

Reasoning about Object Capabilities with Logical Relations and Effect Parametricity

Dominique Devriese
iMinds-DistriNet, KU Leuven, Belgium
Email: dominique.devriese@cs.kuleuven.be

Lars Birkedal
Aarhus University, Denmark
Email: birkedal@cs.au.dk

Frank Piessens
iMinds-DistriNet, KU Leuven, Belgium
Email: frank.piessens@cs.kuleuven.be

Abstract—Object capabilities are a technique for fine-grained privilege separation in programming languages and systems, with important applications in security. However, current formal characterisations do not fully capture capability-safety of a programming language and are not sufficient for verifying typical applications. Using state-of-the-art techniques from programming languages research, we define a logical relation for a core calculus of JavaScript that better characterises capability-safety. The relation is powerful enough to reason about typical capability patterns and supports evolvable invariants on shared data structures, capabilities with restricted authority over them and isolated components with restricted communication channels. We use a novel notion of effect parametricity for deriving properties about effects. Our results imply memory access bounds that have previously been used to characterise capability-safety.

1. Introduction

Privilege separation between components and the Principle Of Least Authority (POLA) are key ingredients for constructing secure systems. Since decades, systems with *object capabilities* [1], [2] have supported a very fine-grained separation of privileges. Lately, a renewed research interest in the technique has produced implementations at the level of the OS [3], the processor [4] and the programming language [2], [5]–[8]. Applications of the technique include sandboxing untrusted code [6], [9], efficient security auditing [5], fault isolation [3]. The technique also has potential applications outside security (e.g. enabling testing in simulated environments, enforcing architectural choices), but these have not drawn as much attention so far.

In the object capabilities model, effects can only be produced by sending messages to *objects*. These can be objects in the sense of object-oriented programming, but not necessarily so. *Device objects* are primitive objects that model resources in the outside world and produce effects in the outside world when receiving a message. *Instance objects* are programmer-defined objects that may hold private state (data or references to other objects). They execute programmer-defined code upon receipt of a message and this may include sending messages to other objects. References

to objects are called *capabilities*, as they represent the authority to invoke methods on those objects.

A *capability-safe* programming language implements certain restrictions to enforce the object-capability model. Such a language guarantees that sending messages to objects is indeed the *only* way to produce effects and that capabilities cannot be forged. Additionally, the language excludes globally accessible mutable state to provide control over the capabilities that a component starts out with. These restrictions enable low-cost fine-grained privilege separation; a programmer defines the authority of a component by controlling the capabilities it holds.

To make this more concrete, consider a web page embedding an untrusted advertisement. The advertisement’s initialisation requires access to the DOM for at least the ad’s designated location on the page. However, it should be prevented from accessing or modifying other parts of the DOM. In a capability-safe JavaScript-like language, the web page could use the following function *rnode* to construct a restricted capability for accessing a subtree of the DOM:

$$rnode \stackrel{\text{def}}{=} \text{func}(node, d) \left\{ \begin{array}{l} \text{getChild} = \text{func}(id) \{ rnode(node.getChild(id), d + 1) \} \\ \text{parent} = \text{func}() \left\{ \begin{array}{l} \text{if } (d \leq 0) \{ \text{error} \} \text{ else} \\ \{ rnode(node.parent(), d - 1) \} \end{array} \right\} \\ \text{getProp} = \text{func}(id) \{ node.getProp(id) \} \\ \text{setProp} = \text{func}(id, v) \{ node.setProp(id, v) \} \\ \text{addChild} = \text{func}(id) \{ rnode(node.addChild(id), d + 1) \} \\ \text{delChild} = \text{func}(id) \{ node.delChild(id) \} \end{array} \right.$$

The function wraps a DOM node and forwards invocations to hypothetical *getChild*, *parent*... methods, but prevents access to nodes more than *d* levels higher in the tree.

Now assume the next web page initialisation function:

$$\text{initWebPage} \stackrel{\text{def}}{=} \text{func}(document, ad) \left\{ \begin{array}{l} \text{document.setProp}(\text{"someProperty"}, 42) \\ \text{let } (adNode = \text{document.addChild}(\text{"ad_div"})) \\ \text{let } (rAdNode = rnode(adNode, 0)) \\ \text{ad.initialize}(rAdNode) \\ \text{document.getProp}(\text{"someProperty"}) == 42 \end{array} \right.$$

A capability-safe language enforces the encapsulation of the $rAdNode$ object, prevents direct access to $document$ through other channels and rules out mutable global state which could (have been made to) contain a reference to $document$. If we assume sane behaviour of $document$ and that ad has no capabilities besides $rAdNode$ to begin with, then we can convince ourselves that if $initWebPage$ terminates, it must return `true`. But can we make this reasoning precise?

Formal reasoning about code in object capability languages requires a good understanding of the model. What exactly does it mean that a language is *capability-safe*? What guarantees can we provide about code written in it and how can we formally prove properties for maximum assurance? Several researchers have worked on this [2], [9], [10], but to this day, our understanding of the model is not satisfactory.

The problem is that previous research (except Spiessens [10], [11], but see Section 7) focuses on *reference graph* dynamics. For a state in a program’s execution, the reference graph is the graph with the allocated objects as nodes, and the references they hold to each other as edges. Properties like *No Authority Amplification* and *Only Connectivity Begets Connectivity* [2], [9] restrict how the reference graph can evolve, depending on the references held by the executing code.

Unfortunately, these properties offer only a conservative bound on components’ authority: the *topology-only bound*. It is based on syntactic structure of objects (whether or not they contain a reference to each other) and ignores their *behaviour*. This limitation is important, because a key quality of the model is the ability to define *custom capabilities* as instance objects. Such capabilities (like our $rAdNode$) restrict other capabilities, make them conditional or revocable or otherwise modify and combine them. The object capabilities community has studied many patterns for defining custom capabilities, but their soundness typically does not follow from topology-only bounds on memory access.

Contributions and Outline In this paper, we propose a novel, more semantic approach for reasoning about object capability languages. The key contribution is a step-indexed Kripke logical relation [12], [13] with two key features: support for reasoning about custom effect properties using *effect parametricity* (EP) and support for reasoning about shared data structures with evolvable invariants and authority using a novel form of possible worlds. Because we use many concepts that may be new to a security audience, we provide a gradual introduction in Section 2, where we introduce (first) logical relations and effect parametricity and (second) step-indexed logical relations for two simple languages without a shared mutable store. Section 3 repeats the definition of λ_{JS} , a pre-existing core calculus for JavaScript [14] and presents the logical relation capturing its capability-safety. It uses ideas and techniques from previous work [15]–[17] but also novel ideas, in particular a variation on Dreyer et al.’s public/private transitions [16] for modelling authority over shared data and our notion of *effect parametricity* to support custom properties about effects. We demonstrate in Section 4 that the logical relation can be used to verify non-trivial examples involving evolvable invariants on shared

data structures, restricted authority over them (like the web-page-with-an-ad example above) and isolating components with restricted communication channels. These examples show that, contrary to previous work on capability-safety, our results capture more than just the topology-only bound on authority. To demonstrate that our results encompass previous results, Section 6 derives such standard syntactical bounds after first repeating a formulation of the properties by Maffeis et al. in Section 5. Section 7 discusses related work and Section 8 concludes. A companion technical report (TR) provides more details about the relational generalisation and contains technical details and proofs of our results [18].

2. A Simpler Setting

In order to explain our work to a wide audience, this section gradually introduces concepts and techniques for two very simple languages, as preparation for the subset of λ_{JS} that we work with later on. The first is $\lambda_{out,FO}$, a simple first-order calculus (no lambdas) which we then extend to the higher-order $\lambda_{out,HO}$. Both contain only a single primitive capability, for producing textual output.

Contrary to λ_{JS} , neither calculus has a mutable store. This reduces the technicalities for reasoning, but makes the results less interesting. Without shared mutable data, there is no point in defining invariants on it or authority over it and less ways to pass capabilities around. What remains is *effect parametricity*: a general property that allows proving custom properties about untrusted expressions’ effects.

Although $\lambda_{out,FO}$ allows introducing our techniques in a very simple setting, it lacks the ability to define custom capabilities, limiting the value of effect parametricity. This limitation is removed in $\lambda_{out,HO}$, where we define and prove correct a simple custom capability that restricts output to upper-case. This simple custom capability models practical scenarios where untrusted code in the browser might receive a capability for injecting only valid HTML in the web page or evaluating a safe subset of JavaScript.

The definitions in this section are an instance of a general technique for reasoning about higher-order programming languages known as *logical relations* (LRs; see e.g. [12], [13], [15], [16]). Our LRs are *unary*, because we use predicates rather than relations. We do not assume prior knowledge about LRs and motivate and introduce the techniques gradually throughout the text.

2.1. An output capability in a first-order calculus

We define $\lambda_{out,FO}$ in Figure 1. It is a simple imperative untyped programming language with expressions e and values v . The (small-step) operational semantics is split in two parts: pure evaluations $e \rightarrow e'$ and impure evaluations $e \rightarrow_o e'$. The latter may produce textual output in the list of strings o . Multi-step impure evaluations $e \rightarrow_o^* e'$, concatenate the outputs of the individual steps. The pure evaluation is defined in terms of primitive pure evaluations $e \mapsto e'$ and evaluation contexts E that produce a strict evaluation order. There are standard `if`, `while`, `let` and sequence $(e; e')$

Syntax:

$$\begin{aligned}
e \in \text{Expr} &::= x \mid v \mid \text{let}(x = e) e \mid \text{if}(e)\{e\} \text{ else } \{e\} \\
&\mid e; e \mid \text{while}(e)\{e\} \mid e.\text{print}(e) \\
v \in \text{Val} &::= \text{num} \mid \text{str} \mid \text{bool} \mid \text{undef} \mid \text{null} \mid \text{out} \\
E &::= \cdot \mid \text{let}(x = E) e \mid \text{if}(E)\{e\} \text{ else } \{e\} \\
&\mid E; e \mid E.\text{print}(e) \mid v.\text{print}(E)
\end{aligned}$$

Pure evaluations:

$$\begin{aligned}
&\text{if}(\text{true})\{e_1\} \text{ else } \{e_2\} \hookrightarrow e_1 \quad (\text{E-IFTRUE}) \\
&\text{if}(\text{false})\{e_1\} \text{ else } \{e_2\} \hookrightarrow e_2 \quad (\text{E-IFFALSE}) \\
v; e \hookrightarrow e \quad (\text{E-BEGIN-DISCARD}) \quad \text{let}(x = v) e \hookrightarrow e[x/v] \quad (\text{E-LET}) \\
&\text{while}(e_1)\{e_2\} \hookrightarrow \\
&\quad \text{if}(e_1)\{e_2; \text{while}(e_1)\{e_2\}\} \text{ else } \{\text{undef}\} \quad (\text{E-WHILE}) \\
&\quad \frac{e_1 \hookrightarrow e_2}{E\langle e_1 \rangle \rightarrow E\langle e_2 \rangle} \quad (\text{E-CXT})
\end{aligned}$$

Impure evaluations:

$$\begin{aligned}
&\frac{e_1 \rightarrow e_2}{e_1 \rightarrow \parallel e_2} \quad (\text{E-PURE}) \\
&E\langle \text{out}.\text{print}(\text{str}) \rangle \rightarrow_{[\text{str}]} E\langle \text{undef} \rangle \quad (\text{E-OUT})
\end{aligned}$$

Figure 1. $\lambda_{\text{out},FO}$, a first-order calculus with an output capability.

expressions. Output is produced by invoking a `print` method on the primitive capability `out`. Note that `out` is an internal value; similar to an object reference (`0x1234`) in a language like Java, programmers cannot write it but the language runtime may give access to it during execution, for example as an argument to a program’s `main`. Note also that $\lambda_{\text{out},FO}$ provides plenty of opportunity for stuck terms, e.g. an `if` branching on a non-boolean or invoking `print` on `true`. We could also throw exceptions in such cases, but stuck expressions are simpler.

Let us now formulate our Fundamental Theorem, guaranteeing effect parametricity. Intuitively, the theorem states that if an *acceptability* property (e.g. only producing upper-case output) holds for the possible effects of the values that an expression has access to (e.g. the mentioned custom capability), then it must also hold for the expression as a whole. These acceptability properties can be chosen freely, as long as they are *admissible*. To understand the conditions for admissibility, one should understand that not every property is enforceable using the general approach of taking arbitrary untrusted code and controlling the capabilities that it has access to. For example, no matter what capabilities a piece of code has access to, it cannot be forced to return a specific value, so admissible effect properties should not distinguish expressions by their return values.

To formalise this, we need to put in place some technical machinery: we (a) define the properties on expressions that we work with, (b) define when such a property is “admissible” and (c) formulate the actual theorem, in terms of an expression’s scope. Let us consider each of the steps.

Semantic properties First, we limit the properties we work with to *semantic* properties, which only consider the

meaning of an expression and not syntactic artefacts. We do this by treating as equivalent all expressions that purely evaluate to the same expression and we only consider properties on *commands*: a syntactic class of expressions that are the intended end result of *pure* evaluations (like *values* are the intended end result of impure evaluations). In addition to values, this also includes expressions like `out.print(“abc”)` which are normal w.r.t. pure but not w.r.t. impure evaluation. Formally, we define a command as either a value or an expression $E\langle \text{cmd}_0 \rangle$, blocked on a print call cmd_0 :

$$\text{cmd}_0 ::= v.\text{print}(v) \quad \text{cmd} \in \text{Cmd} ::= E\langle \text{cmd}_0 \rangle \mid v$$

We will work with *command predicates*, i.e. subsets of commands $P \in \mathcal{P}(\text{Cmd})$. The *expression extension* $\mathcal{E}[P]$ of such a P accepts an expression e if commands that it evaluates to are in P .

$$\mathcal{E}[P] \stackrel{\text{def}}{=} \{e \mid \text{If } e \rightarrow^* \text{cmd}, \text{ then } \text{cmd} \in P\}$$

A first command predicate is *PureVal*, containing only pure values, i.e. numbers, strings, booleans, `undef` or `null`:

$$\text{PureVal} \stackrel{\text{def}}{=} \text{Num} \cup \text{Str} \cup \text{Bool} \cup \{\text{undef}, \text{null}\}.$$

Properties about effects As explained above, effect parametricity is concerned with *admissible* effect properties that are in principle enforceable using the object capability approach. We formalise this in a very general way by taking inspiration from *monads* (used to model effectful computations in pure functional languages like Haskell [19], [20]). We define an *effect interpretation* as a couple (μ, ρ) , where $\mu \in \mathcal{P}(\text{Val}) \rightarrow \mathcal{P}(\text{Cmd})$ and $\rho \in \mathcal{P}(\text{Cap})$ with $\text{Cap} = \{\text{out}\}$ the set of primitive capabilities. Intuitively, for a given value predicate P of acceptable result values, the command predicate μP defines a set of expressions that impurely evaluate to values in P . Typically, μP will not contain all such expressions, but only those that produce effects deemed (in some way) acceptable. ρ defines a set of acceptable primitive capabilities. Defined in terms of the effect interpretation (μ, ρ) , the value predicate $\text{Val}_{\mu, \rho}$ defines the full set of acceptable λ_{JS} values: pure values and primitive capabilities in ρ . $\text{Val}_{\mu, \rho} \stackrel{\text{def}}{=} \text{PureVal} \cup \rho$.¹

The couple (μ, ρ) is *admissible* or *valid* if it satisfies three conditions or *axioms*. We list them and explain below:²

- A-PURE: For a value $v \in P$, v must also be in μP .
- A-BIND: If $\text{cmd} \in \mu P$ and $E\langle v \rangle \in \mathcal{E}[\mu P']$ for all values $v \in P$, then $E\langle \text{cmd} \rangle \in \mathcal{E}[\mu P']$.
- A-PRINT: If $v \in \text{Val}_{\mu, \rho}$ then $v.\text{print}(v') \in \mu \text{Val}_{\mu, \rho}$.

Axiom A-PURE states that it must be acceptable to produce no effects: a value in P must also be in μP . For the second axiom, we consider a command cmd in μP (i.e. producing acceptable effects and a result in P) and an execution context E such that $E\langle v \rangle$ is in $\mu P'$ (i.e. produces acceptable effects and a result in P') for any $v \in P$ (i.e. any possible result of cmd). Then Axiom A-BIND requires that the composed

1. Note that $\text{Val}_{\mu, \rho}$ does not depend on μ , but that will change later.
2. Readers familiar with monads may recognise a correspondence between A-PURE and A-BIND and the monad operations *return* and *bind*.

$$\begin{array}{c}
\Gamma ::= \emptyset \mid \Gamma, x \\
\frac{v \neq \text{out}}{\Gamma \vdash v} \quad \frac{\Gamma \vdash e_1 \quad \Gamma, x \vdash e_2}{\Gamma \vdash \text{let}(x = e_1) e_2} \quad \frac{x \in \Gamma}{\Gamma \vdash x} \\
\frac{\Gamma \vdash e_1 \quad \Gamma \vdash e_2 \quad \Gamma \vdash e_3}{\Gamma \vdash \text{if}(e_1)\{e_2\} \text{ else } \{e_3\}} \quad \frac{\Gamma \vdash e_1 \quad \Gamma \vdash e_2}{\Gamma \vdash e_1; e_2} \\
\frac{\Gamma \vdash e_1 \quad \Gamma \vdash e_2}{\Gamma \vdash \text{while}(e_1)\{e_2\}} \quad \frac{\Gamma \vdash e_1 \quad \Gamma \vdash e_2}{\Gamma \vdash e_1. \text{print}(e_2)}
\end{array}$$

Figure 2. Well-scopedness of expressions in $\lambda_{\text{out},FO}$.

expression $E\langle \text{cmd} \rangle$ should be in $\mathcal{E}[\mu \ P']$. In other words, producing the acceptable effects of cmd , next the acceptable effects of $E\langle v \rangle$ (with v the result value of cmd) and a result in P' should be acceptable in $\mu \ P'$. The third and final axiom A-PRINT makes the link between μ and ρ by requiring that an effect produced using acceptable values (including primitive capabilities in ρ) should be acceptable: if $v \in \text{Val}_{\mu,\rho}$, then $v. \text{print}(v')$ is in $\mu \ \text{Val}_{\mu,\rho}$ for any value v' .

These admissibility axioms impose natural conditions for properties to be (in principle) enforceable on untrusted code using the object capability approach. Axiom A-PURE models the fact that we cannot prevent untrusted code from returning any valid result value that it chooses. Similarly, Axiom A-BIND means that we cannot prevent untrusted code from composing two expressions that we deem acceptable individually. Finally, Axiom A-PRINT means that we cannot prevent the code from exercising the primitive capabilities that we allow it to access.

We point out that Axiom A-BIND should not be interpreted to mean that we cannot enforce policies that are stateful. For example, when we add mutable state in Section 3.1, a policy that only one output may happen can be modelled by using a flag heap variable and requiring that output can only happen when the flag is 0 and after outputting anything, it should be set to 1.

Scope Our third and final step towards formulating effect parametricity, requires formalising the values that an expression has access to. We use a notion of well-scopedness of expressions. Figure 2 defines contexts Γ as lists of variables and a judgement $\Gamma \vdash e$ expressing that e only uses variables in Γ . We define $[\Gamma]_{\mu,\rho}$ as the set of substitutions γ that map the variables in Γ to acceptable values in $\text{Val}_{\mu,\rho}$:

$$[\Gamma]_{\mu,\rho} \stackrel{\text{def}}{=} \{\gamma \mid \gamma(x) \in \text{Val}_{\mu,\rho} \text{ for all } x \in \Gamma\}.$$

We can now formulate our Fundamental Theorem.

Theorem 1 (Fundamental Theorem for $\lambda_{\text{out},FO}$). *For a valid effect interpretation (μ, ρ) , $\Gamma \vdash e$ implies that for any $\gamma \in [\Gamma]_{\mu,\rho}$, $\gamma(e)$ is in $\mathcal{E}[\mu \ \text{Val}_{\mu,\rho}]$.*

The theorem states that for an arbitrary expression e with access to variables in Γ only (i.e. $\Gamma \vdash e$), instantiating those variables with acceptable values (i.e. $\gamma \in [\Gamma]_{\mu,\rho}$) produces an expression $\gamma(e)$ with acceptable effects and result value: $\gamma(e) \in \mathcal{E}[\mu \ \text{Val}_{\mu,\rho}]$. A proof is in the TR, by structural induction on the well-scopedness judgement $\Gamma \vdash e$.

$$\begin{array}{c}
e ::= \dots \mid \lambda x. e \mid e \ e \quad v ::= \dots \mid \lambda x. e \quad E ::= \dots \mid E \ e \mid v \ E \\
(\lambda x. e_1) v_2 \leftrightarrow e_1[x/v_2] \quad \text{(E-APP)}
\end{array}$$

Figure 3. $\lambda_{\text{out},HO}$, a higher-order calculus with an output capability.

Unfortunately, $\lambda_{\text{out},FO}$ is too simple for the theorem to imply many useful results. Without function values or objects, there is nothing to play the role of an instance object: a value that reacts in a programmer-defined way to some form of (method) invocation. This rules out custom capabilities, and the only result to prove is that code without a reference to primitive capability `out` cannot produce output. It is still instructive to see how that follows.

Example 1. *If $\Gamma \vdash e$, $\gamma(x) \neq \text{out}$ for all x in Γ and $\gamma(e) \rightarrow_o^* v$. Then o must be empty.*

Proof. Instantiate Theorem 1 with effect interpretation $(\mu, \rho) = (IO_{\text{triv}}, \text{Ref}_{\text{triv}})$, where $\text{Ref}_{\text{triv}} \stackrel{\text{def}}{=} \emptyset$ and

$$IO_{\text{triv}} \ P \stackrel{\text{def}}{=} \{\text{cmd} \mid \text{If } \text{cmd} \rightarrow_o^* v \text{ then } o = [] \wedge v \in P\}$$

Ref_{triv} defines that there are *no* acceptable primitive capabilities and $IO_{\text{triv}} \ P$ declares a command acceptable if evaluating it produces *no* output and a result satisfying P . In the TR, we show that $(IO_{\text{triv}}, \text{Ref}_{\text{triv}})$ is a valid effect interpretation. With $\Gamma \vdash e$ and $\gamma \in [\Gamma]_{IO_{\text{triv}}, \text{Ref}_{\text{triv}}}$, Theorem 1 implies that $\gamma(e) \in \mathcal{E}[IO_{\text{triv}} \ \text{Val}_{IO_{\text{triv}}, \text{Ref}_{\text{triv}}}]$. $\gamma(e) \rightarrow_o^* v$ then implies (by a simple lemma) that $o = []$. \square

This result is not useless, but probably easier to prove without our machinery. In a higher-order calculus, *custom* capabilities make the Fundamental Theorem more useful.

2.2. An output capability in a higher-order calculus

To obtain a higher-order calculus, we add standard first-class functions in Figure 3, producing calculus $\lambda_{\text{out},HO}$.

The definition of commands does not change and we reuse command predicates and $\mathcal{E}[P]$ as defined before. An effect interpretation is still a couple (μ, ρ) , with $\mu \in \mathcal{P}(\text{Val}) \rightarrow \mathcal{P}(\text{Cmd})$ and $\rho \in \mathcal{P}(\text{Cap})$. However, the set of acceptable values w.r.t. effect interpretation (μ, ρ) becomes more complicated, since lambdas are values too. We define a command predicate $P_1 \rightarrow P_2$ containing lambdas that map values in predicate P_1 to expressions in $\mathcal{E}[P_2]$.

$$P_1 \rightarrow P_2 \stackrel{\text{def}}{=} \{\lambda x. e \mid \text{For all } v \in P_1, e[x/v] \text{ is in } \mathcal{E}[P_2]\}$$

It is natural to extend $\text{Val}_{\mu,\rho}$ with those lambdas that produce acceptable effects and acceptable results when invoked with an acceptable argument, i.e. lambdas in $\text{Val}_{\mu,\rho} \rightarrow \mu \ \text{Val}_{\mu,\rho}$:

$$\text{Val}_{\mu,\rho} \stackrel{\text{def}}{=} \text{Pure Val} \cup \rho \cup (\text{Val}_{\mu,\rho} \rightarrow \mu \ \text{Val}_{\mu,\rho}) \quad \text{(Oops!)}$$

Unfortunately, the resulting definition is recursive in an unacceptable way. It defines $\text{Val}_{\mu,\rho}$ recursively in terms of

itself and this is mathematically unsound³. Formally, the recursive equation may not have a solution (a value predicate that satisfies the equation). Luckily, we are not the first to encounter this problem. A well-studied technique called *step-indexing* can be used to solve it [12].

The idea is to no longer work with command predicates but use *step-indexed* command predicates instead. These are predicates on couples of natural numbers and commands: $P \in \mathcal{P}(\mathbb{N} \times \text{Cmd})$. The fact that a couple (n, cmd) is in a step-indexed predicate P can be understood as saying that cmd will satisfy P when we inspect it for a maximum of n steps. We say that cmd is n -acceptable in P . For example, for noticing that the expression $e = \text{let}(y = 3) \text{ let}(z = 4) \text{ true}$ produces a boolean, we need to perform 2 computational steps. Therefore, $(0, e)$ and $(1, e)$ could be in $\mathcal{E}[P]$ regardless of P but $(2, e)$ only if true is considered acceptable by P . Typically, we require for step-indexed predicates that they are *uniform*; inspecting expressions for more steps should only make more expressions unacceptable. More formally, we define $UPred(A)$, the set of uniform step-indexed predicates over a set A as follows. For technical reasons, we sometimes also use normal predicates in $Pred(A)$.

$$\begin{aligned} Pred(A) &\stackrel{\text{def}}{=} \mathcal{P}(\mathbb{N} \times A) \\ UPred(A) &\stackrel{\text{def}}{=} \left\{ p \in \mathcal{P}(\mathbb{N} \times A) \mid \right. \\ &\quad \left. (n, e) \in p \text{ and } k \leq n \text{ implies that } (k, e) \in p \right\} \end{aligned}$$

We adapt our definitions to follow suit:

$$\mathcal{E}[P] \stackrel{\text{def}}{=} \left\{ (k, e) \mid \text{For all } \text{cmd} \text{ and } i \leq k \text{ s.t. } e \rightarrow^i \text{cmd}, \right. \\ \left. \text{we have that } (k - i, \text{cmd}) \in P \right\}$$

That is: an expression is k -acceptable in the expression extension of command predicate P if a command that it evaluates to in $i \leq k$ steps is $(k - i)$ -acceptable in P .

Effect interpretations become couples (μ, ρ) with $\mu \in UPred(\text{Val}) \rightarrow Pred(\text{Cmd})$ and $\rho \in UPred(\text{Cap})$ and we can now correct the failed definition of $Val_{\mu, \rho}$:

$$\begin{aligned} P_1 \rightarrow P_2 &\stackrel{\text{def}}{=} \left\{ (n, \lambda x.e) \mid \text{For all } i < n, \text{ and } (i, v) \in P_1, \right. \\ &\quad \left. \text{we have that } (i, e[x/v]) \in \mathcal{E}[P_2] \right\} \\ Val_{\mu, \rho} &\stackrel{\text{def}}{=} (\mathbb{N} \times \text{Pure Val}) \cup \rho \cup (Val_{\mu, \rho} \rightarrow \mu \text{ Val}_{\mu, \rho}) \end{aligned}$$

The predicate of functions $P_1 \rightarrow P_2$ defines as n -acceptable those lambdas that, for $i < n$, can be given an i -acceptable argument to obtain an i -acceptable result.

Thanks to step-indexing, this definition of $Val_{\mu, \rho}$ is mathematically valid. Intuitively, this is because deciding whether or not $\lambda x.e$ is n -accepted by $P_1 \rightarrow P_2$, only requires knowing P_1 and P_2 up to $n - 1$ steps. Hence, whether an expression is n -acceptable in $Val_{\mu, \rho}$ only depends on the $(n - 1)$ -acceptable expressions in $Val_{\mu, \rho}$. The more mathematical story is that with a certain *metric* (a function that defines a “distance” between two predicates), $UPred(A)$ can be seen as a *complete metric space* and the Banach fixpoint theorem guarantees that unique fixpoints exist for *contractive* functions (i.e. applying the function to two predicates

3. The contravariant occurrence of $Val_{\mu, \rho}$ precludes induction.

produces new predicates that are strictly closer together than the original two). The TR explains this more formally.

Valid effect interpretations (μ, ρ) must satisfy natural adaptations of the axioms we saw before.

- A-PURE: For $(n, v) \in P$, (n, v) must be in μP .
- A-BIND: If $(n, \text{cmd}) \in \mu P$ and $(i, E\langle v \rangle) \in \mathcal{E}[\mu P']$ for all $i \leq n$ and values $(i, v) \in P$, then $(n, E\langle \text{cmd} \rangle) \in \mathcal{E}[\mu P']$.
- A-PRINT: If $(n, v) \in Val_{\mu, \rho}$, then $(n, v.\text{print}(v')) \in \mu Val_{\mu, \rho}$.

Context interpretations now become step-indexed as well: $(n, \gamma) \in \llbracket \Gamma \rrbracket_{\mu, \rho}$ iff $(n, \gamma(x)) \in Val_{\mu, \rho}$ for all $x \in \Gamma$. Well-scopedness is easily extended to lambdas and applications:

$$\frac{\Gamma, x \vdash e}{\Gamma \vdash \lambda x.e} \quad \frac{\Gamma \vdash e_1 \quad \Gamma \vdash e_2}{\Gamma \vdash e_1 e_2}$$

We can now state the Fundamental Theorem for $\lambda_{\text{out}, HO}$.

Theorem 2 (Fundamental Theorem for $\lambda_{\text{out}, HO}$). *For a valid effect interpretation (μ, ρ) , $\Gamma \vdash e$ implies for any n and $(n, \gamma) \in \llbracket \Gamma \rrbracket_{\mu, \rho}$ that $(n, \gamma(e))$ is in $\mathcal{E}[\mu Val_{\mu, \rho}]$.*

This theorem naturally adapts the previous theorem to step-indexing and as before, it formalises our intuitive notion of effect parametricity; if an expression e only has access to variables in Γ and these are instantiated by values that produce acceptable effects by effect interpretation (μ, ρ) , then the resulting expression $\gamma(e)$ produces effects acceptable by (μ, ρ) and a result that can only produce acceptable effects.

The lambdas in $\lambda_{\text{out}, HO}$ can be used as custom capabilities. Assume a primitive function toUpperCase that converts strings to uppercase and consider the following expression:

$$\text{upp} \stackrel{\text{def}}{=} \lambda s. \text{out}.\text{print}(\text{toUpperCase}(s))$$

The function upp converts a string to upper case and prints it using primitive capability out . The expression is a custom capability, restricting out to allow printing only uppercase text. But can we prove that this restriction actually holds?

Example 2. *For expression e and variable u , if $\emptyset, u \vdash e$ and $e[u \mapsto \text{upp}] \rightarrow_o^* v$, then the output o is in uppercase.*

Proof. Define effect interpretation $(IO_{\text{upp}}, Ref_{\text{upp}})$:

$$\begin{aligned} Ref_{\text{upp}} &\stackrel{\text{def}}{=} \emptyset \\ IO_{\text{upp}} P &\stackrel{\text{def}}{=} \left\{ (n, \text{cmd}) \mid \text{For } i \leq n, \text{ if } \text{cmd} \rightarrow_o^i v \text{ then } o \text{ is} \right. \\ &\quad \left. \text{in upper case and } (n - i, v) \in P \right\} \end{aligned}$$

Ref_{upp} defines that out is not an acceptable primitive capability (upp is only useful when no direct access to out is available). IO_{upp} defines as n -acceptable those expressions that, when evaluated to a value in $i \leq n$ steps, produce only uppercase output and an $n - i$ -acceptable result. In the TR, we show that $(IO_{\text{upp}}, Ref_{\text{upp}})$ is a valid effect interpretation.

Although out is not in ρ , (n, upp) is in $Val_{IO_{\text{upp}}, Ref_{\text{upp}}}$ for any n , so that $(n, [u \mapsto \text{upp}])$ is in $\llbracket u, \emptyset \rrbracket_{IO_{\text{upp}}, Ref_{\text{upp}}}$. This follows because (n, upp) is in

$$(Val_{IO_{\text{upp}}, Ref_{\text{upp}}} \rightarrow IO_{\text{upp}} Val_{IO_{\text{upp}}, Ref_{\text{upp}}}) \subseteq Val_{IO_{\text{upp}}, Ref_{\text{upp}}}.$$

For $i < n$ and $(i, v) \in Val_{IO_{\text{upp}}, Ref_{\text{upp}}}$, we show

$$(i, \text{out}.\text{print}(\text{toUpperCase}(v))) \in \mathcal{E}[IO_{\text{upp}} Val_{IO_{\text{upp}}, Ref_{\text{upp}}}]$$

So, take $i' \leq i$ and $\text{out.print}(toUpperCase(v)) \rightarrow^{i'} \text{cmd}$. Then i' must be 1 and cmd must be $\text{out.print}(v')$ for v' the uppercased version of v . We show that $\text{out.print}(v')$ is $i-1$ -acceptable in $IO_{upp} \text{ Val}_{IO_{upp}, Ref_{upp}}$, so take $i'' \leq i-1$ and $\text{out.print}(v') \rightarrow_o^{i''} v''$. Then $i'' = 1$, $v'' = \text{undef}$ and $o = [v']$. $(i-2, \text{undef})$ is in $Val_{IO_{upp}, Ref_{upp}}$ and o is in upper case. \square

The Fundamental Theorem then implies the safety of upp . For arbitrary n , (n, upp) is in $Val_{IO_{upp}, Ref_{upp}}$ and thus $(n, [u \mapsto upp])$ is in $[[\emptyset, u]]_{IO_{upp}, Ref_{upp}} \cdot (IO_{upp}, Ref_{upp})$ is a valid effect interpretation, so effect parametricity implies that $(n, e[u \mapsto upp])$ is in $\mathcal{E}[IO_{upp} \text{ Val}_{IO_{upp}, Ref_{upp}}]$. By a simple lemma in the TR, this gives the required result. \square

2.3. Dealing with ambient authority

To understand effect parametricity, it is instructive to consider how our proofs would fail for non-capability-safe languages. Remember, for example, that out is *internal* syntax, not *surface* syntax, i.e. the programmer is not allowed to write it. Formally, our well-scopedness judgement $\Gamma \vdash e$ does not permit e to mention out , so that our Fundamental Theorem does not apply to expressions that do mention out .

If we drop this restriction, i.e. make out part of the surface syntax, all programs gain implicit or *ambient authority* to produce arbitrary output. The custom capability upp no longer works: programs can just use the stronger capability out instead of upp to produce arbitrary output. Formally, if we make out part of the surface syntax by adding the rule $\Gamma \vdash \text{out}$ to the well-scopedness judgement, then the Fundamental Theorem should additionally require (as an axiom) that the effect interpretation considers out acceptable, i.e. $(n, \text{out}) \in \rho$ for any n . This excludes the interpretation used to prove safety of upp , but since upp is no longer safe, this makes sense.

We emphasise however, that the additional ambient authority *can* be accommodated with some changes to our Fundamental Theorem. The benefit is small for $\lambda_{\text{out}, HO}$, since an effect interpretation allowing out cannot impose any restriction on effects at all, but imagine that there were primitive capabilities beside out . With a non-globally accessible net capability for accessing the network, for example, the modified Fundamental Theorem still implies properties about how the network can be accessed and custom capabilities restricting network access can still work, despite the ambient authority of the globally accessible out .

3. Capability-safety in λ_{JS}

In this paper, we present our results for λ_{JS} (although we believe they can be adapted to other object capability languages, both typed and untyped, both low-level and high-level). The higher-order store in λ_{JS} adds significant complexity but makes the results more realistic and interesting.

3.1. LambdaJS

Figure 4 shows the syntax and operational semantics of λ_{JS} , as defined by Guha et al. [14], but omitting exceptions

and object prototypes. The calculus is a fairly standard untyped lambda calculus with numbers, strings, booleans, undef and null values. There are n -ary lambdas as func expressions and string-indexed records with field projection $e[e]$, field update $e[e] = e$, field deletion $\text{delete } e[e]$ and record literals $\{\text{str} : e\}$. The record operations are pure; for example, a field update $r[\text{"fld"}] = 5$ does not modify r , but returns a modified copy. Furthermore, λ_{JS} has mutable references with update ($e_1 = e_2$), allocation ($\text{ref } e$; allocates a memory cell with initial value e) and dereference ($\text{deref } e$) expressions. Finally, there are normal if , sequencing ($e; e$) and while expressions and unspecified primitive operators op_n .

The operational semantics of λ_{JS} is defined in two parts. The primitive pure evaluation judgement $e_1 \hookrightarrow e_2$ defines the evaluation of func expressions, field projections, field updates, field deletions, in addition to rules for let , if , sequencing expressions and while expressions in Figure 1. Primitive operations op_n evaluate in terms of an unspecified δ_n function. The actual small-step pure evaluation judgement $e_1 \rightarrow e_2$ is defined in terms of $e_1 \hookrightarrow e_2$ and evaluation contexts, obtaining a strict, left-to-right evaluation order. Finally, the impure evaluation judgement $(\sigma_1, e_1) \rightarrow (\sigma_2, e_2)$ for stores σ embeds pure evaluations, leaving the store unmodified, and defines the behaviour of allocation ($\text{ref } e$), dereference ($\text{deref } e$) and assignment ($e = e$) expressions.

Like $\lambda_{\text{out}, FO}$ and $\lambda_{\text{out}, HO}$, our subset of λ_{JS} has many stuck terms, e.g., an if branching on a non-boolean or invoking a non- func as a function. The original λ_{JS} produced exceptions in those cases, but we omit those for simplicity.

3.2. The logical relation

Let us now extend our previous results to λ_{JS} . As for $\lambda_{\text{out}, HO}$, we use step-indexing to construct them in a mathematically sound way. However, the higher-order store that is present in λ_{JS} adds significant complexity. It adds new types of authority (accessing and modifying heap data structures, respecting certain invariants or protocols, in arbitrary or restricted ways) but also new ways for components to communicate and pass capabilities to each other.

To support all of this, we construct additional machinery to track assumptions that components have about shared data structures and the authority they have to modify them. We construct a relatively rich type of *Kripke worlds* to model arbitrary invariants and protocols on shared state [13], [16]. Our notations loosely follow Birkedal et al. [15] and we use their recipe for constructing recursive worlds (see below).

Kripke possible worlds

λ_{JS} features a higher-order mutable store, i.e., one can use references into *heap* memory that may contain higher-order data such as functions or objects. It is important that our logical relations are powerful enough to support typical usage of such references, and this is quite a challenge. Often, references into the store are shared between multiple components and correctness of a program relies upon invariants on their contents. For example, rnode in the introduction needs to know that the DOM is a tree, not a graph. Sometimes

$$\begin{array}{l}
l \in Loc \quad c \in Const ::= num \mid str \mid bool \mid \mathbf{undef} \mid null \quad \sigma \in Store ::= (l, v) \cdots \quad \delta_n : op_n \times v_1 \cdots v_n \rightarrow c \\
v \in Val ::= c \mid \mathbf{func}(x \cdots)\{\mathbf{return} e\} \mid \{str : v\} \mid l \\
e \in Expr ::= x \mid v \mid \mathbf{let}(x = e) e \mid e(e \cdots) \mid e[e] \mid e[e] = e \mid \mathbf{delete} e[e] \mid \{str : e\} \mid e = e \mid \mathbf{ref} e \mid \mathbf{deref} e \\
\quad \mid \mathbf{if}(e)\{e\} \mathbf{else} \{e\} \mid e; e \mid \mathbf{while}(e)\{e\} \mid op_n(e_1 \cdots e_n) \\
E ::= \cdot \mid \mathbf{let}(x = E) e \mid E(e \cdots) \mid v(v \cdots E, e \cdots) \mid \{str : v \cdots, str : E, str : e \cdots\} \mid E[e] \mid v[E] \mid E[e] = e \mid v[E] = e \mid v[v] = E \\
\quad \mid \mathbf{delete} E[e] \mid \mathbf{delete} v[E] \mid E = e \mid v = E \mid \mathbf{ref} E \mid \mathbf{deref} E \mid \mathbf{if}(E)\{e\} \mathbf{else} \{e\} \mid E; e \mid op_n(v \cdots E e \cdots) \\
\mathbf{func}(x_1 \cdots x_n)\{\mathbf{return} e\}(v_1 \cdots v_n) \hookrightarrow e[x_1/v_1 \cdots x_n/v_n] \quad (\text{E-APP}) \quad \frac{str_x \notin (str_1 \cdots)}{\{str_1 : v_1 \cdots\}[str_x] = v_x \hookrightarrow \{str_x : v_x, str_1 : v_1 \cdots\}} \quad (\text{E-CREATEFIELD}) \\
\frac{str_x \notin (str_1 \cdots)}{\mathbf{delete} \{str_1 : v_1 \cdots\}[str_x] \hookrightarrow \{str_1 : v_1 \cdots\}} \quad (\text{E-DELETEFIELD-NOTFND}) \quad \frac{str_x \notin (str_1 \cdots str_n)}{\{str_1 : v_1 \cdots str_n : v_n\}[str_x] \hookrightarrow \mathbf{undef}} \quad (\text{E-GETFIELD-NOTFND}) \\
\frac{}{\{\cdots str : v \cdots\}[str] \hookrightarrow v} \quad (\text{E-GETFIELD}) \quad \frac{}{op_n(v_1 \cdots v_n) \hookrightarrow \delta_n(op_n, v_1 \cdots v_n)} \quad (\text{E-PRIM}) \\
\frac{}{\{str_1 : v_1 \cdots str_i : v_i \cdots str_n : v_n\}[str_i] = v \hookrightarrow \{str_1 : v_1 \cdots str_i : v \cdots str_n : v_n\}} \quad (\text{E-UPDATEFIELD}) \\
\frac{}{\mathbf{delete} \{str_1 : v_1 \cdots str_x : v_x \cdots str_n : v_n\}[str_x] \hookrightarrow \{str_1 : v_1 \cdots str_n : v_n\}} \quad (\text{E-DELETEFIELD}) \\
\text{Define } e_1 \rightarrow e_2 \text{ if } e_1 = E(e'_1), e_2 = E(e'_2) \text{ and } e'_1 \hookrightarrow e'_2. \\
\frac{e_1 \hookrightarrow e_2}{(\sigma, E(e_1)) \rightarrow (\sigma, E(e_2))} \quad (\text{E-PURE}) \quad \frac{l \notin \text{dom}(\sigma)}{(\sigma, E(\mathbf{ref} v)) \rightarrow (\sigma[l \mapsto v], E(l))} \quad (\text{E-REF}) \quad \frac{}{(\sigma, E(\mathbf{deref} l)) \rightarrow (\sigma, E(\sigma(l)))} \quad (\text{E-DEREF}) \\
\frac{}{(\sigma, E(l = v)) \rightarrow (\sigma[l \mapsto v], E(v))} \quad (\text{E-SETREF})
\end{array}$$

Figure 4. Syntax and operational semantics of λ_{JS} , following Guha et al. [14], but omitting exceptions and object prototypes. For brevity, we write \mathbf{undef} instead of $\mathbf{undefined}$ and sometimes omit brackets and the \mathbf{return} in \mathbf{funcs} and we write (σ, e) instead of σe for clarity. The figure defines values v , expressions e and evaluation contexts E and the primitive pure evaluation judgement $e \hookrightarrow e$ (not repeating rules E-WHILE, E-LET, E-IFTRUE, E-IFFALSE and E-BEGINDISCARD from Figure 1), the pure evaluation judgement $e \rightarrow e$ and the impure evaluation judgement $(\sigma, e) \rightarrow (\sigma, e)$.

invariants evolve during execution, e.g., an auction object may contain a list of bids that is modifiable, but not after the auction is marked final. Components may also have partial authority over a shared data structure. For example, the ad in the example from the introduction only has the authority to modify its part of the DOM, while other components may modify the whole DOM. Similarly, Section 4.3 studies two isolated components, one of which can only push values on a stack while the other can only pop.

To support such patterns, we use another well-established solution: *Kripke* logical relations. The idea is to index our predicates by *possible worlds* $w \in W$, which model a set of assumptions about (a) the current state of data structures in the store, (b) invariants and protocols that will be respected over them in the future and (c) the authority that is available to modify those data structures. Intuitively, a value v is n -accepted by an indexed predicate $P \in W \rightarrow UPred(Val)$ in a world w (i.e. $(n, v) \in P w$) when inspecting v for n operational steps in stores satisfying w cannot make it break invariants in w or return an invalid result. For values (which can be stored and used later), we will require that they are valid not only in w , but in any possible future evolution of $w' \sqsupseteq w$ (to be defined later).

We define our Kripke worlds as follows:

$$\begin{array}{l}
\text{IslandName} \stackrel{\text{def}}{=} \mathbb{N} \\
W \stackrel{\text{def}}{=} \{w \in \text{IslandName} \hookrightarrow \text{Island} \mid \text{dom}(w) \text{ finite}\} \\
\text{Island} \stackrel{\text{def}}{=} \left\{ \begin{array}{l} \iota = (s, \phi, \phi^{\text{pub}}, H) \mid s \in \text{State} \wedge \phi \subseteq \text{State}^2 \wedge \\ H \in \text{State} \rightarrow \text{StorePred} \wedge \phi^{\text{pub}} \subseteq \phi \wedge \\ \phi, \phi^{\text{pub}} \text{ reflexive and transitive} \end{array} \right\} \\
\text{StorePred} \stackrel{\text{def}}{=} \{\psi \in \hat{W} \rightarrow_{mon, ne} UPred(\text{Store})\} \\
roll : \frac{1}{2} \cdot W \cong \hat{W}
\end{array}$$

Worlds $w \in W$ contain a finite set of components or *islands*

describing disjoint parts of the store. Islands are identified by `IslandNames` which are in fact just natural numbers. An island ι represents an evolvable invariant about a heap data structure, as well as a lower bound on the available authority to modify the data structure. Concretely, $\iota = (s, \phi, \phi^{\text{pub}}, H)$ represents a state machine with current state s and transition relation $\phi \in \text{State}^2$. We assume a fixed set `State` of possible states, assumed to contain all the states used in this paper.

Every island contains a function $H \in \text{State} \rightarrow \text{StorePred}$ that defines when a store satisfies the requirements for a state. These store requirements are modelled as `StorePreds`: predicates on stores that are themselves again world-indexed. The arrow $\rightarrow_{mon, ne}$ means that H must be *monotone* (to be explained later) as well as non-expansive (a sanity requirement w.r.t. step-indexing, see the TR).

Finally, ϕ^{pub} represents an assumption about the available authority to modify the data structure. It is a sub-relation of ϕ , and it represents a subset of state machine transitions for which authority is available to make. For example, in Section 4.2, the ad example from the introduction will be verified using an island whose states are trees of values that represent the current state of the DOM. We will prove that the *ad.initialize* call is valid w.r.t. a world in which ϕ^{pub} only allows modifications in the ad's part of the tree (since the ad may not modify the DOM outside of its node), while ϕ allows arbitrary modifications to the DOM (since the DOM may be modified arbitrarily by other code).

The careful reader may notice that the above definition of worlds is again recursive: worlds contain store predicates, which are themselves indexed by worlds. This kind of recursive worlds is needed to deal with the higher-order nature of λ_{JS} and we use a general method by Birkedal et al. [15] for constructing them. The factor $\frac{1}{2}$, the set \hat{W} and the isomorphism *roll* play a technical role in this construction, but we recommend the casual reader to ignore

them here and elsewhere in the paper.

Future worlds Kripke worlds represent assumptions that a piece of code holds on the rest of the system. However, because values may be stored and used later, any assumptions that they rely on must not be invalidated by legitimate future evolutions of the system. There are essentially three legitimate ways for a system to evolve. (a) Fresh data structures may be allocated, and invariants or protocols may be established for them. Additionally, (b) existing invariants may evolve according to their established protocols and finally (c) additional authority over data structures may become available. These three types of evolutions are modelled by the future world relation; a world w_2 is a future world of w_1 ($w_2 \sqsupseteq w_1$) if w_2 contains at least the islands of w_1 , but potentially more (a). For existing islands, the state in w_2 must be reachable from the state in w_1 using the state machine transitions in ϕ , i.e. according to the data structure's established protocol (b). Finally, the state machine's set of public transitions ϕ^{pub} is allowed to grow (but not beyond its complete set of transitions ϕ), representing an increase in available authority (c).

$$w_2 \sqsupseteq w_1 \text{ iff } \text{dom}(w_2) \supseteq \text{dom}(w_1) \wedge \forall j \in \text{dom}(w_1). w_2(j) \sqsupseteq w_1(j)$$

$$(s_2, \phi_2, \phi_2^{\text{pub}}, H_2) \sqsupseteq (s_1, \phi_1, \phi_1^{\text{pub}}, H_1) \text{ iff}$$

$$(\phi_2, H_2) = (\phi_1, H_1) \text{ and } \phi_1 \sqsupseteq \phi_2^{\text{pub}} \supseteq \phi_1^{\text{pub}} \text{ and } (s_1, s_2) \in \phi_1$$

The requirement that values that are valid in a predicate remain valid in legitimate future evolutions of a system is captured by the required *monotonicity* of predicates in $W \rightarrow_{\text{mon}, ne} \text{UPred}(A)$ (in addition to non-expansiveness, a sanity requirement w.r.t. step-indexing). Formally, monotonicity requires that $P w_2 \sqsupseteq P w_1$ whenever $w_2 \sqsupseteq w_1$, i.e. values/stores/... that are valid in a world w_1 must remain valid in future worlds $w_2 \sqsupseteq w_1$. In what follows, predicates on values and stores will typically need to be monotone (because they contain values that may be stored and used later) while predicates on commands and expressions typically need not be, as they cannot be stored (except as part of a `func` value).

The future world relation $w_2 \sqsupseteq w_1$ represents future evolutions of a world that a piece of code should be able to cope with. However, it is not necessarily allowed to make those changes itself. A more restricted *public* future-world relation $w_2 \sqsupseteq^{\text{pub}} w_1$ defines the set of future worlds that one has the authority to transition to:

$$w_2 \sqsupseteq^{\text{pub}} w_1 \text{ iff } \left\{ \begin{array}{l} \text{dom}(w_2) \supseteq \text{dom}(w_1) \wedge \\ \forall j \in \text{dom}(w_1). w_2(j) \sqsupseteq^{\text{pub}} w_1(j) \end{array} \right.$$

$$(s_2, \phi_2, \phi_2^{\text{pub}}, H_2) \sqsupseteq^{\text{pub}} (s_1, \phi_1, \phi_1^{\text{pub}}, H_1) \text{ iff}$$

$$(\phi_2, \phi_2^{\text{pub}}, H_2) = (\phi_1, \phi_1^{\text{pub}}, H_1) \text{ and } (s_1, s_2) \in \phi_1^{\text{pub}}$$

A public future world $w_2 \sqsupseteq^{\text{pub}} w_1$ must also contain at least the islands of w_1 and may contain additional ones. However, the new state of islands in w_2 must now be reachable through the state machine's public transitions in ϕ_1^{pub} , i.e. transitions that can be made using the available authority. Finally, ϕ^{pub} is not allowed to grow in public future worlds, i.e. one does not have the authority to increase one's own authority. Note

that the latter precludes capabilities appearing out of thin air, somewhat similar to the *No Authority Amplification* property that we will discuss in Section 5.

Disjoint worlds (with assumptions about disjoint data structures in the heap) can be combined with the \oplus operator:

$$w \oplus w' = w'' \text{ iff } \text{dom}(w'') = \text{dom}(w) \uplus \text{dom}(w') \wedge \forall j \in \text{dom}(w).$$

$$w''(j) = w(j) \wedge \forall j \in \text{dom}(w'). w''(j) = w'(j)$$

Finally, a store σ n -satisfies a world w if it can be partitioned into parts that n -satisfy the store predicate for the current state of all islands (we write $\iota.H$ for projecting out the store predicate of an island):

$$\sigma :_n w \text{ iff } \left\{ \begin{array}{l} \exists \sigma_j. \sigma = \uplus_{j \in \text{dom}(w)} \sigma_j \text{ and } \forall j \in \text{dom}(w), \\ \forall n' < n. (n', \sigma_j) \in w(j). H(w(j).s) \text{ (roll } w) \end{array} \right.$$

It is worthwhile at this point to take a step back and build an intuitive understanding for our worlds. Generally, they should be interpreted as a set of assumptions with respect to which a value or expression is valid. An expression e that holds capabilities for making certain modifications to a shared data structure, will only be valid w.r.t. a Kripke world with an island governing the data structure. Private transitions for the island would include all possible modifications to the data structure that *any* subject in the system has the authority to make, so that e will be required to tolerate such modifications. Public transitions for the island will include at least the modifications that e itself has the authority to make. Using security terminology, this perspective means that any expression may be seen as a subject and the world in which it is valid as an upper bound on its required authority. Entities in our approach are simply represented by code executing on their behalf.

Command predicates As for $\lambda_{\text{out}, FO}$ and $\lambda_{\text{out}, HO}$, we define a syntactic class of *commands*: either values or expressions blocked on an impure operation.

$$\text{cmd}_0 ::= \text{deref } v \mid v = v \mid \text{ref } v \quad \text{cmd} \in \text{Cmd} ::= E(\text{cmd}_0) \mid v$$

Predicates on commands P are again extended to expressions as $\mathcal{E}[P]$, which closes P under pure evaluation:

$$\mathcal{E}[P] : W \rightarrow_{ne} \text{Pred}(\text{Expr})$$

$$\mathcal{E}[P] w \stackrel{\text{def}}{=} \left\{ (n, e) \mid \begin{array}{l} \text{for all } i \leq n, e'. \text{ if } e \rightarrow^i \text{cmd}, \\ \text{then } (n - i, \text{cmd}) \in P w \end{array} \right\}$$

A few predicates are used as building blocks further on:

$$\text{Cnst} : W \rightarrow_{\text{mon}, ne} \text{UPred}(\text{Val})$$

$$\text{Cnst } w \stackrel{\text{def}}{=} \mathbb{N} \times \text{Const}$$

$$\{P\} : W \rightarrow_{\text{mon}, ne} \text{UPred}(\text{Val})$$

$$\{P\} w \stackrel{\text{def}}{=} \left\{ (n, \overline{\{str : v\}}) \mid \text{for all } i < n, \overline{(i, v)} \in P w \right\}$$

$$P \cup P' : W \rightarrow_{\text{mon}, ne} \text{UPred}(\text{Val})$$

$$(P \cup P') w \stackrel{\text{def}}{=} P w \cup P' w$$

$$([P] \rightarrow P') : W \rightarrow_{\text{mon}, ne} \text{UPred}(\text{Val})$$

$$([P] \rightarrow P') w \stackrel{\text{def}}{=} \left\{ (n, \text{func}(x_1 \dots x_k) \{ \text{return } e \}) \mid \right.$$

$$\left. \text{for all } v_1 \dots v_k, w' \sqsupseteq w, i < n. \overline{(i, v_j)} \in P w' \Rightarrow \right.$$

$$\left. (i, e[x_1/v_1, \dots, x_n/v_k]) \in \mathcal{E}[P'] w' \right\}$$

$Cnst$ n -accepts all constant values. For value predicates P and P' and command predicate P'' , $\{P\}$ accepts records with P -acceptable fields and $P \cup P'$ accepts values from P or P' . $[P] \rightarrow P''$ n -accepts func expressions producing i -acceptable results in P'' when applied to i -acceptable arguments in P , in any future world $w' \sqsupseteq w$ and for any $i < n$. As a technical detail, the quantification over $w' \sqsupseteq w$ and $i < n$ makes $[P] \rightarrow P''$ monotone and uniform even when P'' is not.

Effect interpretations and acceptable values Like for $\lambda_{\text{out},FO}$ and $\lambda_{\text{out},HO}$, we parameterise our LR over an effect interpretation (μ, ρ) with $\rho : W \rightarrow_{\text{mon},ne} \text{UPred}(\text{Loc})$ a predicate of references that are valid in a given world and μ a function that maps a predicate on values to a predicate on commands:

$$\mu : (W \rightarrow_{\text{mon},ne} \text{UPred}(\text{Val})) \rightarrow_{ne} (W \rightarrow_{ne} \text{Pred}(\text{Cmd})).$$

As before, μP accepts commands producing acceptable effects and results acceptable by a value predicate P .

Given an effect interpretation (μ, ρ) that defines acceptable references and effectful expressions in a given world, we define a predicate $\text{JSVal}_{\mu,\rho}$ of all acceptable λ_{JS} values⁴:

$$\begin{aligned} \text{JSVal}_{\mu,\rho} : W \rightarrow_{\text{mon},ne} \text{UPred}(\text{Val}) \\ \text{JSVal}_{\mu,\rho} \stackrel{\text{def}}{=} Cnst \cup \rho \cup \{\text{JSVal}_{\mu,\rho}\} \cup ([\text{JSVal}_{\mu,\rho}] \rightarrow \mu \text{JSVal}_{\mu,\rho}) \end{aligned}$$

The predicate accepts constants, references in ρ , records of acceptable values and functions mapping acceptable values to expressions with acceptable effects and results.

The following axioms are required to hold for a *valid* effect interpretation.

- A-PURE: If $(n, v) \in P w$ then $(n, v) \in \mu P w$
- A-BIND: If $(n, \text{cmd}) \in \mu P w$ and $(n', E\langle v \rangle) \in \mathcal{E}[\mu P'] w'$ for all $n' \leq n$, $w' \sqsupseteq w$ and $(n', v) \in P w'$, then $(n, E\langle \text{cmd} \rangle) \in \mathcal{E}[\mu P'] w$.
- A-ASSIGN: If $(n, v_1) \in \text{JSVal}_{\mu,\rho} w$ and $(n, v_2) \in \text{JSVal}_{\mu,\rho} w$, then $(n, v_1 = v_2) \in \mu \text{JSVal}_{\mu,\rho} w$.
- A-DEREF: If $(n, v) \in \text{JSVal}_{\mu,\rho} w$, $(n, \text{deref } v)$ must be in $\mu \text{JSVal}_{\mu,\rho} w$.
- A-REF: If $(n, v) \in \text{JSVal}_{\mu,\rho} w$, then $(n, \text{ref } v) \in \mu \text{JSVal}_{\mu,\rho} w$.

Axioms A-PURE and A-BIND are as before (except for the quantification over Kripke worlds) and we don't re-explain them for brevity. Axioms A-ASSIGN and A-DEREF essentially require a compatibility between μ and ρ . They require that if a value is accepted by $\text{JSVal}_{\mu,\rho}$ (e.g. references in ρ), then dereferencing it or assigning an acceptable value must be accepted by μ . Finally, Axiom A-REF requires that allocating a new mutable reference (a primitive effect left unrestricted by the language) with an acceptable initial value must be accepted by the effect interpretation.

Fundamental Theorem We conclude this section with the Fundamental Theorem for λ_{JS} : the formal statement of its capability-safety. As before, it informally states that well-formed λ_{JS} terms respect the restrictions on effects imposed

4. Readers with a background in logical relations can see this as the semantic interpretation of the *unitype* of λ_{JS} values.

by a valid effect interpretation, and now additionally that they respect the invariants and protocols of a Kripke world.

The well-scopedness judgement $\Gamma; \Sigma \vdash e$ states that e is syntactically well-formed in context Γ (a list of free variables) and *store shape* Σ (a list of allocated references). Its definition is unsurprising and relegated to the TR. The predicate $[\Sigma]_{\mu,\rho}$ accepts worlds in which the references in Σ are accepted by ρ , and the predicate $[\Gamma]_{\mu,\rho} w$ accepts substitutions for Γ with acceptable values:

$$\begin{aligned} [\Sigma]_{\mu,\rho} : \text{UPred}(W) \\ [\Sigma]_{\mu,\rho} \stackrel{\text{def}}{=} \{(n, w) \mid \text{for all } l \in \Sigma. (n, l) \in \rho w\} \\ [\Gamma]_{\mu,\rho} w : \text{UPred}(\text{Val}^\Gamma) \\ [\Gamma]_{\mu,\rho} w \stackrel{\text{def}}{=} \{(n, \gamma) \mid \forall x \in \Gamma. (n, \gamma(x)) \in \text{JSVal}_{\mu,\rho} w\} \end{aligned}$$

Theorem 3 (Fundamental Theorem for λ_{JS}). *If $\Gamma, \Sigma \vdash e$ then for a valid effect interpretation (μ, ρ) and for all n, γ and w with $(n, w) \in [\Sigma]_{\mu,\rho}$ and $(n, \gamma) \in [\Gamma]_{\mu,\rho} w$, we have that $(n, \gamma(e))$ must be in $\mathcal{E}[\mu \text{JSVal}_{\mu,\rho}] w$.*

The theorem states that substituting acceptable values for an expression's free variables and considering it in a world where its references are acceptable by effect interpretation (μ, ρ) , produces a μ -acceptable expression with a result in $\text{JSVal}_{\mu,\rho}$. The TR contains a proof by induction on the well-scopedness judgement.

Note that our Fundamental Theorem does not offer termination guarantees about untrusted code. This limitation follows from the higher-order untyped language, in which untrusted code cannot be prevented from diverging.

4. Local state abstraction

A special feature of our logical relation is the quantification over effect interpretations. It allows proving properties about primitive effects, as we saw in Section 2 and we will use it again in Section 6. However, in many cases, standard encapsulation of local state (e.g. instance variables of objects) is all we need. In this section, we define an effect interpretation $(IO^{\text{std}}, Ref^{\text{std}})$ for such “standard” reasoning about local state abstraction.

By the definition in Figure 5, an expression is n -accepted by $IO^{\text{std}} P w$ if (1) it is n -accepted by $P w$ if it is already a value and (2) evaluating it to a value in $0 < i \leq n$ steps in a store σ_r that is n -accepted by world w implies that the resulting store σ_r' is $n - i$ -accepted by a public future world $w' \sqsupseteq^{\text{pub}} w$ and the resulting value is $n - i$ accepted by $P w''$. The original store is in fact allowed to additionally contain a *frame* part σ_f which must not be modified by the evaluation. This definition corresponds roughly to what one would typically find in the \mathcal{E} relation of a Kripke logical relation, except that it only provides guarantees after reduction to a value, i.e. the evaluation is allowed to get stuck, and when it does, nothing is guaranteed about σ' and the resulting expression.⁵

5. This is appropriate for our untyped setting with stuck terms, but unusual as Kripke logical relations are most often used for static type systems that rule out stuck terms.

$$\begin{aligned}
IO^{std} &: (W \rightarrow_{mon,ne} UPred(Val)) \rightarrow_{ne} (W \rightarrow_{ne} Pred(Cmd)) \\
IO^{std} P w &\stackrel{\text{def}}{=} \left\{ \begin{array}{l} (n, cmd) \in P w \text{ if } cmd \in Val \text{ and,} \\ \text{for all } \sigma_r, \sigma_f, \sigma', 0 < i \leq n, v, \sigma_r :_n w \wedge \\ (\sigma_r \uplus \sigma_f, cmd) \rightarrow^i (\sigma', v) \Rightarrow \exists \sigma'_r, w' \sqsupseteq^{\text{pub}} w. \\ \sigma' = \sigma'_r \uplus \sigma_f \wedge \sigma'_r :_{n-i} w' \wedge (n-i, v) \in P w' \end{array} \right\} \\
\iota_l^{std} &\stackrel{\text{def}}{=} (l, =, =, H^{std}) \\
H^{std} l w &\stackrel{\text{def}}{=} \left\{ (n, \{l \mapsto v\}) \mid \begin{array}{l} n = 0 \text{ OR } (n-1, v) \in \\ \text{JSVal}_{IO^{std}, Ref^{std}}(\text{roll}^{-1} w) \end{array} \right\} \\
Ref^{std} &: W \rightarrow_{mon,ne} UPred(Loc) \\
Ref^{std} w &\stackrel{\text{def}}{=} \{(n, l) \mid \exists j. w(j) =_{n+1} \iota_l^{std}\}
\end{aligned}$$

Figure 5. An effect interpretation capturing standard local state abstraction.

We instantiate ρ to Ref^{std} , defined in Figure 5 w.r.t. an island ι_l^{std} . The island ι_l^{std} takes ownership of a location l and requires that its value satisfies $\text{JSVal}_{IO^{std}, Ref^{std}}$. Ref^{std} defines as n -acceptable all locations l that are owned by such an island ι_l^{std} , or at least one that $n+1$ -approximates it. The TR shows that effect interpretation (IO^{std}, Ref^{std}) is valid.

Instantiating our Fundamental Theorem with the effect interpretation (IO^{std}, Ref^{std}) produces a logical relation that captures the encapsulation of local state in λ_{JS} and can be used to reason about non-trivial capability patterns. The next sections demonstrate interesting examples of this: (1) a ticket dispenser function featuring a non-trivial protocol on private state, (2) the ad example from the introduction featuring a capability with restricted authority on shared data (the DOM) and (3) two isolated components with different authority on a shared LIFO communication channel.

4.1. Invariants and Protocols

Consider the following expression:

$$\begin{aligned}
\text{ticketDispenser} &\stackrel{\text{def}}{=} \text{func}(\text{attacker}) \\
&\left\{ \begin{array}{l} \text{let}(o = \text{ref } 0) \\ \text{let}(\text{dispTkt} = \text{func}(\{\text{let}(v = \text{deref } o)\{o := v + 2; v\}\}) \\ \text{attacker}(\text{dispTkt}); \text{deref } o \end{array} \right\}
\end{aligned}$$

The expression takes an (untrusted) function argument attacker . The code allocates a new mutable reference o , initially 0 and constructs a function dispTkt . When called, the function increases o 's value by 2 and returns the old value. If we assume for simplicity an infinite range of primitive integers, dispTkt respects a protocol in its usage of o : its value will always remain even and will only ever *increase*. The attacker code is invoked and receives access to dispTkt . After the attacker code returns, the value of o is returned. Since the attacker has access to dispTkt but not o and by inspecting dispTkt , we can expect the following to hold:

Lemma 1. *Take a store σ , $\Sigma = \text{dom}(\sigma)$ and a value $\emptyset; \Sigma \vdash \text{attacker}$. If $(\sigma, \text{ticketDispenser } \text{attacker}) \rightarrow (\sigma', v)$, then v is even and ≥ 0 .*

With our Fundamental Theorem and effect interpretation (IO^{std}, Ref^{std}) , we can prove that this holds.

Proof sketch. If the evaluation terminates, then for some $l \notin \Sigma$, it must factor as follows (omitting the body of dispTkt for brevity) with $v = \sigma'(l)$:

$$\begin{aligned}
(\sigma, \text{ticketDispenser } \text{attacker}) &\rightarrow^* \\
(\sigma[l \mapsto 0], \text{attacker } (\text{func}(\{\dots\}); \text{deref } l)) &\rightarrow^* \\
(\sigma', (v'; \text{deref } l)) &\rightarrow^* (\sigma', \sigma'(l))
\end{aligned}$$

Define a world w with one island ι_l^{std} for every $l \in \Sigma$. Then for any n and $w' \sqsupseteq w$, we have $(n, w') \in \llbracket \Sigma \rrbracket_{IO^{std}, Ref^{std}}$. The next island $\iota_{\text{tk}, l, k}$ captures l 's intended protocol:

$$\begin{aligned}
\iota_{\text{tk}, l, k} &\stackrel{\text{def}}{=} ((l, k), \sqsubseteq_{\text{tk}}, \sqsubseteq_{\text{tk}}, H_{\text{tk}}) \quad \text{for } k \text{ even} \\
(l, k) \sqsubseteq_{\text{tk}} (l', k') &\text{ iff } l = l' \wedge k' \geq k \wedge k', k \text{ even} \\
H_{\text{tk}}(l, k) w &\stackrel{\text{def}}{=} \{(n, \{l \mapsto k\}) \mid n \in \mathbb{N}\}
\end{aligned}$$

Define $w' = w[j \mapsto \iota_{\text{tk}, l, 0}]$ for $j \notin \text{dom}(w)$. Clearly, $w' \sqsupseteq w$.

We show in the TR that for any n :

$$\begin{aligned}
(n, \text{func}(\{\text{return } (\text{let}(v = \text{deref } l)\{l := v + 2; v\}\})) &\in \\
([\text{JSVal}_{IO^{std}, Ref^{std}}] \rightarrow IO^{std} \text{ JSVal}_{IO^{std}, Ref^{std}}) w' &\subseteq \\
\text{JSVal}_{IO^{std}, Ref^{std}} w' &
\end{aligned}$$

The proof takes about two paragraphs. Essentially, it shows that executing the function in a store satisfying a world $w'' \sqsupseteq w'$ produces a new store satisfying a world $w''' \sqsupseteq^{\text{pub}} w''$. Specifically, the island $w''(j) = \iota_{\text{tk}, l, k}$ will take a public transition to $w'''(j) = \iota_{\text{tk}, l, k+2}$. The function's result value in such a store is a number, so satisfies $\text{JSVal}_{IO^{std}, Ref^{std}}$.

By the Fundamental Theorem, attacker n -satisfies $\mathcal{E}[IO^{std} \text{ JSVal}_{IO^{std}, Ref^{std}}] w'$. Then, $\text{attacker } (\text{func}(\{\dots\}))$ must also n -satisfy $\mathcal{E}[IO^{std} \text{ JSVal}_{IO^{std}, Ref^{std}}] w'$ (by a standard lemma) and we can deduce that σ' must n -satisfy some $w''' \sqsupseteq^{\text{pub}} w'$, so that $\sigma'(l)$ must be even and ≥ 0 . \square

4.2. Restricted capabilities

Another important object capability pattern is to give a component *restricted* access to a resource, while other code keeps full authority over it. This can be implemented by giving the component access to a trusted object that has full access to the resource, but whose methods allow only restricted interaction with it. For reasoning about such a component with restricted authority, we use a world that carries two separate protocols for how the shared state can evolve. A first protocol dictates changes to the shared resource that the component itself has the authority to make while the second specifies changes that other code (with potentially greater authority) may make, and which the component under scrutiny should be able to deal with. These two protocols are given by, respectively, the public

and private transitions in an island: a component is itself allowed to make public transitions, but must be able to cope with private transitions made by other code. In this section, we demonstrate this for the example from the introduction which restricts the authority of an untrusted advertisement in a client-side web page.

We repeat the example web page initialisation function:

$$\text{initWebPage} \stackrel{\text{def}}{=} \text{func}(\text{document}, \text{ad}) \left\{ \begin{array}{l} \text{document.setProp}(\text{"someProperty"}, 42) \\ \text{let } (\text{adNode} = \text{document.addChild}(\text{"ad_div"})) \\ \text{let } (\text{rAdNode} = \text{rnode}(\text{adNode}, 0)) \\ \text{ad.initialize}(\text{rAdNode}) \\ \text{document.getProp}(\text{"someProperty"}) == 42 \end{array} \right.$$

The idea is that *initWebPage* receives access to a *document* object that represents the entire web page and carries the authority to modify it. It uses a function *rnode* to construct a restricted capability *rAdNode* for accessing and modifying only the ad’s part of the page and passes that to the untrusted ad’s initialisation code. We do not repeat the definition of *rnode*, which constructs an object that forwards method invocations to the underlying DOM node but returns `null` when asked for a node outside the ad’s turf. If everything works correctly, we should be able to prove that if *initWebPage* terminates, it must return `true`.

To formalise our assumptions about the behaviour of *document* and its child nodes, we use a form of trees defined by the following grammar:

$$\text{tree} ::= v \mid (\text{id} \mapsto \text{tree})^* \quad \text{id} \in \text{String}$$

We define a notion of plugging a subtree in a tree at a certain path (a list of ids) as follows:

$$\begin{aligned} t'[\bar{p} \mapsto t] &\stackrel{\text{def}}{=} t \\ (\text{id}_1 \mapsto c_1, \dots, \text{id}_n \mapsto c_n)[\bar{p}, \bar{p}] \mapsto t &\stackrel{\text{def}}{=} \\ \left\{ \begin{array}{ll} (\text{id}_1 \mapsto c_1 \dots \text{id}_{j-1} \mapsto c_{j-1}, & \text{if } p = \text{id}_j \wedge \\ \text{id}_j \mapsto t', \text{id}_{j+1} \mapsto c_{j+1} \dots \text{id}_n \mapsto c_n) & t' = c_j[\bar{p} \mapsto t] \\ \text{undefined} & \text{otherwise} \end{array} \right. \end{aligned}$$

We define an island $\iota_{l, \text{tree}, P}^{\text{dom}}$ to govern the state of the DOM. It is parameterised by a function $P \in \text{String}^* \rightarrow W \rightarrow \text{mon}, \text{ne} \text{UPred}(\text{Val})$, which defines, for every path in the DOM, a predicate that the DOM property at that path should satisfy.

$$\begin{aligned} \iota_{l, \text{tree}, P}^{\text{dom}} &\stackrel{\text{def}}{=} ((l, \text{tree}), \sqsubseteq^{\text{dom}}, \sqsupseteq^{\text{dom}}, H_P^{\text{dom}}) \\ (l', \text{tree}') \sqsupseteq^{\text{dom}} (l, \text{tree}) &\text{ iff } l = l' \end{aligned}$$

The store predicate H_P^{dom} is defined in the TR such that $H_P^{\text{dom}}(l, t)$ *w* accepts stores containing a representation of DOM tree *t* and for every property in the tree at path *p*, the value satisfies $P \ p \ w$. Note that the island’s transition relation \sqsupseteq^{dom} does not restrict the evolution of the tree.

For our ad example, we define a restricted transition relation $\sqsupseteq_p^{\text{r-dom}}$ expressing the restricted authority of the ad, i.e. only allowing changes to the DOM under path *p*:

$$\begin{aligned} (l_1, t_1) \sqsupseteq_p^{\text{r-dom}} (l_2, t_2) &\text{ iff} \\ l_1 = l_2 \wedge (\forall t'_1, t'_f. t_1 = t_f[p \mapsto t'_1] \Rightarrow \exists t'_2. t_2 = t_f[p \mapsto t'_2]) \end{aligned}$$

Lemma 2. *Assume that document’s methods getChild, parent, getProp, setProp, addChild and delChild behave in the “obvious” way (to be defined formally in the TR) in stores that contain the representation of a state of the DOM. If $w(j) = \iota_{l, t, P}^{\text{dom}}$ for some *j*, *t* and *P*, $\sigma :_n w$, and *ad* is a closed expression and $(\sigma, \text{initWebPage}(\text{document}, \text{ad})) \rightarrow^i (\sigma', v)$ for $i \leq n$, then $v = \text{true}$.*

The proof of this lemma is essentially based on the restricted transition relation $\sqsupseteq^{\text{r-dom}}$ mentioned above. However, it is complicated by the fact that DOM properties can be higher-order, i.e. functions that carry and use capabilities themselves. In the proof, we have to express that DOM properties under “*ad_div*” have the same authority restriction as the ad. But what about the authority of trusted code (including DOM properties outside “*ad_div*”)? Naturally, this other code may modify the rest of the DOM, but can it modify the tree under “*ad_div*”? Crucially, the code must not (by accident or malice) store capabilities with authority greater than the ad’s (e.g. *document* itself) inside the ad’s reach. The most general solution is to require that they do not do this, i.e. that they preserve the authority bound of the ad. However, because *initWebPage* terminates after the *ad.initialize* call, and no trusted code gets a chance to run, this requirement is not in fact necessary for the lemma above. The proof in the TR relies on this simplification and does not restrict the DOM properties outside the ad’s territory. However, we think it can be generalised if needed.

4.3. Isolated but communicating components

Another interesting pattern is similar to the previous example, but features *multiple* untrusted components that have restricted and *different* authority to a shared resource. In such a scenario, it is necessary to prevent the resource from being used as a communication channel for passing capabilities to the other component that it does not hold.

Consider a mashup page that embeds two disjoint pieces of untrusted code: *attacker₁* and *attacker₂*. The mashup sets up a restricted communication channel: a stack that *attacker₁* and *attacker₂* can respectively push to and pop from. Figure 6 shows what the mashup’s code could look like. If we suppose that *attacker₁* and *attacker₂* do not share a communication channel to begin with, then evaluating *mashup* in an arbitrary heap should always produce a non-negative number. Note that this example only works if the stack rejects non-constant values: if not, *attacker₁* could push the *push* function itself on the stack, for *attacker₂* to retrieve and use for stepping outside its pop-only authority.

Formally, the expected result is as follows:

Lemma 3. *If $(n, \text{attacker}_i)$ are in $\text{JSVal}_{\text{IOstd}, \text{Refstd}} w_i$ for $i = 1, 2$ and disjoint worlds w_1 and w_2 and if a store σ *n*-satisfies $w_1 \oplus w_2$ for *n* sufficiently large, then executing $\text{mashup}(\text{attacker}_1, \text{attacker}_2)$ in σ , produces a result ≥ 0 .*

The premise that the *attacker_i* are in $\text{JSVal}_{\text{IOstd}, \text{Refstd}}$ for disjoint worlds expresses the informal requirement that they do not share a communication channel to begin with.

$$\begin{array}{l}
\text{mashup} \stackrel{\text{def}}{=} \text{func}(\text{attacker}_1, \text{attacker}_2) \\
\left. \begin{array}{l}
\text{let } (stk = \text{ref null}) \\
\text{let } (push = \text{func}(v) \\
\left. \begin{array}{l}
\text{if } \neg \text{isConstant}(v) \text{ then } \text{undef} \text{ else} \\
\text{stk} = \text{ref}(\{val = v, rest = \text{deref } stk\}); \text{undef}
\end{array} \right\} \\
\text{let } (pop = \text{func}() \\
\left. \begin{array}{l}
\text{let } (top = \text{deref } stk) \\
\text{if } (top == \text{null}) \{ \text{undef} \} \\
\text{else } \{ \text{stk} = (\text{deref } top).rest; (\text{deref } top).val \}
\end{array} \right\} \\
\text{let } (size = \text{func}() \\
\left. \begin{array}{l}
\text{let } (c = \text{ref } 0) \\
\text{let } (top = \text{ref } (\text{deref } stk)) \\
\text{while } ((\text{deref } top) \neq \text{null}) \{ \\
\quad c = \text{deref } c + 1; \text{top} = (\text{deref } top).rest \\
\} \\
\text{deref } c
\end{array} \right\} \\
\text{attacker}_1(push) \\
\text{let } (s_1 = \text{size}()) \\
\text{attacker}_2(pop) \\
s_1 - \text{size}()
\end{array} \right\}
\end{array}$$

Figure 6. A mashup application embedding two untrusted components, isolated from each other but with a restricted communication channel.

If they did, then attacker_1 could pass the push capability to attacker_2 and break our result. In practice, the requirement should be easily provable thanks to the absence of global mutable state in our capability-safe language. According to the Fundamental Theorem, the premise is satisfied, for example, if we know that the two pieces of code are well-formed in an empty context and empty store typing.

Proof sketch. In this proof sketch, all statements should be interpreted with respect to a sufficiently large step-index n .

We first define an island $\iota_{(s, \bar{v})}^{stack}$ capturing the stack's invariant; its states are the list of values in the stack and its private transition relation \sqsupseteq^{stack} allows arbitrary modifications. We also define two restricted transition relations $\sqsupseteq^{stack\uparrow}$ and $\sqsupseteq^{stack\downarrow}$ that allow the stack to only grow or shrink respectively:

$$\begin{aligned}
\iota_{(s, \bar{v})}^{stack} &\stackrel{\text{def}}{=} ((s, \bar{v}), \sqsupseteq^{stack}, \sqsupseteq^{stack}, H^{stack}) \\
(l', \bar{v}') &\sqsupseteq^{stack} (l, \bar{v}) \text{ iff } l = l' \\
(l', \bar{v}') &\sqsupseteq^{stack\uparrow} (l, \bar{v}) \text{ iff } l = l' \wedge \exists \bar{v}'' . \bar{v}' = \bar{v}'' \bar{v} \\
(l', \bar{v}') &\sqsupseteq^{stack\downarrow} (l, \bar{v}) \text{ iff } l = l' \wedge \exists \bar{v}'' . \bar{v} = \bar{v}'' \bar{v}' \\
H^{stack} &(l, (v_1 \dots v_n)) \ w \stackrel{\text{def}}{=} \\
&\left\{ (n, \sigma) \mid \begin{array}{l} \exists l_1, \dots, l_n, l_{n+1}, v_1 \dots v_n \in \text{Const}. \sigma(l) = l_1 \wedge \\ \text{dom}(\sigma) = \{l, l_1, \dots, l_n\} \wedge l_{n+1} = \text{null} \wedge \\ \forall i \in 1..n. \sigma(l_i) = \{val = v_i, rest = l_{i+1}\} \end{array} \right\}
\end{aligned}$$

After allocating a location l for stk , we define a world w_3 , with $\text{dom}(w_3) = \{j\}$, $w_3(j) = \iota_{(l, \cdot)}^{stack}$ for some j and the empty list \cdot . Next, we show that push is in $\text{JSval}_{IO^{std}, Ref^{std}}$ for world $w_3[j, \phi^{\text{pub}} \mapsto \sqsupseteq^{stack\uparrow}]$ and pop in $\text{JSval}_{IO^{std}, Ref^{std}}$ for world $w_3[j, \phi^{\text{pub}} \mapsto \sqsupseteq^{stack\downarrow}]$. We can derive that $\text{attacker}_1(\text{push})$ is accepted by $\mathcal{E}[IO^{std} \text{ JSval}_{IO^{std}, Ref^{std}}]$

in world $w_4 \stackrel{\text{def}}{=} w_1 \oplus (w_3[j, \phi^{\text{pub}} \mapsto \sqsupseteq^{stack\uparrow}])$ and we get that the resulting store is valid in a world $w'_1 \oplus w'_3$ with $\text{dom}(w'_3) = \{j\}$ and $\text{dom}(w'_1) \supseteq \text{dom}(w_1)$. From the fact that $w'_1 \oplus w'_3 \sqsupseteq^{\text{pub}} w_4$, we know that the stack can only have grown. We can also derive that $\text{attacker}_2(\text{pop})$ is accepted by $\mathcal{E}[IO^{std} \text{ JSval}_{IO^{std}, Ref^{std}}]$ in a world $w_5 \stackrel{\text{def}}{=} w_2 \oplus (w'_3[j, \phi^{\text{pub}} \mapsto \sqsupseteq^{stack\downarrow}])$ and notice that the part of the store satisfying $H^{stack}(w'_3(j).s)$ ($w'_1 \oplus w'_3$) will also satisfy $H^{stack}(w'_3(j).s)$ w_5 because (crucially) the definition of H^{stack} ignores its world argument, which it can because the content of the stack is restricted to first-order values. We then obtain a resulting world $w'_2 \oplus w''_3$ with $\text{dom}(w''_3) = \{j\}$ and $\text{dom}(w'_2) \supseteq \text{dom}(w_2)$. We conclude from the fact that $w'_2 \oplus w''_3 \sqsupseteq^{\text{pub}} w_5$, that the stack can only have shrunk. Since size is called in a store satisfying $w'_1 \oplus w'_2 \oplus w''_3$, $s_1 - \text{size}()$ must give a non-negative result. \square

5. Reference graph dynamics

The previous section shows how the effect interpretation (IO^{std}, Ref^{std}) can be used to verify results that rely on local state abstraction. However, it does not suffice for proving a set of reference graph properties that have previously been proposed as characteristic for object-capability languages. What is lacking is a way to restrict the primitive effects that an expression is allowed to produce. In this section, we introduce those properties, and in the next section, we discuss how they follow from effect parametricity.

The intended properties are standard in the object capability literature and have previously been formalised and proposed as a characterisation of capability-safety by Maffeis et al. [9]. For a programming language with a certain type of operational semantics, they characterise capability-safety through a formalisation of properties about the evolution of the reference graph. Here, we instantiate their formalism for λ_{JS} and explain that their properties are not sufficient to characterise capability-safety. Sometimes, Maffeis et al.'s definitions are more general than our instantiation of it, but not fundamentally stronger.

5.1. Capability safety as reference graph dynamics

We first introduce some notations used by Maffeis et al.: $e \sqsubseteq e'$ if e is a syntactic subterm of e' . Sets $\mathbb{D} \stackrel{\text{def}}{=} \{\mathbf{r}, \mathbf{w}\}$ and $\mathbb{A} \stackrel{\text{def}}{=} \text{Loc} \times \mathbb{D}$ represent read and write permissions and actions on references. For example, (l, \mathbf{r}) denotes reading location l . Allocating and initialising a new memory location is considered a combined reading (\mathbf{r}) and writing (\mathbf{w}) action. The *can influence*-relation on actions $(l, d) \triangleright (l', d')$ holds when $l = l'$, $d = \mathbf{w}$ and $d' = \mathbf{r}$. A set of actions \mathcal{A}_1 can influence another \mathcal{A}_2 ($\mathcal{A}_1 \triangleright \mathcal{A}_2$) when $a_1 \triangleright a_2$ for some $a_1 \in \mathcal{A}_1$, $a_2 \in \mathcal{A}_2$.

We define a labeled version of λ_{JS} 's impure evaluation judgement in Figure 7: $(\sigma, e) \rightarrow_A (\sigma, e)$ with $A \subseteq \mathbb{A}$. We further instantiate Maffeis et al.'s framework by defining $\text{tCap}(e) \stackrel{\text{def}}{=} \{l \mid l \sqsubseteq e\}$, i.e. the capabilities of an expression are the references that it syntactically contains, $\text{priv}(l) \stackrel{\text{def}}{=} \{\mathbf{r}, \mathbf{w}\}$ (i.e. a primitive reference provides reading and writing

$$\frac{e_1 \hookrightarrow e_2}{(\sigma, E\langle e_1 \rangle) \rightarrow_{\emptyset} (\sigma, E\langle e_2 \rangle)} \quad (\text{E-PURE}) \quad \frac{l \notin \text{dom}(\sigma) \quad A = \{(l, \mathbf{r}), (l, \mathbf{w})\}}{(\sigma, E\langle \mathbf{ref} \ v \rangle) \rightarrow_A (\sigma[l \mapsto v], E\langle l \rangle)} \quad (\text{E-REF}) \quad \frac{(\sigma, E\langle \mathbf{deref} \ l \rangle) \rightarrow_{\{(l, \mathbf{r})\}} (\sigma, E\langle \sigma(l) \rangle)}{(\sigma, E\langle l = v \rangle) \rightarrow_{\{(l, \mathbf{w})\}} (\sigma[l \mapsto v], E\langle v \rangle)} \quad (\text{E-DEREF}) \quad (\text{E-SETREF})$$

Figure 7. Action-labeled version \rightarrow_A of the impure λ_{JS} evaluation judgement. Multi-step versions \rightarrow_A^i and \rightarrow_A^* accumulate labels of substeps.

authority) and $\text{cAuth}(\sigma, l)$ is the least set of actions A such that $\{(l, \mathbf{r}), (l, \mathbf{w})\} \subseteq A$ and $(\bigcup_{(l, \mathbf{r}) \in A} \text{tCap}(\sigma(l)) \times \mathbb{D}) \subseteq A$. In other words, a reference l in store σ carries the authority to read and write values at location l and, recursively, the authority of those values.

Maffeis et al. characterise capability-safety by the following reference graph dynamics properties. For brevity, we define $\text{nauth}(\sigma', \sigma) \stackrel{\text{def}}{=} (\text{dom}(\sigma') \setminus \text{dom}(\sigma)) \times \{\mathbf{r}, \mathbf{w}\}$. The map $\text{auth}(\sigma, e) \stackrel{\text{def}}{=} \bigcup_{c \in \text{tCap}(e)} \text{cAuth}(\sigma, c)$ must be a *valid authority map* for the language. This is by definition true iff $(\sigma, e) \rightarrow_A (\sigma', e')$ implies that

- **RG-AUTH1**: $A \subseteq \text{auth}(\sigma, e) \cup \text{nauth}(\sigma', \sigma)$
- **RG-AUTH2**: $\text{auth}(\sigma', e') \subseteq \text{auth}(\sigma, e) \cup \text{nauth}(\sigma', \sigma)$

Additionally, if $(\sigma, e) \rightarrow_A^* (\sigma', v) \not\rightarrow$ and $v \in \text{Val}$, then for any location l , we must have:

- **RG-CONN**: $A \not\vdash \text{cAuth}(\sigma, l)$ implies that $\text{cAuth}(\sigma', l) = \text{cAuth}(\sigma, l) \cup \{(l, \mathbf{r}), (l, \mathbf{w})\}$
- **RG-NOAMPL**: $A \triangleright \text{cAuth}(\sigma, l)$ implies that

$$\text{cAuth}(\sigma', l) \subseteq \text{cAuth}(\sigma, l) \cup \{(l, \mathbf{r}), (l, \mathbf{w})\} \cup \text{auth}(\sigma, e) \cup \text{nauth}(\sigma', \sigma)$$

Property **RG-CONN** is known as “Only Connectivity Begets Connectivity” and **RG-NOAMPL** as “No Authority Amplification”.

5.2. Reference graph dynamics are not enough

The above properties are *necessary but not sufficient* to characterise object-capability languages. They constitute an *over-approximation* of the authority of a term, known as the *topology-only bound on authority*. The approximation is imprecise because it considers indirect references equivalent to direct references and as such ignores objects’ behaviour.

Consider, for example, the ticket dispenser example from Section 4.1, where attacker code was given access to a function $\text{dispTkt} = \text{func}(\{\text{let}(v = \mathbf{deref} \ o)\{o := v + 2; v\}\}$ but no direct access to reference o . The topology-only bound does not distinguish an expression with a reference to dispTkt or a direct reference to o , so the safety of our ticket dispenser cannot follow. In fact, none of the results from Section 4 can be proven using just the topology-only bound on authority.

To be clear, the problem is not just that the soundness of those examples is *hard* to prove using the topology-only bound. Rather, their soundness does not follow because the bound is not strong enough. One can prove this by constructing a language that satisfies the bound but invalidates the examples, for example by adding a `deepInspect` primitive that returns the set of references held by an arbitrary value, i.e. $\text{deepInspect}(\text{dispTkt})$ evaluates to the singleton list $[o]$. More details about this argument are deferred to the TR.

$$\begin{aligned} \iota_{j,L}^{\text{rgn}} &\stackrel{\text{def}}{=} (L, \subseteq, \subseteq, H_j^{\text{rgn}}) \\ \text{Ref}_j^{\text{rgn}} &: W \rightarrow_{\text{mon,ne}} \text{UPred}(\text{Loc}) \\ \text{Ref}_j^{\text{rgn}} \ w &\stackrel{\text{def}}{=} \{(n, l) \mid w(j) = (L, \subseteq, \subseteq, H), l \in L \wedge H =_{n+1} H_j^{\text{rgn}}\} \\ \text{IO}_j^{\text{rgn}} &: (W \rightarrow_{\text{mon,ne}} \text{UPred}(\text{Val})) \rightarrow W \rightarrow_{\text{ne}} \text{Pred}(\text{Cmd}) \\ \text{IO}_j^{\text{rgn}} \ P \ w &\stackrel{\text{def}}{=} \left\{ \begin{array}{l} (n, \text{cmd}) \mid \\ (\text{cmd} \in \text{Val} \Rightarrow (n, \text{cmd}) \in P \ w) \wedge \\ \forall L, H. \text{ if } w(j) = (L, \subseteq, \subseteq, H) \wedge H =_{n+1} H_j^{\text{rgn}} \text{ then} \\ \forall 0 < i \leq n, \sigma_r :_n w, \sigma_f. (\sigma_r \uplus \sigma_f, \text{cmd}) \rightarrow_A^i (\sigma', e') \Rightarrow \\ \exists \sigma'_r, w' \sqsupseteq^{\text{pub}} w, \text{ for } L' = L \cup (\text{dom}(\sigma'_r) \setminus \text{dom}(\sigma_r)). \\ w'(j) = (L', \subseteq, \subseteq, H) \wedge \sigma' = \sigma'_r \uplus \sigma_f \wedge A \subseteq (L' \times \mathbb{D}) \wedge \\ (n - i, e') \in \mathcal{E}[\text{IO}_j^{\text{rgn}} \ P] \ w' \wedge \sigma'_r :_{n-i} w' \end{array} \right\} \\ H_j^{\text{rgn}} &: \mathcal{P}(\text{Loc}) \rightarrow \hat{W} \rightarrow_{\text{mon,ne}} \text{UPred}(\text{Store}) \\ H_j^{\text{rgn}} \ L \ w &\stackrel{\text{def}}{=} \left\{ (n, \sigma) \mid \begin{array}{l} \text{dom}(\sigma) = L \text{ and for all } l \in L. \ n = 0 \text{ or} \\ (n - 1, \sigma(l)) \in \text{JSVal}_{\text{IO}_j^{\text{rgn}}, \text{Ref}_j^{\text{rgn}}}(\text{roll}^{-1} \ w) \end{array} \right\} \end{aligned}$$

Figure 8. Region-based effect interpretation for memory access bounds.

6. Reference Graph Dynamics from Effect Parametricity

Although the reference graph dynamics properties from the previous section are not sufficient to reason about examples like those in Section 4 or to characterise capability safety, they may still be of interest in some applications. In this section, we explain how our Fundamental Theorem implies results that are analogous but a bit more semantic.

First, we need to explain that we cannot derive the properties themselves because of the more semantic nature of our logical relations. Consider, for example, an expression e and a location l that e does not mention (i.e. $l \not\sqsubseteq e$). Now take x fresh and consider $e' \stackrel{\text{def}}{=} \text{let}(x = l) \ e$. Syntactically, e' holds a reference to capability l and Maffeis et al. would consider it to have greater authority, despite the semantic fact that e' never uses l . However, e' purely evaluates to e , so they are equivalent for our logical relations and any results we prove will not distinguish them. Nevertheless, we can prove properties that are analogous but more semantic.

6.1. An Effect Interpretation for Memory Regions

We start from an effect interpretation that formulates a region memory discipline, defined in Figure 8. It uses a special-purpose island $\iota_{j,L}^{\text{rgn}}$ to define the current set of

addresses in a memory region. This set may grow over time (as expressed by the island's transition relation), with newly allocated references or with existing references whose ownership is passed to the region. j is the index of the island in the world. We require that regions are isolated from one another and the rest of the world, i.e. values stored in region j must themselves never access memory outside the region. Formally, the heap invariant H_j^{rgn} enforces this by requiring that values in the region are accepted by $\text{JSVal}_{IO_j^{rgn}, \text{Ref}_j^{rgn}}$.

We instantiate ρ as Ref_j^{rgn} , accepting only region j 's current addresses, and μ as IO_j^{rgn} , which essentially accepts expressions that only access memory within region j and otherwise respect the world invariants and authority bounds. More formally, $IO_j^{rgn} P w$ n -accepts expressions e that are n -accepted by P if they are values and such that executing them in i steps in a w -acceptable store σ_r will only access memory locations inside region j . The resulting expression must $n - i$ -satisfy $\mathcal{E}[IO_j^{rgn} P]$ in a publicly accessible extension w' of world w . The resulting world must extend the region with all freshly allocated references. Additionally, the store is allowed to contain an additional *frame* part σ_f not governed by w which must be left unmodified.

Contrary to IO^{std} , the definition of IO_j^{rgn} does not assume that the evaluation $(\sigma_r \uplus \sigma_f, cmd) \rightarrow_A^i (\sigma', e')$ produces an e' that is a value. This means that IO_j^{rgn} provides guarantees about arbitrary computations, not just those that terminate successfully. Such definitions are typically only found in logical relations for *concurrent* programming languages, but in our case we want memory access bounds even for evaluations that do not terminate or end up in a stuck state. Formally, it does complicate the definition of IO^{rgn} because if e' is potentially not a value, it makes no sense to require that it is accepted by $P w'$. The solution is to require recursively that it is accepted by $\mathcal{E}[IO_j^{rgn} P]$.

6.2. Memory access bounds

With this effect interpretation, our Fundamental Theorem implies memory access bounds similar to the reference graph dynamics properties. To emphasise the correspondence, we define $\text{memBound}(\sigma, e, L)$ as a judgement analogous to (but more semantic than) the statement $\text{auth}(\sigma, e) = L \times \mathbb{D}$.

Definition 1. We define $\text{memBound}(\sigma, e, L)$ for a store σ and $L \subseteq \text{dom}(\sigma)$ iff for $j = 1$, $w = [j \mapsto \iota_{j,L}^{rgn}]$ and any n , we have:

- $(n, e) \in \mathcal{E}[IO_j^{rgn} \text{JSVal}_{IO_j^{rgn}, \text{Ref}_j^{rgn}}] w$
- $\exists \sigma_f, \sigma_r. \sigma = \sigma_r \uplus \sigma_f$ and $\sigma_r \vdash_n w$

The Fundamental Theorem implies that $\text{memBound}(\sigma, e, L)$ holds whenever $\text{auth}(\sigma, e) = L \times \mathbb{D}$ (proof in TR):

Lemma 4. If $\text{auth}(\sigma, e) = L \times \mathbb{D}$, then $\text{memBound}(\sigma, e, L)$.

Furthermore, the memBound judgement satisfies properties akin to the reference graph dynamics properties. The following property corresponds to properties RG-AUTH1 and RG-AUTH2 in the previous section. It specifies that memBound bounds an expression's memory access and that the property is preserved by evaluation.

Lemma 5. If $\text{memBound}(\sigma_1, e_1, L)$ and $(\sigma_1, e_1) \rightarrow_A^i (\sigma_2, e_2)$, then for $L' = L \cup (\text{dom}(\sigma_2) \setminus \text{dom}(\sigma_1))$

- $A \subseteq L' \times \mathbb{D}$ and $\sigma_2|_{\text{dom}(\sigma_1) \setminus L} = \sigma_1|_{\text{dom}(\sigma_1) \setminus L}$.
- $\text{memBound}(\sigma_2, e_2, L')$.

The next two properties correspond to RG-CONN and RG-NOAMPL. The first states that evaluating an expression whose authority *cannot influence* another expression's, leaves that other expression's memBound unaffected. The latter specifies that evaluating an expression whose authority *can influence* another expression's, may increase that other expression's memBound , but not beyond the union of the two expressions' original bounds and newly allocated locations.

Lemma 6. If $(\sigma, e) \rightarrow_A^* (\sigma', v)$, $A \not\triangleright L \times \mathbb{D}$ and $\text{memBound}(\sigma, e', L)$, then still $\text{memBound}(\sigma', e', L)$.

Lemma 7. If $\text{memBound}(\sigma, e_1, L_1)$, $\text{memBound}(\sigma, e_2, L_2)$, $(\sigma, e_1) \rightarrow_A^* (\sigma', v)$ and $A \triangleright L_2 \times \mathbb{D}$, then we have $\text{memBound}(\sigma', e_2, L')$ for $L' = L_1 \cup L_2 \cup (\text{dom}(\sigma') \setminus \text{dom}(\sigma))$.

7. Related Work

Object capability research has a long history, possibly starting with Dennis and Van Horn's proposed hardware protection primitives [1]. Later work on operating systems and processor hardware with capability security primitives is surveyed by Levy [21]. More recent work includes Capicum (primitive capabilities offered by FreeBSD kernel primitives) [3] and CHERI (processor-level object capabilities combined with virtual memory) [4]. A thorough overview of the object capability model and the capability-safe programming language E is in Miller's PhD thesis [2]. Other capability-safe languages include Joe-E [5], Emily [7], W7 [22], Newspeak [8] and Google Caja [6].

Two papers formalise reference graph properties for capability systems using an operational semantics. The first is Maffeis et al.'s paper, presented in Section 5. The second is Shapiro and Weber's verification of the confinement properties of EROS [23]. They prove a property related to one of the properties in Section 5, which also ignores the behaviour of processes/objects.

Dimoulas et al. formalise capability safety in terms of a notion of component boundaries and the ownership of code by security principals [24], in a simple language without mutable state. The formalisation of capability-safety does not seem intended for reasoning about code. Instead, they propose to extend the object capability language with alternative security mechanisms, that enable specific types of reasoning: built-in dynamic access control and an information flow type system. Similarly, Drossopoulou and Noble argue to extend object capability languages with a kind of declarative policies [25]. Generally, we are not convinced that modular enforcement of typical policies requires additional security mechanisms in the language. Rather, they can be implemented cleanly and modularly using standard patterns. Proving that such implementations enforce a declarative policy can be done using techniques as developed here.

Spiessens studied the safety of capability patterns using Knowledge Behaviour Models and a logic called SCOLL [10], [11]. His goal of validating the safety of capability patterns is the same as ours, but he reasons at a higher level of abstraction, viewing executable code as interacting abstract entities with a behaviour specification. As a result, his results do not directly apply to concrete code in a concrete language, but both the implementations and the language must be separately verified to satisfy their specification. Spiessens’ automatic approximation of future behaviour imposes some restrictions in the logic, like the absence of non-monotonic authority changes.

Taly et al. automatically analyse security-critical APIs in a secure subset of JavaScript to guarantee that API implementations do not leak references to objects marked as internal [26]. A limitation is that they only deal with leaking of direct references. Establishing security requires separate verification of objects *not* marked as internal.

Garg et al. [27] and Jia et al. [28] have proposed and studied powerful program logics for *interface-confined* code in, respectively, a first-order and higher-order setting. Like us, they prove properties of arbitrary, untyped attacker code and can reason about memory as well as other effects. They do this for *interface-confined* code, a syntactic requirement that appears similar to the restrictions in an object-capability language. However, the work does not intend to cover object-capability languages, does not look into previous notions of capability-safety (see Section 5), and the authors state that they cannot model object-capability languages in general. Specifically, there is no direct way to reason about untrusted code that gets access to some primitive pointers, but not all, although this can be modeled using additional indirection. It is not clear to what extent the work supports patterns with separate authority over shared data structures as discussed in Section 4.

A second category of related work is on logical relations for proving encapsulation properties in higher-order languages. The logical relations we use are (unary) step-indexed Kripke logical relations [12], [29]. We have generally followed the notations used by Birkedal et al. [15] and we used their recipe for defining recursive worlds using ultrametric space theory. Our worlds are inspired by Dreyer et al.’s [16]. The region-based effect interpretation used in Section 6 produces a logical relation related to one used by Thamsborg et al. for proving relational results about a region-based type-and-effect system [17].

There is a long line of work on using logical relations for proving local state abstraction results [16], [30]–[35]. The work covers increasingly complex languages (e.g. stack variables vs. first-order heap vs. higher-order heap) and uses increasingly complex Kripke worlds for specifying invariants and protocols on the evolution of acceptable/related stores. Recursive types and higher-order heaps are dealt with using either step-indexing like us or more complex (but perhaps more elegant) solutions based on domain theory.

Our work is closely related to this line of work, but there are some differentiating aspects. First, we work for an untyped language. Although it is known that step-indexing

enables logical relations for untyped languages [36], we are (to our knowledge) the first to demonstrate this for a language with a mutable store. Our Kripke worlds are based on and similar to those of Dreyer et al. [16], but some points are novel. Our use of public/private transitions to model restricted authority over a shared resource, and the fact that the public transition relation in our islands can grow are both (to the best of our knowledge) novel. The idea enables a kind of rely-guarantee reasoning, as demonstrated in Section 4.3 and the approach bears a resemblance to permission assertions on a shared region, as in, for example, Dinsdale-Young et al.’s *Concurrent Abstract Predicates* [37]. Also novel is our notion of effect interpretations as a way to define custom world-indexed restrictions on primitive effects, as well as the axioms prescribing compatibility of the effect interpretation to the publicly accessible effects of the language. Finally, perhaps most importantly, the link to object capabilities and our characterisation of capability-safety are novel.

8. Conclusion

In summary, this paper presents a novel approach for formal reasoning about a capability-safe programming language. We use state-of-the-art techniques from programming language research for capturing the encapsulation of the language, supplemented with some additional ideas for capturing specificities of the model. Our demonstration of three typical but non-trivial patterns in Section 4 and our derivation of reference graph dynamics results in Section 6 shows that our approach is significantly more powerful than previous approaches and sufficiently powerful for realistic examples. We have presented our technique for a subset of λ_{JS} , a relatively simple core calculus for JavaScript, but we expect that it generalises to other settings including assembly languages with primitive capabilities [4], [38], capability-safe subsets of JavaScript and typed capability-safe languages [5], [7], [8]. In the TR, we discuss how our techniques generalise to *relational* rather than unary properties.

Our work also offers new insight into the nature of capability-safe programming languages. Summarily, the property groups three characteristics, all captured by our model: (a) encapsulation of local state: a quite common feature in programming languages like Java, ML, JavaScript etc. (b) absence of global mutable state: a less common feature which is nevertheless crucial for isolating components from each other (like the example in Section 4.3) and, finally, (c) primitive effects only available through primitive capabilities, so authority to produce primitive effects can be controlled by giving components access to the capability or not.

With the better understanding of capability-safety, this work creates the potential for a program logic (perhaps with automated tool support) that can be used to conveniently reason about code in a capability-safe language. Specifically, our Fundamental Theorem tells us what semantic properties such a logic or tool can soundly offer as axioms

over untrusted code. As such, this paper contributes a key semantic understanding, but the design of a program logic or automated tool remains future work.

Acknowledgements

Dominique Devriese holds a postdoctoral mandate from the Research Foundation Flanders (FWO). This research is partially funded by project grants from the Research Fund KU Leuven, and from the Research Foundation Flanders (FWO). This research was supported in part by the ModuRes Sapere Aude Advanced Grant from The Danish Council for Independent Research for the Natural Sciences (FNU).

References

- [1] J. B. Dennis and E. C. Van Horn, "Programming semantics for multiprogrammed computations," *Commun. ACM*, vol. 9, no. 3, pp. 143–155, Mar. 1966.
- [2] M. S. Miller, "Robust composition: Towards a unified approach to access control and concurrency control," Ph.D. dissertation, Johns Hopkins University, 2006.
- [3] R. N. M. Watson, J. Anderson, B. Laurie, and K. Kennaway, "Capsicum: Practical capabilities for UNIX," in *USENIX Security*. USENIX, 2010.
- [4] J. Woodruff, R. N. Watson, D. Chisnall, S. W. Moore, J. Anderson, B. Davis, B. Laurie, P. G. Neumann, R. Norton, and M. Roe, "The CHERI capability model: Revisiting RISC in an age of risk," in *ISCA*. IEEE, 2014, pp. 457–468.
- [5] A. Mettler, D. Wagner, and T. Close, "Joe-E: A security-oriented subset of Java," in *NDSS*, 2010.
- [6] M. S. Miller, M. Samuel, B. Laurie, I. Awad, and M. Stay, "Safe active content in sanitized JavaScript," Google, Tech. Rep., 2008.
- [7] M. Stiegler, "Emily: A high performance language for enabling secure cooperation," in *C5*. IEEE, 2007, pp. 163–169.
- [8] G. Bracha, P. von der Ah, V. Bykov, Y. Kashai, W. Maddox, and E. Miranda, "Modules as objects in Newspeak," in *ECOOP*, ser. LNCS. Springer Berlin Heidelberg, 2010, vol. 6183, pp. 405–428.
- [9] S. Maffei, J. Mitchell, and A. Taly, "Object capabilities and isolation of untrusted web applications," in *S&P*. IEEE, 2010, pp. 125–140.
- [10] F. Spiessens, "Patterns of safe collaboration," Ph.D. dissertation, Université Catholique de Louvain, 2007.
- [11] F. Spiessens and P. Van Roy, "A practical formal model for safety analysis in capability-based systems," in *TGC*, ser. LNCS. Springer Berlin Heidelberg, 2005, vol. 3705, pp. 248–278.
- [12] A. W. Appel and D. McAllester, "An indexed model of recursive types for foundational proof-carrying code," *ACM Trans. Program. Lang. Syst.*, vol. 23, no. 5, pp. 657–683, Sep. 2001.
- [13] A. Ahmed, A. Appel, and R. Virga, "A stratified semantics of general references embeddable in higher-order logic," in *LICS*. IEEE, 2002, pp. 75–86.
- [14] A. Guha, C. Saftoiu, and S. Krishnamurthi, "The essence of JavaScript," in *ECOOP*. Springer, 2010, pp. 126–150.
- [15] L. Birkedal, B. Reus, J. Schwinghammer, K. Støvring, J. Thamsborg, and H. Yang, "Step-indexed Kripke models over recursive worlds," in *POPL*. ACM, 2011, pp. 119–132.
- [16] D. Dreyer, G. Neis, and L. Birkedal, "The impact of higher-order state and control effects on local relational reasoning," *J. Funct. Program.*, vol. 22, no. 4–5, pp. 477–528, 2012.
- [17] J. Thamsborg and L. Birkedal, "A Kripke logical relation for effect-based program transformations," in *ICFP*. ACM, 2011, pp. 445–456.
- [18] D. Devriese, L. Birkedal, and F. Piessens, "Reasoning about object capabilities with logical relations and effect parametricity - technical report including proofs and details," Dept. of Computer Science, KU Leuven, Tech. Rep. CW690, 2016. [Online]. Available: <http://www.cs.kuleuven.be/publicaties/rapporten/cw/CW690.abs.html>
- [19] E. Moggi, "Notions of computation and monads," *Inf. Comput.*, vol. 93, no. 1, pp. 55–92, 1991.
- [20] P. Wadler, "Comprehending monads," *Math. Struct. Comp. Sci.*, vol. 2, pp. 461–493, 12 1992.
- [21] H. M. Levy, *Capability-based computer systems*. Digital Press Bedford, 1984, vol. 12.
- [22] J. A. Rees, "A security kernel based on the lambda-calculus," Ph.D. dissertation, MIT, 1995.
- [23] J. Shapiro and S. Weber, "Verifying the EROS confinement mechanism," in *S&P*. IEEE, 2000, pp. 166–176.
- [24] C. Dimoulas, S. Moore, A. Askarov, and S. Chong, "Declarative policies for capability control," in *CSF*. IEEE, 2014, pp. 3–17.
- [25] S. Drossopoulou and J. Noble, "The need for capability policies," in *FTfJP*. ACM, 2013, pp. 6:1–6:7.
- [26] A. Taly, Û. Erlingsson, J. Mitchell, M. Miller, and J. Nagra, "Automated analysis of security-critical JavaScript APIs," in *S&P*. IEEE, May 2011, pp. 363–378.
- [27] D. Garg, J. Franklin, D. Kaynar, and A. Datta, "Compositional system security with interface-confined adversaries," in *MFPS*, ser. ENTCS. Elsevier, 2010, vol. 265, pp. 49 – 71.
- [28] L. Jia, S. Sen, D. Garg, and A. Datta, "A logic of programs with interface-confined code," in *CSF*. IEEE, July 2015, pp. 512–525.
- [29] A. J. Ahmed, "Semantics of types for mutable state," Ph.D. dissertation, Princeton University, 2004.
- [30] P. W. O'Hearn and R. D. Tennent, "Parametricity and local variables," *J. ACM*, vol. 42, no. 3, pp. 658–709, 1995.
- [31] U. Reddy and H. Yang, "Correctness of data representations involving heap data structures," *Sci. Comput. Program.*, vol. 50, no. 1/3, p. 129/160, 2004.
- [32] A. Banerjee and D. A. Naumann, "Ownership confinement ensures representation independence for object-oriented programs," *J. ACM*, vol. 52, no. 6, pp. 894–960, Nov. 2005.
- [33] N. Bohr and L. Birkedal, "Relational reasoning for recursive types and references," in *PLaS*, ser. LNCS. Springer Berlin Heidelberg, 2006, vol. 4279, pp. 79–96.
- [34] L. Birkedal, K. Støvring, and J. Thamsborg, "Relational parametricity for references and recursive types," in *TLDI*. ACM, 2009, pp. 91–104.
- [35] A. Ahmed, D. Dreyer, and A. Rossberg, "State-dependent representation independence," in *POPL*. ACM, 2009, pp. 340–353.
- [36] A. M. Pitts, "Step-indexed biorthogonality: a tutorial example," in *Modelling, Controlling and Reasoning About State*, ser. Dagstuhl Seminar Proceedings. LZfI, 2010, no. 10351.
- [37] T. Dinsdale-Young, M. Dodds, P. Gardner, M. Parkinson, and V. Vafeiadis, "Concurrent abstract predicates," in *ECOOP*, ser. LNCS. Springer Berlin Heidelberg, 2010, vol. 6183, pp. 504–528.
- [38] N. P. Carter, S. W. Keckler, and W. J. Dally, "Hardware support for fast capability-based addressing," in *ASPLOS*. ACM, 1994, pp. 319–327.