

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), 27 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; Dominique Devriese, iMinds-DistriNet, KU Leuven, Street1 Address1, -, -, Belgium, dominique.devriese@cs.kuleuven.be; Lars Birkedal, Department of Computer Science, Aarhus University, Street1 Address1, Aarhus, State1, 8210, Denmark, 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.

© 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.

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 (2).
- We define a novel calling convention enforcing control-flow correctness and encapsulation of stack frames (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; Maffeis 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 logical relation and the fundamental theorem expressing capability safety are presented in 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 (5). The examples demonstrate both the power of our formulation of capability safety and our calling convention.

For reasons of space, some details and all proofs have been omitted; please refer to the technical appendix [Skorstengaard et al. 2018] for those.

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 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 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 and their ordering are depicted in Figure 3: the permissions include 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. **Suggestion: add the locality hierarchy and a short description. Maybe introduce the \sqsubseteq relation here.**

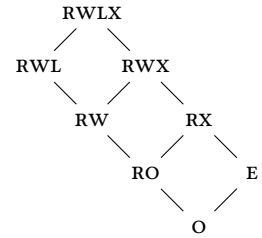


Fig. 3. Permission hierarchy

We assume a finite set of register names RegName . We define register files reg and memories ms as functions mapping register names resp. addresses to words. The state of the entire machine is 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 addition 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 **Suggestion: extend this figure with more (all?) of the operational semantics. In general, should we add all the missing definitions and make this document closer to being self-contained.** shows an excerpt of the operational semantics for a few representative instructions. Essentially, a configuration Φ 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 Φ . 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 .

$$\begin{aligned}
a &\in \text{Addr} \stackrel{\text{def}}{=} \mathbb{N} & r &\in \text{RegName} ::= \text{pc} \mid r_0 \mid r_1 \mid \dots \\
w &\in \text{Word} \stackrel{\text{def}}{=} \mathbb{Z} + \text{Cap} & \text{reg} &\in \text{Reg} \stackrel{\text{def}}{=} \text{RegName} \rightarrow \text{Word} \\
\text{perm} &\in \text{Perm} ::= \text{O} \mid \text{RO} \mid \text{RW} \mid \text{RWL} \mid & m &\in \text{Mem} \stackrel{\text{def}}{=} \text{Addr} \rightarrow \text{Word} \\
&\quad \text{RX} \mid \text{E} \mid \text{RWX} \mid \text{RWLX} & \Phi &\in \text{ExecConf} \stackrel{\text{def}}{=} \text{Reg} \times \text{Mem} \\
g &\in \text{Global} ::= \text{global} \mid \text{local} & ms &\in \text{MemSeg} ::= \text{Addr} \rightarrow \text{Word} \\
\text{Conf} &::= \text{ExecConf} + \{\text{failed}\} + \{\text{halted}\} \times \text{Mem} \\
\text{Cap} &::= \{((\text{perm}, g), b, e, a) \mid b, a \in \text{Addr}, e \in \text{Addr} \cup \{\infty\}\} \\
r &\in \mathbb{Z} + \text{RegName} \\
i &::= \text{jmp } r \mid \text{jnz } r r \mid \text{move } r r \mid \text{load } r r \mid \text{store } r r \mid \text{plus } r r r \mid \text{minus } r r r \mid \\
&\quad \text{lt } r r r \mid \text{lea } r r \mid \text{restrict } r r \mid \text{subseg } r r r \mid \text{isptr } r r \mid \text{getl } r r \mid \\
&\quad \text{getp } r r \mid \text{getb } r r \mid \text{gete } r r \mid \text{geta } r r \mid \text{fail} \mid \text{halt}
\end{aligned}$$

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

$$\begin{aligned}
\Phi &\rightarrow \begin{cases} \llbracket \text{decode}(n) \rrbracket (\Phi) & \text{if } \Phi.\text{reg}(\text{pc}) = ((\text{perm}, g), b, e, a) \text{ and } b \leq a \leq e \\ & \text{and } \text{perm} \in \{\text{RX}, \text{RWX}, \text{RWLX}\} \text{ and } \Phi.\text{mem}(a) = n \\ \text{failed} & \text{otherwise} \end{cases} \\
\text{updPc}(\Phi) &= \begin{cases} \Phi[\text{reg.pc} \mapsto \text{newPc}] & \text{if } \Phi.\text{reg}(\text{pc}) = ((\text{perm}, g), b, e, a) \\ & \text{and } \text{newPc} = ((\text{perm}, g), b, e, a + 1) \\ \text{failed} & \text{otherwise} \end{cases}
\end{aligned}$$

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 $\text{perm} \in \{\text{RWX}, \text{RWLX}, \text{RX}, \text{RW}, \text{RWL}, \text{RO}\}$
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 $w = ((\text{perm}', g'), b, e, a)$
geta $r_1 r_2$	$\text{updPc}(\Phi[\text{reg}.r_1 \mapsto a])$	$\Phi.\text{reg}(r_2) = ((_, _), _, _, a)$
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 $\text{newPc} = \Phi.\text{reg}(r)$
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 $\text{perm} \in \{\text{RWLX}, \text{RWL}\}$
...		
–	<i>failed</i>	otherwise

Fig. 2. An excerpt from the operational semantics.

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. Reviewer E, popl, would like \sqsubseteq and *decodePermPair* explained. If we explain all of these details, this section becomes a bit long and clunky. (see tex for suggestion) 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.

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.

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, but we refer to the technical appendix [Skorstengaard et al. 2018] for full details. Suggestion: add the malloc specification (maybe later as it uses the LR) and possibly a description of the linking tables. (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.

Reviewer C, popl, would have liked a discussion of the relationship between this capability machines and old capability machines from the literature (Cambridge CAP & IBM System 39, see tex comment).

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: Make it more clear that this CC can be used to protect multiple components (reviewer A, ESOP), see tex comment. 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. 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). 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.

However, if we are returning from a function that has itself invoked adversary code, 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

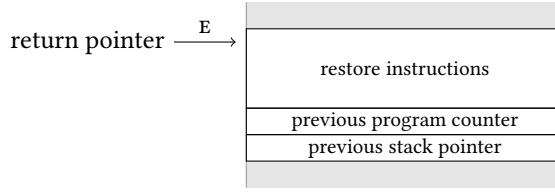


Fig. 4. Structure of an activation record

could use the unused part of the stack to store local pointers and load them from there after they get control back. 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. 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 6 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 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.

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 4. 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.¹ *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.)*

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

¹Note 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.

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 `E`-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`.

Reusable macro instructions Add examples of macros to this section to make it clear what we mean by macros. Move this subsection forward so it is just before the macros are used, i.e. the example section. Reviewer C, ESOP, asks for high-level specifications for the macros in the style of `malloc`. I am not sure we want to actually include this in the submission, so it is maybe not here we should put our attention. (see review in tex comment) We define a number of reusable macros capturing the calling convention and other conveniences. All macros that use the stack assume a stack pointer in register r_{stk} . The macro `fetch r name` fetches the capability related to *name* from the linking table and stores it in register r . The macros `push r` and `pop r` add and remove elements from the stack. The macro `prepstk r` is used when a callback is invoked by the adversary and prepares the received stack pointer by checking that it has permission `RWLX`. The macro `scall $r(\overline{r_{args}}, \overline{r_{priv}})$` jumps to the capability in register r in the manner described above. That is, it pushes local state (the contents of registers $\overline{r_{priv}}$) and the activation record (return code, return pointer, stack pointer) to the stack, creates an `E` return pointer, restricts the stack pointer, clears the unused part of the stack, clears the necessary registers and jumps to r . Upon return, the private state is restored. The macro `mclear r` clears all the memory the capability in register r has authority over. The macro `rclear $regSet$` clears all the registers in $regSet$. The macro `reqglob r` checks whether the word in register r is a global capability. The macro `crtcls $(\overline{x_i}, \overline{r_i})$ r` allocates a closure where r points to the closure's code and a new environment is allocated (using `malloc`) where the contents of $\overline{r_i}$ is stored. In the code referred to by r , an implicit fetch happens when an instruction refers to x_i .

The technical appendix [Skorstengaard et al. 2018] contains detailed descriptions of all the macros. **Suggestion:** Add the call implementation (and possibly other macros?) in the above?

4 LOGICAL RELATION

Suggestions for things to add heres

- Double monotonicity lemma (lemma 77 and 79) with a pointer back to local capability section of popl intro to this section.
- A section with the lemmas we use for reasoning about programs
 - lemma 58, `scall` works. Include description of how it handles the common things in call and leaves the argumentation for callee code and continuation.
 - lemma 59, `malloc` works. If this is included, we should probably also include the `malloc` specification.
 - lemma 60, create closure works.

- Include a description about how the FTLR is used to reason about unknown code.
- Mention that this is a step on the way towards having a program logic, and then in the discussion discuss why a program logic for this language is not what we need.
- 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.

Look this introduction and the following section through to check for overlap. Also consider making backwards references to this introduction from the places where we use advanced techniques, so it becomes clear that they are necessary.

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, 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, in the next section.

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 time². 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).

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

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.

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

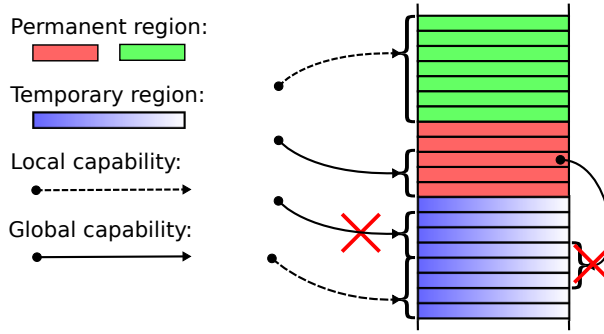


Fig. 5. The relation between local/global capabilities and temporary/permanent regions. The colored fields are regions governing parts of memory. Global capabilities cannot depend on temporary regions.

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.³ Therefore, we can allow the second invocation to happen in *private* future worlds ($W' \sqsubseteq^{\text{priv}} W$), 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 ($W' \sqsubseteq^{\text{pub}} W$, 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 6). 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 capabilities⁴⁵).

4.2 Worlds

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.

³We 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.

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

⁵For 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.

In other words, they are protocols for the region's memory. These are similar to what has been used in logical relations for high-level languages [Ahmed et al. 2009; Devriese et al. 2016; Dreyer et al. 2012]. 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 5.

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 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 Wor and two relations $\sqsubseteq^{\text{priv}}$ and \sqsubseteq^{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) Wor and preorders $\sqsubseteq^{\text{priv}}$ and \sqsubseteq^{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 ξ such that*

$$\begin{aligned} \xi : \text{Wor} &\cong \blacktriangleright(\mathbb{N} \xrightarrow{\text{fin}} \text{Region}) \\ \text{Region} &= \{\text{revoked}\} \uplus \\ &\quad \{\text{temp}\} \times \text{State} \times \text{Rels} \times (\text{State} \rightarrow (\text{Wor} \xrightarrow[\sqsubseteq^{\text{pub}}]{\text{mon, ne}} \text{UPred}(\text{MemSeg}))) \uplus \\ &\quad \{\text{perm}\} \times \text{State} \times \text{Rels} \times (\text{State} \rightarrow (\text{Wor} \xrightarrow[\sqsubseteq^{\text{priv}}]{\text{mon, ne}} \text{UPred}(\text{MemSeg}))) \end{aligned}$$

and for $W, W' \in \text{Wor}$.

$$\begin{aligned} W' \sqsubseteq^{\text{priv}} W &\Leftrightarrow \xi(W') \sqsubseteq^{\text{priv}} \xi(W) \\ W' \sqsubseteq^{\text{pub}} W &\Leftrightarrow \xi(W') \sqsubseteq^{\text{pub}} \xi(W) \end{aligned}$$

In the above theorem, $\text{State} \times \text{Rels}$ corresponds to the aforementioned state transition system where Rels contains pairs of relations corresponding to the public and private transitions, and 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 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* $W' \sqsupseteq^{pub} W$ requires that $\text{dom}(W') \supseteq \text{dom}(W)$ and $\forall r \in \text{dom}(W). W'(r) \sqsupseteq^{pub} W(r)$. 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 $W' \sqsupseteq^{priv} W$ is defined similarly, using a private future region relation. The *public future region* relation is the simplest. It satisfies the following properties:

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

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

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.

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 6 and 7 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 technical appendix [Skorstengaard et al. 2018].

First, the *observation relation* 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

$$\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})) \left| \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. \right\} \\
\mathcal{R} &: \text{World} \xrightarrow[\sqsupset^{\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}) \left| \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. \right\} \\
\mathcal{V} &: \text{World} \xrightarrow[\sqsupset^{\text{pub}}]{\text{mon}, ne} \text{UPred}(\text{Word}) \\
\mathcal{V}(W) &\stackrel{\text{def}}{=} \{(n, i) \mid i \in \mathbb{Z}\} \cup \{(n, ((o, g), b, e, a))\} \cup \\
&\quad \left\{ (n, ((\text{RW}, g), b, e, a)) \left| \begin{array}{l} (n, (b, e)) \in \text{readCond}(g)(W) \text{ and} \\ (n, (b, e)) \in \text{writeCond}(\iota^{\text{nw}}, g)(W) \end{array} \right. \right\} \cup \\
&\quad \{(n, ((\text{E}, g), b, e, a)) \mid (n, (b, e, a)) \in \text{enterCond}(g)(W)\} \cup \\
&\quad \left\{ (n, ((\text{RWLX}, g), b, e, a)) \left| \begin{array}{l} (n, (b, e)) \in \text{readCond}(g)(W) \text{ and} \\ (n, (b, e)) \in \text{writeCond}(\iota^{\text{pw}}, g)(W) \text{ and} \\ (n, (\{\text{RWLX}, \text{RWX}, \text{RX}\}, b, e)) \in \text{execCond}(g)(W) \end{array} \right. \right\} \\
&\quad \cup \dots \text{ and so on for permissions } \text{RO}, \text{RWL}, \text{RX}, \text{ and } \text{RWX}.
\end{aligned}$$

Fig. 6. The logical relation.

$$\begin{aligned}
\text{readCond}(g)(W) &= \left\{ (n, (b, e)) \left| \begin{array}{l} \exists r \in \text{localityReg}(g, W). \\ \exists [b', e'] \supseteq [b, e]. W(r) \stackrel{n}{\sqsubseteq} \iota_{b', e'}^{\text{pw}} \end{array} \right. \right\} \\
\text{writeCond}(\iota, g)(W) &= \left\{ (n, (b, e)) \left| \begin{array}{l} \exists r \in \text{localityReg}(g, W). \\ W(r) \text{ is address-stratified and} \\ \exists [b', e'] \supseteq [b, e]. W(r) \stackrel{n-1}{\sqsupseteq} \iota_{b', e'} \end{array} \right. \right\} \\
\text{execCond}(g)(W) &= \left\{ (n, (P, b, e)) \left| \begin{array}{l} \forall n' < n, W' \sqsupseteq W, a \in [b', e'] \subseteq [b, e], \text{perm} \in P. \\ (n', ((\text{perm}, g), b', e', a)) \in \mathcal{E}(W') \end{array} \right. \right\} \\
\text{enterCond}(g)(W) &= \left\{ (n, (b, e, a)) \left| \begin{array}{l} \forall n' < n. \forall W' \sqsupseteq W. \\ (n', ((\text{RX}, g), b, e, a)) \in \mathcal{E}(W') \end{array} \right. \right\} \\
&\quad \text{where } g = \text{local} \Rightarrow \sqsupseteq = \sqsupseteq^{\text{pub}} \text{ and } g = \text{global} \Rightarrow \sqsupseteq = \sqsupseteq^{\text{priv}}
\end{aligned}$$

Fig. 7. Permission-based conditions

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 E capabilities) are considered safe program counters because 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 7, we define such a condition for each kind of permission a capability can have.

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 $\iota_{b,e}^{pwl}$. The definition of that standard region is omitted for space reasons, but 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 [Skorstengaard et al. 2018]; 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 3, is necessary for any safe use of local capabilities.

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 updated, 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.

As enter capabilities are the only capabilities that have to use the address they point to, it is natural that the *enterCond* is also the only permission-based condition that depends on the current address of the capability. 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 $(n, (b, e)) \in writeCond(\iota^{nwl}, g)(W)$
- $perm = RWLX$ and $(n, (b, e)) \in readCond(g)(W)$ and $(n, (b, e)) \in writeCond(\iota^{pwl}, g)(W)$,

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

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 5 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 5 as the examples depend on these properties for correctness. See the technical appendix [Skorstengaard et al. 2018] for a detailed proof (by induction over the step-index n) of the theorem.

<pre> f1: push 1 fetch r_1 <i>adv</i> scall r_1([],[]) pop r_1 assert r_1 1 halt </pre>	<pre> f2: malloc r_l 1 store r_l 1 fetch r_1 <i>adv</i> call r_1([], [r_l]) assert r_l 1 halt </pre>
---	--

Fig. 8. Two example programs that rely on local-state encapsulation. f_1 uses our stack-based calling convention. f_2 does not rely on a stack.

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 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 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 3. 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].

5.1 Encapsulation of Local State

f_1 and f_2 in Figure 8 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 f_1 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, f_2 uses *malloc* to allocate memory for its local state and uses an activation-record based calling convention (described in the technical appendix) to run 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 5.1 and 5.2 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 5.1. *Let*

$$\begin{aligned} c_{adv} &\stackrel{\text{def}}{=} ((E, \text{global}), \dots) & c_{stk} &\stackrel{\text{def}}{=} ((RWLX, \text{local}), \dots) \\ c_{f1} &\stackrel{\text{def}}{=} ((RWX, \text{global}), \dots) & c_{link} &\stackrel{\text{def}}{=} ((RO, \text{global}), \dots) \\ c_{malloc} &\stackrel{\text{def}}{=} ((E, \text{global}), \dots) & \text{reg} &\in \text{Reg} \\ m &\stackrel{\text{def}}{=} ms_{f1} \uplus ms_{flag} \uplus ms_{link} \uplus ms_{adv} \uplus ms_{malloc} \uplus ms_{stk} \uplus ms_{frame} \end{aligned}$$

where each of the capabilities have an appropriate range of authority and pointer⁶. Furthermore

- ms_{f1} contains c_{link} , c_{flag} and the code of $f1$
- $ms_{flag}(flag) = 0$
- ms_{link} contains c_{adv} and c_{malloc}
- ms_{adv} contains c_{link} and otherwise only instructions.

If $(\text{reg}[pc \mapsto c_{f1}][r_{stk} \mapsto c_{stk}], m) \rightarrow^* (\text{halted}, m')$, then $m'(flag) = 0$

To prove Lemma 5.1, 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 an anti-reduction lemma, it suffices to show that the configuration is safe after some reduction steps. We then use a general lemma for reasoning about `scall`, 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, 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 $f2$:

LEMMA 5.2. *Making similar assumptions about capabilities and linking as in Lemma 5.1 but assuming no stack pointer, if $(\text{reg}[pc \mapsto c_{f2}], m) \rightarrow^* (\text{halted}, m')$, then $m'(flag) = 0$.*

5.2 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 $f3$ and $g1$ in Figure 9.

In $f3$ 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 $f3$ would have 2 on top of its stack and the assertion would fail. However, the security measures in 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:

LEMMA 5.3. *Making similar assumptions about capabilities and linking as in Lemma 5.1 if $(\text{reg}[pc \mapsto c_{f3}][r_{stk} \mapsto c_{stk}], m) \rightarrow^* (\text{halted}, m')$, then $m'(flag) = 0$.*

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

⁶These assumptions are kept intentionally vague for brevity. Full statements are in the technical appendix [Skorstengaard et al. 2018].

<pre> g1: malloc r₂ 1 store r₂ 0 move pc r₃ lea r₃ offset crtcls [(x, r₂)] r₃ rclear RegName \ {pc, r₀, r₁} jmp r₀ f4: reqglob r₁ prepstk r_{stk} (continues in next column) </pre>	<p>(continued from previous column)</p> <pre> store x 0 scall r₁([], [r₀, r₁, r_{env}]) store x 1 scall r₁([], [r₀, r_{env}]) load r₁ x assert r₁ 1 mclear r_{stk} rclear RegName \ {r₀, pc} jmp r₀ </pre>	<pre> f3: push 1 fetch r₁ adv scall r₁([], [r₁]) pop r₂ assert r₂ 1 push 2 scall r₁([], []) halt </pre>
---	---	--

Fig. 9. Two programs that rely on well-bracketedness of scalls to function correctly. *offset* is the offset to f4.

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 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 5.4. *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 5.1. Then if $(\text{reg}_0[\text{pc} \mapsto c_{adv}][r_{stk} \mapsto c_{stk}][r_1 \mapsto c_{g1}], m) \rightarrow^ (\text{halted}, m')$, then $m'(\text{flag}) = 0$.*

As explained in 3, the macro-instruction `reqglob r1` 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 5.4, 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.

Suggestion: *Add an actual proof sketch* In Lemma 5.4 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

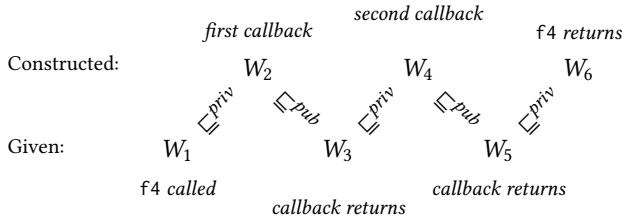


Fig. 10. Illustration of the worlds in the proof of Lemma 5.4. In the proof, the top row of worlds are constructed by us, while the bottom row of worlds are given. W_6 is constructed such that $W_6 \sqsupseteq^{priv} W_5$ and $W_6 \sqsupseteq^{pub} W_1$.

to guarantee that the local state is unchanged. This is important because the return pointer that f_4 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 f_4 returns.

In Figure 10, we illustrate how we accommodate this in the proof of Lemma 5.4 by constructing appropriate worlds for all the situations where control is passed to the adversary. The potentially local return pointer received by f_4 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 ι^{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.

6 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. Reviewer C, popl, would like to know how local capabilities relate to borrowing (see tex comment). 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 g_1 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 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.

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 safe arguments. 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 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 6). 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 6 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.

7 RELATED WORK

Finally, we summarize how our work relates to previous work. We do not repeat the work we discussed in [6](#).

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 [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 papers 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 [6](#), 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.

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 [6](#) 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. 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 ?. 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
- 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.
- 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 – 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>
- 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. Maffei, 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>
- Leo Osvald, Grégory Essertel, Xilun Wu, Lilliam I. González Alayón, 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 and Return Pointer Management - Technical Appendix Including Proofs and Details*. 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 Compositionanl 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.