

1 Reasoning About Foreign Function Interfaces 2 Without Modelling the Foreign Language

3 **Alexi Turcotte**

4 Northeastern University

5 Boston, MA, USA

6 **Ellen Arteca**

7 Northeastern University

8 Boston, MA, USA

9 **Gregor Richards**

10 University of Waterloo

11 Waterloo, Ontario, Canada

12 — Abstract —

13 Foreign function interfaces (FFIs) allow programs written in one language (called the *host* lan-
14 guage) to call functions written in another language (called the *guest* language), and are wide-
15 spread throughout modern programming languages, with C FFIs being the most prevalent. Un-
16 fortunately, reasoning about C FFIs can be very challenging, particularly when using traditional
17 methods which necessitate a full model of the guest language in order to guarantee anything
18 about the whole language. To address this, we propose a framework for defining whole language
19 semantics of FFIs without needing to model the guest language, which makes reasoning about C
20 FFIs feasible. We show that with such a semantics, one can guarantee some form of soundness
21 of the overall language, as well as attribute errors in well-typed host language programs to the
22 guest language. We also present an implementation of this scheme, Poseidon Lua, which shows
23 a speedup over a traditional Lua C FFI.

24 **2012 ACM Subject Classification** .

25 **Keywords and phrases** .

26 **Digital Object Identifier** 10.4230/LIPIcs...

27 **Acknowledgements** The authors would like to thank Rafi Turas for writing the implementation
28 of these techniques in Poseidon Lua. We'd also like to thank Hugo Musso Gualandi for his
29 valuable discussions/feedback. This work was partially funded by NSERC.

30 **1 Introduction**

31 Often, programming languages are designed with a specific purpose or task in mind. For
32 example, domain specific languages (DSLs) exist for a variety of domains (e.g., querying
33 databases), and a programmer will often choose a DSL when solving a problem that falls in
34 its domain. But when a programmer wants to write code which touches on several domains,
35 they turn to more general purpose languages (e.g., Java) to give them the tools they need
36 to do everything they need to do, even though the language might be worse at any one
37 given task as compared to a DSL written specifically for it. With so many programming
38 languages to choose from, not only is picking the right language non-trivial, picking the
39 “wrong” language may come back to haunt you.

40 To make choosing a language easier, many programming languages are equipped to
41 interoperate with other languages, and one of the most common forms of interoperation is



© .;
licensed under Creative Commons License CC-BY

Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

42 the foreign function interface (FFI). FFIs allow code written in one language (called the
 43 *host* language) to call functions written in another language (called the *guest* language), and
 44 also interface with data from the guest language, typically accomplished with wrapper code
 45 surrounding guest language values and regulating access to them. By and large the most
 46 common form of language interoperation is the C FFI since C is so fast; C FFI’s are available
 47 for Python, Lua and many other dynamic languages.

48 Semantically, interfacing with C exposes one to all of C’s foibles and irregularities: Memory
 49 accesses can fail, return values of an incorrect type, or cause system-specific undefined
 50 behavior. As such, FFI’s are usually avoided in language semantics, and assumed to be
 51 either benign or absent. Unfortunately, proving properties of the behavior of a C FFI using
 52 conventional techniques is challenging: Of the existing body of work on formal specification
 53 of language interoperation, some are designed with a very specific use case in mind [6][1],
 54 and others propose general frameworks [16] which are difficult to use when reasoning about
 55 interoperation with C; these general approaches rely on fully defined semantics for all
 56 interoperating languages, which is infeasible when one of those languages is C.

57 In this paper, we aim to describe what behavioral guarantees remain true in the presence
 58 of an FFI, how a language hosting an FFI can guarantee its own type correctness at the
 59 interface, and how that can motivate the implementation of an FFI. We propose a framework
 60 which allows typed languages with a C FFI to be formalized and easily reasoned about *without*
 61 *a full model of C*. Our approach relies on a merger of the guest and host language’s type
 62 systems, which allows us reason statically about the whole language and the host language’s
 63 use of the FFI. Additionally, without a model of C, our semantics are *nondeterministic*—as
 64 there’s no telling what an arbitrary C function might do—and we develop a novel method to
 65 reason about these nondeterministic semantics. In principle, this approach works well with
 66 other languages too, though our model of C’s memory and C’s types in the host language
 67 make languages with similar memory behavior to C’s most suitable.

68 As an example of our framework in action, we also present both the semantics and
 69 implementation of Poseidon Lua, a Typed Lua C FFI. In Poseidon Lua, Typed Lua interfaces
 70 with C by holding direct pointers to C data, and is equipped to dereference these pointers,
 71 cast them, allocate C data directly, as well as call arbitrary C functions. We prove conditional
 72 soundness of Poseidon Lua, and prove that if anything “goes wrong” in well-typed Poseidon
 73 Lua programs, C code is at fault for the error. Interestingly, merging the type systems of the
 74 constituent languages eliminates the need for wrapper code around guest language values,
 75 which contributes to improved overall performance.

76 The main contributions of this paper are:

- 77 ■ a framework for merging type systems of guest and host language to allow interoperation
 78 that can be easily reasoned about;
- 79 ■ a semantics for Poseidon Lua, a Typed Lua C FFI, implemented with our framework;
- 80 ■ an actual implementation of Poseidon Lua;
- 81 ■ improved performance results over the previously existing Lua C FFI.

82 **2 Background**

83 In this section, we will provide requisite background for understanding our proposed frame-
 84 work, as well as our prototype implementation, Poseidon Lua. We will begin with an overview
 85 of foreign function interfaces, as we are describing a framework for reasoning about them.
 86 We will also discuss taint analysis, since the concept of taint features prominently in our
 87 semantics. We will then discuss Lua, Typed Lua, and Featherweight Lua, as all are crucial

88 to understanding our language Poseidon Lua. We end the section with a quick highlight of
89 some related work.

90 2.1 Foreign Function Interfaces

91 A foreign function interface (FFI) is a framework in which code written in one language (called
92 the *host* language) may call code written in another language (called the *guest* language) as
93 well as interface with data from that guest language. In an FFI, the guest language typically
94 exports an API of available functions to the host language, and the host language calls
95 said functions through the function interface. In addition to this function interface, a *data*
96 *interface* is required to manage the use of one language's data in the other language.

97 FFIs are prevalent in modern programming. They date back to Common Lisp [11], which
98 first introduced the concept of calling functions written in another language. Many dynamic
99 languages, such as Python [22] and Perl [19], have easy-to-use C FFIs, allowing programmers
100 to quickly and easily call functions written in C, a language known for its speed. In fact,
101 C FFIs are very common, particularly in systems where performance is critical: Scientific
102 computing environments, such as MATLAB [15] and Julia [9], carry out intensive numeric
103 computations and simulations, and often programmers turn to external C functions available
104 through an FFI to speed up the running time of their computationally intensive programs.
105 This provides the user with an easy-to-use scripting language front end which may not be
106 very performant, but with the ability to call fast functions when speed becomes an issue.

107 Most C FFIs interface with C in environments where C has access to all memory, including
108 that of the host language, but there are exceptions where C is an embedded language with
109 restricted access. A popular such system is Emscripten [27]. Emscripten is a source-to-source
110 compiler from LLVM to JavaScript; its goal is to provide a way to run code on the web which
111 can be compiled with LLVM but not natively run in browsers. Since JavaScript can run in
112 essentially any web setting, compiling a language such as C to JavaScript would enable it to
113 run reliably on a browser. With Emscripten, this can be done by first compiling the original
114 source code down to LLVM, and then translating this to JavaScript. In terms of semantics,
115 C is isolated to its own heap, and cannot interfere with JavaScript's; this is the same setup
116 we have in our semantics, where C is isolated from the host.

117 Idiomatic usage of FFIs is to minimize the data interface between the languages to the
118 point where only primitive, scalar values are passed between the languages, as sharing actual
119 structured data has unfortunate behavior: Often, if the FFI even has the capability to allow
120 the host language to store pointers to guest structures, they are mediated through a wrapper.
121 This wrapper problem is insidious: Consider, for example a list. With each access to the
122 next element of a list, a new wrapper must be allocated, and the old wrapper discarded, so a
123 series of simple accesses instead becomes a series of allocations. If the FFI has no capability
124 to access structured guest data, as in Lua's inbuilt C FFI, the programmer has to write a
125 C accessor for every member they want to access. While the definition of these accessors
126 can be automated, they still incur the FFI to actually access the data, as the accessors are
127 written in C.

128 Formally specifying FFIs (and language interoperation in general) is not unknown to
129 the research community. One example is early work by M. Abadi and coauthors [1], which
130 explores dynamic typing in a statically typed language, a mixing of two very different language
131 paradigms. Other work by K. Gray [6] tackles the problem of multi-language object extension,
132 and presents a sound calculus modelling the language interoperability and the semantics of
133 objects written in one language being extended in another. Additional work by J. Matthews
134 and R. B. Findler [16] realizes whole language semantics by defining full semantics for host

135 and guest languages, and uses *boundaries* to explicitly regulate value conversions. For our
 136 purposes, these approaches are either too specific [1][6], or do not generalize to reasoning
 137 about languages with a C FFI [16]. One particular work has a similar motivation to ours
 138 and has a fairly generalizable approach: linking types presented by Patterson and Ahmed
 139 [21]. This is discussed below in Section 2.4.

140 2.2 Dynamic Taint Analysis

141 Introduced by Newsome and Song in their paper [18], dynamic taint analysis is a technique
 142 initially developed for tracing potential error propagation through a system, in order to
 143 detect exploits on commodity software. The idea is that some data sources are considered
 144 *untrusted*, and data which originates from these sources is labelled with taint. This allows for
 145 the tracking of potential errors, and also can be used to restrict what the tainted data can
 146 be used for. In addition, if there is an error in the program that involves some of the tainted
 147 data, information on what potentially caused the error is all available as taint information.

148 The idea of dynamic taint analysis can be generalized to the tracking or propagation
 149 of any tagged (tainted) data in a program. In this work, we adapt the concept of taint to
 150 reasoning about a C FFI without modelling C: when a C call occurs, we cannot say what
 151 *will* happen, but we can reason about what *could* happen. We can model arbitrary C calls
 152 by tagging any data which could have been modified by the call with taint information
 153 identifying it, and should an error occur involving any of this data, the taint can point to
 154 the call which tampered with the data. Note that this is a property of the semantics for the
 155 purpose of proofs; we do not demand that an implementation track dynamic taint. This is
 156 explained in detail in Sections 3.2 and 4.

157 2.3 The Base for Poseidon Lua

158 Later in this work, we will be presenting Poseidon Lua, a Typed Lua C FFI. In this section,
 159 we present variants of Lua, the host language in Poseidon Lua. First we discuss the Lua
 160 language itself, before turning our attention to its variants and extensions.

161 2.3.1 Lua

162 Lua is a lightweight dynamic imperative scripting language with lexical scoping and first class
 163 functions. Lua is extensible, and offers many metaprogramming mechanisms to facilitate
 164 adaptation of the language. Its main data structure is an associative array known as the
 165 table, which can stand in for most common data structures, such as arrays, records, and
 166 objects. The functionality of tables can be further augmented through metamethods, which
 167 are essentially hooks for the Lua compiler. Classic object-oriented programming patterns,
 168 such as methods and constructors, can be easily encoded in Lua with these table extensions.
 169 A C FFI was developed for Lua by Facebook [3]: called `luaiffib`, it is a standard C FFI
 170 which wraps C data for use by Lua. Note that we did not implement Poseidon Lua on top of
 171 LuaJIT [20], as the implementation merely serves as a demonstration of the semantics, and
 172 JIT compilers are less amenable to such modifications. Also, LuaJIT offers the same sort of
 173 data interface that we do, but without types and with boxed references to C structures—our
 174 techniques would thus apply to it for better performance.

175 Our approach to reasoning about FFIs involves embedding the type system of the guest
 176 into the host language, but Lua has no type system to embed into! For this reason, Lua is
 177 not the host language in Poseidon Lua—as we need a type system, we chose Typed Lua as a
 178 base.

179 2.3.2 Typed Lua

180 Lua is a dynamic language, and as is often the case with these languages (see TypeScript [17]
181 and Typed Racket [25]), there have been a few attempts at adding types in some form. One
182 such example with Lua specifically is Tidal Lock [12], a static analyzer relying on simple
183 type annotations. Another is Typed Lua, an optional type system for Lua [14].

184 In their design of Typed Lua, Maidl et al. performed an automated analysis of existing
185 Lua programs to obtain a clear picture of how programmers use the language; they paid
186 close attention to idiomatic Lua code to ensure that their design aligned with conventional
187 language use. Typed Lua is optionally typed, which means that the type annotations are
188 removed when code is compiled. Typed Lua accounts for a large subset of Lua, but a few
189 parts are omitted, namely polymorphic functions and table types, and certain uses of the
190 `setmetatable` function. The type system of Poseidon Lua largely matches Typed Lua's, and
191 a full discussion will appear in Section 4.1.

192 Like other optionally and gradually typed languages, a program written in Typed Lua has
193 an initial stage of *type compilation*. First, the Typed Lua code gets translated (i.e., compiled)
194 to its corresponding Lua program, and it's during this first phase of compilation that the
195 type information is used. At “type compile” time, typed code can be checked statically for
196 type errors before being translated. The type information has no effect on the generated Lua
197 code; Typed Lua programs are type checked by the compiler, and if they are well-typed, the
198 compiler simply erases the types, generating plain Lua. Then, this Lua code is compiled to
199 bytecode and run on the Lua virtual machine.

200 This multistage process means that there are two distinct versions of Lua involved in
201 running a Typed Lua program. For clarity, in our discussion of Poseidon Lua we will use
202 the following terminology: Typed Lua will be referred to as the *typed language* or the *user*
203 *language*, since this is the language in which the programmer will be writing programs. Then,
204 the *untyped language* or the *runtime language* refers to the subset of Lua resulting from the
205 compilation of user language programs and additional expressions needed to deal with C.
206 Both of these languages' grammar and operational semantics are given in Section 4.

207 In giving a prototype using our framework we needed to develop a formal representation
208 of Poseidon Lua. Poseidon Lua is formalized using a *core calculus* based on Featherweight
209 Lua (FWLua) [10], itself a core calculus of Lua. Details on FWLua are given in the next
210 section.

211 2.3.3 Featherweight Lua

212 There have been a few formal specifications of Lua. First, a semantics was developed by
213 M. Soldevila and coauthors [24] to gain a deeper understanding of Lua programs; it was
214 mechanized in PLT Redex [4] using reduction semantics with evaluation contexts. Another
215 semantics, not unlike Featherweight Java [8] and LambdaJS [7], proposes a core calculus for
216 Lua. Called Featherweight Lua (FWLua) [10], this semantics focuses on formalizing what
217 authors deem to be the essential features of Lua: first-class functions, tables, and metatables.
218 Remaining Lua features, including expression sequencing and control structures, are shown to
219 reduce into FWLua through an extensive desugaring process. The FWLua specification [10]
220 also provides a reference interpreter written in Haskell.

221 The principle goal of FWLua is to capture core Lua idioms, and a crucial aspect of the
222 Lua language is its table construct. Under the hood, Lua handles table access and table
223 write with `rawset` and `rawget` functions, respectively; these are not typically written by the
224 programmer, but are part of how Lua drives table functionality. In their design of FWLua,

XX:6 Reasoning About Foreign Function Interfaces

225 the authors modelled table access and table write wholly with these **rawget** and **rawset**
226 operations, and together with other basic semantic constructs (e.g., functions and binary
227 operations) propose functions which mimic the semantics of full-fledged Lua. For example,
228 to capture Lua’s scoping rules, FWLua reserves certain tables to be so-called “scope tables”:
229 the `_local` table is one such example and is always accessible, and changes whenever a new
230 scope is entered while keeping a reference to its outer scope in its `_outer` member. This way,
231 variable access (say, of `x`) is desugared into a function which first searches through `_local`,
232 and if `x` is not present in `_local`, then it searches recursively through `_local._outer`, and
233 so on until `x` is located, producing `nil` if `x` is not found. This proved challenging to reason
234 about, so we chose to promote variables to first-class language members.

235 To contrast Lua and FWLua, consider the following, which illustrates table construction
236 in Lua:

```
237 local t = {}  
238 t.x -- nil, uninitialized table members are nil  
239 t.x = 42 -- t.x is now 42  
240 t[0] = "hello" -- tables may be indexed like arrays  
241 t["hi"] = 3.14 -- equivalent to t.hi  
242  
243
```

244 As you can see, tables can be accessed in a variety of ways in Lua, and have syntax
245 which specifically supports different access styles, be it array-style or record-style. Tables
246 are incrementally constructed, and can be extended at any time, much like dynamic object
247 extension in JavaScript or other dynamic languages. In FWLua, the above translates to:

```
248 rawset(_local, "t", {})  
249 rawget(rawget(_local, "t"), "x")  
250 rawset(rawget(_local, "t"), "x", 42)  
251 rawset(rawget(_local, "t"), 0, "hello")  
252 rawset(rawget(_local, "t"), "hi", "hello")  
253  
254
```

255 As you can see, the **rawset** and **rawget** functions are used to write and read from a table,
256 respectively. As we mentioned earlier, FWLua desugars variables into special table members:
257 The table `_local` deals with local variables, and the table `_ENV` deals with global variables.

258 2.4 Related Work: Linking Types

259 *Linking types*, presented by Patterson and Ahmed [21], consider a different approach to
260 reasoning about language interoperation. This work considers the languages working together
261 as components within a larger language, which itself encompasses behavior of one language as
262 well as the added behavior of making calls to the other language. Linking types themselves
263 are designed to allow programmers to express and reason about one language’s features in
264 another (possibly) less expressive language which has no concept of those features. With
265 linking types, the programmer can annotate a program to indicate where it interfaces with
266 more expressive code in the linked language. Then, with these types, reasoning about the
267 behavior of the whole program becomes possible.

268 Although both their and our work are motivated by the same essential problem, they
269 both require modelling of both languages and focus more on the language of types than
270 on semantics or proofs. In our work, we take a notably different approach in deciding not
271 to model the behavior of the guest language, and instead work with the semantics of the
272 point of intersection (i.e. the boundary between host and guest), using nondeterminism to
273 consider the potential outcomes of the guest language calls. We believe that our types could

274 be expressed in terms of linking types with no meaningful change to our semantics or proofs,
275 but have not investigated this.

276 **3 The Problem**

277 FFIs are ubiquitous in programming languages, and out of these C FFIs are by far the most
278 popular. Unfortunately, if one wants to make any guarantees about programs using a C
279 FFI, using traditional methods of reasoning is challenging. These necessitate a full semantic
280 model of the guest language to show anything about the overall system, and defining a formal
281 semantics for C is very involved. Further, any such semantics will be compiler-dependent.
282 For example, while the CompCert [2] project was groundbreaking in their implementation of
283 a formally verified C compiler, their guarantees are limited to C programs compiled with
284 this compiler, and do not hold for C programs compiled on other compilers (such as gcc).

285 Hypothetically, if we had a whole language semantics for a system with a C FFI, what
286 might we be interested to show? One result of interest would be some form of type soundness
287 for the host language, to ensure that the inclusion of the FFI in the semantics didn't cause any
288 strange issues. Additionally, we might like to show that if any failures occur in a well-typed
289 program calling a C FFI, then C is in some way *at fault* for the failure. In this work, we
290 show that we can get these results *even without a full model of C!*

291 To achieve this, we will need to be able to reason statically about *use* of the FFI (i.e., the
292 host's interface with the guest). The function interface of an FFI exports function handles, so
293 we can at least check that functions are being called and used correctly, even if we don't know
294 exactly what they do. However, the data interface of FFIs is typically built up *dynamically*,
295 and cannot be reasoned about statically. Indeed, in a conventional FFI, wrappers are built
296 up at runtime as values flow from one language to another, and dynamically regulate access
297 to underlying data.

298 In order to fully guarantee the host language's use of the C FFI correct, we need the data
299 interface to be static, and we can achieve this by embedding C's type system into the type
300 system of the host language. This way, the host language can express C types and statically
301 check its own use of C data instead of relying on runtime wrapper code like in traditional
302 approaches. As it happens, with this scheme wrappers are no longer necessary, and their
303 removal results in improved performance; this is discussed further in Section 5.

304 It's not enough to have a system in place to statically reason about the host language's
305 use of the C FFI, as we still need to consider how we can model calls to C when we have no
306 model of the C code, and how we can reason about the resulting semantics. The mechanisms
307 which enable this are *taint* and *nondeterminism*, discussed next.

308 **3.1 Taint and Nondeterminism**

309 Put simply, without a model of C code, C calls are *nondeterministic*: In this scheme, a
310 well-typed call to a C function could arbitrarily fail or succeed, as there's no telling exactly
311 *what* the function does (e.g., a C function could dereference a null pointer or otherwise crash
312 the program) or *what* the function returns. To account for this, at least two semantic rules
313 for guest language calls are required: one modelling a *successful* call where the function
314 didn't crash and returned correctly, and another modelling *failure*, where the function failed
315 to do so (or, more generally, failed to successfully pass execution back to the host language
316 program). Note that the rule for failure must have strictly more permissive preconditions
317 that any rule modelling a successful call, as failure must always be an option.

318 Unfortunately, this simple model of nondeterministic success and failure does not fully
 319 account for all effects that C can have. For instance, execution a C function could free some
 320 memory that the host program has access to while still terminating and returning successfully,
 321 and the next dereference of a pointer to that memory would fail or return unexpected values.
 322 To fully account for this case where a successful C call has detrimental side effects, we need
 323 some additional mechanism to indicate to subsequent reductions that the function may have
 324 tampered with some data.

325 To model the fact that C code may unexpectedly modify data, we use the concept of *taint*;
 326 here, even successful calls to C functions will taint the memory locations which may have
 327 been modified, indicating the possible presence of a modification which could cause issues
 328 for the next access to this data (e.g., if C deallocated the memory at this location). The
 329 presence of taint at a memory location indicates that use of the location is nondeterministic:
 330 the next use of the location could either succeed, indicating that no fatal modification was
 331 made, or fail, indicating that the C call which prompted the taint to be added modified the
 332 location in such a way as to cause an error on access—either way, the effect will only be
 333 observed on the *next* access. Success in accessing a tainted location does *not* mean that the
 334 value at that location is the value that was there before it became tainted, it just means
 335 that the access did not crash; C could still have changed the value in a way that was not
 336 fatal to the program. Crucially, successfully using a tainted location will *clean* or *remove*
 337 the taint, as from that moment until the next C call we are sure that the location is not
 338 somehow broken, and that its value will not change (unless overwritten by Lua).

339 In summary, nondeterminism and taint together enable us to express the effects that
 340 C may have on the host language program *without* modelling C. Note that since we use a
 341 nondeterministic semantics for C and thus avoid modelling its behavior, in principle this
 342 approach works well with other languages. However, our model of C’s memory and C’s types
 343 in the host language make languages with similar memory behavior to C’s most suitable.

344 To demonstrate this framework, we will present the semantics of Poseidon Lua, a Typed
 345 Lua C FFI. A high-level description of Poseidon Lua will be given in the next section.

346 3.2 Overview of Poseidon Lua

347 Essentially, Poseidon Lua is Typed Lua with a C FFI. It is fine-grained relative to standard
 348 FFIs: Unlike traditional FFIs, in Poseidon Lua the type systems of Lua and C are merged
 349 through a Lua pointer type, and the language has syntax with which the Lua programmer
 350 can allocate and manipulate these pointers. Specifically, Poseidon Lua allows you to: allocate
 351 and use C data, cast said pointers, and call C functions. The formal semantics are discussed
 352 fully in Section 4.

353 In our semantics of Poseidon Lua, Lua directly holds C values through a *pointer* to some
 354 location in a C store, which is separate from Lua’s store. Structs are laid out in the C
 355 store as they would be in C, taking up space proportional to the number of struct members;
 356 these members can then be accessed with an offset equal to its position in the list of struct
 357 members (like accessing elements in an array). As explained, with no model of C, C function
 358 calls are nondeterministic, and successful calls taint everything in the C store, since the
 359 function could have modified any memory C has access to in a way which breaks a later
 360 access—for this reason, our formalization includes optional taint information in the C store.
 361 Access to clean (i.e., taint-free) locations in the C store are deterministic, while accesses to
 362 tainted locations are not, and in the event of successful access to a tainted location the taint
 363 can be removed and future accesses to that same location become deterministic (at least,
 364 until the next call to a C function).

365 In addition to modelling possibly errant C calls, taint allows us to model C’s undefined
366 behavior. One classic example of this is casting pointers in C. In Poseidon Lua, as in C,
367 pointers to C values may be downcast. To model this in our formal semantics, we include
368 types in the C store, alongside taint and the values themselves—the C store is thus a list of
369 triples of (*value*, *type*, *optional taint*). This way, we can model the cast of a Lua pointer
370 (to a C value) to some type *T* by changing the type held at the pointer’s location in the C
371 store to *T*. But that’s not quite enough, as casting pointers is undefined behavior in C, and
372 we can use taint to cleanly capture this: Once cast, the location becomes tainted, and the
373 next access to that location is nondeterministic. In this scenario, taint indicates the cast
374 location’s potential for undefined behavior when it is accessed.

375 Another use of taint in Poseidon Lua is in our modelling of allocation of C pointers. In
376 C, the `calloc` function initializes the allocated memory with 0s, so in allocating a pointer to
377 a pointer, one is actually allocating a pointer to a 0 (which is to be treated as a pointer)!
378 Indeed, if one were to dereference the second pointer, one would be dereferencing 0 which
379 leads to a segmentation fault in most circumstances (0, of course, is NULL in C). To achieve
380 this in our semantics, we taint the allocated memory location when a (Lua pointer to a) C
381 pointer is being allocated, to indicate the potential failure of the next access to this location.

382 Even though we don’t model C, we do make some assumptions about C’s behavior: For
383 one, we assume that C does not touch Lua’s memory, and that its effects are contained to an
384 explicitly defined C store: in other words, the shared memory has clearly defined bounds.
385 This mirrors reality in most other FFIs, where guest code and data is not aware of host code
386 and data. However, it is technically possible for C code to violate this assumption. We also
387 make a simplifying assumption that all allocation and access is by word, which reduces the
388 complexity of C data accesses without loss of generality. We require that C doesn’t write new
389 or mutate existing Lua code, otherwise we would have to scrutinize existing expressions that
390 have yet to be reduced and would be unable to prove anything. We additionally make no
391 explicit mention of the stack pointer, which would needlessly complicate function calls and
392 returns for no real benefit. Further, C functions cannot call Lua functions in our formalization,
393 so as to package all of C’s effects into one black box; this is possible through callbacks, but
394 would again be very complex without meaningfully improving the semantics. Finally, we
395 disregard threads, which avoids needing to reason about the effects of concurrency on top of
396 the effect of C, a layer of complexity which is outside of the scope of this project.

397 4 Semantics

398 Poseidon Lua is our proof-of-concept for the ideas discussed in Section 3. Having highlighted
399 some of the stranger corners of our formal specification of Poseidon Lua in Section 3.2, we
400 will now discuss the C FFI in its entirety.

401 In Poseidon Lua, Lua primarily interacts with C by calling C functions, and our merger of
402 the two languages necessitates that C values be a part of the broader language. To represent
403 these C values, Typed Lua has a concept of a Lua pointer to a C value, which is Lua’s
404 window to accessing C data. This means that Lua never deals directly with C values per
405 se, and instead deals with pointers to these values. With pointers to C values as first-class
406 citizens in Poseidon Lua, we implement the additional functionality of allocating C data as
407 well as downcasting C pointers, both directly from Lua code without needing to call C.

408 We start by describing the type system in detail, and follow with a presentation of a core
409 calculus which models the language. Then, we discuss the typing and reduction relations
410 before concluding with a discussion of soundness and other interesting proven results.

T	$::=$ nil $ $ value $ $ ref T $ $ $T_1 \cup T_2$ $ $ L $ $ B $ $ $T_1 \rightarrow_L T_2$ $ $ $\{f_1, \dots, f_n\}$ $ $ ptr_L T_C	<i>nil type</i> <i>top type</i> <i>reference type</i> <i>union type</i> <i>literal type</i> <i>base type</i> <i>function type</i> <i>table type</i> <i>Lua pointer type</i>	f	$::=$ $s : T$ $ $ const $s : T$	<i>fields</i> <i>const fields</i>
T_C	$::=$ int $ $ $T_C^1 \rightarrow_C T_C^2$ $ $ ptr_C T_C $ $ $\{s_1 : T_C^1, \dots, s_n : T_C^n\}$	<i>C integer type</i> <i>C function type</i> <i>C pointer type</i> <i>C struct type</i>	L	$::=$ $\langle \text{booleans} \rangle$ $ $ $\langle \text{numbers} \rangle$ $ $ $\langle \text{strings} \rangle$	<i>literals</i>
			B	$::=$ boolean $ $ number $ $ string	<i>base types</i>

■ **Figure 1** The Poseidon Lua type system.

4.1 Type Systems

Poseidon Lua’s type system is a combination of Typed Lua’s [14] and C’s type systems. For illustrative purposes, we chose a subset of C’s type system which highlights some of C’s interesting features without getting bogged down in the low-level details; we only formalized integers, pointers, structs, and functions. These are not limitations of the concept, merely simplifications made to the formalization. The story is similar with Typed Lua’s type system; our function type only has a single argument type, and multivariate functions are curried to repeated application of single variable functions, by which a single argument function type suffices. In fleshing out this type system for our core calculus, we found no need for Typed Lua’s type variables, recursive types, and projection types, and were able to greatly simplify their table type. Further, to simplify reasoning about Lua, we only allow string indexing in tables. Again, these are not limitations of the language, and are only simplifications for the purposes of formalization.

Our types are given in Figure 1, and explained in detail throughout this section. Type ordering is as follows:

- **value** is a supertype of all types;
- **nil** is the type of Lua’s `nil` value, and is a subtype of all base types;
- union types are supertypes of their members;
- literal types are the types of literals (e.g. the literal type of `5` is `5`), and base types are the more general typical types of these literals (e.g. the base type of `5` is *numeric*)—that said, literal types are subtypes of their corresponding base types;
- function types are contravariant in their argument types, and covariant in their return types;
- table types have **width subtyping**: A table type T is a supertype of a table type T' which has a superset of all of the fields of T (in other words, adding extra fields preserves the subtyping relationship);
- table types have **depth subtyping** only on **const** fields: If a table type T has a **const** field x with type T_x , and a table type T' has all the same fields as T except that field x has type T'_x , where $T'_x <: T_x$, then $T' <: T$ (in other words, **const** field types may be specialized while preserving the subtyping relationship)

C’s types are included in the Typed Lua type system (and made accessible to the user)

442 via the “Lua pointer” C type $\text{ptr}_L T_C$; here, ptr_L denotes a Lua pointer type, and T_C is the
 443 C type being pointed to (e.g., $\text{ptr}_L \text{int}$ is a Lua pointer to a C integer). As explained above,
 444 Lua only ever deals with *pointers* to C values, and not C values themselves: the only access
 445 to C values is through this pointer. C’s type system is consequently entirely self contained,
 446 and is a strict subset of Lua’s with no ability to reference Lua types. In some sense, C is
 447 “plugged” in to Lua through the $\text{ptr}_L T_C$ type.

448 While we don’t formally model C, we do need some information on C functions in order
 449 to ensure that everything shakes out properly at runtime. For example, in our semantics
 450 we model C functions as black boxes with no function body, and we ask for parameter
 451 and return types for these functions to ensure that they are called with correctly-typed
 452 arguments, even though the function bodies themselves are not modeled. What this means
 453 is that we can make sure that the functions are called correctly, but are not responsible for
 454 their internal behavior. Indeed, FFIs typically export function types as part of their API
 455 and may not always export their code—this is the situation modeled by our semantics. This
 456 is also analogous to a user calling a library for which the source code is not provided, even
 457 when the library is written in the same language as the “library host” language.

458 4.2 The Language

459 In this section, we present a core calculus modelling Poseidon Lua, akin to FWLua [10]. We
 460 will discuss the language of expressions, both typed and untyped, before moving on to the
 461 typing judgment and reduction relation.

462 We present *two* languages (in the same manner as Typed Lua, recall from Section 2.3.2):
 463 The language of *untyped expressions* E , also known as the language of *runtime expressions*,
 464 is the language that will actually reduce at runtime, and the language of *typed expressions*
 465 TE is the language that programmers will interface with and program in, with a few minor
 466 caveats which will be discussed in time. Roughly, the typed language corresponds to Typed
 467 Lua with our added C FFI, and the untyped language corresponds to a subset of Lua with
 468 additional expressions for C interoperation. We begin with the typed language TE .

469 4.2.1 Typed Language

470 Figure 2 presents the language of typed expressions, representing the language that the
 471 programmer will be interfacing with, with some notable exceptions. The *Lua dereference*
 472 and *location update* expressions, and the *Lua location* value are not explicitly written by the
 473 programmer; they are artifacts of our typing judgment which will be presented in Section 4.3.
 474 We sometimes refer to the aforementioned expressions as *intermediate expressions*; the typed
 475 language without these is the *user language*.

476 These expressions largely describe a core calculus of Typed Lua, with the exception of
 477 the following C expressions:

- 478 ■ *C downcast* denotes the cast of expression te to C type T_C ;
- 479 ■ *C allocation* allocates a *C pointer* to a value of C type T_C ;
- 480 ■ *C deref* is used to dereference the C pointer expression te ;
- 481 ■ *C function* describes a C function with type signature T_C . The type T_C is required by
 482 the type transformation to type these functions, as it cannot leverage the function body
 483 (as is the case with traditional functions);
- 484 ■ *C pointer* is a pointer to location n in the C store, with expected C type T_C ;

$te ::= v_t$	<i>value</i>	$v_t ::= \mathbf{nil}$	<i>nil value</i>
$\{s_1 = v_1, \dots, s_n = v_n\}$	<i>table</i>	r	<i>register</i>
$\mathbf{let} x : T := te_1 \mathbf{in} te_2$	<i>let binding</i>	c	<i>constant</i>
$x := te$	<i>variable update</i>	$\mathbf{loc} n$	<i>Lua location</i>
$\mathbf{loc} n := te$	<i>location update</i>	$\lambda x : T. te$	<i>Lua function</i>
$\mathbf{deref} te$	<i>Lua dereference</i>	$\mathbf{cfun} T_C$	<i>C function</i>
$te_1 \mathit{opt} te_2$	<i>binary operation</i>	$\mathbf{ptr} n T_C$	<i>C pointer</i>
$te_1(te_2)$	<i>function call</i>	$r ::= \mathbf{reg} n$	<i>table store loc</i>
x	<i>variable</i>	$c ::= n$	<i>number</i>
$te_1.te_2$	<i>dot access</i>	b	<i>boolean</i>
$te_1.te_2 := te_3$	<i>dot update</i>	s	<i>string</i>
$\mathbf{cast} te T_C$	<i>C downcast</i>	$op ::= +, -, *, /$	<i>arithmetic</i>
$\mathbf{calloc} T_C$	<i>C allocation</i>	$\leq, <, \geq, >$	<i>order</i>
$\mathbf{deref}_C te$	<i>C deref</i>	\wedge, \vee	<i>boolean</i>
$te_1; te_2$	<i>sequence</i>	\dots	<i>concatenation</i>
		$==$	<i>equality</i>

■ **Figure 2** The language of typed expressions.

485 ■ Access to C structs is done through the *dot access* and *dot update* expressions (so long as
 486 te_1 is a C struct), and calling C functions is done through the *function call* expression
 487 (so long as te_1 is a C function).

488 Besides the C expressions, the typed language is standard or otherwise directly analogous to
 489 some untyped expression, which we will discuss in more detail shortly.

490 Typed expressions will all compile into equivalent runtime expressions where the types
 491 have been erased. We explore this runtime language next.

492 4.2.2 Untyped Language

493 The untyped language describes the expressions which will reduce/evaluate at runtime.
 494 Generally speaking, they are analogous to some equivalent typed expression where the types
 495 have been erased. This language essentially describes a core calculus of Lua, based on
 496 FWLua (described in Section 2.3.3), though we added sequencing, let bindings, variables,
 497 table literals, and of course C interoperability. The full language can be found in Figure 3.

498 FWLua is a core calculus of Lua, and a number of minor modifications were required
 499 when adapting FWLua to describe Typed Lua, particularly with tables. Recall that tables
 500 are the principle data structure in Lua; as discussed previously, FWLua desugars all of
 501 Lua's table manipulation into the dual **rawget** and **rawset** constructs. For the purposes of
 502 formalization, we needed to relax FWLua's extreme desugaring; one example of this being
 503 the table literal (*table*) expression. FWLua handles table construction incrementally: an
 504 empty table is first created and stored, and then it is populated with the values at the
 505 programmer's discretion. Unfortunately, this scheme fails in typed languages, as the empty
 506 table is not a subtype of any non-empty tables, so we include a table literal to allow the
 507 expression of a full table when needed for assignments.

508 Our function expression is unchanged from FWLua, though we must include a new C
 509 function expression to allow FFI calls. Unlike the Lua function, which is a traditional lambda
 510 expression, the C function has far less information in it—indeed, it has no function body!

$e ::= v$	<i>value</i>		
$\{s_1 = v_1, \dots, s_n = v_n\}$	<i>table</i>		
rawget $e_1 e_2$	<i>table select</i>		
rawset $e_1 e_2 e_3$	<i>table update</i>		
$e_1 \text{ op } e_2$	<i>binary operation</i>	$v ::= \mathbf{nil}_L$	<i>nil value</i>
$e_1(e_2)$	<i>Lua fun. appl.</i>	r	<i>register</i>
x	<i>variable</i>	c	<i>constant</i>
$x := e$	<i>var. assignment</i>	loc n	<i>Lua store loc.</i>
loc $n := e$	<i>location update</i>	ptr_L $n T_C$	<i>C store pointer</i>
deref e	<i>Lua dereference</i>	$\lambda x.e$	<i>Lua function</i>
let $x := e_1$ in e_2	<i>let binding</i>	cfun	<i>C function</i>
cget $e n T_C$	<i>C store access</i>	$v_C ::= \mathbf{ptr}_C n$	<i>C store pointer</i>
cset $e_1 n e_2 T_C$	<i>C store update</i>	n	<i>C number literal</i>
ccall $e_1 e_2 T_C \beta$	<i>C function call</i>		
calloc $T_C \beta$	<i>C allocation</i>		
cast $e T_C \beta$	<i>C downcast</i>		
$e_1; e_2$	<i>sequence</i>		
err β	<i>error expression</i>		

■ **Figure 3** The language of untyped, runtime expressions.

511 Most of the information needed for a C call is stored in the C function call expression itself.

512 For accesses into C structs, we have the **cget** and **cset** expressions, analogous to **rawget**
 513 and **rawset**. **cget** and **cset** are also used for accessing and writing to C pointers, which will
 514 be discussed in more detail in Section 4.4. In **cget** $e n T_C$, e is a pointer into the C store, n
 515 is the offset of the access, and T_C is the type that the **cget** is expecting to read. Similarly
 516 in **cset** $e_1 n e_2 T_C$, e_1 is a pointer into the C store, n is an offset, e_2 is the value to write,
 517 and T_C is the type that the **cset** is expecting the store to contain at the referenced pointer
 518 (recall that we store type information for each pointer in the C store).

519 To call functions, programmers may write a standard function application as $te_1(te_2)$ in
 520 the typed language of Figure 2. The type transformation can, depending on the type of te_1 ,
 521 transform the application into either a Lua function application or a C call. The Lua function
 522 call expression $e_1(e_2)$ is straightforward, so let us focus on the C call: In **ccall** $e_1 e_2 T_C \beta$, e_1
 523 is the C function being called, e_2 is the argument to that function, T_C is the function's type,
 524 and β is an identifier associated with the call (its line of code). The type is necessary since
 525 C calls exhibit nondeterministic behavior, and we can leverage T_C to reason about the value
 526 that is returned from the function. The line of code information β is related to taint, which
 527 we will describe fully when giving the semantics of the calls.

528 There are also a few expressions for functionality unique to C. As one might expect,
 529 **calloc** $T_C \beta$ allocates something of C type T_C , and β is the identifier uniquely associated with
 530 the allocation, which allows a trace-back if a runtime error occurs. **cast** $e T_C \beta$ downcasts
 531 the pointer e to type T_C , and again β is a unique identifier associated with the cast.

532 4.3 Typing Judgment

533 Making a distinction between typed and untyped languages (or user and runtime languages)
 534 makes sense in many optionally or gradually typed languages, where a typed language is
 535 compiled into an untyped language which will be the one executing at runtime (recall the

XX:14 Reasoning About Foreign Function Interfaces

536 two stage compilation process described in the context of Typed Lua in Section 2.3.2). In
 537 these settings the typing judgment often needs to be modified to connect the languages
 538 together. We define a *type transformation* relation, a modification of the standard *typing*
 539 *judgment* relation, which transforms/compiles a typed expression into its corresponding
 540 untyped expression:

$$541 \quad \Gamma, K \vdash te : T \rightsquigarrow e \quad (1)$$

542 Here, Γ is the typing environment, which assigns types to variables, and K is the typing
 543 context, containing information about the various store typings. Our runtime environment
 544 contains three stores: a table store for Lua tables, a C store for C values, and a variable
 545 store for variables. K can thus be broken up into three store typings: Σ_T describing the
 546 table store, Σ_C for the C store, and Σ_V for the variable store. Roughly speaking, the type
 547 transformation takes a typed expression te and “compiles” it into an untyped expression e ,
 548 assigning to it type T in the context of Γ and K .

549 In the following typing rules, some auxiliary functions will appear in the preconditions to
 550 simplify the notation. They are as follows:

- 551 ■ *goodLayout*(n, T_C, Σ_C) checks to see if location n in the C store typing Σ_C represents
 552 type T_C . If T_C is a primitive type or a pointer type, this succeeds if $\Sigma_C(n) = T_C$. As for
 553 structs, recall that they are laid out contiguously in the store: If T_C is a struct type (for
 554 example, $\{s_1 : T_C^1, \dots, s_n : T_C^n\}$), then each of the fields must be present in Σ_C with the
 555 correct type, i.e. for all fields s_i we must have $\Sigma_C(n + i) = T_C^i$.
- 556 ■ *offsetForType*(s, T_C) computes the offset of member s in structure type T_C . Our formal-
 557 ization of the C store lays out structs according to their type, and this function relates
 558 their type (T_C) to their layout in the store.

559 As we mentioned, in Poseidon Lua, Lua can interact with C in the following ways:
 560 allocation and access of C data, C function calls, and casting of C pointers. In this section
 561 we will focus on the typing rules for the expressions describing this FFI. The full typing rules
 562 are given in Appendix A.1.

563 We will first consider the rule for allocation of C data.

$$\frac{\text{validType}(T_C) \quad \beta \text{ unused}}{\Gamma, K \vdash \mathbf{calloc} T_C : \mathbf{ptr}_L T_C \rightsquigarrow \mathbf{calloc} T_C \beta} \quad (\text{TT_CALLOC})$$

564 In Poseidon Lua, programmers can allocate Lua pointers to C data types (here, T_C),
 565 provided that the type is *valid* for allocation. For this to be the case, T_C must either be
 566 a primitive type, pointer type, or struct (itself recursively made up of valid types). This
 567 prevents programmers from making nonsensical statements, such as allocating C functions
 568 in Lua. The β here is needed when allocating C pointers: In C, allocating a pointer to a
 569 pointer can cause issues if the innermost pointer is not properly initialized, due to the default
 570 values that C inserts (pointer values are often initialized to 0, which is an invalid memory
 571 address for C to access). This semantics will be dealt with in due course, and the inclusion
 572 of β in the **calloc** expression is crucial to achieving the desired behavior—this will be further
 573 discussed in Section 4.4.

574 Having seen C allocation, we turn our attention to typing (Lua pointers to) C values:

$$\frac{n < \text{length}(\Sigma_C) \quad \text{goodLayout}(n, T, \Sigma_C)}{\Gamma, (\Sigma_T, \Sigma_C, \Sigma_V) \vdash \mathbf{ptr}_L n T_C : \mathbf{ptr}_L T_C \rightsquigarrow \mathbf{ptr}_L n T_C} \quad (\text{TT_LUA_PTR})$$

575 C values are always “hidden behind” a Lua pointer in Poseidon Lua, and so from Lua’s
 576 point of view all C values have some \mathbf{ptr}_L type. In the expression $\mathbf{ptr}_L n T_C$, n is the location
 577 referenced by the pointer, and T_C specifies the type that the location is intended to have.
 578 The type information is required since structures do not directly inhabit the C store, and
 579 so accessing a structure would be impossible with a simpler rule, since $\Sigma_C(n)$ will never
 580 have a struct type; the type information allows us to check to see if location n does in fact
 581 correspond to T_C using the *goodLayout* auxiliary function, and only allow the pointer to
 582 type if it does. The typing rule for dereferencing these pointers follows.

$$\frac{\Gamma, K \vdash te : \mathbf{ptr}_L T_C \rightsquigarrow e \quad \text{validForCDeref}(T_C) \quad T_L = \text{coerceCType}(T_C)}{\Gamma, K \vdash \mathbf{deref}_C te : T_L \rightsquigarrow \mathbf{cget} e 0 T_C} \text{ (TT_VAR_C_DEREF)}$$

583 Here, beyond ensuring that te is in fact a Lua pointer, we need to ensure that it is a
 584 pointer to a type that we can dereference. The C store is made up entirely of primitives
 585 and pointers, so we disallow dereferencing of things of another type (for example, we cannot
 586 dereference a C function pointer). Because our type transformation deals with Lua types
 587 only, we need to coerce T_C into a Lua type to type this expression: Indeed, at runtime the
 588 dereference will coerce the value it obtains from the C store, and the coercion at this level
 589 allows such an expression to type. Note also the untyped expression corresponding to the
 590 dereference: \mathbf{cget} can play the part of either simple dereferencing and also struct field access,
 591 depending on the value of its offset parameter (here, 0). An offset of 0 indicates that we are
 592 either getting the first member in a struct, or simply dereferencing a pointer to non-struct
 593 data.

594 We consider C functions next.

$$\frac{}{\Gamma, K \vdash \mathbf{cfun} (ct_1 \rightarrow_C ct_2) : (ct_1 \rightarrow_C ct_2) \rightsquigarrow \mathbf{cfun}} \text{ (TT_C_FUNCTION)}$$

595 Here, note that the C function expression contains the whole type of the function, and
 596 without a body the function trivially types. Type information is necessary because we don’t
 597 model C’s semantics: In typical typing rules for functions, the return type can be determined
 598 thanks to the function body, and we have no such body to rely on here. In some sense, this
 599 is in line with what one would expect when dealing with FFIs, since part of their API is the
 600 full type of the exported functions.

601 Let us consider how one calls these functions:

$$\frac{\Gamma, K \vdash te_1 : (T \rightarrow_C T') \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : T \rightsquigarrow e_2 \quad \beta \text{ unused}}{\Gamma, K \vdash te_1(te_2) : T' \rightsquigarrow \mathbf{ccall} e_1 e_2 T' \beta} \text{ (TT_C_FUN_APPL)}$$

602 In rule TT_C_FUN_APPL, we type the function application according to its return type.
 603 Note the T' in the compiled (on the right of the \rightsquigarrow) C call: The untyped call requires the
 604 return type for reduction to be possible, and we will discuss this in more detail in Section 4.4.
 605 Since C calls are sources of taint, we include β as an identifier uniquely associated with the
 606 call, which corresponds to the line of code occupied by the call. In the event of a failure, we
 607 can determine which call (and, thus, which function handle) is to blame.

608 We will now consider reading from and writing to C structs. First, reading:

$$\frac{\Gamma, K \vdash te_1 : \mathbf{ptr}_L T_1 \rightsquigarrow e_1 \quad \mathit{structType}(T_1) \quad \Gamma, K \vdash te_2 : s \rightsquigarrow e_2 \quad s \in T_1 \quad n = \mathit{offsetForType}(s, T_1)}{\Gamma, K \vdash te_1.te_2 : T_1(s) \rightsquigarrow \mathbf{cget} \ e_1 \ n \ T_1(s)} \text{ (TT_C_DOT_ACCESS)}$$

609 Here, if te_1 types to $\mathbf{ptr}_L T_1$, T_1 is a struct type, and te_2 types to a string literal s which
 610 is a field name in struct T_1 , then the C struct member access types. Note that te_1 must
 611 be a Lua pointer to a C struct, as C structs themselves are not allowed in Poseidon Lua
 612 unless they are behind a Lua pointer. Also, the resulting \mathbf{cget} is given the offset of field s in
 613 T_1 (determined with the $\mathit{offsetForType}$ auxiliary function), since the C store lays out struct
 614 members linearly in an array form.

615 Second, C struct member update:

$$\frac{\Gamma, K \vdash te_1 : \mathbf{ptr}_L T_1 \rightsquigarrow e_1 \quad \mathit{structType}(T_1) \quad \Gamma, K \vdash te_2 : s \rightsquigarrow e_2 \quad \Gamma, K \vdash te_3 : T_1(s) \rightsquigarrow e_3 \quad s \in T_1 \quad n = \mathit{offsetForType}(s, T_1)}{\Gamma, K \vdash te_1.te_2 := te_3 : \mathbf{value} \rightsquigarrow \mathbf{cset} \ e_1 \ n \ e_2 \ T_1(s)} \text{ (TT_C_DOT_UPDATE)}$$

616 As before, if te_1 is a Lua pointer to a C struct type T_1 , and te_2 is a string s which is a
 617 member of that struct, and te_3 is appropriately typed, we can type the C struct update. We
 618 again emit an offset (in place of te_2), which the \mathbf{cset} will use when writing to the C store.

619 Finally, Poseidon Lua allows C values to be downcast, and they type as follows:

$$\frac{\Gamma, K \vdash te : \mathbf{ptr}_L T'_C \rightsquigarrow e \quad \beta \text{ unused}}{\Gamma, K \vdash \mathbf{ccast} \ te \ T_C : T_C \rightsquigarrow \mathbf{ccast} \ e \ T_C \ \beta} \text{ (TT_C_CAST)}$$

620 Here, we notice that casting must be done through the Lua pointer, and so long as T_C is a
 621 C type we allow the cast to go through. There is no mention of T_C and T'_C being compatible
 622 types, as C freely allows casting of pointers, and the cast merely changes the way that the
 623 bits referred to by the pointer are read. As with previous mentions of β , it features here to
 624 allow errors caused by the cast to be easily traced back to the cast.

625 At this point, we have explored each of the typing rules associated with Poseidon Lua's
 626 C FFI. In many cases, such as in TT_C_FUN_APPL, these rules transferred some type
 627 information to their analogous runtime expressions in order to drive the runtime functionality
 628 of the system. We discuss reduction of runtime expressions next.

629 4.4 Operational Semantics

630 The *reduction relation* on untyped expressions, describing the execution of programs, is:

$$631 \quad e / \sigma_T / \sigma_C / \sigma_V \rightarrow e' / \sigma'_T / \sigma'_C / \sigma'_V \quad (2)$$

632 Here, e and e' are expressions in the untyped language, σ_T and σ'_T are *table stores*, σ_C
 633 and σ'_C are *C stores*, σ_V and σ'_V are *variable stores*. At a high level, the table store σ_T is
 634 a list of Lua tables, the variable store σ_V is a list of values, and finally the C store σ_C is
 635 a list of $(v, T_C, \beta?)$ triples, where v is a C value, T_C is its type, and $\beta?$ is optional taint
 636 information (\emptyset represents no taint, or a clean location). As we mentioned in Section 3.2, the
 637 unusual inclusion of type information in the runtime C store is required to properly model C
 638 downcast semantics.

639 To simplify notation, we sometimes write the reduction relation as:

$$640 \quad e / \mathcal{S} \rightarrow e' / \mathcal{S}' \quad (3)$$

641 We refer to \mathcal{S} and \mathcal{S}' above as the *runtime environment*; the set of all the stores making
642 up the state/context of the reduction.

643 It will be necessary to differentiate between C stores based on whether or not they are
644 tainted; for this purpose, we say that a C store is *clean* if none of the elements of the store
645 are themselves tainted. To simplify discussion of tainted environments, we say that a runtime
646 environment is clean if its C store is also clean.

647 At the very highest level, we are formalizing a system wherein Lua code can interface
648 with C in the following manner: allocating C data, reading from and writing to some shared
649 memory with C, downcasting C values, and calling C functions.

650 Our formalization of Lua is based on FWLua [10], and we adapted their big-step semantics
651 to a more standard small-step equivalent. For our discussion of FWLua, see Section 2.3.3.
652 In order to mechanize our formalization, some simplifying modifications to FWLua were
653 required, namely the promotion of variables from syntactic sugar to full-fledged language
654 members. Of course, Lua allows you to declare and use variables, but FWLua desugars
655 variables into access to a special store carried around at runtime. Poseidon Lua requires that
656 FFI calls be made only from well-typed code, and so we adapted the type system of Typed
657 Lua [14], with some modifications made possible by our simplified semantics for Lua.

658 Notable in Poseidon Lua is the merger of Typed Lua's and C's type systems through the
659 Lua pointer type, and consequently the intermixing of values from both Lua and C. Lua
660 makes reference to C values through the *Lua pointer* expression, and can both access and
661 change the data contained in these pointers, as well as cast them to some C type. Lua may
662 also allocate Lua pointers to C values through the **calloc** expression, without needing to
663 make a **ccall**.

664 We will now turn our attention to the operational semantics of Poseidon Lua, with a
665 focus on the C FFI, mirroring discussion of the typing judgment in Section 4.3. The full
666 reduction rules are given in Appendix A.2. We start with the semantics of allocating C data.
667 Consider:

$$\frac{n = \text{length}(\sigma_C) \quad \sigma'_C = \sigma_C + \text{layoutTypeAndTaint}(T_C, \beta)}{\mathbf{calloc} T_C \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{ptr} n T_C / \sigma_T / \sigma'_C / \sigma_V} \quad (\mathbf{R_CALLOC})$$

668 The **calloc** $T_C \beta$ expression allocates enough memory in the C store σ_C to accommodate
669 a value of type T_C . The function *layoutTypeAndTaint* lays out type T_C and taints pointer
670 members (as per our earlier discussion in Section 3.2). If T_C either is or contains a C pointer
671 type, then we taint that location (with taint information β) to indicate to our system that
672 its behavior is undefined until it is successfully accessed or written to. If T_C is a primitive
673 or pointer type, then we simply produce a triplet containing a default value (this is 0 for
674 pointers), the type T_C , and taint if T_C is a pointer type, and if T_C is a struct, we lay out
675 each of its members in a similar fashion. Following allocation, a C pointer with the location
676 of the beginning of the newly allocated memory is produced.

677 Compared with C allocation, C calls have intricate semantics as we do not attempt to
678 model the bodies of arbitrary C functions. Instead, we treat the C functions like black boxes,
679 and consequently C function calls exhibit *nondeterministic* semantics, as any well-typed C
680 call can either succeed or fail if the function body is made up of arbitrary C code (recall that
681 we consider a call successful if it returns to executing the host language with some value
682 of the expected type). In the event of successful execution, we concern ourselves with the
683 return value and the call's potential effects on the rest of the C data. Recall our discussion
684 that even if a call is successful, the function code might have altered the C store in a variety

XX:18 Reasoning About Foreign Function Interfaces

685 of ways (such as freeing some existing memory), and we must account for this possibility.
 686 We will first consider the reduction rule for a successful C call.

$$\frac{\text{value}(v_2) \quad v = \text{makeValueOfType}(ct_2) \quad \sigma'_C = \text{taintCStore}(\sigma_C, \beta)}{\mathbf{ccall\ cfun}\ v_2\ ct_2\ \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow v / \sigma_T / \sigma'_C / \sigma_V} \text{(R_CCALL_WORKED)}$$

687 Here, $\mathbf{ccall\ cfun}\ v_2\ ct_2\ \beta$ calls a C function \mathbf{cfun} with argument v_2 . In this case, the call
 688 succeeds, and $\text{makeValueOfType}(ct_2)$ gives us v , something of type ct_2 . Of course, since it's
 689 possible that the call tampered with the C store, we taint the store with taint information
 690 β , corresponding to the line of code of this function call. This notifies subsequent accesses
 691 to these memory locations of potential tampering, which modifies the semantics of those
 692 accesses. C function calls can also fail:

$$\frac{\text{value}(v_2) \quad \sigma'_C = \text{taintCStore}(\sigma_C, \beta)}{\mathbf{ccall\ cfun}\ v_2\ ct_2\ \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{err}\ \beta / \sigma_T / \sigma'_C / \sigma_V} \text{(R_CCALL_FAILED)}$$

693 To capture that both success and failure are possible outcomes, we ensure that the
 694 premises of both rules are simultaneously satisfied: When all of R_CCALL_WORKED 's
 695 preconditions are met, so are R_CCALL_FAILED 's (and vice-versa). The $\mathbf{err}\ \beta$ expression
 696 is the result of the failing call, and indicates through taint information β which call is to
 697 blame for the failure.

698 Having seen the intricacies of C calls, we will turn our attention to the semantics of
 699 casting C pointers, another source of taint. For brevity, we only present the rule for casting
 700 a clean location (the other rule is not notably different). Consider:

$$\frac{n < \text{length}(\sigma_C) \quad \sigma_C(n) = (v, T_C, \emptyset) \quad \sigma'_C = \text{update}(\sigma_C, n, (v, T'_C, \beta))}{\mathbf{ccast}\ (\mathbf{ptr}_L\ n\ T_C)\ T'_C\ \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{ptr}_L\ n\ T'_C / \sigma_T / \sigma'_C / \sigma_V} \text{(R_CCAST)}$$

701 Here, the location n in σ_C is updated with the new type T'_C and taint information
 702 associated with the cast (thanks to the update auxiliary function— $\text{update}(s, l, v)$ reads as
 703 “update s at location l with value v ”). In C, casting a pointer merely changes how the bits
 704 being pointed to are read, and the cast may even cause an error; we achieve similar semantics
 705 with taint. When attempting to read location n in σ_C after it was cast, taint indicates that
 706 the access should be nondeterministic. To keep our system as general as possible, we don't
 707 attempt to model the cast per se, and the next read will replace v with a new value of type
 708 T'_C if successful, or fail with an error. We discuss the semantics of accesses next.

709 Thus far, we focused on the introduction of taint and fairly direct sources of nondetermin-
 710 ism, and we will turn our attention to taint's effect on the semantics of our system, as well as
 711 how it can be removed from the runtime environment. As an example, recall our semantics
 712 for C casts: When casting a location to some type T_C , the location becomes tainted. Now,
 713 imagine that the next use of the location is to store something of type T_C in it; if this write
 714 succeeds, from then on we are sure about the value present at the location. Such an operation
 715 is said to *clean* the taint from the location; in our formalization, taint represents uncertainty
 716 about a C value, and once we become certain of it (e.g., we have accessed the value and no
 717 errors have occurred) we can safely remove the taint.

718 In more formal terms, the presence of taint at a location in σ_C indicates that accessing
 719 that location yields nondeterministic results. To capture this, we ensure that a read or

720 write to a tainted location can reduce to more than one expression *under the same premises*;
 721 namely, said read or write can succeed or fail.

722 Consider the following semantics for accessing a clean location in σ_C :

$$\frac{\sigma_C(n+o) = (v_C, T_C, \emptyset) \quad v_{out} = \text{coerceToLua}(v_C)}{\mathbf{cget}(\mathbf{ptr}_L n T'_C) o T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_{out} / \sigma_T / \sigma_C / \sigma_V} \text{ (R_CGET_NO_TAINT)}$$

723 Here, the expression $\mathbf{cget}(\mathbf{ptr}_L n T'_C) o T_C$ accesses σ_C at location n with offset o , and
 724 is expecting something of type T_C . In this reduction rule, location $n+o$ in σ_C is clean,
 725 and so the (well-typed) store access cannot fail. The access steps to v_{out} , which is the Lua
 726 equivalent of the C value contained in σ_C , determined through the *coerceToLua* auxiliary
 727 function. Note that the pointer's type (T'_C) does not necessarily need to match the expected
 728 type of the access (T_C); this is because **cgets** can be used for struct member access, where
 729 T'_C would be a struct type and T_C would be the type of the member.

730 *coerceToLua*(v_C) is a function which takes a C value v and coerces it to a Lua value. If
 731 v_C is a C integer, then it is coerced to a Lua constant with the same numeric value. If v_C
 732 is a C pointer $\mathbf{ptr}_C m ct$, then it is coerced into a Lua pointer $\mathbf{ptr}_{\text{Lua}} m ct$ (to the same
 733 location). Otherwise, the coercion fails.

734 Note the presence of a type T_C in the **cget** expression. A condition of reading (and
 735 writing) from σ_C is that the type specified for the read must match the type held in σ_C .
 736 This allows us to enforce the correct use of downcast locations, as the cast changes the type
 737 in σ_C , and future reads (and writes) must specify the new type.

738 We will now consider accesses to tainted locations, which can either fail or succeed. First,
 739 consider a successful access:

$$\frac{\sigma_C(n+o) = (v, T_C, \beta) \quad v' = \text{makeValueOfType}(T_C) \quad \sigma'_C = \text{update}(\sigma_C, n+o, (v', T_C, \emptyset) \quad v_{out} = \text{coerceToLua}(v'))}{\mathbf{cget}(\mathbf{ptr} n T'_C) o T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_{out} / \sigma_T / \sigma_C / \sigma_V} \text{ (R_CGET_TAINT_WORKS)}$$

740 Here, we access σ_C at location n with offset o , and are expecting something of type T_C
 741 as before. However, $\sigma_C(n+o)$ is tainted, resulting in nondeterminism (i.e. we do not know
 742 whether an access to this value will fail or succeed). In this reduction rule, we deal with
 743 the case of a successful access to tainted locations. Here, a successful access returns some
 744 value of the appropriate type (thanks to the *makeValueOfType* auxiliary function). The C
 745 store at $n+o$ is cleaned and updated with the new value; from this moment on, use of
 746 this location is deterministic. Note that the value was observed to be *something* of type
 747 T_C , though not necessarily the same value that was in that location before the C call which
 748 initially necessitated the addition of the taint.

749 The following reduction rule deals with failing access:

$$\frac{\sigma_C(n+o) = (v, T_C, \beta)}{\mathbf{cget}(\mathbf{ptr} n T'_C) o T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{err} \beta / \sigma_T / \sigma_C / \sigma_V} \text{ (R_CGET_TAINT_FAILS)}$$

750 Here, the access fails, reporting the taint information identifying the call which tampered
 751 with this data. Note that satisfaction of rule R_CGET_TAINT_WORKS's preconditions
 752 implies satisfaction of this rule's preconditions—this ensures that access to tainted locations
 753 can fail in any situation that it can succeed.

754 Similar to **cget**, **cset** has nondeterministic semantics when dealing with tainted locations.
 755 First, consider writes to clean locations:

$$\frac{\sigma_C(n+o) = (v, T_C, \emptyset) \quad \text{value}(v_2) \quad v_{put} = \text{coerceToC}(v_2) \quad \sigma'_C = \text{update}(\sigma_C, n+o, (v_{put}, T_C, \emptyset))}{\mathbf{cset}(\mathbf{ptr}_L n T'_C) o v_2 T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_2 / \sigma_T / \sigma'_C / \sigma_V} \text{(R_CSET_NO_TAINT)}$$

756 In the expression $\mathbf{cset}(\mathbf{ptr}_L n T'_C) o v_2 T_C$, we write v_2 to location n with offset o in σ_C ,
 757 and we expect the location to have type T_C . Since location $n+o$ in σ_C is clean, the store
 758 update cannot fail.

759 Note that we must first coerce v_2 to a C value v_{put} to store it in σ_C . $\text{coerceToC}(v_2)$ is
 760 similar to the coerceToLua function, though it coerces Lua values to C instead. For example,
 761 if v_2 is a numeric constant, the function produces a C integer with the same numeric value,
 762 and if v_2 is a Lua pointer $\mathbf{ptr}_L m ct$, an equivalent C pointer $\mathbf{ptr}_C m ct$ is produced.

763 The rule for **csets** on tainted locations is given below:

$$\frac{\sigma_C(n+o) = (v, T_C, \beta) \quad \text{value}(v_2) \quad v_{put} = \text{coerceToC}(v_2) \quad \sigma'_C = \text{update}(\sigma_C, n+o, (v_{put}, T_C, \emptyset))}{\mathbf{cset}(\mathbf{ptr}_L n T'_C) o v_2 T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_2 / \sigma_T / \sigma'_C / \sigma_V} \text{(R_CSET_TAINT_WORKS)}$$

764 Here, we again coerce v_2 to a C value v_{put} to location n with offset o in σ_C , and we
 765 expect the location to have type T_C . However, $\sigma_C(n+o)$ is tainted, and so we are in a
 766 state of nondeterminism. In rule **R_CSET_TAINT_WORKS**, the write succeeds: We update
 767 $\sigma_C(n+o)$ with the new value v_{put} and clean the taint. Of course, failure is always an option:

$$\frac{\sigma_C(n+o) = (v, T_C, \beta)}{\mathbf{cset}(\mathbf{ptr}_L n T'_C) o v_2 T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{err} \beta / \sigma_T / \sigma_C / \sigma_V} \text{(R_CSET_TAINT_FAILS)}$$

768 In this parallel case to **T_CSET_TAINT_WORKS**, the write fails, and reports the taint
 769 information stored at $\sigma_C(n+o)$.

770 By now, we have explored each of the reduction relations related to Poseidon Lua's C
 771 FFI. In Section 3, we claimed that even without a model of C, as is the case in our system,
 772 the merger of the type systems of C and Typed Lua allows us to prove meaningful and
 773 interesting results about the language as a whole. The next section presents the results which
 774 we have proved, and sketches the proofs.

775 4.5 Proofs

776 There are two major results that we would like to prove about our semantics of Poseidon
 777 Lua. First, we would like to show some form of *soundness*, though clearly we can't have
 778 traditional type safety due to interoperation with C. Even so, we designed our semantics in
 779 such a way as to *track* C's effect on the overall system, and we can leverage that to show
 780 (conditional) soundness of the host language. Note that our proofs are mechanized in Coq,
 781 and this code is included in the artifact; a brief sketch of each proof is given here, but for
 782 the full details refer to the code.

783 We start with a sketch of preservation.

784 **► Theorem 1 (Preservation).** *For all $K, te, T, e, S,$ and S' such that $\{\}, K \vdash te : T \rightsquigarrow e,$
 785 $e / S \rightarrow e' / S', S$ is well-typed with respect to $K,$ and both environments S and S' are clean,
 786 then there exists store typing $K',$ typed expression $te',$ and type T' such that $\{\}, K' \vdash te' : T' \rightsquigarrow e'$
 787 $T' \rightsquigarrow e'$ with $T' <: T, S'$ is well-typed with respect to $K',$ and K' extends $K.$*

788 **Proof sketch.** Standard proof by induction on the typing derivation $\{\}, K \vdash te : T \rightsquigarrow e$.
 789 Any case where the error expression is reached is in violation of the runtime environments S
 790 and S' being clean, as taint is required in order to get an error. ◀

791 Essentially, the statement of preservation for Poseidon Lua differs from traditional
 792 statements of preservation in the stipulation that the runtime environments S and S' be
 793 clean. Clean environments ensure that the C error expression cannot be reached, and that
 794 the semantics are deterministic, as it's the presence of taint which begets nondeterminism.
 795 We can similarly show progress.

796 ▶ **Theorem 2 (Progress).** *For all K, te, T, e , and S such that $\{\}, K \vdash te : T \rightsquigarrow e$, S is
 797 well-typed with respect to K , and S is clean, then either e is a value, or there exists clean
 798 environment S' , and expression e' such that $e / S \rightarrow e' / S'$.*

799 **Proof sketch.** Another standard proof by induction on the typing derivation $\{\}, K \vdash te :$
 800 $T \rightsquigarrow e$. As with preservation, any case where the error expression is reached is in violation
 801 of the runtime environments S and S' being clean. ◀

802 As was the case in preservation, the statement of progress here is distinguished by the
 803 requirement that runtime environment S be clean. With a clean S' , progress connects
 804 cleanly with preservation, allowing us to show soundness of Poseidon Lua contingent on clean
 805 environments. A sketch of soundness follows.

806 ▶ **Theorem 3 (Soundness).** *For all K, te, T, e , and S such that $\{\}, K \vdash te : T \rightsquigarrow e$, either
 807 e diverges, or there exists clean environment S' , and value v such that $e / S \rightarrow^* v / S'$ and
 808 all intermediate environments are clean.*

809 **Proof sketch.** A standard proof, which basically amounts to applications of progress and
 810 preservation, and the intermediate environment of each step in the chain of reductions is
 811 guaranteed to be clean by construction (in a sense, progress generates clean environments). ◀

812 Roughly speaking, Theorem 3 states that Poseidon Lua programs in clean environments
 813 do not get stuck. The restriction to clean environments is due to the guest language, C,
 814 potentially interfering with the host language: C calls taint the environment, and accessing
 815 tainted values can lead to a stuck state even in well-typed programs. This isn't to say
 816 that you can't use C at all, as allocating simple pointers and structs does not taint the
 817 environment, and it is equally valid if some taint was once present and had been cleaned by
 818 successful accesses or writes.

819 Unfortunately, our statement of soundness doesn't say much for the realistic use case
 820 of Poseidon Lua (and C FFIs in general), as these systems are designed to call C code.
 821 That said, we are not without options: as before, our inclusion of taint allows us to reason
 822 about C's effects on the overall language. Crucially, failing C reductions result in the error
 823 expression **err** β , and the taint information β can be used to identify the true culprit for the
 824 crash, even if that culprit was some earlier, seemingly unrelated expression. In short, we can
 825 show that C is to blame for failures in well-typed Poseidon Lua programs.

826 ▶ **Theorem 4 (Always Blame C).** *If the error expression **err** β is reached, then there exists
 827 some C expression which is to blame.*

828 **Proof sketch.** Effectively, this can be shown by construction of our semantics. **err** β can
 829 only be reached through reduction from a C expression, and the only way that such a
 830 reduction can occur is if there was some taint in the runtime environment. In **err** β , β is

1	p = calloc Point	p = calloc Point
2	cCall1(p)	cCall1(p)
3	cCall2(p)	print(p.x)
4	cCall3(p)	cCall2(p)
5	print(p.x)	print(p.x)
6		cCall3(p)
7		print(p.x)

■ **Figure 4** Illustrative example.

831 taint information which identifies some C call, cast, or allocation (as those are the only
832 expressions which can taint), and it's the identified expression that will be blamed. ◀

833 At a high level, Theorem 4 indicates that runtime errors in well-typed Poseidon Lua
834 are attributable to C. This signifies that our interoperation scheme does not allow for any
835 additional errors which are the fault of the host language, and any errors introduced by the
836 C FFI can be traced back to C.

837 Taken together, Theorems 3 and 4 are analogous to soundness of static code and the
838 gradual guarantee in gradually typed languages [26][23], though the context is otherwise quite
839 different. This similarity betrays a certain connection between gradual typing and language
840 interoperation, a connection equally noted by aforementioned work on linking types [21].

841 As we know, program execution in a tainted environment is nondeterministic. In this
842 state, many executions are possible, and they can be categorized as follows: the program
843 either terminates successfully, terminates unsuccessfully, or it executes until the environment
844 is cleaned of taint. Interestingly, executions which clean the taint actually *reclaim* soundness,
845 and are deterministic at least until the next C call.

846 We can show one other interesting result about Poseidon Lua programs which call C.
847 First, recall that only clean locations gain taint when a C call occurred; this ensures proper
848 error tracking in the event of multiple C calls possibly tainting the same data. For an
849 illustrative example, consider the code in Figure 4.

850 Assume the leftmost program fails at the access to `p.x`, blaming `cCall1` and identifying
851 it as the start of our search; here, we cannot say for sure which of `cCall1`, `cCall2`, or `cCall3`
852 mucked with `p.x`. However, we can generate a modified program which can isolate the
853 faulty C call. Consider the snippet on the right. If `cCall1` was the culprit of the failure,
854 then the access immediately following it will fail. If not, and `cCall2` was at fault, then the
855 access immediately after `cCall2` will fail. If neither of these are true, then `cCall3` is at fault,
856 causing the final access to `p.x` to fail. This amounts to fault localization: When we are
857 uncertain about which of a number of unsafe operations are at fault for a runtime failure, we
858 can generate a new program which isolates the faulty operation.

859 5 Poseidon Lua: Implementation

860 As a demonstration of the practicality of these semantics, they have been implemented as
861 modifications to Lua 5.3.3 [13] and Typed Lua [14]. Lua is extended to provide low-level
862 interfaces, and Typed Lua is extended to make use of them with C types. The extensions to

863 Lua have no guarantees of safety or correctness on their own, and are treated as an internal
 864 implementation language for the modifications of Typed Lua. Typed Lua is extended with C
 865 types, through the addition of a C pointer in Lua which refers to C data (as explained in
 866 Section 4.1).

867 Typed Lua’s grammar is extended as follows:

```

868 T ::= (all existing Typed Lua types) | PtrType
869
870 PtrType ::= ptr ptr* PtrTargetType
871
872 PtrTargetType ::= CVoidType | CPrimitiveType | Name
873
874 CType ::= CPrimitiveType | PtrType
875
876 CVoidType ::= void
877
878 CPrimitiveType ::= char | int | double
879
880 Statement ::= (all existing Typed Lua statements) | StructDeclaration
881
882 StructDeclaration ::= struct Name StructIdDecList end
883
884 StructIdDecList ::= StructIdDec StructIdDec*
885
886 StructIdDec ::= Id : CType
887
888 Expression ::= (all existing Typed Lua expressions) | CallocExpr
889
890 CallocExpr ::= calloc ( PtrTargetType )

```

891 T, in particular, is the existing Typed Lua non-terminal for types. As a consequence, any
 892 variable, parameter or field in Poseidon Lua may contain a *pointer* to a C value, but may not
 893 contain a C value directly. All other types are unmodified, and behave as they do in Typed
 894 Lua. As in C, the Poseidon Lua compiler assures that every type named in a C pointer type
 895 has a corresponding struct declaration, and that no name corresponds to multiple structure
 896 declarations, and as in C, the struct declaration defines the memory layout of objects of that
 897 type. Unlike in C, declarations are not required to precede uses of the type they declare. A
 898 simple wrapper for `calloc` is provided to assure that allocations are always of the correct
 899 size. For this prototype, we implemented only `chars`, `ints` and `doubles`, but there is no
 900 conceptual limitation on implementing any other primitive type. For convenience, Poseidon
 901 Lua also provides syntax and semantics for C arrays, but they are not discussed in this work.

902 This modified Typed Lua compiles to Lua, extended with intrinsics to manipulate memory
 903 directly. Typed Lua code which doesn’t use C features is unchanged: That is, if C `ptrs`
 904 are not used, `calloc` is not used, and the code passes type checking, then it compiles into
 905 identical Lua code without type annotations or declarations (i.e. the types are erased).
 906 Lua already provides a datatype, “light user data”, intended for storing pointers to C data,
 907 and this datatype is used for all `ptr`-typed variables and fields. This is why Lua was used
 908 for this prototype. However, Lua’s light user data is completely opaque to Lua code: In
 909 order to use it, one must implement a C interface, from which the underlying pointers are

```

struct House
    num_rooms : int
end
local house_1 : ptr House = calloc(House)
house_1.num_rooms = 6

```

■ **Figure 5** Simple Poseidon Lua code example

```

local house_1 = CS_calloc(4)
CS_storeInt(house_1, 0, 6)

```

■ **Figure 6** Simple Lua code example compiled from Poseidon Lua

910 exposed. Our principle extensions to Lua are low-level operators to directly manipulate
 911 memory through these pointers: `CS_loadChar`, `CS_storeChar`, and similar for ints, doubles
 912 and pointers. In addition, `CS_calloc` and `CS_free` are provided to give direct access to
 913 C’s `calloc` and `free`, a literal `CS_NULL` corresponding to C’s `NULL` is provided to check for
 914 errors, and `CS_loadString` and `CS_storeString` are provided to convert between C strings
 915 (0-terminated `char` arrays) and Lua strings. “CS” in this context is an abbreviation of “C
 916 Semantics”.

917 Each of these low-level operators converts data between Lua’s native data types and C’s,
 918 given a C pointer stored in a Lua light user data, and an offset. The conversions themselves
 919 are trivial. None of these operators are intended for direct use by end users. Instead, Poseidon
 920 Lua’s Typed Lua implementation compiles code which uses C types—that is, code which
 921 accesses members of `ptr`-typed variables or fields—to Lua which uses the correct operators.
 922 Internally, each low-level operator is compiled to its own opcode in Lua’s bytecode.

923 As a simple example, the Poseidon Lua in Figure 5 compiles to the Lua in Figure 6.

924 As the changes in our semantics are concerned principally with C data, rather than C
 925 functions, we use a modified `luaiffb` for the function component of the interface. Poseidon
 926 Lua’s modified `luaiffb` is changed only by replacing their wrapper objects with Lua’s light
 927 user data, which can then be handled by Typed Lua types. The jump between C and
 928 Lua code incurs much less overhead than wrapping C data for use in Lua, so no further
 929 modifications are necessary.

930 5.1 Performance

931 Poseidon Lua code which doesn’t use C types is just regular Typed Lua: when compiled into
 932 Lua code this will be identical to the equivalent Typed Lua program being compiled into Lua,
 933 and so will not display any performance difference. Thus, to compare the performance of
 934 Poseidon Lua against `luaiffb`, we need benchmarks which particularly measure the access to
 935 structured data. Unfortunately, we know of no benchmark suite intended specifically for this
 936 purpose, so instead we ported four benchmarks from the Computer Language Benchmarks
 937 Game [5]. The subset of benchmarks from CLBG were selected because they had Lua versions
 938 which used structured data types. In each case, they were rewritten so that every structured
 939 datatype used a C struct, the shape of which was taken from the C version of each benchmark.
 940 In Poseidon Lua, these structs were represented as `struct` declarations, and in `luaiffb`, as
 941 their dynamic declarations. In both cases, no actual C calls are made: The data is stored in
 942 C-compatible structures, and accessed through them, but the benchmark code is entirely

Benchmark	Poseidon Lua		luaffib		Lua	
	Time (s)	Std. Dev.	Time (s)	Std. Dev.	Time (s)	Std. Dev.
binary-trees	18.8	0.447	202.4	2.97	22.0	0.707
n-body	4.0	0	40.6	1.14	4.0	0.707
spectral-norm	108.2	0	270.8	2.59	105.6	0.894
fannkuch-redux	66.8	2.95	528.8	9.68	55.0	0

■ **Figure 7** Comparison of performance results over various benchmarks

943 Lua. We compare the performance of luaffib, which uses wrappers, to Poseidon Lua, which
 944 does not. We also include the original Lua benchmark, which does not use C structured
 945 data, for reference, although we expect no significant performance difference with respect
 946 to it. The results and standard deviations are shown in Figure 7. As expected, Poseidon
 947 Lua shows a substantial speedup over luaffib, due to the absence of allocated wrappers at
 948 runtime. Our performance is close to original Lua, though in some benchmarks the cost of
 949 converting between C’s primitive types and Lua’s overwhelms other benefits.

950 The benchmarks were performed on Lua 5.3.3 as well as our modified version thereof, on
 951 a quad-core 1.8GHz 64-bit Intel desktop PC running Ubuntu 14.04.3LTS.

952 6 Conclusions

953 In this paper, we presented a framework for reasoning about C FFI’s without fully modelling
 954 the guest language. This framework relies on making the data interface of the FFI static by
 955 combining the type systems of the host and guest languages, and doesn’t require a model of
 956 the guest language beyond its direct interactions with the host. We also saw how making the
 957 data interface static eliminates the need for burdensome wrappers in FFI implementations,
 958 as the host language can statically check its own use of the FFI instead of needing to rely on
 959 the dynamic checks in the wrappers.

960 To showcase our framework, we presented Poseidon Lua, a Typed Lua C FFI. We gave the
 961 formal semantics of the C FFI in Poseidon Lua, and even without modelling C were able to
 962 guarantee some level of soundness of the host language, as well as prove that well-typed host
 963 language code is not to blame for errors that occur. We also presented an implementation of
 964 Poseidon Lua, and confirmed that making the data interface static does indeed improve the
 965 performance of the FFI.

966 While we focus on a C FFI, in principle our approach also works for other choices of guest
 967 language, as we deliberately avoid modelling C. That said, our model of C’s memory and
 968 C’s types in the host language make languages with similar memory behavior to C’s most
 969 suitable, though one could plug in any type system and model memory differently if they are
 970 so inclined. We focused on a C FFI because they are very common, and prove particularly
 971 challenging to reason about with traditional methods.

972 ——— References ———

- 973 1 M. Abadi, L. Cardelli, B. Pierce, and G. Plotkin. Dynamic typing in a statically-typed
 974 language. In *Proceedings of the 16th ACM SIGPLAN-SIGACT Symposium on Principles
 975 of Programming Languages*, POPL ’89, pages 213–227, New York, NY, USA, 1989. ACM.
 976 URL: <http://doi.acm.org/10.1145/75277.75296>, doi:10.1145/75277.75296.
- 977 2 CompCert. CompCert Main Page. <http://compcert.inria.fr>. Accessed: 2018-07-23.

- 978 3 Facebook. luaffib. <https://github.com/facebookarchive/luaffib>. Accessed: 2019-
979 01-10.
- 980 4 Matthias Felleisen, Robert Bruce Findler, and Matthew Flatt. *Semantics engineering with*
981 *PLT Redex*. Mit Press, 2009.
- 982 5 Isaac Gouy. The Computer Language Benchmarks Game. [https://benchmarksgame-team.](https://benchmarksgame-team.pages.debian.net/benchmarksgame/)
983 [pages.debian.net/benchmarksgame/](https://benchmarksgame-team.pages.debian.net/benchmarksgame/). Accessed: 2019-01-10.
- 984 6 Kathryn E Gray. Safe cross-language inheritance. In *European Conference on Object-*
985 *Oriented Programming*, pages 52–75. Springer, 2008.
- 986 7 Arjun Guha, Claudiu Saftoiu, and Shriram Krishnamurthi. The essence of JavaScript. In
987 *European Conference on Object-Oriented Programming*, pages 126–150. Springer, 2010.
- 988 8 Atsushi Igarashi, Benjamin C Pierce, and Philip Wadler. Featherweight Java: a minimal
989 core calculus for Java and GJ. *ACM Transactions on Programming Languages and Systems*
990 (*TOPLAS*), 23(3):396–450, 2001.
- 991 9 Julia. Calling C and Fortran Code. [https://docs.julialang.org/en/stable/manual/](https://docs.julialang.org/en/stable/manual/calling-c-and-fortran-code/index.html)
992 [calling-c-and-fortran-code/index.html](https://docs.julialang.org/en/stable/manual/calling-c-and-fortran-code/index.html). Accessed: 2018-07-14.
- 993 10 Hanshu Lin. Operational semantics for Featherweight Lua. *Master’s Projects*, page 387,
994 2015.
- 995 11 Lisp. CFFI The Common Foreign Function Interface. [https://common-lisp.net/](https://common-lisp.net/project/cffi/)
996 [project/cffi/](https://common-lisp.net/project/cffi/). Accessed: 2018-07-25.
- 997 12 Tidal Lock. Tidal Lock Gradual Static Typing for Lua. [https://github.com/fab13n/](https://github.com/fab13n/metalua/tree/tilo/src/tilo)
998 [metalua/tree/tilo/src/tilo](https://github.com/fab13n/metalua/tree/tilo/src/tilo). Accessed: 2018-06-20.
- 999 13 Lua. Lua 5.3 documentation. <https://www.lua.org/manual/5.3/>. Accessed: 2018-06-20.
- 1000 14 André Murbach Maidl, Fabio Mascarenhas, and Roberto Ierusalimsky. Typed Lua: An
1001 optional type system for lua. In *Proceedings of the Workshop on Dynamic Languages and*
1002 *Applications*, pages 1–10. ACM, 2014.
- 1003 15 MathWorks. Matlab Calling C Shared Libraries. [https://www.mathworks.com/help/](https://www.mathworks.com/help/matlab/using-c-shared-library-functions-in-matlab-.html)
1004 [matlab/using-c-shared-library-functions-in-matlab-](https://www.mathworks.com/help/matlab/using-c-shared-library-functions-in-matlab-.html).html. Accessed: 2018-07-04.
- 1005 16 Jacob Matthews and Robert Bruce Findler. Operational semantics for multi-language pro-
1006 grams. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 31(3):12,
1007 2009.
- 1008 17 Microsoft. Typescript – language specification version 1.8. Technical report, Microsoft,
1009 January 2016.
- 1010 18 James Newsome and Dawn Xiaodong Song. Dynamic taint analysis for automatic detection,
1011 analysis, and signature generation of exploits on commodity software. In *NDSS*, volume 5,
1012 pages 3–4. Citeseer, 2005.
- 1013 19 Graham Ollis. Perl Perl Foreign Function Interface based on GNU fcall. [https:](https://metacpan.org/pod/FFI)
1014 [//metacpan.org/pod/FFI](https://metacpan.org/pod/FFI). Accessed: 2018-07-06.
- 1015 20 Mike Pall. The LuaJIT project. *Web site: http://luajit.org*, 2008.
- 1016 21 Daniel Patterson and Amal Ahmed. Linking Types for Multi-Language Software: Have
1017 Your Cake and Eat It Too. In *2nd Summit on Advances in Programming Languages (SNAPL*
1018 *2017)*, volume 71 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 12:1–
1019 12:15, 2017. doi:10.4230/LIPIcs.SNAPL.2017.12.
- 1020 22 Python. CFFI Documentation. <https://cffi.readthedocs.io/en/latest/>. Accessed:
1021 2018-07-06.
- 1022 23 Jeremy G Siek, Michael M Vitousek, Matteo Cimini, and John Tang Boyland. Refined
1023 criteria for gradual typing. In *LIPIcs-Leibniz International Proceedings in Informatics*,
1024 volume 32. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2015.
- 1025 24 Mallku Soldevila, Beta Ziliani, Bruno Silvestre, Daniel Fridlender, and Fabio Mascarenhas.
1026 Decoding Lua: Formal semantics for the developer and the semanticist. In *Proceedings*
1027 *of the 13th ACM SIGPLAN International Symposium on Dynamic Languages, DLS 2017*,

- 1028 pages 75–86, New York, NY, USA, 2017. ACM. URL: <http://doi.acm.org/10.1145/3133841.3133848>, doi:10.1145/3133841.3133848.
- 1029
- 1030 **25** Sam Tobin-Hochstadt, Vincent St-Amour, Eric Dobson, and Asumu Takikawa. Typed
- 1031 Racket Documentation. <https://docs.racket-lang.org/ts-guide/>. Accessed: 2018-08-
- 1032 01.
- 1033 **26** Philip Wadler and Robert Bruce Findler. Well-typed programs can't be blamed. In
- 1034 *European Symposium on Programming*, pages 1–16. Springer, 2009.
- 1035 **27** Alon Zakai. Emscripten: an LLVM-to-JavaScript compiler. In *Proceedings of the ACM*
- 1036 *international conference companion on Object oriented programming systems languages and*
- 1037 *applications companion*, pages 301–312. ACM, 2011.

1038 **A** Appendix

1039 **A.1** Full Typing Rules

$$\frac{\text{validType}(T_C) \quad \beta \text{ unused}}{\Gamma, K \vdash \mathbf{calloc} T_C : \mathbf{ptr}_L T_C \rightsquigarrow \mathbf{calloc} T_C \beta} \quad (\text{TT_CALLOC})$$

$$\frac{n < \text{length}(\Sigma_C) \quad \text{goodLayout}(n, T, \Sigma_C)}{\Gamma, (\Sigma_T, \Sigma_C, \Sigma_V) \vdash \mathbf{ptr}_L n T_C : \mathbf{ptr}_L T_C \rightsquigarrow \mathbf{ptr}_L n T_C} \quad (\text{TT_LUA_PTR})$$

$$\frac{\Gamma, K \vdash te : \mathbf{ptr}_L T_C \rightsquigarrow e \quad \text{validForCDeref}(T_C) \quad T_L = \text{coerceCType}(T_C)}{\Gamma, K \vdash \mathbf{deref}_C te : T_L \rightsquigarrow \mathbf{cget} e 0 T_C} \quad (\text{TT_VAR_C_DEREF})$$

$$\frac{}{\Gamma, K \vdash \mathbf{cfun} (ct_1 \rightarrow_C ct_2) : (ct_1 \rightarrow_C ct_2) \rightsquigarrow \mathbf{cfun}} \quad (\text{TT_C_FUNCTION})$$

$$\frac{\Gamma, K \vdash te_1 : (T \rightarrow_C T') \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : T \rightsquigarrow e_2 \quad \beta \text{ unused}}{\Gamma, K \vdash te_1(te_2) : T' \rightsquigarrow \mathbf{ccall} e_1 e_2 T' \beta} \quad (\text{TT_C_FUN_APPL})$$

$$\frac{\Gamma, K \vdash te_1 : \mathbf{ptr}_L T_1 \rightsquigarrow e_1 \quad \text{structType}(T_1) \quad \Gamma, K \vdash te_2 : s \rightsquigarrow e_2 \quad s \in T_1 \quad n = \text{offsetForType}(s, T_1)}{\Gamma, K \vdash te_1.te_2 : T_1(s) \rightsquigarrow \mathbf{cget} e_1 n T_1(s)} \quad (\text{TT_C_DOT_ACCESS})$$

$$\frac{\Gamma, K \vdash te_1 : \mathbf{ptr}_L T_1 \rightsquigarrow e_1 \quad \text{structType}(T_1) \quad \Gamma, K \vdash te_2 : s \rightsquigarrow e_2 \quad \Gamma, K \vdash te_3 : T_1(s) \rightsquigarrow e_3 \quad s \in T_1 \quad n = \text{offsetForType}(s, T_1)}{\Gamma, K \vdash te_1.te_2 := te_3 : \mathbf{value} \rightsquigarrow \mathbf{cset} e_1 n e_2 T_1(s)} \quad (\text{TT_C_DOT_UPDATE})$$

$$\frac{\Gamma, K \vdash te : \mathbf{ptr}_L T'_C \rightsquigarrow e \quad \beta \text{ unused}}{\Gamma, K \vdash \mathbf{ccast} te T_C : T_C \rightsquigarrow \mathbf{ccast} e T_C \beta} \quad (\text{TT_C_CAST})$$

$$\frac{\forall i, f_i = s_i : T_i \vee f_i = \mathbf{const} s_i : T_i \quad \forall i, \Gamma, K \vdash tv_i : T_i \rightsquigarrow v_i}{\Gamma, K \vdash \{s_1 = tv_1, \dots, s_n = tv_n\} : \{f_1, \dots, f_n\} \rightsquigarrow \{s_1 = v_1, \dots, s_n = v_n\}} \quad (\text{TT_TABLE})$$

$$\frac{\Gamma, K \vdash te_1 : T \rightsquigarrow e_1 \quad \Gamma + \{x \mapsto T\}, K \vdash te_2 : T' \rightsquigarrow e_2}{\Gamma, K \vdash \mathbf{let} \ x : T := te_1 \ \mathbf{in} \ te_2 : T' \rightsquigarrow \mathbf{let} \ x := e_1 \ \mathbf{in} \ e_2} \quad (\text{TT_LET})$$

$$\frac{x \in \Gamma}{\Gamma, K \vdash x : \Gamma(x) \rightsquigarrow x} \quad (\text{TT_VAR})$$

$$\frac{n < \text{length}(\Sigma_T)}{\Gamma, (\Sigma_T, \Sigma_C, \Sigma_V) \vdash \mathbf{reg} \ n : \Sigma_T(n) \rightsquigarrow \mathbf{reg} \ n} \quad (\text{TT_REG})$$

$$\frac{\Gamma, K \vdash te : \mathbf{ref} \ T \rightsquigarrow e}{\Gamma, K \vdash \mathbf{deref} \ te : T \rightsquigarrow \mathbf{deref} \ e} \quad (\text{TT_VAR_DEREF})$$

$$\frac{x \in \Gamma \quad \Gamma, K \vdash te : \Gamma(x) \rightsquigarrow e}{\Gamma, K \vdash x := te : T \rightsquigarrow x := e} \quad (\text{TT_VAR_ASSIGN})$$

$$\frac{\Gamma, (\Sigma_T, \Sigma_C, \Sigma_V) \vdash te : \Sigma_V(n) \rightsquigarrow e}{\Gamma, (\Sigma_T, \Sigma_C, \Sigma_V) \vdash \mathbf{loc} \ n := te : T \rightsquigarrow \mathbf{loc} \ n := e} \quad (\text{TT_LOC_UPDATE})$$

$$\frac{\Gamma + \{x \mapsto T\}, K \vdash te : T' \rightsquigarrow e}{\Gamma, K \vdash \lambda x : T. te : (T \rightarrow_L T') \rightsquigarrow \lambda x. e} \quad (\text{TT_FUNCTION})$$

$$\frac{\Gamma, K \vdash te_1 : (T \rightarrow_L T') \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : T \rightsquigarrow e_2}{\Gamma, K \vdash te_1(te_2) : T' \rightsquigarrow e_1(e_2)} \quad (\text{TT_LUA_FUN_APPL})$$

$$\frac{\Gamma, K \vdash te_1 : T_1 \rightsquigarrow e_1 \quad \text{tableType}(T_1) \quad \Gamma, K \vdash te_2 : s \rightsquigarrow e_2 \quad s \in T_1}{\Gamma, K \vdash te_1.te_2 : T_1(s) \rightsquigarrow \mathbf{rawget} \ e_1 \ e_2} \quad (\text{TT_DOT_ACCESS})$$

$$\frac{\Gamma, K \vdash te_1 : T_1 \rightsquigarrow e_1 \quad \text{tableType}(T_1) \quad \Gamma, K \vdash te_2 : s \rightsquigarrow e_2 \quad s \in T_1 \quad \Gamma, K \vdash te_3 : T_1(s) \rightsquigarrow e_3}{\Gamma, K \vdash te_1.te_2 := te_3 : \mathbf{value} \rightsquigarrow \mathbf{rawset} \ e_1 \ e_2 \ e_3} \quad (\text{TT_DOT_UPDATE})$$

$$\frac{\Gamma, K \vdash te : T \rightsquigarrow e \quad T <: T'}{\Gamma, K \vdash te : T' \rightsquigarrow e} \quad (\text{TT_SUBSUMPTION})$$

$$\frac{c \text{ constant}}{\Gamma, K \vdash c : c \rightsquigarrow c} \quad (\text{TT_CONST})$$

$$\frac{\Gamma, K \vdash te_1 : \mathbf{number} \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : \mathbf{number} \rightsquigarrow e_2 \quad op \in \{+, -, *, /\}}{\Gamma, K \vdash te_1 \ op \ te_2 : \mathbf{number} \rightsquigarrow e_1 \ op \ e_2} \quad (\text{TT_BINOP_ARITH})$$

$$\frac{\Gamma, K \vdash te_1 : \mathbf{number} \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : \mathbf{number} \rightsquigarrow e_2}{\Gamma, K \vdash te_1 \text{ op } te_2 : \mathbf{boolean} \rightsquigarrow e_1 \text{ op } e_2} \quad (\text{TT_BINOP_ORDER})$$

$$\frac{\Gamma, K \vdash te_1 : \mathbf{boolean} \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : \mathbf{boolean} \rightsquigarrow e_2}{\Gamma, K \vdash te_1 \text{ op } te_2 : \mathbf{boolean} \rightsquigarrow e_1 \text{ op } e_2} \quad (\text{TT_BINOP_BOOLS})$$

$$\frac{\Gamma, K \vdash te_1 : \mathbf{string} \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : T_2 \rightsquigarrow e_2 \quad T_2 \in \{\mathbf{string}, \mathbf{number}\}}{\Gamma, K \vdash te_1 .. te_2 : \mathbf{string} \rightsquigarrow e_1 .. e_2} \quad (\text{TT_BINOP_STRING})$$

$$\frac{\Gamma, K \vdash te_1 : T_1 \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : T_2 \rightsquigarrow e_2}{\Gamma, K \vdash te_1 == te_2 : \mathbf{boolean} \rightsquigarrow e_1 == e_2} \quad (\text{TT_BINOP_EQ})$$

$$\frac{\Gamma, K \vdash te_1 : T_1 \rightsquigarrow e_1 \quad \Gamma, K \vdash te_2 : T_2 \rightsquigarrow e_2}{\Gamma, K \vdash te_1; te_2 : T_2 \rightsquigarrow e_1; e_2} \quad (\text{TT_SEQUENCE})$$

1040 A.2 Full Reduction Rules

$$\frac{n = \text{length}(\sigma_C) \quad \sigma'_C = \sigma_C + \text{layoutTypeAndTaint}(T_C, \beta)}{\mathbf{calloc} T_C \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{ptr} n T_C / \sigma_T / \sigma'_C / \sigma_V} \quad (\text{R_CALLOC})$$

$$\frac{\text{value}(v_2) \quad v = \text{makeValueOfType}(ct_2) \quad \sigma'_C = \text{taintCStore}(\sigma_C, \beta)}{\mathbf{ccall cfun} v_2 ct_2 \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow v / \sigma_T / \sigma'_C / \sigma_V} \quad (\text{R_CCALL_WORKED})$$

$$\frac{\text{value}(v_2) \quad \sigma'_C = \text{taintCStore}(\sigma_C, \beta)}{\mathbf{ccall cfun} v_2 ct_2 \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{err} \beta / \sigma_T / \sigma'_C / \sigma_V} \quad (\text{R_CCALL_FAILED})$$

$$\frac{n < \text{length}(\sigma_C) \quad \sigma_C(n) = (v, T_C, \emptyset) \quad \sigma'_C = \text{update}(\sigma_C, n, (v, T'_C, \beta))}{\mathbf{ccast}(\mathbf{ptr}_L n T_C) T'_C \beta / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{ptr}_L n T'_C / \sigma_T / \sigma'_C / \sigma_V} \quad (\text{R_CCAST})$$

$$\frac{\sigma_C(n+o) = (v_C, T_C, \emptyset) \quad v_{out} = \text{coerceToLua}(v_C)}{\mathbf{cget}(\mathbf{ptr}_L n T'_C) o T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_{out} / \sigma_T / \sigma_C / \sigma_V} \quad (\text{R_CGET_NO_TAINT})$$

$$\frac{\sigma_C(n+o) = (v, T_C, \beta) \quad v' = \text{makeValueOfType}(T_C) \quad \sigma'_C = \text{update}(\sigma_C, n+o, (v', T_C, \emptyset)) \quad v_{out} = \text{coerceToLua}(v')}{\mathbf{cget}(\mathbf{ptr} n T'_C) o T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_{out} / \sigma_T / \sigma_C / \sigma_V} \quad (\text{R_CGET_TAINT_WORKS})$$

$$\frac{\sigma_C(n+o) = (v, T_C, \beta)}{\mathbf{cget}(\mathbf{ptr} n T'_C) o T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow \mathbf{err} \beta / \sigma_T / \sigma_C / \sigma_V} \quad (\text{R_CGET_TAINT_FAILS})$$

XX:30 Reasoning About Foreign Function Interfaces

$$\frac{\sigma_C(n+o) = (v, T_C, \emptyset) \quad \text{value}(v_2) \quad v_{\text{put}} = \text{coerceToC}(v_2) \quad \sigma'_C = \text{update}(\sigma_C, n+o, (v_{\text{put}}, T_C, \emptyset))}{\text{cset}(\text{ptr}_{\mathbf{L}} n T'_C) o v_2 T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_2 / \sigma_T / \sigma'_C / \sigma_V} \text{(R_CSET_NO_TAINT)}$$

$$\frac{\sigma_C(n+o) = (v, T_C, \beta) \quad \text{value}(v_2) \quad v_{\text{put}} = \text{coerceToC}(v_2) \quad \sigma'_C = \text{update}(\sigma_C, n+o, (v_{\text{put}}, T_C, \emptyset))}{\text{cset}(\text{ptr}_{\mathbf{L}} n T'_C) o v_2 T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow v_2 / \sigma_T / \sigma'_C / \sigma_V} \text{(R_CSET_TAINT_WORKS)}$$

$$\frac{\sigma_C(n+o) = (v, T_C, \beta)}{\text{cset}(\text{ptr}_{\mathbf{L}} n T'_C) o v_2 T_C / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{err } \beta / \sigma_T / \sigma_C / \sigma_V} \text{(R_CSET_TAINT_FAILS)}$$

$$\frac{n = \text{length}(\sigma_T) \quad t_n = \text{buildTable}(\{s_1 = v_1, \dots, s_n = v_n\})}{\{s_1 = v_1, \dots, s_n = v_n\} / \sigma_T / \sigma_C / \sigma_V \rightarrow (\text{reg } n) / \sigma_T + t_n / \sigma_C / \sigma_V} \text{(R_TABLE)}$$

$$\frac{\text{value}(e_1) \quad l = \text{length}(\sigma_V)}{\text{let } x := e_1 \text{ in } e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow [x \leftarrow l] e_2 / \sigma_T / \sigma_C / \sigma_V + e_1} \text{(R_LET)}$$

$$\frac{\text{value}(e_2) \quad l = \text{length}(\sigma_V)}{(\lambda x. e)(e_2) / \sigma_T / \sigma_C / \sigma_V \rightarrow [x \leftarrow l] e / \sigma_T / \sigma_C / \sigma_V + e_2} \text{(R_FUN_APP)}$$

$$\frac{\sigma_V(l) = v \quad \text{value}(v)}{\text{deref}(\text{loc } l) / \sigma_T / \sigma_C / \sigma_V \rightarrow v / \sigma_T / \sigma_C / \sigma_V} \text{(R_LOC_DEREF)}$$

$$\frac{\text{value}(e) \quad \sigma'_V = \text{update}(\sigma_V, l, e)}{\text{loc } l := e / \sigma_T / \sigma_C / \sigma_V \rightarrow e / \sigma_T / \sigma_C / \sigma'_V} \text{(R_LOC_UPDATE)}$$

$$\frac{\sigma_T(n) = T \quad T(s) = v}{\text{rawget}(\text{reg } n) s / \sigma_T / \sigma_C / \sigma_V \rightarrow v / \sigma_T / \sigma_C / \sigma_V} \text{(R_RAWGET)}$$

$$\frac{\text{value}(e_3) \quad \sigma_T(n) = T \quad s \in T \quad T' = \text{update}(T, s, e_3) \quad \sigma'_T = \text{update}(\sigma_T, n, T')}{\text{rawset}(\text{reg } n) s e_3 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{reg } n / \sigma'_T / \sigma_C / \sigma_V} \text{(R_RAWSET)}$$

$$\frac{e / \sigma_T / \sigma_C / \sigma_V \rightarrow e' / \sigma'_T / \sigma'_C / \sigma'_V}{x := e / \sigma_T / \sigma_C / \sigma_V \rightarrow x := e' / \sigma'_T / \sigma'_C / \sigma'_V} \text{(R_VAR_ASSIGN_STEP_1)}$$

$$\frac{e / \sigma_T / \sigma_C / \sigma_V \rightarrow e' / \sigma'_T / \sigma'_C / \sigma'_V}{\text{loc } l := e / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{loc } l := e' / \sigma'_T / \sigma'_C / \sigma'_V} \text{(R_LOC_UPDATE_STEP_1)}$$

$$\frac{e_1 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 / \sigma'_T / \sigma'_C / \sigma'_V}{\text{let } x := e_1 \text{ in } e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{let } x := e'_1 \text{ in } e_2 / \sigma'_T / \sigma'_C / \sigma'_V} \text{(R_LET_STEP)}$$

$$\frac{e_1 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 / \sigma'_T / \sigma'_C / \sigma'_V}{\text{rawget } e_1 e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{rawget } e'_1 e_2 / \sigma'_T / \sigma'_C / \sigma'_V} \text{(R_RAWGET_STEP_1)}$$

$$\frac{\text{value}(e_1)}{e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_2 / \sigma'_T / \sigma'_C / \sigma'_V} \text{rawget } e_1 e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{rawget } e_1 e'_2 / \sigma'_T / \sigma'_C / \sigma'_V \quad (\text{R_RAWGET_STEP_2})$$

$$\frac{e_1 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 / \sigma'_T / \sigma'_C / \sigma'_V}{\text{rawset } e_1 e_2 e_3 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{rawset } e'_1 e_2 e_3 / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_RAWSET_STEP_1})$$

$$\frac{\text{value}(e_1)}{e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_2 / \sigma'_T / \sigma'_C / \sigma'_V} \text{rawset } e_1 e_2 e_3 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{rawset } e_1 e'_2 e_3 / \sigma'_T / \sigma'_C / \sigma'_V \quad (\text{R_RAWSET_STEP_2})$$

$$\frac{\text{value}(e_1) \quad \text{value}(e_2)}{e_3 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_3 / \sigma'_T / \sigma'_C / \sigma'_V} \text{rawset } e_1 e_2 e_3 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{rawset } e_1 e_2 e'_3 / \sigma'_T / \sigma'_C / \sigma'_V \quad (\text{R_RAWSET_STEP_3})$$

$$\frac{e_1 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 / \sigma'_T / \sigma'_C / \sigma'_V}{e_1(e_2) / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1(e_2) / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_FUN_APP_STEP_1})$$

$$\frac{\text{value}(e_1)}{e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_2 / \sigma'_T / \sigma'_C / \sigma'_V} \frac{e_1(e_2) / \sigma_T / \sigma_C / \sigma_V \rightarrow e_1(e'_2) / \sigma'_T / \sigma'_C / \sigma'_V}{e_1(e_2) / \sigma_T / \sigma_C / \sigma_V \rightarrow e_1(e'_2) / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_FUN_APP_STEP_2})$$

$$\frac{e_1 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 / \sigma'_T / \sigma'_C / \sigma'_V}{e_1 \text{ op } e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 \text{ op } e_2 / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_BINOP_STEP_1})$$

$$\frac{\text{value}(e_1)}{e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_2 / \sigma'_T / \sigma'_C / \sigma'_V} \frac{e_1 \text{ op } e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e_1 \text{ op } e'_2 / \sigma'_T / \sigma'_C / \sigma'_V}{e_1 \text{ op } e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e_1 \text{ op } e'_2 / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_BINOP_STEP_2})$$

$$\frac{\text{value}(e_1) \quad \text{value}(e_2)}{\text{validL}(e_1) \quad \text{validR}(e_2)} \frac{e_1 \text{ op } e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{evalOp}(e_1, e_2, \text{op}) / \sigma_T / \sigma_C / \sigma_V}{e_1 \text{ op } e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{evalOp}(e_1, e_2, \text{op}) / \sigma_T / \sigma_C / \sigma_V} \quad (\text{R_BINOP})$$

$$\frac{e / \sigma_T / \sigma_C / \sigma_V \rightarrow e' / \sigma'_T / \sigma'_C / \sigma'_V}{\text{cget } e \text{ o } T / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{cget } e' \text{ o } T / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_CGET_STEP})$$

$$\frac{e_1 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 / \sigma'_T / \sigma'_C / \sigma'_V}{\text{cset } e_1 \text{ o } e_2 T / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{cset } e'_1 \text{ o } e_2 T / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_CSET_STEP_1})$$

$$\frac{\text{value}(e_1)}{e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_2 / \sigma'_T / \sigma'_C / \sigma'_V} \text{cset } e_1 \text{ o } e_2 T / \sigma_T / \sigma_C / \sigma_V \rightarrow \text{cset } e_1 \text{ o } e'_2 T / \sigma'_T / \sigma'_C / \sigma'_V \quad (\text{R_CSET_STEP_2})$$

$$\frac{e_1 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1 / \sigma'_T / \sigma'_C / \sigma'_V}{e_1; e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e'_1; e_2 / \sigma'_T / \sigma'_C / \sigma'_V} \quad (\text{R_SEQ_STEP_1})$$

$$\frac{\text{value}(e_1)}{e_1; e_2 / \sigma_T / \sigma_C / \sigma_V \rightarrow e_2 / \sigma_T / \sigma_C / \sigma_V} \quad (\text{R_SEQ_STEP_THROUGH})$$