

Language Engineering - A nice set of notes

Josh Felmeden

November 15, 2019

Contents

1	Introduction to Semantics	3
2	Structural Operational Semantics	4
2.1	Termination and looping	4
2.2	Determinism and Equivalence	5
2.3	Provably correct implementation	6
2.3.1	Arithmetic code	7
2.3.2	State changing code	7
2.3.3	Computation sequences	7
2.4	The execution function	8
2.5	Code Translation of Expressions	8
2.6	Semantic function	9
2.6.1	Correctness of translation	9
3	Denotational Semantics	10
3.1	Conditional functions	10
3.2	Least fixpoint	11
4	Chain-Complete Partial Orders	11
4.1	Definitions of PO-Set (partially ordered set)	11
5	The language Exc	12
5.1	Continuation-style denotational semantics	13

1 Introduction to Semantics

Semantics are really complex and they actually exist in the real world as problems that can arise when the semantics are unclear. In the example of the Derek Bentley case, Bentley tells Chris (who is holding a gun, and a policeman standing in front of him to 'let him have it!'. Here, it appears that he could be talking about the gun, or to kill him. The same kind of thing can happen in computing when we are unsure of the references of certain objects.

Here are some examples learned from natural languages:

- Syntactic complexity
 - Jack built the house the malt the rat the cat killed ate lay in
- Syntactic ambiguity
 - Let him have it, Chris!
- Semantic Complexity
 - It depends on what the meaning of the word 'is' is!
- Semantic ambiguity
 - I haven't slept for ten days
- Semantic undefinedness
 - Colourless green ideas sleep furiously
- Interaction of syntax and semantics
 - Time flies like an arrow, fruit flies like a banana.

We can apply these things to computing terms, too.

- Syntactic complexity

```
x-=y = (x=x+y) - y      //switches variables x and y
```

- Syntactic ambiguity

```
if (...) if (...) ..; else ..      //dangling else
```

- Semantic Complexity

```
y = x++ + x++      //sequence points
```

- Semantic ambiguity

```
(x%2=1) ? "odd" : "even"      //unspecified in C89 if x<0
```

- Semantic undefinedness

```
while(x/x)    //division error or infinite loop
```

- Interaction of syntax and semantics

```
A * B    //lever hack
```

To put this another way:

- **Syntax:** concerned with the form of expressions and whether or not the program actually *compiles*
- **Semantics:** concerned with the meaning of expressions and what the program does when it *runs*
- **Pragmatics:** concerned with issues like design patterns, program style, industry standards, etc.

2 Structural Operational Semantics

We're going to look at doing some compilation (of the *while*) language.

2.1 Termination and looping

The execution of the statement S in state σ terminates iff there exists a finite derivation sequence from $\langle S, \sigma \rangle$. The derivation sequence looks like:

$$\langle S, \theta \rangle \Rightarrow \gamma_1 \Rightarrow \dots \Rightarrow \gamma_n \text{ where } \gamma_n \text{ is terminal } \sigma' \text{ or stuck } \langle S', \sigma' \rangle$$

The while language never gets stuck, but some language might if we try to divide by zero because we don't know how to process this.

The execution of the statement S in a state θ loops iff there exists an infinite derivation sequence from $\langle S, \sigma \rangle$

$$\langle S, \sigma \rangle \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots$$

S always terminates iff its execution terminates in all states σ .

S always loops if the execution loops in all states σ .

The execution of statement S in state σ terminates successfully iff it ends with a terminal configuration.

Note while has no stuck configurations, so termination implies successful termination!

2.2 Determinism and Equivalence

The structural operation semantics is (strongly) **deterministic** iff $\langle S, \sigma \rangle \Rightarrow \gamma$ and $\langle S, \sigma \rangle \Rightarrow \gamma'$ imply that $\gamma = \gamma'$ for all $S, \sigma, \gamma, \gamma'$

It is **weakly deterministic** iff $\langle S, \sigma \rangle \Rightarrow^* \sigma'$ and $\langle S, \sigma \rangle \Rightarrow^* \sigma''$ imply that $\sigma' = \sigma''$ for all $S, \sigma, \sigma', \sigma''$. This is different from the strong determinism above because it says that for every successfully terminating branch, (it doesn't matter how we get there) we get to the same final state.

Two statements are **semantically equivalent** whenever it holds that for *all states* σ

$$\langle S_1, \sigma \rangle \Rightarrow^* \gamma \text{ iff } \langle S_2, \sigma \rangle \Rightarrow^* \gamma \text{ whenever } \gamma \text{ is terminal or stuck}$$

This means that there is an infinite derivation sequence from $\langle S_1, \sigma \rangle$ iff there is an infinite derivation from $\langle S_2, \sigma \rangle$.

Note! The length of these could be different (because of the * again.)

For a deterministic structural operational semantics, we can define a semantic function as follows:

- $S_{sos}[[\cdot]] : \text{Stm} \rightarrow (\text{State} \hookrightarrow \text{State})$
- $S_{sos}[[S]]\sigma = \sigma'$ if $\langle S, \sigma \rangle \Rightarrow^* \sigma'$ and **undefined** otherwise
- Note that the semantic function is only guaranteed to return a partial function between states due to the existence of statements whose execution loops in one or more states
- $S_{sos}[[\text{while true do skip}]] = \{\}$
 - If we apply this on any state, we get undefined back BUT it is not in and of itself undefined. It is simply the empty set.
 - What could we do if the semantics is not deterministic?
 - * The problem is that depending on what choice we made, we might get a different answer. But the definition says that we only return one function. So therefore, we need to be able to collect them up into a list.
 - * One way of doing it is: $S'_{sos}[[\cdot]] : \text{Stm} \hookrightarrow (\text{State} \hookrightarrow \text{State})$ to ignore the ambiguous cases
 - * Another way is to allow a set of final states: $S''_{sos}[[\cdot]] : \text{Stm} \rightarrow (\text{State} \rightarrow \mathcal{P} \text{State})$. This is bad because we get a set of states.
 - * $S'''_{sos}[[\cdot]] : \text{Stm} \rightarrow (\mathcal{P}(\text{State}) \rightarrow \mathcal{P}(\text{State}))$ now facilitates function composition. Basically, we pass a load of states, and the function returns a list of all functions that can be reached from any of those functions.

Theorem. For all statements S of **While**, it holds that $S_{ns}[[S]] = S_{sos}[[S]]$. Basically, for all statements, then:

$$\{(\sigma, \sigma') \in \text{State}^2 \mid \langle S, \sigma \rangle \rightarrow \sigma'\} = \{(\sigma, \sigma') \in \text{State}^2 \mid \langle S, \sigma \rangle \Rightarrow^* \sigma'\}$$

This mess can be decomposed into two different facts:

$$\langle S, \sigma \rangle \Rightarrow^* \sigma' \text{ implies } \langle S, \sigma \rangle \rightarrow \sigma'$$

And

$$\langle S, \sigma \rangle \rightarrow \sigma' \text{ implies } \langle S, \sigma \rangle \Rightarrow^* \sigma'$$

Very subtle, right? This can also be decomposed further into some cool stuff but I don't think it's helpful. See lemma 2.28 in the book for the derivation sequence.

2.3 Provably correct implementation

We're now going to look at the correctness of a compiler from **While** to an abstract machine **AM**. Initially, we will consider a simple **stack** machine with a set of abstract instructions. Later on, we'll refine it to use memory addresses. Let's formalise some aspects of the abstract machine.

The configurations in the machine are going to be a triple: $\langle c, e, s \rangle$:

- c is the code to be executed $c \in \text{Code} = \text{inst}^*$
- e is the evaluation stack (of expressions) $e \in \text{Stack} = (Z \cup T)^*$
- s is the storage (for variables) $s \in \text{State} = \text{Var} \rightarrow Z$

The instructions will be:

```
inst ::= push-n | add | mult | sub
      | true | false | eq | le | and | neg
      | fetch-x | store-x | noop | branch(c,c) | loop(c,c)

c ::= empty | inst:c
```

noop is basically a skip. Also, we'll be passing around code in this example, but later on we'll replace the 'code' by memory addresses where the code is stored.

2.3.1 Arithmetic code

At this point, we might have the following code:

$\langle \text{PUSH-n}:c, e, s \rangle$	\triangleright	$\langle c, \mathcal{N}[[n]]:e, s \rangle$	
$\langle \text{ADD}:c, z_1:z_2:e, s \rangle$	\triangleright	$\langle c, (z_1 * z_2):e, s \rangle$	if $z_1, z_2 \in \mathbb{Z}$
$\langle \text{TRUE}:c, e, s \rangle$	\triangleright	$\langle c, \mathbf{tt}:e, s \rangle$	
$\langle \text{EQ}:c, z_1:z_2:e, s \rangle$	\triangleright	$\langle c, (z_1 = z_2):e, s \rangle$	if $z_1, z_2 \in \mathbb{Z}$

Here, the ‘.’ is much like the ‘cons’ function from Haskell, in that if we take ADD for example; $\text{ADD}:c$ means that we have the statement ADD, and then more code following it. In the same way, with the arguments of ADD, we need two integers z_1, z_2 on the stack, represented by $z_1 : z_2 : e$.

Obviously, there are more keywords, but they follow the same format as these existing ones.

2.3.2 State changing code

Now, let's look at some of the state rules:

$\langle \text{FETCH-x}:c, e, s \rangle$	\triangleright	$\langle c, s(sx):e, s \rangle$	
$\langle \text{STORE-x}:c, z:e, s \rangle$	\triangleright	$\langle c, e, s[x \mapsto z] \rangle$	if $z \in \mathbb{Z}$
$\langle \text{NOOP}:c, e, s \rangle$	\triangleright	$\langle c, e, s \rangle$	
$\langle \text{BRANCH}(c_1, c_2):c, t:e, s \rangle$	\triangleright	$\begin{cases} \langle c_1:c, e, s \rangle \\ \langle c_2:c, e, s \rangle \end{cases}$	$\begin{matrix} \text{if } t = \mathbf{tt} \\ \text{if } t = \mathbf{ff} \end{matrix}$
$\langle \text{LOOP}(c_2, c_2):c, e, s \rangle$	\triangleright	$\langle c_1:\text{BRANCH}(c_2:\text{LOOP}(c_1, c_2), \text{NOOP}):c, e, s \rangle$	

2.3.3 Computation sequences

- A configuration γ can have one of two forms. It can either be **incomplete** or **terminal**.
- An incomplete configuration be either **stuck** if there is no γ' such that $\gamma \triangleright \gamma'$, OR it is **unstuck** if the opposite is true.
- A computation sequence from $\langle c, \epsilon, \sigma \rangle$ is either a **finite sequence** such that all γ is a terminal or stuck configuration, or it is **infinite**, such that $\gamma_0 = \langle c, \epsilon, \sigma \rangle$ and $\gamma_i \triangleright \gamma_{i+1}$ for all $0 \leq i$.
- **Note!** $\gamma \triangleright^k \gamma'$ means that γ' can be obtained from γ in exactly k steps of \triangleright .
- **Note!** $\gamma \triangleright^* \gamma'$ means that γ' can be obtained from γ in a *finite* number of steps.

Termination and looping is pretty basic and expected, so I won't cover that here.

2.4 The execution function

We can define an execution function for our abstract machine (AM) as follows:

- $\mathcal{M}[\cdot]: \text{Code} \rightarrow (\text{State} \hookrightarrow \text{State})$
- $\mathcal{M}[c]\sigma =$

$$\begin{cases} \sigma' & \text{if } \langle c, \epsilon, \sigma' \rangle \triangleright^* \langle \epsilon, e, \sigma' \rangle \\ \text{(Undefined)} & \text{otherwise} \end{cases}$$

2.5 Code Translation of Expressions

Now, we're looking at a function that can translate from **while** into this AM code. So, $\mathcal{CA}[\cdot]: \text{Aexp} \rightarrow \text{Code}$.

$$\begin{aligned} \mathcal{CA}[n] &= \text{PUSH-}n \\ \mathcal{CA}[x] &= \text{FETCH-}n \\ \mathcal{CA}[a_1 + a_2] &= \mathcal{CA}[a_2] : \mathcal{CA}[a_1] : \text{ADD} \\ \mathcal{CA}[a_1 * a_2] &= \mathcal{CA}[a_2] : \mathcal{CA}[a_1] : \text{MULT} \\ \mathcal{CA}[a_1 - a_2] &= \mathcal{CA}[a_2] : \mathcal{CA}[a_1] : \text{SUB} \end{aligned}$$

For the arithmetic ones, we need to take care because with subtraction, the order matters, hence a_2 being pushed onto the stack before a_1

Now, we can look at $\mathcal{CB}[\cdot]: \text{Bexp} \rightarrow \text{Code}$ (the binary ones).

$$\begin{aligned} \mathcal{CB}[\text{true}] &= \text{TRUE} \\ \mathcal{CB}[\text{false}] &= \text{FALSE} \\ \mathcal{CB}[a_1 = a_2] &= \mathcal{CA}[a_2] : \mathcal{CA}[a_1] : \text{EQ} \\ \mathcal{CB}[a_1 \leq a_2] &= \mathcal{CA}[a_2] : \mathcal{CA}[a_1] : \text{LE} \\ \mathcal{CB}[\neg b] &= \mathcal{CB}[b] : \text{NEG} \\ \mathcal{CB}[b_1 \wedge b_2] &= \mathcal{CB}[b_2] : \mathcal{CB}[b_1] : \text{AND} \end{aligned}$$

We have to use the **stack** rather than work directly with the results.

Finally, we need to translate statements: $\mathcal{CS}[\cdot]: \text{Stm} \rightarrow \text{Code}$.

$$\begin{aligned}
\mathcal{CS}[\![x := a]\!] &= \mathcal{CA}[\![a]\!] : \text{STORE-}x \\
\mathcal{CS}[\![\text{skip}]\!] &= \text{NOOP} \\
\mathcal{CS}[\![S_1; S_2]\!] &= \mathcal{CS}[\![S_1]\!] : \mathcal{CS}[\![S_2]\!] \\
\mathcal{CS}[\![\text{if } b \text{ then } S_1 \text{ else } S_2]\!] &= \mathcal{CB}[\![b]\!] : \text{BRANCH}(\mathcal{CS}[\![S_1]\!], \mathcal{CS}[\![S_2]\!]) \\
\mathcal{CS}[\![\text{while } b \text{ do } S]\!] &= \text{LOOP}(\mathcal{CB}[\![b]\!], \mathcal{CS}[\![S]\!])
\end{aligned}$$

2.6 Semantic function

We can now define a semantic function for **While** (by translating and executing the program on our AM), and it will be called \mathcal{S}_{am} .

$$\begin{aligned}
\mathcal{S}_{am}[\![\cdot]\!] &:\rightarrow (\text{State} \hookrightarrow \text{State}) \\
\mathcal{S}_{am}[\![\cdot]\!] &= \sigma' \text{ if } \mathcal{M}[\![\mathcal{CS}]\!]\sigma = \sigma', \text{ Undefined otherwise} \\
\mathcal{S}_{am}[\![S]\!]\sigma &= (\mathcal{M}^o \mathcal{CS})[\![S]\!]\sigma \\
\mathcal{S}_{am}[\![S]\!] &= (\mathcal{M}^o \mathcal{CS})[\![S]\!]
\end{aligned}$$

Here's an example using factorials:

$$\begin{aligned}
&\mathcal{CS}[\![y:=1, \text{ while } \neg(x=1) \text{ do } (y := y * x; x := x - 1)]\!] \\
&= \mathcal{CS}[\![y := 1]\!] : \\
&\quad \mathcal{CS}[\![\text{while } \neg(x=1) \text{ do } (y := y * x; x := x - 1)]\!] \\
&= \mathcal{CA}[\![1]\!] : \text{store-}y : \\
&\quad \text{loop}(\mathcal{CB}[\![y := y * x; x := x - 1]\!]) \\
&= \text{push-1 : store-}y : \\
&\quad \text{loop}(\mathcal{CB}[\![x = 1]\!] : \text{neg}, \\
&\quad \quad \mathcal{CS}[\![y := y * x]\!] : \mathcal{CS}[\![x := x - 1]\!]) \\
&\quad \dots \\
&= \text{push-1 : store-}y : \\
&\quad \text{loop}(\text{push-1 : fetch-}x : \text{eq} : \text{neg}, \\
&\quad \quad \text{fetch-}x : \text{fetch-}y : \text{mult} : \text{store-}y : \\
&\quad \quad \text{push-1 : fetch-}x : \text{sub} : \text{store-}x)
\end{aligned}$$

Because we have specified the functions, they are allowed to be undefined (such as when they do infinite loops).

2.6.1 Correctness of translation

We can say that for any arithmetic expression, all intermediate configurations have a non-empty stack.

Similarly, for all boolean expressions b , all intermediate configurations have a non-empty stack.

Note! We don't have to always change state because we defined the program. In fact, if we wanted to allow $\mathcal{A}[\cdot]$ to allow state changes, we'd have to change the definition of \mathcal{A} . That is to say, if we want to allow 'side-effects', we'd need to change the definition to allow it to return both an integer and a state S .

3 Denotational Semantics

A denotational semantics defines the meaning of a program by a partial function called a **state-transformer** from (initial) states to (final) states. We're gonna be operating at a level that takes functions as arguments, and gives arguments.

A denotational semantics must be **compositional**. This means we must define the meaning of expressions in terms of their *sub-expressions*.

There are two flavours of denotational semantics: **direct-style** and **continuation-style**. We'll mostly look at the first flavour.

We need two more notions to formalise the meaning of if-then-else, and while statements. For this, we'll use **conditional functions** and **fixpoint operators**.

The denotational semantics of **While** statements is specified in a similar way to the functions for arithmetic and boolean operations.

So, we define the direct-style semantic function to be:

$$S_{ds} : \text{Stm} \rightarrow (\text{State} \hookrightarrow \text{State})$$

Here is what the semantics looks like:

$$\begin{aligned} S_{ds}[x := a]s &= s[x \mapsto \mathcal{A}[a]s] \\ S_{ds}[\text{skip}] &= id \\ S_{ds}[S_1; S_2] &= S_{ds}[S_2] \circ S_{ds}[S_1] \\ S_{ds}[\text{if } b \text{ then } S_1 \text{ else } S_2] &= \text{cond}(\mathcal{B}[b], S_{ds}[S_1], S_{ds}[S_2]) \\ S_{ds}[\text{while } b \text{ do } S] &= \text{FIX } F \text{ where } Fg = \text{cond}(\mathcal{B}[b], g \circ S_{ds}[S], id) \end{aligned}$$

'FIX F' is the (least) fixpoint of the **functional** F, which is the functional of the loop. The g is an implicit argument of the functional F. We'll look at this a bit more, cause it's quite rough.

3.1 Conditional functions

The idea of a **conditional function** is closely related to the denotational semantics of conditionals.

We want to use one of two functions c or d to map some inputs x to some outputs y ; and we have a boolean test b for, that will determine which function to apply in each case x .

$$\text{cond} : (X \rightarrow T) \times (X \hookrightarrow Y) \times (X \hookrightarrow Y) \rightarrow (X \hookrightarrow Y)$$

We give in three functions, and we get out a function that takes X to Y . We can rewrite this as:

$$\text{cond}(b, c, d)x = \begin{cases} c(x) & \text{if } b(x) = tt \\ d(x) & \text{otherwise} \end{cases}$$

Where b, c, d correspond to the functions we supply. We will be interested when X and Y are both states. So:

$$\text{cond} : (\text{State} \rightarrow T) \times (\text{State} \hookrightarrow \text{State}) \times (\text{State} \hookrightarrow \text{State}) \times (\text{State} \hookrightarrow \text{State})$$

3.2 Least fixpoint

The idea of the **least fixpoint** of a function is actually quite heavily studied in maths (particularly in lambda calculus).

There are also a lot of heavy duty theorems, but we'll be using this one: **Definition:** For any unary operator $f : X \rightarrow X$ on some domain X with a partial order \leq :

4 Chain-Complete Partial Orders

Sets can have upper and lower bounds depending on where they come in the order of chains. For example, those at the top of the set are the upper bound of every element since they are unable to have an upper bound themselves, while the ones at the bottom are the lower bound, since they have no lower bound themselves.

4.1 Definitions of PO-Set (partially ordered set)

A PO-Set (D, \sqsubseteq) is called a *chain-complete* partially ordered set (ccpo) whenever the least upper bound $\sqcup Y$ exists for all chains $Y \subseteq D$.

- To be a CCPO, each chain must have an upper bound

Furthermore, a PO-Set (D, \sqsubseteq) is called a *complete lattice* whenever the least upper bound $\sqcup Y$ exists for all subsets $Y \subseteq D$.

- To be a lattice, each chain must have a ******

Note that every CCPO (D, \sqsubseteq) has a (necessarily unique) element denoted $\perp = \sqcup \emptyset$. This means that it is given by the least upper bound of the empty chain. First observe \emptyset is a chain, since we know that $\emptyset \subseteq D$ by the basic set properties and $d \sqsubseteq e$ vacuously holds for all $d, e \in \emptyset$ (of which there are none). So, the least upper bound \perp ******

The question is: is our relation a chain-complete partial order (from the slides)? The answer, simply, is no.

- This is because we have a whole bunch of items at the end. There is no least element of the ordering, so it cannot be a CCPO (this is an important part of the CCPO).

What we can do is the **lifted** relation, which we obtain by adding a least element \perp . But, does this now make it a CCPO? Yeah, it does. Nice. HOWEVER, it does not make it a **complete lattice**.

5 The language Exc

Exc is an extension of the **while** language whose statements are defined like as follows. Exc stands for extension:

$$S ::= x := a \mid \text{skip} \mid \text{if } b \text{ then } S_1 \text{ else } S_2 \mid \text{while } b \text{ do } S \mid \text{begin } S_1 \text{ handle } e : S_2 \text{ end} \mid \text{raise } e$$

The idea is that whenever a **raise** exception instruction is encountered, then the execution of the current (encapsulating blocks of) code (such as S_1) is aborted and control is passed to the (most recently defined) handler (such as S_2) for the exception (e).

Consider the following **Exc** statement.

```
begin
  while true do
    if x <= 0
    then raise exit
    else x := x-1
  and exit:
    y:=7
end
```

If this program is run when x is negative then it will terminate after setting y to 7 (leaving x unchanged).

The meaning of an exception needs to capture the result of executing the relevant handler and following by any remaining code after the definition of that handler.

But this requires a more complicated semantic definition which allows, for example, for the fact that (differently from the **while** language) we don't necessarily have to continue running S_3 in the following program after running either S_1 or S_2 (if either of them raises an exception) ******

5.1 Continuation-style denotational semantics

We can define the set **Cont** of continuations as the set of all partial functions between states so that **Cont** = **State** \hookrightarrow **State**

Intuitively, a continuation is simply a state transformer that describes the input-output behaviour of a (part of) program.

This concept allows us to define the behaviour of a statement by the effect that it has (c') on a continuation defining the behaviour of the code following that statement (c).

Recall the **direct style denotational semantics** (3). Continuation style denotational semantics are as follows:

$$\begin{aligned}
 S'_{cs} : \text{Stm} &\rightarrow (\text{Cont} \rightarrow \text{Cont}) \\
 S'_{cs}[[x := a]]cs &= cs[x \mapsto \mathcal{A}[[a]]s] \\
 S'_{cs}[[\text{skip}]] &= \text{id} \\
 S'_{cs}[[S_1; S_2]] &= S'_{cs}[[S_1]]^o S'_{cs}[[S_2]]
 \end{aligned}$$