

Agenda

Simple arithmetic expressions

Syntax

Semantics

Implementation



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]

Compiler Construction

WWW: <http://www.cs.uu.nl/wiki/Cco>

Edition 2011/2012

2



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]



Introduction

§3

3. Simple arithmetic expressions

To develop tools and techniques for reasoning about the syntax and semantics of high-level programming languages, we will now consider a small “programming language” for simple arithmetic expressions.

3



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]



4



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]



3.1 Syntax



Styles of syntax

§3.1

- ▶ Inductive definitions.
- ▶ Inference rules.
- ▶ Generation procedures.
- ▶ BNF-notation.



Syntactic forms

§3.1

Our language will consist of just a few syntactic forms: natural numbers, addition, and multiplication.

For example:

12

$7 + 21$

$18 * 2$

$3 * (17 + 5)$



Inductive definition

§3.1

One way of formally defining the syntax of a language is by means of an **inductive definition**.

Let \mathbb{N} be the set of natural numbers $\{0, 1, 2, \dots\}$.

Definition: The set of *terms* is the smallest set \mathbf{Tm} such that:

1. If $n \in \mathbb{N}$, then $n \in \mathbf{Tm}$.
2. If $t_1 \in \mathbf{Tm}$ and $t_2 \in \mathbf{Tm}$, then $t_1 + t_2 \in \mathbf{Tm}$.
3. If $t_1 \in \mathbf{Tm}$ and $t_2 \in \mathbf{Tm}$, then $t_1 * t_2 \in \mathbf{Tm}$.

Note the use of the word “smallest”: \mathbf{Tm} has no elements other than the ones required by the three clauses of the definition.



The given inductive definition does not say anything about the use of parentheses and precedence rules to disambiguate phrases, such as $3 * (17 + 5)$.

Formally, we are defining *trees* and the structure of a term is immediate from the shape of the tree. The rules that prescribe which trees form terms define the **abstract syntax** of our language.

However, if we write down terms in programs etc., we use a string representation, i.e., a **concrete syntax**, in which we do use parentheses and precedence rules.

In an implementation, the translation from concrete into abstract syntax is carried out by a **parser**.



Inference rules (cont'd)

To define our abstract syntax, we employ the following rules:

$$\frac{n \in \mathbb{N}}{n \in \mathbf{Tm}} [num]$$

$$\frac{t_1 \in \mathbf{Tm} \quad t_2 \in \mathbf{Tm}}{t_1 + t_2 \in \mathbf{Tm}} [add]$$

$$\frac{t_1 \in \mathbf{Tm} \quad t_2 \in \mathbf{Tm}}{t_1 * t_2 \in \mathbf{Tm}} [mul]$$



An alternative approach to defining abstract syntax, is in **natural deduction style**, i.e., by means of so-called **inference rules**:

$$\frac{premise_1 \quad \dots \quad premise_n}{conclusion} [name]$$

Each inference rule consists of zero or more premises, a conclusion, and, optionally, a name: if we have established all premises, then we may derive the conclusion.

An inference rule without premises, is called an **axiom**.



Generation procedure

Yet another means to defining abstract syntax is by giving an explicit **generation procedure**:

$$\begin{aligned} \mathbf{Tm}_0 &= \emptyset \\ \mathbf{Tm}_{i+1} &= \mathbb{N} \cup \{t_1 + t_2, t_1 * t_2 \mid t_1, t_2 \in \mathbf{Tm}_i\} \end{aligned}$$

And then let

$$\mathbf{Tm} = \bigcup_i \mathbf{Tm}_i$$



During the course, we will mostly use **Backus-Naur Form** (BNF) to define abstract syntax.

First we introduce metavariables:

$n \in \mathbf{Num} = \mathbb{N}$	numerals
$t \in \mathbf{Tm}$	terms

The **metavariable** n ranges over the set \mathbf{Num} of numerals, while the metavariable t (or variations such as t_1 and t_2) ranges over the set \mathbf{Tm} of terms.

The abstract syntax of terms is then given by:

$t ::= n \mid t_1 + t_2 \mid t_1 * t_2$



The meaning of terms

§3.2

Now we have established the precise syntax of terms, we need to rigorously define their meaning, i.e., we need to formally define the **semantics** of our language.

For example:

- ▶ The meaning of the program 12 is the natural number 12.
- ▶ The meaning of the program $7 + 21$ is the natural number 28.
- ▶ The program $18 * 2$ has the same meaning as the program $4 * 9$.
- ▶ The program $3 * (17 + 5)$ has the same meaning as the program $3 * 22$.



3.2 Semantics



Styles of semantics

§3.2

- ▶ Axiomatic semantics.
- ▶ Denotational semantics.
- ▶ Operational semantics.
 - ▶ Small-step (structural) operational semantics.
 - ▶ Big-step (natural) operational semantics.



Before we can formalise the semantics of our language, we need some notation.

For $n_1, n_2 \in \mathbb{N}$, we write $n_1 \underline{+} n_2$ to denote the “normal” addition of n_1 and n_2 .

Hence, we distinguish between addition in the **object language**, i.e., $t_1 + t_2$, for which we have yet to provide a meaning, and addition in the **metalanguage**, of which the meaning is assumed to be well-understood.

Similarly, we write $n_1 \underline{*} n_2$ for the metalanguage multiplication of two natural numbers.

$\underline{+}$ and $\underline{*}$ only work on natural numbers, while $+$ and $*$ can be used on arbitrary terms to form new syntactic terms.



Axiomatic semantics (cont'd)

An axiomatic semantics is undirected, rendering it less suitable for a direct mapping to an implementation:

$$\begin{aligned}
 &3 * (17 + 5) \\
 &= \{1 + 2 =_{\text{ax}} 3\} \\
 &(1 + 2) * (17 + 5) \\
 &= \{11 + 6 =_{\text{ax}} 17\} \\
 &(1 + 2) * (11 + 6 + 5)
 \end{aligned}$$

Intuitively, a semantics should provide a means to “simplify” a term.



An **axiomatic semantics** for a language is given by supplying a set of equalities between terms:

$$\begin{aligned}
 n_1 + n_2 &=_{\text{ax}} n_1 \underline{+} n_2 \\
 n_2 * n_2 &=_{\text{ax}} n_1 \underline{*} n_2
 \end{aligned}$$

Example:

$$\begin{aligned}
 &3 * (17 + 5) \\
 &= \{17 + 5 =_{\text{ax}} 22\} \\
 &3 * 22 \\
 &= \{3 * 22 =_{\text{ax}} 66\} \\
 &66
 \end{aligned}$$



Denotational semantics

A **denotational semantics** provides a mapping from the set of terms to another set:

$$\begin{aligned}
 &[[\cdot]] : \mathbf{Tm} \rightarrow \mathbb{N} \\
 &[[n]] = n \\
 &[[t_1 + t_2]] = [[t_1]] \underline{+} [[t_2]] \\
 &[[t_1 * t_2]] = [[t_1]] \underline{*} [[t_2]]
 \end{aligned}$$



For example:

$$\begin{aligned}
 \llbracket 3 * (17 + 5) \rrbracket &= \llbracket 3 \rrbracket * \llbracket 17 + 5 \rrbracket \\
 &= \llbracket 3 \rrbracket * (\llbracket 17 \rrbracket \pm \llbracket 5 \rrbracket) \\
 &= \llbracket 3 \rrbracket * (\llbracket 17 \rrbracket \pm 5) \\
 &= 3 * (\llbracket 17 \rrbracket \pm 5) \\
 &= 3 * (17 \pm 5) \\
 &= 3 * 22 \\
 &= 66
 \end{aligned}$$

☞ A denotational semantics does not necessarily prescribe an “order of evaluation.”



Small-step operational semantics

The first type of operational semantics that we shall consider is a so-called **small-step** or **structured** operational semantics.

A small-step operational semantics defines a machine that takes a term and performs a single step of computation on it, yielding a simpler term. Every step is referred to as a **reduction**.

Computation halts as soon as the term is transformed into a value.

We define the semantics as a set of inference rules with conclusions of the form

$$t \longrightarrow t'$$



An **operational semantics** describes an abstract “machine” that operates on terms, simplifying them to **values**.

For our language, values are just numerals:

$$v \in \mathbf{Val} \quad \text{values}$$

$$v ::= n$$

Often, but not always, we have that $\mathbf{Val} \subset \mathbf{Tm}$.



Small-step operational semantics (cont'd)

$$\frac{}{n_1 + n_2 \longrightarrow n_1 \pm n_2} [r\text{-add}]$$

$$\frac{t_1 \longrightarrow t'_1}{t_1 + t_2 \longrightarrow t'_1 + t_2} [r\text{-add}_1]$$

$$\frac{t_2 \longrightarrow t'_2}{v_1 + t_2 \longrightarrow v_1 + t'_2} [r\text{-add}_2]$$

Rules $[r\text{-add}_1]$ and $[r\text{-add}_2]$ are examples of so-called **congruence rules**: as opposed to the axiom $[r\text{-add}]$ they do not denote a “real” computation step, but rather guide the reduction of a compound term, effectively dictating an order of evaluation.



For multiplication, we have a similar triple of rules:

$$\frac{}{n_1 * n_2 \longrightarrow n_1 * n_2} [r-mul]$$

$$\frac{t_1 \longrightarrow t'_1}{t_1 * t_2 \longrightarrow t'_1 * t_2} [r-mul_1]$$

$$\frac{t_2 \longrightarrow t'_2}{v_1 * t_2 \longrightarrow v_1 * t'_2} [r-mul_2]$$



A sequence of reduction steps.

$$t_1 \longrightarrow t_2 \longrightarrow \dots \longrightarrow t_n$$

such that $t_i \longrightarrow t_{i+1}$ for $1 \leq i \leq n$ is sometimes also written as $t_1 \longrightarrow^* t_n$.

If t_n is a value, $t_n = v$, we say that t_1 evaluates to v in $(n - 1)$ steps.

For instance, $3 * (17 + 5)$ evaluates to 66 in two steps:

$$3 * (17 + 5) \longrightarrow^* 66$$



Example:

$$3 * (17 + 5) \longrightarrow 3 * 22 \longrightarrow 66$$

As a pair of derivation trees, one for each step:

$$\frac{17 + 5 \longrightarrow 22}{3 * (17 + 5) \longrightarrow 3 * 22}$$

$$\frac{}{3 * 22 \longrightarrow 66}$$



An alternative style of operational semantics, **big-step** or **natural** operational semantics, formalises the concept of evaluation directly.

A big-step operational semantics for our language of arithmetic expressions is defined as a natural deduction system with rules for deriving judgements of the form

$$t \Downarrow v$$

where t is term and v the value that it evaluates to.



$$\frac{}{n \Downarrow n} [e\text{-num}]$$

$$\frac{t_1 \Downarrow n_1 \quad t_2 \Downarrow n_2}{t_1 + t_2 \Downarrow n_1 \pm n_2} [e\text{-add}]$$

$$\frac{t_1 \Downarrow n_1 \quad t_2 \Downarrow n_2}{t_1 * t_2 \Downarrow n_1 * n_2} [e\text{-mul}]$$

- ☞ The big-step operational semantics is quite similar to the denotational semantics for our language.
- ☞ The big-step operational semantics is not as precise as the small-step operational semantics with respect to the order of evaluation.



Language = syntax + semantics

Syntax:

$n \in$	Num	numerals
$t \in$	Tm	terms
$v \in$	Val	values

$t ::=$	$n \mid t_1 + t_2 \mid t_1 * t_2$
$v ::=$	n

Semantics:

$t \longrightarrow t'$	reduction
------------------------	-----------

or

$t \Downarrow v$	evaluation
------------------	------------



For example:

$$3 * (17 + 5) \Downarrow 66$$

Derivation tree:

$$\frac{\frac{3 \Downarrow 3}{3 \Downarrow 3} \quad \frac{\frac{17 \Downarrow 17 \quad 5 \Downarrow 5}{17 + 5 \Downarrow 22}}{3 * (17 + 5) \Downarrow 66}}$$



3.3 Implementation




```
type Num_ = Integer
```

- ☞ *Integer* is the Haskell type of arbitrary sized integers. In contrast to the the set `Num` it also contains negative numbers, so some carefulness is in order.
- ☞ The name *Num_* (with the underscore postfix) is used to avoid a name conflict with the type class *Num* of numeric types.



Syntax: terms and values (cont'd)

Some helper functions related to the observation that $\text{Val} \subseteq \text{Tm}$:

```
isValue :: Tm → Bool
isValue (Num _) = True
isValue _       = False
```

```
fromTm :: Tm → Val
fromTm (Num n) = VNum n
fromTm (Add _ _) = error "fromTm: Add"
fromTm (Mul _ _) = error "fromTm: Mul"
```



```
data Tm = Num Num_ | Add Tm Tm | Mul Tm Tm
data Val = VNum Num_
```

```
instance Tree Tm where
```

```
fromTree (Num n) = App "Num" [fromTree n ]
fromTree (Add t1 t2) = App "Add" [fromTree t1 t2]
fromTree (Mul t1 t2) = App "Mul" [fromTree t1 t2]

toTree = parseTree [ app "Num" (Num <$> arg
                             , app "Add" (Add <$> arg <*> arg)
                             , app "Mul" (Mul <$> arg <*> arg)
                   ]
```

```
instance Tree Val where
```

```
fromTree (VNum n) = App "VNum" [fromTree n]
toTree = parseTree [ app "VNum" (VNum <$> arg) ]
```



Small-step operational semantics

```
reduce :: Tm → Maybe Tm -- produces Nothing for irreducible
                        -- terms
```

(exercise)

```
eval :: Tm → Val
```

```
eval t = case reduce t of
```

```
Nothing | isValue t → fromTm t
         | otherwise → ⊥ -- impossible
Just t'      → eval t'
```



$eval :: Tm \rightarrow Val$ $eval (Num\ n) = VNum\ n$ $eval (Add\ t_1\ t_2) = \mathbf{let}\ VNum\ n_1 = eval\ t_1$
 $\qquad\qquad\qquad VNum\ n_2 = eval\ t_2$
 $\qquad\qquad\qquad \mathbf{in}\ VNum\ (n_1 + n_2)$ $eval (Mul\ t_1\ t_2) = \mathbf{let}\ VNum\ n_1 = eval\ t_1$
 $\qquad\qquad\qquad VNum\ n_2 = eval\ t_2$
 $\qquad\qquad\qquad \mathbf{in}\ VNum\ (n_1 * n_2)$

- ☞ Big-step semantics is implemented more directly.
- ☞ Evaluation order is enforced by Haskell's + and *.

