories and Subject Descriptors F.3.1 [Logics and Meanings of Programs]: Specifying and Verifying and Reasoning about Programs; D.2.4 [Software Engineering]: Software/Program Verification — correctness proofs, formal methods; D.3.1 [Programming Languages]: Formal Definitions and Theory—semantics; D.4.5 [Operating Systems]: Reliability—verification

General Terms Languages, Verification

Keywords Concurrency Verification, Proof-Carrying Code, Rely-Guarantee, Dynamic Thread Creation

1. Introduction

Proof-carrying code (PCC) [28] is a general framework that can, in principle, verify safety properties of arbitrary machine-language programs. Existing PCC systems [29, 6, 2] and typed assembly languages (TAL) [27, 26], however, can only handle *sequential* programs. This severely limits their applicability since most real-world

systems use some form of concurrency in their core software. Certifying low-level concurrent programs is an important task because it helps increase the reliability of software infrastructure and is crucial for scaling the PCC and TAL technologies to realistic systems.

As an important first step, Yu and Shao [42]—at last year’s ICFP—proposed a certified formal framework (known as CCAP) for specifying and reasoning about general properties of concurrent programs at the assembly level. They applied the “invariance proof” technique for verifying general safety properties and the rely-guarantee methodology [23] for decomposition. They introduced a notion of “local guarantee” for supporting thread-modular verification even inside the middle of an atomic instruction sequence. Their thread model, however, is rather restrictive in that no threads can be created or terminated dynamically and no sharing of code is allowed between threads; both of these features are widely supported and used in mainstream programming languages such as C, Java, and Concurrent ML [34].

Certifying dynamic thread creation and termination turns out to be a much harder problem than we had originally anticipated [42]. Dynamic thread creation and termination imply a changing thread *environment* (*i.e.*, the collection of all live threads in the system other than the thread under concern). Such dynamic environment cannot be tracked during static verification, yet we still must somehow reason about it. For example, we must ensure that a newly created thread does not interfere with existing live threads, but at the same time we do not want to enforce non-interference for threads that have no overlap in their lifetime. Using one copy of code to create multiple threads also complicates program specification.

Existing work on the verification of concurrent programs almost exclusively uses high-level calculi (*e.g.*, CSP [21], CCS [25], TLA [24]). Also, existing work on the rely-guarantee methodology for shared-memory concurrency only supports properly nested concurrent code in the form of $\text{cobegin } P_1 \parallel \dots \parallel P_n \text{ coend}$ (which is a language construct for parallel composition where code blocks P_1, \dots, P_n execute in parallel and all terminate at the coend point). They do not support dynamic thread creation and termination.

Modularity is also needed to make verification scale. Existing work on the rely-guarantee methodology supports thread modularity, *i.e.*, different threads can be verified separately without looking into other threads’ code. However, they do not support code reuse very well. In CCAP, if a procedure is called in more than one thread, it must be verified multiple times using different specifications, one for each calling thread. We want a procedure to be specified and verified once so it can be reused for different threads.

In this paper, we propose a new framework for supporting certified multi-threaded assembly programming (CMAP). CMAP is based on a realistic abstract machine which supports dynamic thread creation with argument passing (the “fork” operation) as well as termination (the “exit” operation). Thread “join” can also be

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

ICFP’05 September 26–28, 2005, Tallinn, Estonia.

Copyright © 2005 ACM 1-59593-064-7/05/0009...\$5.00.

implemented in our language using synchronization. Our approach builds on previous work on type systems and program verification but makes the following important new contributions:

- The “fork/join” thread model is more general than “cobegin/coend” in that it supports unbounded dynamic thread creation, which poses new challenges for verification. To our knowledge, our work is the first to successfully apply the rely-guarantee method to verify concurrent programs with dynamic thread creation and termination. Our CMAP framework provides a foundation for certifying realistic multi-threaded programs and makes an important step toward generating concurrent certified code.
- The presence of dynamic threads makes it impossible to track the actual live threads during verification. This poses great challenge in enforcing the rely-guarantee condition. To solve this, we collect all dynamic threads into a single *environment* (i.e., the dynamic thread queue) and reason about the environment’s assumption and guarantee requirements as a whole. Although the dynamic thread queue cannot be tracked statically, we can update and approximate the environment’s assumption and guarantee at each program point. In fact, we can unify the concepts of the current running thread’s assumption/guarantee with its environment’s guarantee/assumption. As we will demonstrate in Sections 3 and 4, making this work in a formal framework (i.e., CMAP) is not trivial and it constitutes our main technical contribution.
- To ensure that the dynamic thread environment is well-formed, we enforce the invariant that the active threads in the system never interfere with each other. We maintain this invariant by following the approach used for type checking the dynamic data heap [27]. By combining the type-based proof techniques with the rely-guarantee based reasoning, we get a simple, extensible, and expressive framework for reasoning about the flexible “fork/join” thread model.
- We allow one copy of thread code to be activated multiple times at different places. Different “incarnations” may have different behavior, depending on the value of the thread argument. This allows us to support unbounded thread creation.
- We show how to maintain thread-modular reasoning even in the presence of dynamic thread creation and termination. Unlike CCAP, we allow each code segment to be specified independently of threads. Our work provides great support for code- and verification sharing between threads.
- We have also solved some practical issues such as thread argument passing and the saving and restoring of thread-private data at context switches. These issues are important for realistic multi-threaded programming but as far as we know have never been discussed in existing work.

We have developed CMAP and proved its soundness using the Coq proof assistant [37]. The implementation in Coq is available for download [10]. Our work makes an important advance toward building a complete PCC system for multi-threaded programs. Without formal systems such as CMAP, we cannot formally reason about concurrent assembly code. Still, more work must be done before we can construct a fully practical system. For example, a highly desirable goal is a high-level language with concise human-readable annotations that can be automatically compiled into CMAP programs and proofs. We leave this as future work.

In the rest of this paper, we first give an overview of the rely-guarantee-based reasoning and a detailed explanation of the key challenges in verifying multi-threaded assembly code (Section 2). We then give an informal description of our approach to address these problems in Section 3, and present our work on CMAP with

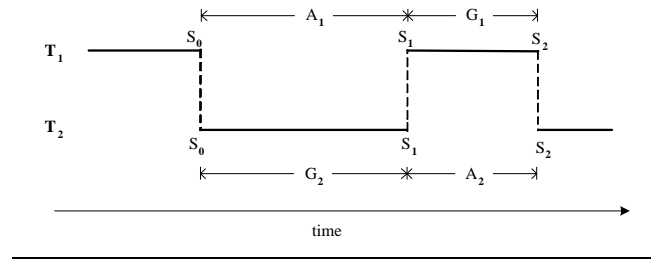


Figure 1. Rely-guarantee-based reasoning

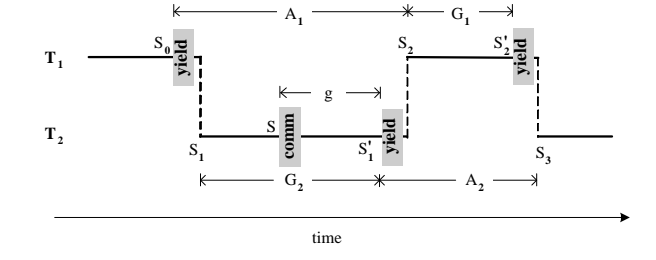


Figure 2. R-G in a non-preemptive setting

formal semantics in Section 4. We use a few examples to illustrate CMAP-based program verification in Section 5. Finally we discuss related work and conclude.

2. Background and Challenges

2.1 Rely-Guarantee-Based Reasoning

The rely-guarantee (R-G) proof method [23] is one of the best-studied approaches to the compositional verification of shared-memory concurrent programs. Under the R-G paradigm, every thread is associated with a pair (A, G) , with the meaning that if the environment (i.e., the collection of all of the rest threads) satisfies the assumption A , the thread will meet its guarantee G to the environment. In the shared-memory model, the assumption A of a thread describes what atomic transitions may be performed by other threads, while the guarantee G of a thread must hold on every atomic transition of the thread. They are typically modeled as predicates on a pair of states, which are often called *actions*.

For instance, in Figure 1 we have two interleaving threads T_1 and T_2 . T_1 ’s assumption A_1 adds constraints on the transition (S_0, S_1) made by the environment (T_2 in this case), while G_1 describes the transition (S_1, S_2) , assuming the environment’s transition satisfies A_1 . Similarly A_2 describes (S_1, S_2) and G_2 describes (S_0, S_1) .

We need two steps to reason about a concurrent program consisting of T_1 and T_2 . First, we check that there is no interference between threads, i.e., that each thread’s assumption can be satisfied by its environment. In our example, non-interference is satisfied as long as $G_1 \Rightarrow A_2$ (a shorthand for $\forall S, S'. G_1(S, S') \Rightarrow A_2(S, S')$), and $G_2 \Rightarrow A_1$. Second, we check that T_1 and T_2 do not lie, that is, they satisfy their guarantee as long as their assumption is satisfied. As we can see, the first step only uses the specification of each thread, while the second step can be carried out independently without looking at other threads’ code. This is how the R-G paradigm achieves thread-modularity.

2.2 R-G in Non-Preemptive Thread Model

CMAP adopts a non-preemptive thread model, in which threads yield control voluntarily with a `yield` instruction, as shown in Figure 2. The preemptive model can be regarded as a special case of

```

Variables:          Initially:
nat[100] data;     data[i] = ni, 0 ≤ i < 100

main1 :
  data[0] := f(0);
  fork(chld, 0);
  ⋮
  data[99] := f(99);
  fork(chld, 99);
  ⋮

main2 :
  nat i := 0;
  while(i < 100){
    data[i] := f(i);
    fork(chld, i);
    i := i + 1;
  }
  ⋮

void chld(int x){
  data[x] := g(x, data[x]);
}

```

Figure 3. Loop: high-level program

the non-preemptive one, in which an explicit `yield` is used at every program point. Also, on real machines, programs might run in both preemptive and non-preemptive settings: preemption is usually implemented using interrupts; a program can disable the interrupt to get into non-preemptive setting.

An “atomic” transition in a non-preemptive setting then corresponds to a sequence of instructions between two yields. For instance, in Figure 2 the state pair (S_2, S'_2) corresponds to an atomic transition of thread T_1 . A difficulty in modeling concurrency in such a setting is that the effect of an “atomic” transition cannot be completely captured until the end. For example, in Figure 2, the transition (S_1, S'_1) should satisfy G_2 . But when we reach the intermediate state S , we have no idea of what the whole transition (*i.e.*, (S_1, S'_1)) will be. At this point, neither (S_1, S) nor (S, S'_1) need satisfy G_2 . Instead, it may rely on the remaining commands (the commands between `comm` and `yield`, including `comm`) to complete an adequate state transition. In CCAP [42], a “local guarantee” g is introduced for every program point to capture further state changes that must be made by the following commands before it is safe for the current thread to yield control. For instance, the local guarantee g attached to `comm` in Figure 2 describes the transition (S, S'_1) .

2.3 Challenges for Dynamic Thread Creation

To prove safety properties of multi-threaded programs, the key problem is to enforce the invariant that *all executing threads must not interfere with each other*. As mentioned in Section 2.1, threads *do not interfere* (or they satisfy the *non-interference* or *interference-free* property) only if each thread’s assumption is implied by the guarantee of all other threads¹. For languages that do not support dynamic thread creation, the code for each thread corresponds to exactly one executing thread. Using the rely-guarantee method, we can assign an assumption and guarantee to each thread code and enforce non-interference by checking all of these assumptions and guarantees, as is done in [42] and [13]. However, the following example shows that this simple approach cannot support dynamic thread creation and multiple “incarnation” of the thread code.

In Figure 3, the high-level pseudo code (using C-like syntax) shows the code for (two versions of) a main thread and child threads. The main thread initializes 100 pieces of data using some function f , and distributes them to 100 child threads that will work on their own data (by applying a function g) in parallel. The `fork` function creates a child thread that will execute the function pointed

¹We will formalize the Non-Interference property in Section 4.4.

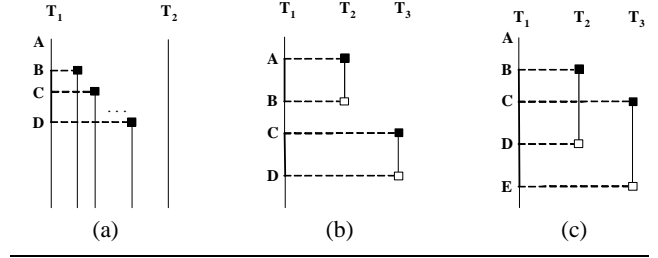


Figure 4. Interleaving of threads

to by the first argument. The second argument of `fork` is passed to the function as argument. The thread `main1` does this in sequential code while `main2` uses code with a loop. We assume that the high level code runs in preemptive mode. In other words, there is an implicit `yield` at any program point.

It is easy to see that both versions are “well-behaved” as long as the function g has no side effects, and all other threads in the rest of the system do not update the array of `data`. However, the simple approach used in [42] and [13] even cannot provide a specification for such trivial code.

1. Figure 4 (a) illustrates the execution of `main1` (time goes downwards). When doing data initialization (at stage A-B, meaning from point A to point B), the main thread needs to assume that no other threads in the environment (say, T_2) can change the array of `data`. However, the composition of the main thread’s environment changes after a child thread is created. The assumption used at A-B is no longer appropriate for this new environment since the first child thread will write to `data[0]`. And the environment will keep changing with the execution of the main thread. How can we specify the `main1` thread to support such a dynamic thread environment?

One possible approach is that the main thread relaxes its assumption to make exceptions for its child threads. However, it is hard to specify the parent-child relationship. Another approach is to use something like the program counter in the assumption and guarantee to indicate the phase of computation. This means the specification of the main thread is sensitive to the implementation. Also the program structure of the main thread has to be exposed to the specification of the child threads, which compromises modularity. The worst thing is that this approach simply won’t work for the version `main2`.

2. Since multiple child threads are created, we must make sure that there is no interference between these children. It is easy for the above example since we can let the assumption and guarantee of the `chld` code be parameterized by its argument, and require $G_i \Rightarrow A_j$ given $i \neq j$. However, this approach cannot be generalized for threads that have dummy arguments and their behavior does not depend on their arguments at all. In this case $G_i \equiv G_j$ and $A_i \equiv A_j$ for any i and j . Then requiring $G_i \Rightarrow A_j$ is equivalent to requiring $G_i \Rightarrow A_i$, which cannot be true in general, given the meaning of assumptions and guarantees described in section 2.1. Do we need to distinguish these two kinds of threads and treat them differently? And how do we distinguish them?
3. Another issue introduced by dynamic thread creation, but not shown in this example program, is that the lifetimes of some threads may not overlap. In the case shown in Figure 4 (b), the lifetimes of T_2 and T_3 do not overlap and we should not statically enforce non-interference between them. Again, how can we specify and check the interleaving of threads, which can be as complex as shown in Figure 4 (c)?

In the next section we'll show how these issues are resolved in our development of CMAP.

3. Our Approach

In the rest of this paper, to distinguish the executing thread and the thread code, we call the dynamically running thread the “dynamic thread” and the thread code the “static thread”. In Figure 3 the function `child` is the static child thread, from which 100 dynamic child threads are activated.

As explained in Section 2.3, the approach that requiring non-interference of static threads is too rigid to support dynamic thread creation. Our approach, instead, enforces the thread non-interference in a “lazy” way. We maintain a dynamic thread queue which contains all of the active threads in the system. When a new thread is created, it is added to the dynamic thread queue. A thread is removed from the queue when its execution terminates. We also require that, when specifying the program, each static thread be assigned an assumption/guarantee pair. However, we do not check for non-interference between static thread specifications. Instead, each dynamic thread is also assigned an assumption and guarantee at the time of creation, which is an instantiation of the corresponding static thread specification with the thread argument. We require that dynamic threads do not interfere with each other, which can be checked by inspecting their specifications.

Our approach is very flexible in that each dynamic thread does not have to stick to one specification during its lifetime. When its environment changes, its specification can be changed accordingly. As long as the new specification does not introduce interference with other existing dynamic threads, and the subsequent behavior of this thread satisfies the new specification, the whole system is still interference-free. In this way, we can deal with the changing environment resulting from dynamic thread creation and termination. Problem 1 in Section 2.3 can be solved now.

If the lifetimes of two threads do not overlap they will not show up in the system at the same time. Therefore we do not need to check for interference at all. Also, since each dynamic thread has its own specification, we no longer care about the specification of the corresponding static thread. Therefore problems 2 and 3 shown in Section 2.3 are no longer an issue in our approach.

3.1 Typing The Dynamic Thread Queue

We define the dynamic thread queue \mathbb{Q} as a set of thread identifiers τ_i , and the assignment Θ of assumption/guarantee pairs to dynamic threads as a partial mapping from τ_i to (A_i, G_i) ². The queue \mathbb{Q} is “well-typed” with regard to Θ if:

- $\mathbb{Q} = \text{dom}(\Theta)$, where $\text{dom}(\Theta)$ is the domain of Θ ;
- threads in \mathbb{Q} do not interfere, i.e., $\forall \tau_i, \tau_j. \tau_i \neq \tau_j \Rightarrow (G_i \Rightarrow A_j)$; and
- each dynamic thread τ_i is “well-behaved” with regard to (A_i, G_i) , i.e., if A_i is satisfied by the environment, τ_i 's execution does not get stuck and satisfies G_i .

Therefore, the invariant we need to maintain is that during the execution of the program, for the queue \mathbb{Q} at each step there exists a Θ such that \mathbb{Q} is well-typed with regard to Θ . In fact, we do not require Θ to be part of the program specification. We only need to ensure that there *exists* such a Θ at each step, which may be changing.

The content of the thread queue keeps changing, so how can we track the set of threads in the queue by a static inspection of the program? Here we follow the approach used for type-checking the

²This is a temporary formulation to illustrate our basic idea. We will use different definitions in our formal development of CMAP in Section 4.

dynamic data heap [27], which is dynamically updated by the store instruction and extended by the alloc instruction. We can ensure our invariant holds as long as the following conditions are satisfied:

- At the initial state (when the program starts to run) we can find a Θ to type-check the initial \mathbb{Q} . Usually the initial \mathbb{Q} only contains the main thread, which will start to execute its first instruction, so we can simply assign the assumption/guarantee in the specification of the static main thread to the dynamic main thread.
- For each instruction in the program, assume that before the execution of the instruction there is a Θ such that \mathbb{Q} is well typed. Then as long as certain constraints are satisfied to execute the instruction, there must exist a Θ' that can type check the resulting \mathbb{Q}' . For most instructions which do not change the content of the thread queue, this condition can be trivially satisfied. We are only interested in the “fork” and “exit” operation which will change the content of \mathbb{Q} .

For the “exit” instruction, the second condition can also be satisfied by the following lemma which can be trivially proven.

Lemma 3.1 (Thread Deletion)

If \mathbb{Q} is well-typed with regard to Θ , then for all $\tau \in \mathbb{Q}$ we know $\mathbb{Q} \setminus \{\tau\}$ is well-typed with regard to $\Theta \setminus \{\tau\}$.

For the “fork” instruction, things are trickier. We need to ensure that the new child thread does not interfere with threads in the parent thread's environment. We also require that the parent thread does not interfere with the child thread. The following lemma ensures the first requirement.

Lemma 3.2 (Queue Extension I)

Suppose \mathbb{Q} is well-typed with regard to Θ and the current executing thread is τ . If

- $\Theta(\tau) = (A, G)$;
- a new thread τ' is created by τ ;
- (A', G') is the instantiation of the corresponding static thread specification by the thread argument;
- $A \Rightarrow A'$ and $G' \Rightarrow G$;

then $\mathbb{Q}' \cup \{\tau'\}$ is well-typed with regard to $\Theta' \{\tau' \rightsquigarrow (A', G')\}$, where $\mathbb{Q}' = \mathbb{Q} \setminus \{\tau\}$ and $\Theta' = \Theta \setminus \{\tau\}$.

Here \mathbb{Q}' is the environment of the current thread τ . Since τ does not interfere with its environment (because \mathbb{Q} is well-typed), we know that its assumption A is an approximation of what the environment can guarantee (G_e), and similarly that G is an approximation of the environment's assumption (A_e). By this interpretation, we can unify the concepts of the current running thread's assumption/guarantee with its environment's guarantee/assumption. To ensure the new thread τ' does not interfere with τ 's environment, we need $G' \Rightarrow A_e$ and $G_e \Rightarrow A'$, which can be derived from $G' \Rightarrow G$ and $A \Rightarrow A'$.

Still, we need to ensure that thread τ does not interfere with τ' . As mentioned above, A and G are approximations of G_e and A_e , respectively. Since the environment is extended with the child thread, the guarantee G'_e for the new environment is $G_e \vee G'$ and the assumption for the new environment A'_e is $A_e \wedge A'$. We want to change A and G correspondingly to reflect the environment change. First, the following lemma says that the specification of a dynamic thread can be changed during its lifetime.

Lemma 3.3 (Queue Update)

Suppose \mathbb{Q} is well-typed with regard to Θ and that the current executing thread is τ . If

- $\Theta(\tau) = (A, G)$;
- $G'' \Rightarrow G$ and $A \Rightarrow A''$;
- the subsequent behavior of the current thread satisfies (A'', G'') ;

then \mathbb{Q} is well-typed with regard to $\Theta\{\tau \rightsquigarrow (A'', G'')\}$.

Now we can change the specification of the parent thread τ to let it reflect the change of the environment.

Lemma 3.4 (Queue Extension II)

Suppose \mathbb{Q} is well-typed with regard to Θ and the current executing thread is τ . If

- $\Theta(\tau) = (A, G)$;
- a new thread τ' is created by τ ;
- (A', G') is the instantiation of the corresponding static thread specification by the thread argument;
- $(A \Rightarrow A') \wedge (G' \Rightarrow G)$;
- the remainder behavior of the thread τ also satisfies $(A \vee G', G \wedge A')$;

then $\mathbb{Q} \cup \{\tau'\}$ is well-typed with regard to $\Theta\{\tau' \rightsquigarrow (A', G'), \tau \rightsquigarrow (A \vee G', G \wedge A')\}$.

If τ later creates another thread τ'' , because the specification of τ already reflects the existence of τ' , by Lemma 3.2 we know that τ'' will not interfere with τ' as long as its specification satisfies the constraints. Therefore we do not need to explicitly check that τ' and τ'' are activated from the same static thread or that multiple activations of a static thread do not interfere with each other.

These lemmas are used to prove the soundness of CMAP. They are somewhat similar to the heap update and heap extension lemmas used in TAL's soundness proof [27]. People familiar with the traditional rely-guarantee method may feel this is nothing but the parallel composition rule used to support nested `cobegin/coend`. However, by combining the invariant-based proof technique used by type systems and the traditional rely-guarantee method, we can now verify multi-threaded assembly program with a more flexible program structure than the `cobegin/coend` structure. In particular, programs which are not properly nested, as shown in Figure 4(c) and the `main2` program in Figure 3, can be supported in our system. This is one of the most important contributions of this paper.

3.2 Parameterized Assumption/Guarantee

The assumptions and guarantees are interfaces between threads, which should only talk about shared resources. As we allow multiple activations of static threads, the behavior of a dynamic thread may depend on its arguments, which is the thread's private data. Therefore, to specify a static thread, the assumption and guarantee need to be parameterized over the thread argument.

In our thread model, the flat memory space is shared by all threads, and as in most operating systems, the register file is saved at the moment of context switch. Therefore the register file is thread-private data. The thread argument is stored in a dedicated register.

Rather than letting the assumption and guarantee be parameterized over the thread argument, we let them be parameterized over the whole register file. This makes our specification language very expressive. For instance, we allow the dynamic thread to change its specification during its lifetime to reflect change in the environment. If the thread has private data that tracks the composition of the environment, and its specification is parameterized by such data, then its specification automatically changes with the change of the data, which in turn results from the change of the thread environment. This is the key technique we use to support unbounded dynamic thread creation, as shown in the program `main2` in Figure 3.

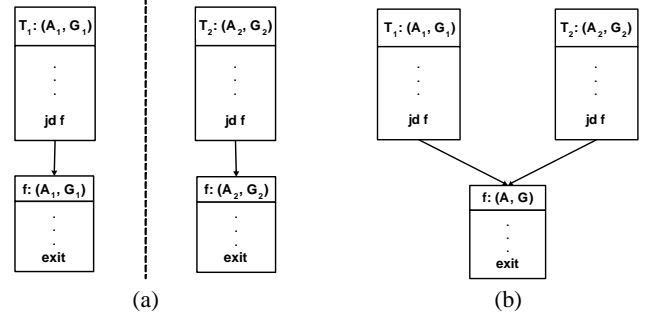


Figure 5. Code sharing between different threads

3.3 Support of Modular Verification

The rely-guarantee method supports thread modularity well, *i.e.*, code of one thread can be verified independently without inspecting other threads' code. However, it does not have good support for code reuse. In CCAP, each thread has its own code heap and there is no sharing of code between threads. As Figure 5(a) shows, if an instruction sequence is used by multiple threads, it has multiple copies in different threads' code heaps, each copy verified with regard to the specifications of these threads.

Based on the rely-guarantee method, the thread modularity is also supported in our system. In addition, using our "lazy" checking of thread non-interference, and by the queue update lemma, we can allow instruction sequences to be specified independently of their calling thread, thus achieving better modularity.

As shown in Figure 5(b), we assign an assumption/guarantee pair to the specification of each instruction sequence, and require the instruction sequence be well-behaved with regard to its own assumption/guarantee. Similar to threads, the instruction sequence is well-behaved if, when it is executed by a dynamic thread, its execution is safe and satisfies its guarantee, as long as other dynamic threads in the environment satisfy the assumption. The instruction sequence only needs to be verified once with respect to its own specification, and can be executed by different threads as long as certain constraints are satisfied.

Intuitively, it is safe for a dynamic thread τ with specification (A_i, G_i) to execute the instruction sequence labeled by f as long as executing it does not require a stronger assumption than A_i , nor does it violate the guarantee G_i . Therefore, if the specification of f is (A, G) , τ can call f as long as $A_i \Rightarrow A$ and $G \Rightarrow G_i$. The intuition is backed up by our queue update lemma.

4. CMAP

The language CMAP is based on an "untyped" low-level abstract machine supporting multi-threaded programs with dynamic thread creation and argument passing. The "type" system of CMAP uses the calculus of inductive constructions (CiC) [32] to essentially support reasoning in higher-order predicate logic.

4.1 The Abstract Machine

Figure 6 shows the definition of our abstract machine.

A CMAP program (corresponding to a complete machine state) is made up of an updatable state \mathbb{S} (which is made up of the shared memory \mathbb{M} and the register file \mathbb{R}), a dynamic thread queue \mathbb{Q} , two shared code heaps \mathbb{C} (for basic code blocks) and \mathbb{T} (for thread entries), and the current instruction sequence \mathbb{I} of the currently executing thread. Here \mathbb{C} and \mathbb{T} can be merged, but conceptually it is cleaner to have them separated because the specification of \mathbb{T} and is different from that of \mathbb{C} (See Section 4.3: for \mathbb{T} we do not need to specify a local guarantee). Memory is a partial mapping from

$((M, \mathbb{R}), \mathbb{Q}, T, C, I) \mapsto \mathbb{P}$	
if $I =$	then $\mathbb{P} =$
fork $h, r; I''$	$((M, \mathbb{R}), \mathbb{Q}\{t \rightsquigarrow (\mathbb{R}', I')\}, T, C, I'')$ where $I' = T(h), t \notin \text{dom}(\mathbb{Q}), t \neq \mathbb{R}(rt)$, and $\mathbb{R}' = \{r_0 \rightsquigarrow _, \dots, r_{15} \rightsquigarrow _, rt \rightsquigarrow t, ra \rightsquigarrow \mathbb{R}(r)\}$
yield; I''	$((M, \mathbb{R}'), (\mathbb{Q}\{\mathbb{R}(rt) \rightsquigarrow (\mathbb{R}, I'')\}) \setminus \{t\}, T, C, I'')$ where $t \in \text{dom}(\mathbb{Q})$ and $(\mathbb{R}', I') = \mathbb{Q}(t)$ or $t = \mathbb{R}(rt)$ and $(\mathbb{R}', I') = (\mathbb{R}, I'')$
exit	$((M, \mathbb{R}'), \mathbb{Q} \setminus \{t\}, T, C, I'')$ where $t \in \text{dom}(\mathbb{Q})$ and $(\mathbb{R}', I') = \mathbb{Q}(t)$
jd f	$((M, \mathbb{R}), \mathbb{Q}, T, C, I')$ where $I' = C(f)$
bgt $r_s, r_t, f; I''$	$((M, \mathbb{R}), \mathbb{Q}, T, C, I'')$ if $\mathbb{R}(r_s) \leq \mathbb{R}(r_t)$, $((M, \mathbb{R}), \mathbb{Q}, T, C, I')$ otherwise, where $I' = C(f)$
beq $r_s, r_t, f; I''$	$((M, \mathbb{R}), \mathbb{Q}, T, C, I'')$ if $\mathbb{R}(r_s) \neq \mathbb{R}(r_t)$, $((M, \mathbb{R}), \mathbb{Q}, T, C, I')$ otherwise, where $I' = C(f)$
c; I'' for remaining cases of c	$(\text{Next}(c, (M, \mathbb{R})), \mathbb{Q}, T, C, I'')$

Figure 8. Operational semantics of CMAP

(Program)	$\mathbb{P} ::= (S, \mathbb{Q}, T, C, I)$
(State)	$S ::= (M, \mathbb{R})$
(Memory)	$M ::= \{l \rightsquigarrow w\}^*$
(RegFile)	$\mathbb{R} ::= \{r \rightsquigarrow w\}^*$
(Register)	$r ::= r_0 \mid r_1 \mid \dots \mid r_{15} \mid rt \mid ra$
(CdHeap)	$C ::= \{f \rightsquigarrow I\}^*$
(Labels)	$f, l ::= n \text{ (nat nums)}$
(WordVal)	$w ::= n \text{ (nat nums)}$
(TEnties)	$T ::= \{h \rightsquigarrow I\}^*$
(TQueue)	$\mathbb{Q} ::= \{t \rightsquigarrow (\mathbb{R}, I)\}^*$
(THandles)	$h ::= n \text{ (nat nums)}$
(ThrdID)	$t ::= n \text{ (nat nums)}$
(InstrSeq)	$I ::= c; I \mid \text{jd } f \mid \text{exit}$
(Cmmd)	$c ::= \text{yield} \mid \text{fork } h, r \mid \text{add } r_d, r_s, r_t \mid \text{sub } r_d, r_s, r_t$ $\mid \text{movi } r_d, w \mid \text{bgt } r_s, r_t, f \mid \text{beq } r_s, r_t, f$ $\mid \text{ld } r_d, r_s(w) \mid \text{st } r_d(w), r_s$

Figure 6. The abstract machine

if $c =$	then $\text{Next}(c, (M, \mathbb{R})) =$
add r_d, r_s, r_t	$(M, \mathbb{R}\{r_d \rightsquigarrow \mathbb{R}(r_s) + \mathbb{R}(r_t)\})$
sub r_d, r_s, r_t	$(M, \mathbb{R}\{r_d \rightsquigarrow \mathbb{R}(r_s) - \mathbb{R}(r_t)\})$
movi r_d, w	$(M, \mathbb{R}\{r_d \rightsquigarrow w\})$
ld $r_d, r_s(w)$	$(M, \mathbb{R}\{r_d \rightsquigarrow M(\mathbb{R}(r_s) + w)\})$ where $(\mathbb{R}(r_s) + w) \in \text{dom}(M)$
st $r_d(w), r_s$	$(M\{(\mathbb{R}(r_d) + w) \rightsquigarrow \mathbb{R}(r_s)\}, \mathbb{R})$ where $(\mathbb{R}(r_d) + w) \in \text{dom}(M)$

Figure 7. Auxiliary state update function

memory locations to word-sized values. The register file \mathbb{R} maps registers to word-sized values. In our machine, there are 16 general purpose registers ($r_0 - r_{15}$) and two special registers (rt and ra) that hold the current thread id and the thread argument. Code heaps map code labels to instruction sequences, which are lists of instructions terminated by a `jd` or `exit` instruction. Code labels pointing to thread entries are also called *thread handles*. T maps thread handles to thread entries (instruction sequences from which a thread starts to execute). Thread entries are also called *static threads*. The current instruction sequence I plays the role of the program counter of the current executing thread.

For simplicity, we just model the queue of ready threads, which is the dynamic thread queue \mathbb{Q} that maps the dynamic thread id to an execution context of a thread. The dynamic thread id t is a natural number generated randomly at run time. The thread execution context includes the snapshot of the register file and the program point where the thread will resume its execution. Note that \mathbb{Q} does not contain the current executing thread, which is different from the dynamic thread queue used in Section 3.

The instruction set of CMAP just contains the most basic and common assembly instructions. It also includes primitives `fork`, `exit` and `yield` to support multi-threaded programming which can be viewed as system calls to a thread library. We do not have a `join` instruction because thread join can be implemented using synchronization. Readers who are eager to see a CMAP program can take a quick look at Figure 10 in Section 5.1 (ignore the specifications and annotations in program for the time being), which is the CMAP implementation of programs `main2` and `ch1d` shown in Figure 3.

The execution of CMAP programs is modeled as small-step transitions from one program to another, *i.e.*, $\mathbb{P} \mapsto \mathbb{P}'$. Figure 8 defines the program transition function.

The primitive `fork` creates a new thread using the static thread h , and passes the value $\mathbb{R}(r)$ to it as the argument. The new thread will be assigned a fresh thread id and placed in the dynamic thread queue waiting for execution. The current thread continues with the subsequent instructions.

At a `yield` instruction, the current thread will give up the control of the machine. Its execution context is stored in \mathbb{Q} . The scheduler will pick one thread non-deterministically from the thread queue (which might be the yielding thread itself), restore its execution context, and execute it.

The `exit` instruction terminates the execution of the current thread and non-deterministically selects a thread from the thread queue. Here we have the implicit assumption that there is always an idle thread in the thread queue that never changes the state or terminates, ensuring that the thread queue will never be empty.

Semantics for the rest of the instructions are standard. The “next state” function defined in Figure 7 describes the effects of some instructions on the state.

4.2 The Meta-Logic

To encode the specification and proofs, we use the calculus of inductive constructions (CiC) [37, 32], which is an extension of the calculus of constructions (CC) [7] with inductive definitions. CC corresponds to Church’s higher-order predicate logic via the Curry-Howard isomorphism.

CiC has been shown strongly normalizing [38], hence the corresponding logic is consistent. It is supported by the Coq proof assistant [37], which we have used to implement CMAP.

In the remainder of this paper, we will mainly use the more familiar mathematical and logical notations, instead of strict CiC or Coq representation. We use *Prop* to denote the type of all propositions. No knowledge of CiC and Coq is required to understand them.

4.3 Program Specifications

The verification constructs of CMAP are defined in Figure 9. The program specification Φ is a global invariant (*Inv*), a static thread

(ProgSpec)	$\Phi ::= (Inv, \Delta, \Psi)$
(ThrdSpec)	$\Delta ::= \{h \rightsquigarrow \theta\}^*$
(ThrdType)	$\theta ::= (p, A, G)$
(CdHpSpec)	$\Psi ::= \{f \rightsquigarrow (p, g, A, G)\}^*$
(ActTSpec)	$\Theta ::= \{t \rightsquigarrow (p, A, G)\}^*$
(Invariant)	$Inv \in Mem \rightarrow Prop$
(Assertion)	$p \in State \rightarrow Prop$
(Assumption)	$A \in RegFile \rightarrow Mem \rightarrow Mem \rightarrow Prop$
(Guarantee)	$G, g \in RegFile \rightarrow Mem \rightarrow Mem \rightarrow Prop$

Figure 9. Verification constructs of CMAP

specification Δ , and a code heap specification Ψ . The invariant Inv is a programmer specified predicate, which implies a safety property of concern. It must hold throughout the execution of the program. The static thread specification Δ contains a specification θ for each static thread in \mathbb{T} . Each θ contains a precondition p to invoke this thread, and an assumption A and guarantee G for this thread with regard to the environment at its creation time.

A code heap specification Ψ assigns a quadruple (p, g, A, G) to each instruction sequence. The assertion p is the precondition to execute the code sequence. The local guarantee g , as introduced in Section 2.2, describes a valid state transition – it is safe for the current thread to yield control only after making a state transition described by g . As explained in Section 3.3., we assign a pair of A and G to each instruction sequence as part of its specification. The instruction sequence can be verified with regard to its own specification without knowing which thread executes it. Here the A and G reflect knowledge of the dynamic thread environment at the time the instruction sequence is executed.

The global invariant Inv is a CiC term of type $Mem \rightarrow Prop$, *i.e.*, a predicate over memory. Inv does not specify the register file, which contains thread-private data and keeps changing. In contrast to Inv , assertions (p) are predicates over the whole state. Assumptions and guarantees (*e.g.*, A , G and g) are CiC terms with type $RegFile \rightarrow Mem \rightarrow Mem \rightarrow Prop$, which means predicates over a register file and two instances of memory. Assumptions and guarantees specify the behavior of threads by describing the change of shared memory. As mentioned in Section 3.2, they are parameterized over the register file, which contains the private data of threads.

We also define the specification Θ of active threads in the thread queue \mathbb{Q} . Within Θ , each triple (p, A, G) describes a dynamic thread at its yield point (or at the beginning if it has just forked). The assertion p gives the constraint of the state when the thread gets control back to execute its remaining instructions. The assumption and guarantee used by the thread at the yield point are given by A and G . As we said in Section 3.1, the A and G of each dynamic thread may change during the lifetime of the thread. Notice that Θ is not part of the program specification. It is used only in the soundness proof.

4.4 Inference Rules

We use the following judgement forms to define the inference rules:

$\Phi; \Theta; (p, g, A, G) \vdash \mathbb{P}$	(well-formed program)
$\Phi; \Theta; (g, S) \vdash \mathbb{Q}$	(well-formed dynamic threads)
$\Phi \vdash \mathbb{T}$	(well-formed static threads)
$\Phi \vdash \mathbb{C}$	(well-formed code heap)
$\Phi; (p, g, A, G) \vdash \mathbb{I}$	(well-formed instr. sequence)

Before introducing the inference rules, we first define some shorthands in Table 1 to simplify our presentation.

Well-formed programs. The PROG rule shows the invariants that need to be maintained during program transitions.

Representation	Definition
$\mathbb{S}(x)$	$\mathbb{R}(x)$ where $\mathbb{S} = (M, \mathbb{R})$
$Inv \wedge p$	$\lambda(M, \mathbb{R}). Inv M \wedge p (M, \mathbb{R})$
$g \mathbb{S} M'$	$g \mathbb{R} M M'$, where $\mathbb{S} = (M, \mathbb{R})$
$A \mathbb{S} M', G \mathbb{S} M'$	similar to $g \mathbb{S} M'$
$p \circ c$	$\lambda \mathbb{S}. p (\text{Next}(c, \mathbb{S}))$
$g \circ c$	$\lambda \mathbb{S}. \lambda M'. g (\text{Next}(c, \mathbb{S})) M'$
$A \circ c, G \circ c$	similar to $g \circ c$
$p \Rightarrow p'$	$\forall \mathbb{S}. p \mathbb{S} \Rightarrow p' \mathbb{S}$
$p \Rightarrow g$	$\forall \mathbb{R}, M. p (M, \mathbb{R}) \Rightarrow g \mathbb{R} M M$
$p \Rightarrow g \Rightarrow g'$	$\forall \mathbb{S}, M'. p \mathbb{S} \Rightarrow g \mathbb{S} M' \Rightarrow g' \mathbb{S} M'$
$A \Rightarrow A'$	$\forall \mathbb{R}, M, M'. A \mathbb{R} M M' \Rightarrow A' \mathbb{R} M M'$
$G \Rightarrow G'$	$\forall \mathbb{R}, M, M'. G \mathbb{R} M M' \Rightarrow G' \mathbb{R} M M'$
$p \xrightarrow{\mathbb{R}} g \xrightarrow{\mathbb{R}'} p'$	$\forall M, M'. p (M, \mathbb{R}) \Rightarrow g \mathbb{R} M M' \Rightarrow p (M', \mathbb{R}')$
$(Inv, p \circ c, A, \mathbb{R})$	$(Inv \wedge p) \xrightarrow{\mathbb{R}} A \xrightarrow{\mathbb{R}} p$

Table 1. Assertion definitions and syntactic sugar

$$\begin{array}{l}
(Inv, \Delta, \Psi) = \Phi \quad (M, \mathbb{R}) = \mathbb{S} \quad t = \mathbb{R}(t) \\
\Phi \vdash \mathbb{T} \quad \Phi \vdash \mathbb{C} \quad (Inv \wedge p) \mathbb{S} \quad \Phi; (p, g, A, G) \vdash \mathbb{I} \\
\Phi; \Theta; (g, \mathbb{S}) \vdash \mathbb{Q} \quad \text{NI}(\Theta \{t \rightsquigarrow (p, A, G)\}; \mathbb{Q} \{t \rightsquigarrow (\mathbb{R}, \mathbb{I})\}) \\
\hline
\Phi; \Theta; (p, g, A, G) \vdash (\mathbb{S}, \mathbb{Q}, \mathbb{T}, \mathbb{C}, \mathbb{I}) \quad (\text{PROG})
\end{array}$$

The well-formedness of a program is judged with respect to the program specification Φ , the dynamic thread specification Θ , and the specification of the current executing thread (p, g, A, G) . Compared with the triples in Θ , we need the local guarantee g here to specify the transition which the current thread must make before it yields control.

In the first line, we give the composition of the program specification, the current state and the current thread id. Here we use a pattern match representation, which will be used in the rest of this paper.

We require the code, including the thread entry points \mathbb{T} and the code heap \mathbb{C} , always be well-formed with respect to the program specification. Since Φ , \mathbb{T} and \mathbb{C} do not change during program transitions, the check for the first two premises in line 2 can be done once and for all.

The next premise shows the constraints on the current state \mathbb{S} : it must satisfy both the global invariant Inv and the assertion p of the current thread. The last premise in line 2 essentially requires it be safe for the current thread to execute the remainder instruction sequence \mathbb{I} .

Premises in line 3 require the well-formedness of dynamic threads in \mathbb{Q} , which is checked by the rule DTHRDS, and the non-interference between all the live threads.

Non-interference. The non-interference macro $\text{NI}(\Theta, \mathbb{Q})$ requires that each dynamic thread be compatible with all the other. It is formally defined as:

$$\begin{array}{l}
\forall t_i, t_j \in \text{dom}(\Theta). \forall M, M'. \\
t_i \neq t_j \Rightarrow G_i \mathbb{R}_i M M' \Rightarrow A_j \mathbb{R}_j M M',
\end{array}$$

where $(-, A_i, G_i) = \Theta(t_i)$, $(-, A_j, G_j) = \Theta(t_j)$, $\mathbb{R}_i = \mathbb{Q}(t_i)$ and $\mathbb{R}_j = \mathbb{Q}(t_j)$.

As explained in Section 3, here we enforce the non-interference by a lazy check of specifications Θ of dynamic threads in \mathbb{Q} , instead of checking the specification Δ of all static threads in \mathbb{T} .

Well-formed dynamic threads. The rule DTHRDS ensures each dynamic thread in \mathbb{Q} is in good shape with respect to the specification Θ , the current program state \mathbb{S} , and the transition g that the current thread need do before other threads can take control.

$$\begin{array}{c}
(\mathbb{R}_k, \mathbb{I}_k) = \mathbb{Q}(t_k), (\mathbb{P}_k, \mathbb{A}_k, \mathbb{G}_k) = \Theta(t_k), \forall t_k \in \text{dom}(\mathbb{Q}) \\
(Inv, \mathbb{P}_k \circ \mathbb{A}_k, \mathbb{R}_k) \quad \forall M'.g \mathbb{R} M M' \Rightarrow \mathbb{P}_k(M', \mathbb{R}_k) \\
(Inv, \Delta, \Psi); (\mathbb{P}_k, \mathbb{G}_k, \mathbb{A}_k, \mathbb{G}_k) \vdash \mathbb{I}_k \\
\hline
(Inv, \Delta, \Psi); \Theta; (g, (M, \mathbb{R})) \vdash \mathbb{Q}
\end{array}
\quad (\text{DTHRDS})$$

The first line gives the specification of each thread when it reaches a **yield** point, and its execution context.

The first premise in line 2 requires that if it is safe for a thread to take control at certain state, it should be also safe to do so after any state transition satisfying its assumption. Note that the state transition will not affect the thread-private data in \mathbb{R}_k .

The second premise describes the relationship between the local guarantee of the current thread and the preconditions of other threads. For any transitions starting from the current state (M, \mathbb{R}) , as long as it satisfies g , it should be safe for other threads to take control at the result state.

Line 3 requires that for each thread its remainder instruction sequence must be well formed. Note we use \mathbb{G}_k as the local guarantee because after **yield**, a thread starts a new “atomic transition” described by its global guarantee.

Well-formed static threads. This rule checks the static thread entry is well formed with respect to its specification in Δ .

$$\begin{array}{c}
(\mathbb{P}_k, \mathbb{A}_k, \mathbb{G}_k) = \Delta(h_k), \mathbb{I}_k = \mathbb{T}(h_k), \forall h_k \in \text{dom}(\mathbb{T}) \\
\forall \mathbb{R}. (Inv, \mathbb{P}_k \circ \mathbb{A}_k, \mathbb{R}) \quad (Inv, \Delta, \Psi); (\mathbb{P}_k, \mathbb{G}_k, \mathbb{A}_k, \mathbb{G}_k) \vdash \mathbb{I}_k \\
\hline
(Inv, \Delta, \Psi) \vdash \mathbb{T}
\end{array}
\quad (\text{THRDS})$$

The first line gives the specification for each static thread in \mathbb{T} . We implicitly require $\text{dom}(\Delta) = \text{dom}(\mathbb{T})$.

The first premise in line 2 says that if it is safe to invoke the new thread at certain state, it should also be safe to delay the invocation after any transition satisfying the assumption of this thread.

The initial instruction sequence of each thread must be well-formed with respect to the thread specification. This is required by the last premise.

Well-formed code heap. A code heap is well formed if every instruction sequence is well-formed with respect to its corresponding specification in Ψ .

$$\begin{array}{c}
\text{dom}(\Psi) = \text{dom}(\mathbb{C}) \quad (Inv, \Delta, \Psi); \Psi(\mathbb{f}) \vdash \mathbb{C}(\mathbb{f}), \forall \mathbb{f} \in \text{dom}(\Psi) \\
\hline
(Inv, \Delta, \Psi) \vdash \mathbb{C}
\end{array}
\quad (\text{CDHP})$$

Thread creation. The **FORK** rule describes constraints on new thread creation, which enforces the non-interference between the new thread and existing threads.

$$\begin{array}{c}
(\mathbb{P}', \mathbb{A}', \mathbb{G}') = \Delta(h) \\
\mathbb{A} \Rightarrow \mathbb{A}'' \quad \mathbb{G}'' \Rightarrow \mathbb{G} \quad (\mathbb{A} \vee \mathbb{G}'') \Rightarrow \mathbb{A}' \quad \mathbb{G}' \Rightarrow (\mathbb{G} \wedge \mathbb{A}'') \\
\forall \mathbb{R}, \mathbb{R}'. (\mathbb{R}(\text{rt}) \neq \mathbb{R}'(\text{rt}) \wedge \mathbb{R}(x) = \mathbb{R}'(ra)) \Rightarrow ((Inv \wedge \mathbb{P}) \xrightarrow{\mathbb{R}} g \xrightarrow{\mathbb{R}'} \mathbb{P}') \\
(Inv, \Delta, \Psi); (\mathbb{P}, g, \mathbb{A}', \mathbb{G}'') \vdash \mathbb{I} \\
\hline
(Inv, \Delta, \Psi); (\mathbb{P}, g, \mathbb{A}, \mathbb{G}) \vdash \text{fork } h, r; \mathbb{I}
\end{array}
\quad (\text{FORK})$$

As explained in Section 3.1, the parent thread can change its specification to reflect the change of the environment. To maintain the non-interference invariant, constraints between specifications of the parent and child threads have to be satisfied, as described in Lemma 3.4. Here we enforce these constraints by premises in line 2, where $(\mathbb{A} \vee \mathbb{G}'') \Rightarrow \mathbb{A}'$ is the shorthand for:

$$\begin{array}{c}
\forall (M, \mathbb{R}). \forall M', M''. \forall \mathbb{R}'. (Inv \wedge \mathbb{P})(M, \mathbb{R}) \Rightarrow \mathbb{R}(\text{rt}) \neq \mathbb{R}'(\text{rt}) \\
\Rightarrow \mathbb{R}(x) = \mathbb{R}'(ra) \Rightarrow (\mathbb{A} \vee \mathbb{G}'') \mathbb{R} M M' M'' \Rightarrow \mathbb{A}' \mathbb{R}' M' M'',
\end{array}$$

and $\mathbb{G}' \Rightarrow (\mathbb{G} \wedge \mathbb{A}'')$ for:

$$\begin{array}{c}
\forall (M, \mathbb{R}). \forall M', M''. \forall \mathbb{R}'. (Inv \wedge \mathbb{P})(M, \mathbb{R}) \Rightarrow \mathbb{R}(\text{rt}) \neq \mathbb{R}'(\text{rt}) \\
\Rightarrow \mathbb{R}(x) = \mathbb{R}'(ra) \Rightarrow \mathbb{G}' \mathbb{R}' M' M'' \Rightarrow (\mathbb{G} \wedge \mathbb{A}'') \mathbb{R} M' M''.
\end{array}$$

Above non-interference checks use the extra knowledge that:

- the new thread id is different with its parent's, *i.e.*, $\mathbb{R}(\text{rt}) \neq \mathbb{R}'(\text{rt})$;
- the argument of the new thread comes from the parent's register x , *i.e.*, $\mathbb{R}(x) = \mathbb{R}'(ra)$;
- the parent's register file satisfies the precondition, *i.e.*, $(Inv \wedge \mathbb{P})(M, \mathbb{R})$.

In most cases, the programmer can just pick $(\mathbb{A} \vee \widehat{\mathbb{G}'})$ and $(\mathbb{G} \wedge \widehat{\mathbb{A}'})$ as \mathbb{A}'' and \mathbb{G}'' respectively, where $\widehat{\mathbb{G}'}$ and $\widehat{\mathbb{A}'}$ are instantiations of \mathbb{G}' and \mathbb{A}' using the value of the child's argument.

The premise in line 3 says that after the current thread completes the transition described by g , it should be safe for the new thread to take control with its new register file (\mathbb{R}') , whose relationship between the parents register file \mathbb{R} is satisfied.

The last premise checks the well-formedness of the remainder instruction sequence. Since the **fork** instruction does not change states, we need not change the precondition \mathbb{P} and g .

Yielding and termination.

$$\begin{array}{c}
\forall \mathbb{R}. (Inv, \mathbb{P} \circ \mathbb{A}, \mathbb{R}) \quad (Inv \wedge \mathbb{P}) \Rightarrow g \quad (Inv, \Delta, \Psi); (\mathbb{P}, \mathbb{G}, \mathbb{A}, \mathbb{G}) \vdash \mathbb{I} \\
\hline
(Inv, \Delta, \Psi); (\mathbb{P}, g, \mathbb{A}, \mathbb{G}) \vdash \text{yield}; \mathbb{I}
\end{array}
\quad (\text{YIELD})$$

The **YIELD** rule requires that it is safe for the yielding thread to take back control after any state transition satisfying the assumption \mathbb{A} . Also the current thread cannot yield until it completes the required state transition, *i.e.*, an identity transition satisfies the local guarantee g . Lastly, one must verify the remainder instruction sequence with the local guarantee reset to \mathbb{G} .

$$\begin{array}{c}
(Inv \wedge \mathbb{P}) \Rightarrow g \\
\hline
(Inv, \Delta, \Psi); (\mathbb{P}, g, \mathbb{A}, \mathbb{G}) \vdash \text{exit}
\end{array}
\quad (\text{EXIT})$$

The rule **EXIT** is simple: it is safe for the current thread to terminate its execution only after it finishes the required transition described by g , which is an identity transition.

Type-checking other instructions.

$$\begin{array}{c}
c \in \{\text{add } r_d, r_s, r_t, \text{sub } r_d, r_s, r_t, \text{movi } r_d, w\}, r_d \notin \{\text{rt}, ra\} \\
(Inv \wedge \mathbb{P}) \Rightarrow (Inv \wedge \mathbb{P}') \circ c \quad (Inv \wedge \mathbb{P}) \Rightarrow (g' \circ c) \Rightarrow g \\
\mathbb{A} \Rightarrow \mathbb{A}' \circ c \quad \mathbb{G}' \circ c \Rightarrow \mathbb{G} \quad (Inv, \Delta, \Psi); (\mathbb{P}', g', \mathbb{A}', \mathbb{G}') \vdash \mathbb{I} \\
\hline
(Inv, \Delta, \Psi); (\mathbb{P}, g, \mathbb{A}, \mathbb{G}) \vdash c; \mathbb{I}
\end{array}
\quad (\text{SIMP})$$

The rule **SIMP** covers the verification of instruction sequences starting with a simple command such as **add**, **sub** or **movi**. We require that the program not update registers rt and ra . In these cases, one must find an intermediate precondition $(\mathbb{P}', g', \mathbb{A}', \mathbb{G}')$ under which the remainder instruction sequence \mathbb{I} is well-formed.

The global invariant Inv and the intermediate assertion \mathbb{P}' must hold on the updated machine state, and the intermediate guarantee g' applied to the updated machine state must be no weaker than the current guarantee g applied to the current state.

Since \mathbb{A} and \mathbb{G} are parameterized over \mathbb{R} , which will be changed by the instruction c , one may change \mathbb{A} and \mathbb{G} to ensure the assumption does not become stronger and the guarantee does not become weaker.

$$\begin{array}{c}
(\mathbb{P}', g', \mathbb{A}', \mathbb{G}') = \Psi(\mathbb{f}) \\
(Inv \wedge \mathbb{P}) \Rightarrow \mathbb{P}' \quad (Inv \wedge \mathbb{P}) \Rightarrow g' \Rightarrow g \quad \mathbb{A} \Rightarrow \mathbb{A}' \quad \mathbb{G}' \Rightarrow \mathbb{G} \\
\hline
(Inv, \Delta, \Psi); (\mathbb{P}, g, \mathbb{A}, \mathbb{G}) \vdash \text{jd } \mathbb{f}
\end{array}
\quad (\text{JD})$$

The **JD** rule checks the specification of the target instruction sequence. As mentioned before, each instruction sequence can have

its own specification, independent of the thread that will jump to it. It is safe for a thread to execute an instruction sequence as long as executing the instruction sequence does not require a stronger assumption than the thread's assumption \mathbb{A} , nor does it break the guarantee \mathbb{G} of the thread.

Inference rules for memory operations are quite similar to the SIMP rule. Here we also need ensure the memory address is in the domain of the data heap.

$$\frac{\begin{array}{l} c = \text{ld } r_d, r_s(w) \quad r_d \notin \{\text{rt}, \text{ra}\} \\ \forall M, \forall R. (Inv \wedge p) (M, R) \Rightarrow ((R(r_s) + w) \in \text{dom}(M)) \\ (Inv \wedge p) \Rightarrow (Inv \wedge p') \circ c \quad (Inv \wedge p) \Rightarrow (g' \circ c) \Rightarrow g \\ \forall S, M, M'. \mathbb{A} (S, R) M M' \Rightarrow \mathbb{A}' (\text{Next}(c, S), R) M M' \\ \forall S, M, M'. G' (\text{Next}(c, S), R) M M' \Rightarrow G (S, R) M M' \\ (Inv, \Delta, \Psi); (p', g', \mathbb{A}', G') \vdash I \end{array}}{(Inv, \Delta, \Psi); (p, g, \mathbb{A}, G) \vdash c; I} \quad (\text{LD})$$

$$\frac{\begin{array}{l} c = \text{st } r_d(w), r_s \\ \forall M, \forall R. (Inv \wedge p) (M, R) \Rightarrow ((R(r_d) + w) \in \text{dom}(M)) \\ (Inv \wedge p) \Rightarrow (Inv \wedge p') \circ c \quad (Inv \wedge p) \Rightarrow (g' \circ c) \Rightarrow g \\ (Inv, \Delta, \Psi); (p', g', \mathbb{A}, G) \vdash I \end{array}}{(Inv, \Delta, \Psi); (p, g, \mathbb{A}, G) \vdash c; I} \quad (\text{ST})$$

Rules for conditional branching instructions are similar to the JD rule, which are straightforward to understand.

$$\frac{\begin{array}{l} (p', g', \mathbb{A}', G') = \Psi(f) \quad \mathbb{A} \Rightarrow \mathbb{A}' \quad G' \Rightarrow G \\ \forall S. (Inv \wedge p) S \Rightarrow (S(r_s) > S(r_t)) \Rightarrow (p' S) \\ \forall S. \forall M'. (Inv \wedge p) S \Rightarrow (S(r_s) > S(r_t)) \Rightarrow (g' S M') \Rightarrow (g S M') \\ \forall S. (Inv \wedge p) S \Rightarrow (S(r_s) \leq S(r_t)) \Rightarrow (p'' S) \\ \forall S. \forall M'. (Inv \wedge p) S \Rightarrow (S(r_s) \leq S(r_t)) \Rightarrow (g'' S M') \Rightarrow (g S M') \\ (Inv, \Delta, \Psi); (p'', g'', \mathbb{A}, G) \vdash I \end{array}}{(Inv, \Delta, \Psi); (p, g, \mathbb{A}, G) \vdash \text{bgt } r_s, r_t, f; I} \quad (\text{BGT})$$

$$\frac{\begin{array}{l} (p', g', \mathbb{A}', G') = \Psi(f) \quad \mathbb{A} \Rightarrow \mathbb{A}' \quad G' \Rightarrow G \\ \forall S. (Inv \wedge p) S \Rightarrow (S(r_s) = S(r_t)) \Rightarrow (p' S) \\ \forall S. \forall M'. (Inv \wedge p) S \Rightarrow (S(r_s) = S(r_t)) \Rightarrow (g' S M') \Rightarrow (g S M') \\ \forall S. (Inv \wedge p) S \Rightarrow (S(r_s) \neq S(r_t)) \Rightarrow (p'' S) \\ \forall S. \forall M'. (Inv \wedge p) S \Rightarrow (S(r_s) \neq S(r_t)) \Rightarrow (g'' S M') \Rightarrow (g S M') \\ (Inv, \Delta, \Psi); (p'', g'', \mathbb{A}, G) \vdash I \end{array}}{(Inv, \Delta, \Psi); (p, g, \mathbb{A}, G) \vdash \text{beq } r_s, r_t, f; I} \quad (\text{BEQ})$$

4.5 Soundness of CMAP

The soundness of CMAP inference rules with respect to the operational semantics of the machine is established following the syntactic approach of proving type soundness [39]. From the “progress” and “preservation” lemmas, we can guarantee that given a well-formed program under compatible assumptions and guarantees, the current instruction sequence will be able to execute without getting “stuck”. Furthermore, any safety property derivable from the global invariant will hold throughout the execution. We define $\mathbb{P} \xrightarrow{n} \mathbb{P}'$ as the relation of n -step ($n \geq 0$) program transitions. The soundness of CMAP is formally stated as Theorem 4.3.

Lemma 4.1 (Progress)

$\Phi = (Inv, \Delta, \Psi)$. If there exist Θ, p, g, \mathbb{A} and \mathbb{G} such that $\Phi; \Theta; (p, g, \mathbb{A}, G) \vdash ((M, R), Q, T, C, I)$, then $(Inv M)$, and there exists a program $\tilde{\mathbb{P}}$ such that $((M, R), Q, T, C, I) \mapsto \tilde{\mathbb{P}}$.

Lemma 4.2 (Preservation)

If $\Phi; \Theta; (p, g, \mathbb{A}, G) \vdash \mathbb{P}$ and $\mathbb{P} \mapsto \tilde{\mathbb{P}}$, where $\mathbb{P} = (S, Q, T, C, I)$ and $\tilde{\mathbb{P}} = (\tilde{S}, \tilde{Q}, \tilde{T}, \tilde{C}, \tilde{I})$, then there exist $\tilde{\Theta}, \tilde{p}, \tilde{g}, \tilde{\mathbb{A}}$ and $\tilde{\mathbb{G}}$ such that $\Phi; \tilde{\Theta}; (\tilde{p}, \tilde{g}, \tilde{\mathbb{A}}, \tilde{G}) \vdash \tilde{\mathbb{P}}$.

Theorem 4.3 (Soundness)

$\Phi = (Inv, \Delta, \Psi)$. If there exist Θ, p, g, \mathbb{A} and \mathbb{G} such that $\Phi; \Theta; (p, g, \mathbb{A}, G) \vdash \mathbb{P}_0$, then for any $n \geq 0$, there exist M, R, Q, T, C and I such that $\mathbb{P}_0 \xrightarrow{n} ((M, R), Q, T, C, I)$ and $(Inv M)$.

The proofs for these two lemmas and the soundness theorem are given in Appendix A. We have also implemented the complete CMAP system [10] in the Coq proof assistant so we are confident that CMAP is indeed sound and can be used to certify general multi-threaded programs.

5. Examples

5.1 Unbounded Dynamic Thread Creation

In Figure 3 we showed a small program `main2` which spawns child threads within a `while` loop. This kind of unbounded dynamic thread creation cannot be supported using the `cobegin/coend` structure. We show how such a program is specified and verified using our logic. To simplify the specification, we trivialize the function `f` and `g` and let $f(i) = i$ and $g(x, _) = x + 1$.

We assume the high-level program works in a preemptive mode. Figure 10 shows the CMAP implementation, where `yield` instructions are inserted to simulate the preemption. This also illustrates that our logic is general enough to simulate preemptive thread model.

To verify the safety property of the program, the programmer need find a global invariant and specifications for each static thread and the code heap. In Figure 10 we show definitions of assertions that are used to specify the program. Proof sketches are also inserted in the program. For ease of reading, we use named variables as short-hands for their values in memory. The primed variables represent the value of the variable after state transition. We also introduce the shorthand $[r]$ for $R(r)$.

The following formulae show the specifications of static threads, and the initial memory and instruction sequence.

$$\begin{array}{l} Inv \equiv \text{True} \\ \Delta \equiv \{\text{main} \rightsquigarrow (\text{True}, \mathbb{A}_0, \mathbb{G}_0), \text{chld} \rightsquigarrow (p, \mathbb{A}, G)\} \\ \Psi \equiv \{\text{loop} \rightsquigarrow (p', G_1, \mathbb{A}_1, G_1), \text{cont} \rightsquigarrow (p_3, G_3, \mathbb{A}_3, G_3)\} \\ \text{Initial } M \equiv \{\text{data} \rightsquigarrow _, \dots, \text{data} + 99 \rightsquigarrow _ \} \\ \text{Initial } I \equiv \text{movi } r_0, 0; \text{movi } r_1, 1; \text{movi } r_2, 100; \text{jd loop} \end{array}$$

For each child thread, it is natural to assume that no other threads will touch its share of the data entry and guarantee that other threads' data entry will not be changed, as specified by \mathbb{A} and \mathbb{G} . For the main thread, it assumes at the beginning that no other threads in the environment will change any of the data entries (\mathbb{A}_0).

Specifying the loop body is not easy. We need find a loop invariant $(p', G_1, \mathbb{A}_1, G_1)$ to attach to the code label `loop`. At first glance, \mathbb{A}_1 and \mathbb{G}_1 can be defined the same as \mathbb{A}_3 and \mathbb{G}_3 , respectively. However, this does not work because our FORK rule requires $G \Rightarrow G_1$, which cannot be satisfied. Instead, our \mathbb{A}_1 and \mathbb{G}_1 are polymorphic over the loop index r_0 , which reflects the composition of the changing thread environments. At the point that a new child thread is forked but the value of r_0 has not been changed to reflect the environment change, we explicitly change the assumption and guarantee to \mathbb{A}_2 and \mathbb{G}_2 . When the value of r_0 is increased, we cast \mathbb{A}_2 and \mathbb{G}_2 back to \mathbb{A}_1 and \mathbb{G}_1 .

5.2 The Readers-Writers Problem

Our logic is general enough to specify and verify general properties of concurrent programs. In this section, we give a simple solution of the readers-writers problem and show that there is no race conditions. This example also shows how P/V operations and lock primitives can be implemented and specified in CMAP.

Figure 11 shows the C-like pseudo code for readers and writers. Note that this simple code just ensures that there is no race con-

```

p ≡ 0 ≤ [ra] < 100    A ≡ data[ra] = data'[ra]
G ≡ ∀i.(0 ≤ i < 100 ∧ i ≠ [ra]) ⇒ (data[i] = data'[i])
A₀ ≡ ∀i.0 ≤ i < 100 ⇒ (data[i] = data'[i])    G₀ ≡ True
A₁ ≡ ∀i.(0 ≤ i < 100 ∧ i ≥ [r₀]) ⇒ (data[i] = data'[i])
G₁ ≡ ∀i.(0 ≤ i < 100 ∧ i < [r₀]) ⇒ (data[i] = data'[i])
A₂ ≡ ∀i.(0 ≤ i < 100 ∧ i > [r₀]) ⇒ (data[i] = data'[i])
G₂ ≡ ∀i.(0 ≤ i < 100 ∧ i ≤ [r₀]) ⇒ (data[i] = data'[i])
A₃ ≡ True    G₃ ≡ ∀i.0 ≤ i < 100 ⇒ (data[i] = data'[i])
p' ≡ 0 ≤ [r₀] < 100 ∧ [r₁] = 1 ∧ [r₂] = 100    p₃ ≡ [r₀] = 100

```

```

main : -{(True, A₀, G₀)}
  movi r₀, 0
  movi r₁, 1
  movi r₂, 100
  jd loop
loop : -{(p', G₁, A₁, G₁)}
  beq r₀, r₂, cont
  yield
  st r₀(data), r₀
  yield
  fork chld, r₀
  -{(p', G₁, A₂, G₂)}
  yield
  add r₀, r₀, r₁
  -{(p', G₁, A₁, G₁)}
  yield
  jd loop
cont : -{(p₃, G₃, A₃, G₃)}
  ...

```

Figure 10. Loop: the CMAP program

```

Variables :
  int[100] rf, wf;
  int cnt, writ, l, v;
Initially :
  rf[i] = wf[i] = 0, 0 ≤ i < 100;
  cnt = 0 ∧ writ = 1 ∧ l = 0;

writer(int x){
  while(true){
    P(writ);
    wf[x] := 1;
    write v ...
    wf[x] := 0;
    V(writ);
  }
}

reader(int x){
  while(true){
    lock_acq(1);
    if (cnt = 0){
      P(writ);
    }
    cnt := cnt + 1;
    rf[x] := 1;
    lock_rel(1);
    read v ...
    lock_acq(1);
    rf[x] := 0;
    cnt := cnt - 1;
    if(cnt = 0){
      V(writ);
    }
    lock_rel(1);
  }
}

```

Figure 11. Readers & writers : the high-level program

ditions. It does not ensure fairness. Verification of liveness properties is part of our future work. We assume that 100 readers and writers will be created by a main thread. The main thread and its specification will be very similar to the main program shown in the previous section, so we omit it here and just focus on the readers and writers code. The array of `rf` and `wf` are not necessary for the implementation. They are introduced as auxiliary variables just for specification and verification purpose.

Figure 13 shows the CMAP implementation of the high-level pseudo code. Yielding is inserted at all the intervals of the atomic operations of the high-level program. The lock primitives and P/V operations, as shown in Figure 12, are implemented as program macros parameterized by the return label. They will be instantiated

```

acq(f) : -{}
  yield
  movi r₀, 0
  movi r₁, 1
  ld r₂, r₁(0)
  bgt r₂, r₀, acq
  st r₁(0), rt
  jd f

rel(f) : -{}
  yield
  movi r₀, 0
  st r₀(1), r₀
  jd f

p_writ(f) : -{}
  yield
  movi r₀, 0
  movi r₁, writ
  ld r₂, r₁(0)
  beq r₂, r₀, p_writ
  st r₁(0), r₀
  jd f

v_writ(f) : -{}
  movi r₁, 1
  movi r₂, writ
  st r₂(0), r₁
  jd f

```

Figure 12. Lock & semaphore primitives

```

reader : -{(p₁, g, Aᵣ, Gᵣ)}
  JD acq(cont_1)

cont_1 : -{(p₂, g, Aᵣ, Gᵣ)}
  yield
  movi r₀, 0
  movi r₁, cnt
  ld r₂, r₁(0)
  yield
  beq r₂, r₀, getw
  jd inc_cnt

getw : -{(p₃, g, Aᵣ, Gᵣ)}
  JD p_writ(inc_cnt)

inc_cnt : -{(p₄, g, Aᵣ, Gᵣ)}
  yield
  movi r₁, cnt
  ld r₂, r₁(0)
  yield
  movi r₃, 1
  add r₂, r₂, r₃
  st r₁(0), r₂
  yield
  movi r₁, 1
  st ra(rf), r₁
  -{(p₅, g, Aᵣ, Gᵣ)}
  JD rel(cont_2)

cont_2 : -{(p₆, g, Aᵣ, Gᵣ)}
  yield
  ...
  yield
  ...
  yield
  -{(p₆, g, Aᵣ, Gᵣ)}
  JD acq(cont_3)

cont_3 : -{(p₇, g, Aᵣ, Gᵣ)}
  yield
  movi r₀, 0
  st ra(rf), r₀
  yield
  movi r₃, 1
  sub r₂, r₂, r₃
  st r₁(0), r₂
  yield
  beq r₂, r₀, relw
  -{(p₈, g, Aᵣ, Gᵣ)}
  JD rel(reader)

relw : -{(p₉, g, Aᵣ, Gᵣ)}
  yield
  JD v_writ(cont_4)

cont_4 : -{(p₈, g, Aᵣ, Gᵣ)}
  JD rel(reader)

writer : -{(p₁₀, g, Aᵂ, Gᵂ)}
  JD p_writ(cont_5)

cont_5 : -{(p₁₁, g₁, Aᵂ, Gᵂ)}
  movi r₁, 1
  st ra(wf), r₁
  yield
  ...
  yield
  movi r₀, 0
  st ra(wf), r₀
  -{(p₁₁, g₂, Aᵂ, Gᵂ)}
  JD v_writ(writer)

```

Figure 13. Readers & writers : the CMAP program

and inlined in the proper position of the code. We introduce a pseudo instruction `JD` to represent the inlining of the macro.

We define the global invariant and reader/writer's assumptions and guarantees in Figure 14. The code heap specifications are embedded in the code as annotations, as shown in Figure 13. Specifications of lock primitives and P/V operations are given at places they are inlined. Definitions of assertions and local guarantees used in code heap specifications are shown in Figure 14.

The assertion inv_1 says that out of the critical section protected by the lock, the value of the counter `cnt` is always consistent and reflects the number of the readers that can read the value of `v`; inv_2 says at one time there is at most one writer that can change the value of `v`; while inv_3 and inv_4 states the relationship between the counter `cnt`, the semaphore variable `writ` and the actual number

$$\begin{aligned}
inv_1 &\equiv 1 = 0 \Rightarrow \sum_i \mathbf{rf}[i] = \mathbf{cnt} & inv_2 &\equiv \sum_i \mathbf{wf}[i] \leq 1 \\
inv_3 &\equiv \sum_i \mathbf{wf}[i] = 1 \Rightarrow (\mathbf{cnt} = 0 \wedge \mathbf{writ} = 0) \\
inv_4 &\equiv \mathbf{writ} = 1 \Rightarrow \mathbf{cnt} = 0 & Inv &\equiv inv_1 \wedge inv_2 \wedge inv_3 \wedge inv_4 \\
idr &\equiv \forall i. \mathbf{rf}[i] = \mathbf{rf}'[i] \\
idr_1(i) &\equiv \forall j. i \neq j \Rightarrow \mathbf{rf}[j] = \mathbf{rf}'[j] & idr_2(i) &\equiv \mathbf{rf}[i] = \mathbf{rf}'[i] \\
idw &\equiv \forall i. \mathbf{wf}[i] = \mathbf{wf}'[i] \\
idw_1(i) &\equiv \forall j. i \neq j \Rightarrow \mathbf{wf}[j] = \mathbf{wf}'[j] & idw_2(i) &\equiv \mathbf{wf}[i] = \mathbf{wf}'[i] \\
A_r &\equiv idr_2(\mathbf{ra}) \wedge (\mathbf{rf}[\mathbf{ra}] = 1 \Rightarrow v = v') \wedge (1 = [\mathbf{rt}] \Rightarrow \mathbf{cnt} = \mathbf{cnt}') \\
G_r &\equiv idw \wedge idr_1(\mathbf{ra}) \wedge v = v' \wedge (1 \notin \{[\mathbf{rt}], 0\} \Rightarrow \mathbf{cnt} = \mathbf{cnt}') \\
A_w &\equiv idw_2(\mathbf{ra}) \wedge (\mathbf{wf}[\mathbf{ra}] = 1 \Rightarrow v = v') \\
G_w &\equiv idr \wedge idw_1(\mathbf{ra}) \wedge ((\mathbf{writ} = 0 \wedge \mathbf{wf}[\mathbf{ra}] = 0) \Rightarrow v = v') \\
g &\equiv \lambda R, M, M'. M = M' & g_1 &\equiv \lambda R, M, M'. M\{\mathbf{wf}[\mathbf{ra}] \sim 1\} = M' \\
g_2 &\equiv \lambda R, M, M'. M\{\mathbf{writ} \sim 1\} = M' & p_1 &\equiv \mathbf{rf}[\mathbf{ra}] = 0 \wedge 1 \neq [\mathbf{rt}] \\
p_2 &\equiv \mathbf{rf}[\mathbf{ra}] = 0 \wedge 1 = [\mathbf{rt}] & p_3 &\equiv p_2 \wedge \mathbf{cnt} = 0 \\
p_4 &\equiv p_2 \wedge \mathbf{writ} = 0 & p_5 &\equiv \mathbf{rf}[\mathbf{ra}] = 1 \wedge 1 = [\mathbf{rt}] \\
p_6 &\equiv \mathbf{rf}[\mathbf{ra}] = 1 \wedge 1 \neq [\mathbf{rt}] \wedge \mathbf{cnt} > 0 \\
p_7 &\equiv \mathbf{rf}[\mathbf{ra}] = 1 \wedge 1 = [\mathbf{rt}] \wedge \mathbf{cnt} > 0 & p_8 &\equiv 1 = [\mathbf{rt}] \wedge \mathbf{rf}[\mathbf{ra}] = 0 \\
p_9 &\equiv 1 = [\mathbf{rt}] \wedge \mathbf{writ} = 0 \wedge \mathbf{cnt} = 0 \wedge \mathbf{rf}[\mathbf{ra}] = 0 \\
p_{10} &\equiv \mathbf{wf}[\mathbf{ra}] = 0 & p_{11} &\equiv \mathbf{wf}[\mathbf{ra}] = 0 \wedge \mathbf{writ} = 0 \wedge \mathbf{cnt} = 0
\end{aligned}$$

Figure 14. Program specifications

of writers that can write the data, which must be maintained to ensure the mutual exclusive access between readers and writers. The program invariant Inv , which is the conjunction of the four, must be satisfied at any step of the execution.

The assumptions and guarantees of readers and writers are parameterized by their thread id (in \mathbf{rt}) and thread argument (in \mathbf{ra}). The reader assumes (A_r) that if it has the right to read the data v (i.e., $\mathbf{rf}[\mathbf{ra}] = 1$), nobody can change the data. Also, if it owns the lock, nobody else can change the value of the counter. Finally, it assumes its auxiliary variable $\mathbf{rf}[\mathbf{ra}]$ will never be changed by others. Readers guarantee (G_r) that they will never change the shared data v and auxiliary variables of other threads, and it will not revise the counter unless it owns the lock. Writers' assumption (A_w) and guarantee (G_w) are defined in a similar way to ensure the non-interference.

5.3 More Examples

Yu and Shao [42] have shown by examples that CCAP is expressive enough to verify general properties of concurrent programs, like mutual exclusion, deadlock freedom, and partial correctness. CMAP, as an extension of CCAP with dynamic thread creation, is more expressive than CCAP (it is straightforward to translate the CCAP code and verification to PCC packages in CMAP). Therefore, CMAP also applies for those CCAP examples. In Appendix B, we give a variation of the GCD program shown in [42], where collaborating threads are forked by a main thread. We also show how thread join can be implemented and reasoned in CMAP.

6. Related Work and Conclusion

We have presented a certified programming framework for verifying multi-threaded assembly code with unbounded dynamic thread creation. Our work is related to two directions of research: concurrency verification and PCC. The rely-guarantee method [23, 35, 36, 1, 13, 42] is one of the best studied technique for compositional concurrent program verification. However, most of the work on it are based on high-level languages or calculi, and none of them support unbounded dynamic thread creation. On the other hand, many PCC frameworks [28, 27, 2, 18, 41, 8, 42] have been proposed for machine/assembly code verification, but most of them only support sequential code. The only intersection point of these two directions is the work on CCAP [42], which applies the R-G method at the assembly level to verify general concurrency properties, like mu-

tual exclusion and deadlock-freedom. Unfortunately, CCAP does not support dynamic thread creation either.

Recently, O'Hearn and others [31, 4, 3] applied separation logic to reason about concurrent programs. Bornat et al. [3] showed how to verify the race-free property of a program solving the readers and writers problem, which is similar to the program presented in this paper. However, their work heavily depends on the higher-level language features, such as resources and conditional critical regions. As other traditional work on concurrency verification, they only support nested cobegin/coend structure. It is not clear how their technique can be applied to support assembly code verification with unbounded dynamic thread creation.

A number of model checkers have been developed for concurrent software verification, but most of them only check programs with a fixed finite number of threads [22, 9, 19, 5, 16, 14]. The CIRC algorithm [20] supports unbounded number of threads, but it doesn't model the dynamic thread creation and the changing thread environment. Qadeer and Rehof [33] pointed out that verification of concurrent boolean program with unbounded parallelism is decidable if the number of context switches is bounded, but they do not directly verify dynamic thread creation either. Given a context bound k , they reduce a dynamic concurrent program P to a program Q with $k+1$ threads and verify Q instead. 3VMC [40] supports both unbounded number of threads and dynamic thread creation, but it is not based on the rely-guarantee method and does not support compositional verification.

Many type systems are also proposed to reason about concurrent programs [11, 12, 15, 17]. Unlike CMAP, which uses high-order predicate logic to verify general program properties, they are designed to automatically reason about specific properties of programs, like races, deadlocks and atomicity. Also, they do not directly generate proofs about program properties. Instead, proofs are implicitly embedded in their soundness proofs.

CMAP extends previous work on R-G method and CCAP with dynamic thread creation. We unify the concepts of a thread's assumption/guarantee and its environment's guarantee/assumption, and allow a thread to change its assumption and guarantee to track the changing environment caused by dynamic thread creation. Code segments in CMAP can be specified and verified once and used in multiple threads with different assumptions and guarantees, therefore CMAP achieves better modularity than CCAP. Some practical issues, such as argument passing at thread creation, thread local data, and multiple invocation of one copy of thread code, are also discussed to support practical multi-threaded programming. CMAP has been developed using the Coq proof assistant, along with a formal soundness proof and the verified example programs.

The goal of our work is to provide an explicit and general framework (with smaller TCB) such that code and specifications at the higher level can be compiled down to the assembly level. Although directly specifying and proving CMAP programs may be daunting, it is simpler at the higher level. Also, there are common concurrent idioms that would permit the assumption/guarantee to be generated automatically during compilation. For example, for critical sections protected by lock, the assumption is always like "nobody else will update the protected data if I am holding the lock..." (see Figure 14). Another scenario is illustrated by the example shown in Figure 3 and 10, where each thread is exclusively responsible for a single piece of global data. In these scenarios, the assumption is always like "nobody else will touch my share of data" and the guarantee is like "I will not update other threads' data". Automatically generating assumptions/guarantees and proofs for common idioms, and compiling higher level concurrent programs and specifications down to CMAP, will be our future work.

The thread primitives in CMAP are higher-level pseudo instructions. They can be replaced by system calls to certified thread li-

barries, which is part of our ongoing work. Also we separate issues in multi-threaded programming from the embedded code pointer problem which is addressed in a companion paper [30]. Applying that framework to our thread model will be part of the future work.

7. Acknowledgment

We thank Andrew McCreight and anonymous referees for suggestions and comments on an earlier version of this paper. Rodrigo Ferreira helped implement the CMAP language and its soundness proof in the Coq proof system. This research is based on work supported in part by grants from Intel and Microsoft, and NSF grant CCR-0208618. Any opinions, findings, and conclusions contained in this document are those of the authors and do not reflect the views of these agencies.

References

- [1] M. Abadi and L. Lamport. Conjoining specifications. *ACM Trans. on Programming Languages and Systems*, 17(3):507–535, 1995.
- [2] A. W. Appel. Foundational proof-carrying code. In *Proc. 16th Annual IEEE Symposium on Logic in Computer Science*, pages 247–258. IEEE Computer Society, June 2001.
- [3] R. Bornat, C. Calcagno, P. W. O’Hearn, and M. J. Parkinson. Permission accounting in separation logic. In *Proc. 32th ACM Symp. on Principles of Prog. Lang.*, pages 259–270, 2005.
- [4] S. D. Brookes. A semantics for concurrent separation logic. In *Proc. 15th International Conference on Concurrency Theory (CONCUR’04)*, volume 3170 of *LNCS*, pages 16–34, 2004.
- [5] S. Chaki, E. Clarke, A. Groce, S. Jha, and H. Veith. Modular verification of software components in C. In *ICSE’03: International Conference on Software Engineering*, pages 385–395, 2003.
- [6] C. Colby, P. Lee, G. Necula, F. Blau, M. Plesko, and K. Cline. A certifying compiler for Java. In *Proc. 2000 ACM Conf. on Prog. Lang. Design and Impl.*, pages 95–107, New York, 2000. ACM Press.
- [7] T. Coquand and G. Huet. The calculus of constructions. *Information and Computation*, 76:95–120, 1988.
- [8] K. Crary. Toward a foundational typed assembly language. In *Proc. 30th ACM Symp. on Principles of Prog. Lang.*, pages 198–212, 2003.
- [9] C. Demartini, R. Iosif, and R. Sisto. dSPIN: A dynamic extension of SPIN. In *Proc. 5th and 6th International SPIN Workshops on Theoretical and Practical Aspects of SPIN Model Checking*, pages 261–276, London, UK, 1999. Springer-Verlag.
- [10] R. Ferreira and X. Feng. Coq (v8.0) implementation for CMAP language and the soundness proof. <http://flint.cs.yale.edu/publications/cmap.html>, Mar. 2005.
- [11] C. Flanagan and M. Abadi. Object types against races. In *CONCUR’99 – Concurrency Theory, volume 1664 of LNCS*, pages 288–303. Springer-Verlag, 1999.
- [12] C. Flanagan and M. Abadi. Types for safe locking. In *Proceedings of the 8th European Symposium on Programming Languages and Systems*, pages 91–108. Springer-Verlag, 1999.
- [13] C. Flanagan, S. N. Freund, and S. Qadeer. Thread-modular verification for shared-memory programs. In *Proc. 2002 European Symposium on Programming*, pages 262–277. Springer-Verlag, 2002.
- [14] C. Flanagan and S. Qadeer. Thread-modular model checking. In *Proc. 10th SPIN workshop*, pages 213–224, May 2003.
- [15] C. Flanagan and S. Qadeer. A type and effect system for atomicity. In *Proceedings of the ACM SIGPLAN 2003 conference on Programming Language Design and Implementation*, pages 338–349, 2003.
- [16] D. Giannakopoulou, C. S. Pasareanu, and H. Barringer. Assumption generation for software component verification. In *ASE’02: Automated Software Engineering*, pages 3–12, 2002.
- [17] D. Grossman. Type-safe multithreading in cyclone. In *Proceedings of the 2003 ACM SIGPLAN International workshop on Types in Languages Design and Implementation*, pages 13–25, 2003.
- [18] N. A. Hamid, Z. Shao, V. Trifonov, S. Monnier, and Z. Ni. A syntactic approach to foundational proof-carrying code. In *Proc. Seventeenth Annual IEEE Symposium on Logic in Computer Science (LICS’02)*, pages 89–100. IEEE Computer Society, July 2002.
- [19] K. Havelund and T. Pressburger. Model checking Java programs using Java pathfinder. *Software Tools for Technology Transfer (STTT)*, 2(4):72–84, 2000.
- [20] T. A. Henzinger, R. Jhala, and R. Majumdar. Race checking by context inference. In *PLDI’04: Programming Language Design and Implementation*, pages 1–13, 2004.
- [21] C. A. R. Hoare. Communicating sequential processes. *Communications of the ACM*, 21(8):666–677, 1978.
- [22] G. J. Holzmann. Proving properties of concurrent systems with SPIN. In *Proc. 6th International Conference on Concurrency Theory (CONCUR’95)*, volume 962 of *LNCS*, pages 453–455. Springer, 1995.
- [23] C. B. Jones. Tentative steps toward a development method for interfering programs. *ACM Trans. on Programming Languages and Systems*, 5(4):596–619, 1983.
- [24] L. Lamport. The temporal logic of actions. *ACM Trans. on Programming Languages and Systems*, 16(3):872–923, May 1994.
- [25] R. Milner. *A Calculus of Communicating Systems*. Springer-Verlag New York, Inc., 1982.
- [26] G. Morrisett, K. Crary, N. Glew, D. Grossman, R. Samuels, F. Smith, D. Walker, S. Weirich, and S. Zdancewic. TALx86: a realistic typed assembly language. In *1999 ACM SIGPLAN Workshop on Compiler Support for System Software*, pages 25–35, Atlanta, GA, May 1999.
- [27] G. Morrisett, D. Walker, K. Crary, and N. Glew. From System F to typed assembly language. In *Proc. 25th ACM Symp. on Principles of Prog. Lang.*, pages 85–97. ACM Press, Jan. 1998.
- [28] G. Necula. Proof-carrying code. In *Proc. 24th ACM Symp. on Principles of Prog. Lang.*, pages 106–119. ACM Press, Jan. 1997.
- [29] G. Necula and P. Lee. The design and implementation of a certifying compiler. In *Proc. 1998 ACM Conf. on Prog. Lang. Design and Impl.*, pages 333–344, New York, 1998.
- [30] Z. Ni and Z. Shao. Certified assembly programming with embedded code pointers. Technical report, Dec. 2004. <http://flint.cs.yale.edu/publications/ecap.html>.
- [31] P. W. O’Hearn. Resources, concurrency and local reasoning. In *Proc. 15th International Conference on Concurrency Theory (CONCUR’04)*, volume 3170 of *LNCS*, pages 49–67, 2004.
- [32] C. Paulin-Mohring. Inductive definitions in the system Coq—rules and properties. In *Proc. TLCA*, volume 664 of *LNCS*, 1993.
- [33] S. Qadeer and J. Rehof. Context-bounded model checking of concurrent software. In *TACAS’05: Tools and Algorithms for the Construction and Analysis of Systems*, pages 93–107, 2005.
- [34] J. H. Reppy. CML: A higher concurrent language. In *Proc. 1991 Conference on Programming Language Design and Implementation*, pages 293–305, Toronto, Ontario, Canada, 1991. ACM Press.
- [35] E. W. Stark. A proof technique for rely/guarantee properties. In *Proc. 5th Conf. on Foundations of Software Technology and Theoretical Computer Science*, volume 206 of *LNCS*, pages 369–391, 1985.
- [36] K. Stølen. A method for the development of totally correct shared-state parallel programs. In *Proc. 2nd International Conference on Concurrency Theory (CONCUR’91)*, pages 510–525, 1991.
- [37] The Coq Development Team. The Coq proof assistant reference manual. The Coq release v7.1, Oct. 2001.
- [38] B. Werner. *Une Théorie des Constructions Inductives*. PhD thesis, A L’Université Paris 7, Paris, France, 1994.
- [39] A. K. Wright and M. Felleisen. A syntactic approach to type soundness. *Information and Computation*, 115(1):38–94, 1994.
- [40] E. Yahav. Verifying safety properties of concurrent Java programs using 3-valued logic. In *Proc. 28th ACM Symp. on Principles of Prog. Lang.*, pages 27–40, New York, NY, USA, 2001. ACM Press.
- [41] D. Yu, N. A. Hamid, and Z. Shao. Building certified libraries for PCC: Dynamic storage allocation. In *Proc. 2003 European Symposium on Programming, LNCS Vol. 2618*, pages 363–379, 2003.
- [42] D. Yu and Z. Shao. Verification of safety properties for concurrent assembly code. In *Proc. 2004 ACM SIGPLAN Int’l Conf. on Functional Prog.*, September 2004.

A. The Soundness of CMAP

Lemma A.1 (Progress)

$\Phi = (Inv, \Delta, \Psi)$. If there exists Θ, p, g, A and G such that $\Phi; \Theta; (p, g, A, G) \vdash ((M, R), Q, T, C, I)$, then $(Inv M)$, and there exists a program \tilde{P} such that $((M, R), Q, T, C, I) \mapsto \tilde{P}$.

Proof sketch: By induction over the structure of I . $(Inv M)$ holds by the assumption and the inversion of the rule PROG. In the cases where I starts with `add`, `sub`, `movi`, `fork`, `exit`, or `yield`, the program can always make a step by the definition of the operational semantics (recall in our abstract machine we always assume there is an “idle” thread in Q , therefore the instructions `fork` and `exit` can go through). In the cases where I starts with `ld` and `st`, the side conditions for make a step, as defined by the operational semantics, are established by the rules LD and ST. In the cases where I starts with `bgt` or `beq`, or where I is `jd`, the operational semantics may fetch a code block from the code heap; such a code block exists by the inversion of the rule CDHP. \square

Lemma A.2 (Preservation)

If $\Phi; \Theta; (p, g, A, G) \vdash P$ and $P \mapsto \tilde{P}$, where $P = (S, Q, T, C, I)$ and $\tilde{P} = (\tilde{S}, \tilde{Q}, \tilde{T}, \tilde{C}, \tilde{I})$, then there exist $\tilde{\Theta}, \tilde{p}, \tilde{g}, \tilde{A}$ and \tilde{G} such that $\Phi; \tilde{\Theta}; (\tilde{p}, \tilde{g}, \tilde{A}, \tilde{G}) \vdash \tilde{P}$.

Proof sketch. By the assumption $\Phi; \Theta; (p, g, A, G) \vdash P$ and the inversion of the rule PROG, we know that

1. $(Inv, \Delta, \Psi) = \Phi, (M, R) = S, t = S(rt)$;
2. $\Phi \vdash T$;
3. $\Phi \vdash C$;
4. $(Inv \wedge p) S$;
5. $\Phi; (p, g, A, G) \vdash I$;
6. $\Phi; \Theta; (g, S) \vdash Q$;
7. $NI(\Theta\{t \rightsquigarrow (p, A, G)\}, Q\{t \rightsquigarrow (R, I)\})$ (see Section 4.4 for the definition of NI).

We prove the lemma by induction over the structure of I . Here we only give the detailed proof of cases where I starts with `fork` and `yield`. Proof for the rest cases are trivial.

Case $I = \text{fork } h, r; I''$.

By the operational semantics, we know that $\tilde{S} = S, \tilde{Q} = Q\{t' \rightsquigarrow (R', I')\}$, and $\tilde{I} = I''$, where $R' = \{r_0 \rightsquigarrow \dots, r_{15} \rightsquigarrow \dots, rt \rightsquigarrow t', ra \rightsquigarrow S(r)\}$, $I' = T(h)$, $t' \notin \text{dom}(Q)$ and $t' \neq t$. According to 5 and the inversion of the FORK rule, we know that there exist p', A', G', A'' and G'' such that

- f.1 $(p', A', G') = \Delta(h)$;
- f.2 $A \Rightarrow A'', G'' \Rightarrow G, (A \vee G'') \Rightarrow A', \text{ and } G' \Rightarrow (G \wedge A'')$;
- f.3 $\forall (M, R), \forall R'. \forall M'. (Inv \wedge p) (M, R) \Rightarrow R(rt) \neq R'(rt) \Rightarrow R(x) = R'(ra) \Rightarrow g R M M' \Rightarrow p' (M', R')$;
- f.4 $(Inv, \Delta, \Psi); (p, g, A'', G'') \vdash I''$

Then we let $\tilde{\Theta} = \Theta\{t' \rightsquigarrow (p', A', G')\}$, $\tilde{p} = p, \tilde{g} = g, \tilde{A} = A''$, and $\tilde{G} = G''$.

According to the rule PROG, to prove $\Phi; \tilde{\Theta}; (\tilde{p}, \tilde{g}, \tilde{A}, \tilde{G}) \vdash \tilde{P}$, we need prove the following:

- $\Phi \vdash T$ and $\Phi \vdash C$. They trivially follow 2 and 3. Since these two conditions never change, we will omit the proof of them in the following cases.
- $(Inv \wedge p) S$. By 4.
- $\Phi; (p, g, A'', G'') \vdash I''$. By f.4.
- $\Phi; \tilde{\Theta}; (g, S) \vdash \tilde{Q}$. By 6 and the inversion of the rule DTHRDS, we need check the following for the new thread t' :
 - $\forall M', M''. (Inv \wedge p')(M', R') \Rightarrow A' R' M' M'' \Rightarrow p'(M'', R')$. By 2 and the inversion of the rule THRDS.

- $(Inv, \Delta, \Psi); (p', G', A', G') \vdash I'$. By 2 and the inversion of the rule THRDS.
- $\forall M''. g S M'' \Rightarrow p'(M'', R')$. By 4, and f.3.

- $NI(\tilde{\Theta}\{t \rightsquigarrow (p, A'', G'')\}, \tilde{Q}\{t \rightsquigarrow (R, I'')\})$. By 4, 7, f.2, and Lemma 3.4.

Case $I = \text{yield}; I''$.

By the operational semantics, we know that $\tilde{S} = (M, R')$, $\tilde{Q} = (Q\{R(rt) \rightsquigarrow (R, I'')\}) \setminus \{t\}$, and $\tilde{I} = I''$, where $t \in \text{dom}(Q)$ and $(R', I') = Q(t)$ or $t = R(rt)$ and $(R', I') = (R, I'')$.

If $t = R(rt)$, the `yield` instruction is like a `no-op` instruction. The proof is trivial. So we just consider the case that $t \in \text{dom}(Q)$ and $(R', I') = Q(t)$.

We let $\tilde{\Theta} = (\Theta\{R(rt) \rightsquigarrow (p, A, G)\}) \setminus \{t\}$, $\tilde{p} = p', \tilde{g} = G', \tilde{A} = A'$, and $\tilde{G} = G'$, where $(p', A', G') = \Theta(t)$. According to the rule PROG, to prove $\Phi; \tilde{\Theta}; (\tilde{p}, \tilde{g}, \tilde{A}, \tilde{G}) \vdash \tilde{P}$, we need prove:

- $(Inv \wedge p') (M, R')$. By 4 we know that $Inv (M, R')$ (note that Inv is independent of R).
- $\Phi; (p', G', A', G') \vdash I'$. By 6 and the inversion of DTHRDS.
- $\Phi; \tilde{\Theta}; (G', (M, R')) \vdash \tilde{Q}$. To prove this, we only check the following for the yielding thread:
 - $\forall M', M''. (Inv \wedge p)(M', R) \Rightarrow A R M' M'' \Rightarrow p(M'', R)$. This follows 5 and the inversion of the rule YIELD.
 - $(Inv, \Delta, \Psi); (p, G, A, G) \vdash I''$. This follows 5 and the inversion of the rule YIELD.
 - $\forall M', M''. G' R' M' M'' \Rightarrow p(R, M'')$. By 7 we know $\forall M', M''. G' R' M' M'' \Rightarrow A R M' M''$. Therefore, we only need prove: $\forall M'. A R M' M'' \Rightarrow p(M'', R)$. By 5, the inversion of the rule YIELD, and 4 we get it.
- $NI(\tilde{\Theta}\{t \rightsquigarrow (p', A', G')\}, \tilde{Q}\{t \rightsquigarrow (R', I')\})$, which follows 7. \square

Theorem A.3 (Soundness)

$\Phi = (Inv, \Delta, \Psi)$. If there exists Θ, p, g, A and G such that $\Phi; \Theta; (p, g, A, G) \vdash P_0$, then for any $n \geq 0$, there exist M, R, Q, T, C and I such that $P_0 \mapsto^n ((M, R), Q, T, C, I)$ and $(Inv M)$.

Proof sketch. Given the progress and the preservation lemmas, this theorem can be easily proved by induction over n . However, instead of proving $(Inv M)$, we need strengthen the induction hypothesis and prove “ $\Phi; \tilde{\Theta}; (\tilde{p}, \tilde{g}, \tilde{A}, \tilde{G}) \vdash ((M, R), Q, T, C, I)$ for some $\tilde{\Theta}, \tilde{p}, \tilde{g}, \tilde{A}$ and \tilde{G} ”. \square

B. Partial Correctness of A Lock-Free Program

Figure 16 and 17 give the CMAP implementation of the GCD algorithm shown in Figure 15, where two threads collaborate to compute the greatest common divisor of two numbers. This is a lock-free algorithm, *i.e.*, no synchronization is required to ensure the non-interference even if atomic operations are machine instructions. To show the lock-free property, we insert `yield` instructions at every program point of the `chld` thread (some `yield` instructions are omitted in the `main` thread for clarity).

In Figure 18 we show definitions of assertions used to specify the program. The following formulae show the specifications of static threads, and the initial memory and instruction sequence.

$$\begin{aligned}
Inv &\equiv \text{True} \\
\Delta &\equiv \{\text{main} \rightsquigarrow (\text{True}, A_0, G_0), \text{chld} \rightsquigarrow (p_2 \wedge p_8, A_g, G_g)\} \\
\text{Initial } M &\equiv \{a \rightsquigarrow _, b \rightsquigarrow _, \text{flag} \rightsquigarrow _, (\text{flag} + 1) \rightsquigarrow _ \} \\
\text{Initial } I &\equiv \text{jd begn}
\end{aligned}$$

```

Variables:
  nat a, b;
  nat[2] flag;

Main :
  a :=  $\alpha$ ; b :=  $\beta$ ;
  flag[0] := 0;
  fork(gcd, 0);
  flag[1] := 0;
  fork(gcd, 1);
  while (!flag[0])
    yield;
  while (!flag[1])
    yield;
  post proc...

void gcd(int x){
  while (a  $\neq$  b){
    if ((a > b)&&(x = 0))
      a := a - b;
    else if ((a < b)&&(x = 1))
      b := b - a;
  }
  flag[x] := 1;
}

```

Figure 15. GCD: high-level program

```

main :  $-\{(\text{True}, A_0, G_0)\}$       next :  $-\{p_2, g, A_1, G_1\}$ 
  jd begn                          yield
begn :  $-\{(\text{True}, G_0, A_0, G_0)\}$   st r1(1), r0
  movi r0,  $\alpha$                     yield
  movi r1, a                         movi r2, 1
  st r1(0), r0                       fork chld, r2
  yield                               join1 :  $-\{p_2 \wedge p_3, g, A_2, G_2\}$ 
  movi r0,  $\beta$                       yield
  movi r1, b                         ld r2, r1(0)
  st r1(0), r0                       beq r0, r2, join1
  yield                               join2 :  $-\{p_2 \wedge p_4, g, A_2, G_2\}$ 
  movi r0, 0                          yield
  movi r1, flag                      ld r2, r1(1)
  st r1(0), r0                       beq r0, r2, join2
  yield                               post :  $-\{p_2 \wedge p_5, g, A_2, G_2\}$ 
  movi r2, 0                          yield
  fork chld, r2                       ...

```

Figure 16. GCD: the main thread

Inv is simply set to True . The partial correctness of the gcd algorithm is inferred by the fact that p_2 is established at the entry point and G_g is satisfied.

Two child threads are created using the same static thread chld . They use thread arguments to distinguish their task. Correspondingly, the assumption A_g and G_g of the static thread chld also use the thread argument to distinguish the specification of different dynamic copies. The child thread does not change its assumption and guarantee throughout its lifetime. Therefore we omit A_g and G_g in the code heap specifications. Since we insert yield instructions at every program point, every local guarantee is simply G_g , which is also omitted in the specifications.

At the data initiation stage, the main thread assumes no threads in the environment changes a , b and the flags, and guarantees nothing, as shown in A_0 and G_0 (see Figure 18). After creating child threads, the main thread changes its assumption and guarantee to reflect the changing environment. The new assumption is just the disjunction of the previous assumption and the guarantee of the new thread (instantiated by the thread id and thread argument), similarly the new guarantee is the conjunction of the previous guarantee and the new thread's assumption.

This example also shows how thread join can be implemented in CMAP by synchronization. We use one flag for each thread to indicate if the thread is alive or not. The assumptions and guarantees of the main thread also take advantages of these flags so that it can weaken its guarantee and strengthen its assumption after the child threads die.

```

chld :  $-\{p_2 \wedge p_8, A_g, G_g\}$       calc1 :  $-\{p_2 \wedge p_9 \wedge p_{12}\}$ 
  jd gcd                             yield
gcd :  $-\{p_2 \wedge p_8\}$               beq r0, ra, cld1
  yield                               yield
  movi r0, 0                          jd loop
loop :  $-\{p_2 \wedge p_8 \wedge p_9\}$     cld1 :  $-\{p_2 \wedge p_9 \wedge p_{13}\}$ 
  yield                               yield
  ld r1, r0(a)                        sub r3, r1, r2
  yield                               yield
   $-\{p_2 \wedge p_9 \wedge p_{10}\}$       st r0(a), r3
  ld r2, r0(b)                        yield
  yield                               jd loop
   $-\{p_2 \wedge p_9 \wedge p_{11}\}$     calc2 :  $-\{p_2 \wedge p_9 \wedge p_{14}\}$ 
  beq r1, r2, done                    yield
  bgt r1, r2, calc1                   movi r3, 1
  yield                               beq r3, ra, cld2
  jd calc2                             yield
done :  $-\{p_1 \wedge p_2 \wedge p_8\}$     cld2 :  $-\{p_2 \wedge p_9 \wedge p_{15}\}$ 
  yield                               yield
  movi r0, 1                          sub r3, r2, r1
  yield                               yield
  st ra(flag), r0                     st r0(b), r3
  exit                               yield
                                       jd loop

```

Figure 17. GCD: the child thread

```

id1  $\equiv a = a'$       id2  $\equiv b = b'$ 
id3  $\equiv (\text{flag}[0] = \text{flag}'[0]) \wedge (\text{flag}[1] = \text{flag}'[1])$ 
G0  $\equiv \text{True}$       A0  $\equiv \text{id}_1 \wedge \text{id}_2 \wedge \text{id}_3$ 
p  $\equiv \text{id}_2 \wedge (a > b \Rightarrow (\text{GCD}(a, b) = \text{GCD}(a', b'))) \wedge (a \leq b \Rightarrow \text{id}_1)$ 
p'  $\equiv \text{id}_1 \wedge (a < b \Rightarrow (\text{GCD}(a, b) = \text{GCD}(a', b'))) \wedge (a \geq b \Rightarrow \text{id}_2)$ 
Gg  $\equiv (\text{flag}'[ra] = 0 \Rightarrow (([ra] = 0 \wedge p) \vee ([ra] = 1 \wedge p'))) \wedge (\text{flag}'[ra] = 1 \Rightarrow (\text{id}_1 \wedge \text{id}_2 \wedge a' = b')) \wedge (\text{flag}[ra] = 1 \Rightarrow S = S') \wedge (\text{flag}[1 - ra] = \text{flag}'[1 - ra])$ 
Ag  $\equiv (\text{flag}'[ra] = 0) \Rightarrow ((\text{flag}[ra] = \text{flag}'[ra]) \wedge (([ra] = 0 \wedge p') \vee ([ra] = 1 \wedge p)))$ 
G1  $\equiv G_0 \wedge (\text{flag}'[0] = 0 \Rightarrow ((\text{flag}[0] = \text{flag}'[0]) \wedge p'))$ 
A1  $\equiv A_0 \vee ((\text{flag}[1] = \text{flag}'[1]) \wedge (\text{flag}'[0] = 0 \Rightarrow p) \wedge (\text{flag}'[0] = 1 \Rightarrow (\text{id}_1 \wedge \text{id}_2 \wedge a' = b'))) \wedge (\text{flag}[0] = 1 \Rightarrow S = S')$ 
G2  $\equiv G_1 \wedge (\text{flag}'[1] = 0 \Rightarrow ((\text{flag}[1] = \text{flag}'[1]) \wedge p))$ 
A2  $\equiv A_1 \vee ((\text{flag}[0] = \text{flag}'[0]) \wedge (\text{flag}'[1] = 0 \Rightarrow p') \wedge (\text{flag}'[1] = 1 \Rightarrow (\text{id}_1 \wedge \text{id}_2 \wedge a' = b'))) \wedge (\text{flag}[1] = 1 \Rightarrow S = S')$ 
g  $\equiv S = S'$       p1  $\equiv a = b$       p2  $\equiv \text{GCD}(a, b) = \text{GCD}(a', b)$ 
p3  $\equiv \text{flag}[0] = 0 \vee (\text{flag}[0] = 1 \wedge p_1)$ 
p4  $\equiv (\text{flag}[1] = 0 \vee \text{flag}[1] = 1) \wedge \text{flag}[0] = 1 \wedge p_1$ 
p5  $\equiv \text{flag}[0] = 1 \wedge \text{flag}[1] = 1 \wedge p_1$ 
p6  $\equiv [ra] = 0 \wedge \text{flag}[0] = 0$       p7  $\equiv [ra] = 1 \wedge \text{flag}[1] = 0$ 
p8  $\equiv p_6 \vee p_7$       p9  $\equiv [x_0] = 0$ 
p10  $\equiv (p_6 \wedge [x_1] = a) \vee (p_7 \wedge ([x_1] \leq b \Rightarrow [x_1] = a))$ 
p11  $\equiv (p_6 \wedge [x_1] = a \wedge ([x_1] \geq [x_2] \Rightarrow [x_2] = b)) \vee (p_7 \wedge [x_2] = b \wedge ([x_1] \leq [x_2] \Rightarrow [x_1] = a))$ 
p12  $\equiv [x_1] > [x_2] \wedge ([ra] = 0 \Rightarrow (\text{flag}[0] = 0 \wedge [x_1] = a \wedge [x_2] = b))$ 
p13  $\equiv p_6 \wedge [x_1] > [x_2] \wedge [x_1] = a \wedge [x_2] = b$ 
p14  $\equiv [x_1] < [x_2] \wedge ([ra] = 1 \Rightarrow (\text{flag}[1] = 0 \wedge [x_1] = a \wedge [x_2] = b))$ 
p15  $\equiv p_7 \wedge [x_1] < [x_2] \wedge [x_1] = a \wedge [x_2] = b$ 

```

Figure 18. GCD: Assertion Definitions