

A Typed C11 Semantics for Interactive Theorem Proving

Robbert Krebbers Freek Wiedijk

ICIS, Radboud University Nijmegen, The Netherlands

mail@robbertkrebbers.nl freek@cs.ru.nl

Abstract

We present a semantics of a significant fragment of the C programming language as described by the C11 standard. It consists of a small step semantics of a core language, which uses a structured memory model to capture subtleties of C11, such as strict-aliasing restrictions related to unions, that have not yet been addressed by others. The semantics of actual C programs is defined by translation into this core language. We have an explicit type system for the core language, and prove type preservation and progress, as well as type correctness of the translation.

Due to unspecified order of evaluation, our operational semantics is non-deterministic. To explore all defined and undefined behaviors, we present an executable semantics that computes a stream of finite sets of reachable states. It is proved sound and complete with respect to the operational semantics.

Both the translation into the core language and the executable semantics are defined as Coq programs. Extraction to OCaml is used to obtain a C interpreter to run and test the semantics on actual C programs. All proofs are fully formalized in Coq.

Categories and Subject Descriptors D.3.1 [Programming Languages]: Formal Definitions and Theory; F.3.1 [Logics and Meanings of Programs]: Semantics of Programming Languages

Keywords ISO C11 Standard, Operational Semantics, Executable Semantics, Interactive Theorem Proving, Coq

1. Introduction

Many programs need a high runtime performance, close control of the underlying hardware, to be able to run with a minimal runtime environment, or to be very portable. All these reasons explain the ongoing popularity of the C programming language, especially for embedded systems and systems programming in general. Of course one pays for these benefits by using a low-level language in which it is very easy to make mistakes with potentially disastrous consequences. This is especially shown by the proliferation of malware on computers all around the world. Without the widespread use of C and its derivative languages, malware would have a much harder time tainting systems everywhere.

An interesting approach to remedy this situation is to use *proofs* to establish the safety of C programs. That way one gets all performance, control and portability benefits of C, without the dangers

caused by it. Such approaches range from light weight methods like static analysis (which is by nature incomplete) to systems where a user can interactively reason about programs.

Still with systems of the latter kind, it is often unclear whether the system matches the compiler that is used, or might be used in the future, because most of these systems are about a version of C that is quite specific. Moreover, the semantics of C that such a system embodies is generally not made explicit, which makes it hard to establish that it does not contain semantic errors.

For this reason, the CH₂O project [13–18] aims at developing an *explicit* formal semantics that should match the official description of C, namely the C11 standard [11], as closely as possible. Our goal is that if one proves something about a program with respect to the CH₂O semantics, it will behave that way with *any* C11 compliant compiler. There are two projects that are very close:

- The CompCert project by Leroy *et al.* [20] has a formal semantics for a C-like language called CompCert C in Coq. This language can be compiled with a compiler written in Coq, which has been proved correct with respect to this formal semantics.
- Ellison and Rosu [6] have developed a formal semantics of C in the \mathbb{K} -framework. In this work the goal is also to explicitly formally model the C11 standard as closely as possible.

In both projects, as well as in ours, there is a *formal* description of a significant part of C close to the C11 standard and an *executable* interpreter that matches the semantics precisely. However, both projects also differ from our work:

- **Proof infrastructure.** Although Ellison and Rosu support more C features than we do, they do not have proof infrastructure around it. Their semantics, despite being written in a formal framework, should more be seen as a debugger or state space search tool. Our entire semantics, as well as proofs of metatheoretical properties about it, has been formalized in Coq.
- **Explicit typing.** Neither of these projects have an explicit notion of typing, whereas we have defined type judgments. Therefore we have been able to prove properties of our language like type preservation and progress.
- **Formal translation from abstract syntax.** We process the abstract syntax of C to an intermediate language *inside* Coq. We prove that this translation only yields well-typed programs. Although CompCert’s parser has been formally verified [12], the majority of the translation from abstract syntax to CompCert C is performed using unverified OCaml code. This translation involves semantic transformations such as evaluation of constant expressions, that we have proved sound and complete.
- **Core language.** Ellison and Rosu do not have an intermediate language: the semantics operates on the abstract syntax of a program itself. We claim that by first going to a *core* language, the description of various semantic features becomes more princi-

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

CPP ’15, January 13–14 2015, Mumbai, India.

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-3300-9/15/01...\$15.00.

<http://dx.doi.org/10.1145/2676724.2693571>

pled and simpler. For example, this way we do not have duplication in the semantics of various looping constructs.

- **Closer to the C11 standard.** In the CompCert semantics the goal is *not* to make sure that what can be proved according to the semantics will hold for *any* C11 compiler, it just has to hold for one specific compiler: the CompCert compiler. To that end, CompCert is allowed to make choices for implementation defined behavior (*e.g.* integer representations) and gives a semantics to various undefined behaviors (such as aliasing violations).

Like in our work, Ellison and Rosu’s goal is to *exactly* describe the C semantics. However, for some programs their semantics is less precise than ours (see Section 2.4). We can be more precise because in our memory model data objects are structured like trees, while in theirs they consist of flat arrays of bytes.

Approach. To be able to test our semantics, we have developed an ‘interpreter’. This is not an ordinary interpreter that executes a program according to *one* interpretation of the C standard, and that has arbitrary behavior in the case of undefined behavior. Instead, our interpreter calculates *all* behaviors of the program that are allowed by the standard, and when the program has undefined behavior, it will explicitly state this undefinedness.

Our interpreter is also different from compilers or interpreters that insert tests for undefined behaviors as a protection. Those compilers or interpreters generally only follow one possible execution order. Instead, our interpreter is not primarily meant to be a debugging tool, but instead is an *exploration* tool, intended to explore the implications of the C standard [5, 30].

We aim to follow the C11 standard [11] as closely as possible. Unfortunately, the standard is often ambiguous due to its use of natural language (*e.g.*, see the message [24] on the standard’s committee’s mailing list, and Defect Reports #260 and #236 [10]). In the case of ambiguities, we err on the side of caution and favor undefined behavior. Since our interests lie in proving the correctness of C programs with respect to arbitrary C11 compilers, assigning undefined behavior in the case of doubt is the safe choice.

The organization of this paper follows the structure of our C interpreter as shown in Figure 1. We have two languages:

- **CH₂O core C** is described in Section 3. This language is quite abstract, and not very close to actual C.
- **CH₂O abstract C** is described in Section 6. This language is very close to the abstract syntax trees of actual C.

The interpreter then consists of four passes to get from C source code to the behavior of the program.

Related work. Apart from the related work discussed above, there is an abundance of related work on (executable) semantics for large imperative programming languages in proof assistants.

Norrish has defined a semantics of a large part of the C89 type system and semantics in the interactive theorem prover HOL [26]. His main focus was to capture non-determinism in expressions, and our expression semantics is based on his. However, many of the language features we consider are more recent than C89.

Campbell has defined an executable semantics for CompCert C [4] and proved soundness and completeness with respect to the deterministic (leftmost-innermost) operational semantics. The resulting interpreter, that can be used to effectively test the CompCert C semantics, has been reimplemented by Leroy as part of the official CompCert distribution. This reimplementaion is also able to handle non-determinism and explore a program’s whole state space. Since the amount of non-determinism in CompCert C is more restrained than in our semantics (it is at most finitary branching, whereas we have infinitary branching), our completeness proof with respect to the operational semantics is more involved.

Compared to the CompCert C executable semantics, we handle a larger part of the interpreter inside Coq. In particular, our translation from abstract syntax is implemented in Coq, and the machinery to compute the reachable state set is implemented using a verified implementation of hash-sets in Coq.

There have been many previous efforts on creating sound and complete executable semantics. For example, Zhao *et al.* does this for (a deterministic version of) LLVM [31], Bodin *et al.* for Javascript [2], and Lochbihler and Bulwahn for Java [23]. The latter uses code extraction to automatically obtain an executable semantics from an inductive definition. Due to excessive non-determinism in C, our executable semantics has been written by hand.

There has been a lot of related work on tool support for writing definitions of programming languages which can be automatically translated to different proof assistants. Most notably, there are the OTT [29] and Lem [27] tools. Since we wish to use advanced features of Coq, we have not used such tools.

Contribution. In [14, 18] we have introduced an operational and axiomatic semantics for subtle features of a C-like language. Since the focus of these papers was mostly on separation logic, we used a basic memory model that merely supported integers and pointers. In this paper, we extend the operational semantics as follows:

- We integrate our new memory model that supports array, struct, and union types [13, 15]. As usual, such integration suffers from feature interaction leading to many changes to the operational semantics (like introducing new language constructs), as well as the memory model (like extending the permission system).
- We extend the language with a type system and prove that the operational semantics enjoys type presentations and progress.
- We extend the supported C fragment significantly (*e.g.* type-defs, enums, static and global variables, initializers, constant expressions) by defining a translation from a larger fragment of C. This translation is proved type sound.
- We define an executable semantics that can be used to test the semantics against actual C programs. This executable semantics is proved sound and complete. The completeness proof is non-trivial due to excessive non-determinism in our semantics.
- All definitions and proofs have been formalized using Coq. The development can be extracted to OCaml, resulting in an interpreter that can be used to explore all behaviors of C programs.

This paper describes a large formalization effort, and we only have space to show representative parts of definitions. All details can be found online at:

<http://robbertkrebbers.nl/research/ch2o/>.

2. Challenges

The C11 standard gives compilers considerable freedom in what behaviors a program may have [11, 3.4]. It uses the following notions of under-specification:

- **Unspecified behavior:** two or more behaviors are allowed. For example: the execution order of expressions. The choice may vary for each use of the construct.
- **Implementation defined behavior:** like unspecified behavior, but the compiler has to document its choice. For example: size and endianness of integers.
- **Undefined behavior:** the standard imposes no requirements at all, the program is even allowed to crash. For example: dereferencing a NULL pointer or signed integer overflow.

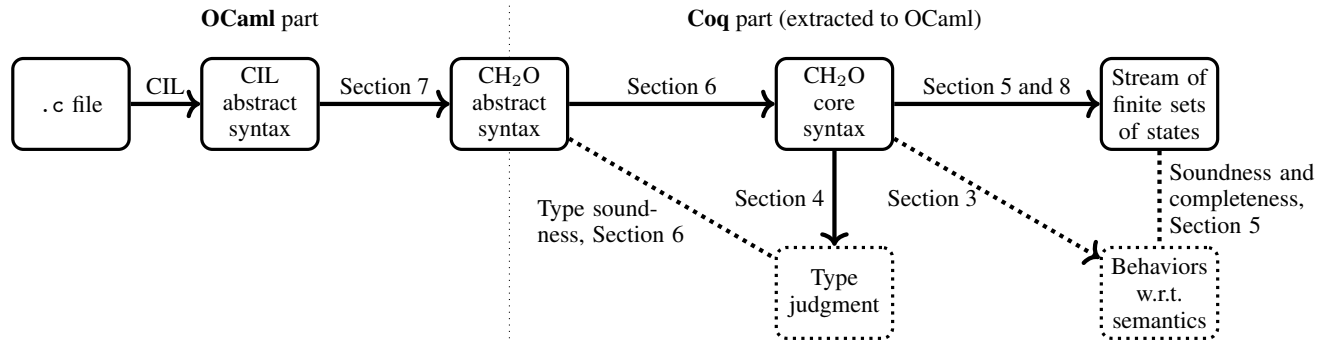


Figure 1. Overview of the architecture of the interpreter.

Under-specification is used extensively to make C portable, and to allow compilers to generate fast code. For example, when dereferencing a pointer, no code has to be generated to check whether the pointer is valid or not. If the pointer is invalid (NULL or a dangling pointer), the compiled program may do something arbitrary instead of having to exit with a nice error message.

Since the CH₂O semantics intends to be a formal version of the C11 standard, it has to capture the behavior of *any* C compiler, and has to take *all* under-specification seriously (even if that makes the semantics complex). In this section, we will describe a number of subtle forms of underspecification and give examples of bizarre behaviors exhibited by actual compilers.

2.1 Sequence point violations and non-determinism

Instead of having to follow a specific execution order, in C the execution order of expressions is unspecified. This is a common cause of portability and maintenance problems because a compiler may use an arbitrary execution order for each expression.

To make more effective optimizations possible (*e.g.* delaying of side-effects and interleaving), the C standard requires the programmer to ensure that all execution orders satisfy a certain condition. If this condition is not met, the program has undefined behavior. Let us consider an example where this is the case:

```
int x;
int y = (x = 3) + (x = 4);
printf("%d %d\n", x, y);
```

By considering all possible execution orders, one would naively expect this program to print 4 7 or 3 7, depending on whether the assignment $x = 3$ or $x = 4$ is executed first. However, the *sequence point restriction* does not allow an object to be modified more than once (or being read after being modified) between two *sequence points* [11, 6.5p2]. A sequence point occurs for example at the end of a full expression, before a function call, and after the first operand of the conditional `? :` operator [11, Annex C]. Hence, both execution orders lead to a sequence point violation, and therefore result in undefined behavior. Indeed, when compiled with `gcc -O2` (version 4.6.4), the above program prints 4 8, which does not correspond to any of the execution orders.

Non-determinism in C is even more unrestrained than some may think. The execution order in $e_1 + e_2$ is unspecified, but this does not mean that either e_1 or e_2 will be executed entirely before the other. Instead, it means that execution can be interleaved; first a part of e_1 , then a part of e_2 , then another part of e_1 , and so on. Hence, the following expression is allowed to print bac:

```
printf("a") + (printf("b") + printf("c"));
```

Our approach to handling non-determinism and sequence points is inspired by Ellison and Rosu [6] and Norrish [26]. Each object in memory carries a permission that is changed into a locked variant whenever a write occurs. Having a locked permission, the memory model prohibits any access (read or write) to this object. At the next sequence point, the permission is changed back into the unlocked variant, making future accesses possible again.

2.2 Non-local control and block scope variables

C allows unrestricted `gotos` which (unlike `break` and `continue`) may not only jump out of blocks, but can also jump into blocks. Blocks may contain local variables, called *block scope variables*, which can be “taken out of their block” by keeping a pointer to them. Leaving a block results in the memory of these variables being freed, and thus making pointers to them invalid. Consider:

```
int *p = NULL;
1: if (p) { return (*p); }
   else { int j = 10; p = &j; goto 1; }
```

When the label 1 is passed for the first time, the variable `p` is NULL. Hence, execution continues in the block where `p` is assigned a pointer to the block scope variable `j`. After the `goto 1`, the block containing `j` is left, and the memory of `j` is freed. Hence, the conditional `if (p)` on a dangling pointer after the `goto` yields undefined behavior.

2.3 Indeterminate memory and pointers

A pointer value becomes indeterminate when the object it points to reaches the end of its lifetime [11, 6.2.4] (goes out of scope, or has been deallocated). Not only dereferencing indeterminate pointers has undefined behavior, but also using such pointers in pointer arithmetic and pointer comparisons. For example, consider:

```
int *p = malloc(sizeof(int));
free(p);
int *q = malloc(sizeof(int));
printf("%d\n", p == q);
```

If both calls to `malloc` succeed¹ the program still has undefined behavior because `p` has become indeterminate.

When a pointer has been deallocated, not just the argument of `free` becomes indeterminate, but also all other copies of that pointer. In our memory model we represent pointer values symbolically, and keep track of the memory areas that have been deallocated. The behavior of operations like `==` depends on the memory state, which allows us to accurately capture undefined behaviors.

¹ The size of the CH₂O memory is unbounded, hence `malloc` cannot fail. Addressing finiteness of the memory is left for future work.

2.4 Effective types and aliasing restrictions

Aliasing means that multiple pointers refer to the same object in memory. Consider:

```
int f(int *p, int *q) {
    int x = *q; *p = 10; return x;
}
```

When `f` is called with aliased pointers `p` and `q`, the assignment to `*p` also affects `*q`. As a result, a compiler cannot transform the function body of `f` into the shorter `*p = 10; return (*q);`.

Unlike this example, there are many situations in which pointers can be assumed *not* to alias. It is essential for an optimizing compiler to determine when aliasing cannot occur, and use this information to generate faster code. The technique of determining whether pointers can alias or not is called *alias analysis*.

In *type-based alias analysis*, type information is used to determine if pointers are aliased. Consider the following example:

```
short g(int *p, short *q) {
    short x = *q; *p = 10; return x;
}
```

Here, a compiler should be able to assume that `p` and `q` are not aliased because they point to objects with different types. However, the static type system of C is too weak to ensure that property because a union type can be used to call `g` with aliased pointers.

```
union INT_SHRT { int x; short y; };
union INT_SHRT u = { .y = 3 };
g(&u.x, &u.y);
```

A union is C's version of a *sum* type. Contrary to sum types in many other programming languages, unions are *untagged* instead of *tagged*, which means that the current variant of a union is not stored. As a result, unions destroy the property that each memory area has a unique type that is known at compile time. The *effective type* [11, 6.5p6-7] of a memory area thus depends on the *run time behavior* of the program.

The *strict-aliasing restrictions* [11, 6.5p6-7] imply that a pointer to a variant of a union type (not to the whole union itself) can only be used for an access (read or write) if the union is in that particular variant. Calling `g` with aliased pointers (as in the example where `u` is in the `y` variant, and is accessed through a pointer `p` to the `x` variant) results in undefined behavior.

Contrary to existing formalizations (Leroy *et al.* [21] and Ellison and Rosu [6]) where each object consists of an array of bytes, we use structured trees [13]. By giving the memory more structure, we capture the strict-aliasing restrictions: effective types are modeled by the state of the trees in the memory model.

2.5 Type-punning

Despite the aliasing restrictions of C, it is under certain conditions allowed to access a union through another variant than the current one. This is called *type-punning* [11, 6.5.2.3]. Since the C standard is ambiguous about these conditions², we follow the interpretation by the GCC documentation [8]. It states: "type-punning is allowed, provided the memory is accessed through the union type".

For example, according to this interpretation the following program has implementation defined behavior (on architectures where `shorts` do not have trap values):

```
union INT_SHRT { int x; short y; };
```

²The term *type-punning* is merely used in a footnote, but for the related *common initial segment* rule, it uses the notion of *visible*, which is not clearly defined either.

```
union INT_SHRT u; u.x = 3;
printf("%d", u.y);
```

Type-punning may only be performed directly via an l-value of union type. Indirect type-punning via a pointer to a particular variant of a union type yields undefined behavior. For example:

```
union INT_SHRT u; short *p = &u.y; u.x = 3;
printf("%d", *p);
```

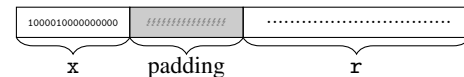
We formalize the informal semantics of GCC by decorating the formal references to subobjects with annotations [13]. Whenever a pointer to a variant of some union is stored in memory, or used as the argument of a function, the annotations are changed to ensure that type-punning is no longer possible via that pointer. In [13], we have shown that the GCC interpretation is correct by proving that a compiler can indeed perform type-based alias analysis [13].

2.6 Byte-level operations

Apart from *high-level* access to objects in memory by means of typed expressions, C also allows *low-level* access by means of byte-wise manipulation. Each object of type τ can be interpreted as an *unsigned char* array of length `sizeof(τ)`, which is called the *object representation* [11, 6.2.6.1p4]. An object can thus be copied by copying its object representation. For example:

```
struct { short x, *r; } s1 = { 10, &s1.x }, s2;
for (size_t i = 0; i < sizeof(s1); i++)
    ((unsigned char*)&s2)[i] =
        ((unsigned char*)&s1)[i];
```

The alignment requirements put restrictions on the addresses at which objects may be allocated [11, 6.2.8]. For each type τ , there is an implementation defined integer constant `_Alignof(τ)`. Objects of type τ should be allocated at addresses that are a multiple of `_Alignof(τ)`. On 32-bits architectures with `_Alignof(short*) = 4` (e.g. x86), the object representation of `s1` might be:



Due to the alignment requirements, the object representation of the above struct contains a hole. The bytes belonging to these holes are called *padding bytes*, and their values remain indeterminate even after initialization of the whole struct. To facilitate byte-wise copying of structs, a memory model should allow one to read and write all bytes, even those that are indeterminate, using an expression of type `unsigned char` [11, 6.2.6.1p5]. Note that in general, reading indeterminate memory has undefined behavior (for example, reading from an uninitialized `int` variable).

Our memory model uses a symbolic representation of bits in order to distinguish determinate and indeterminate memory. This way, we can precisely keep track of the situations in which reading and writing indeterminate memory is permitted.

The following excerpt from the C11 standard points out another challenge with respect to padding bytes [11, 6.2.6.1p6]:

```
When a value is stored in an object of structure or union type, including in a member object, the bytes of the object representation that correspond to any padding bytes take unspecified values
```

Let us illustrate this difficulty by an example:

```
struct { short x, *r; } s1 = { 10, NULL };
((unsigned char*)&s1)[3] = 1;
s1.x = 11;
printf("%d\n", ((unsigned char*)&s1)[3]);
```

We assign 1 to the third byte `((unsigned char*)&s1)[3]` of the struct `s1`, and then assign 11 to the field `x` of the struct `s1`. The last assignment makes all padding bytes of `s1` indeterminate, and consequently the argument of `printf` will be indeterminate. In our memory model based on trees, we enforce that padding bytes are always indeterminate, and therefore have this behavior implicitly.

2.7 Arithmetic conversions and overflow

In order to make C portable, the C standard gives compilers a lot of freedom to represent integers and to perform integer arithmetic. In particular, the sizes of integer types are implementation defined. For example, `int` does not necessarily have to be 32 bits and be able to exactly hold values between -2^{31} and $2^{31} - 1$. Only some minimum limits are described [11, 5.2.4.2.1]. In order to capture different architectures, our memory model is parametrized by an abstract interface of *integer implementations*.

Overflow of signed integers yields undefined behavior, whereas overflow has defined behavior (and wraps around modulo) for unsigned integers. For example, consider the following function:

```
int f(int x) { return x < x + 1; }
```

A compiler is thus allowed to optimize the function `f` to always return 1. Indeed, when compiled with `gcc -O2` (version 4.6.3), the following program prints `INT_MAX < INT_MAX+1 = 1`.

```
printf("INT_MAX < INT_MAX+1 = %d\n", f(INT_MAX));
```

The *integer promotions* [11, 6.3.1.1p2] and *usual arithmetic conversions* [11, 6.3.1.8p1] make the situations in which signed integer overflow occurs more subtle. For example, the following program has defined behavior (provided that `SHRT_MAX < INT_MAX`) because the arguments of `<` are promoted to `int` type.

```
printf("SHRT_MAX < SHRT_MAX+1 = %d\n",
      SHRT_MAX < SHRT_MAX + 1);
```

2.8 End-of-array pointers

The way the C standard deals with pointer equality is subtle. Consider the following excerpt [11, 6.5.9p6]:

Two pointers compare equal if and only if [...] or one is a pointer to one past the end of one array object and the other is a pointer to the start of a different array object that happens to immediately follow the first array object in the address space.

End-of-array pointers are peculiar because they cannot be dereferenced. Nonetheless, their use is common programming practice when looping through arrays.

```
void inc_array(int *p, int n) {
    int *end = p + n;
    while (p < end) (*p++)++;
}
```

End-of-array pointers can also be used in a way that is not stable under compilation. In the example below, the `printf` is executed only if `x` and `y` are allocated adjacently in memory (typically the stack).

```
int x, y;
if (&x + 1 == &y)
    printf("x and y are allocated adjacently\n");
```

Inspired by the CompCert memory model of Leroy *et al.* [21], we represent pointer values symbolically. The use of symbolic

pointer representations restricts the situations in which pointer arithmetic and pointer comparisons have defined behavior. This way, we assign undefined behavior to questionable uses of end-of-array pointers while assigning the correct defined behavior to pointer comparisons as in the first example above.

3. Operational semantics of CH₂O core C

In this section we describe the memory model and the operational semantics of our fragment of C11, and show how these address the challenges of Section 2. Since the memory model and operational semantics are an extension of previous work [13–15, 18], we do not describe these in detail. We refer to Figure 2 for the syntax and to the Coq formalization for other details.

3.1 Notations

We let *option B* denote the *option type over B*, whose elements are inductively defined as either \perp or x for some $x \in B$. Elements of the option type are denoted as $x^? \in \text{option } B$. A *partial function* $f : A \rightarrow \text{option } B$ is called *finite* if its domain is finite. The set of these functions is denoted as $A \rightarrow_{\text{fin}} B$.

We let *list B* denote the *list type over B*, whose elements are inductively defined as either ϵ or $x\vec{x}$ for some $x \in B$ and $\vec{x} \in \text{list } B$. We let $|\vec{x}| \in \mathbb{N}$ denote the length \vec{x} .

3.2 Memory model

Significant existing formal versions of a semantics of C, in particular those by Leroy *et al.* [21] and Ellison and Rosu [6], model the memory as a finite partial function from *object identifiers* to *objects*, where each object consists of a single array of bytes. Every block scope, global and static variable, and invocation of `malloc` is associated a unique object identifier with its own object in memory. This approach separates unrelated objects by construction, and is therefore more suitable for reasoning about memory transformations than the more intuitive approach of modeling the memory as an array of bytes (which is closer to an actual machine).

However, since no information about dynamic types is stored using this approach, not all undefined behaviors related to the aliasing restrictions (see Section 2.4) and interactions between high- and low-level memory accesses (see Section 2.6) can be described accurately. To remedy these shortcomings, we model the memory as a finite partial function from object identifiers to well-typed trees whose structure corresponds to the structure of data types in C. The leafs of these trees consist of arrays of (symbolic) bits that represent base values (integers and pointers).

The key concepts of our memory model are the following.

- *Pointers* contain paths following the structure of types. These paths, called *references*, have annotations \circ and \bullet to restrict the situations in which type-punning has defined behavior (see Section 2.5 for the description of type-punning).
 - The annotation \bullet means that type-punning (accessing another variant of the union than the current one) is allowed.
 - The annotation \circ means that type-punning is forbidden.

Whenever a pointer is stored in memory, its annotations are set to \circ . These annotations give a formal treatment of the strict-aliasing restrictions of C11. Pointers are annotated with the type to which they are cast.

- *Bits* are represented symbolically because each object (including those containing pointers and indeterminate memory) can be interpreted as an `unsigned char` array called the object representation (see Section 2.6). A bit is either a concrete bit 0 or 1, the i th bit $(\text{ptr } p)_i$ of a pointer p , or the indeterminate bit ℓ . Integers are represented using sequences of concrete bits, and pointers as sequences of abstract pointer bits.

Operators

$$\begin{aligned} \odot_c \in \text{compop} &::= == \mid <= \mid < \\ \odot_a \in \text{arithop} &::= + \mid - \mid * \mid / \mid \% \\ \odot_b \in \text{bitop} &::= \& \mid | \mid \wedge \\ \odot_s \in \text{shifto} &::= \ll \mid \gg \\ \odot \in \text{binop} &::= \odot_c \mid \odot_b \mid \odot_a \mid \odot_s \\ \odot_u \in \text{unop} &::= - \mid \sim \mid ! \\ \alpha \in \text{assign} &::= := \mid \odot := \mid := \odot \end{aligned}$$

Types

$$\begin{aligned} k \in K &::= \text{Set of integer ranks} \\ t \in \text{tag} &::= \text{Set of struct/union names} \\ si \in \text{signedness} &::= \text{signed} \mid \text{unsigned} \\ \tau_i \in \text{inttype} &::= si \ k \\ \tau_b \in \text{basetype} &::= \tau_i \mid \tau * \mid \text{void} \\ \tau, \sigma \in \text{type} &::= \tau_b \mid \tau[n] \mid \text{struct } t \mid \text{union } t \\ \Gamma \in \text{env} &::= \text{tag} \rightarrow_{\text{fin}} \text{list type} \end{aligned}$$

Separation algebras and permissions

$$\begin{aligned} \mathcal{L}(A) &::= \{\circ, \bullet\} \times A && \text{(Lockable SA)} \\ \mathcal{C}(A) &::= \mathbb{Q} \times A && \text{(Counting SA)} \\ \mathcal{T}_T^x(A) &::= A \times T \text{ with } x \in T && \text{(Tagged SA)} \\ p \in \text{perm} &::= \mathcal{L}(\mathcal{C}(\mathbb{Q})) \end{aligned}$$

Memory model

$$\begin{aligned} o \in \text{index} &::= \text{Set of memory indexes} \\ r \in \text{refseg} &::= \xrightarrow{\tau[n]} i \mid \xrightarrow{\text{struct } t} i \mid \xrightarrow{\text{union } t} q \ i \text{ with } q \in \{\circ, \bullet\} \\ \vec{r} \in \text{ref} &::= \text{list refseg} \\ a \in \text{addr} &::= (o : \tau, \vec{r}, i)_{\sigma > \sigma'} \\ p \in \text{ptr} &::= \text{NULL } \tau \mid a \\ b \in \text{bit} &::= 0 \mid 1 \mid (\text{ptr } p)_i \mid \ell \\ \mathbf{b} \in \text{pbit} &::= \mathcal{T}_{\text{bit}}^{\ell}(\text{perm}) \\ v_b \in \text{baseval} &::= \text{indet } \tau_b \mid \text{int}_{\tau_i} z \mid \text{ptr } p \mid \text{byte } \vec{b} \text{ with } z \in \mathbb{Z} \\ v \in \text{val} &::= v_b \mid \text{array}_{\tau} \vec{v} \mid \text{struct}_t \vec{v} \\ &\mid \text{union}_t(i, v) \mid \overline{\text{union}_t} \vec{v} \end{aligned}$$

$$\begin{aligned} w \in \text{mtree} &::= \text{base}_{\tau_b} \vec{\mathbf{b}} \mid \text{array}_{\tau} \vec{w} \mid \text{struct}_t \vec{w} \vec{\mathbf{b}} \\ &\mid \text{union}_t(i, w, \vec{\mathbf{b}}) \mid \overline{\text{union}_t} \vec{\mathbf{b}} \\ m \in \text{mem} &::= \text{index} \rightarrow_{\text{fin}} (\text{mtree} \times \text{bool} + \text{type}) \\ \Omega \in \text{lockset} &::= \text{index} \rightarrow_{\text{fin}} \text{list bool} \end{aligned}$$

CH₂O core C

$$\begin{aligned} f \in \text{funname} &::= \text{Set of function names} \\ l \in \text{labelname} &::= \text{Set of label names} \\ e \in \text{expr} &::= x_i^{\tau} \mid [v]_{\Omega} \mid [a]_{\Omega} \mid *e \mid \&e \mid e ._1 r \mid e ._r r \\ &\mid [r := e_1] e_2 \mid e_1 \alpha e_2 \mid f(\vec{e}) \mid \text{abort } \tau \\ &\mid \text{load } e \mid \text{alloc}_{\tau} e \mid \text{free } e \mid \odot_u e \mid e_1 \odot e_2 \\ &\mid e_1 ? e_2 : e_3 \mid (e_1, e_2) \mid (\tau) e \\ s \in \text{stmt} &::= e \mid \text{skip} \mid \text{goto } l \mid \text{throw } n \mid \text{return } e \mid l : \\ &\mid \text{local}_{\tau} s \mid \text{catch } s \mid s_1 ; s_2 \\ &\mid \text{loop } s \mid \text{if } (e) s_1 \text{ else } s_2 \\ \delta \in \text{funenv} &::= \text{funname} \rightarrow_{\text{fin}} \text{stmt} \end{aligned}$$

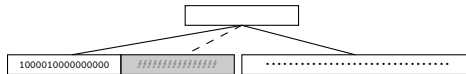
CH₂O core C states

$$\begin{aligned} \mathcal{E}_s \in \text{ectx}_s &::= * \square \mid \& \square \mid \square ._1 r \mid \square ._r r \mid [r := \square] e_2 \\ &\mid [r := e_1] \square \mid \square \alpha e_2 \mid e_1 \alpha \square \\ &\mid f(\vec{e}_1, \square, \vec{e}_2) \mid \text{load } \square \mid \text{alloc}_{\tau} \square \mid \text{free } \square \\ &\mid \odot_u \square \mid \square \odot e_2 \mid e_1 \odot \square \\ &\mid \square ? e_2 : e_3 \mid (\square, e_2) \mid (\tau) \square \\ \mathcal{E} \in \text{ectx} &::= \text{list ectx}_s \\ \mathcal{S}_s \in \text{sctx}_s &::= \text{catch } \square \mid \square ; s_2 \mid s_1 ; \square \mid \text{loop } \square \\ &\mid \text{if } (e) \square \text{ else } s_2 \mid \text{if } (e) s_1 \text{ else } \square \\ \mathcal{S}_e \in \text{sctx}_e &::= \square \mid \text{return } \square \mid \text{if } (\square) s_1 \text{ else } s_2 \\ \mathcal{P}_s \in \text{ctx}_s &::= \mathcal{S}_s \mid \text{local}_{o:\tau} \square \mid (\mathcal{S}_e, e) \\ &\mid \text{resume } \mathcal{E} \mid \text{params } f \vec{o}^{\tau} \\ \mathcal{P} \in \text{ctx} &::= \text{list ctx}_s \\ d \in \text{direction} &::= \searrow \mid \nearrow \mid \curvearrowright \mid \uparrow n \mid \uparrow v \\ \phi_U \in \text{undef} &::= \ell \mathcal{E}(e) \mid \ell \mathcal{S}_e([v]_{\Omega}) \\ \phi \in \text{focus} &::= (d, s) \mid e \mid \text{call } f \vec{v} \mid \overline{\text{return}} f v \mid \overline{\text{undef}} \phi_U \\ S \in \text{state} &::= \mathbf{S}(\mathcal{P}, \phi, m) \end{aligned}$$

Figure 2. Syntax of CH₂O core C. In these definitions we let $i, n \in \mathbb{N}$.

By representing bits symbolically, we ensure that additional information that is stored to describe certain undefined behaviors is preserved while performing byte-wise operations on object representations.

- *Memory trees* are abstract trees whose structure corresponds to the shape of data types. These are used to represent each object in memory. For example, the memory tree corresponding to `struct { short x, *r; } s = { 33; &s.x }` might be:

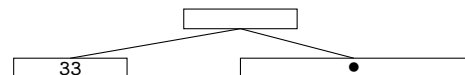


The leaves of memory trees contain symbolic bits: the integer 33 is represented by its binary representation 1000010000000000,

the padding is represented by symbolic indeterminate bits ℓ , and the pointer p by a sequence $(\text{ptr } p)_0 \dots (\text{ptr } p)_{31}$. These bits are annotated with permissions. Permissions are built modularly using a telescope of separation algebras [15].

The *memory* is a forest of memory trees. Compared to modeling each object in memory as an array of bits, memory trees store the variant of each union, and are explicit about padding bits.

- *Abstract values* are similar to memory trees, but have base values (integers and pointers) on their leaves. The abstract value of `struct { short x, *r; } s = { 33; &s.x }` might be:



Abstract values are used in the external interface of the memory model (which consists of functions to perform stores and loads) and hide internal details of the memory (permissions, padding, and object representations).

The external interface of the CH₂O memory model consists of operations with the following types:

- (-)(-)Γ : mem → addr → option val (1)
- forceΓ : addr → mem → mem (2)
- (-)[- := -]Γ : mem → addr → val → mem (3)
- writableΓ : addr → mem → Prop (4)
- lockΓ : addr → mem → mem (5)
- unlock : lockset → mem → mem (6)
- allocΓ : index → type → bool → mem → mem (7)
- freeable : addr → mem → Prop (8)
- free : index → mem → mem (9)

The operation $m\langle a \rangle_\Gamma$ yields the value stored at address a in the memory m . This operation yields \perp in case the permissions are not sufficient for a read access, effective types are violated, or in case a is an end-of-array address. Reading from (the abstract) memory is not a pure operation, it may affect the effective types [11, 6.5p6-7]. That means, whenever a union is accessed using an address to a different variant (provided the annotations \circ and \bullet permit it), the variant of the union should be changed in memory. This impurity is factored out by the operation $\text{force}_\Gamma a m$.

The operation $m[a := v]_\Gamma$ stores the value v at address a in m . A memory store is only permitted in case permissions are sufficient, effective types are not violated, and the address is not end-of-array. These side-conditions are described by $\text{writable}_\Gamma a m$.

After a successful store, the operation $\text{lock}_\Gamma a m$ is used to lock the object at address a in m . The lock operation temporarily reduces the permissions of a to prohibit all accesses to a . Locking yields a formal treatment of the sequence point restriction (which states that modifying an object more than once between two sequence points is undefined behavior, see Section 2.1). At a subsequent sequence point (for example, at the end ; of a full expression, or after the first operand of the conditional ? : operator has been evaluated), the operation $\text{unlock } \Omega m$ is used to unlock previously locked objects Ω in memory m . Unlocking has the intended effect of making subsequent accesses possible again.

The operation $\text{alloc}_\Gamma o \tau \beta m$ allocates a new object of type τ in the memory m . The new object has object identifier $o \notin \text{dom } m$ which is non-deterministically chosen by the operation semantics (see Section 3.3). The boolean β expresses whether the new object is allocated by `malloc` (the expression $\text{alloc}_\tau e$ in our language).

Conversely, the operation $\text{free } o m$ deallocates an object with object identifier o in the memory m . For the case of `malloced` memory, the side-condition $\text{freeable } a m$ ensures that the permissions are sufficient, and that a indeed points to a whole `malloced` object with object identifier $\text{index } a$.

As part of the Coq development, we have proven many properties about the interaction of memory operations. For example that stores to non-overlapping addresses commute.

3.3 Expression semantics

We define the semantics of expressions and statements by a small step operational semantics. That means, computation is defined as the reflexive transitive closure of a reduction relation. For expressions, we first define head reduction \rightarrow_h , and use *evaluation contexts* (as introduced by Felleisen *et al.* [7]) to select a redex in a whole expression. Given a stack $\rho \in \text{list index}$, the rules for head reduction are as follows:

1. $(x_i^\tau, m) \rightarrow_h ([\text{top}_\tau o]_\emptyset, m)$ if $\rho(i) = o$
2. $([*\text{ptr } a]_\Omega, m) \rightarrow_h ([a]_\Omega, m)$ if a is not a dangling address
3. $(\&[a]_\Omega, m) \rightarrow_h ([\text{ptr } a]_\Omega, m)$
4. $([a]_\Omega .i r, m) \rightarrow_h ([a']_\Omega, m)$
where a' is obtained by appending r to a
5. $([v]_\Omega .r r, m) \rightarrow_h ([v']_\Omega, m)$
where v' is the sub-value at location r in v
6. $([r := [v_1]_{\Omega_1}][v_2]_{\Omega_2}, m) \rightarrow_h ([v']_{\Omega_1 \cup \Omega_2}, m)$
where v' is obtained by inserting v_1 at location r in v_2
7. $([a]_{\Omega_1} := [v]_{\Omega_2}, m) \rightarrow_h ([v]_{\{a\} \cup \Omega_1 \cup \Omega_2}, \text{lock}_\Gamma a (m[a := v]_\Gamma))$
provided that $\text{writable}_\Gamma a m$
8. $(\text{load } [a]_\Omega, m) \rightarrow_h ([v]_\Omega, \text{force}_\Gamma a m)$ if $m\langle a \rangle_\Gamma = v$
9. $(\text{alloc}_\tau [\text{int}_{\tau_i} n]_\Omega, m) \rightarrow_h ([\text{ptr } a]_\Omega, \text{alloc}_\Gamma o (\tau[n]) \text{ true } m)$
for any o satisfying $o \notin \text{dom } m$, where a is the 0-th element of $\text{top}_{\tau[n]} o$, and provided that $n \neq 0$, and $\text{sizeof}_\Gamma \tau \cdot n$ is representable by an integer of type signed `ptr_rank`
10. $(\text{free } [\text{ptr } a]_\Omega, m) \rightarrow_h ([\text{void}]_\Omega, \text{free } (\text{index } a) m)$
provided that $\text{freeable } a m$
11. $(\odot_u [v]_\Omega, m) \rightarrow_h ([\odot_u v]_\Omega, m)$
provided that $\Gamma; m \vdash (\odot_u v)$ defined
12. $([v_1]_{\Omega_1} \odot [v_2]_{\Omega_2}, m) \rightarrow_h ([v_1 \odot v_2]_{\Omega_1 \cup \Omega_2}, m)$
provided that $\Gamma; m \vdash (v_1 \odot v_2)$ defined
13. $(([v]_\Omega, e), m) \rightarrow_h (e, \text{unlock } \Omega m)$
14. $([v]_\Omega ? e_2 : e_3, m) \rightarrow_h (e_2, \text{unlock } \Omega m)$ if $\text{istrue}_m v$
15. $([v]_\Omega ? e_2 : e_3, m) \rightarrow_h (e_3, \text{unlock } \Omega m)$ if $\text{isfalse } v$
16. $(\tau [v]_\Omega, m) \rightarrow_h ((\tau)v_\Omega, m)$
provided that $\Gamma; m \vdash (\tau)v$ defined

(Rule 7 just includes the case of a normal assignment $:=$, the rules for assignment operators $\odot :=$ and $:= \odot$ are similar).

The leafs of expressions, values $[v]_\Omega$ and addresses $[a]_\Omega$, are annotated with a finite set $\Omega \in \text{lockset}$ of *locked* addresses. This set is initially empty, but whenever a write is performed, the written object is locked in memory and its address is added to Ω (see rule 7 for assignment above). Locking a enforces the sequence point restriction because the permission model [14, 15] enforces that consecutive reads and writes to a will fail. Whenever a sequence point occurs, for example in rule 14 for the conditional expression, the locked objects in Ω are unlocked. This makes subsequent accesses to these objects possible again.

The C standard distinguishes *l-values* and *r-values* [11, 6.3.2.1]. This distinction originates from the assignment $e_1 = e_2$ where the left-hand side does not denote an actual value, but an address. L-values (that do not have array type) are implicitly converted to r-values by *l-value conversion* [11, 6.3.2.1p2].

In CH₂O core C we make this distinction, as well as the conversion between the two, explicit. L-values reduce to addresses $[a]_\Omega$ while r-values reduce to abstract values $[v]_\Omega$. The expressions $*e$, $\&e$ and $\text{load } e$ are used to map between the two. CH₂O abstract C has implicit l-value conversion which is being disambiguated during the translation to core C (see Section 6).

CH₂O core C uses De Bruijn indices for local variables, which means that a local variable x_i^τ refers to the i -th item on the stack $\rho \in \text{list index}$. Accessing the value of a local variable has an extra level of indirection: the stack ρ contains a reference to the value in memory instead of the value of the variable itself. This way, pointers to both local and allocated storage are treated uniformly. Evaluation of a variable x_i^τ consists of looking up its object identifier o in the stack ρ , and returning the address $\text{top}_\rho o := (o : \tau, \epsilon, 0)_{\tau > \tau}$

the top of the memory tree at position o in memory (rule 1). Using a load, the actual value can be obtained (rule 8).

Due to explicit l-value conversion, CH₂O core C has two operators for field indexing of structs and unions: $e \cdot_1 r$ for l-values, and $e \cdot_r r$ for r-values. The semantics of these operators are different: the first refines an address to refer to the corresponding field (rule 4), whereas the latter refines the actual value (rule 5).

In the reduction of whole programs (see Section 3.4), we allow $\mathcal{E}[e_1]$ to reduce to $\mathcal{E}[e_2]$ provided that $(e_1, m_1) \rightarrow_h (e_2, m_2)$. To make expression evaluation non-deterministic (see Section 2.1), we include both the contexts $\square + e_2$ and $e_1 + \square$. However, to enforce that the first operand of the conditional $e_1 ? e_2 : e_3$ is executed entirely before the others, it is essential that we omit the contexts $e_1 ? \square : e_3$ and $e_1 ? e_2 : \square$. This approach is also used by Norrish [26], Leroy [20] and Ellison and Rosu [6]. Also, like Ellison and Rosu [6], we implicitly use non-determinism to capture undefined behavior due to sequence point violations. For example, in $x + (x = 10)$ only one execution order (performing the read after the assignment) leads to a sequence point violation.

The side-condition $\Gamma; m \vdash (v_1 \odot v_2)$ defined for binary operators depends on the memory m (rule 12). This is needed because operations on pointers (for example comparing a pointer, see Section 2.3) only have defined behavior in case the pointers in question have not been deallocated.

The abort expression `abort τ` does not have a corresponding rule for head reduction and therefore has undefined behavior as explained in Section 3.4. The translation from CH₂O abstract C inserts abort τ expressions in branches of the program that should not be reached to describe all undefined behaviors.

Contrary to our previous work [14], we take the *integer promotions* [11, 6.3.1.1p2] and the *usual arithmetic conversions* [11, 6.3.1.8p1] for operators into account (see Section 2.7). Naively, it seems best to make these promotions and conversions explicit in core C and let the translation of CH₂O abstract C into core C insert appropriate casts. However, for assignment operators it is unclear how this could be handled by translation. For example, consider:

```
int x = INT_MAX;
x += INT_MAX + 2U;
```

This program has defined behavior, and x has the value 0 after its execution. Since `INT_MAX + 2U` has type `unsigned int`, the addition $x + (\text{INT_MAX} + 2U)$ performed by the assignment operator `+=` promotes both arguments to `unsigned int`, and the result will thus wrap around modulo instead of exhibiting undefined behavior. However, if we would have homogeneous assignment operators, then we would have to insert the following `int` cast:

```
int x = INT_MAX;
x += (int)(INT_MAX + 2U);
```

In this case, `(int)(INT_MAX + 2U)` results in signed integer overflow which causes undefined behavior.

3.4 Whole programs semantics

The small step reduction \rightarrow of *program states* $\mathbf{S}(\mathcal{P}, \phi, m)$ uses a zipper-like data structure \mathcal{P} , called a *program context* [18], which describes the location of the part of the program ϕ , called the *focus*, that is being executed. Program contexts extend the traditional zipper [9] by annotating each block scope variable with its associated object identifier, and furthermore contain the full call stack of the program. We have five kinds of focuses ϕ :

- The focus (d, s) describes execution of a *statement* s in *direction* d . The directions correspond to the different forms of (non-local) control: *down* \searrow , *up* \nearrow , *goto* \curvearrowright , *throw* \uparrow , or *top* \dagger .

When a `goto` l statement is executed, the direction is changed, and a small step traversal through the zipper is performed until the label l has been reached where normal control flow is resumed. This is described by the following rules:

$$\begin{aligned} \mathbf{S}(\mathcal{P}, (\searrow, \text{goto } l), m) &\rightarrow \mathbf{S}(\mathcal{P}, (\curvearrowright l, \text{goto } l), m) \\ \mathbf{S}(\mathcal{P}, (\curvearrowright l, \mathcal{S}_s[s]), m) &\rightarrow \mathbf{S}(\mathcal{S}_s \mathcal{P}, (\curvearrowright l, s), m) \\ &\quad \text{if } l \in \text{labels } s \\ \mathbf{S}(\mathcal{S}_s \mathcal{P}, (\curvearrowright l, s), m) &\rightarrow \mathbf{S}(\mathcal{P}, (\curvearrowright l, \mathcal{S}_s[s]), m) \\ &\quad \text{if } l \notin \text{labels } s \\ \mathbf{S}(\mathcal{P}, (\curvearrowright l, l :), m) &\rightarrow \mathbf{S}(\mathcal{P}, (\nearrow, l :), m) \end{aligned}$$

The point of this traversal is not so much to *search* for the label, but much more to incrementally *calculate* the required allocations and deallocations. Instead, the point of this traversal is to accurately describe the interaction between the `goto/break/continue/return` statement and block scope variables (see Section 2.2 for its subtleties).

When entering a scope with a local variable `local τ s` (this construct is nameless because we use De Bruijn indices for variables), the singular context `local o, τ \square` is appended to the head of the program context. The rule for `goto` is:

$$\begin{aligned} \mathbf{S}(\mathcal{P}, (\curvearrowright l, \text{local}_{\tau} s), m) \\ \rightarrow \mathbf{S}((\text{local}_{\tau, o} \square) \mathcal{P}, (\curvearrowright l, s), \text{alloc}_{\Gamma} o \tau \text{ false } m) \end{aligned}$$

Here, we should have $o \notin \text{dom } m$ and $l \in \text{labels } s$. After execution of this rule, the program context associates the block scope variable with its object identifier o . The *corresponding stack* `getstack \mathcal{P}` of a program context \mathcal{P} (which is used by the expression semantics) is obtained by an iteration over \mathcal{P} .

- The focus $\overline{\text{undef}} \phi_U$ describes undefined behavior. The undefined state ϕ_U indicates how the undefined behavior is caused.
- The focus e describes execution of an expression e . For example, provided $(e_1, m_1) \rightarrow_h (e_2, m_2)$, we have:

$$\mathbf{S}(\mathcal{P}, \mathcal{E}[e_1], m_1) \rightarrow \mathbf{S}(\mathcal{P}, \mathcal{E}[e_2], m_2)$$

The non-deterministic decomposition in an evaluation context \mathcal{E} and subexpression e_1 models unspecified execution order of expressions.

In case the expression e_1 is unsafe, *i.e.* it is a redex that cannot be \rightarrow_h -reduced (due to undefined behavior), we have:

$$\mathbf{S}(\mathcal{P}, \mathcal{E}[e_1], m) \rightarrow \mathbf{S}(\mathcal{P}, \overline{\text{undef}} (\mathcal{E} \langle e_1 \rangle), m)$$

For function calls we have the following rule:

$$\begin{aligned} \mathbf{S}(\mathcal{P}, \mathcal{E}[f([v_0]_{\Omega_0}, \dots, [v_n]_{\Omega_n})], m) \\ \rightarrow \mathbf{S}((\text{resume } \mathcal{E}) \mathcal{P}, \overline{\text{call}} f \vec{v}, \text{unlock} (\bigcup \vec{\Omega}) m) \end{aligned}$$

This rule is not part of the expression head reduction \rightarrow_h because it changes the focus to `call $f \vec{v}$` .

- The focus `call $f \vec{v}$` describes calling a function f with arguments v . The rule for calling a function is:

$$\mathbf{S}(\mathcal{P}, \overline{\text{call}} f \vec{v}, m_1) \rightarrow \mathbf{S}((\text{params } f \vec{\sigma}) \mathcal{P}, (\searrow, s), m_2)$$

Here, m_2 is obtained by allocating the function arguments v at object identifiers $\vec{\sigma}$ with types $\vec{\tau}$. The reduction relation \rightarrow is parametrized by an environment $\delta \in \text{funname} \rightarrow_{\text{fin}} \text{stmt}$ that maps function names to function bodies. For the above rule, we thus should have $\delta f = s$. The program context `params $f \vec{\sigma}$` delimits the stack of the callee, and behaves like a sequence `local $o_0: \tau_0$ $\square, \dots, \text{local}_{o_1: \tau_1} \square$` .

- The focus $\overline{\text{return}} f v$ describes returning from a function f with return value v . Some of the corresponding rules are:

$$\begin{aligned} \mathbf{S}((\text{params } f \overline{\sigma\vec{\tau}}) \mathcal{P}, (\uparrow v, s), m) &\rightarrow \mathbf{S}(\mathcal{P}, \overline{\text{return}} v, \text{free } \vec{\sigma} m) \\ \mathbf{S}((\text{resume } \mathcal{E}) \mathcal{P}, \overline{\text{return}} v, m) &\rightarrow \mathbf{S}(\mathcal{P}, \mathcal{E}[[v]_{\emptyset}], m) \end{aligned}$$

Distinguishing different kinds of states for different forms of program execution is inspired by Leroy’s CompCert C [19]. However, our treatment of statement execution is entirely different.

The C standard describes while, for, and do-while looping statements [11, 6.8.5]. Adding all of these looping statements as primitives to CH₂O core C would cause duplication (since the semantics of these loops are similar). To that end, we only have a statement loop s for an infinite loop, which combined with the throw n statement that jumps to the n th surrounding catch s block, can encode all C loops (see Section 6 for the translation).

4. Type system of CH₂O core C

Contrary to strongly typed programming languages like Java, ML or Safe Haskell, the C programming language does not enjoy type safety. That means, well-typedness according to the C type system does not ensure the absence of crashes. The reason for this is that even well-typed programs may exhibit executions with undefined behaviors such as dereferences of the NULL pointer, or aliasing or sequence point violations.

Nonetheless, the C type system still ensures the absence of basic mistakes like using an integer where a pointer is expected (without explicit cast) or using an undeclared variable. Clearly, the following program is not valid and should be ruled out by the type system:

```
int f(int x) {
  if(0) { return &x; } else { return 10; }
}
```

We will present a type system for CH₂O core C and prove that it enjoys type preservation and a weak form of type safety. This weak form of type safety guarantees that a well-typed program can keep on reducing, while possibly eventually reaching a final state or an undef ϕ_U state. From a programmer’s point of view, weak type safety gives few guarantees because a program may still crash, but from a formalist point of view, it ensures that our semantics behaves well and that we have not forgotten any reductions.

For each syntax category we introduce a corresponding typing judgment. For example, for types we have the judgment $\Gamma \vdash \tau$ that ensures that τ does not contain any cyclic structs and unions. The typing judgment for expressions is $\Gamma, \Gamma_f, \Gamma_m, \vec{\tau} \vdash e : \tau_r$, where:

- The map $\Gamma \in \text{tag} \rightarrow_{\text{fin}} \text{list type}$ associates a list of field types to each struct/union name.
- The map $\Gamma_f \in \text{funname} \rightarrow_{\text{fin}} (\text{list type} \times \text{type})$ associates a list of argument types and a return type to each function name.
- The map $\Gamma_m \in \text{index} \rightarrow_{\text{fin}} (\text{type} \times \text{bool})$ associates a type and boolean to each object identifier. The boolean indicates whether the object has been deallocated or is still alive. For technical reasons, we keep track of the types of deallocated objects.
- The list $\vec{\tau} \in \text{list type}$ gives the types of the De Bruijn variables.

The type $\tau_l, \tau_r \in \text{type} + \text{type}$ indicates whether the expression yields an l-value (address) or r-value (abstract value). We list some typing rules (with typing environments omitted from now on):

$$\frac{\vec{\tau}(i) = \tau}{x_i^{\vec{\tau}} : \tau_l} \quad \frac{e : \tau_l}{\&e : (\tau^*)_r} \quad \frac{\Gamma_f(f) = (\vec{\tau}, \sigma) \quad \vec{e} : \vec{\tau}_r}{f(\vec{e}) : \tau_r}$$

The typing judgment for expression contexts $\mathcal{E} : \tau_r \mapsto \sigma_r$ states that given an expression e of type τ_r , the expression $\mathcal{E}[e]$ obtained by substituting e for the hole \square in \mathcal{E} has type σ_r .

The typing judgment for statements is of the shape $s : (\beta, \tau^?)$. Here, $\tau^? \in \text{option type}$ denotes the type of the return statements in s , and is \perp if s does not contain any returns, and $\beta \in \text{bool}$ denotes whether s is statically known to never reach the end because it always executes a return statement. We list some typing rules:

$$\frac{}{\text{skip} : (\text{false}, \perp)} \quad \frac{e : \tau_r}{\text{return } e : (\text{true}, \tau)} \quad \frac{}{\text{goto } l : (\text{true}, \perp)}$$

$$\frac{s_1 : (\beta_1, \tau_1^?) \quad s_2 : (\beta_2, \tau_2^?) \quad \tau_1^? \cup \tau_2^? = \sigma^?}{s_1 ; s_2 : (\beta_2, \sigma^?)}$$

Here, $\cup : \text{option type} \rightarrow \text{option type} \rightarrow \text{option (option type)}$ is defined as $\tau \cup \tau := \tau$, $\perp \cup \tau_r := \tau_r$, and $\tau_l \cup \perp := \tau_l$.

The judgment $\delta : \Gamma_f$ for function environments ensures that each function body in $\delta \in \text{funname} \rightarrow_{\text{fin}} \text{stmt}$ is well-typed with respect to its prototype in $\Gamma_f \in \text{funname} \rightarrow_{\text{fin}} (\text{list type} \times \text{type})$. Finally, the judgment $S : g$ ensures that all components of the state S are typed with respect to a main function $g \in \text{funname}$ (in CH₂O core C, `main` can be any function g of any type, so we need to ensure that the argument and return types match).

Lemma 4.1. *All typing judgments satisfy uniqueness of typing and have a corresponding type inference function.*

Proof. Since all type judgments are defined in a syntax directed fashion, this follows trivially. \square

Fact 4.2. *All typing judgments are closed under weakening.*

Lemma 4.3 (Type preservation). *If a state S_1 with $S_1 : g$ and $S_1 \rightarrow S_2$, then for resulting state S_2 we have $S_2 : g$.*

Proof. By case analysis on the derivation of $S_1 \rightarrow S_2$, using the fact that all memory operations preserve typing. \square

Lemma 4.4 (Weak progress). *If a state is typed $S_1 : g$, then either:*

1. *It can reduce further, that is $S_1 \rightarrow S_2$ for some S_2 .*
2. *It is an undefined state, that is $S_1 = \mathbf{S}(\mathcal{P}, \overline{\text{undef}} \phi_U, m)$ for some \mathcal{P}, ϕ_U and m .*
3. *It is a final state, that is $S_1 = \mathbf{S}(\epsilon, \overline{\text{return}} g v, m)$ for some v and m .*
4. *The labels are incorrect, that is $S_1 = \mathbf{S}(\mathcal{P}, (\wedge l, s), m)$ for some \mathcal{P}, l, s and m with $l \notin \text{labels } s \cup \text{labels } \mathcal{P}$.*
5. *The throws are incorrect, that is $S_1 = \mathbf{S}(\mathcal{P}, (\uparrow n, s), m)$ for some \mathcal{P}, n, s and m with $n \not\prec \text{num_catches } \mathcal{P}$.*

Proof. By case analysis on the structure of S_1 . \square

For technical reasons, we do not let the judgment $S_1 : g$ ensure correctness of gotos and throws. To that end, the progress theorem includes two cases to account for stuck non-local control. However, the judgment $\delta : \Gamma_f$ for function environments ensures that all gotos and throws are correct. To that end, whenever we start from an initial state, these cases of stuck non-local control cannot happen.

Theorem 4.5 (Weak type safety). *Given an initial state S_1 for main g with arguments \vec{v} , then if $S_1 \rightarrow^* S_2$ we have either:*

1. *$S_2 \rightarrow S_3$ for some S_3 .*
2. *$S_2 = \mathbf{S}(\mathcal{P}, \overline{\text{undef}} \phi_U, m)$ for some \mathcal{P}, ϕ_U and m .*
3. *$S_2 = \mathbf{S}(\epsilon, \overline{\text{return}} g v, m)$ for some v and m .*

Proof. Using Lemmas 4.4 and 4.3, and the fact that \rightarrow preserves validity of gotos and throws. \square

5. Executable semantics of CH₂O core C

The semantics of CH₂O core C (Section 3) is defined as an inductively defined reduction relation $_ \rightarrow _ : \text{state} \rightarrow \text{state} \rightarrow \text{Prop}$. Since this relation is not executable, we also define a corresponding function $\text{exec} : \text{state} \rightarrow \mathcal{P}_{\text{fin}}(\text{state})$ that *computes* the finite set of consecutive states.

In previous work [13], we already defined basic operations, such as accessing the memory or performing a cast, as Coq functions that are effectively executable. Hence, creating an executable version of our semantics seems straightforward. However, non-determinism makes the situation more complicated.

- Execution of expressions is non-deterministic. For example, in $(\ast p = 1) + (\ast q = 2)$, the assignments may be executed in any order. This occurs in the following rule:

$$\mathbf{S}(\mathcal{P}, \mathcal{E}[e_1], m_1) \rightarrow \mathbf{S}(\mathcal{P}, \mathcal{E}[e_2], m_2)$$

Here, an expression is decomposed non-deterministically in an evaluation context \mathcal{E} and a subexpression e_1 .

- The choice of picking an object identifier for newly allocated memory is non-deterministic. For example, this occurs in the rule for entering a block scope with a local variable:

$$\begin{aligned} & \mathbf{S}(\mathcal{P}, (\backslash, \text{local } \tau s), m) \\ \rightarrow & \mathbf{S}((\text{local}_{o:\tau} \square) \mathcal{P}, (\backslash, s), \text{alloc}_{\Gamma} o \tau \text{ false } m) \end{aligned}$$

Here, any unused object identifier $o \notin \text{dom } m$ may be chosen.

The first source of non-determinism is finitary because expressions can only be decomposed in a finite number of ways. However, there is an infinite choice of picking fresh object identifiers.

In order to deal with the first source of non-determinism, we define a function $\text{redexes} : \text{expr} \rightarrow \mathcal{P}_{\text{fin}}(\text{ectx} \times \text{expr})$ that decomposes an expression into a finite set of all possible combinations of evaluation contexts and redexes.

Lemma 5.1 (Soundness and completeness of redex splitting). *We have $(\mathcal{E}, e) \in \text{redexes } e'$ iff $e' = \mathcal{E}[e]$ and e is a redex.*

The second source of non-determinism is more difficult to deal with. Since there is an infinite choice of fresh object identifiers, we cannot just try them all. Instead, we will just pick one:

$$\begin{aligned} & \text{exec } (\mathbf{S}(\mathcal{P}, (\backslash, \text{local } \tau s), m)) := \\ & \text{let } o := \text{fresh } m \text{ in} \\ & \{ \mathbf{S}((\text{local}_{\tau.o} \square) \mathcal{P}, (\backslash, s), \text{alloc}_{\Gamma} o \tau \text{ false } m) \} \end{aligned}$$

Using a canonical object identifier for newly allocated memory removes the second source of non-determinism entirely. For the soundness theorem, this choice does not matter.

Theorem 5.2 (Soundness). *If $S_2 \in \text{exec } S_1$, then $S_1 \rightarrow S_2$.*

Proof. By case analysis on S_1 using Lemma 5.1. \square

The converse of the above theorem is not true. Given a reduction $S_1 \rightarrow S_2$, we do not necessarily have $S_2 \in \text{exec } S_1$ because the operational semantics may have used a different object identifier. For completeness, we show that this choice does not matter.

Definition 5.3. *A state S_1 is an f -permutation of state S_2 , notation $S_1 \sim_f S_2$, if S_2 is obtained by renaming the object identifiers in S_1 with respect to $f : \text{index} \rightarrow \text{option index}$.*

Fact 5.4. *We have the following properties:*

1. Identity: $S \sim_f S$ for f with $f x = x$ for each $x \in \text{dom } S$.
2. Composition: If $S_1 \sim_f S_2$ and $S_2 \sim_{f'} S_3$, then $S_1 \sim_{f' \circ f} S_3$.
3. Symmetry: If $S_1 \sim_f S_2$, then $S_2 \sim_{f^{-1}} S_1$.

4. Weakening: If $S_1 \sim_f S_2$ and $f' \supseteq f$, then $S_1 \sim_{f'} S_2$.

Lemma 5.5. *If $S_1 \rightarrow S_2$ and $S'_1 \sim_f S_1$, then there exists an $f' \supseteq f$ and S'_2 such that:*

$$\begin{array}{ccc} S'_1 & \text{-----} \triangleright & S'_2 \\ \downarrow f & & \downarrow f' \\ S_1 & \text{-----} \triangleright & S_2 \end{array}$$

Proof. By case analysis on the derivation of $S_1 \rightarrow S_2$. \square

Lemma 5.6. *If $S_1 \rightarrow S_2$, then there exists an f and S'_2 such that:*

$$\begin{array}{ccc} & & S'_2 \\ & \text{exec} \text{-----} \triangleright & \downarrow f \\ S_1 & \text{-----} \triangleright & S_2 \end{array}$$

Here, $S_1 \xrightarrow{\text{exec}} S'_2$ denotes $S'_2 \in \text{exec } S_1$.

Proof. By case analysis on the derivation of $S_1 \rightarrow S_2$. Lemma 5.1 is used for expression steps, and properties of the memory model for steps that involve allocation of new objects. \square

Theorem 5.7 (Completeness). *If $S_1 \rightarrow^* S_2$, then there exists an f and S'_2 such that:*

$$\begin{array}{ccc} & & S'_2 \\ & \text{exec} \text{-----} \triangleright^* & \downarrow f \\ S_1 & \text{-----} \triangleright^* & S_2 \end{array}$$

Proof. By induction from the right on the derivation of $S_1 \rightarrow^* S_2$. In the inductive case, we have $S_1 \rightarrow^* S_2 \rightarrow S_3$. By the induction hypothesis we obtain an f and S'_2 with $S_1 \xrightarrow{\text{exec}}^* S'_2$ and $S'_2 \sim_f S_2$. Next, we obtain an f_1 and S'_3 with $S'_2 \rightarrow S'_3$ and $S'_3 \sim_{f_1} S_3$ by Lemma 5.5. Using Lemma 5.6, we obtain an f_2 and S''_3 with $S'_2 \xrightarrow{\text{exec}} S''_3$ and $S''_3 \sim_{f_2} S'_3$. Displayed as a diagram:

$$\begin{array}{ccccc} & & & & S''_3 \\ & & & & \downarrow f_2 \\ & & & \text{exec} \text{-----} \triangleright^* & S'_3 \\ & \text{exec} \text{-----} \triangleright^* & S'_2 & \text{-----} \triangleright^* & \downarrow f_1 \\ S_1 & \text{-----} \triangleright^* & S_2 & \text{-----} \triangleright^* & S_3 \end{array}$$

By composition of permutations, we have $S''_3 \sim_{f_1 \circ f_2} S_3$, which concludes the proof. \square

6. CH₂O abstract C

CH₂O abstract C bridges the gap between the abstract syntax tree obtained from the parser (Section 7) and CH₂O core C. Its syntax therefore closely resembles the structure of a `.c` file. In particular, it uses named variables instead of De Bruijn indices, makes type annotations (such as those on variables) implicit, and extends CH₂O core C with various features (such as enums and C-like loops). Its semantics is specified by translation into core C. In this section we will highlight some differences between both languages and prove type soundness of the translation.

Like a `.c` file, the syntax of CH₂O abstract C (see Figure 3) consists of a sequence of declarations:

- The items $(s, \text{struct } \vec{\tau} \vec{x})$ and $(u, \text{union } \vec{\tau} \vec{x})$ declare a struct s or union u with fields \vec{x} of type $\vec{\tau}$.
- The item $(t, \text{typedef } \tau)$ declares a *typedef* t abbreviating τ .
- The item $(u, \text{enum } x := e^? : \tau_i)$ declares an *enumeration type* u . This declares constants \vec{x} of integer type τ_i whose values are specified by the constant expressions $e^?$. The type τ_i is not fixed because it may differ for each enum [11, 6.7.2.2p4].

$x \in \text{string} ::= \text{Set of strings}$ $k \in \text{cintrank} ::= \text{char} \mid \text{short} \mid \text{int} \mid \text{long} \mid \text{long long} \mid \text{ptr}$ $si \in \text{signedness} ::= \text{signed} \mid \text{unsigned}$ $\tau_1 \in \text{cinttype} ::= si^? k$ $\tau \in \text{ctype} ::= \text{void} \mid \text{def } x \mid \tau_1 \mid \tau * \mid \tau[e]$ $\quad \mid \text{struct } x \mid \text{union } x \mid \text{enum } x \mid \text{typeof } e$ $e \in \text{cexpr} ::= x \mid \text{const}_{\tau_1} z \mid \text{sizeof } \tau \mid \tau_1 \text{ min}$ $\quad \mid \tau_1 \text{ max} \mid \tau_1 \text{ bits} \mid \&e \mid *e \mid e_1 \alpha e_2$ $\quad \mid x(\vec{e}) \mid \text{abort} \mid \text{alloc}_{\tau} e \mid \text{free } e$ $\quad \mid \odot_u e \mid e_1 \odot e_2 \mid e_1 \&\& e_2 \mid e_1 \parallel e_2$ $\quad \mid e_1 ? e_2 : e_3 \mid (e_1, e_2) \mid (\tau) I \mid e . x$ $r \in \text{crefseg} ::= [e] \mid .x$	$I \in \text{cinit} ::= e \mid \{\vec{r} := \vec{I}\}$ $sto \in \text{cstorage} ::= \text{static} \mid \text{extern} \mid \text{auto}$ $s \in \text{cstmt} ::= e \mid \text{skip} \mid \text{goto } x \mid \text{break} \mid \text{continue}$ $\quad \mid \text{return } e^? \mid \{s\} \mid \vec{sto} \tau x := I^? ; s$ $\quad \mid \text{typedef } x := \tau ; s \mid s_1 ; s_2 \mid x : s$ $\quad \mid \text{while}(e) s \mid \text{for}(e_1 ; e_2 ; e_3) s$ $\quad \mid \text{do } s \text{ while}(e) \mid \text{if}(e) s_1 \text{ else } s_2$ $d \in \text{decl} ::= \text{struct } \vec{\tau} \vec{x} \mid \text{union } \vec{\tau} \vec{x} \mid \text{typedef } \tau$ $\quad \mid \text{enum } x := e^? : \tau_1 \mid \text{global } I^? : \vec{sto} \tau$ $\quad \mid \text{fun } (\tau x^?) s^? : \vec{sto} \tau$ $\Theta \in \text{decls} ::= \text{list}(\text{string} \times \text{decl})$
--	---

Figure 3. Syntax of CH₂O abstract C. The operators are shared with core C, see Figure 2.

- The item $(x, \text{global } I^? : \vec{sto} \tau)$ declares a global variable x with initializer $I^?$.
- The item $(f, \text{fun } (\tau x^?) s^? : \vec{sto} \tau)$ declares a function f with return type τ and arguments \vec{x} of type $\vec{\tau}$. If the function body $s^?$ is omitted, the function can already be used in consecutive (mutually recursive) function declarations but has to be redeclared with a corresponding function body later.

The translator from CH₂O abstract C into core C turns a sequence of declarations Θ into:

- An environment $\Gamma \in \text{env}$ of struct/union declarations.
- An environment $\Gamma_f \in \text{funenv}$ of function types.
- An initial memory $m \in \text{mem}$ that contains storage for global and static variables declared in Θ .
- A function environment $\delta \in \text{funname} \rightarrow_{\text{fin}} \text{stmt}$ that contains statements for all function declarations in Θ .

An important part of the translation into CH₂O core is evaluation of constant expressions [11, 6.6]. For example, an array type $\tau[e]$ may use an arbitrary expression e to specify the length. The type $\text{int}[\text{sizeof}(x) + 10]$ is thus valid C.

To evaluate such expressions during translation, we have defined an evaluator $\llbracket _ \rrbracket_{\Gamma, \rho, m} : \text{expr} \rightarrow \text{option}(\text{addr} + \text{val})$ for constant expressions. It returns an address or value depending on whether the expression is an l- or r-value. An improvement on work by others is that we have proved soundness and completeness of this evaluator with respect to the operational semantics.

Lemma 6.1. *Constant expression evaluation is sound and complete. That is, given a constant expression e , then:*

$$\llbracket e \rrbracket_{\Gamma, \text{getstack } \mathcal{P}, m} = v \quad \text{iff} \quad \mathbf{S}(\mathcal{P}, e, m) \rightarrow^* \mathbf{S}(\mathcal{P}, v, m).$$

The translator transforms traditional C looping constructs into primitive constructs. For example, $\text{while}(e) s$ is translated into:

catch (loop (if (e') skip else throw 0 ; catch s'))

Here, e' and s' are the translations of e and s , respectively. The break and continue statements in s are translated into throw 1 and throw 0, that jump to the outer, respectively the inner catch.

As noted in Section 3, C makes a distinction between l-values (that denote addresses) and r-values (that denote abstract values). Since this distinction is implicit in CH₂O abstract C, but explicit in core C, we have to disambiguate it. In particular, we have to insert

casts to perform *l-value conversion* [11, 6.3.2.1p2]. For example, the expression $*p := y$ is translated into:

$$*(\text{load } x_i^{\tau^*}) := \text{load } x_j^{\tau}$$

Here, i and j are the De Bruijn indices corresponding to p and y . Also, we disambiguate field indexing $e . i$ into $e' . i$ and $e' . r$ depending on whether e is an l- or r-value.

In abstract C, the integer literal 0 is overloaded to denote both the NULL pointer and the integer constant 0 [11, 6.3.2.3p3]. We have to disambiguate this form of overloading in the translation to core C, because the typing rules of various expressions have special cases for NULL pointers [11, 6.5].

A subtle part is the translation of compound literals. Consider:

```
struct S { int x, y, z; } s = { .y=10, .x=11 };
```

Here, fields can be initialized in any order, and fields that are not initialized receive the value 0. We translate such initializers into sequence of $[_ := _]$ inserts into val_0 (struct S), the struct value initialized with zeros.

Finally, C provides many macros like INT_MIN to obtain information about the sizes of integer types [11, 5.2.4.2.1]. CH₂O abstract C has primitive constructs for these macros, such as $\tau_1 \text{ min}$. Since our translator is parametrized by an environment that describes the integer sizes of the architecture, these constructs will be translated into a value that depends on the architecture.

Theorem 6.2 (Type soundness). *If the translator succeeds with environments Γ , Γ_f , δ , and initial memory m , they are well-typed.*

7. Testing the semantics

In order to compute the behaviors of a C program, we have turned the single step function $\text{exec} : \text{state} \rightarrow \mathcal{P}_{\text{fin}}(\text{state})$ into two functions to compute a stream of reachable states.

- The function $\text{all} : \mathcal{P}_{\text{fin}}(\text{state}) \rightarrow (\mathcal{P}_{\text{fin}}(\text{state}) \times \mathcal{P}_{\text{fin}}(\text{state}))^\omega$ computes a stream of all reachable states. The n th element (\mathbf{S}, \mathbf{N}) of the stream $(\text{all } \mathbf{I})$ contains the intermediate states $\mathbf{S} \subseteq \text{state}$ and normal forms $\mathbf{N} \subseteq \text{state}$ after n steps starting in initial states $\mathbf{I} \subseteq \text{state}$.
- The function $\text{some} : (\mathbb{N} \rightarrow \mathbb{N}) \rightarrow \text{state} \rightarrow (\text{state} + \text{state})^\omega$ computes a specific execution stream. Given a selection function $f : \mathbb{N} \rightarrow \mathbb{N}$ that specifies which redexes should be chosen, the stream $(\text{some } f \mathbf{I})$ represents a trace starting in $\mathbf{I} \in \text{state}$.

9. Conclusion

We have defined an operational and executable semantics for a significant fragment of C11, that have been proved sound and complete with respect to each other in Coq. Contrary to previous work, we have used a typed core language, and specified the semantics of C programs by translation into this core language. This translation is implemented as a Coq program that has been proved sound.

Using the extraction mechanism of Coq this yields an interpreter that can explore the full state-space of a C program. We followed the C11 standard closely, and as a result our interpreter can detect subtle undefined behaviors not yet addressed by others.

Since C is a large and subtle language, there are many directions for future work. First of all, features could be added to the semantics. For example, floats [3], bit-fields, concurrency [28], const qualifiers, register storage, *etc.*

Despite the fact that the biggest part of our interpreter has been implemented in Coq, it still uses some glue that is written in OCaml. In particular, the translation of string literals and `printf` (that is currently just present for debugging purposes) is inaccurate. To implement these features properly, we have to support functions on a variable number of arguments, I/O, and the `const` qualifier in CH₂O core C. Instead of supporting just basic forms of I/O such as `printf`, it may better to support external function calls in a general way. The work of Beringer *et al.* [1] may be useful for that.

Instead of using the CIL parser, which is written in OCaml, it would be interesting to use the parser by Jourdan *et al.* [12] that is part of CompCert and has been implemented and verified in Coq. This way, the entire interpreter could be implemented in Coq.

We have tested our semantics on a small test suite. Considering the bugs that we have found this way were minor and easy to fix, it seems to indicate that using a proof assistant to prove metatheory about the language [13–15, 18] has already provided a good way of debugging the semantics. Nonetheless, it would be useful to test the semantics with a more extensive test suite, like the GCC torture tests [8], or using a tool like CSmith [30].

Lastly, our implementation of the executable semantics implements the operational semantics in a rather naive way, and is using inefficient data structures to represent the memory (for example, it uses lists to represent array objects). In order to make the interpreter more robust, it would be useful to optimize the executable semantics. This could be achieved by using more efficient data structures, or by contracting more redexes in one step.

Acknowledgments

We thank Herman Geuvers and the anonymous reviewers for their helpful comments. This work is financed by the Netherlands Organisation for Scientific Research (NWO).

References

- [1] L. Beringer, G. Stewart, R. Dockins, and A. W. Appel. Verified Compilation for Shared-Memory C. In *ESOP*, volume 8410 of *LNCS*, pages 107–127, 2014.
- [2] M. Bodin, A. Charguéraud, D. Filaretto, P. Gardner, S. Maffei, D. Naudziuniene, A. Schmitt, and G. Smith. A trusted mechanised JavaScript specification. In *POPL*, pages 87–100, 2014.
- [3] S. Boldo, J.-H. Jourdan, X. Leroy, and G. Melquiond. A Formally-Verified C Compiler Supporting Floating-Point Arithmetic. In *ARITH*, pages 107–115, 2013.
- [4] B. Campbell. An Executable Semantics for CompCert C. In *CPP*, volume 7679 of *LNCS*, pages 60–75, 2012.
- [5] W. Dietz, P. Li, J. Regehr, and V. S. Adve. Understanding integer overflow in C/C++. In *ICSE*, pages 760–770, 2012.
- [6] C. Ellison and G. Rosu. An executable formal semantics of C with applications. In *POPL*, pages 533–544, 2012.
- [7] M. Felleisen, D. P. Friedman, E. E. Kohlbecker, and B. F. Duba. A syntactic theory of sequential control. *Theoretical Computer Science*, 52:205–237, 1987.
- [8] GCC, the GNU Compiler Collection. <http://gcc.gnu.org/>.
- [9] G. P. Huet. The Zipper. *Journal of Functional Programming*, 7(5):549–554, 1997.
- [10] ISO. WG14 Defect Report Summary. <http://www.open-std.org/jtc1/sc22/wg14/www/docs/>.
- [11] ISO. *ISO/IEC 9899-2011: Programming languages – C*. ISO Working Group 14, 2012.
- [12] J. Jourdan, F. Pottier, and X. Leroy. Validating LR(1) Parsers. In *ESOP*, volume 7211 of *LNCS*, pages 397–416, 2012.
- [13] R. Krebbers. Aliasing restrictions of C11 formalized in Coq. In *CPP*, volume 8307 of *LNCS*, 2013.
- [14] R. Krebbers. An Operational and Axiomatic Semantics for Non-determinism and Sequence Points in C. In *POPL*, pages 101–112, 2014.
- [15] R. Krebbers. Separation algebras for C verification in Coq. In *VSTTE*, volume 8471 of *LNCS*, 2014.
- [16] R. Krebbers, X. Leroy, and F. Wiedijk. Formal C semantics: CompCert and the C standard. In *ITP*, volume 8558 of *LNAI*, pages 543–548, 2014.
- [17] R. Krebbers and F. Wiedijk. A Formalization of the C99 Standard in HOL, Isabelle and Coq. In *CICM*, volume 6824 of *LNAI*, pages 297–299, 2011.
- [18] R. Krebbers and F. Wiedijk. Separation Logic for Non-local Control Flow and Block Scope Variables. In *FoSSaCS*, volume 7794 of *LNCS*, pages 257–272, 2013.
- [19] X. Leroy. A formally verified compiler back-end. *Journal of Automated Reasoning*, 43(4):363–446, 2009.
- [20] X. Leroy. Formal verification of a realistic compiler. *CACM*, 52(7):107–115, 2009.
- [21] X. Leroy, A. W. Appel, S. Blazy, and G. Stewart. The CompCert Memory Model, Version 2. Research report RR-7987, INRIA, 2012.
- [22] P. Letouzey. A New Extraction for Coq. In *TYPES’02*, volume 2646 of *LNCS*, pages 200–219, 2003.
- [23] A. Lochbihler and L. Bulwahn. Animating the Formalised Semantics of a Java-Like Language. In *ITP*, volume 6898 of *LNCS*, pages 216–232, 2011.
- [24] N. Maclaren. What is an Object in C Terms?, 2001. Mailing list message, <http://www.open-std.org/jtc1/sc22/wg14/9350>.
- [25] G. C. Necula, S. McPeak, S. P. Rahul, and W. Weimer. CIL: Intermediate Language and Tools for Analysis and Transformation of C Programs. In *CC*, volume 2304 of *LNCS*, pages 213–228, 2002.
- [26] M. Norrish. *C formalised in HOL*. PhD thesis, University of Cambridge, 1998.
- [27] S. Owens, P. Böhm, F. Z. Nardelli, and P. Sewell. Lem: A Lightweight Tool for Heavyweight Semantics. In *ITP*, volume 6898 of *LNCS*, pages 363–369, 2011.
- [28] J. Sevcík, V. Vafeiadis, F. Z. Nardelli, S. Jagannathan, and P. Sewell. CompCertTSO: A Verified Compiler for Relaxed-Memory Concurrency. *Journal of the ACM*, 60(3):22, 2013.
- [29] P. Sewell, F. Z. Nardelli, S. Owens, G. Peskine, T. Ridge, S. Sarkar, and R. Strnisa. Ott: effective tool support for the working semanticist. In *ICFP*, pages 1–12, 2007.
- [30] X. Yang, Y. Chen, E. Eide, and J. Regehr. Finding and understanding bugs in C compilers. In *PLDI*, pages 283–294, 2011.
- [31] J. Zhao, S. Nagarakatte, M. M. K. Martin, and S. Zdancewic. Formalizing the LLVM intermediate representation for verified program transformations. In *POPL*, pages 427–440, 2012.