

A CASE STUDY IN SPECIFYING THE DENOTATIONAL SEMANTICS OF C

N. S. PAPASPYROU

National Technical University of Athens

Department of Electrical and Computer Engineering

Division of Computer Science, Software Engineering Laboratory

Polytechnioupoli, 15780 Zografou, Athens, Greece.

Tel. +30-1-7722486, Fax. +30-1-7722519.

1. Introduction

C is a well known and very popular general purpose programming language which represents, together with its descendants, a strong and indisputable status quo in the current software industry. Its semantics is informally defined in the ISO/IEC 9899:1990 standard [1] using natural language. This causes a number of ambiguities and problems of interpretation, clearly manifested in numerous discussions taking place in the newsgroup `comp.std.c`. It is worthwhile noticing that members of the standardization committee and other distinguished researchers participating in the discussions often give contradictory answers when asked about the intended semantics of surprisingly small programs, and that their answers are usually based on different possible interpretations of the standard. With all this in mind, the necessity for a formal description of the semantics of C becomes apparent. Such a description would serve as a precise standard for compiler implementation and would provide a basis for reasoning about properties of C programs.

The semantics of many popular programming languages have been formally specified in literature using various formalisms. However, in most cases these specifications are incomplete, inaccurate or both, in varying degrees. By *incomplete* we mean that they do not specify the semantics of the whole language but that of a subset, often leaving out the most complicated features. By *inaccurate* we mean that the formal descriptions are not entirely correct, either because of intended simplifications or by mistake.

Significant research has been conducted recently concerning semantic aspects of C. In what seems to be the earliest formal approach, Sethi addresses mainly the semantics of pre-ANSI C declarations, using the denotational approach and making several simplifications, e.g. requiring left-to-right evaluation [2]. In a different paper, Sethi addresses the semantics of C's control structures using again a denotational framework [3]. In the work of Gurevich and Higgins a formal semantics for C is given using the formalism of evolving algebras [4]. Again, a number of simplifications are made, e.g. no interleaving is possible in expression evaluation and side effects are assumed to take place at the same

time that they are generated. In the work of Cook and Subramanian an incomplete semantics for C is developed in the theorem prover Nqthm [5]. Cook et al. have also developed a denotational semantics for C based on temporal logic, which again makes a number of simplifying assumptions, mainly concerning evaluation order [6]. An operational semantics for C has been sketched, in terms of a random access machine, as a part of the MATHS project in California State University. Finally, in the work of Norrish a complete operational semantics for C is given using small-step reductions [7]. To the best of our knowledge, this is the only approach that formalizes correctly C's unspecified order of evaluation and sequence points. No similar denotational approach is known to us.

This chapter summarizes the results of our research, aiming at the development of an accurate and complete formal description for the semantics of the C programming language [8]. For this purpose, we have chosen the denotational approach.¹ As a remedy for the most important drawback of classic denotational semantics, its lack of modularity, we have used a number of monads which represent different aspects of computations.² In this way, the developed semantics was significantly improved in terms of modularity and elegance and its development was greatly facilitated.

2. Overview of the semantics

Our denotational specification for the semantics of C can be best understood as part of the *abstract interpreter* depicted in Fig. 1. The left part of the figure is a module diagram of the interpreter, showing the chain of actions performed and the processed data. Each action takes as input the result of previous actions. The initial piece of data is a *source program* and the final result is a representation of this program's *meaning*, i.e. a description of the program's behaviour when it is executed. The right part of the figure presents parts of a small example that will be discussed gradually until the end of this section. It should be noted that the example is simplified and does not correctly specify the semantics of C.

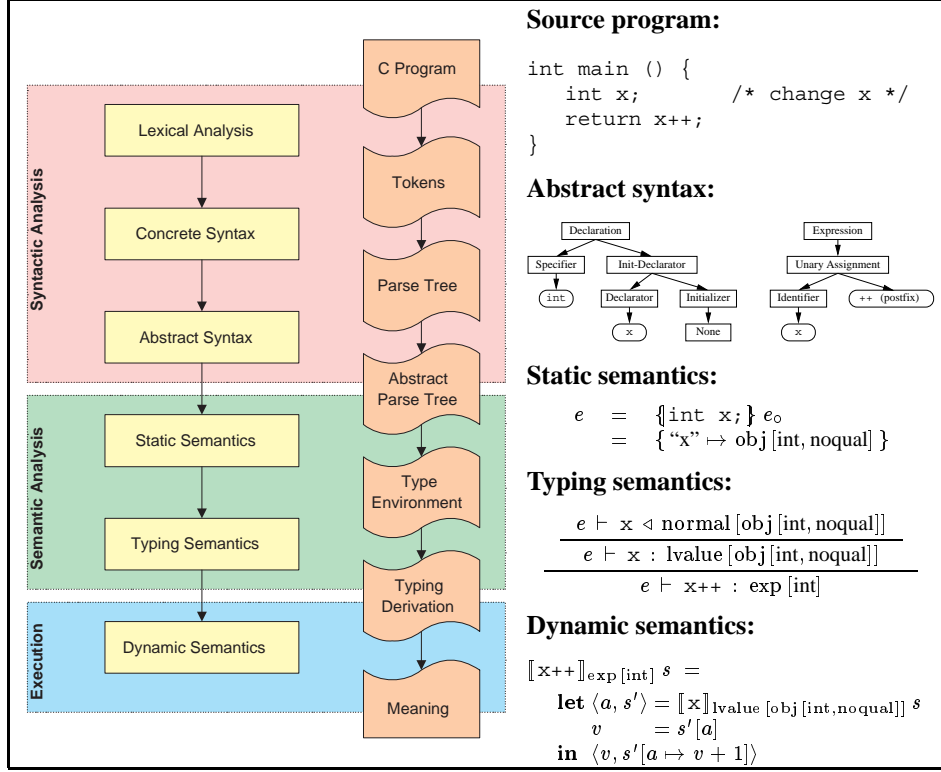
The interpreter consists of three layers, each containing a series of actions. *Syntactic analysis* aims at checking the syntactic validity of the source program and signalling syntax errors. Syntactically correct programs are transformed to abstract parse trees, which represent their structure in detail. *Semantic analysis* is used to define aspects of C that cannot be defined by a context-free grammar. It aims at checking the semantic validity of the program and signalling semantic errors, e.g. use of undeclared identifiers or type mismatches. Finally, *execution* aims at describing the meaning of programs. Our research mainly focuses on the last two layers, containing a total of three actions: *static semantics*, *typing semantics* and *dynamic semantics*. A brief overview of these actions and their collaboration is given in this section. The following three sections present the basics of each action in more detail.

Static semantics keeps track of identifiers that are defined in the source program. It aims at detecting static semantic errors, such as the redefinition of an identifier in the same scope, as well as associating identifiers with appropriate types or values. For each

¹Introductions to denotational semantics, including useful bibliography, can be found in [9] and [10].

²Introductions to monads and their use in denotational semantics can be found in [11] and [12].

Figure 1: An abstract interpreter for C.



syntactically well-formed program phrase P , its static semantic meaning is denoted by $\llbracket P \rrbracket$. Such meanings are typically types, type environments, i.e. associations of identifiers to types, or functions handling these two. Considering the simple C program that is used as an example in Fig. 1, let us isolate the declaration “`int x;`”. The static semantic meaning of this declaration is a function that updates the type environment by declaring an integer variable “`x`”. If e_0 is the empty type environment, containing no declarations, then the result of $\{ \text{int } x; \}$ applied to e_0 is the updated environment e shown in Fig. 1, which contains a declaration for “`x`”.

Typing semantics focuses on program phrases. It aims at the detection of type-mismatch errors, such as assignment to a constant value, and at associating syntactically well-formed phrases with appropriate phrase types. Such associations are given by means of *typing derivations*, i.e. formal proofs that phrases are well-typed. Typing derivations use inference rules to prove *typing judgements*, such as “ $e \vdash P : \theta$ ”, which states that phrase P has type θ in environment e . In the same example program of Fig. 1, let us now consider the expression “`x++`”. Assuming that the static semantic analysis has resulted in type environment e , it is possible to derive that “`x++`” is an expression that computes an integer value, in other words “ $e \vdash x++ : \text{exp}[\text{int}]$ ”. A sketch of such a typing derivation is shown in Fig. 1. The derivation makes use of two inference rules, stating in short that

variables are l-values and how l-values can be used as operands of postfix “++”. The initial assumption means that “x” is declared to be a variable of type int in e , and is not further analyzed here.

Finally, dynamic semantics aims primarily at defining the execution behaviour of well-typed programs. For each well-typed program phrase P of type θ , its dynamic semantic meaning is denoted by $\llbracket P \rrbracket_\theta$. Such a meaning is typically a function describing some aspect of the execution of P . The typing derivation for P is important since it determines the way in which the dynamic semantics will be calculated. Going back once more to the same example, let us consider again the expression “x++”, together with its typing derivation. Assuming a simple direct semantics for the dynamic meaning of expressions, the dynamic semantics for “x++” is given in Fig. 1. It is a function, taking the initial program state s and returning the result of the expression’s evaluation and the final program state. Notice that the typing derivation dictates the types that are used for the dynamic semantics of phrases on both sides of the equation. For the sake of clarity, a represents the address of the object designated by the l-value “x”, s' is the program state after evaluating “x” and v is the value stored in a at the program state s' . The result of the evaluation is v , and the final program state is the same as s' , with the value stored in a incremented by one.

3. Static semantics

The static semantics of C can be thought of as the symbol table in our abstract interpreter. It calculates the environments containing type information for all identifiers defined in the source program and, for this reason, it mainly deals with the program’s declarations. At the same time static semantic errors are detected. Apart from the complicated syntax of declarations that is characteristic of C, the static semantics is further perplexed by the presence of incomplete types in C’s type system. Forward declarations of tags used in the recursive definition of structures and unions are also sources of complexity.

The domains that we use in the static semantics of C are summarized in Fig. 2. Most of them are defined as coalesced sums, by enumeration of their elements. The domain ordering relation \sqsubseteq is crucial in the treatment of incomplete types: $x \sqsubseteq y$ denotes that y is a better approximation of a possibly incomplete element x . Notice the number of different types that are dictated by C’s type system. Among them, data types provide the basis for the type system, representing types that C programs can manipulate as first class elements. Object types are associated with objects in memory and consist of qualified versions of data types and array types. Denotable types are associated with identifiers in type environments, while identifier types are used for the classification of these identifiers. Member types are associated with identifiers defined as members of structures or unions. Finally, phrase types are associated with program phrases by the typing semantics of §4. Most domains for environments can be taken as functions from identifiers to types; their full definition is omitted in this chapter.

The use of a simple error monad in the definition of the static semantics provides an elegant way of generating and propagating errors. In brief, a monad is a tuple $\langle M, \text{unit}, * \rangle$, where M is a domain constructor, $\text{unit} : A \rightarrow M(A)$ and $* : M(A) \times (A \rightarrow M(B)) \rightarrow M(B)$ are polymorphic functions for arbitrary domains A and B , satisfying three monad

Figure 2: Static semantic domains.

Auxiliary domains.		
$I : \mathbf{Ide}$ (identifiers), $t : \mathbf{Tag}$ (tags), $\sigma : \mathbf{TagType}$ (tag types).		
Domains for types.		
$\tau : \mathbf{Type}_{dat}$	= void char signed-char unsigned-char short-int unsigned-short-int int unsigned-int long-int unsigned-long-int float double long-double ptr $[\phi]$ enum $[\epsilon]$ struct $[t, \pi]$ union $[t, \pi]$	(data types)
$q : \mathbf{Qual}$	= noqual const volatile const-volatile	(type qualifiers)
$\alpha : \mathbf{Type}_{obj}$	= obj $[\tau, q]$ array $[\alpha, n]$	(object types)
$f : \mathbf{Type}_{fun}$	= func $[\tau, p]$	(function types)
$\phi : \mathbf{Type}_{den}$	= α f	(denotable types)
$m : \mathbf{Type}_{mem}$	= α bitfield $[\tau, q, n]$	(member types)
$\beta : \mathbf{Type}_{bit}$	= int signed-int unsigned-int	(bit-field types)
$v : \mathbf{Type}_{val}$	= τ f	(value types)
$\delta : \mathbf{Type}_{ide}$	= normal $[\phi]$ typedef $[\phi]$ enum-const $[n]$	(identifier types)
$\theta : \mathbf{Type}_{phr}$	(See Fig. 3)	(phrase types)
Domains for environments.		
$e : \mathbf{Ent}$ (types), $\epsilon : \mathbf{Enum}$ (enumerations), $\pi : \mathbf{Memb}$ (members), $p : \mathbf{Prot}$ (function prototypes).		

laws. Given a domain D , the domain $M(D)$ is a domain of *computations* resulting in values of type D . Function unit converts values to (trivial) computations and the binary operator $*$ is used to extract the results of computations.³ Monad E can be defined by taking $E(D) = D \oplus \mathbf{U}$ where \mathbf{U} is a singleton domain, whose element represents errors.

Using monad E to allow static errors, the static meaning of declarations can be defined as a function of the domain:

$$\blacktriangleright \llbracket \text{declaration} \rrbracket : \mathbf{Ent} \rightarrow E(\mathbf{Ent})$$

that is, a function which takes as argument an initial type environment and returns an updated type environment, which contains information about the declared identifiers. The equations defining such functions are often very complex but can be somewhat simplified by the use of several auxiliary operations, defined separately.

The meaning of recursively defined types requires special treatment. In a way similar to the one suggested in [2], a monadic closure operator is used, defined as:

$$\text{mclo } z \ f = \mathbf{clo } z (\lambda x. x * f) = \bigsqcup_{n=0}^{\infty} (\lambda x. x * f)^n z$$

on condition that f is continuous and $z \sqsubseteq z * f$. This operator is applied to the initial type environment and the static meaning of a list of declarations, in order to obtain a least upper bound for recursively defined types.

³This description, although naïve, is sufficient for the purpose of this chapter.

Figure 3: Some of the phrase types.

Phrase type	Description
$\text{exp } [v]$	Expression, whose result is a non-constant r-value of type v .
$\text{lvalue } [m]$	Expression, whose result is an l-value of type m .
$\text{val } [\tau]$	Expression, whose result is a constant r-value of type τ .
$\text{arg } [p]$	Actual arguments of a function with prototype p .
$\text{stmt } [\tau]$	Statement in a function returning a result of type τ .
decl	Declaration.

Figure 4: Typing judgements.

Main typing relation.	
$e \vdash \text{phrase} : \theta$	The given <i>phrase</i> can be attributed phrase type θ in type environment e .
Predicates as judgements.	
P	Predicate P is <i>true</i> , where P can be any valid predicate over truth values T .
$v := z$	The static semantic valuation $z : E(D)$ produces the (non-error) value $v : D$.
Judgements related to environments.	
$e \vdash I \triangleleft \delta$	Identifier I is associated with identifier type δ in type environment e .
$\pi \vdash I \triangleleft m$	Identifier I is associated with member type m in member environment π .
Judgements related to expressions.	
$e \vdash E \gg \tau$	Expression E can be assigned to an object of data type τ in type environment e .
$e \vdash E = \text{NULL}$	Expression E is a null pointer constant in type environment e .
$e \vdash T \equiv \phi$	Type name T denotes type ϕ in type environment e .

4. Typing semantics

The primary aim of typing semantics is the association of program phrases with phrase types. Fig. 3 shows all the phrase types that we use. Typing rules are inference rules whose premises and conclusion are typing judgements. Several forms of typing judgements are necessary in order to simplify the typing rules. A summary of the most important ones is given in Fig. 4. Approximately 200 typing rules are used in our approach, in order to specify the typing semantics of C. 70% of those deal with expressions, 10% with statements and the remaining 20% with declarations. In the rest of this section we will present some non-trivial examples of typing rules and we will conclude with a discussion on our typing semantics.

The following rules specify the typing semantics of four types of expressions, in accordance with the ANSI C standard. Rule E1 states that decimal constants with no suffix are attributed type $\text{val } [\tau]$, where τ is the first integer type that can represent their value. According to E2, identifiers that have been defined as normal variables in an e are l-values of the appropriate type. Similarly, rules E3 and E4 specify the typing semantics of function calls and the indirection operator.

$$\begin{array}{c}
\text{suffix}(n) = \emptyset \quad \text{isDecimal}(n) \\
\hline
\tau := \text{repType}(n, [\text{int}, \text{long-int}, \text{unsigned-long-int}]) \quad (E1) \\
e \vdash n : \text{val}[\tau]
\end{array}
\quad
\begin{array}{c}
e \vdash I \triangleleft \text{normal}[\alpha] \\
\hline
e \vdash I : \text{lvalue}[\alpha] \quad (E2)
\end{array}$$

$$\begin{array}{c}
e \vdash E : \text{exp}[\text{ptr}[\text{func}[\tau, p]]] \quad e \vdash \text{args} : \text{arg}[p] \\
\hline
e \vdash E(\text{args}) : \text{exp}[\tau] \quad (E3)
\end{array}
\quad
\begin{array}{c}
e \vdash E : \text{exp}[\text{ptr}[\alpha]] \\
\hline
e \vdash *E : \text{lvalue}[\alpha] \quad (E4)
\end{array}$$

The following two rules specify in part the semantics of assignments. According to E5, the expression on the left side of the simple assignment operator must be a modifiable l-value, while the expression on the right side must be assignable to the corresponding data type. Rule A1 states that an expression of arithmetic type τ can be assigned to an object of another arithmetic type τ' .

$$\begin{array}{c}
e \vdash E_1 : \text{lvalue}[m] \quad \text{isModifiable}(m) \\
\tau := \text{datify } m \quad e \vdash E_2 \gg \tau \\
\hline
e \vdash E_1 = E_2 : \text{exp}[\tau] \quad (E5)
\end{array}
\quad
\begin{array}{c}
e \vdash E : \text{exp}[\tau] \\
\text{isArithmetic}(\tau) \quad \text{isArithmetic}(\tau') \\
\hline
e \vdash E \gg \tau' \quad (A1)
\end{array}$$

A small number of typing rules specifies conversions that take place implicitly in the evaluation of expressions. Such conversions are called implicit coercions. For example, rule C1 states that l-values are implicitly converted to the values stored in the designated objects. Even more obviously, according to rule C2 constant values may be treated as normal non-constant expressions.

$$\begin{array}{c}
e \vdash E : \text{lvalue}[\text{obj}[\tau, q]] \quad \text{isComplete}(\tau) \\
\hline
e \vdash E : \text{exp}[\tau] \quad (C1)
\end{array}
\quad
\begin{array}{c}
e \vdash E : \text{val}[\tau] \\
\hline
e \vdash E : \text{exp}[\tau] \quad (C2)
\end{array}$$

The typing semantics of statements is specified in a similar way. Rule S1 deals with compound statements. Notice that a compound statement defines a new scope, containing its declarations, and therefore a new environment e' must be calculated, allowing for recursively defined structures or unions. This new environment e' is used for the typing of the compound statement's body. Rules S2 and S3, specifying the semantics of `while` and `return` statements respectively, are relatively easier. In S3, the type of the returned expression must be assignable to the function's returned type.

$$\begin{array}{c}
e' := \text{rec } \{\text{declaration-list}\} (\uparrow e) \\
e' \vdash \text{declaration-list} : \text{decl} \quad e' \vdash \text{statement-list} : \text{stmt}[\tau] \\
\hline
e \vdash \{\text{declaration-list statement-list}\} : \text{stmt}[\tau] \quad (S1)
\end{array}$$

$$\begin{array}{c}
e \vdash \text{expr} : \text{exp}[\tau'] \quad \text{isScalar}(\tau') \\
\hline
e \vdash \text{stmt} : \text{stmt}[\tau] \quad (S2)
\end{array}
\quad
\begin{array}{c}
e \vdash \text{expr} \gg \tau \\
\hline
e \vdash \text{return expr} ; : \text{stmt}[\tau] \quad (S3)
\end{array}$$

The suggested typing semantics for C leads to two forms of ambiguity problems. The first concerns the uniqueness of typing results: the main typing relation does not always provide a unique phrase type for a given program phrase. (Implicit coercion rules such as C1 and C2 are one source of such ambiguities.) This form of ambiguity is in fact useful. A given program phrase can be attributed different phrase types, depending on its role in the program, and a different dynamic semantic meaning may exist for each phrase type. The second form concerns the uniqueness of typing derivations for a given typing judgement.

Figure 5: Dynamic semantic domains.

Auxiliary domains and environments.
$\mathbf{Addr} = \mathbf{Obj} \times \mathbf{Offset}, \mathbf{Obj}, \mathbf{Fun}, \mathbf{Offset}, \mathbf{BitOfs}, \llbracket e \rrbracket_{Ent}, \llbracket p \rrbracket_{Prot}, \mathbf{Cod}, \mathbf{Lab}_\tau$
Domains for types.
$\llbracket \text{void} \rrbracket_{dat} = \mathbf{U}, \llbracket \text{int} \rrbracket_{dat} = \mathbf{N}, \llbracket \text{ptr}[\alpha] \rrbracket_{dat} = \mathbf{Addr} \oplus \mathbf{U}, \llbracket \text{ptr}[f] \rrbracket_{dat} = \mathbf{Fun} \oplus \mathbf{U}$
$\llbracket \text{obj}[\tau, q] \rrbracket_{obj} = \mathbf{Addr}, \llbracket \text{array}[\alpha, n] \rrbracket_{obj} = \mathbf{N} \rightarrow \llbracket \alpha \rrbracket_{obj}$
$\llbracket \text{func}[\tau, p] \rrbracket_{fun} = \llbracket p \rrbracket_{Prot} \rightarrow \mathbf{G}(\llbracket \tau \rrbracket_{dat})$
$\llbracket \alpha \rrbracket_{mem} = \llbracket \alpha \rrbracket_{obj}, \llbracket \text{bitfield}[\tau, q, n] \rrbracket_{mem} = \mathbf{Addr} \times \mathbf{BitOfs}$

For example, there are two different derivations concluding with the fact that the sum of two integer constants is an integer expression: the first adds the constants and coerces the constant sum using C3, while the second coerces the summands separately using C3 and adds the resulting expressions. It is required that all different possible derivations for a given typing judgement result in the same dynamic semantics for the program phrase.

5. Dynamic semantics

The dynamic semantics of C specify the execution behaviour of well-typed programs. As a useful side-effect, run-time errors and other sources of undefined behaviour are detected. The most important source of complexity in an accurate definition of C's dynamic semantics is the unspecified evaluation order, combined with the fact that expressions generate side effects. In order to disallow undesired ambiguities, the ANSI C standard has introduced restrictions imposed on expression evaluation with the mechanism of sequence points. Additional restrictions are imposed on the access of objects between consecutive sequence points; however, according to our interpretation of the standard this mechanism does not always prevent non-determinism.⁴ The dynamic semantics is further perplexed by pointer arithmetic, complex control statements like `for` and `switch`, and the presence of `goto` in combination with block scopes containing variable declarations.

For each static type, a dynamic semantic domain is defined for representing the dynamic meaning of values of this type. The definitions of some dynamic semantic domains are shown in Fig. 5. Among other things, the domain of integer numbers \mathbf{N} is used to represent values of type `int`, pointers to objects are represented by the objects' address or a special null value, and addresses of objects are offsets in the biggest (possibly aggregate) objects containing them, in order to correctly model pointer arithmetic. The domain for type environments $\llbracket e \rrbracket_{Ent}$ contains the dynamic meanings of all identifiers defined in e , while that for function prototypes contains the values of all parameters. Domain \mathbf{Cod} contains the code of all defined functions and \mathbf{Lab}_τ contains the meanings of labeled statements in a function returning a result of type τ and is used in the semantics of jumps.

A number of monads is used in order to represent various aspects of the computations that are related to the execution of C programs. Brief descriptions of these monads are

⁴This issue has been discussed a lot in `comp.std.c`. Several opinions have been expressed but no conclusion has been reached.

Figure 6: Monads used in the dynamic semantics.

Auxiliary monads and monad transformers.	
$P(D)$	Powerdomain monad, allowing for non-determinism in values of domain D .
$R(M)(T)$	Resumption monad transformer, allowing interleaving in computations of type $M(T)$. Defined as a solution of the equation $R(M)(T) = T \oplus M(R(M)(T))$.
Monads for computations.	
$V(T)$	Computation of a constant value of type T , not accessing the memory (value monad).
$C(T)$	Computation of a non-constant value of type T , possibly accessing the memory (continuation monad). Defined as $C(T) = (T \rightarrow \mathbf{C}) \rightarrow \mathbf{C}$, where $\mathbf{C} = \mathbf{S} \rightarrow P(\mathbf{A})$ is the domain of non-deterministic continuations, \mathbf{S} is the domain of program states and \mathbf{A} the domain of final program answers.
$G(T)$	Computation of a non-constant value of type T in the evaluation of an expression, possibly accessing the memory and allowing for interleaving. Defined as $G(T) = R(C)(T)$.
$K_\tau(T)$	Computation of a non-constant value of type T in the execution of a statement, possibly accessing the memory or terminating the function by returning a result of type τ .

shown in Fig. 6. The powerdomain monad, based on the convex powerdomain, is used to model non-determinism. The resumption monad transformer is used to model the interleaving in the evaluation of expressions, that is required by C 's unspecified order of evaluation.⁵

The definition of dynamic semantics for program phrases is similar to that of static semantics. An important difference, however, is that the typing derivations provide useful information about the semantic meanings of both the defined phrase and its components. Thus, typing derivations control the definition of dynamic semantics, instead of abstract syntax, and there is one dynamic equation for each typing rule. In the rest of this section we will illustrate the definition of dynamic semantics by presenting some small examples.

Let us consider the simple case of typing rule E4 in §4. The dynamic semantics for an expression of the form “ $\star E$ ” under the typing given in E4 can be defined as follows:

$$\begin{aligned}
 \blacktriangleright \quad \llbracket \text{exp } [v] \rrbracket & : e : \mathbf{Ent} \twoheadrightarrow \llbracket e \rrbracket_{Ent} \rightarrow \mathbf{Cod} \rightarrow G(\llbracket v \rrbracket_{val}) \\
 \llbracket \star E \rrbracket_{lvalue [\alpha]} & = \lambda e. \lambda \rho. \lambda \xi. \\
 \quad \llbracket E \rrbracket_{exp [ptr [\alpha]]} e \rho \xi & * (\lambda d_e. \mathbf{case } d_e \mathbf{ of} \\
 \quad \mathbf{inl } a & \Rightarrow \mathbf{unit } a \\
 \quad \mathbf{otherwise} & \Rightarrow \mathbf{error})
 \end{aligned}$$

The first line states that the dynamic semantic meaning for phrases of type $\text{exp } [\tau]$ is a function taking as arguments the static and dynamic environments and the code environment and returning an interleaved computation with a result of type $\llbracket \tau \rrbracket_{dat}$. Notice the dependent function type, denoted here as $x : A \twoheadrightarrow B(x)$, which dictates a connection between the static and the dynamic environments. The equation that follows defines $\llbracket \star E \rrbracket_{lvalue [\alpha]}$ in terms of $\llbracket E \rrbracket_{exp [ptr [\alpha]]}$. If the pointer contains an object's address, this address is used in the resulting l-value. An error occurs if the pointer is null. Dynamic semantic meanings for phrases of type $\text{lvalue } [m]$ are functions of the form:

$$\blacktriangleright \quad \llbracket \text{lvalue } [m] \rrbracket : e : \mathbf{Ent} \twoheadrightarrow \llbracket e \rrbracket_{Ent} \rightarrow \mathbf{Cod} \rightarrow G(\llbracket m \rrbracket_{mem})$$

⁵The resumption monad transformer is defined in detail in a paper that will be submitted for publication in the near future. For an introduction to monad transformers, the reader is referred to [13].

In a similar way, the following equations define the dynamic semantics that correspond to typing rules E3 and C1 of §4 respectively.

$$\begin{aligned}
\llbracket E(\text{args}) \rrbracket_{\text{exp}[\tau]} &= \lambda e. \lambda \rho. \lambda \xi. \\
&\quad \llbracket E \rrbracket_{\text{exp}[\text{ptr}[\text{func}[\tau, p]]]} e \rho \xi \bowtie \llbracket \text{args} \rrbracket_{\text{arg}[p]} e \rho \xi * (\lambda \langle d_e, d_p \rangle. \\
&\quad \text{seqpt} * (\lambda u. \text{case } d_e \text{ of} \\
&\quad \quad \text{inl } d_f \Rightarrow \text{let } \langle f, b_f \rangle = \xi[d_f] \text{ in isCompatible}(f, \text{func}[\tau, p]) \rightarrow b_f d_p, \text{ error} \\
&\quad \quad \text{otherwise} \Rightarrow \text{error})) \\
\llbracket E \rrbracket_{\text{exp}[\tau]} &= \lambda e. \lambda \rho. \lambda \xi. \llbracket E \rrbracket_{\text{lvalue}[\text{obj}[\tau, q]]} e \rho \xi * \text{getValue}_{\text{obj}[\tau, q] \rightarrow \tau}
\end{aligned}$$

The first thing to notice is the use of operator $\cdot \bowtie \cdot : G(A) \times G(B) \rightarrow G(A \times B)$ for the interleaving of two computations, in order to model the unspecified order in which the function's designator and its arguments are evaluated. The second is the use of $\text{seqpt} : G(\mathbf{U})$ to introduce a sequence point just before the function is actually called. The function's dynamic meaning is looked up in the code environment. Furthermore, the function's actual type must be compatible with the type of the designator that is used. In the second equation, the important point is the use of $\text{getValue}_{m \rightarrow \tau} : \llbracket m \rrbracket_{\text{mem}} \rightarrow G(\llbracket \tau \rrbracket_{\text{dat}})$ which retrieves a stored value from memory.

The dynamic semantics of statements is defined in a similar way. We give here two relatively simple examples, corresponding to the typing rules S2 and S3 of §4. Notice that the dynamic meaning of statements uses also a label environment of type \mathbf{Lab}_τ . The calculation of this environment requires a least fixed point, to allow infinite loops implemented by `goto` statements.

$$\begin{aligned}
\blacktriangleright \quad \llbracket \text{stmt}[\tau] \rrbracket &: e : \mathbf{Ent} \rightarrow \llbracket e \rrbracket_{\text{Ent}} \rightarrow \mathbf{Cod} \rightarrow \mathbf{Lab}_\tau \rightarrow K_\tau(\mathbf{U}) \\
\llbracket \text{while } (\text{expression}) \text{ statement} \rrbracket_{\text{stmt}[\tau]} &= \lambda e. \lambda \rho. \lambda \xi. \lambda \ell. \mathbf{fix}(\lambda g. \\
&\quad \text{lift}_{G \rightarrow K}(\llbracket \text{expression} \rrbracket_{\text{exp}[\tau']} e \rho \xi) * (\lambda d. \\
&\quad \quad \text{checkBoolean}_{\tau'}(d) \rightarrow \\
&\quad \quad \quad \text{setBreakContinue}(\text{unit } \mathbf{u}, g) (\llbracket \text{statement} \rrbracket_{\text{stmt}[\tau]} e \rho \xi \ell) * (\lambda u. g), \text{ unit } \mathbf{u})) \\
\llbracket \text{return expression} \rrbracket_{\text{stmt}[\tau]} &= \lambda e. \lambda \rho. \lambda \xi. \lambda \ell. \text{lift}_{G \rightarrow K}(\mathcal{A} \llbracket \text{expression} \rrbracket_{\text{exp}[\tau']} e \rho \xi) * \text{result}
\end{aligned}$$

The first equation uses the least fixed point operator \mathbf{fix} to model the semantics of the `while` statement. The correct continuations that will be used in case of a `break` or `continue` statement are passed to the body of the loop by using `setBreakContinue`. The second equation defines the semantics of the `return` statement. The expression is evaluated and converted as if by assignment to the function's return type. Notice the use of $\text{result} : \llbracket \tau \rrbracket_{\text{dat}} \rightarrow K_\tau(\mathbf{U})$ which signals the termination of this function and specifies the returned result.

6. Evaluation

A significant effort has been made to evaluate our approach in defining a denotational semantics for the C programming language. In this task, the major issue was to assess how complete and accurate the developed semantics is. Unfortunately, there is no systematic way to evaluate our approach and be absolutely certain that the results are correct: there is simply no way to compare a formal system of this complexity against an informal specification. For this reason, we have resorted in testing an interpreter that directly implements our semantics, by using some test suites for C implementations that were available.

An earlier version of our semantics was first implemented using SML as the implementation language. Later, SML was abandoned and Haskell was used instead, mainly because it has a richer type system, more flexible syntax, elegant support for monads and also because lazy evaluation avoids a number of non-termination problems. The current implementation consists of approximately 15,000 lines of Haskell code, which are distributed roughly as follows: 3,000 lines for the static semantics, 3,000 lines for the typing semantics, 5,000 lines for the dynamic semantics, 3,000 lines for parsing and pretty-printing and 1,000 more lines of general code and code related to testing. As it was expected, the implementation is very slow and this presents a serious handicap in our yet unfinished evaluation process, significantly limiting the size of test programs.

Although the evaluation of our semantics is still under way and minor bugs are waiting to be fixed, the results indicate that the developed semantics is complete and accurate to a great extent, with respect to the ANSI C standard. The most important deviations from the standard are that the developed semantics requires function prototypes to exist for all called functions, something already favoured by the current standard, and that storage specifiers other than `typedef` are currently ignored. Static variables may be preprocessed out, but a solution integrated in the semantics is currently investigated. As less important deviations, the developed semantics requires fully bracketed initializations and forbids the declaration of identifiers, other than labels, in expressions or statements.

7. Conclusion and future work

In this chapter, we have presented a summary of our work in developing a formal semantics for the ANSI C programming language, following the denotational approach. The developed semantics is satisfactorily complete and accurate, with respect to the standard. A significant contribution of our research, besides the developed semantics itself, is the application of monads and monad transformers for the specification of a real programming language. Furthermore, interesting results have been achieved in our attempt to model the interleaving of computations and non-determinism using monads and monad transformers, which may be useful in specifying the semantics of programming languages supporting parallelism.

Our research in the near future will focus on the process of evaluating and improving the developed semantics. Beyond that, we would like to study the practical applications that a formal semantics for C may have in the software industry, especially in tools for program transformation, debugging and understanding. The implementation of the developed semantics also gave rise to an interesting question: what are the characteristics of a programming language that make it suitable for implementing denotational specifications, especially using monadic notation? Finally, another direction for future research aims at studying and specifying the semantics of C's object-oriented descendants, C++ and Java.

References

- [1] American National Standards Institute, New York, NY. *ANSI/ISO 9899-1990, American National Standard for Programming Language: C*, 1990. Revision and redesignation of ANSI X3.159-1989.

- [2] R. Sethi. A case study in specifying the semantics of a programming language. In *Proceedings of the 7th Annual ACM Symposium on Principles of Programming Languages*, pages 117–130, January 1980.
- [3] R. Sethi. Control flow aspects of semantics-directed compiling. *ACM Transactions on Programming Languages and Systems*, 5(4):554–595, October 1983.
- [4] Y. Gurevich and J. K. Huggins. The semantics of the C programming language. In E. Börger et al., editors, *Selected Papers from CSL'92 (Computer Science Logic)*, volume 702 of *Lecture Notes in Computer Science*, pages 274–308. Springer Verlag, New York, NY, 1993.
- [5] J. Cook and S. Subramanian. A formal semantics for C in Nqthm. Technical Report 517D, Trusted Information Systems, October 1994.
- [6] J. Cook, E. Cohen, and T. Redmond. A formal denotational semantics for C. Technical Report 409D, Trusted Information Systems, September 1994.
- [7] M. Norrish. An abstract dynamic semantics for C. Technical Report TR-421, University of Cambridge, Computer Laboratory, May 1997.
- [8] N. S. Papaspyrou. *A Formal Semantics for the C Programming Language*. PhD thesis, National Technical University of Athens, Software Engineering Laboratory, February 1998.
- [9] R. D. Tennent. The denotational semantics of programming languages. *Communications of the ACM*, 19(8):437–453, August 1976.
- [10] P. D. Mosses. Denotational semantics. In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science*, volume B, chapter 11, pages 577–631. Elsevier Science Publishers B.V., 1990.
- [11] E. Moggi. An abstract view of programming languages. Technical Report ECS-LFCS-90-113, University of Edinburgh, Laboratory for Foundations of Computer Science, 1990.
- [12] P. Wadler. The essence of functional programming. In *Proceedings of the 19th Annual Symposium on Principles of Programming Languages (POPL'92)*, January 1992.
- [13] S. Liang, P. Hudak, and M. Jones. Monad transformers and modular interpreters. In *Conference Record of the 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL'95)*, San Francisco, CA, January 1995.