

# Reasoning About a Machine with Local Capabilities

Provably Safe Stack and Return Pointer Management

LAU SKORSTENGAARD, Aarhus University, Denmark

DOMINIQUE DEVRIESE, KU Leuven, Belgium

LARS BIRKEDAL, Aarhus University, Denmark

## ACM Reference Format:

Lau Skorstengaard, Dominique Devriese, and Lars Birkedal. 2010. Reasoning About a Machine with Local Capabilities: Provably Safe Stack and Return Pointer Management. *ACM Trans. Web* 9, 4, Article 39 (March 2010), 45 pages. <https://doi.org/0000001.0000001>

## 1 INTRODUCTION

Compromising software security is often based on attacks that break programming language properties relied upon by software authors, such as control-flow correctness, local-state encapsulation, etc. Commodity processors offer little support for defending against such attacks: they offer security primitives with only coarse-grained memory protection and limited compartmentalization scalability. As a result, defenses against attacks on control-flow correctness and local-state encapsulation are either limited to only certain common forms of attacks (leading to an attack-defense arms race) and/or rely on techniques like machine code rewriting [Abadi et al. 2005; Wahbe et al. 1993], machine code verification [Morrisett et al. 1999], virtual machines with a native stack [Lindholm et al. 2014] or randomization [Forrest et al. 1997]. The latter techniques essentially emulate protection techniques on existing hardware, at the cost of performance, system complexity and/or security.

Capability machines are a type of processors that remediate these limitations with a better security model at the hardware level. They are based on old ideas [Carter et al. 1994; Dennis and Van Horn 1966; Shapiro et al. 1999], but have recently received renewed interest; in particular, the CHERI project has proposed new ideas and ways of tackling practical challenges like backwards compatibility and realistic OS support [Watson et al. 2015; Woodruff et al. 2014]. Capability machines tag every word (in the register file and in memory) to enforce a strict separation between numbers and capabilities (a kind of pointers that carry authority). Memory capabilities carry the authority to read and/or write to a range of memory locations. There is also a form of object capabilities, which represent the authority to invoke a piece of code without exposing the code's encapsulated private state (e.g., the M-Machine's enter capabilities or CHERI's sealed code/data pairs).

Unlike commodity processors, capability machines lend themselves well to enforcing local-state encapsulation. Potentially, they will enable compilation schemes that enforce this property in an efficient but also 100% watertight way (ideally evidenced by a mathematical proof, guaranteeing

---

Authors' addresses: Lau Skorstengaard, Department of Computer Science, Aarhus University, Street1 Address1, Aarhus, State1, 8210, Denmark, [lau@cs.au.dk](mailto:lau@cs.au.dk); Dominique Devriese, iMinds-DistriNet, KU Leuven, Street1 Address1, -, -, Belgium, [dominique.devriese@cs.kuleuven.be](mailto:dominique.devriese@cs.kuleuven.be); Lars Birkedal, Department of Computer Science, Aarhus University, Street1 Address1, Aarhus, State1, 8210, Denmark, [birkedal@cs.au.dk](mailto:birkedal@cs.au.dk).

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2009 Copyright held by the owner/author(s). Publication rights licensed to the Association for Computing Machinery.

1559-1131/2010/3-ART39 \$15.00

<https://doi.org/0000001.0000001>

that we do not end up in a new attack-defense arms race). However, a lot needs to happen before we get there. For example, it is far from trivial to devise a compilation scheme adapted to the details of a specific source language's notion of encapsulation (e.g., private member variables in OO languages often behave quite differently than private state in ML-like languages). And even if such a scheme were defined, a formal proof depends on a formalization of the encapsulation provided by the capability machine at hand.

A similar problem is the enforcement of control-flow correctness on capability machines. An interesting approach is taken in CheriBSD [Watson et al. 2015]: the standard contiguous C stack is split into a central, trusted stack, managed by trusted call and return instructions, and disjoint, private, per-compartment stacks. To prevent illegal use of stack references, the approach relies on local capabilities, a type of capabilities offered by CHERI to temporarily relinquish authority, namely for the duration of a function invocation whereafter the capability can be revoked. However, details are scarce (how does it work precisely? what features are supported?) and a lot remains to be investigated (e.g., combining disjoint stacks with cross-domain function pointers seems like it will scale poorly to large numbers of components?). Finally, there is no argument that the approach is watertight and it is not even clear what security property is targeted exactly.

DD: I think we should include a comparison with the conference paper, and explain what the delta is for this journal version.

In this paper, we make two main contributions: (1) an alternative calling convention that uses local capabilities to enforce stack frame encapsulation and well-bracketed control flow, and (2) perhaps more importantly, we adapt and apply the well-studied techniques of step-indexed Kripke logical relations for reasoning about code on a representative capability machine with local capabilities in general and correctness and security of the calling convention in particular. More specifically, we make the following contributions:

- We formalize a simple but representative capability machine featuring local capabilities and its operational semantics (Section 2).
- We define a novel calling convention enforcing control-flow correctness and encapsulation of stack frames (Section 3). It relies solely on local capabilities and does not require OS support (like a trusted stack or call/return instructions). It supports higher-order cross-component calls (e.g., cross-component function pointers) and can be efficient assuming only one additional piece of processor support: an efficient instruction for clearing a range of memory.
- We present a novel step-indexed Kripke logical relation for reasoning about programs on the capability machine. It is an untyped logical relation, inspired by previous work on object capabilities [Devriese et al. 2016]. We prove an analogue of the standard fundamental theorem of logical relations — to the best of our knowledge, our theorem is the most general and powerful formulation of the formal guarantees offered by a capability machine (a form of capability safety [Devriese et al. 2016; Maffei et al. 2010]), including the specific guarantees offered for local capabilities. It is very general and not tied to our calling convention or a specific way of using the system's capabilities. We are the first to apply these techniques for reasoning about capability machines and we believe they will prove useful for many other purposes than our calling convention.
- We introduce two novel technical ideas in the unary, step-indexed Kripke logical relation used to formulate the above theorem: the use of a single orthogonal closure (rather than the earlier used biorthogonal closure) and a variant of Dreyer et al. [2012]'s public and private future worlds [Dreyer et al. 2012] to express the special nature of local capabilities. The

```

99       $a \in \text{Addr} \stackrel{\text{def}}{=} \mathbb{N}$ 
100       $w \in \text{Word} \stackrel{\text{def}}{=} \mathbb{Z} + \text{Cap}$ 
101       $r \in \text{RegName} ::= \text{pc} \mid r_0 \mid r_1 \mid \dots$ 
102       $\text{perm} \in \text{Perm} ::= \text{O} \mid \text{RO} \mid \text{RW} \mid \text{RWL} \mid$ 
103       $\text{RX} \mid \text{E} \mid \text{RWX} \mid \text{RWLX}$ 
104       $g \in \text{Global} ::= \text{global} \mid \text{local}$ 
105       $\text{Conf} ::= \text{ExecConf} + \{\text{failed}\} + \{\text{halted}\} \times \text{Mem}$ 
106       $\text{Cap} ::= \{((\text{perm}, g), b, e, a) \mid b, a \in \text{Addr}, e \in \text{Addr} \cup \{\infty\}\}$ 
107
108       $r \in \mathbb{Z} + \text{RegName}$ 
109       $i ::= \text{jmp } r \mid \text{jnz } r r \mid \text{move } r r \mid \text{load } r r \mid \text{store } r r \mid \text{plus } r r r \mid \text{minus } r r r \mid$ 
110       $\text{lt } r r r \mid \text{lea } r r \mid \text{restrict } r r \mid \text{subseg } r r r \mid \text{isptr } r r \mid \text{getl } r r \mid$ 
111       $\text{getp } r r \mid \text{getb } r r \mid \text{gete } r r \mid \text{geta } r r \mid \text{fail} \mid \text{halt}$ 

```

Fig. 1. The syntax of our capability machine assembly language.

logical relation and the fundamental theorem expressing capability safety are presented in Section 4.

- We demonstrate our results by applying them to challenging examples, specifically constructed to demonstrate local-state encapsulation and control-flow correctness guarantees in the presence of cross-component function pointers (Section 8). The examples demonstrate both the power of our formulation of capability safety and our calling convention.

We have written a technical appendix [Skorstengaard et al. 2018] which contains additional details and proofs left out from this paper.

## 2 A CAPABILITY MACHINE WITH LOCAL CAPABILITIES

In this paper, we work with a formal capability machine with all the characteristics of real capability machines, as well as local capabilities much like CHERI's. Otherwise, it is kept as simple as possible. It is inspired by both the M-Machine [Carter et al. 1994] and CHERI [Watson et al. 2015]. To avoid uninteresting details, we assume an infinite address space and unbounded integers.

We define the syntax of our capability machine in Figure 1. We assume an infinite set of addresses  $\text{Addr}$  and define machine words as either integers or capabilities of the form  $((\text{perm}, g), \text{base}, \text{end}, a)$ . Such a capability represents the authority to execute permissions  $\text{perm}$  on the memory range  $[\text{base}, \text{end}]$ , together with a current address  $a$  and a locality tag  $g$  indicating whether the capability is global or local. There is no notion of pointers other than capabilities, so we will use the terms interchangeably. The available permissions are null permission (O), readonly (RO), read/write (RW), read/execute (RX) and read/write/execute (RWX) permissions. Additionally, there are three special permissions: read/write-local (RWL), read/write-local/execute (RWLX) and enter (E), which we will explain below. Permissions and locality are ordered and both orderings are depicted in Figure 3. We denote the pairwise ordering of permission and locality with  $\sqsubseteq$ .

We assume a finite set of register names  $\text{RegName}$ . We define register files  $\text{reg}$  and memories  $\text{ms}$  as functions mapping register names resp. addresses to words. The state of the entire machine is

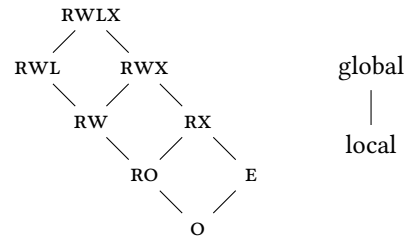


Fig. 3. Permission and locality hierarchy.

$$\Phi \rightarrow \begin{cases} \llbracket decode(n) \rrbracket (\Phi) & \text{if } \Phi.\text{reg}(\text{pc}) = ((\text{perm}, g), b, e, a) \text{ and } b \leq a \leq e \\ & \text{and } \text{perm} \in \{\text{RX}, \text{RWX}, \text{RWLX}\} \text{ and } \Phi.\text{mem}(a) = n \\ \text{failed} & \text{otherwise} \end{cases}$$

$$\text{updPc}(\Phi) = \begin{cases} \Phi[\text{reg.pc} \mapsto \text{newPc}] & \text{if } \Phi.\text{reg}(\text{pc}) = ((\text{perm}, g), b, e, a) \\ & \text{and } \text{newPc} = ((\text{perm}, g), b, e, a + 1) \\ \text{failed} & \text{otherwise} \end{cases}$$

$i$	$\llbracket i \rrbracket (\Phi)$	Conditions
fail	<i>failed</i>	
halt	$(\text{halted}, \Phi.\text{mem})$	
move $r_1 \ r_2$	$\text{updPc}(\Phi[\text{reg}.r_1 \mapsto w])$	$r_2 \in \text{Reg} \Rightarrow w = \Phi.\text{reg}(r_2)$ and $r_2 \in \mathbb{Z} \Rightarrow w = r_2$
load $r_1 \ r_2$	$\text{updPc}(\Phi[\text{reg}.r_1 \mapsto w])$	$\Phi.\text{reg}(r_2) = ((\text{perm}, g), b, e, a)$ and $w = \Phi.\text{mem}(a)$ and $b \leq a \leq e$ and $\text{perm} \in \{\text{RWX}, \text{RWLX}, \text{RX}, \text{RW}, \text{RWL}, \text{RO}\}$
store $r_1 \ r_2$	$\text{updPc}(\Phi[\text{mem}.a \mapsto w])$	$\Phi.\text{reg}(r_1) = ((\text{perm}, g), b, e, a)$ and $\text{perm} \in \{\text{RWX}, \text{RWLX}, \text{RW}, \text{RWL}\}$ and $b \leq a \leq e$ and $w = \Phi.\text{reg}(r_2)$ and if $w = ((\_, \text{local}), \_, \_, \_)$ , then $\text{perm} \in \{\text{RWLX}, \text{RWL}\}$
jmp $r$	$\Phi[\text{reg.pc} \mapsto \text{newPc}]$	if $\Phi.\text{reg}(r) = ((\text{E}, g), b, e, a)$ , then $\text{newPc} = ((\text{RX}, g), b, e, a)$ otherwise $\text{newPc} = \Phi.\text{reg}(r)$
restrict $r_1 \ r_2$	$\text{updPc}(\Phi[\text{reg}.r_1 \mapsto w])$	$\Phi.\text{reg}(r_2) = ((\text{perm}, g), b, e, a)$ and $(\text{perm}', g') = \text{decodePermPair}(\Phi.\text{reg}(r_2))$ and $(\text{perm}', g') \sqsubseteq (\text{perm}, g)$ and $w = ((\text{perm}', g'), b, e, a)$
subseg $r_1 \ r_2 \ r_3$	$\text{updPc}(\Phi[\text{reg}.r_1 \mapsto w])$	$\Phi.\text{reg}(r_1) = ((\text{perm}, g), b, e, a)$ and for $i \in \{2, 3\}$ $n_i = \Phi.\text{reg}(r_i)$ and $n_2 \in \mathbb{N}$ and $b \leq n_2$ and $n_3 \leq e$ where either $n_3 \in \mathbb{N}$ or $(n_3 = -42 \text{ and } e = \infty)$ and $\text{perm} \neq \text{E}$ and $w = ((\text{perm}, g), n_1, n_2, a)$
lea $r_1 \ r_2$	$\text{updPc}(\Phi[\text{reg}.r_1 \mapsto c])$	$\Phi.\text{reg}(r_1) = ((\text{perm}, g), b, e, a)$ and $n = \Phi.\text{reg}(r_2)$ and $n \in \mathbb{Z}$ and $\text{perm} \neq \text{E}$ and $c = ((\text{perm}, g), b, e, a + n)$
geta $r_1 \ r_2$	$\text{updPc}(\Phi[\text{reg}.r_1 \mapsto a])$	$\Phi.\text{reg}(r_2) = ((\_, \_), \_, \_, a)$
		...
–	<i>failed</i>	otherwise

Fig. 2. An excerpt from the operational semantics.

represented as a configuration that is either a running state  $\Phi \in \text{ExecConf}$  containing a memory and a register file, or a failed or halted state, where the latter keeps hold of the final state of memory.

The machine's instruction set is rather basic. Instructions  $i$  include relatively standard jump (jmp), conditional jump (jnz) and move (move, copies words between registers) instructions. Also familiar are load and store instructions for reading from and writing to memory (load and store) and arithmetic operators (lt (less than), plus and minus, operating only on numbers). There are three instructions for modifying capabilities: lea (modifies the current address), restrict (modifies the permission and local/global tag) and subseg (modifies the range of a capability). Importantly, these

instructions take care that the resulting capability always carries less authority than the original (e.g. `restrict` will only weaken a permission). Finally, the instruction `isptr` tests whether a word is a capability or a number and instructions `getp`, `getl`, `getb`, `gete` and `geta` provide access to a capability's permissions, local/global tag, base, end and current address, respectively.

Figure 2 shows the operational semantics for a few representative instructions. Essentially, a configuration  $\Phi$  either decodes and executes the instruction at  $\Phi.\text{reg}(\text{pc})$  if it is executable and its address is in the valid range or otherwise fails. The table in the figure shows for instructions  $i$  the result of executing them in configuration  $\Phi$ . The instructions `fail` and `halt` obviously fail and `halt` respectively. `move` simply modifies the register file as requested and updates the pc to the next instruction using the meta-function *updPc*.

DD: Follow order of figure in discussion?

The load instruction loads the contents of the requested memory location into a register, but only if the capability has appropriate authority (i.e. read permission and an appropriate range). `restrict` updates a capability's permissions and global/local tag in the register file, but only if the new permissions are weaker than the original. It also never turns local capabilities into global ones. `geta` queries the current address of a capability and stores it in a register.

The `jmp` instruction updates the program counter to a requested location, but it is complicated by the presence of `enter` capabilities, modeled after the M-Machine's [Carter et al. 1994]. `Enter` capabilities cannot be used to read, write or execute and their address and range cannot be modified. They can only be used to jump to, but when that happens, their permission changes to `rx`. They can be used to represent a kind of closures: an opaque package containing a piece of code together with local encapsulated state. Such a package can be built as an `enter` capability  $c = ((E, g), b, e, a)$  where the range  $[b, a - 1]$  contains local state (data or capabilities) and  $[a, e]$  contains instructions. The package is opaque to an adversary holding  $c$  but when  $c$  is jumped to, the instructions can start executing and have access to the local data through the updated version of  $c$  that is then in pc. LS: As we talked about what is described here, is not how we make closures in the TR, so we should consider whether we should refrain from using that word here.

Finally, the store instruction updates the memory to the requested value if the capability has write authority for the requested location. However, the instruction is complicated by the presence of `local` capabilities, modeled after the ones in the CHERI processor [Watson et al. 2015]. Basically, local capabilities are special in that they can only be kept in registers, i.e. they cannot be stored to memory. This means that local capabilities can be temporarily given to an adversary, for the duration of an invocation: if we take care to clear the capability from the register file after control is passed back to us, they will not have been able to store the capability. However, there is one exception to the rule above: local capabilities can be stored to memory for which we have a capability with write-local authority (i.e. permission `RWL` or `RWLX`). This is intended to accommodate a stack, where register contents can be stored, including local capabilities. As long as all capabilities with write-local authority are themselves local and the stack is cleared after control is passed back by the adversary, we will see that this does not break the intended behavior of local capabilities.

We point out that our local capabilities capture only a part of the semantics of local capabilities in CHERI. Specifically, in addition to the above, CHERI's default implementation of the CCall exception handler forbids local capabilities from being passed across module boundaries. Such a restriction fundamentally breaks our calling convention, since we pass around local return pointers and stack capabilities. However, CHERI's CCall is not implemented in hardware, but in software, precisely to allow experimenting with alternative models like ours. DD: Update to CHERI's new ccall story.

In order to have a reasonably realistic system, we use a simple model of linking where a program has access to a linking table that contains capabilities for other programs. We also assume malloc to be part of the trusted computing base satisfying a certain specification. Malloc and linking tables are described further in the next section. The specification of malloc is presented in Section 5 as it uses the semantic model we build in Section 4. For full details on the linking table, we refer to the technical appendix [Skorstengaard et al. 2018]. (Reviewer A, esop suggest explaining the linking model)

Suggestion: We also have the flag table used for assertions, maybe we should add a description of it here.

### 3 STACK AND RETURN POINTER MANAGEMENT USING LOCAL CAPABILITIES

Suggestion: say something about arguments in this paragraph instead of deferring it to the discussion only?

Suggestion: Add some of our presentation illustration figures to illustrate some of the issues? May not be feasible as they rely on “animation”.

Suggestion: Include a description of (possibly definition for) well-bracketedness? (requested by reviewer C, popl). Reviewer D, popl, asks for the same but also a definition of local state encapsulation. In the POPL notes, I wrote that I talked with Lars about it and this is something that is not easily defined. Perhaps we should have a discussion of this.

One of the contributions in this paper is a demonstration that local capabilities on a capability machine support a calling convention that enforces control-flow correctness in a way that is provably watertight, potentially efficient, does not rely on a trusted central stack manager and supports higher-order interfaces to an adversary, where an adversary is just some unknown piece of code. In this section, we explain this convention’s high-level approach, the security measures to be taken in a number of situations (motivating each separately with a summary table at the end). After that, we define a number of reusable macro-instructions that can be used to conveniently apply the proposed convention in subsequent examples.

The basic idea of our approach is simple: we stick to a single, rather standard, C stack and register-passed stack and return pointers, much like a standard C calling convention. However, to prevent various ways of misusing this basic scheme, we put local capabilities to work and take a number of not-always-obvious safety measures. The safety measures are presented in terms of what we need to do to protect ourselves against an adversary, but this is only for presentation purposes as our code assumes no special status on the machine. In fact, an adversary can apply the same safety measures to protect themselves against us. In the next paragraphs, we will explain the issues to be considered in all the relevant situations: when (1) starting our program, (2) returning to the adversary, (3) invoking the adversary, (4) returning from the adversary, (5) invoking an adversary callback and (6) having a callback invoked by the adversary.

**Program start-up** We assume that the language runtime initializes the memory as follows: a contiguous array of memory is reserved for the stack, for which we receive a stack pointer in a special register  $r_{stk}$ . We stress that the stack is not built-in, but merely an abstraction we put on this piece of the memory. The stack pointer is local and has RWLX permission. Note that this means that we will be placing and executing instructions on the stack. Crucially, the stack is the only part of memory for which the runtime (including malloc, loading, linking) will ever provide RWLX or RWL capabilities. Additionally, our examples typically also assume some memory to store instructions or static data. Another part of memory (called the heap) is initially governed by malloc and at program start-up, no other code has capabilities for this memory. Malloc hands out RWX



capabilities for allocated regions as requested (no RWLX or RWL permissions). For simplicity, we assume that memory allocated through malloc cannot be freed.

**Returning to the adversary** Perhaps the simplest situation is returning to the adversary after they invoked our code (Figure 4a). In this case, we have received a return pointer from them, and we just need to jump to it as usual. An obvious security measure to take care of is properly clearing the non-return-value registers before we jump (since they may contain data or capabilities that the adversary should not get access to, Figure 4b). Additionally, we may have used the stack for various purposes (register spilling, storing local state when invoking other functions etc.), so we also need to clear that data before returning to the adversary (Figure 4c).

However, if we are returning from a function that has itself invoked adversary code (Figure 4d), then clearing the used part of the stack is not enough. The unused part of the stack may also contain data and capabilities, left there by the adversary, including local capabilities since the stack is write-local. As we will see later, we rely on the fact that the adversary cannot keep hold of local capabilities when they pass control to the trusted code and receive control back. In this case, the adversary could use the unused part of the stack to store local pointers and load them from there after they get control back (Figure 4e). To prevent this, we need to clear (i.e. overwrite with zeros) the entire part of the stack that the adversary has had access to, not just the parts that we have used ourselves (Figure 4f). Since we may be talking about a large part of memory, this requirement is the most problematic aspect of our calling convention for performance, but see Section 9 for how this might be mitigated.

**Invoking the adversary** A slightly more complex case is invoking the adversary. As above, we clear all the non-argument registers, as well as the part of the stack that we are not using (because, as above, it may contain local capabilities from previously executed code that the adversary could exploit in the same way). We leave a copy of the stack pointer in  $r_{stk}$ , but only after we have used the subseg instruction to shrink its authority to the part that we are not using ourselves.

In one of the registers, we also provide a return pointer, which must be a local capability. If it were global, the adversary would be able to store away the return pointer in a global data structure (i.e. there exists a global capability for it), and jump to it later, in circumstances where this should not be possible. For example, they could store the return pointer, legally jump to it a first time, wait to be invoked again and then jump to the old return pointer a second time, instead of the new return pointer received for the second invocation. Similarly, they could store the return pointer, invoke a function in our code, wait for us to invoke them again and then jump to the old return pointer rather than the new one, received for the second invocation. By making the return pointer local, we prevent such attacks: the adversary can only store local capabilities through write-local capabilities, which means (because of our assumptions above): on the stack. Since the stack pointer itself is also local, it can also only be stored on the stack. Because we clear the part of the stack that the adversary has had access to before we pass control back, there is no way for them to recover either of these local capabilities.

Note that storing stack pointers for use during future invocations would also be dangerous in itself, i.e. not just because it can be used to store return pointers. Imagine the adversary stores their stack pointer, invokes trusted code that uses part of the stack to store private data and then invokes the adversary again with a stack pointer restricted to exclude the part containing the private data. If the adversary had a way of keeping hold of their old stack pointer, it could access the private data stored there by the trusted code and break local-state encapsulation.

**Returning from the adversary** So return pointers must be passed as local capabilities. But what should their permissions be, what memory should they point to and what should that memory (the activation record) contain? Let us answer the last question first by considering what should happen when the adversary jumps to a return pointer. In that case, the program counter should be

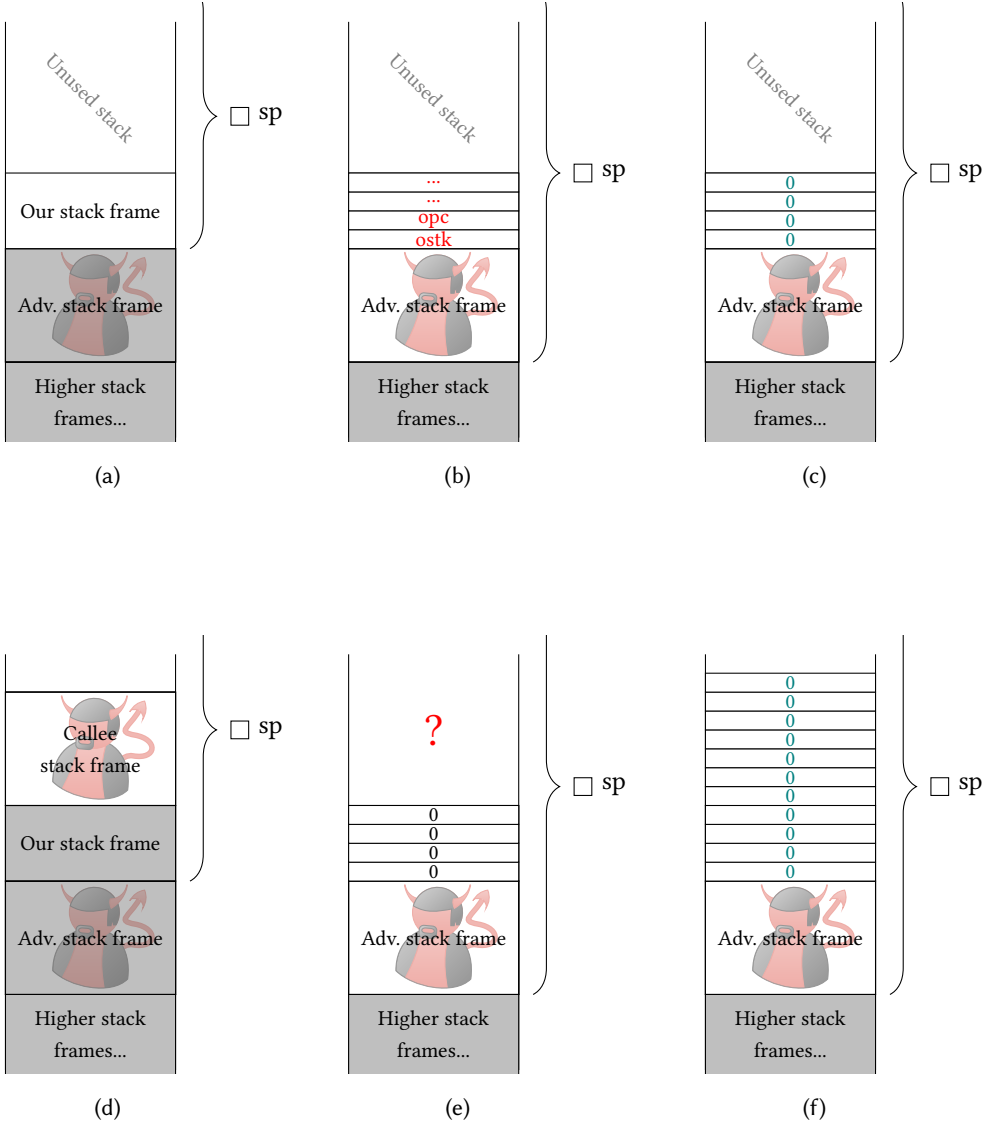


Fig. 4. Illustration of the situations related to stack clearing when returning to an adversary. The greyed out areas are parts of the stack that are not accessible at the time.

restored to the instruction after the jump to the adversary, so the activation record should store this old program counter. Additionally, the stack pointer should also be restored to its original value. Since the adversary has a more restricted authority over the stack than the code making the



call, we cannot hope to reconstruct the original stack pointer from the stack pointer owned by the adversary. Instead, it should be stored as part of the activation record.

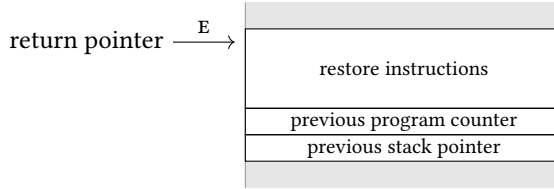


Fig. 5. Structure of an activation record

Clearly, neither of these capabilities should be accessible by the adversary. In other words, the return pointer provided to the adversary must be a capability that they can jump to but not read from, i.e. an enter capability. To make this work, we construct the activation record as depicted in Figure 5. The  $E$  return pointer has authority over the entire activation record (containing the previous return and stack pointer), and its current address points to a number of restore instructions in the record, so that upon invocation, these instructions are executed and can load the old stack pointer and program counter back into the register file. As the return pointer is an enter pointer, the adversary cannot get hold of the activation record's contents, but after invocation, its permission is updated to  $RX$ , so the contents become available to the restore instructions.

The final question that remains is: where should we store this activation record? The attentive reader may already see that there is only one possibility: since the activation record contains the old stack pointer, which is local, the activation record can only be constructed in a part of memory where we have write-local access, i.e. on the stack. Note that this means we will be placing and executing instructions on the stack, i.e. it will not just contain code pointers and data. This means that our calling convention should be combined with protection against stack smashing attacks (i.e. buffer overflows on the stack overwriting activation records' contents). Luckily, the capability machine's fine-grained memory protection should make it reasonably easy for a compiler to implement such protection, by making sure that only appropriately bounded versions of the stack pointer are made available to source language code.

**Invoking an adversary callback** If we have a higher-order interface to the adversary, we may need to invoke an adversary callback. In this case, not so much changes with respect to the situation where we invoke static adversary code. The adversary can provide a callback as a capability for us to jump to, either an  $E$ -capability if they want to protect themselves from us or just an  $RX$  capability if they are not worried about that. However, there is one scenario that we need to prevent: if they construct the callback capability to point into the stack, it may contain local capabilities that they should not have access to upon invocation of the callback. As before, this includes return and stack pointers from previous stack frames that they may be trying to illegally use inside the callback.

To prevent this, we only accept callbacks from the adversary in the form of global capabilities, which we dynamically check before invoking them (and we fail otherwise). This should not be an overly strict requirement: our own callbacks do not contain local data themselves, so there should be no need for the adversary to construct callbacks on the stack.<sup>1</sup> *Maybe a word on arguments passed from the adversary to us that we pass on to a callback given by the adversary (mentioned by reviewer E, popl, see tex comment.)*

<sup>1</sup>Note that it does prevent a legitimate but non-essential scenario where the adversary wants to give us temporary access to a callback not allocated on the stack.

**Having a callback invoked by the adversary** The above leaves us with perhaps the hardest scenario: how to provide a callback to the adversary. The basic idea is that we allocate a block of memory using malloc that we fill with the capabilities and data that the callback needs, as well as some prelude instructions that load the data into registers and jumps to the right code. Note that this implies that no local capabilities can be stored as part of a closure. We can then provide the adversary with an enter-capability covering the allocated block and pointing to the contained prelude instructions. However, the question that remains in this setup is: from where do we get a stack pointer when the callback is invoked?

Our answer is that the adversary should provide it to us, just as we provide them with a stack pointer when we invoke their code. However, it is important that we do not just accept any capability as a stack pointer but check that it is safe to use. Specifically, we check that it is indeed an RWLX capability. Without this check, an adversary could potentially get control over our local stack frame during a subsequent callback by passing us a local rwx capability to a global data structure instead of a proper stack pointer and a global callback for our callback to invoke. If our local state contains no local capabilities, then, otherwise following our calling convention, the callback would not fail and the adversary could use a stored capability for the global data structure to access our local state. To prevent this from happening, we need to make sure the stack capability carries RWLX authority, since the system wide assumption then tells us that the adversary cannot have global capabilities to our local stack.

**Calling convention** With the security measures introduced and motivated, let us summarize our proposed calling convention:

*At program start-up* A local RWLX stack pointer resides in register  $r_{stk}$ . No global write-local capabilities.

*Before returning to the adversary* Clear non-return-value registers. Clear the part of the stack we had access to (not just the part we used).

*Before invoking the adversary* Push activation record to the stack. Create return pointer as local  $\epsilon$ -capability to the instructions in the record. Restrict the stack capability to the unused part and clear it. Clear non-argument registers.

*Before invoking an adversary callback* Make sure callback is global.

*When invoked by an adversary* Make sure received stack pointer has permission RWLX.

**Modularity** The calling convention ensures well-bracketed calls and local-state encapsulation for the caller, but not the callee. In the above presentation to make it easy to distinguish between the parties involved, we present the callee as some adversarial code that we do not trust. In reality, the callee could be well-behaved and wish to ensure well-bracketed calls and local-state encapsulation as well. The calling convention puts no restriction on the callee that the caller itself does not follow, so by following the calling convention, the callee can also obtain those guarantees. In other words, the calling convention is modular and scales to scenarios with multiple distrusting parties invoking each other.

## 4 LOGICAL RELATION

Suggestions for things to add heres

- (The lemma is there, consider pointer) Double monotonicity lemma (lemma 77 and 79) with a pointer back to local capability section of popl intro to this section.
- Reviewer B, ESOP suggests inferring general security properties from the LR. We had a section in the discussion that explained why this was not really possible. I have reintroduced it in the discussion.

Now that we have defined our calling convention, how can we be sure that it works? More concretely, suppose that we have a program that uses the convention in its interaction with untrusted adversary code. Can we formally prove the program's correctness if it relies on well-bracketed control flow and private state encapsulation for the interaction with the adversary? Clearly, such a proof should depend on a formal expression of the guarantees provided by the capability machine, including the specific guarantees for local capabilities.

In this section, we construct such a formalization. We make use of some well-studied and powerful (but non-trivial) machinery from the literature. Specifically, we employ a unary step-indexed Kripke logical relation with recursive worlds, and some additional special characteristics of our own. Step-indexing, Kripke logical relations and recursive worlds are techniques that may be familiar from lambda calculus settings, but it may not be clear to the reader how they apply in this more low-level assembly language. Therefore, in the next section, we do not immediately dive into the details, but we first try to provide some informal intuition about how all of this machinery comes into play in our setting.

Note: even though the calling convention is the main application in this paper, the logical relation we construct is very general and should be regarded as an independent contribution.

#### 4.1 Formalizing the guarantees of the capability machine

What differentiates a capability machine from a more standard assembly language is that we can bound the authority of an executing block of code, based solely on the capabilities it has access to. Specifically, it does not matter which instructions are actually executed, i.e., the bound also applies to untrusted adversary code that has not been inspected or modified in any way.

*Worlds.* But what does a “bound on the authority” of an executing block of code mean? In our setting, there are no externally observable side-effects and the only primitive authority that code may hold is authority over memory. As such, the authority bounds we consider are related to memory, but in a form that is more fine-grained than standard read/write authority: a piece of code's authority can be bounded by arbitrary memory invariants that it is required to respect. Specifically, we will define worlds  $W \in \text{World}$ , which describe a set of memory invariants, and our results will express authority bounds on code as safety with respect to such a world, i.e., the fact that the code respects the invariants registered in the world.

*Safe values.* So, let's say that we have a world  $W$  expressing that the memory must contain value 42 at address 0, may contain arbitrary values at addresses 50-60, a rw capability for address 0 at address 73, and an integer at address 100 that may only increase over time<sup>2</sup>. Our main theorem will state that if the current register file only contains safe words (numbers or capabilities which preserve the invariants in  $W$  under any interaction), then an execution will necessarily also preserve the memory invariants (irrespective of the instructions being executed).

To make this more precise, we need to define the set  $\mathcal{V}(W) \in \mathcal{P}(\text{Word})$  of words that are safe w.r.t.  $W$ . Essentially, the set should only include words that preserve  $W$ 's invariants under any interaction, but should otherwise be as liberal as possible. Numbers are clearly always safe, as they cannot be used to break invariants. Whether a capability is safe depends on the authority that it carries. In the above-described world, a read capability for address 0 is safe, as it can only be used to read the value 42, which is itself safe. However, a write capability for address 0 is not safe: it can be used to overwrite the memory at that address with a value other than 42, breaking the invariant for that address.

<sup>2</sup>Indeed, we will allow a notion of evolvable invariants, aka protocols, that can express such a temporal property.

*Step-indexing.* More generally, we want to define that a read capability for memory range  $[b, e]$  is safe, if the world guarantees that the words at those addresses are themselves safe. However, this definition is cyclic: what if the world guarantees that the memory at address  $a$  will contain a read capability for address  $a$ ? The definition then just says that a read capability for address  $a$  is safe if and only if the same read capability for address  $a$  is safe. This form of cyclic reasoning is related to similar challenges in languages with recursive types or higher-order ML-style references, and a standard solution is to use step-indexing [Appel and McAllester 2001]: essentially, the cycle is broken by defining safety up to a certain number of interaction steps. All words will be considered safe up to 0 steps (since if there is no interaction, nothing unsafe can happen), and, for example, a read capability will be safe up to  $n$  steps if the world guarantees that the words at the corresponding addresses are safe up to  $n - 1$  steps. We can then prove that the above read capability for address  $a$  is safe up to any number of steps.

*Future worlds.* So worlds are defined as a set of invariants on the memory, but what if we allocate fresh memory through malloc? We may want to establish new invariants for this freshly allocated memory and be sure that the adversary will also respect those (if we don't provide them with capabilities through which the new invariants can be broken). To accommodate this, we allow worlds to evolve, for example by adding additional invariants for freshly allocated memory. Formally, we define valid ways for a world  $W$  to evolve into a new world  $W'$  through a future-world relation  $W' \sqsupseteq W$  and we ensure that the set of safe words in world  $W$  must remain safe in any future world  $W'$ . Defining safety w.r.t. a notion of evolvable worlds makes our logical relation into a Kripke logical relation [Pitts and Stark 1998].

*Invariants and Recursive Worlds.* So worlds group a set of memory invariants, but how are they actually defined formally? We represent each invariant by a region  $r \in \text{Region}$ . We will see later that regions contain state machines to support a notion of evolvable invariants, but in every state, they also contain a predicate  $H$  that defines which memory segments are acceptable in the current state of the invariant. Unfortunately, it is not enough to just take  $H \in \mathcal{P}(\text{MemSeg})$ , because sometimes the invariant may itself be world-dependent. For example, we may want to express invariants like “the memory at address 50 contains a value that is safe in the current world”. As explained, worlds may evolve, and the set of safe values may grow in future worlds, and therefore we need to index  $H$  over worlds, i.e., take  $H \in \text{World} \rightarrow \mathcal{P}(\text{MemSeg})$ . We then end up with worlds containing regions with world-indexed predicates, i.e., the set of worlds must be recursively defined. We will see how such a recursive definition can be accommodated using techniques from the literature (essentially an advanced application of step-indexing).

*Local capabilities.* A final thing we provide informal intuition for is how our results take into account local capabilities and their special treatment by the hardware. Normally, when we invoke an untrusted piece of code and provide it with certain global capabilities, it may have stored those capabilities in memory and we will only be able to reinvoke the code if we can guarantee that those values are still valid. Formally, worlds represent the invariants that global capabilities' safety relies on and the reinvocation is only safe in future worlds, where the invariants are respected.

However, if we provide the adversary with local capabilities in that first invocation, then the situation is a bit different. The adversary has no way to store these local capabilities, so if we make sure that there are also no old local capabilities in the register file for the second invocation (including the capability being invoked), then the adversary cannot use them any more.<sup>3</sup> Therefore,

<sup>3</sup>We are ignoring write-local capabilities in this discussion. If the adversary does have access to write-local capabilities in the first and second invocation, then the memory they address must entirely be cleared before the second invocation in order for the reinvocation to remain safe.

we can allow the second invocation to happen in private future worlds ( $W' \sqsupseteq^{priv} W$ ), in which safe global capabilities remain safe, but local ones don't. This private future world relation is more liberal than the standard, public one ( $W' \sqsupseteq^{pub} W$ , in which all safe capabilities remain safe). Concretely, worlds may contain temporary regions, representing invariants that only local capabilities may rely on for their safety, and which may be revoked (disabled) in private future worlds.

Interestingly, this idea is a variant of a notion of public/private future worlds that has been previously used in the literature (see Section 9). However, temporary regions are new in our setting and there is an interesting interplay with the recursiveness of the worlds: for a temporary region, the predicate  $H \in \text{World} \rightarrow \mathcal{P}(\text{MemSeg})$  (which defines the safe memory segments in the current world) is only required to be monotone w.r.t. public future worlds (i.e. safe memory segments remain safe in public future worlds), but for permanent regions, it must be monotone w.r.t. private future worlds. As a consequence, the memory for a permanent region may not contain local capabilities (as their safety would be broken in private future worlds), which in turn implies that only local capabilities may have write-local permission (a general sanity requirement when using local capabilities<sup>4</sup><sup>5</sup>).

## 4.2 Worlds

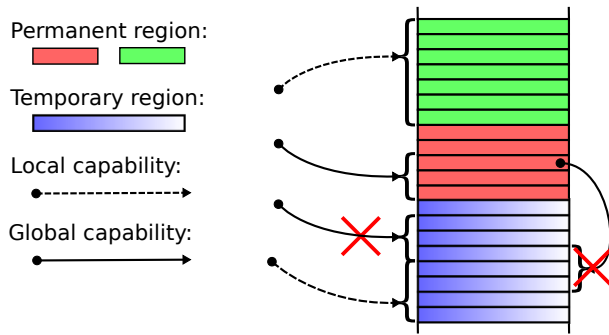


Fig. 6. The relation between local/global capabilities and temporary/permanent regions. The colored fields are regions governing parts of memory. Global capabilities cannot depend on temporary regions.

A world is a finite map from region names, modeled as natural numbers, to regions that each correspond to an invariant of part of the memory. We have three types of regions: permanent, temporary, and revoked. Each permanent and temporary region contains a state transition system, with public and private transitions, to describe how the invariants are allowed to change over time. In other words, they are protocols for the region's memory. Protocols imposed by permanent regions stay in place indefinitely. Any capability, local or global, can depend on these protocols. Protocols imposed by temporary regions can be revoked in private future worlds. Doing this may break the safety of local capabilities but not global ones. This means that local capabilities can safely depend on the protocols imposed by temporary regions, but global capabilities cannot, since a global capability may outlive a temporary region that is revoked. This is illustrated in Figure 6.

For technical reasons, we do not actually remove a revoked temporary region from the world, but we turn it into a special revoked region that exists for this purpose. Such a revoked region

<sup>4</sup>In fact, local capabilities become useless as soon as the adversary has access to a single global, write-local capability.

<sup>5</sup>For explanation purposes, this discussion ignores certain ways to allow for local capabilities in a permanent region, for example, by not requiring that they are valid or requiring that they are local versions of valid global capabilities.

contains no state transition system and puts no requirements on the memory. It simply serves as a mask for a revoked temporary region. Masking a region like this goes back to earlier work of Ahmed [2004] and was also used by Thamsborg and Birkedal [2011].

Regions are used to define safe memory segments, but this set may itself be world-dependent. In other words, our worlds are defined recursively. Recursive worlds are common in Kripke models and the following lemma uses the method of Birkedal and Bizjak [2014]; Birkedal et al. [2011] for constructing them. The formulation of the lemma is technical, so we recommend that non-expert readers ignore the technicalities and accept that there exists a set of worlds  $\text{Wor}$  and two relations  $\sqsubseteq^{\text{priv}}$  and  $\sqsubseteq^{\text{pub}}$  satisfying the (recursive) equations in the theorem (where the  $\blacktriangleright$  operator can be safely ignored).

**THEOREM 4.1.** *There exists a c.o.f.e. (complete ordered family of equivalences)  $\text{Wor}$  and preorders  $\sqsubseteq^{\text{priv}}$  and  $\sqsubseteq^{\text{pub}}$  such that  $(\text{Wor}, \sqsubseteq^{\text{priv}})$  and  $(\text{Wor}, \sqsubseteq^{\text{pub}})$  are preordered c.o.f.e.'s, and there exists an isomorphism  $\xi$  such that*

$$\xi : \text{Wor} \cong \blacktriangleright(\mathbb{N} \xrightarrow{\text{fin}} \text{Region})$$

$$\text{Region} = \{\text{revoked}\} \uplus$$

$$\{\text{temp}\} \times \text{State} \times \text{Rels} \times (\text{State} \rightarrow (\text{Wor} \xrightarrow[\sqsubseteq^{\text{pub}}]{\text{mon, ne}} \text{UPred}(\text{MemSeg}))) \uplus$$

$$\{\text{perm}\} \times \text{State} \times \text{Rels} \times (\text{State} \rightarrow (\text{Wor} \xrightarrow[\sqsubseteq^{\text{priv}}]{\text{mon, ne}} \text{UPred}(\text{MemSeg})))$$

$$\text{and for } W, W' \in \text{Wor.} \quad W' \sqsubseteq^{\text{priv}} W \Leftrightarrow \xi(W') \sqsubseteq^{\text{priv}} \xi(W)$$

$$W' \sqsubseteq^{\text{pub}} W \Leftrightarrow \xi(W') \sqsubseteq^{\text{pub}} \xi(W)$$

In the above theorem,  $\text{State} \times \text{Rels}$  corresponds to the aforementioned state transition system where  $\text{Rels}$  contains pairs of relations corresponding to the public and private transitions, and  $\text{State}$  is an unspecified set that we assume to contain at least the states we use in this paper. The last part of the temporary and permanent regions is a state interpretation function that determines what memory segments the region permits in each state of the state transition system. The different monotonicity requirements in the two interpretation functions reflects how permanent regions rely only on permanent protocols whereas temporary regions can rely on both temporary and permanent protocols.  $\text{UPred}(\text{MemSeg})$  is the set of step-indexed, downwards closed predicates on memory segments:  $\text{UPred}(\text{MemSeg}) = \{A \subseteq \mathbb{N} \times \text{MemSeg} \mid \forall(n, ms) \in A. \forall m \leq n. (m, ms) \in A\}$ .

With the recursive domain equation solved, we could take  $\text{Wor}$  as our notion of worlds, but it is technically more convenient to work with the following definition instead:

$$\text{World} = \mathbb{N} \xrightarrow{\text{fin}} \text{Region}$$

**4.2.1 Future Worlds.** The future world relations model how memory may evolve over time. The public future world relation  $W' \sqsubseteq^{\text{pub}} W$  requires that  $\text{dom}(W') \supseteq \text{dom}(W)$  and  $\forall r \in \text{dom}(W). W'(r) \sqsubseteq^{\text{pub}} W(r)$ . That is, in a public future world, new regions may have been allocated, and existing regions may have evolved according to the public future region relation (defined below). The private future world relation  $W' \sqsubseteq^{\text{priv}} W$  is defined similarly, using a private future region relation. The public future region relation is the simplest. It satisfies the following properties:

$$\frac{(s, s') \in \phi_{\text{pub}}}{(v, s', \phi_{\text{pub}}, \phi, H) \sqsubseteq^{\text{pub}} (v, s, \phi_{\text{pub}}, \phi, H)} \quad \frac{(\text{temp}, s, \phi_{\text{pub}}, \phi, H) \in \text{Region}}{(\text{temp}, s, \phi_{\text{pub}}, \phi, H) \sqsubseteq^{\text{pub}} \text{revoked}}$$

$$\frac{}{\text{revoked} \sqsubseteq^{\text{pub}} \text{revoked}}$$



Both temporary and permanent regions are only allowed to transition according to the public part of their transition system. Additionally, revoked regions must either remain revoked or be replaced by a temporary region. This means that the public future world relations allows us to reinstate a region that has been revoked earlier. The private future region relation satisfies:

$$\frac{(s, s') \in \phi}{(v, s', \phi_{pub}, \phi, H) \sqsupseteq^{priv} (v, s, \phi_{pub}, \phi, H)} \quad \frac{r \in \text{Region}}{r \sqsupseteq^{priv} (\text{temp}, s, \phi_{pub}, \phi, H)} \quad \frac{r \in \text{Region}}{r \sqsupseteq^{priv} \text{revoked}}$$

Here, revocation of temporary regions is allowed. In fact, temporary regions can be replaced by an arbitrary other region, not just the special revoked. Conversely, revoked regions may also be replaced by any other region. On the other hand, permanent regions cannot be masked away. They are only allowed to transition according to the private part of the transition system.

In the future world relations, we must remember all region names in order to keep them extensional which is why we need to use masks. It does, however, not matter what kind of region is used to mask as it is the region name we must retain. This is why we allow revoked regions to be replaced by other regions in the two future world relations. This may seem unnecessary as the future world relation allows new regions to be added, but in private future worlds, for technical reasons, we need to be able to reinstate old regions with the same region names as they had in a past world. Further, we do not want public future worlds to be able to prevent reinstatement in private future worlds. Private future worlds cannot revoke perm regions, so we cannot allow public future worlds to use perm regions as a mask for a revoked region.

Notice that the public future region relation is a subset of the private future region relation.

**4.2.2 World Satisfaction.** A memory satisfies a world, written  $ms :_n W$ , if it can be partitioned into disjoint parts such that each part is accepted by an active (permanent or temporary) region. Revoked regions are not taken into account as their memory protocols are no longer in effect.

$$ms :_n W \text{ iff } \begin{cases} \exists P : \text{active}(W) \rightarrow \text{MemSeg. } ms = \biguplus_{r \in \text{active}(W)} P(r) \text{ and} \\ \forall r \in \text{active}(W). \\ \exists H, s. W(r) = (\_, s, \_, \_, H) \text{ and } (n, P(r)) \in H(s)(\xi^{-1}(W)) \end{cases}$$

### 4.3 Logical Relation

The logical relation defines semantically when values, program counters, and configurations are capability safe. The definition is found in Figures 7 and 8 and we provide some explanations in the following paragraphs. For space reasons, we omit some definitions and explain them only verbally, but precise definitions can be found in the appendix.

First, the observation relation  $\mathcal{O}$  defines what configurations we consider safe. A configuration is safe with respect to a world, when the execution of said configuration does not break the memory protocols of the world. Roughly speaking, this means that when the execution of a configuration halts, then there is a private future world that the resulting memory satisfies. Notice that failing is considered safe behavior. In fact, the machine often resorts to failing when an unauthorized access is attempted, such as loading from a capability without read permission. This is similar to Devriese et al. [2016]’s logical relation for an untyped language, but unlike typical logical relations for typed languages, which require that programs do not fail.

The register-file relation  $\mathcal{R}$  defines safe register-files as those that contain safe words (i.e. words in  $\mathcal{V}$ ) in all registers but pc. The expression relation  $\mathcal{E}$  defines that a word is safe to use as a program counter if it can be plugged into a safe register file (i.e. a register file in  $\mathcal{R}$ ) and paired with a memory satisfying the world to become a safe configuration. Note that integers and non-executable capabilities (e.g. ro and  $\mathbb{E}$  capabilities) are considered safe program counters because

$$\begin{aligned}
O &: \text{World} \xrightarrow{ne} \text{UPred}(\text{Reg} \times \text{MemSeg}) \\
O(W) &\stackrel{\text{def}}{=} \left\{ (n, (\text{reg}, \text{ms})) \mid \begin{array}{l} \forall \text{ms}_f, \text{mem}', i \leq n. (\text{reg}, \text{ms} \uplus \text{ms}_f) \rightarrow_i (\text{halted}, \text{mem}') \Rightarrow \\ \exists W' \sqsupseteq^{\text{priv}} W, \text{ms}_r, \text{ms}'. \\ \text{mem}' = \text{ms}' \uplus \text{ms}_r \uplus \text{ms}_f \text{ and } \text{ms}' :_{n-i} W' \end{array} \right\} \\
\mathcal{R} &: \text{World} \xrightarrow[\sqsupseteq^{\text{pub}}]{\text{mon}, ne} \text{UPred}(\text{Reg}) \\
\mathcal{R}(W) &\stackrel{\text{def}}{=} \{(n, \text{reg}) \mid \forall r \in \text{RegName} \setminus \{\text{pc}\}. (n, \text{reg}(r)) \in \mathcal{V}(W)\} \\
\mathcal{E} &: \text{World} \xrightarrow{ne} \text{UPred}(\text{Word}) \\
\mathcal{E}(W) &\stackrel{\text{def}}{=} \left\{ (n, \text{pc}) \mid \begin{array}{l} \forall n' \leq n, (n', \text{reg}) \in \mathcal{R}(W), \text{ms} :_{n'} W. \\ (n', (\text{reg}[\text{pc} \mapsto \text{pc}], \text{ms})) \in O(W) \end{array} \right\} \\
\mathcal{V} &: \text{World} \xrightarrow[\sqsupseteq^{\text{pub}}]{\text{mon}, ne} \text{UPred}(\text{Word}) \\
\mathcal{V}(W) &\stackrel{\text{def}}{=} \{(n, i) \mid i \in \mathbb{Z}\} \cup \\
&\quad \{(n, ((o, g), b, e, a))\} \cup \\
&\quad \{(n, ((ro, g), b, e, a)) \mid (n, (b, e)) \in \text{readCond}(g)(W)\} \cup \\
&\quad \left\{ (n, ((rw, g), b, e, a)) \mid \begin{array}{l} (n, (b, e)) \in \text{readCond}(g)(W) \text{ and} \\ (n, (b, e)) \in \text{writeCond}(l^{nwl}, g)(W) \end{array} \right\} \cup \\
&\quad \left\{ (n, ((rwl, g), b, e, a)) \mid \begin{array}{l} (n, (b, e)) \in \text{readCond}(g)(W) \text{ and} \\ (n, (b, e)) \in \text{writeCond}(l^{pwl}, g)(W) \end{array} \right\} \cup \\
&\quad \left\{ (n, ((rx, g), b, e, a)) \mid \begin{array}{l} (n, (b, e)) \in \text{readCond}(g)(W) \text{ and} \\ (n, (\{rx\}, b, e)) \in \text{execCond}(g)(W) \end{array} \right\} \cup \\
&\quad \{(n, ((e, g), b, e, a)) \mid (n, (b, e, a)) \in \text{enterCond}(g)(W)\} \cup \\
&\quad \left\{ (n, ((rwx, g), b, e, a)) \mid \begin{array}{l} (n, (b, e)) \in \text{readCond}(g)(W) \text{ and} \\ (n, (b, e)) \in \text{writeCond}(l^{nwl}, g)(W) \text{ and} \\ (n, (\{rwx, rx\}, b, e)) \in \text{execCond}(g)(W) \end{array} \right\} \cup \\
&\quad \left\{ (n, ((rwlx, g), b, e, a)) \mid \begin{array}{l} (n, (b, e)) \in \text{readCond}(g)(W) \text{ and} \\ (n, (b, e)) \in \text{writeCond}(l^{pwl}, g)(W) \text{ and} \\ (n, (\{rwlx, rwx, rx\}, b, e)) \in \text{execCond}(g)(W) \end{array} \right\}
\end{aligned}$$

Fig. 7. The logical relation.

when they are plugged into a register file and paired with a memory, the execution will immediately fail, which is safe.

The value relation  $\mathcal{V}$  defines when words are safe. [Explain why we use UPred\(\) with reference to the introduction of this section.](#) We make the value relation as liberal as possible by considering what is the most we can allow an adversary to use a capability for without breaking the memory protocols. Non-capability data is always safe because it provides no authority. Capabilities give the authority to manipulate memory and potentially break memory protocols, so they need to satisfy certain conditions to be safe. In Figure 8, we define such a condition for each kind of permission a capability can have.

$$\begin{aligned}
\text{readCond}(g)(W) &= \left\{ (n, (b, e)) \mid \begin{array}{l} \exists r \in \text{localityReg}(g, W). \\ \exists [b', e'] \supseteq [b, e]. W(r) \stackrel{n}{\subseteq} \iota_{b', e'}^{pwl} \end{array} \right\} \\
\text{writeCond}(\iota, g)(W) &= \left\{ (n, (b, e)) \mid \begin{array}{l} \exists r \in \text{localityReg}(g, W). \\ W(r) \text{ is address-stratified and} \\ \exists [b', e'] \supseteq [b, e]. W(r) \stackrel{n-1}{\supseteq} \iota_{b', e'} \end{array} \right\} \\
\text{execCond}(g)(W) &= \left\{ (n, (P, b, e)) \mid \begin{array}{l} \forall n' < n, W' \sqsupseteq W, a \in [b', e'] \subseteq [b, e], \text{perm} \in P. \\ (n', ((\text{perm}, g), b', e', a)) \in \mathcal{E}(W') \end{array} \right\} \\
\text{enterCond}(g)(W) &= \left\{ (n, (b, e, a)) \mid \begin{array}{l} \forall n' < n. \forall W' \sqsupseteq W. \\ (n', ((\text{rx}, g), b, e, a)) \in \mathcal{E}(W') \end{array} \right\}
\end{aligned}$$

where  $g = \text{local} \Rightarrow \sqsupseteq = \sqsupseteq^{\text{pub}}$  and  $g = \text{global} \Rightarrow \sqsupseteq = \sqsupseteq^{\text{priv}}$

Fig. 8. Permission-based conditions

$$\begin{aligned}
H_A^{nwl}, H_A^{pwl} &: \text{State} \rightarrow (\text{Wor} \xrightarrow[\sqsupseteq^{\text{pub}}]{\text{mon, ne}} \text{UPred}(\text{MemSeg})) \\
\iota_A^{pwl}, \iota_A^{nwl} &: \text{Region} \\
\iota_A^{pwl} &\stackrel{\text{def}}{=} (\text{temp}, 1, =, =, H_A^{pwl}) \\
H_A^{pwl} \text{ s } \hat{W} &\stackrel{\text{def}}{=} \left\{ (n, ms) \mid \begin{array}{l} \text{dom}(ms) = A \wedge \\ \forall a \in A. (n-1, ms(a)) \in \mathcal{V}(\xi(\hat{W})) \end{array} \right\} \cup \{(0, ms)\} \\
\iota_A^{nwl} &\stackrel{\text{def}}{=} (\text{temp}, 1, =, =, H_A^{nwl}) \\
\iota_A^{nwl, p} &\stackrel{\text{def}}{=} (\text{perm}, 1, =, =, H_A^{nwl}) \\
H_A^{nwl} \text{ s } \hat{W} &\stackrel{\text{def}}{=} \left\{ (n, ms) \mid \begin{array}{l} \text{dom}(ms) = A \wedge \\ \forall a \in A. \\ ms(a) \text{ is non-local} \wedge \\ (n-1, ms(a)) \in \mathcal{V}(\xi(\hat{W})) \end{array} \right\} \cup \{(0, ms)\}
\end{aligned}$$

Fig. 9. The standard permit write-local and no write-local regions.

For capabilities with read permission, the *readCond* ensures that it can only be used to read safe words, i.e. words in the value relation. To guarantee this, we require that the addressed memory is governed by a region  $W(r)$  that imposes safety as a requirement on the values contained. This safety requirement is formulated in terms of a standard region  $\iota_{[b, e]}^{pwl}$ . It simply requires all the words in the range  $[b, e]$  to be safe, i.e. in the value relation. Requiring that  $W(r) \stackrel{n}{\subseteq} \iota_{[b, e]}^{pwl}$  means that  $W(r)$  must accept only safe values like  $\iota_{[b, e]}^{pwl}$ , but can be even more restrictive if desired. The read condition also takes into account the locality of the capability because, generally speaking, global capabilities should only depend on permanent regions. Concretely, we use the function *localityReg*( $g, W$ ), which projects out all active (non-revoked) regions when the locality  $g$  is local,

but only the permanent regions when  $g$  is global. The definition of the standard region  $\iota_{[b,e]}^{pwl}$  can be found in Figure 9; it makes use of the isomorphism from Theorem 4.1.

For a capability with write permission, *writeCond* must be satisfied for the capability's range of authority. An adversary can use such a capability to write any word they can get a hold of, and we can safely assume that they can only get a hold of safe words, so the region governing the relevant memory must allow any safe word to be written there. In order to make the logical relation as liberal as possible, we make this a lower bound of what the region may allow. For write capabilities, we also have to take into account the two flavours of write permissions: write and write-local. In the case of write-local capabilities, the region needs to allow (at least) any safe word to be written, but in the case of write capabilities, the capability cannot be used to write local capabilities, so the region only needs to allow safe non-local values. In the write condition, this is handled by parameterizing it with a region. For the write-local capabilities the write condition is applied with the standard region  $\iota_{[b,e]}^{pwl}$  that we described previously. For the write capabilities we use a different standard region  $\iota_{[b,e]}^{nwl}$  which requires that the words in  $[b, e]$  are non-local and safe. As before, we use *localityReg* to pick an appropriate region based on the capability's locality. Finally, there is a technical requirement that the region must be address-stratified. Intuitively, this means that if a region accepts two memory segments, then it must also accept every memory segment "in between", that is every memory segment where each address contains a value from one of the two accepted memory segments. An interesting property of the write condition is that they prohibit global write-local capabilities which, as discussed in Section 3, is necessary for any safe use of local capabilities. The standard region  $\iota_{[b,e]}^{nwl}$  is defined in Figure 9.

The conditions *enterCond* and *execCond* are very similar. Both require that the capability can be safely jumped to. However, executable capabilities can be updated to point anywhere in their range, so they must be safe as a program counter (in the  $\mathcal{E}$ -relation) no matter the current address. **The range of an executable capability can also be reduced, so they must also be safe as program counter no matter what their range of authority is reduced to.** In contrast, enter capabilities are opaque and can only be used to jump to the address they point to.

**This is why the *enterCond* depends on the current address of the capability, unlike for other types of capabilities.** They also change permission when jumped to, so we require them to be safe as a program counter after the permission is changed to RX. Because the capabilities are not necessarily invoked immediately, this must be true in any future world, but it depends on the capability's locality which future worlds we consider. If it is global, then we require safety as a program counter in private future worlds (where temporary regions may be revoked). For local capabilities, it suffices to be safe in public future worlds, where temporary regions are still present.

In the technical appendix, we prove that safety of all values is preserved in public future worlds, and that safety of global values is also preserved in private future worlds:

LEMMA 4.2 (DOUBLE MONOTONICITY OF VALUE RELATION).

- If  $W' \sqsupseteq^{pub} W$  and  $(n, w) \in \mathcal{V}(W)$ , then  $(n, w) \in \mathcal{V}(W')$ .
- If  $W' \sqsupseteq^{priv} W$  and  $(n, w) \in \mathcal{V}(W)$  and  $w = ((perm, global), b, e, a)$  (i.e.  $w$  is a global capability), then  $(n, w) \in \mathcal{V}(W')$ .

#### 4.4 Safety of the Capability Machine

With the logical relation defined, we can now state the fundamental theorem of our logical relation: a strong theorem that formalizes the guarantees offered by the capability machine. Essentially, it says a capability that only grants safe authority is capability safe as a program counter.

THEOREM 4.3 (FUNDAMENTAL THEOREM). *If one of the following holds:*

- $perm = RX$  and  $(n, (b, e)) \in readCond(g)(W)$
- $perm = RWX$  and  $(n, (b, e)) \in readCond(g)(W)$  and  
 $(n, (b, e)) \in writeCond(l^{nwl}, g)(W)$
- $perm = RWLX$  and  $(n, (b, e)) \in readCond(g)(W)$  and  
 $(n, (b, e)) \in writeCond(l^{pwl}, g)(W),$

then  $(n, ((perm, g), b, e, a)) \in \mathcal{E}(W)$

PROOF SKETCH. Induction over  $n$ . By definition of  $\mathcal{E}(W)$ , show  $(n', (reg[pc \mapsto ((perm, g), b, e, a)], ms)) \in \mathcal{O}(W)$  assuming  $n' \leq n$ ,  $(n', reg) \in \mathcal{R}(W)$ , and  $ms \vdash_{n'} W$ . By definition of  $\mathcal{O}$  let  $ms_f$ ,  $mem'$ , and  $i \leq n'$  be given and for  $\Phi = (reg[pc \mapsto ((perm, g), b, e, a)], ms \uplus ms_f)$  assume  $\Phi \rightarrow_i (halted, mem')$  and show there exists  $W' \sqsubseteq^{priv} W$  that part of  $mem'$  satisfies. First observe that  $i \neq 0$  as  $\Phi$  is a non-halted configuration, so  $\Phi$  takes at least one step, i.e.  $\Phi \rightarrow \Phi' \rightarrow_{i-1} (halted, mem')$ .

The rest of the proof considers the different ways the step  $\Phi \rightarrow \Phi'$  could have occurred, depending on the instruction being executed. For each of these cases, we argue (1) that  $\Phi'$  is consistent with the world in the sense that the register-file and memory still respect the world and (2) that the rest of the execution respects the world. Depending on where the pc in  $\Phi'$  comes from, the second result is proven in one of two ways. If the step to  $\Phi'$  was a jump, then the new pc is one of the safe values in  $\Phi$ 's registers and the value relation can be used to argue that those can be safely jumped to. On the other hand, if the pc was just incremented, then we can apply the induction hypothesis.

In order to argue (1), we consider what configurations the operational semantics allows us to get to from the initial state. If we consider the memory in  $\Phi'$ , then it either (a) remains unchanged or (b) one address has been updated and in that case, the register-file contains an appropriate capability for writing. The latter occurs when the instruction being executed is a store, and otherwise the former. In case (a), we can conclude that the memory still respects the world just by downwards closure of memory satisfaction. In case (b), we use an auxiliary lemma that uses the safety of the write capability used by the store instruction to show that the updated memory satisfies the world. To show that the updated register-file is safe, we consider the changes made to it by all instructions in separate lemmas and show that they all preserve safety of the register file.

The complete proof can be found in the technical appendix [Skorstengaard et al. 2018].  $\square$

The permission-based conditions of Theorem 4.3 make sure that the capability only provides safe authority in which case the capability must be in the  $\mathcal{E}$  relation, i.e. it can safely be used as a program counter in an otherwise safe register-file.

The Fundamental Theorem can be understood as a general expression of the guarantees offered by the capability machine, an instance of a general property called capability safety [Devriese et al. 2016; Maffeis et al. 2010]. To understand this, consider that the theorem says the capability  $((perm, g), b, e, a)$  is safe as a program counter, without any assumption about what instructions it actually points to (the only assumptions we have are about the read or write authority that it carries). As such, the theorem expresses the capability safety of the machine, which guarantees that any instruction is fine and will not be able to go beyond the authority of the values it has access to. We demonstrate this in Section 8 where Theorem 4.3 is used to reason about capabilities that point to arbitrary instructions. The relation between Theorem 4.3 and local-state encapsulation and control-flow correctness, will also be shown by example in Section 8 as the examples depend on these properties for correctness.

Suggestion: Add a small section about reasoning about programs. Add the “scall works” and “malloc works” lemmas and describe how they are used. This section could also point out the things one have to argue about that are hidden in similar high-level proofs (e.g., the return pointer valid when we return corresponding to restoring memory invariants of caller.).

## 5 MALLOC

In the examples we present in Section 8, we will assume the existence of a trusted malloc routine so that both the trusted code and the adversary can allocate new memory. Malloc is considered part of the trusted computing base, as mentioned in Section 2. This is unavoidable: if we cannot trust malloc, then we cannot use the memory it allocates as we have no idea whether it is aliased by some untrusted program.

Our semantic model is not specific to a particular implementation of malloc, so rather than providing the implementation, we provide a specification that malloc must satisfy. The specification expresses standard expectations of how malloc should behave, but making it realistic requires using some of the technical machinery from the logical relation. As such, this section is a bit technical and can safely be skipped on first read. The malloc specification is presented in Figure 10, and in the following, we will connect the description of the specification by referring to the item described.

We require a global capability for invoking malloc (because if the capability were local, then a program with access to malloc would have to give up this access when invoking untrusted code, 1). The capability is assumed to have enter permission (1), so that malloc can protect its internal state even when the capability is shared with untrusted code.

In addition to these syntactic requirements, we also specify standard expectations for how malloc behaves. Intuitively, we require that when malloc is invoked with a length argument, then it will return a capability for a fresh piece of memory of that size. It should be fresh in the sense that malloc has not already given out a capability for any part of that memory before and will not do so in the future. Also, when invoked with a nonsensical length argument such as a negative integer or a capability, malloc should simply fail. However, formulating these requirements is harder than one might expect. The problem is that a realistic implementation of malloc needs to rely on internal state that changes after every invocation, and relies on invariants on that state. We can only expect that malloc behaves according to its specification, if its internal state satisfies the current state of the invariants in an executing system. To express this, we will use the semantic model defined in Section 4.

To allow malloc implementations to rely on internal state and invariants for that state, we assume an initial region for malloc (2). The region is assumed to be permanent (since safety of the global malloc capability will depend on its presence, 2a). Furthermore, we want to express that malloc does not depend on any other memory than its own internal state. This is expressed by a restriction on the malloc region’s state interpretation function, which (as explained in Section 4.2) defines what memory segments it permits in a given world. We require that the accepted memory segments only depend on the presence of the malloc region itself, i.e. that in any world, the same memory segments are accepted if we remove all regions except the malloc region. This property should continue to hold throughout execution, so it must hold true for any private future region of the initial malloc region (2b).

The malloc specification then dictates what malloc should do when invoked in a memory with its internal state valid according to the malloc region (some future evolution of the initial malloc region). If malloc is invoked with an invalid length argument (that is, a negative integer or a capability), then we simply require malloc to fail (2d). This part of the specification does not actually rely on the malloc region: for simplicity, we assume malloc does not need its internal state to check the argument. When malloc is invoked with a valid length (2c), then it should return a fresh



memory segment of the correct length. This segment is required to come from the footprint of malloc, i.e. the memory owned by the malloc region before the call. After malloc returns, we require the malloc region to have evolved (according to the public future world relation) to a new state where the new footprint is disjoint from the allocated memory. This implies that future invocations of malloc can never return previously-allocated memory.

For convenience, we also require that malloc returns the non-argument registers of the register-file unchanged after the call. This allows the caller to keep private capabilities in the register file, without having to protect them by storing them in a private stack frame.

**Definition 5.1 (Malloc Specification).**  $c_{\text{malloc}}$  satisfies the specification for malloc iff the following conditions hold:

- (1)  $c_{\text{malloc}} = ((\mathbb{E}, \text{global}), \_, \_, \_)$
- (2) There exists a  $\iota_{\text{malloc},0}$  such that
  - (a)  $\iota_{\text{malloc},0}.v = \text{perm}$
  - (b) For all  $\iota' \sqsupseteq^{\text{priv}} \iota_{\text{malloc},0}, W, i$  with  $W(i) = \iota'$ , we have that
 
$$\iota'.H(\iota'.s)(\xi^{-1}(W)) = \iota'.H(\iota'.s)(\xi^{-1}([i \mapsto W(i)]))$$
  - (c) For all  $\Phi \in \text{ExecConf}, ms_{\text{footprint}}, ms_{\text{frame}} \in \text{MemSeg}, i, n, \text{size} \in \mathbb{N}, w_{\text{ret}} \in \text{Word}, \iota_{\text{malloc}} \sqsupseteq^{\text{priv}} \iota_{\text{malloc},0}$ , we have that
 

If  $\Phi.\text{mem} = ms_{\text{footprint}} \uplus ms_{\text{frame}} \wedge ms_{\text{footprint}} :_n [i \mapsto \iota_{\text{malloc}}] \wedge$   
 $\Phi.\text{reg}(r_1) = \text{size} \wedge \text{size} \geq 0 \wedge \Phi.\text{reg}(r_0) = w_{\text{ret}} \wedge$   
 $\Phi.\text{reg}(\text{pc}) = \text{updPcPerm}(c_{\text{malloc}})$

Then,

$$\exists \Phi' \in \text{ExecConf}. \exists ms'_{\text{footprint}}, ms_{\text{alloc}} \in \text{MemSeg}.$$

$$\exists j \in \mathbb{N}. j > 0 \wedge \exists b', e' \in \text{Addr}. \exists \iota'_{\text{malloc}} \in \text{Region}.$$

$$\Phi \rightarrow_j \Phi' \wedge$$

$$\Phi'.\text{mem} = ms'_{\text{footprint}} \uplus ms_{\text{alloc}} \uplus ms_{\text{frame}} \wedge$$

$$\iota'_{\text{malloc}} \sqsupseteq^{\text{pub}} \iota_{\text{malloc}} \wedge$$

$$ms'_{\text{footprint}} :_{n-j} [i \mapsto \iota'_{\text{malloc}}] \wedge$$

$$\text{dom}(ms_{\text{alloc}}) = [b', e'] \wedge \forall a \in [b', e']. ms_{\text{alloc}}(a) = 0 \wedge$$

$$\Phi'.\text{reg} = \Phi.\text{reg}[\text{pc} \mapsto \text{updPcPerm}(w_{\text{ret}})][r_1 \mapsto ((\text{RWX}, \text{global}), b', e', b')] \wedge$$

$$\text{size} - 1 = e' - b'$$
  - (d) For all  $\Phi \in \text{ExecConf}$ ,
 

If  $(\Phi.\text{reg}(r_1) \notin \mathbb{Z} \vee \Phi.\text{reg}(r_1) < 0) \wedge \Phi.\text{reg}(\text{pc}) = \text{updPcPerm}(c_{\text{malloc}})$

Then  $\exists j \in \mathbb{N}. \Phi \rightarrow_j \text{failed}$

Fig. 10. The specification of malloc.

As described previously, the specification of malloc ensures that malloc has no capabilities pointing out of malloc. It does, however, not say anything about capabilities that point in to malloc. We obviously cannot allow this if we want to trust malloc. We have chosen to keep the malloc specification simple and let this assumption be in the lemmas that use malloc. It is sufficient for these lemmas to require that there are no outside capabilities for malloc in the initial configuration as capabilities cannot appear out of thin air and the malloc specification makes sure that capabilities are not leaked.

Malloc should not just be available to trusted programs, but also to possibly malicious programs. This is safe as it follows from the specification that the malloc capability is always safe in a world with the malloc region:

LEMMA 5.2 (MALLOC IS SAFE TO PASS TO ADVERSARY). *For all  $c_{\text{malloc}}$  that satisfies the specification for malloc with region  $\iota_{\text{malloc},0}$ , if  $W(r) \sqsupseteq^{\text{priv}} \iota_{\text{malloc},0}$ , then  $(n, c_{\text{malloc}}) \in \mathcal{V}(W)$  for all  $n$ .*

The reason we allow trusted code and adversary to invoke malloc is just to make our work more realistic, but we are otherwise not interested in its details. As such, we do not give a malloc implementation. We are, however, confident that it is possible to make an implementation of malloc that satisfies the malloc specification in Definition 5.1. There are in fact two simplifications in our system that makes things easier: First, we do not consider deallocation of memory which means that the data structure malloc uses to keep track of free memory does not have to handle reclaimed memory. Second, the malloc specification does not permit malloc to run out of memory and thus refuse allocation. This is possible on our simple capability machine because it has an infinite address space. An initial capability with an infinite range of authority would of course need to be part of malloc, but it could also double as the data structure that keeps track of free memory.

## 6 REUSABLE MACRO INSTRUCTIONS

With the calling convention and logical relation defined, we would like to show its usefulness by proving the correctness of a series of examples that rely on well-bracketedness and local-state encapsulation and use the calling convention to enforce these properties. However, the programs that run on our capability machine are assembly programs and even program examples that would be small in a high-level language become big and unintelligible at this low level. Thus to make our program examples intelligible, we introduce a series of low-level abstractions in the form of macros. We define a number of reusable macros capturing the calling convention as well as other conveniences. The macros that utilise the stack assume that it is available in register  $r_{\text{stk}}$ . *Lau, suggestion: should we mention temporary registers? They are used in the macros, but it seems like a minor technical detail.*

The macro `scall  $r(\overline{r_{\text{args}}}, \overline{r_{\text{priv}}})$`  captures those parts of the calling convention related to actually transferring control to some adversarial code. Specifically, it pushes the contents of the private registers,  $\overline{r_{\text{priv}}}$ , to the stack. Then it pushes the “restoration” code to the stack. The restoration code will be executed as the first thing upon return, and restores the stack pointer and the old program counter. After the restoration code is pushed to the stack, a copy of the pc is pushed to the stack after it has been adjusted to point to the first instruction after the jump. From the stack pointer, a protected return pointer is created by adjusting it to point to the return code and encapsulating it by restricting its permission to  $\text{E}$ . Next, the stack pointer is restricted to the unused part, and the unused part of the stack is cleared (as discussed in Section 3). Finally, the non-argument registers are cleared and we jump to  $r$ . Upon return after the restoration code has been executed, the restoration code is popped from the stack and the private words we pushed to the stack before the call are popped to the private registers. The implementation of `scall` is presented in Figure 11, and the restoration code is presented in Figure 12. The implementation of `scall` uses some of the macros we present next.

The macro `mclear  $r$`  clears all the memory the capability in register  $r$  has authority over. It is used by `scall` to clear the unused part of the stack before control is transferred. It should also be used to clear the stack before returning to adversarial code. Similarly, the macro `rclear  $R$`  clears all the registers in the set  $R$ . It is also used by `scall`, and it should also be used before returning to adversarial code.

```

1079 // push private registers to the stack
1080   push  $r_{priv,1}$ 
1081   ...
1082   push  $r_{priv,n}$ 
1083 // push restoration code to the stack
1084   push  $encode(i_1)$ 
1085   push  $encode(i_2)$ 
1086   push  $encode(i_3)$ 
1087   push  $encode(i_4)$ 
1088 // push old pc to the stack
1089   move r_t1 pc
1090   lea r_t1 ofc
1091   push r_t1
1092 // push stack pointer to the stack
1093   push r_stk
1094 // set up protected return pointer
1095   move r_0 r_stk
1096   lea r_0 ofc_rec // -5 is the offset to the first instruction of
1097 the recovery code
1098   restrict r_0  $encodePermPair((local, e))$ 
1099 // restrict stack capability to unused part
1100   geta r_t1 r_stk
1101   plus r_t1 r_t1 1
1102   getb r_t2 r_stk
1103   subseg r_stk r_t1 r_t2
1104 // clear unused part of the stack
1105   mclear r_stk
1106 // clear non-argument registers
1107   rclear R
1108   jmp r
1109 after:
1110 // pop the restore code
1111   pop r_t1
1112   pop r_t1
1113   pop r_t1
1114   pop r_t1
1115 // pop the private state into appropriate registers
1116   pop  $r_{priv,1}$ 
1117   ...
1118   pop  $r_{priv,n}$ 
1119

```

Fig. 11. Implementation of  $scall\ r(r_{args,1}, \dots, r_{args,m}, r_{priv,1}, \dots, r_{priv,n})$ . The restoration code is given in Fig. 12. The variable  $ofc$  is the offset to the label after, and  $ofc_{rec} = -5$  which is the offset to the first instruction of the activation record. The set  $R = RegName \setminus \{pc, r\_stk, r\_0, r, r_{args,1}, \dots, r_{args,m}\}$ .

The macros `prepstack r` and `reqglob r` are the last macros related to the calling convention. The former, `prepstack` should be used when one receives a stack from an unknown source. The `prepstack` macro first ensures the stack has permission `read/write-local/execute`, and then it

```

1128  i1 = move r_t1 pc
1129  i2 = lea r_t1 ofc
1130  i3 = load r_stk r_t1
1131  i4 = pop pc
1132

```

Fig. 12. The restoration code used in `scall`. The variable `ofc` is 5 which is the offset to the address where the old stack pointer is stored on the stack.

adjusts the stack capability, so it follows the convention for the stack<sup>6</sup>. The other macro, `reqglob`, ensures that the capability in register  $r$  is global. This macro should be used on callbacks received from unknown sources to ensure that they cannot be derived from the stack pointer.

The remaining macros are not directly related to the calling convention, but they help making the examples in Section 8 intelligible. The macros `push r` and `pop r` add and remove elements from the stack. The macro `fetch r name` fetches the capability related to `name` from the linking table and stores it in register  $r$ . The macro `malloc r n` invokes `malloc` with size argument  $n$ . The `malloc` macro assumes that a capability for `malloc` resides in the linking table and is basically a fetch of the `malloc` capability followed by a setup of a return pointer. Finally, the macro `crtcls  $\overline{(x_i, r_i)}$  r` allocates a closure where  $r$  points to the closure's code and a new environment is allocated (using `malloc`) and the contents of  $\overline{r_i}$  is stored in the environment. In the code referred to by  $r$ , an implicit load from the environment happens when an instruction refers to  $x_i$ .

The Appendix contains the implementation of all the macros used in `scall`. The technical appendix [Skorstengaard et al. 2018] contains more detailed descriptions of all the macros as well as all implementations. We stress that the macros correspond to series of instructions as seen in Figure 11 and 12, and that they are introduced for intelligibility. The examples of Section 8, the program examples are stated using the macros, but the proofs work on the expanded macros.

## 7 REASONING ABOUT PROGRAMS ON A CAPABILITY MACHINE

There are many details to get right when programming in assembly. These details carry over to proofs about assembly programs, so many of the proofs in the next section about example programs are a bit cumbersome. It is especially annoying, when the same line of reasoning is applied in multiple places. To mitigate this, we have defined a number of lemmas that capture common reasoning patterns in these proofs. We present the most important such lemmas in this section, to give the reader an idea of how they work.

Our capability machine only allows the final memory of an execution to be observed, so naturally the correctness lemmas we prove in Section 8 are statements about the memory in the halted configuration. In order to prove something about some part of the final memory, we establish an invariant for it as a region in a world and show that the initial configuration is in the  $\mathcal{O}$ -relation w.r.t. that world. The  $\mathcal{O}$ -relation says that if the initial configuration halts successfully, then the memory in the halting state must still respect a private future world of the initial world. It is, however, bothersome and error prone to try to argue about the entire execution in one go. Instead we want to reason modularly, in the sense that we only want to reason about parts of the execution at a time. To this end, we prove an anti-reduction lemma that essentially says that if we can show that an initial configuration,  $\Phi$ , steps to a configuration  $\Phi'$  and that  $\Phi'$  is in the  $\mathcal{O}$ -relation, then  $\Phi$  is also in the observation relation.

<sup>6</sup>The stack capability should always point to the topmost word on the stack. A stack received from an unknown source can be treated as empty, so the stack capability should point just outside its range of authority.

LEMMA 7.1 (ANTI-REDUCTION FOR  $O$ ).

$$\begin{aligned}
 & \forall n, n', i, \text{reg}, \text{reg}', \text{ms}, \text{ms}', \text{ms}_r, W, W'. \\
 & n' \geq n - i \wedge W' \sqsubseteq^{\text{priv}} W \wedge \\
 & (\forall \text{ms}_f. (\text{reg}, \text{ms} \uplus \text{ms}_r \uplus \text{ms}_f) \rightarrow_i (\text{reg}', \text{ms}' \uplus \text{ms}_r \uplus \text{ms}_f)) \wedge \\
 & (n', (\text{reg}', \text{ms}')) \in O(W') \\
 & \Rightarrow (n, (\text{reg}, \text{ms} \uplus \text{ms}_r)) \in O(W)
 \end{aligned}$$

The anti-reduction works well when we reason about the execution of known code, because we can execute it. When we want to reason about unknown code, the anti-reduction lemma does not apply because we do not know what instructions are being executed. This is where the FTLR (Theorem 4.3) comes into play. As a reminder, the FTLR says that if a capability only has access to safe values with respect to a world, then it is safe to use it for execution in the same world. The unknown code we consider will be assumed to only have access to safe values, which typically means that we assume it has access to a linking table with a malloc capability and otherwise consists of instructions. This allows us to use the FTLR seeing as malloc is a safe value (cf. Lemma 5.2) and instructions are really just integers and they are always safe values. Here, it is important to remember that the FTLR gives us that the capability for the untrusted code is safe w.r.t. some world, so even though we described the usual assumption on the unknown code's memory, it is the semantic model of the memory, i.e. the world, that the FTLR considers. In order to argue that a specific configuration is safe, we need to argue that the configuration is safe with respect to the world. That is, we need to argue that the memory satisfies the world and the register-file is in the  $\mathcal{R}$ -relation. For instance, in the case where the untrusted code assumes control first, we will use the FTLR to show that the capability for the unknown code can be used for execution. The unknown code will get access to some known, trusted code through a capability in the initial register-file. Because the capability is in the register-file, we will have to argue that it is safe which we do using Lemma 7.1.

Another common pattern in proofs about programs on our capability machine concerns the use of scall. The following lemma captures the commonalities of reasoning about programs using scall. From another point of view, it can be seen as a specification for scall.

LEMMA 7.2 (scall WORKS). *If*

- (1)  $\text{ms} :_n \text{revokeTemp}(W)$
- (2)  $\text{dom}(\text{ms}_f) \cap (\text{dom}(\text{ms}_{\text{stk}} \uplus \text{ms}_{\text{unused}} \uplus \text{ms})) = \emptyset$
- (3)  $(\text{reg}, \text{ms})$  is looking at scall  $r(\overline{r_{\text{arg}}}, \overline{r_{\text{priv}}})$  followed by  $c_{\text{next}}$ <sup>7</sup>
- (4)  $\text{reg}$  points to stack with  $\text{ms}_{\text{stk}}$  used and  $\text{ms}_{\text{unused}}$  unused
- (5) Hyp-Callee *If*
  - $\text{dom}(\text{ms}_{\text{unused}}) = \text{dom}(\text{ms}_{\text{rec}} \uplus \text{ms}'_{\text{unused}})$ ,
  - $W' = \text{revokeTemp}(W)[i^{\text{sta}}(\text{temp}, \text{ms}_{\text{stk}} \uplus \text{ms}_{\text{rec}} \uplus \text{ms}_f), i^{\text{pwl}}(\text{dom}(\text{ms}'_{\text{unused}}))]$ ,
  - $\text{ms}'' :_{n-1} W'$
  - $\text{reg}'$  points to stack with  $\emptyset$  used and  $\text{ms}'_{\text{unused}}$  unused
  - $\text{reg}' = \text{reg}_0[\text{pc} \mapsto \text{updPcPerm}(\text{reg}(r)), \overline{r_{\text{arg}}} \mapsto \text{reg}(\overline{r_{\text{arg}}}), r_0 \mapsto c_{\text{ret}}, r_{\text{stk}} \mapsto c_{\text{stk}}, r \mapsto \text{reg}(r)]$
  - $(n-1, c_{\text{ret}}) \in \mathcal{V}(W')$
  - $(n-1, c_{\text{stk}}) \in \mathcal{V}(W')$

then we have that  $(n-1, (\text{reg}', \text{ms}'')) \in O(W')$
- (6) Hyp-Cont *If*

<sup>7</sup>Defined in the appendix.

- $n' \leq n - 2$
- $W'' \sqsupseteq^{pub} \text{revokeTemp}(W)$
- $ms'' :_{n'} \text{revokeTemp}(W'')$
- for all  $r$ , we have that:

$$\text{reg}'(r) \begin{cases} = c_{next} & \text{if } r = \text{pc} \\ = \text{reg}(r) & \text{if } r \in \overline{r_{priv}} \\ \in \mathcal{V}(\text{revokeTemp}(W'')) & \text{if } \text{reg}'(r) \text{ is a global capability and } r \notin \{\text{pc}, \overline{r_{priv}}, r_{stk}\} \end{cases}$$

- $\text{reg}'$  points to stack with  $ms_{stk}$  used and  $ms''_{unused}$  unused for some  $ms''_{unused}$
- then we have that  $(n', (\text{reg}', ms'' \uplus ms_f \uplus ms_{stk} \uplus ms''_{unused})) \in O(W'')$

Then

- $(n, (\text{reg}, ms \uplus ms_f \uplus ms_{stk} \uplus ms_{unused})) \in O(W)$

Roughly, the scall lemma states that an invocation of scall is safe if scall is executed in a reasonable state (1-3), the callee is safe to execute (5), and returning to the code after the scall is safe (6).

In the common case scall is used to invoke untrusted code, and we will only have very general assumptions about that code, in which case we will use the FTLR (Theorem 4.3) to argue Assumption 5, using what we know about the values it has access to (e.g. linking table, malloc capability etc.). By assumption in Hyp-Callee, the memory is safe, so it suffices to show that the register-file contains safe values, which amounts to showing that the arguments in the call are safe.

By using Lemma 7.2 to reason about scall, the proofs become agnostic to the actual implementation of scall. In other words, should we change the implementation of scall, then we just need to prove that Lemma 7.2 holds for the new implementation in order to gain that all our results still hold true.

Also often used is the malloc instruction, so we also prove Lemma 7.3 to help reason about it. The structure of the malloc lemma is close to that of the scall lemma. It also has requirements on the configuration just before malloc is executed (1-8) as well as requirements on the execution afterwards (9). It does not have any requirements on the callee as the malloc specification defines its behavior.

LEMMA 7.3 (malloc WORKS). *If*

- (1)  $(\text{reg}, ms)$  is looking at malloc  $r$   $k$  followed by  $c_{next}$
- (2)  $k \geq 0$
- (3)  $(\text{reg}, ms)$  links malloc as  $k$  to  $c_{malloc}$
- (4)  $c_{malloc}$  satisfies the malloc specification with  $\iota_{malloc,0}$
- (5)  $W \sqsupseteq^{priv} [i \mapsto \iota_{malloc,0}]$
- (6)  $ms :_n W$
- (7)  $ms = ms' \uplus ms_{footprint}$
- (8)  $ms_{footprint} :_n [i \mapsto W(i)]$
- (9) Hyp-Cont *If*
  - $n' \leq n - 1$
  - $\iota_{malloc} \sqsupseteq^{pub} W(i)$
  - $ms'_{footprint} \uplus ms' :_{n'} W[i \mapsto \iota_{malloc}]$



<pre> 1275      f1: push 1 1276          fetch <math>r_1</math> <i>adv</i> 1277          scall <math>r_1</math>([],[]) 1278          pop <math>r_1</math> 1279          assert <math>r_1</math> 1 1280          halt </pre>	<pre> f2: malloc <math>r_l</math> 1     store <math>r_l</math> 1     fetch <math>r_1</math> <i>adv</i>     call <math>r_1</math>([], [<math>r_l</math>])     assert <math>r_l</math> 1     halt </pre>
---	--

Fig. 13. Two example programs that rely on local-state encapsulation. f1 uses our stack-based calling convention. f2 does not rely on a stack.

$$\begin{aligned}
& \bullet \text{ } ms'_{\text{footprint}} : n' [i \mapsto \iota_{\text{malloc}}] \\
& \text{reg}'(r') = \begin{cases} c_{\text{next}} & r' = \text{pc} \\ ((RWX, \text{global}), b, e, a) & r' = r \\ \text{reg}(r) & r' \notin \text{RegName}_t \cup \{\text{pc}, r, r_1\} \end{cases} \\
& \bullet e - b = k - 1 \\
& \bullet \text{dom}(ms_{\text{alloc}}) = [b, e] \\
& \bullet \forall a \in [b, e]. ms_{\text{alloc}}(a) = 0 \\
& \text{Then we have } (n', (\text{reg}', ms' \uplus ms'_{\text{footprint}} \uplus ms_{\text{alloc}})) \in O(W[\iota_{\text{malloc}}])
\end{aligned}$$

Then

$$(n, (\text{reg}, ms)) \in O(W)$$

In the technical appendix [Skorstengaard et al. 2018], we also define a lemma for the macro used to create closures, `crtcls`.

## 8 EXAMPLES

**Suggestion: add a couple of short proof sketches.** In this section, we demonstrate how our formalization of capability safety allows us to prove local-state encapsulation and control-flow correctness properties for challenging program examples. The security measures of Section 3 are deployed to ensure these properties. Since we are dealing with assembly language, there are many details to the formal treatment, and therefore we necessarily omit some details in the lemma statements. The examples may look deceptively short, but it is because they use the macro instructions described in Section 6. The examples would be unintelligible without the macros, as each macro expands to multiple basic instructions. The interested reader can find all the technical details in the technical appendix [Skorstengaard et al. 2018].

### 8.1 Encapsulation of Local State

f1 and f2 in Figure 13 demonstrate the capability machine's encapsulation of local state. They are very similar: both store some local state, call an untrusted piece of code (*adv*), and then test whether the local state is unchanged. They differ in the way they do this. Program f1 uses our stack-based calling convention (captured by `scall`) to call the adversary, so it can use the available stack to store its local state. On the other hand, f2 uses `malloc` to allocate memory for its local state and uses a calling convention based on heap allocated activation records (described in the technical appendix) to invoke the adversarial code.

For both programs, we can prove that if they are linked with an adversary, *adv*, that is allowed to allocate memory but has no other capabilities, then the assertion will never fail during executing

(see Lemmas 8.1 and 8.2 below). The two examples also illustrate the versatility of the logical relation. The logical relation is not specific to any calling convention, so we can use it to reason about both programs, even though they use different calling conventions.

In order to formulate results about  $f_1$  and  $f_2$ , we need a way to observe whether the assertion fails. To this end, we assume they have access to a flag (an address in memory). If the assertion fails, then the flag is set to 1 and execution halts. The correctness lemma for  $f_1$  then states:

LEMMA 8.1. *Let*

$$\begin{aligned} c_{adv} &\stackrel{\text{def}}{=} ((E, \text{global}), \dots) & c_{stk} &\stackrel{\text{def}}{=} ((RWLX, \text{local}), \dots) \\ c_{f_1} &\stackrel{\text{def}}{=} ((RWX, \text{global}), \dots) & c_{link} &\stackrel{\text{def}}{=} ((RO, \text{global}), \dots) \\ c_{malloc} &\stackrel{\text{def}}{=} ((E, \text{global}), \dots) & reg &\in \text{Reg} \\ m &\stackrel{\text{def}}{=} ms_{f_1} \uplus ms_{flag} \uplus ms_{link} \uplus ms_{adv} \uplus ms_{malloc} \uplus ms_{stk} \uplus ms_{frame} \end{aligned}$$

where each of the capabilities have an appropriate range of authority and pointer<sup>8</sup>. Furthermore

- $c_{malloc}$  satisfies the specification for malloc with  $\iota_{malloc,0}$
- $ms_{malloc} :_n [0 \mapsto \iota_{malloc,0}]$
- $ms_{f_1}$  contains  $c_{link}$ ,  $c_{flag}$  and the code of  $f_1$
- $ms_{flag}(flag) = 0$
- $ms_{link}$  contains  $c_{adv}$  and  $c_{malloc}$
- $ms_{adv}$  contains  $c_{link}$  and otherwise only instructions.

If  $(reg[pc \mapsto c_{f_1}][r_{stk} \mapsto c_{stk}], m) \rightarrow^* (\text{halted}, m')$ , then  $m'(flag) = 0$

To prove Lemma 8.1, it suffices to show that the start configuration is safe (in the  $\mathcal{O}$  relation) for a world with a permanent region that requires the assertion flag to be 0. By Lemma 7.1, it suffices to show that the configuration is safe after some reduction steps. We then use the `scall` lemma (Lemma 7.2), by which it suffices to show that (1) the configuration that `scall` will jump to is safe and (2) that the configuration just after `scall` is done cleaning up is safe. We use the Fundamental Theorem to reason about the unknown adversarial code as described in Section 7, but notice that the adversary capability is an `enter` capability, which the Fundamental Theorem says nothing about. Luckily the `enter` capability becomes `rx` after the jump and then the Fundamental Theorem applies.

We have a similar lemma for  $f_2$ :

LEMMA 8.2. *Making similar assumptions about capabilities and linking as in Lemma 8.1 but assuming no stack pointer, if  $(reg[pc \mapsto c_{f_2}], m) \rightarrow^* (\text{halted}, m')$ , then  $m'(flag) = 0$ .*

## 8.2 Well-Bracketed Control-Flow

Using the stack-based calling convention of `scall`, we get well-bracketed control-flow. To illustrate this, we look at two example programs  $f_3$  and  $g_1$  in Figure 14.

In  $f_3$  there are two calls to an adversary and in order for the assertion in the middle to succeed, they need to be well-bracketed. If the adversary were able to store the return pointer from the first call and invoke it in the second call, then  $f_3$  would have 2 on top of its stack and the assertion would fail. However, the security measures in Section 3 prevent this attack: specifically, the return pointer is local, so it can only be stored on the stack, but the part of the stack that is accessible to the adversary is cleared before the second invocation. In fact, the following lemma shows that there are also no other attacks that can break well-bracketedness of this example, i.e. the assertion never fails. It is similar to the two previous lemmas:

<sup>8</sup>These assumptions are kept intentionally vague for brevity. Full statements are in the Appendix.

1373	g1: malloc $r_2$ 1	(continued from previous column)	f3: push 1
1374	store $r_2$ 0	store $x$ 0	fetch $r_1$ <i>adv</i>
1375	move pc $r_3$	scall $r_1$ ([], [ $r_0$ , $r_1$ , $r_{env}$ ])	scall $r_1$ ([], [ $r_1$ ])
1376	lea $r_3$ <i>offset</i>	store $x$ 1	pop $r_2$
1377	crtcls [( $x$ , $r_2$ )] $r_3$	scall $r_1$ ([], [ $r_0$ , $r_{env}$ ])	assert $r_2$ 1
1378	rclear RegName \ {pc, $r_0$ , $r_1$ }	load $r_1$ $x$	push 2
1379	jmp $r_0$	assert $r_1$ 1	scall $r_1$ ([], [])
1380	f4: reqglob $r_1$	mclean $r_{stk}$	halt
1381	prepstk $r_{stk}$	rclear RegName \ { $r_0$ , pc}	
1382	(continues in next column)	jmp $r_0$	

Fig. 14. Two programs that rely on well-bracketedness of scalls to function correctly. *offset* is the offset to f4.

LEMMA 8.3. *Making similar assumptions about capabilities and linking as in Lemma 8.1 if ( $\text{reg}[pc \mapsto c_{f3}][r_{stk} \mapsto c_{stk}], m$ )  $\rightarrow^*$  (halted,  $m'$ ), then  $m'(\text{flag}) = 0$ .*

The final example, g1 with f4, is a faithful translation of a tricky example known from the literature (known as the awkward example) [Dreyer et al. 2012; Pitts and Stark 1998]. It consists of two parts, g1 and f4. g1 is a closure generator that generates closures with one variable  $x$  set to 0 in its environment and f4 as the program (note we can omit some calling convention security measures because the stack is not used in the closure generator). f4 expects one argument, a callback. It sets  $x$  to 0 and calls the callback. When it returns, it sets  $x$  to 1 and calls the callback a second time. When it returns again, it asserts  $x$  is 1 and returns. This example is more complicated than the previous ones because it involves a closure invoked by the adversary and an adversary callback invoked by us. *Suggestion: Maybe give some more intuition as to why it is difficult.* As explained in Section 3, this means that we need to check (1) that the stack pointer that the closure receives from the adversary has write-local permission and (2) that the adversary callback is global.

To illustrate how subtle this program is, consider how an adversary could try to make the assertion fail. In the second callback an adversary can get to the first callback by invoking the closure one more time. If there were any way for the adversary to transfer the return pointer from the point where it reinvokes the closure to where the closure reinvokes the callback, then the assertion could be made to fail. Similarly, if there were any way for the adversary to store a stack pointer or trick the trusted code into preserving it across an invocation, the assertion can likely be made to fail too. However, our calling convention prevents any of this from happening, as we prove in the following lemma.

LEMMA 8.4. *Let*

$$c_{adv} \stackrel{\text{def}}{=} ((RWX, \text{global}), \dots) \quad c_{g1} \stackrel{\text{def}}{=} ((E, \text{global}), \dots)$$

*and otherwise make assumptions about capabilities and linking similar to Lemma 8.1. Then if ( $\text{reg}_0[pc \mapsto c_{adv}][r_{stk} \mapsto c_{stk}][r_1 \mapsto c_{g1}], m$ )  $\rightarrow^*$  (halted,  $m'$ ), then  $m'(\text{flag}) = 0$ .*

As explained in Section 3, the macro-instruction reqglob  $r_1$  checks that the callback is global, essentially to make sure it is not allocated on the stack where it might contain old stack pointers or return pointers. Otherwise, the encapsulation of our local stack frame could be broken. In the proof of Lemma 8.4, this requirement shows up because we invoke the callback in a world that is only a private future world of the one where we received the callback, precisely because we have invalidated the adversary's local state (particularly their old stack and return capabilities). The callback is still valid in this private future world, but only because we know that it is global.

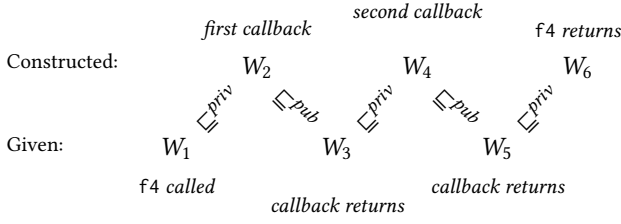


Fig. 15. Illustration of the worlds in the proof of Lemma 8.4. In the proof, the top row of worlds are constructed by us, while the bottom row of worlds are given.  $W_6$  is constructed such that  $W_6 \sqsupseteq^{\text{priv}} W_5$  and  $W_6 \sqsupseteq^{\text{pub}} W_1$ .

**Suggestion: Add an actual proof sketch** In Lemma 8.4 the order of control has been inverted compared to the previous lemmas. In this lemma, the adversary assumes control first with a capability for the closure creator  $g1$ . Consequently, we need to check that all arguments are safe to use and that we clean up before returning in the end. The inversion of control poses an interesting challenge when it comes to reasoning about the adversary’s local state during the execution of  $f4$  and the callbacks where the adversary should not rely on the local state from before the call of  $f4$ . This is easily done by revoking all the temporary regions of the world given at the start of  $f4$ . However, when  $f4$  returns, the adversary is again allowed to rely on its old local state so we need to guarantee that the local state is unchanged. This is important because the return pointer that  $f4$  receives may be local, and the adversary is allowed to allocate the activation record on the stack (just like we do) so they can store and recover their old stack pointer after  $f4$  returns.

In Figure 15, we illustrate how we accommodate this in the proof of Lemma 8.4 by constructing appropriate worlds for all the situations where control is passed to the adversary. The potentially local return pointer received by  $f4$  from the adversary is safe in  $W_1$ , so it can only be used in public future worlds. Let us see how  $W_6$  is constructed as a public future world of  $W_1$ , as well as a private future world of  $W_5$ . The worlds  $W_2$  and  $W_4$  are constructed by revoking all temporary regions and adding a  $\iota^{\text{pwl}}$  region for the stack we pass away and static temporary regions (a region that accepts exactly one memory) for our local stack as well as the adversary’s local stack. The given worlds  $W_3$  and  $W_5$  are public future worlds in which the adversary chooses to return to us. None of these worlds can have masked the temporary regions of  $W_1$  with permanent regions.  $W_6$  is constructed by revoking the temporary regions of  $W_5$  and reinstating the temporary regions of  $W_1$ . This makes  $W_6$  a public future world of  $W_1$ . It is therefore safe to use the return pointer to return to the adversary, since we have restored validity of any local state they might have stashed away.

## 9 DISCUSSION

Reviewer C, esop, would have liked a discussion of (a) what could be done on a machine with a smaller set of capabilities, e.g., no local capability (we could probably not have done anything as the stack pointer could be stored on the heap and reused at unintended times.) and (b) what could be done with a larger/stronger set of available capabilities (we could include a discussion of linear capabilities). Reviewer C, esop, asks whether our attacker model is reasonable (see tex comment). Maybe we should include a short discussion of the attacker model. Should we include a discussion of how this “scales” to other things like a multi-core setting or if we needed tail calls? **Calling convention**

**Formulating control flow correctness** While we claim that our calling convention enforces control-flow correctness, we do not prove a general theorem that shows this, because it is not clear what such a theorem should look like. Formulations in terms of a control-flow graph, like the one by

Abadi et al. [2005], do not take into account temporal properties, like the well-bracketedness that Example `g1` relies on. In fact, our examples show that our logical relation imply a stronger form of control-flow correctness than such formulations, although this is not made very explicit. As future work, we consider looking at a more explicit and useful way to formalize control-flow correctness. The idea would be to define a variant of our capability machine with call and return instructions and well-bracketed control flow built-in to the operational semantics, and then prove that compiling such code to our machine using our calling convention is fully abstract [Abadi 1998]. Maybe elaborate on why full-abstractness and not weaker properties like “robust safety”.

Performance and the requirement for stack clearing The additional security measures of the calling convention described in Section 3 impose an overhead compared to a calling convention without security guarantees. However, most of our security measures require only a few atomic checks or register clearings on boundary crossings between trusted code and adversary, which should produce an acceptable performance overhead. The only exception are the requirements for stack clearing that we have in two situations: when returning to the adversary and when invoking an adversary callback. As we have explained, we need to clear all of the stack that we are not using ourselves, not just the part that we have actually used. In other words, on every boundary cross between trusted code and adversary code, a potentially large region of memory must be cleared. We believe this is actually a common requirement for typical usage scenarios of local capabilities and capability machines like CHERI should consider to provide special support for this requirement, in the form of a highly-optimized instruction for erasing a large block of memory. Nevertheless, from a discussion with the designers of the CHERI capability machine, we gather that it is not immediately clear whether and how such a primitive could be implemented efficiently in the CHERI context. Maybe refer to some of the papers reviewer A, esop mentions for stack clearing. See comment in tex file.

The need for stack clearing in our calling convention is an instance of a general caveat when using CHERI’s local capabilities as a restricted form of capability revocation. Consider how our use of local capabilities can be interpreted as temporarily delegating the stack and return capability to callees and then revoking the granted authority after the callee returns. From this perspective, local capabilities are a general feature enabling this temporary delegation of authority for the duration of an invocation and this is also how they are described by the authors [Watson et al. 2015]. However, our requirement for stack clearing on boundary crossings is also general. Revoking authority that was granted temporarily using local capabilities, requires clearing all memory for which the invokee had write-local authority (or at least erasing all local capabilities from that memory). Without micro-architectural support for efficiently clearing large ranges of memory, local capabilities can only be used for revocation in scenarios where the duration of a revocation is unimportant or the adversary only has write-local access to small amounts of memory. In future work, we are planning to investigate an alternative notion of linear capabilities that offer a limited form of revocation without such limitations.

Note, by the way, that CheriBSD’s use of local capabilities in `ccall` does not actually involve a form of revocation. CheriBSD’s model involves a trusted stack manager that gives every compartment access to its own private stack using a local, write-local capability [Watson et al. 2015]. The locality of the stack capability allows the trusted stack manager to prevent compartments from leaking their stack pointer in a boundary crossing, but those capabilities are never actually revoked. In fact, a compartment can easily store away such local capabilities in its own private stack and recover them there during future invocations.

Modularity It is important that our calling convention is modular, i.e. we do not assume that our code is specially privileged w.r.t. the adversary, and they can apply the same measures to protect themselves from us as we do to protect ourselves from them. More concretely, the requirements we

have on callbacks and return pointers received from the adversary are also satisfied by callbacks and return pointers that we pass to them. For example, our return pointers are local capabilities because they must point to memory where we can store the old stack pointer, but the adversary's return pointers are also allowed to be local. Adversary callbacks are required to be global but the callbacks we construct are allocated on the heap and also global.

Arguments and local capabilities Local capabilities are a central part of the calling convention as they are used to construct stack and return pointers. The use of local capabilities for the calling convention unfortunately limits the extent to which local capabilities can be used for other things. Say we are using the calling convention and receive a local capability other than the stack and return pointer, then we need to be careful if we want to use it because it may be an alias to the stack pointer. That is, if we first push something to the stack and then write to the local capability, then we may be (tricked into) overwriting our own local state. The logical relation helps by telling us what we need to ascertain or check in such scenarios to guarantee safety and preserve our invariants, but such checks may be costly and it is not clear to us whether there are practical scenarios where this might be realistic.

We also need to be careful when we receive a capability from an adversary that we want to pass on to a different (instance of the) adversary. It turns out that the logical relation again tells us when this is safe. Namely, the logical relation says that we can only pass on arguments that are safe in the world we invoke the adversary in. For instance, when we receive a stack pointer from an adversary, then we may at some point want to pass on part of this stack pointer to, say, a callback. In order to do so, we need to make sure the stack pointer is safe which means that, if we have revoked temporary invariants, the stack must not directly or indirectly allow access to local values that we cannot guarantee safety of. When received from an adversary, we have to consider the contents of the stack unsafe, so before we pass it on, we have to clear it, or perform a dynamic safety analysis of the stack contents and anything it points to. Clearing everything is not always desirable and a dynamic safety analysis is hard to get right and potentially expensive.

In summary, the use of local capabilities for other things than stack and return pointers is likely only possible in very specific scenarios when using our calling convention. While this is unfortunate, it is not unheard of that processors have built-in constructs that are exclusively used for handling control flow, such as, for example, the call and return instructions that exist in some instruction sets.

We could also include a brief discussion of concurrency and the fact that this probably wouldn't scale to it (but well-bracketedness is not the same in that case anyway.) Reviewer C, *esop*, wants to know what happens when there are several stack blocks (see *tex* for comment). We could here emphasize that while we talk about a single stack there could be multiple local read/write-local/execute-capabilities. If we get the stack from an untrusted source, which is the case in the awkward example, then we will use anything that looks like the stack as the stack. However, due to the operational semantics for local capabilities all other "stack blocks" would have to be stored on the stack or in the register file. As the calling convention dictates that we clear these, the only place the alternative stack blocks will not be cleared is on a part of the stack that we do not have access to. The return capability would basically have to encapsulate this part of memory (perhaps indirectly) which means that from the moment we have been called, we will make sure there is only one stack available until our program returns from the initial call to it. Single stack A single stack is a good choice for the simple capability machine presented here, because it works well with higher-order functions. An alternative to a single stack would be to have a separate stack per component. The trouble with this approach is that, with multiple stacks and local stack pointers, it is not clear how components would retrieve their stack pointer upon invocation without compromising safety. A safe approach could be to have stack pointers stored by a central, trusted stack management



component, but it is not clear how that could scale to large numbers of separate components. Handling large numbers of components is a requirement if we want to use capability machines to enforce encapsulation of, for example, every object in an object-oriented program or every closure in a functional program.

### Reasoning about capability machine programs

Semantic, but not syntactic properties The logical relation defined in Section 4 allows us to reason about capability machine programs. A limitation w.r.t. previous work is that the logical relation is tailored exclusively towards semantic properties, not syntactic ones.

Imagine, for example, that we invoke a block of adversary code in such a way that it only ever receives capabilities within a specific range of memory. After the code returns, we may try to prove that any capabilities passed back to us in the registers are still confined to that range of memory. The property of falling in a certain address range is syntactic in the sense that it talks about the specific implementation of a higher-order value rather than its behavior, like the invariants that are required/preserved when we use it.

Such syntactic properties are hard to prove in our system. For the example cited, it would be easy to conclude that the returned values are in the value relation (see Figure 7). This gives us a lot of semantic information, like conditions under which they are safe to use and invariants that will be preserved when we do, but it does not tell us much about the address of the capability. As a very concrete example, capabilities with permission `o` are always in the value relation, irrespective of their address. Semantically, this makes perfect sense, since they are always safe since they cannot be used for anything anyway. However, it also means we do not get syntactic information about them.

For our purposes, this restriction is unproblematic, since we are only interested in proving semantic properties (e.g., an assertion will never fail). However, in other situations, we may be interested in proving more syntactic properties like the ones that are often considered in object capability literature: confinement, no authority amplification etc. Although such properties are more restrictive and less easy to use for reasoning, Devriese et al. [2016] have demonstrated how a logical relation like ours can be adapted to also support them, by quantifying the logical relation over a custom interpretation of effectful computations and the type of references. We expect their solution can be readily adapted to our setting, modulo some details (like the fact that we do not just have read-write capabilities, but also others).

### Logical relation

Single orthogonal closure The definitions of  $\mathcal{E}$  and  $\mathcal{V}$  in Figure 7 apply a single orthogonal closure, a new variant of an existing pattern called biorthogonality. Biorthogonality is a pattern for defining logical relations [Krivine 1994; Pitts and Stark 1998] in terms of an observation relation of safe configurations (like we do). The idea is to define safe evaluation contexts as the set of contexts that produce safe observations when plugging safe values and define safe terms as the set of terms that can be plugged into safe evaluation contexts to produce safe observations. This is an alternative to more direct definitions where safe terms are defined as terms that evaluate to safe values. An advantage of biorthogonality is that it scales better to languages with control effects like call/cc. Our definitions can be seen as a variant of biorthogonality, where we take only a single orthogonal closure: we do not define safe evaluation contexts but immediately define safe terms as those that produce safe observations when plugged with safe values. This is natural because we model arbitrary assembly code that does not necessarily respect a particular calling convention: return

pointers are in principle values like all others and there is no reason to treat them specially in the logical relation.

Interestingly, Hur and Dreyer [2011] also use a step-indexed, Kripke logical relation for an assembly language (for reasoning about correct compilation from ML to assembly), but because they only model non-adversarial code that treats return pointers according to a particular calling convention, they can use standard biorthogonality rather than a single orthogonal closure like us.

Public/private future worlds A novel aspect of our logical relation is how we model the temporary, revokable nature of local capabilities using public/private future worlds. The main insight is that this special nature generalizes that of the syntactically-enforced unstorable status of evaluation contexts in lambda calculi without control effects (of which well-bracketed control flow is a consequence). To reason about code that relies on this (particularly, the original awkward example), Dreyer et al. [2012] (DNB) formally capture the special status of evaluation contexts using Kripke worlds with public and private future world relations. Essentially, they allow relatedness of evaluation contexts to be monotone with respect to a weaker future world relation (public) than relatedness of values, formalizing the idea that it is safe to make temporary internal state modifications (private world transitions, which invalidate the continuation, but not other values) while an expression is performing internal steps, as long as the code returns to a stable state (i.e. transitions to a public future world of the original) before returning. We generalize this idea to reason about local capabilities: validity of local capabilities is allowed to be monotone with respect to a weaker future-world relation than other values, which we can exploit to distinguish between state changes that are always safe (public future worlds) and changes that are only valid if we clear all local capabilities (private future worlds). Our future world relations are similar to DNB's (for example, our proof of the awkward example uses exactly the same state transition system), but they turn up in an entirely different place in the logical relation: rather than using public future worlds for the special syntactic category of evaluation contexts, they are used in the value relation depending on the locality of the capability at hand. Additionally, our worlds are a bit more complex because, to allow local memory capabilities and write-local capabilities, they can contain (revokable) temporary regions that are only monotonous w.r.t. public future worlds, while DNB's worlds are entirely permanent.

Local capabilities in high-level languages We point out that local capabilities are quite similar to a feature proposed for the high-level language Scala: Osvald et al. [2016]'s second-class or local values. They are a kind of values that can be provided to other code for immediate use without allowing them to be stored in a closure or reference for later use. We believe reasoning about such values will require techniques similar to what we provide for local capabilities.

Reviewer C, popl, would like to know how local capabilities relate to borrowing (see tex comment).

## 10 RELATED WORK

Finally, we summarize how our work relates to previous work. We do not repeat the work we discussed in Section 9.

Capability machines originate with Dennis and Van Horn [1966] and we refer to Levy [1984] and Watson et al. [2015] for an overview of previous work. The capability machine formalized in Section 2 is a simple but representative model, modeled mainly after the M-Machine [Carter et al. 1994] (the enter pointers resemble the M-Machine's) and CHERI [Watson et al. 2015; Woodruff et al. 2014] (the memory and local capabilities resemble CHERI's). The latter is a recent and relatively mature capability machine, which combines capabilities with a virtual memory approach, in the interest of backwards compatibility and gradual adoption. As discussed, our local capabilities can cross module boundaries, contrary to what is enforced by CHERI's default CCall implementation.

Plenty of other papersLS: If there are plenty, then I guess we should cite more than one? enforce well-bracketed control flow at a low level, but most are restricted to preventing particular types of attacks and enforce only partial correctness of control flow. This includes particularly the line of work on control-flow integrity [Abadi et al. 2005]. Those use a quite different attacker model than us: they assume an attacker that is not able to execute code, but can overwrite arbitrary data at any time during execution (to model buffer overflows). By checking the address of every indirect jump and using memory access control to prevent overwriting code, this work enforces what they call control-flow integrity, formalized as the property that every jump will follow a legal path in the control-flow graph. As discussed in Section 9, such a property ignores temporal properties and seems hard to use for reasoning.

More closely related to our work are papers that use a trusted stack manager and some form of memory isolation to enforce control-flow correctness as part of a secure compilation result [Juglaret et al. 2016; Patrignani et al. 2016]. Our work differs from theirs in that we use a different form of low-level security primitive (a capability machine with local capabilities rather than a machine with a primitive notion of compartments) and we do not use a trusted stack manager, but a decentralized calling convention based on local capabilities. Also, both prove a secure compilation result from a high-level language, which clearly implies a general form of control-flow correctness, while we define a logical relation that can be used to reason about specific programs that rely on well-bracketed control flow.

LS: Reviewer C, esop, would like a better comparison with Devriese et al. [2016] (but also admits that they did not follow the below discussion).

Our logical relation is a unary, step-indexed Kripke logical relation with recursive worlds [Ahmed 2004; Appel and McAllester 2001; Birkedal et al. 2011; Pitts and Stark 1998], closely related to the one used by Devriese et al. [2016] to formulate capability safety in a high-level JavaScript-like lambda calculus. Our Fundamental Theorem is similar to theirs and expresses capability safety of the capability machine. Because we are not interested in externally observable side-effects (like console output or memory access traces), we do not require their notion of effect parametricity. Our logical relation uses several ideas from previous work, like Kripke worlds with regions containing state transition systems [Ahmed et al. 2009], public/private future worlds [Dreyer et al. 2012] (see Section 9 for a discussion), and biorthogonality [Benton and Hur 2009; Hur and Dreyer 2011; Pitts and Stark 1998].

Swasey et al. [2017] have recently developed a logic, OCPL, for verification of object capability patterns. The logic is based on Iris [Jung et al. 2016, 2015; Krebbers et al. 2017a], a state of the art higher-order concurrent separation logic and is formalized in Coq, building on the Iris Proof Mode for Coq [Krebbers et al. 2017b]. OCPL gives a more abstract and modular way of proving capability safety for a lambda-calculus (with concurrency) compared to the earlier work by Devriese et al. [2016].

In the future we would also like to investigate a new program logic for reasoning about capability safety for our capability machine model. In fact, we think the lemmas in Section 7 are suggestive of the style of results that could be written in such a logic. We think Iris would also be a natural starting point for such an endeavour, since Iris is really a framework, which can be instantiated to different programming languages. OCPL was able to leverage existing Iris specifications for the high-level language; for our capability machine model, however, it would be necessary to devise new kinds of specifications for our low-level programs with unstructured control-flow. It is likely that we could get inspiration from earlier work on logics for assembly programming languages, such as XCAP [Ni and Shao 2006]. We could further add that this is not the level of abstraction that one would like to reason. One would like to reason about the programs we actually write - that is programs written in high-level languages. If we had a fully-abstract compiler from a nice

high-level language to this low-level machine, then we could reason about the program in the high-level language and then because the compiler is fully-abstract, we would retain all the nice properties from the high-level language. How does this work fit in to this picture? The calling convention ensures two properties that we expect our nice high-level languages to have, namely well-bracketedness and local-state encapsulation. If we want to do fully-abstract compilation, then we need some way to guarantee this after compilation. Further, logical relations are often used in full-abstractness proofs. This work expands our knowledge about logical relations should be defined for low-languages and in particular capability machines. (maybe also mention that enforcement necessary for fully-abstract compilation. Could be capability, but it could also be something else - we just need something which we do not have now.)

El-Korashy also defined a formal model of a capability machine, namely CHERI, and uses it to prove a compartmentalization result [El-Korashy 2016] (not implying control-flow correctness). He also adapts control-flow integrity (see above) to the machine and shows soundness, seemingly without relying on capabilities.

## REFERENCES

- Martín Abadi. 1998. Protection in Programming-Language Translations: Mobile Object Systems. In *European Conference on Object-Oriented Programming (Lecture Notes in Computer Science)*. Springer Berlin Heidelberg, 291–291. [https://doi.org/10.1007/3-540-49255-0\\_70](https://doi.org/10.1007/3-540-49255-0_70)
- Martín Abadi, Mihai Budiu, Úlfar Erlingsson, and Jay Ligatti. 2005. Control-Flow Integrity. In *Conference on Computer and Communications Security*. ACM, 340–353. <https://doi.org/10.1145/1102120.1102165>
- Amal Ahmed, Derek Dreyer, and Andreas Rossberg. 2009. State-Dependent Representation Independence. In *POPL*. ACM, 340–353.
- Amal Jamil Ahmed. 2004. *Semantics of types for mutable state*. Ph.D. Dissertation. Princeton University.
- P. America and J. J. M. M. Rutten. 1989. Solving Reflexive Domain Equations in a Category of Complete Metric Spaces. *J. Comput. Syst. Sci.* 39, 3 (1989), 343–375.
- Andrew W. Appel and David McAllester. 2001. An Indexed Model of Recursive Types for Foundational Proof-carrying Code. *ACM Trans. Program. Lang. Syst.* 23, 5 (Sept. 2001), 657–683. <https://doi.org/10.1145/504709.504712>
- Nick Benton and Chung-Kil Hur. 2009. Biorthogonality, Step-Indexing and Compiler Correctness. In *International Conference on Functional Programming*. ACM, 97–108. <https://doi.org/10.1145/1596550.1596567>
- Lars Birkedal and AleÅa Bizjak. 2014. A Taste of Categorical Logic – ÅT Tutorial Notes. <http://cs.au.dk/~birke/modules/tutorial/categorical-logic-tutorial-notes.pdf>. (2014).
- Lars Birkedal, Bernhard Reus, Jan Schwinghammer, Kristian Støvring, Jacob Thamsborg, and Hongseok Yang. 2011. Step-indexed Kripke Models over Recursive Worlds. In *POPL*. ACM, 119–132. <https://doi.org/10.1145/1926385.1926401>
- L. Birkedal, K. Støvring, and J. Thamsborg. 2010. The category-theoretic solution of recursive metric-space equations. *Theor. Comput. Sci.* 411, 47 (2010), 4102–4122.
- Nicholas P. Carter, Stephen W. Keckler, and William J. Dally. 1994. Hardware Support for Fast Capability-based Addressing. In *Architectural Support for Programming Languages and Operating Systems*. ACM, 319–327. <https://doi.org/10.1145/195473.195579>
- Jack B. Dennis and Earl C. Van Horn. 1966. Programming Semantics for Multiprogrammed Computations. *Commun. ACM* 9, 3 (March 1966), 143–155. <https://doi.org/10.1145/365230.365252>
- Dominique Devriese, Lars Birkedal, and Frank Piessens. 2016. Reasoning about Object Capabilities using Logical Relations and Effect Parametricity. In *IEEE European Symposium on Security and Privacy*. IEEE.
- Derek Dreyer, Georg Neis, and Lars Birkedal. 2012. The Impact of Higher-Order State and Control Effects on Local Relational Reasoning. *J. Funct. Program.* 22, 4–5 (2012), 477–528.
- Akram El-Korashy. 2016. *A Formal Model for Capability Machines: An Illustrative Case Study towards Secure Compilation to CHERI*. Master’s thesis. Saarland University.
- S. Forrest, A. Somayaji, and D. H. Ackley. 1997. Building Diverse Computer Systems. In *Hot Topics in Operating Systems*. 67–72. <https://doi.org/10.1109/HOTOS.1997.595185>
- Chung-Kil Hur and Derek Dreyer. 2011. A Kripke Logical Relation Between ML and Assembly. In *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 133–146. <https://doi.org/10.1145/1926385.1926402>
- Y. Juglaret, C. Hritcu, A. A. D. Amorim, B. Eng, and B. C. Pierce. 2016. Beyond Good and Evil: Formalizing the Security Guarantees of Compartmentalizing Compilation. In *Computer Security Foundations Symposium (CSF)*. 45–60. <https://doi.org/10.1109/CSF.2016.11>

- Ralf Jung, Robbert Krebbers, Lars Birkedal, and Derek Dreyer. 2016. Higher-order ghost state. In *ICFP*. 256–269.
- Ralf Jung, David Swasey, Filip Sieczkowski, Kasper Svendsen, Aaron Turon, Lars Birkedal, and Derek Dreyer. 2015. Iris: Monoids and Invariants as an Orthogonal Basis for Concurrent Reasoning. In *POPL*. 637–650.
- Robbert Krebbers, Ralf Jung, AleÅa Bizjak, Jacques-Henri Jourdan, Derek Dreyer, and Lars Birkedal. 2017a. The essence of higher-order concurrent separation logic. In *European Symposium on Programming (ESOP)*.
- Robbert Krebbers, Amin Timany, and Lars Birkedal. 2017b. Interactive Proofs in Higher-Order Concurrent Separation Logic. In *POPL*.
- Jean-Louis Krivine. 1994. Classical Logic, Storage Operators and Second-Order Lambda-Calculus. *Annals of Pure and Applied Logic* 68, 1 (June 1994), 53–78. [https://doi.org/10.1016/0168-0072\(94\)90047-7](https://doi.org/10.1016/0168-0072(94)90047-7)
- Henry M Levy. 1984. *Capability-based computer systems*. Vol. 12. Digital Press Bedford.
- Tim Lindholm, Frank Yellin, Gilad Bracha, and Alex Buckley. 2014. *The Java Virtual Machine Specification*. Pearson Education.
- S. Maffeis, J.C. Mitchell, and A. Taly. 2010. Object Capabilities and Isolation of Untrusted Web Applications. In *S&P*. IEEE, 125–140. <https://doi.org/10.1109/SP.2010.16>
- Greg Morrisett, David Walker, Karl Crary, and Neal Glew. 1999. From System F to Typed Assembly Language. *ACM Trans. Program. Lang. Syst.* 21, 3 (May 1999), 527–568. <https://doi.org/10.1145/319301.319345>
- Z. Ni and Z. Shao. 2006. Certified Assembly Programming with Embedded Code Pointers. In *POPL*.
- Leo Osvald, Gr  gory Essertel, Xilun Wu, Lilliam I. Gonz  lez Alay  n, and Tiark Rompf. 2016. Gentrification Gone Too Far? Affordable 2Nd-Class Values for Fun and (Co-)Effect. In *Object-Oriented Programming, Systems, Languages, and Applications*. ACM, 234–251. <https://doi.org/10.1145/2983990.2984009>
- M. Patrignani, D. Devriese, and F. Piessens. 2016. On Modular and Fully-Abstract Compilation. In *Computer Security Foundations Symposium (CSF)*. 17–30. <https://doi.org/10.1109/CSF.2016.9>
- A. M. Pitts and I. D. B. Stark. 1998. Operational Reasoning for Functions with Local State. In *Higher Order Operational Techniques in Semantics*, Andrew D. Gordon and Andrew M. Pitts (Eds.). Cambridge University Press, New York, NY, USA, 227–274.
- Jonathan S. Shapiro, Jonathan M. Smith, and David J. Farber. 1999. EROS: A Fast Capability System. In *Symposium on Operating Systems Principles (SOSP ’99)*. ACM, 170–185. <https://doi.org/10.1145/319151.319163>
- Lau Skorstengaard, Dominique Devriese, and Lars Birkedal. 2018. *Reasoning About a Machine with Local Capabilities: Provably Safe Stack*. Technical Report. Dept. of Computer Science, Aarhus University. <https://cs.au.dk/~birke/papers/local-capabilities-conf-tr.pdf>. Available online: <https://cs.au.dk/~birke/papers/local-capabilities-conf-tr.pdf>.
- D. Swasey, D. Garg, and D. Dreyer. 2017. Robust and Compositionl Verification of Object Capability Patterns. (2017). <http://www.mpi-sws.org/~dreyer/papers/ocpl/paper.pdf> Submitted for publication.
- Jacob Thamsborg and Lars Birkedal. 2011. A Kripke Logical Relation for Effect-based Program Transformations. In *ICFP*. ACM, 445–456. <https://doi.org/10.1145/2034773.2034831>
- Robert Wahbe, Steven Lucco, Thomas E. Anderson, and Susan L. Graham. 1993. Efficient Software-Based Fault Isolation. In *Symposium on Operating Systems Principles*. ACM, 203–216. <https://doi.org/10.1145/168619.168635>
- R. N. M. Watson, J. Woodruff, P. G. Neumann, S. W. Moore, J. Anderson, D. Chisnall, N. Dave, B. Davis, K. Gudka, B. Laurie, S. J. Murdoch, R. Norton, M. Roe, S. Son, and M. Vadera. 2015. CHERI: A Hybrid Capability-System Architecture for Scalable Software Compartmentalization. In *IEEE Symposium on Security and Privacy*. 20–37. <https://doi.org/10.1109/SP.2015.9>
- Jonathan Woodruff, Robert N.M. Watson, David Chisnall, Simon W. Moore, Jonathan Anderson, Brooks Davis, Ben Laurie, Peter G. Neumann, Robert Norton, and Michael Roe. 2014. The CHERI Capability Model: Revisiting RISC in an Age of Risk. In *International Symposium on Computer Architecture*. IEEE Press, 457–468.

## A APPENDIX

In this appendix, we give more precise formulations of lemmas that were mentioned in the paper, and the most important supporting definitions and lemmas. The goal is to provide details that can help to understand in more detail what we discuss in the paper. Full details and proofs are not given here, but for those we refer to the technical appendix [Skorstengaard et al. 2018].

### A.1 Logical relation

#### *n*-subset simulation

$$\frac{(s, \phi_{pub}, \phi) = (s', \phi'_{pub}, \phi') \quad \forall \hat{W}. H \ s \ \hat{W} \stackrel{n}{\subseteq} H' \ s' \ \hat{W}}{(v, s, \phi_{pub}, \phi, H) \stackrel{n}{\subseteq} (v', s', \phi'_{pub}, \phi', H')}$$



## Erasure

$$\lfloor W \rfloor_S \stackrel{\text{def}}{=} \lambda r. \begin{cases} W(r) & W(r).v \in S \\ \perp & \text{otherwise} \end{cases}$$

## Active region projection

$$\text{active} : \text{World} \rightarrow 2^{\text{RegionName}}$$

$$\text{active}(W) \stackrel{\text{def}}{=} \text{dom}(\lfloor W \rfloor_{\{\text{perm}, \text{temp}\}})$$

## Projection of regions based on locality

$$\text{localityReg}(g, W) \stackrel{\text{def}}{=} \begin{cases} \text{dom}(\lfloor W \rfloor_{\{\text{perm}, \text{temp}\}}) & \text{if } g = \text{local} \\ \text{dom}(\lfloor W \rfloor_{\{\text{perm}\}}) & \text{if } g = \text{global} \end{cases}$$

## Address stratification

$\iota = (v, s, \phi_{\text{pub}}, \phi, H)$  is address-stratified iff

$$\forall s', \hat{W}, n, ms, ms'.$$

$$(n, ms), (n, ms') \in H \ s' \ \hat{W} \Rightarrow$$

$$\text{dom}(ms) = \text{dom}(ms') \wedge$$

$$\forall a \in \text{dom}(ms). (n, ms[a \mapsto ms'(a)]) \in H \ s' \ \hat{W}$$

## A.2 Complete ordered family of equivalences (c.o.f.e)

This is an excerpt from [Birkedal and Bizjak \[2014\]](#) about c.o.f.e.'s.

*Definition A.1 (o.f.e.).* An ordered family of equivalence (o.f.e) is a set and a family of equivalences  $(X, (\stackrel{n}{=})_{n=0}^{\infty})$  that satisfy the following properties:

- $\stackrel{0}{=}$  is the total relation on  $X$
- $\forall n. \forall x, y \in X. x \stackrel{n+1}{=} y \Rightarrow x \stackrel{n}{=} y$
- $\forall x, y \in X. (\forall n. x \stackrel{n}{=} y) \Rightarrow x = y$

We say that an o.f.e.  $(X, (\stackrel{n}{=})_{n=0}^{\infty})$  is inhabited if there exists an element  $x \in X$ .

If you are familiar with metric spaces observe that o.f.e.'s are but a different presentation of bisected 1-bounded ultrametric spaces.

*Definition A.2 (Cauchy sequences and limits).* Let  $(X, (\stackrel{n}{=})_{n=0}^{\infty})$  be an o.f.e. and  $\{x_n\}_{n=0}^{\infty}$  be a sequence of elements of  $X$ . Then  $\{x_n\}_{n=0}^{\infty}$  is a Cauchy sequence if

$$\forall k \in \mathbb{N}, \exists j \in \mathbb{N}, \forall n \geq j, x_j \stackrel{n}{=} kx_n$$

or in words, the elements of the chain get arbitrarily close.

An element  $x \in X$  is the limit of the sequence  $\{x_n\}_{n=0}^{\infty}$  if

$$\forall k \in \mathbb{N}, \exists j \in \mathbb{N}, \forall n \geq j, x \stackrel{n}{=} kx_n.$$

A sequence may or may not have a limit. If it has we say that the sequence converges. The limit is necessarily unique in this case and we write  $\lim_{n \rightarrow \infty} x_n$  for it.

*Definition A.3 (c.o.f.e.).* A complete ordered family of equivalences (c.o.f.e) is an o.f.e  $(X, (\stackrel{n}{=})_{n=0}^{\infty})$  where all Cauchy sequences have a limit.



*Definition A.4.* Let  $\left(X, \left(\overset{n}{\equiv} n_X\right)_{n=0}^{\infty}\right)$  and  $\left(Y, \left(\overset{n}{\equiv} n_Y\right)_{n=0}^{\infty}\right)$  be two ordered families of equivalences and  $f$  a function from the set  $X$  to the set  $Y$ . The function  $f$  is

- non-expansive if for any  $x, x' \in X$ , and any  $n \in \mathbb{N}$ ,

$$x \overset{n}{\equiv} n_X x' \implies f(x) \overset{n}{\equiv} n_Y f(x')$$

- contractive if for any  $x, x' \in X$ , and any  $n \in \mathbb{N}$ ,

$$x \overset{n}{\equiv} n_X x' \implies f(x) \overset{n}{\equiv} n + 1_Y f(x')$$

**THEOREM A.5 (BANACH'S FIXED POINT THEOREM).** Let  $\left(X, \left(\overset{n}{\equiv} n\right)_{n=0}^{\infty}\right)$  be a inhabited c.o.f.e. and  $f : X \rightarrow X$  a contractive function. Then  $f$  has a unique fixed point.

*Definition A.6 (The category  $\mathcal{U}$ ).* The category  $\mathcal{U}$  of complete ordered families of equivalences has as objects complete ordered families of equivalences and as morphisms non-expansive functions.

From now on, we often use the underlying set  $X$  to denote a (complete) o.f.e.  $\left(X, \left(\overset{n}{\equiv} n_X\right)_{n=0}^{\infty}\right)$ , leaving the family of equivalence relations implicit.

*Definition A.7.* The functor  $\blacktriangleright$  is a functor on  $\mathcal{U}$  defined as

$$\begin{aligned} \blacktriangleright \left(X, \left(\overset{n}{\equiv} n\right)_{n=0}^{\infty}\right) &= \left(X, \left(\overset{n}{\equiv} \right)_{n=0}^{\infty}\right) \\ \blacktriangleright (f) &= f \end{aligned}$$

where  $\overset{0}{\equiv}$  is the total relation and  $x \overset{n+1}{\equiv} x'$  iff  $x \overset{n}{\equiv} nx'$

*Definition A.8.* A functor  $F : \mathcal{U}^{\text{op}} \times \mathcal{U} \rightarrow \mathcal{U}$  is locally non-expansive if for all objects  $X, X', Y$ , and  $Y'$  in  $\mathcal{U}$  and  $f, f' \in \text{hom } \mathcal{U} X X'$  and  $g, g' \in \text{hom } \mathcal{U} Y' Y$  we have

$$f \overset{n}{\equiv} f' \wedge g \overset{n}{\equiv} g' \implies F(f, g) \overset{n}{\equiv} F(f', g').$$

It is locally contractive if the stronger implication

$$f \overset{n}{\equiv} f' \wedge g \overset{n}{\equiv} g' \implies F(f, g) \overset{n+1}{\equiv} F(f', g').$$

holds. Note that the equalities are equalities on function spaces.

**PROPOSITION A.9.** If  $F$  is a locally non-expansive functor then  $\blacktriangleright \circ F$  and  $F \circ (\blacktriangleright^{\text{op}} \times \blacktriangleright)$  are locally contractive. Here, the functor  $F \circ (\blacktriangleright^{\text{op}} \times \blacktriangleright)$  works as

$$(F \circ (\blacktriangleright^{\text{op}} \times \blacktriangleright))(X, Y) = F(\blacktriangleright^{\text{op}}(X), \blacktriangleright(Y))$$

on objects and analogously on morphisms and  $\blacktriangleright^{\text{op}} : \mathcal{U}^{\text{op}} \rightarrow \mathcal{U}^{\text{op}}$  is just  $\blacktriangleright$  working on  $\mathcal{U}^{\text{op}}$  (i.e., its definition is the same).

*Definition A.10.* A fixed point of a locally contractive functor  $F$  is an object  $X \in \mathcal{U}$ , such that  $F(X, X) \cong X$ .

The following is America and Rutten's fixed point theorem [America and Rutten 1989].

**THEOREM A.11.** Every locally contractive functor  $F$  such that  $F(1, 1)$  is inhabited has a unique fixed point. The fixed point is unique among inhabited c.o.f.e.'s. If in addition  $F(\emptyset, \emptyset)$  is inhabited then the fixed point of  $F$  is unique.

In Birkedal et al. [2010] one can find a category-theoretic generalization, which shows how to obtain fixed points of locally contractive functors on categories enriched in  $\mathcal{U}$ , in particular on the category of preordered c.o.f.e.'s.

### A.3 Load instruction sufficiency lemma

LEMMA A.12 (CONDITIONS FOR STORE INSTRUCTION ARE SUFFICIENT). *If*

- $ms = ms' \uplus ms_f$
- $ms' :_n W$
- $((perm, g), b, e, a) = c$
- $(n, c) \in \mathcal{V}(W)$
- $writeAllowed(perm)$
- $withinBounds(c)$
- $(n, w) \in \mathcal{V}(W)$
- *if*  $w = ((\_, local), \_, \_, \_)$ , *then*  $perm \in \{RWLX, RWL\}$

*then*  $a \in \text{dom}(ms')$  (i.e.  $ms[a \mapsto w] = ms'[a \mapsto w] \uplus ms_f$ ) *and*  $ms'[a \mapsto w] :_n W$

### A.4 Macros

Implementation of the macros used in `scall`. Implementations of the macros not presented here can be found in the technical appendix [Skorstengaard et al. 2018].

`push r`

```
    lea r_stk 1
    store r_stk r
```

`pop r`

```
    load r r_stk
    minus r_t1 0 1
    lea r_stk r_t1
```

`rclear  $r_1, \dots, r_n$`

```
move r_1 0
move r_2 0
...
move r_n 0
```

`mclean r`

```
    move r_t r
    getb r_t1 r_t
    geta r_t2 r_t
    minus r_t2 r_t1 r_t2
    lea r_t r_t2
    gete r_t2
    minus r_t1 r_t2 r_t1
    plus r_t1 r_t1 1
    move r_t2 pc
    lea r_t2 ofc_end
    move r_t3 pc
    lea r_t3 ofc_iter
```

`iter:`

```
    jnz r_t2 r_t1
    store r_t 0
    lea r_t 1
    plus r_t1 r_t1 1
    jmp r_t3
```

`end:`

```

1961      move r_t 0
1962      move r_t1 0
1963      move r_t2 0
1964      move r_t3 0

```

Where  $ofc_{end}$  and  $ofc_{iter}$  are the offsets to the label end and iter, respectively.

## A.5 Example correctness lemmas

LEMMA A.13 (CORRECTNESS LEMMA FOR  $f1$ , COPY OF LEMMA 8.1).

For all  $n \in \mathbb{N}$  let

$$\begin{aligned}
 c_{adv} &\stackrel{\text{def}}{=} ((E, \text{global}), b_{adv}, e_{adv}, b_{adv} + \text{offsetLinkFlag}) \\
 c_{f1} &\stackrel{\text{def}}{=} ((RWX, \text{global}), f1 - \text{offsetLinkFlag}, 1f, f1) \\
 c_{malloc} &\stackrel{\text{def}}{=} ((E, \text{global}), b_{malloc}, e_{malloc}, b_{malloc} + \text{offsetLinkFlag}) \\
 m &\stackrel{\text{def}}{=} ms_{f1} \uplus ms_{flag} \uplus ms_{link} \uplus ms_{adv} \uplus ms_{malloc} \uplus ms_{frame}
 \end{aligned}$$

and

- $c_{malloc}$  satisfies the specification for malloc and  $\iota_{malloc,0}$  is the region from the specification.

where

$$\begin{aligned}
 \text{dom}(ms_{f1}) &= [f1 - \text{offsetLinkFlag}, 1f] \\
 \text{dom}(ms_{flag}) &= [flag, flag] \\
 \text{dom}(ms_{link}) &= [link, link + 1] \\
 \text{dom}(ms_{adv}) &= [b_{adv}, e_{adv}] \\
 ms_{malloc} :_n [0 \mapsto \iota_{malloc,0}]
 \end{aligned}$$

and

- $ms_{f1}(f1 - \text{offsetLinkFlag}) = ((RO, \text{global}), link, link + 1, link)$ ,  $ms_{f1}(f1 - \text{offsetLinkFlag} + 1) = ((RW, \text{global}), flag, flag, flag)$ , the rest of  $ms_{f1}$  contains the code of  $f1$ .
- $ms_{flag} = [flag \mapsto 0]$
- $ms_{link} = [link \mapsto c_{malloc}, link + 1 \mapsto c_{adv}]$
- $ms_{adv}$  contains a global read-only capability for  $ms_{link}$  on its first address. The remaining cells of the memory segment only contain instructions.

if

$$(\text{reg}[pc \mapsto c_{f1}], m) \rightarrow_n (\text{halted}, m'),$$

then

$$m'(flag) = 0$$

LEMMA A.14 (CORRECTNESS LEMMA FOR f2, DETAILED VERSION OF LEMMA 8.2). *let*

$$c_{adv} \stackrel{\text{def}}{=} ((E, \text{global}), b_{adv}, e_{adv}, b_{adv} + \text{offsetLinkFlag})$$

$$c_{f2} \stackrel{\text{def}}{=} ((RWX, \text{global}), f2 - \text{offsetLinkFlag}, 2f, f2)$$

$$c_{malloc} \stackrel{\text{def}}{=} ((E, \text{global}), b_{malloc}, e_{malloc}, b_{malloc} + \text{offsetLinkFlag})$$

$$c_{stk} \stackrel{\text{def}}{=} ((RWLX, \text{local}), b_{stk}, e_{stk}, b_{stk} - 1)$$

$$c_{link} \stackrel{\text{def}}{=} ((RO, \text{global}), link, link + 1, link)$$

$$reg \in \text{Reg}$$

$$m \stackrel{\text{def}}{=} ms_{f2} \uplus ms_{flag} \uplus ms_{link} \uplus ms_{adv} \uplus ms_{malloc} \uplus ms_{stk} \uplus ms_{frame}$$

and

- $c_{malloc}$  satisfies the specification for malloc and  $\iota_{malloc,0}$  is the region from the specification.

where

$$\text{dom}(ms_{f2}) = [f2 - \text{offsetLinkFlag}, 2f]$$

$$\text{dom}(ms_{flag}) = [flag, flag]$$

$$\text{dom}(ms_{link}) = [link, link + 1]$$

$$\text{dom}(ms_{stk}) = [b_{stk}, e_{stk}]$$

$$\text{dom}(ms_{adv}) = [b_{adv}, e_{adv}]$$

$$ms_{malloc} :_n [0 \mapsto \iota_{malloc,0}] \quad \text{for all } n \in \mathbb{N}$$

and

- $ms_{f2}(f2 - \text{offsetLinkFlag}) = ((RO, \text{global}), link, link + 1, link)$ ,  $ms_{f2}(f2 - \text{offsetLinkFlag} + 1) = ((RW, \text{global}), flag, flag, flag)$ , the rest of  $ms_{f2}$  contains the code of  $f2$ .
- $ms_{flag} = [flag \mapsto 0]$
- $ms_{link} = [link \mapsto c_{malloc}, link + 1 \mapsto c_{adv}]$
- $ms_{adv}(b_{adv}) = c_{link}$  and  $\forall a \in [b_{adv} + 1, e]$ .  $ms_{adv}(a) \in \mathbb{Z}$

if

$$(reg[pc \mapsto c_{f2}][r_{stk} \mapsto c_{stk}], m) \rightarrow_n (\text{halted}, m'),$$

then

$$m'(flag) = 0$$

LEMMA A.15 (CORRECTNESS LEMMA FOR f3, DETAILED VERSION OF LEMMA 8.3). *For all*  $n \in \mathbb{N}$  *let*

$$c_{adv} \stackrel{\text{def}}{=} ((E, \text{global}), b_{adv}, e_{adv}, b_{adv} + \text{offsetLinkFlag})$$

$$c_{f3} \stackrel{\text{def}}{=} ((RWX, \text{global}), f3 - \text{offsetLinkFlag}, 3f, f3)$$

$$c_{stk} \stackrel{\text{def}}{=} ((RWLX, \text{local}), b_{stk}, e_{stk}, b_{stk} - 1)$$

$$c_{malloc} \stackrel{\text{def}}{=} ((E, \text{global}), b_{malloc}, e_{malloc}, b_{malloc} + \text{offsetLinkFlag})$$

$$c_{link} \stackrel{\text{def}}{=} ((RO, \text{global}), link, link + 1, link)$$

$$reg \in \text{Reg}$$

$$m \stackrel{\text{def}}{=} ms_{f3} \uplus ms_{flag} \uplus ms_{link} \uplus ms_{adv} \uplus ms_{malloc} \uplus ms_{stk} \uplus ms_{frame}$$

and

- $c_{\text{malloc}}$  satisfies the specification for malloc.

where

$$\text{dom}(ms_{f3}) = [f3 - \text{offsetLinkFlag}, 3f]$$

$$\text{dom}(ms_{\text{flag}}) = [\text{flag}, \text{flag}]$$

$$\text{dom}(ms_{\text{link}}) = [\text{link}, \text{link} + 1]$$

$$\text{dom}(ms_{\text{stk}}) = [b_{\text{stk}}, e_{\text{stk}}]$$

$$\text{dom}(ms_{\text{adv}}) = [b_{\text{adv}}, e_{\text{adv}}]$$

$$ms_{\text{malloc}} \cdot n \ [0 \mapsto \iota_{\text{malloc},0}]$$

and

- $ms_{f3}(f3 - \text{offsetLinkFlag}) = ((\text{RO}, \text{global}), \text{link}, \text{link} + 1, \text{link}), ms_{f3}(f3 - \text{offsetLinkFlag} + 1) = ((\text{RW}, \text{global}), \text{flag}, \text{flag}, \text{flag}),$  the rest of  $ms_{f3}$  contains the code of  $f3$ .
- $ms_{\text{flag}} = [\text{flag} \mapsto 0]$
- $ms_{\text{link}} = [\text{link} \mapsto c_{\text{malloc}}, \text{link} + 1 \mapsto c_{\text{adv}}]$
- $ms_{\text{adv}}(b_{\text{adv}}) = c_{\text{link}}$  and all other addresses of  $ms_{\text{adv}}$  contain instructions.

if

$$(\text{reg}[\text{pc} \mapsto c_{f3}][r_{\text{stk}} \mapsto c_{\text{stk}}], m) \rightarrow_n (\text{halted}, m'),$$

then

$$m'(\text{flag}) = 0$$

LEMMA A.16 (CORRECTNESS OF  $g1$ , DETAILED VERSION OF LEMMA 8.4). For all  $n \in \mathbb{N}$  let

$$c_{\text{adv}} \stackrel{\text{def}}{=} ((\text{RWX}, \text{global}), b_{\text{adv}}, e_{\text{adv}}, b_{\text{adv}} + \text{offsetLinkFlag})$$

$$c_{g1} \stackrel{\text{def}}{=} ((\text{E}, \text{global}), g1 - \text{offsetLinkFlag}, 4f, g1)$$

$$c_{\text{stk}} \stackrel{\text{def}}{=} ((\text{RWLX}, \text{local}), b_{\text{stk}}, e_{\text{stk}}, b_{\text{stk}} - 1)$$

$$c_{\text{malloc}} \stackrel{\text{def}}{=} ((\text{E}, \text{global}), b_{\text{malloc}}, e_{\text{malloc}}, b_{\text{malloc}} + \text{offsetLinkFlag})$$

$$c_{\text{link}} \stackrel{\text{def}}{=} ((\text{RO}, \text{global}), \text{link}, \text{link}, \text{link})$$

$$m \stackrel{\text{def}}{=} ms_{g1} \uplus ms_{\text{flag}} \uplus ms_{\text{link}} \uplus ms_{\text{adv}} \uplus ms_{\text{malloc}} \uplus ms_{\text{stk}} \uplus ms_{\text{frame}}$$

where

- $c_{\text{malloc}}$  satisfies the specification for malloc with  $\iota_{\text{malloc},0}$

$$\text{dom}(ms_{g1}) = [g1 - \text{offsetLinkFlag}, 4f]$$

$$\text{dom}(ms_{\text{flag}}) = [\text{flag}, \text{flag}]$$

$$\text{dom}(ms_{\text{link}}) = [\text{link}, \text{link}]$$

$$\text{dom}(ms_{\text{stk}}) = [b_{\text{stk}}, e_{\text{stk}}]$$

$$\text{dom}(ms_{\text{adv}}) = [b_{\text{adv}}, e_{\text{adv}}]$$

$$ms_{\text{malloc}} \cdot n \ [0 \mapsto \iota_{\text{malloc},0}]$$

and

- $ms_{g1}(g1 - \text{offsetLinkFlag}) = ((\text{RO}, \text{global}), \text{link}, \text{link}, \text{link}), ms_{g1}(g1 - \text{offsetLinkFlag} + 1) = ((\text{RW}, \text{global}), \text{flag}, \text{flag}, \text{flag}),$  the rest of  $ms_{g1}$  contains the code of  $g1$  immediately followed by the code of  $f4$ .
- $ms_{\text{flag}} = [\text{flag} \mapsto 0]$

- $ms_{link} = [link \mapsto c_{malloc}]$
- $ms_{adv}(b_{adv}) = c_{link}$  and all other addresses of  $ms_{adv}$  contain instructions.
- $\forall a \in \text{dom}(ms_{stk}). ms_{stk}(a) = 0$

if

$$(reg_0[pc \mapsto c_{adv}][r_{stk} \mapsto c_{stk}][r_1 \mapsto c_{g1}], m) \rightarrow_n (\text{halted}, m'),$$

then

$$m'(\text{flag}) = 0$$

## A.6 Reasoning about programs definitions

*Definition A.17.* We say that  $(reg, ms)$  is looking at  $[i_0, \dots, i_n]$  followed by  $c_{next}$  iff

- $reg(pc) = ((p, g), b, e, a)$
- $p = \text{RWX}, p = \text{RX}, \text{ or } p = \text{RWLX}$
- $a + n \leq e, b \leq a \leq e$
- $ms(a + 0, \dots, a + n) = [i_0, \dots, i_n]$
- $c_{next} = ((p, g), b, e, a + n + 1)$

*Definition A.18.* We say that “ $(reg, ms)$  links *key* as *j* to  $c_{malloc}$ ” iff

- $reg(pc) = ((perm, g), b, e, a)$
- $ms(b) = ((\_, \_), b_{link}, \_, \_)$
- $ms(b_{link} + j) = c$

*Definition A.19.* We say that  $reg$  points to stack with  $ms_{stk}$  used and  $ms_{unused}$  unused iff

- $reg(r_{stk}) = ((\text{RWLX}, \text{local}), b_{stk}, e_{stk}, a_{stk})$
- $\text{dom}(ms_{unused}) = [a_{stk} + 1, \dots, e_{stk}]$
- $\text{dom}(ms_{stk}) = [b_{stk}, \dots, a_{stk}]$  LS: Maybe make it clear what happens when  $ms_{stk}$  is empty
- $b_{stk} - 1 \leq a_{stk}$