

A Methodology for the Definition of Programming Languages

NIKOLAOS S. PAPASPYROU

Department of Electrical and Computer Engineering
National Technical University of Athens
Polytechnioupoli, 15780 Zografou, Athens
GREECE
E-mail: nickie@softlab.ntua.gr

Abstract: The formal definition of a programming language is a valuable tool for the study, design, evaluation and even for the implementation of the language. A methodology for developing formal definitions of programming languages is proposed in this paper. We follow the denotational approach and use monads, in order to improve the modularity and elegance of the result. The definition of semantics is divided in three distinct consecutive phases: static, typing and dynamic semantics, each using the results of the previous two. Emphasis has been given on the straightforward translation of the developed definitions into abstract interpreters, in order to obtain verifiably correct prototype implementations for the studied languages. The methodology is described through a complete example, defining a small imperative programming language. Its use for the definition of ANSI C has resulted in considerable success.

Key-Words: Programming languages, formal definition, denotational semantics, monads.

1 Introduction

The study of programming languages invariably distinguishes between two fundamental features: *syntax* and *semantics*. Although the two cannot always be clearly separated, syntax refers to the appearance and structure of the well-formed sentences of the language, whereas semantics refers to the meanings of these sentences. Every person who uses a programming language must understand both its syntax and semantics at some level of abstraction. Programmers usually understand a programming language by means of examples, intuition and descriptions in natural language. Such descriptions are informal, typically based on a set of assumptions about the reader's knowledge and understanding, and therefore inherently ambiguous.

In the definition of programming languages, the syntax usually formally specified. On the other hand, semantics is most commonly specified in an informal way. Research in the area of formal programming language semantics started in the 1960s. Since then, the product of more than three decades of research has been the development and thorough study of numerous methods and formalisms, usually classified for histori-

cal reasons as following one of three main approaches: *operational*, *denotational* or *axiomatic* semantics.

Denotational semantics is a formalism introduced by Scott and Strachey in the late 1960s [1, 2, 3]. According to the denotational approach, semantics is described by attributing mathematical *denotations* to programs and program segments. Denotations are elements of appropriate mathematical structures called *domains*. Variations of the λ -calculus over domains are commonly used as metalanguages for expressing denotations.

The main objective of our research is to develop and evaluate a methodology for the formal definition of programming languages, with the main effort lying on the definition of semantics. Our methodology is described in this paper through a simple, yet not trivial example of a programming language called \mathcal{L} . Three distinct consecutive phases for the definition of semantics are proposed, each phase using the results of the previous ones. The methodology has been used with considerable success for the definition of ANSI C [4].

The rest of the paper is structured as follows. In Section 2, an overview of the proposed methodology

is given and the various stages in the formal definition of a programming language are briefly presented. Section 3 contains a short introduction to monads and their correspondence with computations. In the following sections 4 through 7 the various stages of the proposed methodology are presented in detail, using a small imperative programming language called \mathcal{L} as an example. Section 8 discusses the implementation of programming language definitions obtained by the proposed methodology, in the form of prototypical interpreters. Section 9 attempts an evaluation of the proposed methodology, discussing the results of an ambitious project concerning the definition of ANSI C. Finally, Section 10 concludes with a few remarks and directions for future research.

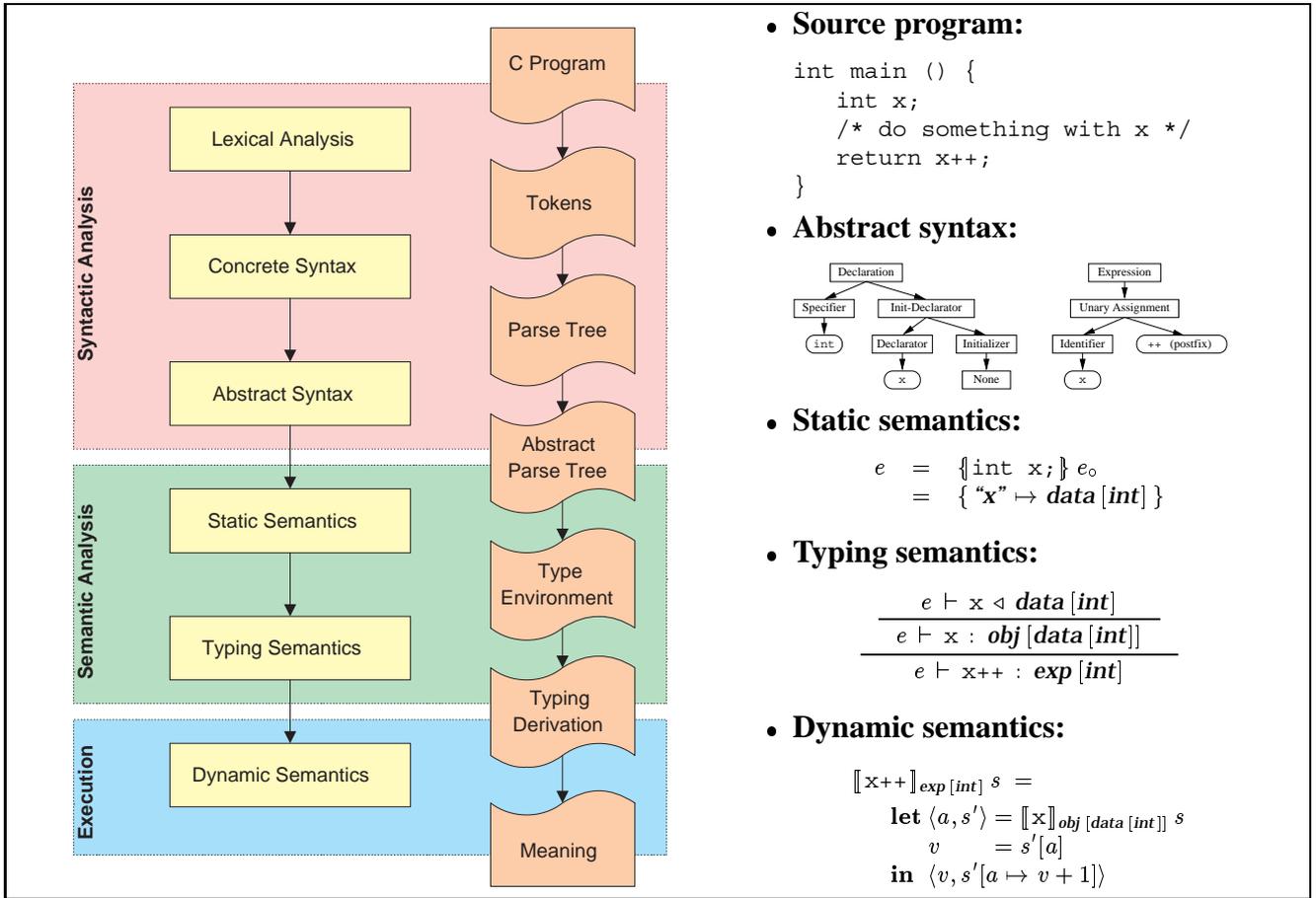
2 Overview of the methodology

According to the proposed methodology, the formal definition of a given programming language \mathcal{L} can be best understood as part of an *abstract interpreter* for programs written in \mathcal{L} . Such an interpreter is depicted in Fig. 1 as a data process. The left part of the figure is a module diagram of the interpreter. It shows the chain of actions performed by the interpreter as well as the data that is processed by these actions. Each action uses the results of the 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 what the program does when executed. The right part of the figure gives a small example illustrating the function of the interpreter, where the source language is taken to be C.

The interpreter consists of the following three layers, represented in Fig. 1 as big boxes. Each layer contains 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 syntactic structure in detail. This layer consists of three actions:
 - *Lexical analysis*: transforms the source program to a sequence of lexical units, also called *tokens*.
 - *Concrete syntax*: checks the syntactic validity of the source program by grouping tokens together in more complex syntax entities.
 - *Abstract syntax*: simplifies the parse tree by extracting superfluous information and results in an abstract parse tree.
- *Semantic analysis*: aims at checking the semantic validity of the program, as represented by the abstract parse tree, and signalling semantic errors, e.g. use of undeclared identifiers or type mismatches. It contains two actions:
 - *Static semantics*: aims at detecting static semantic errors, e.g. redefinition of an identifier in the same scope, as well as associating identifiers with appropriate types or values. Static semantics is based solely on the abstract syntax. For each syntactically well-formed program phrase P , its static semantic meaning is denoted by $\llbracket P \rrbracket$. Static semantic meanings are mathematical objects, typically type environments associating identifiers to types.
 - *Typing Semantics*: aims at detecting type mismatch errors, e.g. 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, that is, formal proofs that phrases are well-typed. Typing semantics is based not only on the abstract syntax of the program, but also on static semantic environments.
- *Execution*: the main part of the interpreter. It contains only one action:
 - *Dynamic Semantics*: aims at defining the execution behaviour of well-typed programs. As a useful side effect, run-time errors and other sources of undefined behaviour are detected at the same time. Dynamic semantics is based on the abstract syntax, the static semantic environments and the typing derivations. For each well-typed program phrase P of type θ , its dynamic semantic meaning is denoted by $\llbracket P \rrbracket_\theta$. Such meanings are also mathematical objects, typically functions describing aspects of the execution of the corresponding program phrases. The typing derivation for P is important, as the same phrase is allowed to have different meanings when attributed different types.

Fig. 1: An abstract interpreter for \mathcal{L} .



3 Computations and monads

The concept of monads comes from category theory. Monads have been proposed by Moggi as a useful structuring tool for denotational semantics [5, 6]. Moggi demonstrated their use for representing different aspects of computations and defined monads for programming language features such as state, exceptions and continuations. In a short time, the idea of monads became very popular in the functional programming community as a way of structuring functional programs and simulating non-functional features. The work of Wadler [7] played a very important role in this direction.

In the last few years, research related to the application of monads in denotational semantics has focused on the combination of monads to structure semantic interpreters. Monad transformers, which were also first proposed by Moggi, have attracted the attention of many researchers. In the work of Liang, Hudak

and Jones [8], monad transformers are demonstrated to successfully modularize semantic interpreters and the lifting of monad operations is investigated.

The notion of computation is very significant in the semantics of programming languages. If D is a domain of values, then $M(D)$ can be taken to denote a domain of computations which return values from D . In this sense, M can be thought of as a domain constructor which specifies the characteristics of computations, e.g. whether they are deterministic, whether they require access to the program state, or whether they can cause run-time errors. This is the principal concept behind the use of monads in denotational semantics. In brief, a monad is a triple $\langle M, \text{unit}_M, *_M \rangle$, where M is a domain constructor, $\text{unit}_M : A \rightarrow M(A)$ and $\cdot *_M \cdot : M(A) \times (A \rightarrow M(B)) \rightarrow M(B)$ are polymorphic functions for arbitrary domains A and B , satisfying three monad laws.

The definition of M reflects our notion of computa-

tion. The computation denoted by $\mathit{unit}_M v$ is usually the simplest computation returning the value v . If m is a computation returning a value v in D , then the computation denoted by $m *_M f$ is the combined computation of m followed by computation $f v$. Many different types of computations are implicit in the semantics of programs and a number of monads is required to represent them. Typical examples of such computations and related monads are:

- *Error monad*: Computations allowing errors.
- *Value monad*: Constant computations.
- *State monad*: Computations accessing the state
- *Continuation monad*: Computations accessing the state and allowing unrestricted jumps.
- *Powerdomain monad*: Non-deterministic computations.

4 Syntax

The formal definition of a programming language begins with the definition of its syntax. There are various standard formalisms for this purpose, the most widely known and used being context-free grammars, usually expressed in the Backus Naur Form (BNF) and its variations. This formalism is used in the proposed methodology. As a convention, non-terminal symbols are written in italics. The symbol $::=$ separates the two sides of production rules, whereas the symbol $|$ separates alternative productions. Finally, the symbol ϵ denotes the empty string.

Since the area of formal syntax specification has been thoroughly studied and is not so interesting from a researcher's point of view, our definition of the syntax of the example language \mathcal{L} will be slightly informal. We will only focus on the abstract syntax of \mathcal{L} , which is given in Fig. 2. The intended semantics of \mathcal{L} can be easily deduced from the abstract syntax. The language supports variables of just one data type (integer), functions and procedures to which parameters can be passed by value or by reference, a variety of statements including a loop statement and a variety of expression operators.

5 Static semantics

Static semantics 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. The definition of static semantics can be very complicated in the case of programming languages with rich type systems. It consists of two steps. At the first step, the static semantic domains must be defined, together with the basic operations that are allowed on the elements of each domain. For this reason, a careful analysis of the type system of \mathcal{L} is necessary.

Computations performed during this step are compile-time computations and can typically be represented by a simple error monad E , defined below. An additional operation on the monad, used for the generation of compile-time errors is given and it is easy to prove that errors are correctly propagated through the use of the standard monad operations.

- $E(D) = D \oplus U$
- $\mathit{unit}_E : D \rightarrow E(D)$
 $\mathit{unit}_E = \mathit{inl}$
- $\cdot *_E \cdot : E(A) \times (A \rightarrow E(B)) \rightarrow E(B)$
 $m *_E f = \mathbf{case\ } m \mathbf{ of}$
 $\mathit{inl}\ v \Rightarrow f\ v$
 $\mathit{inr}\ u \Rightarrow \mathit{inr}\ u$
- $\mathit{error}_E : E(D)$
 $\mathit{error}_E = \mathit{inr}\ u$

From the abstract syntax of \mathcal{L} , it is apparent that at least two classes of types exist. A *data type* is a basic type of the language: one can assign values of such a type to variables, pass them as parameters to functions, use them in expressions, etc. The only data type supported by \mathcal{L} is *integer*. On the other hand, a *return-type* is a type that can be returned from a function. Apart from the *integer* data type, the empty type *nothing* can also be returned from a function.

At a closer look, one can identify more classes of types. Another important class in the study of programming languages is the class of *denotable types*, that is, types that can be denoted by identifiers but whose values cannot always be assigned to variables or passed to functions as parameters. Functions in \mathcal{L} , are typical such values. They can only be used in expressions and statements in a limited way, i.e. in function calls. Identifiers in \mathcal{L} can denote three kinds of things: data objects (variables or parameters passed by-value), reference objects (parameters passed by-reference) and function names.

Fig. 2: Abstract syntax of \mathcal{L} .

```

program ::= body
body ::= definition-list statement-list
definition-list ::=  $\epsilon$  | definition definition-list
definition ::= var  $I$  as data-type | fun  $I$  ( formal-parameter-list ) as return-type body end
data-type ::= integer
return-type ::= data-type | nothing
formal-parameter-list ::=  $\epsilon$  | formal-parameter formal-parameter-list
formal-parameter ::=  $I$  as data-type byval |  $I$  as data-type byref
statement-list ::=  $\epsilon$  | statement statement-list
statement ::= skip |  $I$  := expression | begin statement-list end | if condition then statement else statement
| while condition do statement | call  $I$  ( expression-list ) | return | return expression
expression ::=  $n$  |  $I$  | call  $I$  ( expression-list ) | un-op expression | expression bin-op expression
condition ::= true | false | not condition | expression rel-op expression | condition log-op condition
expression-list ::=  $\epsilon$  | expression expression-list
un-op ::= + | -
bin-op ::= + | - | * | / | mod
rel-op ::= = | <> | > | < | >= | <=
log-op ::= and | or

```

Expression types form yet another type class. Apart from data types, expressions can also be of **boolean** type which can only be used in conditions. Finally, the class of *phrase types* contains one type for each kind of program phrase. All the type classes that were mentioned are defined below in the form of static semantic domains. The values of the domain for phrase types are further explained in Fig. 3.

- $I : \mathbf{Ide}$ (domain of identifiers)
- $\tau : \mathbf{Type}_{dat} = \mathit{integer}$
- $r : \mathbf{Type}_{ret} = \tau \mid \mathit{nothing}$
- $\phi : \mathbf{Type}_{den} = \mathit{data}[\tau] \mid \mathit{ref}[\tau] \mid \mathit{func}[r, p]$
- $v : \mathbf{Type}_{exp} = \tau \mid \mathit{boolean}$
- $\theta : \mathbf{Type}_{phr} = \mathit{prog} \mid \mathit{body}[r] \mid \mathit{stmt}[r]$
 $\mid \mathit{decl} \mid \mathit{exp}[v] \mid \mathit{obj}[\phi]$
 $\mid \mathit{arg}[p] \mid \mathit{prot}[p] \mid \mathit{par}[\tau, m]$
- $m : \mathbf{Mode} = \mathit{byval} \mid \mathit{byref}$

As a particular case of static domains, *static environments* are mappings from identifiers to appropriate static type domains. In the case of \mathcal{L} , static type environments are elements of domain \mathbf{Env} and map identifiers of a program scope to denotable types. The domain \mathbf{Env} and the required operations on its elements are shown below as an abstract data type. Its

Fig. 3: Static semantics, domain of phrase types.

Type	Description
prog	Complete program
body [r]	Body of function returning a result of type r
stmt [τ]	Statement in a function returning a result of type r .
decl	Declaration.
exp [v]	Expression, whose result is an r-value of type v .
obj [ϕ]	Expression, whose result is an l-value of type ϕ .
arg [p]	Actual parameters to a function with prototype p .
prot [p]	Formal parameters to a function with prototype p .
par [τ, m]	Formal parameter of type τ and passing mode m .

implementation is slightly perplexed because of nested scopes and is not included in this paper.

- $e : \mathbf{Env}$ (domain of type environments)
- ▶ $e_o : \mathbf{Env}$ (empty environment)
- ▶ $\cdot[\cdot] : \mathbf{Env} \times \mathbf{Ide} \rightarrow \mathbf{E}(\mathbf{Type}_{den})$ (lookup function)
- ▶ $\cdot[\cdot \mapsto \cdot] : \mathbf{Env} \times \mathbf{Ide} \times \mathbf{Type}_{den} \rightarrow \mathbf{E}(\mathbf{Env})$ (update function)
- ▶ $\uparrow \cdot : \mathbf{Env} \rightarrow \mathbf{Env}$ (new enclosed scope)
- ▶ $\downarrow \cdot : \mathbf{Env} \rightarrow \mathbf{Env}$ (enclosing scope)

A second kind of static environment is also shown below as an abstract data type. It represents function

Fig. 4: Static semantics for \mathcal{L} , functions and equations.

▶	$\{\{program\}\} : \mathbf{E}(\mathbf{Env})$ $\{\{body\}\} = \{\{body\}\} e_o$
▶	$\{\{body\}\} : \mathbf{Env} \rightarrow \mathbf{E}(\mathbf{Env})$ $\{\{declaration-list\} \text{ statement-list}\} = \{\{declaration-list\}\}$
▶	$\{\{declaration-list\}\} : \mathbf{Env} \rightarrow \mathbf{E}(\mathbf{Env})$ $\{\{\epsilon\}\} = \mathbf{unit}$ $\{\{declaration\} \text{ declaration-list}\} = \{\{declaration\}\} ; \{\{declaration-list\}\}$
▶	$\{\{declaration\}\} : \mathbf{Env} \rightarrow \mathbf{E}(\mathbf{Env})$ $\{\{\text{var } I \text{ as } data\text{-type}\}\} = \lambda e. \{\{data\text{-type}\}\} * (\lambda \tau. e[I \mapsto \mathbf{data}[\tau]])$ $\{\{\text{fun } I \text{ (} formal\text{-parameter-list) \text{ as } return\text{-type } body \text{ end}\}\} = \lambda e.$ $\{\{return\text{-type}\}\} * (\lambda r. \{\{formal\text{-parameter-list}\}\} * (\lambda p. e[I \mapsto \mathbf{func}[r, p]]))$
▶	$\{\{formal\text{-parameter-list}\}\} : \mathbf{E}(\mathbf{Prot})$ $\{\{\epsilon\}\} = \mathbf{unit} p_o$ $\{\{formal\text{-parameter-list}\} \text{ formal\text{-parameter-list}\} =$ $\{\{formal\text{-parameter}\}\} * (\lambda \langle \tau, m \rangle. \{\{formal\text{-parameter-list}\}\} * (\lambda p. \langle \tau, m \rangle \leq p))$
▶	$\{\{formal\text{-parameter}\}\} : \mathbf{E}(\mathbf{Type}_{dat} \times \mathbf{Mode})$ $\{\{I \text{ as } data\text{-type} \text{ byval}\}\} = \{\{data\text{-type}\}\} * (\lambda \tau. \mathbf{unit} \langle \tau, \mathbf{byval} \rangle)$ $\{\{I \text{ as } data\text{-type} \text{ byref}\}\} = \{\{data\text{-type}\}\} * (\lambda \tau. \mathbf{unit} \langle \tau, \mathbf{byref} \rangle)$
▶	$\mathcal{F}\{\{formal\text{-parameter-list}\}\} : \mathbf{Env} \rightarrow \mathbf{E}(\mathbf{Env})$ $\mathcal{F}\{\{\epsilon\}\} = \mathbf{unit}$ $\mathcal{F}\{\{formal\text{-parameter}\} \text{ formal\text{-parameter-list}\} = \mathcal{F}\{\{formal\text{-parameter}\}\} ; \mathcal{F}\{\{formal\text{-parameter-list}\}\}$
▶	$\mathcal{F}\{\{formal\text{-parameter}\}\} : \mathbf{Env} \rightarrow \mathbf{E}(\mathbf{Env})$ $\mathcal{F}\{\{I \text{ as } data\text{-type} \text{ byval}\}\} = \lambda e. \{\{data\text{-type}\}\} * (\lambda \tau. e[I \mapsto \mathbf{data}[\tau]])$ $\mathcal{F}\{\{I \text{ as } data\text{-type} \text{ byref}\}\} = \lambda e. \{\{data\text{-type}\}\} * (\lambda \tau. e[I \mapsto \mathbf{ref}[\tau]])$
▶	$\{\{data\text{-type}\}\} : \mathbf{E}(\mathbf{Type}_{dat})$ $\{\{\text{integer}\}\} = \mathbf{unit} \text{integer}$
▶	$\{\{return\text{-type}\}\} : \mathbf{E}(\mathbf{Type}_{ret})$ $\{\{data\text{-type}\}\} = \{\{data\text{-type}\}\}$ $\{\{\text{nothing}\}\} = \mathbf{unit} \text{nothing}$

prototypes, i.e. mappings from the formal parameters of a function to their type and passing mode. The order of formal parameters is important and this is why they are represented by integer numbers instead of names.

- $p : \mathbf{Prot}$ (domain of function prototypes)
- ▶ $p_o : \mathbf{Prot}$ (empty prototype)
- ▶ $\cdot[\cdot] : \mathbf{Prot} \times \mathbf{N} \rightarrow \mathbf{E}(\mathbf{Type}_{dat} \times \mathbf{Mode})$ (lookup function)
- ▶ $\cdot \leq \cdot : (\mathbf{Type}_{dat} \times \mathbf{Mode}) \times \mathbf{Prot} \rightarrow \mathbf{E}(\mathbf{Prot})$ (prepend parameter)
- ▶ $\Downarrow \cdot : \mathbf{Prot} \rightarrow \mathbf{E}(\mathbf{Type}_{dat} \times \mathbf{Mode} \times \mathbf{Prot})$ (extract 1st parameter)

At the second step, for each non-terminal symbol involved in the static semantics, a *semantic function* must be defined. This function represents the static meaning of phrases generated by this non-terminal symbol and its definition is given by induction on the

set of production rules. For the case of \mathcal{L} , the complete set of static semantic functions is defined in Fig. 4. There is one equation for each production rule. The type of each function reflects the static meaning of the corresponding program phrase, e.g. the meaning of declarations is typically a function from some initial type environment to a compile-time computation of an updated type environment, which contains the newly declared identifiers.

In some cases, program phrases may have more than one static meaning, depending on their use. For example, formal parameters determine the prototype of a defined function, but also modified the static type environment when this function is called. The multiple meanings are distinguished by prefixed calligraphic letters. Furthermore, most programming languages support recursively defined types, requiring the use of the least fixed point operator in their static semantics.

6 Typing semantics

The primary aim of typing semantics is the association of program phrases with phrase types. The relation between program phrases and phrase types is formally established by means of *typing judgements*. The most common form of judgement is the main typing relation $e \vdash \textit{phrase} : \theta$. However, in specifying the typing semantics of a complex programming language, additional forms of typing judgements may be necessary. Fig. 5 shows the forms of judgements that are needed for the typing semantics of \mathcal{L} .

Fig. 5: Typing semantics, judgements.

Judgement	Description
$e \vdash \textit{phrase} : \theta$	The given <i>phrase</i> can be attributed phrase type θ in type environment e .
$e \vdash I \triangleleft \delta$	Identifier I is associated with the denotable type ϕ in type environment e .
$v := z$	The compile-time computation $z : \mathbf{E}(D)$ produces the (non-error) value $v : D$.

Typing rules are inference rules whose premises and conclusion are typing judgements. The proof of a typing judgement using a set of typing rules is called *typing derivation*. A typing derivation combines a number of typing rules in a tree-like structure, in such a way that the conclusion of one rule becomes a premise of some other rule.

A large subset of the rules that define the typing semantics of \mathcal{L} is given in Fig. 6. Rules I1 and I2 specify conversions that take place implicitly in the evaluation of expressions, stating that l-values are implicitly converted to the values stored in the designated objects. Such conversions are called *implicit coercions*.

7 Dynamic semantics

Dynamic semantics specifies the execution behaviour of well-typed programs. As a useful side-effect, run-time errors and other sources of undefined behaviour are detected. As in the case of static semantics, dynamic semantics is specified in two steps. The first step aims at the definition of dynamic semantic domains, whereas the second aims at the definition of dynamic semantic functions.

The definition of dynamic semantic domains begins

with a mapping from static types, i.e. elements of static type domains, to dynamic domains representing the values of these types. In the case of \mathcal{L} , this mapping is done below.

- $\llbracket \textit{integer} \rrbracket_{dat} = \mathbf{N}$
- $\llbracket \tau \rrbracket_{ret} = \llbracket \tau \rrbracket_{dat}$
- $\llbracket \textit{nothing} \rrbracket_{ret} = \mathbf{U}$
- $\llbracket \textit{data}[\tau] \rrbracket_{den} = \mathbf{Obj}$
- $\llbracket \textit{ref}[\tau] \rrbracket_{den} = \mathbf{Obj}$
- $\llbracket \textit{func}[r, p] \rrbracket_{den} = \llbracket p \rrbracket_{Prot} \rightarrow \mathbf{G}(\llbracket r \rrbracket_{ret})$
- $\llbracket \tau \rrbracket_{exp} = \llbracket \tau \rrbracket_{dat}$
- $\llbracket \textit{boolean} \rrbracket_{exp} = \mathbf{T}$

A few additional domains are necessary, such as the domain \mathbf{Obj} of object descriptors, i.e. object locations in memory, \mathbf{S} of program states and the domain \mathbf{C} of program continuations. Complete definitions are omitted for brevity.

- $o : \mathbf{Obj}$ (domain of object descriptors)
- $s : \mathbf{S}$ (domain of program states)
- ▶ $s_o : \mathbf{S}$ (initial program state)
- ▶ $\textit{write} : \mathbf{Obj} \rightarrow \mathbf{N} \rightarrow \mathbf{S} \rightarrow \mathbf{S}$
- ▶ $\textit{read} : \mathbf{Obj} \rightarrow \mathbf{S} \rightarrow \mathbf{N}$
- \mathbf{A} (domain of program answers)
- $\mathbf{C} = \mathbf{S} \rightarrow \mathbf{A}$ (domain of continuations)
- ▶ $\textit{error}_c : \mathbf{C}$

The domain ordering relation in dynamic semantic domains is usually very important. Bottom elements can model non-termination and top elements, if they exist, can model the occurrence of run-time errors. Intermediate values represent results of computations which produce at least some non-erroneous results. Non-termination and errors must be appropriately propagated by various operations of the dynamic domains and by monad operations.

The dynamic semantics of programming languages can be greatly simplified by proper use of monads and monad transformers. A monad must be defined for each kind of computation that is inherent in the given programming language. In the case of \mathcal{L} , there are two kinds of computations that may be expressed by monads, in addition to monad \mathbf{E} representing compile-time (i.e. constant) computations. *Expression computations* are represented by monad \mathbf{G} , the standard continuation monad; they may affect the program state, i.e. read and

Fig. 6: Typing semantics, rules.

<u>Declarations</u>	
$\frac{e := \{\{body\}\} e_0}{e \vdash body : \mathbf{body}[\mathbf{nothing}]} \quad (B1)$	$\frac{e \vdash declaration\text{-list} : \mathbf{decl} \quad e' := \{\{declaration\text{-list}\}\} e}{e' \vdash statement\text{-list} : \mathbf{stmt}[r]} \quad (B2)$
$\frac{}{e \vdash \text{var } I \text{ as } data\text{-type} : \mathbf{decl}} \quad (D3)$	$\frac{e \vdash formal\text{-parameter}\text{-list} : \mathbf{prot}[p] \quad r := \{\{return\text{-type}\}\} e' := \mathcal{F}\{\{formal\text{-parameter}\text{-list}\}\}(\uparrow e) \quad e' \vdash body : \mathbf{body}[r]}{e \vdash \text{fun } I (formal\text{-parameter}\text{-list}) \text{ as } return\text{-type } body \text{ end} : \mathbf{decl}} \quad (D4)$
$\frac{\tau := \{\{data\text{-type}\}\}}{e \vdash I \text{ as } data\text{-type} \text{ byval} : \mathbf{par}[\tau, \mathbf{byval}]} \quad (P3)$	$\frac{\tau := \{\{data\text{-type}\}\}}{e \vdash I \text{ as } data\text{-type} \text{ byref} : \mathbf{par}[\tau, \mathbf{byref}]} \quad (P4)$
<u>Statements</u>	
$\frac{}{e \vdash \epsilon : \mathbf{stmt}[r]} \quad (S1)$	$\frac{e \vdash statement : \mathbf{stmt}[r] \quad e \vdash statement\text{-list} : \mathbf{stmt}[r]}{e \vdash statement\ statement\text{-list} : \mathbf{stmt}[r]} \quad (S2)$
$\frac{}{e \vdash \text{skip} : \mathbf{stmt}[r]} \quad (S3)$	$\frac{e \vdash I : \mathbf{obj}[\mathbf{func}[\mathbf{nothing}, p]] \quad e \vdash expression\text{-list} : \mathbf{arg}[p]}{e \vdash \text{call } I (expression\text{-list}) : \mathbf{stmt}[r]} \quad (S9)$
$\frac{e \vdash I : \mathbf{obj}[\mathbf{data}[\tau]] \quad e \vdash expression : \mathbf{exp}[\tau]}{e \vdash I := expression : \mathbf{stmt}[r]} \quad (S4)$	$\frac{e \vdash I : \mathbf{obj}[\mathbf{ref}[\tau]] \quad e \vdash expression : \mathbf{exp}[\tau]}{e \vdash I := expression : \mathbf{stmt}[r]} \quad (S5)$
$\frac{e \vdash statement\text{-list} : \mathbf{stmt}[r]}{e \vdash \text{begin } statement\text{-list} \text{ end} : \mathbf{stmt}[r]} \quad (S6)$	$\frac{}{e \vdash \text{return} : \mathbf{stmt}[\mathbf{nothing}]} \quad (S10)$
$\frac{e \vdash condition : \mathbf{exp}[\mathbf{boolean}] \quad e \vdash statement : \mathbf{stmt}[r]}{e \vdash \text{while } condition \text{ do } statement : \mathbf{stmt}[r]} \quad (S7)$	$\frac{e \vdash expression : \mathbf{exp}[\tau]}{e \vdash \text{return } expression : \mathbf{stmt}[\tau]} \quad (S11)$
$\frac{e \vdash condition : \mathbf{exp}[\mathbf{boolean}] \quad e \vdash statement_1 : \mathbf{stmt}[r] \quad e \vdash statement_2 : \mathbf{stmt}[r]}{e \vdash \text{if } condition \text{ then } statement_1 \text{ else } statement_2 : \mathbf{stmt}[r]} \quad (S8)$	
<u>Expressions</u>	
$\frac{}{e \vdash n : \mathbf{exp}[\mathbf{integer}]} \quad (E1)$	$\frac{e \vdash I \triangleleft \phi}{e \vdash I : \mathbf{obj}[\phi]} \quad (E2)$
$\frac{}{e \vdash expression : \mathbf{obj}[\mathbf{data}[\tau]]} \quad (I1)$	$\frac{e \vdash expression : \mathbf{obj}[\mathbf{data}[\tau]]}{e \vdash expression : \mathbf{exp}[\tau]} \quad (I1)$
$\frac{e \vdash expression : \mathbf{obj}[\mathbf{ref}[\tau]]}{e \vdash expression : \mathbf{exp}[\tau]} \quad (I2)$	$\frac{e \vdash I : \mathbf{obj}[\mathbf{func}[\tau, p]] \quad e \vdash expression\text{-list} : \mathbf{arg}[p]}{e \vdash \text{call } I (expression\text{-list}) : \mathbf{exp}[\tau]} \quad (E3)$
$\frac{e \vdash expression : \mathbf{exp}[\mathbf{integer}]}{e \vdash \text{un-op } expression : \mathbf{exp}[\mathbf{integer}]} \quad (E4)$	$\frac{e \vdash expression_1 : \mathbf{exp}[\mathbf{integer}] \quad e \vdash expression_2 : \mathbf{exp}[\mathbf{integer}]}{e \vdash expression_1 \text{ bin-op } expression_2 : \mathbf{exp}[\mathbf{integer}]} \quad (E5)$
$\frac{}{e \vdash \epsilon : \mathbf{arg}[p_0]} \quad (A1)$	$\frac{e \vdash expression : \mathbf{exp}[\tau] \quad e \vdash expression\text{-list} : \mathbf{arg}[p'] \quad p := \langle \tau, \mathbf{byval} \rangle \triangleleft p'}{e \vdash expression\ expression\text{-list} : \mathbf{arg}[p]} \quad (A2)$
$\frac{e \vdash expression : \mathbf{obj}[\mathbf{data}[\tau]] \quad e \vdash expression\text{-list} : \mathbf{arg}[p'] \quad p := \langle \tau, \mathbf{byref} \rangle \triangleleft p'}{e \vdash expression\ expression\text{-list} : \mathbf{arg}[p]} \quad (A3)$	$\frac{e \vdash expression : \mathbf{obj}[\mathbf{ref}[\tau]] \quad e \vdash expression\text{-list} : \mathbf{arg}[p'] \quad p := \langle \tau, \mathbf{byref} \rangle \triangleleft p'}{e \vdash expression\ expression\text{-list} : \mathbf{arg}[p]} \quad (A4)$

alter the values of variables. However, they may not terminate the execution of the program or the current function.

- $G(D) = (D \rightarrow C) \rightarrow C$
- ▶ $unit_G : D \rightarrow G(D)$
 $unit_G = \lambda v. \lambda \kappa. \kappa v$
- ▶ $\cdot *_G \cdot : G(A) \times (A \rightarrow G(B)) \rightarrow G(B)$
 $m *_G f = \lambda \kappa. m(\lambda v. f v \kappa)$
- ▶ $error_G : G(D)$
 $error_G = \lambda \kappa. error_C$
- ▶ $lift_{E \rightarrow G} : E(D) \rightarrow G(D)$
 $lift_{E \rightarrow G} = \lambda m. \text{case } m \text{ of}$
 $inl v \Rightarrow unit_G v$
 $inr u \Rightarrow error_G$
- ▶ $setState : (S \rightarrow S) \rightarrow G(S)$
 $setState = \lambda f. \lambda \kappa. \lambda s. \kappa s (f s)$
- ▶ $getState : G(S)$
 $getState = setState id$

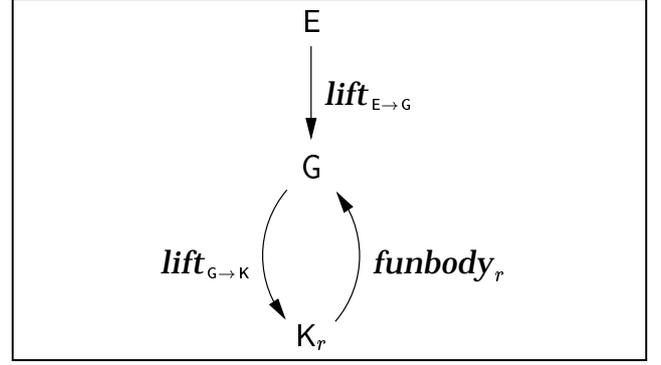
On the other hand, *statement* computations are represented by monad K_r ; they may affect the program state but may also terminate the program or the current function by returning a result of type r .

- $K_r(D) = ([r]_{ret} \rightarrow C) \rightarrow G(D)$
- ▶ $unit_K : D \rightarrow K_r(D)$
 $unit_K = \lambda v. \lambda \kappa_r. unit_G v$
- ▶ $\cdot *_K \cdot : K_r(A) \times (A \rightarrow K_r(B)) \rightarrow K_r(B)$
 $m *_K f = \lambda \kappa_r. m \kappa_r *_G (\lambda v. f v \kappa_r)$
- ▶ $error_K : K_r(D)$
 $error_K = \lambda \kappa_r. error_G$
- ▶ $lift_{G \rightarrow K} : G(D) \rightarrow K_r(D)$
 $lift_{G \rightarrow K} = \lambda m. \lambda \kappa_r. m$
- ▶ $funbody_r : K_r(\mathbf{U}) \rightarrow G([r]_{ret})$
 $funbody_{nothing} = \lambda k. \lambda \kappa. k \kappa \kappa$
 $funbody_\tau = \lambda k. \lambda \kappa. k \kappa (\lambda u. error_C)$
- ▶ $result_r : [r]_{ret} \rightarrow K_r(\mathbf{U})$
 $result_r = \lambda d. \lambda \kappa_r. \lambda \kappa. \kappa_r d$

The relation between the three monads and the conversion functions are shown in Fig. 7.

Dynamic environments are a special case of dynamic semantic domains. They correspond directly to static environments and associate static program elements, such as identifiers or formal parameters, with the dynamic entities that these elements represent during the program's execution. In the case of \mathcal{L} , two kinds of dynamic environments are needed: the domain $\llbracket e \rrbracket_{Env}$ of dynamic type environments corresponding to the static type environment e , and the domain $\llbracket p \rrbracket_{Prot}$ of dynamic function prototypes corresponding to the static function prototype p . They are

Fig. 7: Relations between the defined monads.



both defined below as abstract data types. In brief, the first associates variables with objects in memory, whereas the second associates formal with actual parameters.¹

- $\rho : \llbracket e \rrbracket_{Env}$ (domain of dynamic type environments)
- $\llbracket e \Downarrow I \rrbracket = \begin{cases} \llbracket \phi \rrbracket_{den} & , \text{ if } e[I] = unit_E \phi \\ \mathbf{U} & , \text{ if } e[I] = error_E \end{cases}$
- ▶ $\rho_o : \llbracket e_o \rrbracket_{Env}$ (empty dynamic environment)
- ▶ $\cdot \uparrow \cdot : e : \mathbf{Env} \times \llbracket \downarrow e \rrbracket_{Env} \rightarrow \llbracket e \rrbracket_{Env}$
- ▶ $lookup : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{Env} \rightarrow \mathbf{Ide} \rightarrow E(\llbracket e \Downarrow I \rrbracket)$
- ▶ $create : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{Env} \rightarrow \mathbf{Ide} \rightarrow E(\llbracket e \rrbracket_{Env})$
- ▶ $assign : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{Env} \rightarrow I : \mathbf{Ide} \rightarrow \llbracket e \Downarrow I \rrbracket \rightarrow E(\llbracket e \rrbracket_{Env})$
- $\llbracket p \rrbracket_{Prot}$ (domain of dynamic function prototypes)
- $\llbracket p \Downarrow i \rrbracket = \begin{cases} \llbracket \tau \rrbracket_{dat} & , \text{ if } p[i] = unit_E \langle \tau, byval \rangle \\ \llbracket ref[\tau] \rrbracket_{den} & , \text{ if } p[i] = unit_E \langle \tau, byref \rangle \\ \mathbf{U} & , \text{ if } p[i] = error_E \end{cases}$
- $\llbracket \tau, byval \rrbracket_{Par} = \llbracket \tau \rrbracket_{dat}$
- $\llbracket \tau, byref \rrbracket_{Par} = \mathbf{Obj}$
- ▶ $noargs : \llbracket p_o \rrbracket_{Prot}$
- ▶ $prepend_{p, \tau, m, p'} : \llbracket \tau, m \rrbracket_{Par} \rightarrow \llbracket p' \rrbracket_{Prot} \rightarrow E(\llbracket p \rrbracket_{Prot})$
- ▶ $split_{p, \tau, m, p'} : \llbracket p \rrbracket_{Prot} \rightarrow E(\llbracket \tau, m \rrbracket_{Par} \times \llbracket p' \rrbracket_{Prot})$

During the second step in the definition of dynamic semantics, dynamic functions and equations must be defined. In general, a dynamic semantic function maps well-typed program phrases to dynamic semantic domains. A program phrase is well-typed if there exists a typing derivation for it, and the conclusion of this typing derivation specifies the phrase type. At least one

¹In the notation that we use, the domains $x : A \multimap B(x)$ and $x : A \times B(x)$ are respectively dependent function and dependent product domains.

Fig. 8: Dynamic semantics for \mathcal{L} , functions and equations for declarations.

▶	$\llbracket \text{prog} \rrbracket : \mathbf{K}_{\text{nothing}}(\mathbf{U})$
(B1)	$\llbracket \text{body} \rrbracket_{\text{prog}} = \text{lift}_{\mathbf{E} \rightarrow \mathbf{K}} (\{\text{body}\} e_o) * (\lambda e. \llbracket \text{body} \rrbracket_{\text{body} [r]} e (e \uparrow \rho_o))$
▶	$\llbracket \text{body} [r] \rrbracket : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{K}_r(\mathbf{U})$
(B2)	$\llbracket \text{declaration-list statement-list} \rrbracket_{\text{body} [r]} = \lambda e. \lambda \rho.$ $\text{lift}_{\mathbf{E} \rightarrow \mathbf{K}} (\{\text{declaration-list}\} e) * (\lambda e'.$ $\text{lift}_{\mathbf{E} \rightarrow \mathbf{K}} (\mathbf{mcl}_o (\llbracket \text{declaration-list} \rrbracket_{\text{decl}} e' \rho) (\mathcal{F} \llbracket \text{declaration-list} \rrbracket_{\text{decl}} e')) * (\lambda \rho'.$ $\llbracket \text{statement-list} \rrbracket_{\text{stmt} [r]} e' \rho'))$
▶	$\llbracket \text{decl} \rrbracket : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{E}(\llbracket e \rrbracket_{\mathbf{Env}})$
(D3)	$\llbracket \text{var } I \text{ as data-type} \rrbracket_{\text{decl}} = \lambda e. \lambda \rho. \mathbf{create} e \rho I$
(D4)	$\llbracket \text{fun } I (\text{formal-parameter-list}) \text{ as return-type body end} \rrbracket_{\text{decl}} = \lambda e. \mathbf{unit}$
▶	$\mathcal{F} \llbracket \text{decl} \rrbracket : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{E}(\llbracket e \rrbracket_{\mathbf{Env}})$
(D3)	$\mathcal{F} \llbracket \text{var } I \text{ as data-type} \rrbracket_{\text{decl}} = \lambda e. \mathbf{unit}$
(D4)	$\mathcal{F} \llbracket \text{fun } I (\text{formal-parameter-list}) \text{ as return-type body end} \rrbracket_{\text{decl}} = \lambda e. \lambda \rho.$ $\mathcal{F} \{\text{formal-parameter-list}\} (\uparrow e) * (\lambda e'. \{\text{return-type}\} * (\lambda r.$ $\mathbf{let} f = \lambda d_p. \llbracket \text{formal-parameter-list} \rrbracket_{\text{prot} [p]} d_p e' (e' \uparrow \rho) * (\lambda \rho'. \mathbf{funbody}_r (\llbracket \text{body} \rrbracket_{\text{body} [r]} e' \rho'))$ $\mathbf{in} \mathbf{assign} e \rho I f))$
▶	$\llbracket \text{prot} [p] \rrbracket : \llbracket p \rrbracket_{\text{prot}} \rightarrow e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{G}(\llbracket e \rrbracket_{\mathbf{Env}})$
▶	$\llbracket \text{par} [\tau, m] \rrbracket : \llbracket \tau, m \rrbracket_{\text{par}} \rightarrow e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{G}(\llbracket e \rrbracket_{\mathbf{Env}})$
(P3)	$\llbracket I \text{ as data-type byval} \rrbracket_{\text{par} [\tau, \text{byval}]} = \lambda d_a. \lambda e. \lambda \rho.$ $\text{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\mathbf{create} e \rho I) * (\lambda \rho'.$ $\text{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\mathbf{lookup} e \rho' I) * (\lambda o.$ $\mathbf{setState} (\mathbf{write} o d_a) * (\lambda s.$ $\mathbf{unit} \rho'))$
(P4)	$\llbracket I \text{ as data-type byref} \rrbracket_{\text{par} [\tau, \text{byref}]} = \lambda d_a. \lambda e. \lambda \rho.$ $\text{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\mathbf{assign} e \rho I d_a)$

function is required for each phrase type, representing the dynamic meaning of phrases that can be attributed this type. More than one functions may be necessary if there are more than one possible meanings for a program phrase. Subsequently, one equation is required for each typing rule whose conclusion is a main typing relation. A large subset of the dynamic semantic functions and equations for \mathcal{L} , corresponding to the typing rules given in Fig. 6 is given in Fig. 8, 9 and 10.

8 Implementation

Easily produced and verifiably correct implementations, directly based on formal definitions, are valuable tools for the semantic analysis, design and evaluation of programming languages. Although not in the spirit of operational semantics, a denotational description of a programming language defines directly an execution model in the form of the abstract interpreter discussed in Section 2. This execution model can be implemented by translating the denotational semantics

to a program, written in some target language.

Several languages have been suggested and used for this purpose, with more or less success. It seems that typed functional programming languages, such as ML or Haskell, are more suitable for this purpose. The latter is a better choice for implementing the semantics obtained by the proposed methodology, mainly because of its richer type system, more flexible syntax, elegant support for monads and also because lazy evaluation avoids a number of non-termination problems.

A direct implementation of the definition of \mathcal{L} using Haskell as the target language consists of approximately 700 lines of code: 15 lines were needed for the implementation of the abstract syntax, 120 for the static semantics, 200 for the typing semantics and 275 for the dynamic semantics. The translation from the source programs in \mathcal{L} to the abstract syntax was implemented as 700 additional lines of Flex/Bison code²

²For the complete implementation the reader is referred to [9].

Fig. 9: Dynamic semantics for \mathcal{L} , functions and equations for statements.

►	$\llbracket \text{stmt } [r] \rrbracket : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{K}_r(\mathbf{U})$
(S3)	$\llbracket \text{skip} \rrbracket_{\text{stmt } [r]} = \lambda e. \lambda \rho. \mathbf{unit } u$
(S4, 5)	$\llbracket I := \text{expression} \rrbracket_{\text{stmt } [r]} = \lambda e. \lambda \rho. \llbracket I \rrbracket_{\text{obj } [\phi]} e \rho * (\lambda o. \mathbf{lift}_{\mathbf{G} \rightarrow \mathbf{K}} (\llbracket \text{expression} \rrbracket_{\text{exp } [\tau]} e \rho) * (\lambda v. \mathbf{lift}_{\mathbf{G} \rightarrow \mathbf{K}} (\mathbf{setState } (\mathbf{write } o v)) * (\lambda s. \mathbf{unit } u)))$
(S6)	$\llbracket \text{begin statement-list end} \rrbracket_{\text{stmt } [r]} = \llbracket \text{statement-list} \rrbracket_{\text{stmt } [r]}$
(S7)	$\llbracket \text{if condition then statement}_1 \text{ else statement}_2 \rrbracket_{\text{stmt } [r]} = \lambda e. \lambda \rho. \mathbf{lift}_{\mathbf{G} \rightarrow \mathbf{K}} (\llbracket \text{condition} \rrbracket_{\text{exp } [\text{boolean}]} e \rho) * (\lambda v. \mathbf{if } v \text{ then } \llbracket \text{statement}_1 \rrbracket_{\text{exp } [\text{stmt } [r]]} e \rho \text{ else } \llbracket \text{statement}_2 \rrbracket_{\text{exp } [\text{stmt } [r]]} e \rho)$
(S8)	$\llbracket \text{while condition do statement} \rrbracket_{\text{stmt } [r]} = \lambda e. \lambda \rho. \mathbf{mfix } (\lambda k. \mathbf{lift}_{\mathbf{G} \rightarrow \mathbf{K}} (\llbracket \text{condition} \rrbracket_{\text{exp } [\text{boolean}]} e \rho) * (\lambda v. \mathbf{if } v \text{ then } \llbracket \text{statement} \rrbracket_{\text{exp } [\text{stmt } [r]]} e \rho * (\lambda u. k) \text{ else } \mathbf{unit } u))$
(S9)	$\llbracket \text{call } I (\text{expression-list}) \rrbracket_{\text{stmt } [r]} = \lambda e. \lambda \rho. \llbracket I \rrbracket_{\text{obj } [\text{func } [\text{nothing}, p]]} e \rho * (\lambda f. \mathbf{lift}_{\mathbf{G} \rightarrow \mathbf{K}} (\llbracket \text{expression-list} \rrbracket_{\text{arg } [p]} e \rho) * (\lambda d_p. \mathbf{lift}_{\mathbf{G} \rightarrow \mathbf{K}} (f d_p)))$
(S10)	$\llbracket \text{return} \rrbracket_{\text{stmt } [r]} = \lambda e. \lambda \rho. \mathbf{result } u$
(S11)	$\llbracket \text{return expression} \rrbracket_{\text{stmt } [r]} = \lambda e. \lambda \rho. \mathbf{lift}_{\mathbf{G} \rightarrow \mathbf{K}} (\llbracket \text{expression} \rrbracket_{\text{exp } [\tau]} e \rho) * (\lambda v. \mathbf{result } v)$

Fig. 10: Dynamic semantics for \mathcal{L} , functions and equations for expressions.

►	$\llbracket \text{exp } [v] \rrbracket : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{G}(\llbracket v \rrbracket_{\text{exp}})$
►	$\llbracket \text{obj } [\phi] \rrbracket : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{G}(\llbracket \phi \rrbracket_{\text{den}})$
(E1)	$\llbracket n \rrbracket_{\text{exp } [\text{integer}]} = \lambda e. \lambda \rho. \mathbf{unit } \mathcal{N}[n]$
(E2)	$\llbracket I \rrbracket_{\text{obj } [\phi]} = \lambda e. \lambda \rho. \mathbf{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\mathbf{lookup } e \rho I)$
(E3)	$\llbracket \text{call } I (\text{expression-list}) \rrbracket_{\text{exp } [\tau]} = \lambda e. \lambda \rho. \llbracket I \rrbracket_{\text{obj } [\text{func } [\tau, r]]} e \rho * (\lambda f. \llbracket \text{expression-list} \rrbracket_{\text{arg } [p]} e \rho * (\lambda d_p. f d_p))$
(E4)	$\llbracket \text{un-op expression} \rrbracket_{\text{exp } [\text{integer}]} = \lambda e. \lambda \rho. \llbracket \text{expression} \rrbracket_{\text{exp } [\text{integer}]} e \rho * (\lambda n. \mathbf{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\mathcal{UO}[\text{un-op}] n))$
(E5)	$\llbracket \text{expression}_1 \text{ bin-op expression}_2 \rrbracket_{\text{exp } [\text{integer}]} = \lambda e. \lambda \rho. \llbracket \text{expression}_1 \rrbracket_{\text{exp } [\text{integer}]} e \rho * (\lambda n_1. \llbracket \text{expression}_2 \rrbracket_{\text{exp } [\text{integer}]} e \rho * (\lambda n_2. \mathbf{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\mathcal{BO}[\text{bin-op}] \langle n_1, n_2 \rangle)))$
(II, 2)	$\llbracket \text{expression} \rrbracket_{\text{exp } [\tau]} = \lambda e. \lambda \rho. \llbracket \text{expression} \rrbracket_{\text{obj } [\phi]} e \rho * (\lambda \phi. \mathbf{getState } * (\lambda s. \mathbf{unit } (\mathbf{read } s o)))$
►	$\llbracket \text{arg } [p] \rrbracket : e : \mathbf{Env} \rightarrow \llbracket e \rrbracket_{\mathbf{Env}} \rightarrow \mathbf{G}(\llbracket p \rrbracket_{\text{Prot}})$
(A1)	$\llbracket \epsilon \rrbracket_{\text{arg } [p_o]} = \lambda e. \lambda \rho. \mathbf{unit } \mathbf{noargs}$
(A2)	$\llbracket \text{expression expression-list} \rrbracket_{\text{arg } [p]} = \lambda e. \lambda \rho. \mathbf{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\Downarrow p) * (\lambda \langle \tau, \text{byval}, p' \rangle. \llbracket \text{expression} \rrbracket_{\text{exp } [\tau]} e \rho * (\lambda d. \llbracket \text{expression-list} \rrbracket_{\text{arg } [p']} e \rho * (\lambda d_{p'}. \mathbf{prepend}_{p, \tau, \text{byval}, p'} d d_{p'})))$
(A3, 4)	$\llbracket \text{expression expression-list} \rrbracket_{\text{arg } [p]} = \lambda e. \lambda \rho. \mathbf{lift}_{\mathbf{E} \rightarrow \mathbf{G}} (\Downarrow p) * (\lambda \langle \tau, \text{byref}, p' \rangle. \llbracket \text{expression} \rrbracket_{\text{obj } [\phi]} e \rho * (\lambda d. \llbracket \text{expression-list} \rrbracket_{\text{arg } [p']} e \rho * (\lambda d_{p'}. \mathbf{prepend}_{p, \tau, \text{byref}, p'} d d_{p'})))$
►	$\mathcal{UO}[\text{un-op}] : \mathbf{N} \rightarrow \mathbf{E}(\mathbf{N})$
►	$\mathcal{BO}[\text{bin-op}] : \mathbf{N} \times \mathbf{N} \rightarrow \mathbf{E}(\mathbf{N})$

9 Evaluation

The proposed methodology has been used for the formal definition of the ANSI C programming language, as described in [4]. The definition was based on the standard for the language and is as complete and accurate as possible. The definition of abstract syntax consists of 39 production rules, whereas the definition of static semantics involves 17 domains, one monad, 33 semantic functions and 93 semantic equations. Furthermore, approximately 200 typing rules are used for the typing semantics. The dynamic semantics is the most complicated part, because of C's inherent complexity: 41 domains are used, 6 monads, one monad

transformer, 45 semantic functions and more than 250 dynamic equations.

A structured implementation of the semantics of ANSI C in Haskell has also been developed. It 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. The implementation was used for the evaluation of the completeness and accuracy of the developed semantics, with very positive results.

10 Conclusion and future work

In this paper, we have presented a methodology for the formal definition of programming languages, following the denotational approach and using monads to improve modularity and elegance of notation. The main contribution of the proposed methodology are the three distinct consecutive phases that it proposes, for the definition of the semantics. The obtained programming language definitions are well structured and modular and can be implemented in the form of abstract interpreters in a relatively straightforward way.

Our research in the near future will focus on the process of evaluating and improving the proposed methodology, by applying it to programming languages of different nature and complexity. A secondary aim is to study the practical applications that formal definitions of programming languages may have in the software industry, especially in tools for program transformation, debugging and understanding, and how our methodology can be enriched to support these applications.

References:

- [1] R. D. Tennent, "The denotational semantics of programming languages," *Communications of the ACM*, vol. 19, pp. 437–453, Aug. 1976.
- [2] P. D. Mosses, "Denotational semantics," in *Handbook of Theoretical Computer Science* (J. van Leeuwen, ed.), vol. B, ch. 11, pp. 577–631, Elsevier Science Publishers B.V., 1990.
- [3] R. D. Tennent, *Semantics of Programming Languages*. Englewood Cliffs, NJ: Prentice Hall, 1991.
- [4] N. S. Papaspyrou, *A Formal Semantics for the C Programming Language*. PhD thesis, National Technical University of Athens, Software Engineering Laboratory, Feb. 1998.
- [5] E. Moggi, "Computational lambda calculus and monads," in *IEEE Symposium on Logic in Computer Science*, pp. 14–23, 1989.
- [6] E. Moggi, "An abstract view of programming languages," Tech. Rep. ECS-LFCS-90-113, University of Edinburgh, Laboratory for Foundations of Computer Science, 1990.
- [7] P. Wadler, "The essence of functional programming," in *Proceedings of the 19th Annual Symposium on Principles of Programming Languages (POPL'92)*, Jan. 1992.
- [8] 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), Jan. 1995.
- [9] N. Papaspyrou, "A methodology for the definition of programming languages," Tech. Rep. CSD-SW-TR-1-99, National Technical University of Athens, Software Engineering Laboratory, Apr. 1999.