

Introduction: Languages, Semantics, Interpreters, Compilers

Dmitry Boulytchev

1 Language and semantics

A language is a collection of programs. A program is an abstract syntax tree (AST), which describes the hierarchy of constructs. An abstract syntax of a programming language describes the format of abstract syntax trees of programs in this language. Thus, a language is a set of constructive objects, each of which can be constructively manipulated.

The semantics of a language \mathcal{L} is a total map

$$\llbracket \bullet \rrbracket_{\mathcal{L}} : \mathcal{L} \rightarrow \mathcal{D}$$

where \mathcal{D} is some semantic domain. The choice of the domain is at our command; for example, for Turing-complete languages \mathcal{D} can be the set of all partially-recursive (computable) functions.

2 Interpreters

In reality, the semantics often is described using interpreters:

$$eval : \mathcal{L} \rightarrow \text{Input} \rightarrow \text{Output}$$

where Input and Output are sets of (all possible) inputs and outputs for the programs in the language \mathcal{L} . We claim *eval* to possess the following property

$$\forall p \in \mathcal{L}, \forall x \in \text{Input} : \llbracket p \rrbracket_{\mathcal{L}} x = eval\ p\ x$$

In other words, an interpreter takes a program and its input as arguments, and returns what the program would return, being run on that argument. The equality in the definitional property of an interpreter has to be read “if the right hand side is defined, then the left hand side is defined, too, and their values coincide”, and vice-versa.

Why interpreters are so important? Because they can be written as programs in a meta-language, or a language of implementation. For example, if we take *ocaml* as a language of implementation, then an interpreter of a language \mathcal{L} is some *ocaml* program *eval*, such that

$$\forall p \in \mathcal{L}, \forall x \in \text{Input} : \llbracket p \rrbracket_{\mathcal{L}} x = \llbracket eval \rrbracket_{\text{ocaml}} p\ x$$

How to define $\llbracket \bullet \rrbracket_{\text{ocaml}}$? We can write an interpreter in some other language. Thus, a tower of meta-languages and interpreters comes into consideration. When to stop? When the meta-language is simple enough for intuitive understanding (in reality: some math-based frameworks like operational, denotational or game semantics, etc.)

Pragmatically: if you have a good implementation of a good programming language you trust, you can write interpreters of other languages.

3 Compilers

A compiler is just a language transformer

$$\text{comp} : \mathcal{L} \rightarrow \mathcal{M}$$

for two languages \mathcal{L} and \mathcal{M} ; we expect a compiler to be total and to possess the following property:

$$\forall p \in \mathcal{L} \quad \llbracket p \rrbracket_{\mathcal{L}} = \llbracket \text{comp } p \rrbracket_{\mathcal{M}}$$

Again, the equality in this definition is understood functionally. The property itself is called a complete (or full) correctness. In reality compilers are partially correct, which means, that the domain of compiled programs can be wider.

And, again, we expect compilers to be defined in terms of some implementation language. Thus, a compiler is a program (in, say, ocaml), such, that its semantics in ocaml possesses the following property (fill the rest yourself).

4 The first example: language of expressions

Abstract syntax:

$$\begin{array}{ll} \mathcal{X} &= \{x, y, z, \dots\} & \text{(variables)} \\ \otimes &= \{+, -, *, /, \%, <, <=, >, >=, ==, !=, !!, \&\&\} & \text{(binary operators)} \\ \mathcal{E} &= \mathcal{X} & \text{(expressions)} \\ &\quad \mathbb{N} \\ &\quad \mathcal{E} \otimes \mathcal{E} \end{array}$$

Semantics of expressions:

- state $\sigma : \mathcal{X} \rightarrow \mathbb{Z}$ assigns values to (some) variables;
- semantics $\llbracket \bullet \rrbracket$ assigns each expression a total map $\Sigma \rightarrow \mathbb{Z}$, where Σ is the set of all states.

Empty state \perp : undefined for any variable.
Denotational style of semantic description:

$$\begin{aligned} \llbracket n \rrbracket &= \lambda\sigma. n & , \quad n \in \mathbb{N} \\ \llbracket x \rrbracket &= \lambda\sigma. \sigma x & , \quad x \in \mathcal{X} \\ \llbracket A \otimes B \rrbracket &= \lambda\sigma. (\llbracket A \rrbracket \sigma \oplus \llbracket B \rrbracket \sigma) & , \quad A, B \in \mathcal{E} \end{aligned}$$

\otimes	\oplus in ocaml	
+	+	
-	-	
*	*	
/	/	
%	mod	
<	<	} see note 1 below
>	>	
<=	<=	
>=	>=	
==	=	
!=	<>	} see note 2 below
&&	&&	
!!		

Note 1: the result is converted into integers (true \rightarrow 1, false \rightarrow 0).

Note 2: the arguments are converted to booleans (0 \rightarrow false, not 0 \rightarrow true),
the result is converted to integers as in the previous note.

Important observations:

1. $\llbracket \bullet \rrbracket$ is defined compositionally: the meaning of an expression is defined in terms of meanings of its proper subexpressions. This is an important property of denotational style.
2. $\llbracket \bullet \rrbracket$ is total, since it takes into account all possible ways to deconstruct any expression.
3. $\llbracket \bullet \rrbracket$ is deterministic: there is no way to assign different meanings to the same expression, since we deconstruct each expression unambiguously.
4. \otimes is an element of language syntax, while \oplus is its interpretation in the meta-language of semantic description (simpler: in the language of interpreter implementation).
5. This concrete semantics is strict: for a binary operator both its arguments are evaluated unconditionally; thus, for example, $1 \text{ !! } x$ is undefined in empty state.