# Introduction to Programming Languages, Compilers and Tools

### Dmitry Boulytchev

### November 19, 2019

## 1  Introduction: Languages, Semantics, Interpreters, Compilers

### 1.1  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 $\mathscr{L}$ is a total map

$$[\![\bullet]\!]_{\mathscr{L}} : \mathscr{L} \to \mathscr{D}$$

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

### 1.2  Interpreters

In reality, the semantics often is described using *interpreters*:

$$eval : \mathscr{L} \to \texttt{Input} \to \texttt{Output}$$

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

$$\forall p \in \mathscr{L}, \forall x \in \texttt{Input} : [\![p]\!]_{\mathscr{L}} \, x = eval \ p \ x$$

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 coinside", and vice-versa.

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

$$\forall p \in \mathscr{L}, \forall x \in \texttt{Input} : [\![p]\!]_{\mathscr{L}} \, x = [\![eval]\!]_{\mathrm{ocaml}} \, p \, x$$

How to define $[\![\bullet]\!]_{\mathrm{ocaml}}$? 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.

## 1.3  Compilers

A compiler is just a language transformer

$$comp : \mathscr{L} \to \mathscr{M}$$

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

$$\forall p \in \mathscr{L} \;\; [\![p]\!]_{\mathscr{L}} = [\![comp \; p]\!]_{\mathscr{M}}$$

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).

## 1.4  The first example: language of expressions

Abstract syntax:

$$
\begin{array}{rcll}
\mathscr{X} & = & \{x, y, z, \ldots\} & \text{(variables)} \\
\otimes & = & \{+, -, \times, /, \%, <, \leq, >, \geq, =, \neq, \vee, \wedge\} & \text{(binary operators)} \\
\mathscr{E} & = & \mathscr{X} & \text{(expressions)} \\
& & \mathbb{N} & \\
& & \mathscr{E} \otimes \mathscr{E} &
\end{array}
$$

Semantics of expressions:

- state $\sigma : \mathscr{X} \to \mathbb{Z}$ assigns values to (some) variables;

- semantics $[\![\bullet]\!]$ assigns each expression a total map $\Sigma \to \mathbb{Z}$, where $\Sigma$ is the set of all states.

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

$$\begin{array}{rcll}
[\![n]\!] & = & \lambda\sigma.n & , \quad n \in \mathbb{N} \\
[\![x]\!] & = & \lambda\sigma.\sigma x & , \quad x \in \mathscr{X} \\
[\![A \otimes B]\!] & = & \lambda\sigma.([\![A]\!]\sigma \oplus [\![B]\!]\sigma) & , \quad A, B \in \mathscr{E}
\end{array}$$

| $\otimes$ | $\oplus$ in ocaml | |
|---|---|---|
| $+$ | $+$ | |
| $-$ | $-$ | |
| $\times$ | $*$ | |
| $/$ | $/$ | |
| $\%$ | **mod** | |
| $<$ | $<$ | |
| $>$ | $>$ | |
| $\leq$ | $<=$ | |
| $\geq$ | $>=$ | see note 1 below |
| $=$ | $=$ | |
| $\neq$ | $<>$ | |
| $\wedge$ | $\&\&$ | |
| $\vee$ | $\|$ | see note 2 below |

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

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

Important observations:

1. $[\![\bullet]\!]$ 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. $[\![\bullet]\!]$ is total, since it takes into account all possible ways to deconstruct any expression.

3. $[\![\bullet]\!]$ 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, $1 \vee x$ is undefined in empty state.

# 2 Statements, Stack Machine, Stack Machine Compiler

## 2.1 Statements

More interesting language — a language of simple statements:

$$\mathscr{S} \quad = \quad \begin{aligned}&\mathscr{X} := \mathscr{E}\\&\mathbf{read}\,(\mathscr{X})\\&\mathbf{write}\,(\mathscr{E})\\&\mathscr{S}\,;\mathscr{S}\end{aligned}$$

Here $\mathscr{E}, \mathscr{X}$ stand for the sets of expressions and variables, as in the previous lecture. Again, we define the semantics for this language

$$[\![\bullet]\!]_{\mathscr{S}} : \mathscr{S} \mapsto \mathbb{Z}^* \to \mathbb{Z}^*$$

with the semantic domain of partial functions from integer strings to integer strings. This time we will use *big-step operational semantics*: we define a ternary relation "$\Rightarrow$"

$$\Rightarrow \subseteq \mathscr{C} \times \mathscr{S} \times \mathscr{C}$$

where $\mathscr{C} = \Sigma \times \mathbb{Z}^* \times \mathbb{Z}^*$ — a set of all configurations during a program execution. We will write $c_1 \overset{S}{\Rightarrow} c_2$ instead of $(c_1, S, c_2) \in \Rightarrow$ and informally interpret the former as "the execution of a statement $S$ in a configuration $c_1$ completes with the configuration $c_2$". The components of a configuration are state, which binds (some) variables to their values, and input and output streams, represented as (finite) strings of integers.

The relation "$\Rightarrow$" is defined by the following deductive system (see Fig. 1). The first three rules are *axioms* as they do not have any premises. Note, according to these rules sometimes a program cannot do a step in a given configuration: a value of an expression can be undefined in a given state in rules ASSIGN and WRITE, and there can be no input value in rule READ. This style of a semantics description is called big-step operational semantics, since the results of a computation are immediately observable at the right hand side of "$\Rightarrow$" and, thus, the computation is performed in a single "big" step. And, again, this style of a semantic description can be used to easily implement a reference interpreter.

With the relation "$\Rightarrow$" defined we can abbreviate the "surface" semantics for the language of statements:

$$\forall S \in \mathscr{S}, \forall \iota \in \mathbb{Z}^* \ : \ [\![S]\!]_{\mathscr{S}}\iota = o \Leftrightarrow \langle \Lambda, i, \varepsilon \rangle \overset{S}{\Longrightarrow} \langle \_, \_, o \rangle$$

# 3 Stack Machine

Stack machine is a simple abstract computational device, which can be used as a convenient model to constructively describe the compilation process.

$$\langle \sigma, \iota, o \rangle \xrightarrow{\text{X} := e} \langle \sigma[x \leftarrow [\![e]\!]_{\mathscr{E}} \, \sigma], \iota, o \rangle \qquad\qquad [\textsc{Assign}]$$

$$\langle \sigma, z\iota, o \rangle \xrightarrow{\textbf{read}\,(x)} \langle \sigma[x \leftarrow z], \iota, o \rangle \qquad\qquad [\textsc{Read}]$$

$$\langle \sigma, \iota, o \rangle \xrightarrow{\textbf{write}\,(e)} \langle \sigma, \iota, o([\![e]\!]_{\mathscr{E}} \, \sigma) \rangle \qquad\qquad [\textsc{Write}]$$

$$\frac{c_1 \xRightarrow{S_1} c' \qquad c' \xRightarrow{S_2} c_2}{c_1 \xRightarrow{S_1 \, ; S_2} c_2} \qquad\qquad [\textsc{Seq}]$$

Figure 1: Big-step operational semantics for statements

In short, stack machine operates on the same configurations, as the language of statements, plus a stack of integers. The computation, performed by the stack machine, is controlled by a program, which is described as follows:

$$
\begin{aligned}
\mathscr{I} \quad &= \quad \text{BINOP} \otimes \\
& \qquad \text{CONST}\,\mathbb{N} \\
& \qquad \text{READ} \\
& \qquad \text{WRITE} \\
& \qquad \text{LD}\,\mathscr{X} \\
& \qquad \text{ST}\,\mathscr{X} \\
\mathscr{P} \quad &= \quad \varepsilon \\
& \qquad \mathscr{I}\,\mathscr{P}
\end{aligned}
$$

Here the syntax category $\mathscr{I}$ stands for *instructions*, $\mathscr{P}$ — for *programs*; thus, a program is a finite string of instructions.

The semantics of stack machine program can be described, again, in the form of big-step operational semantics. This time the set of stack machine configurations is

$$\mathscr{C}_{SM} = \mathbb{Z}^* \times \mathscr{C}$$

where the first component is a stack, and the second — a configuration as in the semantics of statement language. The rules are shown on Fig. 2; note, now we have one axiom and six inference rules (one per instruction).

As for the statement, with the aid of the relation "$\Rightarrow$" we can define the surface semantics of stack machine:

$$\forall p \in \mathscr{P}, \forall i \in \mathbb{Z}^* \; : \; [\![p]\!]_{SM} \, i = o \Leftrightarrow \langle \varepsilon, \langle \Lambda, i, \varepsilon \rangle \rangle \xRightarrow{p} \langle \_, \langle \_, \_, o \rangle \rangle$$

## 3.1 A Compiler for the Stack Machine

A compiler of the statement language into the stack machine is a total mapping

$$[\![\bullet]\!]_{comp} : \mathscr{S} \mapsto \mathscr{P}$$

$$c \overset{\varepsilon}{\Longrightarrow} c \qquad\qquad\qquad \left[\text{S\scriptsize TOP}_{SM}\right]$$

$$\cfrac{\langle (x \oplus y) :: st, c \rangle \overset{p}{\Longrightarrow} c'}{\langle y :: x :: st, c \rangle \overset{(\texttt{BINOP}\, \otimes)p}{=\!\!=\!\!=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{B\scriptsize INOP}_{SM}\right]$$

$$\cfrac{\langle z :: st, c \rangle \overset{p}{\Longrightarrow} c'}{\langle st, c \rangle \overset{(\texttt{CONST}\, z)p}{=\!\!=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{C\scriptsize ONST}_{SM}\right]$$

$$\cfrac{\langle z :: st, \langle s, i, o \rangle \rangle \overset{p}{\Longrightarrow} c'}{\langle st, \langle s, z :: i, o \rangle \rangle \overset{(\texttt{READ})p}{=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{R\scriptsize EAD}_{SM}\right]$$

$$\cfrac{\langle st, \langle s, i, o@z \rangle \rangle \overset{p}{\Longrightarrow} c'}{\langle z :: st, \langle s, i, o \rangle \rangle \overset{(\texttt{WRITE})p}{=\!\!=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{W\scriptsize RITE}_{SM}\right]$$

$$\cfrac{\langle (s\, x) :: st, \langle s, i, o \rangle \rangle \overset{p}{\Longrightarrow} c'}{\langle st, \langle s, i, o \rangle \rangle \overset{(\texttt{LD}\, x)p}{=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{L\scriptsize D}_{SM}\right]$$

$$\cfrac{\langle st, \langle s[x \leftarrow z], i, o \rangle \rangle \overset{p}{\Longrightarrow} c'}{\langle z :: st, \langle s, i, o \rangle \rangle \overset{(\texttt{ST}\, x)p}{=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{S\scriptsize T}_{SM}\right]$$

Figure 2: Big-step operational semantics for stack machine

We can describe the compiler in the form of denotational semantics for the source language. In fact, we can treat the compiler as a *static* semantics, which maps each program into its stack machine equivalent.

As the source language consists of two syntactic categories (expressions and statments), the compiler has to be "bootstrapped" from the compiler for expressions $\llbracket \bullet \rrbracket_{comp}^{\mathscr{E}}$:

$$
\begin{aligned}
\llbracket x \rrbracket_{comp}^{\mathscr{E}} &= [\texttt{LD}\, x] \\
\llbracket n \rrbracket_{comp}^{\mathscr{E}} &= [\texttt{CONST}\, n] \\
\llbracket A \otimes B \rrbracket_{comp}^{\mathscr{E}} &= \llbracket A \rrbracket_{comp}^{\mathscr{E}} @ \llbracket B \rrbracket_{comp}^{\mathscr{E}} @ (\texttt{BINOP}\, \otimes)
\end{aligned}
$$

And now the main dish:

$$
\begin{aligned}
\llbracket x := e \rrbracket_{comp} &= \llbracket e \rrbracket_{comp}^{\mathscr{E}} @ [\texttt{ST}\, x] \\
\llbracket \mathbf{read}\, (x) \rrbracket_{comp} &= [\texttt{READ};\ \texttt{ST}\, x] \\
\llbracket \mathbf{write}\, (e) \rrbracket_{comp} &= \llbracket e \rrbracket_{comp}^{\mathscr{E}} @ [\texttt{WRITE}] \\
\llbracket S_1;\, S_2 \rrbracket_{comp} &= \llbracket S_1 \rrbracket_{comp} @ \llbracket S_2 \rrbracket_{comp}
\end{aligned}
$$

6

$$
\begin{array}{rcll}
[\![n]\!] & = & \lambda\sigma.n & \text{[CONST]} \\
[\![x]\!] & = & \lambda\sigma.\sigma x & \text{[VAR]} \\
[\![A \otimes B]\!] & = & \lambda\sigma.([\![A]\!]\sigma \oplus [\![B]\!]\sigma) & \text{[BINOP]}
\end{array}
$$

(a) Denotational semantics for expressions

$$
c \overset{\varepsilon}{\Longrightarrow} c \qquad\qquad\qquad \left[\text{STOP}_{SM}\right]
$$

$$
\frac{\langle (x \oplus y) :: st, s \rangle \overset{p}{\Longrightarrow} c'}{\langle y :: x :: st, s \rangle \overset{(\texttt{BINOP } \otimes)p}{=\!\!=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{BINOP}_{SM}\right]
$$

$$
\frac{\langle z :: st, s \rangle \overset{p}{\Longrightarrow} c'}{\langle st, s \rangle \overset{(\texttt{CONST } z)p}{=\!\!=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{CONST}_{SM}\right]
$$

$$
\frac{\langle (s\ x) :: st, s \rangle \overset{p}{\Longrightarrow} c'}{\langle st, s \rangle \overset{(\texttt{LD } x)p}{=\!\!=\!\!=\!\!=\!\!\Longrightarrow} c'} \qquad\qquad \left[\text{LD}_{SM}\right]
$$

(b) Big-step operational semantics for stack machine

$$
\begin{array}{rcll}
[\![x]\!]^{\mathscr{E}}_{comp} & = & [\,\texttt{LD } x\,] & \text{[VAR}_{comp}\text{]} \\
[\![n]\!]^{\mathscr{E}}_{comp} & = & [\,\texttt{CONST } n\,] & \text{[CONST}_{comp}\text{]} \\
[\![A \otimes B]\!]^{\mathscr{E}}_{comp} & = & [\![A]\!]^{\mathscr{E}}_{comp} @ [\![B]\!]^{\mathscr{E}}_{comp} @ [\,\texttt{BINOP}\otimes\,]) & \text{[BINOP}_{comp}\text{]}
\end{array}
$$

(c) Compilation

Figure 3: All relevant definitions

# 4  Structural Induction

We have considered two languages (a language of expressions $\mathscr{E}$ and a language of stack machine programs $\mathscr{P}$), and a compiler from the former to the latter. It can be formally proven, that the compiler is (fully) correct in the sense, given in the lecture 1. Due to the simplicity of the languages, the proof technique — *structural induction* — is simple as well.

First, we collect all needed definitions in one place (see Fig. 3). We simplified the description of stack machine semantics a little bit: first, we dropped off all instructions, which cannot be generated by the expression compiler, and then, we removed the input and output streams from the stack machine configurations, since they are never affected by the remaining instructions.

**Lemma 1.** *(Determinism) Let p be an arbitrary stack machine program, and let c, $c_1$ and $c_2$ be arbitrary configurations. Then*

$$c \xRightarrow{p} c_1 \land c \xRightarrow{p} c_2 \Rightarrow c_1 = c_2$$

*Proof.* Induction on the structure of $p$.

**Base case**. If $p = \varepsilon$, then, by the rule $\text{STOP}_{SM}$, $c_1 = c$ and $c_2 = c$. Since no other rule can be applied, we're done.

**Induction step**. If $p = \iota p'$, then, by condition, we have

$$\frac{c' \xRightarrow{p'} c_1}{c \xRightarrow{\iota p'} c_1}$$

and

$$\frac{c'' \xRightarrow{p'} c_2}{c \xRightarrow{\iota p'} c_2}$$

where $c'$ and $c''$ depend only on $c$ and $\iota$. By the case analysis on $\iota$ we conclude, that $c' = c''$. Since $p'$ is shorter, than $p$, we can apply the induction hypothesis, which gives us $c_1 = c_2$. $\qquad\square$

**Lemma 2.** *(Compositionality) Let $p = p_1 p_2$ be an arbitrary stack machine program, subdivided into arbitrary subprograms $p_1$ and $p_2$. Then,*

$$\forall c_1, c_2 : c_1 \xRightarrow{p} c_2 \Leftrightarrow \exists c' : c_1 \xRightarrow{p_1} c' \land c' \xRightarrow{p_2} c_2$$

*Proof.* Induction on the structure of $p$.

**Base case**. The base case $p = \varepsilon$ is trivial: use the rule $\text{STOP}_{SM}$ and get $c' = c_2 = c_1$.

**Induction step**. When $p = \iota p'$, then there are two cases:

- Either $p_1 = \varepsilon$, then $c' = c_1$ trivially by the rule $\text{STOP}_{SM}$, and we're done.

- Otherwise $p_1 = \iota p_1'$, and, thus, $p = \iota p_1' p_2$. In order to prove the lemma, we need to prove two implications:

  1. Let $c_1 \xRightarrow{p = \iota p_1' p_2} c_2$. Technically, we need here to consider three cases (one for each type of the instruction $\iota$), but in all cases the outcome would be the same: we have the picture

  $$\frac{c'' \xRightarrow{p_1' p_2} c_2}{c_1 \xRightarrow{p = \iota p_1' p_2} c_2}$$

  where $c''$ depends only on $\iota$ and $c_1$. Since $p_1' p_2$ is shorter, than $p$, we can apply the induction hypothesis, which gives us a configuration $c'$, such, that $c'' \xRightarrow{p_1'} c'$ and $c' \xRightarrow{p_2} c_2$. The observation $c_1 \xRightarrow{\iota p_1'} c'$ concludes the proof (note, we implicitly use determinism here).

8

2. Let there exists $c'$, such that $c_1 \xrightarrow{\iota p'_1} c'$ and $c' \xrightarrow{p_2} c_2$. From the first relation we have

$$\frac{c'' \xrightarrow{p'_1} c'}{c_1 \xrightarrow{\iota p'_1} c'}$$

where $c''$ depends only on $\iota$ and $c_1$. Since $p'_1 p_2$ is shorter, than $p$, we can apply the induction hypothesis, which gives us $c'' \xrightarrow{p'_1 p_2} c_2$, and, thus, $c_1 \xrightarrow{\iota p'_1 p_2} c_2$ (again, we implicitly use determinism here).

$\square$

**Theorem 1.** *(Correctness of compilation) Let $e \in \mathcal{E}$ be arbitrary expression, $s$ — arbitrary state, and $st$ — arbitrary stack. Then*

$$\langle st, s \rangle \xrightarrow{[\![e]\!]^{\mathcal{E}}_{comp}} \langle ([\![e]\!]\, s) :: st, s \rangle \text{ iff } ([\![e]\!]\, s) \text{ is defined}$$

*Proof.* Induction on the structure of $e$.

**Base case**. There are two subcases:

1. $e$ is a constant $z$. Then:

   - $[\![e]\!]\, s = z$ for each state $s$;
   - $[\![e]\!]^{\mathcal{E}}_{comp} = [\texttt{CONST z}]$;
   - $\langle st, s \rangle \xrightarrow{[\texttt{CONST z}]} \langle z :: st, s \rangle$ for arbitrary $st$ and $s$.

   This concludes the first base case.

2. $e$ is a variable $x$. Then:

   - $[\![s]\!]\, s = s\,x$ for each state $s$, such that $s\,x$ is defined;
   - $[\![e]\!]^{\mathcal{E}}_{comp} = [\texttt{LD x}]$;
   - $\langle st, s \rangle \xrightarrow{[\texttt{CONST z}]} \langle (s\,x) :: st, s \rangle$ for arbitrary $st$ and arbitrary $s$, such that $s\,x$ is defined.

   This concludes the second base case.

**Induction step**. Let $e$ be $A \otimes B$. Then:

- $[\![A \otimes B]\!]\, s = [\![A]\!]\, s \oplus [\![B]\!]\, s$ for each state $s$, such that both $[\![A]\!]\, s$ and $[\![B]\!]\, s$ are defined;
- $[\![A \otimes B]\!]^{\mathcal{E}}_{comp} = [\![A]\!]^{\mathcal{E}}_{comp} @ [\![B]\!]^{\mathcal{E}}_{comp} @ [\texttt{BINOP} \oplus]$;

9

- by the inductive hypothesis, for arbitrary *st* and *s*

$$\langle st, s \rangle \xrightarrow{\;[\![A]\!]^{\mathscr{E}}_{comp}\;} \langle ([\![A]\!]\, s) :: st, s \rangle \text{iff } ([\![A]\!]\, s) \text{ is defined}$$

and

$$\langle ([\![A]\!]\, s) :: st, s \rangle \xrightarrow{\;[\![B]\!]^{\mathscr{E}}_{comp}\;} \langle ([\![B]\!]\, s) :: ([\![A]\!]\, s) :: st, s \rangle \text{iff } ([\![A]\!]\, s) \text{ and } ([\![A]\!]\, s) \text{ are defined}$$

Taking into account the semantics of BINOP$\otimes$ and applying the compositionality lemma, the theorem follows.

$\square$

Control Flow Constructs

## 4.1  Structural Control Flow

We add a few structural control flow constructs to the language:

$$\mathscr{S} \;\; +\!= \;\; \textbf{skip}$$
$$\textbf{if } \mathscr{E} \textbf{ then } \mathscr{S} \textbf{ else } \mathscr{S}$$
$$\textbf{while } \mathscr{E} \textbf{ do } \mathscr{S}$$

The big-step operational semantics is straightforward and is shown on Fig. 4.

In the concrete syntax for the constructs we add the closing keywords "**if**" and "**od**" as follows:

$$\textbf{if } e \textbf{ then } s_1 \textbf{ else } s_2 \textbf{ fi}$$
$$\textbf{while } e \textbf{ do } s \textbf{ od}$$

## 4.2  Extended Stack Machine

In order to compile the extended language into a program for the stack machine the latter has to be extended. First, we introduce a set of label names

$$\mathscr{L} = \{l_1, l_2, \dots\}$$

Then, we add three extra control flow instructions:

$$\mathscr{I} \;\; +\!= \;\; \text{LABEL } \mathscr{L}$$
$$\text{JMP } \mathscr{L}$$
$$\text{CJMP}_x\, \mathscr{L}, \text{ where } x \in \{\text{nz}, \text{z}\}$$

In order to give the semantics to these instructions, we need to extend the syntactic form of rules, used in the description of big-step operational smeantics. Instead of the rules in the form

$$c \xRightarrow{\;\texttt{skip}\;} c \qquad\qquad\qquad [\textsc{Skip}]$$

$$\frac{[\![e]\!]\,\sigma \neq 0 \qquad c \xRightarrow{\;S_1\;} c'}{c = \langle \sigma, \_, \_ \rangle \xRightarrow{\;\texttt{if } e \texttt{ then } S_1 \texttt{ else } S_2\;} c'} \qquad [\textsc{If-True}]$$

$$\frac{[\![e]\!]\,\sigma = 0 \qquad c \xRightarrow{\;S_1\;} c'}{c = \langle \sigma, \_, \_ \rangle \xRightarrow{\;\texttt{if } e \texttt{ then } S_1 \texttt{ else } S_2\;} c'} \qquad [\textsc{If-False}]$$

$$\frac{[\![e]\!]\,\sigma \neq 0 \qquad c \xRightarrow{\;S\;} c' \qquad c' \xRightarrow{\;\texttt{while } e \texttt{ do } S\;} c''}{c = \langle \sigma, \_, \_ \rangle \xRightarrow{\;\texttt{while } e \texttt{ do } S\;} c''} \qquad [\textsc{While-True}]$$

$$\frac{[\![e]\!]\,\sigma = 0}{c = \langle \sigma, \_, \_ \rangle \xRightarrow{\;\texttt{while } e \texttt{ do } S\;} c} \qquad [\textsc{While-False}]$$

Figure 4: Big-step operational semantics for control flow statements

$$\frac{c \xRightarrow{\;p\;} c'}{c' \xRightarrow{\;p'\;} c''}$$

we use the following form

$$\frac{\Gamma \vdash c \xRightarrow{\;p\;} c'}{\Gamma' \vdash c' \xRightarrow{\;p'\;} c''}$$

where $\Gamma, \Gamma'$ — *environments*. The structure of environments can be different in different cases; for now environment is just a program. Informally, the semantics of control flow instructions can not be described in terms of just a current instruction and current configuration — we need to take the whole program into account. Thus, the enclosing program is used as an environment.

Additionally, for a program $P$ and a label $l$ we define a subprogram $P[l]$, such that $P$ is uniquely represented as $p'(\texttt{LABEL } l)P[l]$. In other words $P[l]$ is a unique suffix of $P$, immediately following the label $l$ (if there are multiple (or no) occurrences of label $l$ in $P$, then $P[l]$ is undefined).

All existing rules have to be rewritten — we need to formally add a $P \vdash \ldots$ part everywhere. For the new instructions the rules are given on Fig. 5.

Finally, the top-level semantics for the extended stack machine can be redefined as follows:

$$\frac{P \vdash c \overset{p}{\Longrightarrow} c'}{P \vdash c \xxrightarrow{(\text{LABEL }l)p} c'} \qquad\qquad \left[\text{L\tiny ABEL}_{SM}\right]$$

$$\frac{P \vdash c \xxrightarrow{P[l]} c'}{P \vdash c \xxrightarrow{(\text{JMP }l)p} c'} \qquad\qquad \left[\text{J\tiny MP}_{SM}\right]$$

$$\frac{P \vdash c \xxrightarrow{P[l]} c'}{P \vdash \langle z :: st, c \rangle \xxrightarrow{(\text{CJMP}_{nz}\,l)p} c'}, \; z \neq 0 \qquad \left[\text{C\tiny JMP}^{+}_{nz\,SM}\right]$$

$$\frac{P \vdash c \overset{p}{\Longrightarrow} c'}{P \vdash \langle z :: st, c \rangle \xxrightarrow{(\text{CJMP}_{nz}\,l)p} c'}, \; z = 0 \qquad \left[\text{C\tiny JMP}^{-}_{nz\,SM}\right]$$

$$\frac{P \vdash c \xxrightarrow{P[l]} c'}{P \vdash \langle z :: st, c \rangle \xxrightarrow{(\text{CJMP}_{z}\,l)p} c'}, \; z = 0 \qquad \left[\text{C\tiny JMP}^{+}_{z\,SM}\right]$$

$$\frac{P \vdash c \overset{p}{\Longrightarrow} c'}{P \vdash \langle z :: st, c \rangle \xxrightarrow{(\text{CJMP}_{z}\,l)p} c'}, \; z \neq 0 \qquad \left[\text{C\tiny JMP}^{-}_{z\,SM}\right]$$

Figure 5: Big-step operational semantics for extended stack machine

$$\forall p \in \mathscr{P}, \forall i \in \mathbb{Z}^* \;:\; [\![p]\!]_{SM}\, i = o \Leftrightarrow p \vdash \langle \varepsilon, \langle \lambda, i, \varepsilon \rangle \rangle \overset{p}{\Longrightarrow} \langle _{\text{-}}, \langle _{\text{-}}, _{\text{-}}, o \rangle \rangle$$

## 4.3   Syntax Extensions

With the structural control flow constructs already implemented, it is rather simple to "saturate" the language with more elaborated control constructs, using the method of *syntactic extension*. Namely, we may introduce the following constructs

```
   if e₁ then s₁
elif e₂ then s₂
...
elif eₖ then sₖ
[ else sₖ₊₁ ]
  fi
```

and

```
   for s₁, e, s₂ do s₃ od
```

only at the syntactic level, directly parsing these constructs into the original abstract syntax tree, using the following conversions:

```
   if e_1 then s_1                      if e_1 then s_1
elif e_2 then s_2                  else if e_2 then s_2
...                        ⤳       ...
elif e_k then s_k                  else if e_k then s_k
else s_{k+1}                       else s_{k+1}
fi                                      fi
                                        ...
                                        fi



   if e_1 then s_1                      if e_1 then s_1
elif e_2 then s_2                  else if e_2 then s_2
...                        ⤳       ...
elif e_k then s_k                  else if e_k then s_k
fi                                 else skip
                                        fi
                                        ...
                                        fi



for s_1, e, s_2 do s_3 od        ⤳        s_1;
                                          while e do
                                            s_3;
                                            s_2
                                          od
```

The advantage of syntax extension method is that it makes it possible to add certain constructs with almost zero cost — indeed, no steps have to be made in order to implement the extended constructs (besides parsing). Note, the semantics of extended constructs is provided for free as well (which is not always desirable). Another potential problem with syntax extensions is that they can easily provide unreasonable results. For example, one may be tempted to implement a post-condition loop using syntax extension:

```
repeat s until e                 ⤳        s;
                                          while e == 0 do
                                            s
                                          od
```

However, for nested **repeat** constructs the size of extended program is exponential w.r.t. the nesting depth, which makes the whole idea unreasonable.

# 5 Procedures

Procedures are unit-returning functions. We consider adding procedures as a separate step, since introducing full-fledged functions would require essential redefinition of the semantics for expressions; at the same time the code generation for funictions is a little trickier, than for procedures, so it is reasonable to split the implementation in two steps.

At the source level procedures are added as a separate syntactic category — definition $\mathscr{D}$:

$$\begin{aligned} \mathscr{D} \quad = \quad & \varepsilon \\ & (\textbf{fun}\ \mathscr{X}\ (\mathscr{X}^*)\ \textbf{local}\ \mathscr{X}^*\ \{\mathscr{S}\})\mathscr{D} \end{aligned}$$

In other words, a definition is a (possibly empty) sequence of procedure descriptions. Each description consists of a name for the procedure, a list of names for its arguments and local variables, and a body (statement). In concrete syntax a single definition looks like

```
fun name (a₁, a₂, …, aₖ)
  local l₁, l₂, …, lₙ {
  s
}
```

where *name* — a name for the procedure, $a_i$ — its arguments, $l_i$ — local variables, $s$ — body.

We also need to add a call statement to the language:

$$\mathscr{S}+=\mathscr{X}(\mathscr{E}^*)$$

In a concrete syntax a call to a procedure $f$ with arguments $e_1,\ldots,e_k$ looks like

$$f\,(e_1,\ldots,e_k)$$

Finally, we have to redefine the syntax category for programs at the top level:

$$\mathscr{L}=\mathscr{D}\mathscr{S}$$

In other words, we extend a statement with a set of definitions.

With procedures, we need to introduce the notion of *scope*. When a procedure is called, its arguments are associated with actual parameter values. A procedure is also in "posession" of its local variables. So, in principle, the context of a procedure execution is a set of arguments, local variables and, possibly, some other variables, for example, the global ones. However, the exact details of procedure context manipilation can differ essentially from language to language.

In our case, we choose a static scoping — each procedure, besides its arguments and local variables, has an access only to global variables. To describe this semantics, we need to change the definition of a state we've used so far:

$$\Sigma = (\mathscr{X} \to \mathbb{Z}) \times 2^{\mathscr{X}} \times (\mathscr{X} \to \mathbb{Z})$$

Now the state is a triple: a *global* state, a set of variables, and a *local* state. Informally, in a new state $\langle \sigma_g, S, \sigma_l \rangle$ $S$ describes a set of local variables, $\sigma_l$ — their values, and $\sigma_g$ — the values of all other accessible variables.

We need to redefine all state-manipulation primitives; first, the valuation of variables:

$$\langle \sigma_g, S, \sigma_l \rangle \ x = \left\{ \begin{array}{ll} \sigma_g \ x & , \quad x \notin S \\ \sigma_l \ x & , \quad x \in S \end{array} \right.$$

Then, updating the state:

$$\langle \sigma_g, S, \sigma_l \rangle \ [x \leftarrow z] = \left\{ \begin{array}{ll} \langle \sigma_g[x \leftarrow z], S, \sigma_l \rangle & , \quad x \notin S \\ \langle \sigma_g, S, \sigma_l[x \leftarrow z] \rangle & , \quad x \in S \end{array} \right.$$

As an empty state, we take the following triple:

$$\Lambda = \langle \Lambda, \varnothing, \Lambda \rangle$$

Finally, we need two transformations for states:

$$\begin{array}{rcl} \mathbf{enter} \ \langle \sigma_g, \_, \_ \rangle \ S & = & \langle \sigma_g, S, \Lambda \rangle \\ \mathbf{leave} \ \langle \sigma_g, \_, \_ \rangle \ \langle \_, S, \sigma_l \rangle & = & \langle \sigma_g, S, \sigma_l \rangle \end{array}$$

The first one simulates entering the new scope with a set of local variables $S$; the second one simulates leaving from an inner scope (described by the first state) to an outer one (described by the second).

All exising rules for big-step operational semantics have to be enriched by *functional environment* $\Gamma$, which binds procedure names to their definitions. As this binding never changes during the program interpretation, we need only to propagate this environment, adding $\Gamma \vdash ...$ for each transition "$... \Longrightarrow ...$". The only thing we need now is to describe the rule for procedure calls:

$$\frac{\mathbf{fun} \ f \ (\bar{a}) \ \mathtt{local} \ \bar{l} \ \{S\} = \Gamma f \qquad \Gamma \vdash \left\langle \mathbf{enter} \ \sigma \ (\bar{a}\bar{l})\overline{[a \leftarrow [\![e]\!]\sigma]}, \iota, o \right\rangle \overset{S}{\Longrightarrow} \langle \sigma', \iota', o' \rangle}{\Gamma \vdash \langle \sigma, \iota, o \rangle \overset{f(\bar{e})}{\Longrightarrow} \langle \mathbf{leave} \ \sigma' \sigma, \iota', o' \rangle} \quad [\text{Call}]$$

where $\Gamma f = \mathbf{fun} \ f \ (\bar{a}) \ \mathtt{local} \ \bar{l} \ \{S\}$.

# 6 Extended Stack Machine

In order to support procedures and calls, we enrich the stack machine with three following instructions:

$$\begin{array}{rcl} \mathscr{I} & += & \mathtt{BEGIN} \ \mathscr{X}^* \ \mathscr{X}^* \\ & & \mathtt{CALL} \ \mathscr{X} \\ & & \mathtt{END} \end{array}$$

Informally speaking, instruction `BEGIN` performs entering into the scope of a procedure; its operands are the lists of argument names and local variables; `END` leaves the scope and returns to the call site, and `CALL` performs the call itself.

We need to enrich the configurations for the stack machine as well:

$$\mathscr{C}_{SM} = (\mathscr{P} \times \Sigma) \times \mathbb{Z}^* \times \mathscr{C}$$

Here we added a *control stack* — a stack of pairs of programs and states. Informally, when performing the `CALL` instruction, we put the following program and current state on a stack to use them later on, when corresponding `END` instruction will be encountered. As all other instructions does not affect the control stack, it gets threaded through all rules of operational semantics unchanged.

Now we specify additional rules for the new instructions:

$$\frac{P \vdash \left\langle cs, st, \left\langle \mathbf{enter}\, \sigma\, (\bar{a}@\bar{l}) \overline{[a \leftarrow z]}, i, o \right\rangle \right\rangle \overset{p}{\Longrightarrow} c'}{P \vdash \langle cs, \bar{z}@st, \langle \sigma, i, o \rangle \rangle \xRightarrow{(\texttt{BEGIN}\, \bar{a}\, \bar{l})p} c'} \qquad \left[\text{Begin}_{SM}\right]$$

$$\frac{P \vdash \langle (p, \sigma) :: cs, st, \langle \sigma, i, o \rangle \rangle \xRightarrow{P[f]} c'}{P \vdash \langle cs, st, \langle \sigma, i, o \rangle \rangle \xRightarrow{(\texttt{CALL}\, f)p} c'} \qquad \left[\text{Call}_{SM}\right]$$

$$\frac{P \vdash \langle cs, st, \langle \mathbf{leave}\, \sigma\, \sigma', i, o \rangle \rangle \xRightarrow{p'} c'}{P \vdash \langle (p', \sigma') :: cs, st, \langle \sigma, i, o \rangle \rangle \xRightarrow{\texttt{END}\, p} c'} \qquad \left[\text{EndRet}_{SM}\right]$$

$$P \vdash \langle \varepsilon, st, c \rangle \xRightarrow{\texttt{END}\, p} \langle \varepsilon, st, c \rangle \qquad \left[\text{EndStop}_{SM}\right]$$

# 7 Functions

## 7.1 Functions in Expressions

At the syntax level function calls are introduced with the following construct:

$$\mathscr{E} \quad += \quad \mathscr{X}\,\mathscr{E}^*$$

In the concrete syntax function call looks conventional:

$f\ (e_1, e_2, \ldots, e_k)$

where $f$ — a function name, $e_i$ — its actual parameters.

Surpisingly, such a mild extension results in a complete redefinition of the semantics. Indeed, as function body can perform arbitrary actions on state, input and output streams, expressions with function calls can now have side effects. In order to express these side effects we need to redefine the semantics completely, this time using big-step operational style.

We extend the configuration, used in the semantic description for statements, with a fourth component — an optional integer value:

$$\Phi \vdash \langle \sigma, i, o, \_ \rangle \xrightarrow{\;n\;}_{\mathscr{E}} \langle \sigma, i, o, n \rangle \qquad\qquad [\text{Const}]$$

$$\Phi \vdash \langle \sigma, i, o, \_ \rangle \xrightarrow{\;x\;}_{\mathscr{E}} \langle \sigma, i, o, \sigma\, x \rangle \qquad\qquad [\text{Var}]$$

$$\frac{\Phi \vdash c \xrightarrow{\;A\;}_{\mathscr{E}} c' \qquad \Phi \vdash c' \xrightarrow{\;B\;}_{\mathscr{E}} c''}{\Phi \vdash c \xrightarrow{\;A \otimes B\;}_{\mathscr{E}} \mathbf{ret}\,(\mathbf{val}\,c' \oplus \mathbf{val}\,c'')} \qquad [\text{Binop}]$$

$$\frac{\begin{array}{c} \Phi \vdash c_{j-1} \xrightarrow{\;e_j\;}_{\mathscr{E}} c_j = \langle \sigma_j, i_j, o_j, v_j \rangle \\ \Phi\, f = \mathbf{fun}\, f\,(\bar{a})\,\mathbf{local}\,\bar{l}\,\{s\} \\ \mathbf{skip}, \Phi \vdash \langle \mathbf{enter}\, \sigma_k\,(\bar{a}@\bar{l})\,[\overline{a_j \leftarrow v_j}], i_k, o_k, {-\!\!\!-} \rangle \xrightarrow{\;s\;} \langle \sigma', i', o', n \rangle \end{array}}{\Phi \vdash c_0 = \langle \sigma_0, \_, \_, \_ \rangle \xrightarrow{\;f(\overline{e_k})\;}_{\mathscr{E}} \langle \mathbf{leave}\,\sigma'\,\sigma_0, i', o', n \rangle} \qquad [\text{Call}]$$

Figure 6: Big-step Operational Semantics for Expressions

$$\mathscr{C} = \Sigma \times \mathbb{Z}^* \times \mathbb{Z}^* \times \mathbb{Z}^?$$

This component will correspond to an optional return value for function/procedure calls (either integer value $n$ or "—", nothing).

We introduce two primitives to make the semantic description shorter:

$$\mathbf{val}\,\langle s, i, o, v \rangle = v$$
$$\mathbf{ret}\,\langle s, i, o, \_ \rangle\,v = \langle s, i, o, v \rangle$$

The rules themselves are summarized in the Fig. 6. Note the use of double environment for evaluating the body of a function in the rule CALL; note also, that now semantics of expressions and statements are mutually recursive.

## 7.2 Return Statement

In order to make it possible to return values from procedures we add a return statement to the language of statements:

$$\mathscr{S} \mathrel{+}= \mathbf{return}\,\mathscr{E}^?$$

And, again, this small addition leads to redefinition of the semantics in *continuation-passing style* (CPS).

First, we define the following meta-operator "$\diamond$" on statements:

$$\begin{aligned} s \diamond \mathbf{skip} &= s \\ s_1 \diamond s_2 &= s_1; s_2 \end{aligned}$$

Then, we add another environment component, $K$, in the description of semantic relation "$... \vdash ... \overset{...}{\Longrightarrow} ...$". Informally speaking, now

$$K, \Phi \vdash c \overset{s}{\Longrightarrow} c'$$

is read "the execution of $s$, immediately followed by $K$, in the configuration $c$ results in the configuration $c'$". Statement $K$ is called *continuation*. The rules themselves are shown on Fig. 7. Note, the rule for the call statement is exactly the same, as for the call expression.

# 8 Arrays and strings

An array can be represented as a pair: the length of the array and a mapping from indices to elements. If we denote $\mathscr{E}$ the set of elements then the set of all arrays $\mathscr{A}(\mathscr{E})$ can be defined as follows:

$$\mathscr{A}(\mathscr{E}) = \mathbb{N} \times (\mathbb{N} \to \mathscr{E})$$

For an array $(n, f)$ we assume $\mathtt{dom}\, f = [0..n-1]$. An element selection function:

$$\bullet[\bullet] : \mathscr{A}(\mathscr{E}) \to \mathbb{N} \to \mathscr{E}$$

$$(n, f)[i] = \begin{cases} f\, i & , \quad i < n \\ \bot & , \quad \text{otherwise} \end{cases}$$

We represent arrays by *references*. Thus, we introduce a (linearly) ordered set of locations

$$\mathscr{L} = \{l_0, l_1, \dots\}$$

Now, the set of all values the programs operate on can be described as follows:

$$\mathscr{V} = \mathbb{Z} \uplus \mathscr{L}$$

Here, every value is either an integer, or a reference (some location). The disjoint union "$\uplus$" makes it possible to unambiguously discriminate between the shapes of each value. To access arrays, we introduce an abstraction of memory:

$$\mathscr{M} = \mathscr{L} \to \mathscr{A}(\mathscr{V})$$

We now add two more components to the configurations: a memory function $\mu$ and the first free memory location $l_m$, and define the following primitive:

$$\mathbf{mem}\, \langle s, \mu, l_m, i, o, v \rangle = \mu$$

which gives a memory function from a configuration.

$$\textbf{skip}, \Phi \vdash c \xRightarrow{\textbf{skip}} c \qquad \left[\text{SkipSkip}\right]$$

$$\frac{\textbf{skip}, \Phi \vdash c \xRightarrow{K} c'}{K, \Phi \vdash c \xRightarrow{\textbf{skip}} c'} \qquad \left[\text{Skip}\right]$$

$$\frac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} \langle \sigma, i, o, n \rangle \quad \textbf{skip}, \Phi \vdash \langle \sigma[x \leftarrow n], i, o, - \rangle \xRightarrow{K} c'}{K, \Phi \vdash c \xRightarrow{x := e} c'} \qquad \left[\text{Assign}\right]$$

$$\frac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} \langle \sigma, i, o, n \rangle \quad \textbf{skip}, \Phi \vdash \langle \sigma, i, o@[n], - \rangle \xRightarrow{K} c'}{K, \Phi \vdash c \xRightarrow{\textbf{write}\,(e)} c'} \qquad \left[\text{Write}\right]$$

$$\frac{\textbf{skip}, \Phi \vdash \langle \sigma[x \leftarrow z], i, o, - \rangle \xRightarrow{K} c'}{K, \Phi \vdash \langle \sigma, z : i, o, - \rangle \xRightarrow{\textbf{read}\,(x)} c'} \qquad \left[\text{Read}\right]$$

$$\frac{s_2 \diamond K, \Phi \vdash c \xRightarrow{s_1} c'}{K, \Phi \vdash c \xRightarrow{s_1;\, s_2} c'} \qquad \left[\text{Seq}\right]$$

$$\frac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} \langle \sigma, i, o, n \rangle \quad K, \Phi \vdash \langle \sigma, i, o, - \rangle \xRightarrow{s_1} c' \quad n \neq 0}{K, \Phi \vdash c \xRightarrow{\textbf{if}\ e\ \textbf{then}\ s_1\ \textbf{else}\ s_2} c'} \qquad \left[\text{IfTrue}\right]$$

$$\frac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} \langle \sigma, i, o, 0 \rangle \quad K, \Phi \vdash \langle \sigma, i, o, - \rangle \xRightarrow{s_2} c'}{K, \Phi \vdash c \xRightarrow{\textbf{if}\ e\ \textbf{then}\ s_1\ \textbf{else}\ s_2} c'} \qquad \left[\text{IfFalse}\right]$$

$$\frac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} \langle \sigma, i, o, n \rangle \quad \textbf{while}\ e\ \textbf{do}\ s \diamond K, \Phi \vdash \langle \sigma, i, o, - \rangle \xRightarrow{s} c' \quad n \neq 0}{K, \Phi \vdash c \xRightarrow{\textbf{while}\ e\ \textbf{do}\ s} c'} \qquad \left[\text{WhileTrue}\right]$$

$$\frac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} \langle \sigma, i, o, 0 \rangle \quad \textbf{skip}, \Phi \vdash \langle \sigma, i, o, - \rangle \xRightarrow{K} c'}{K, \Phi \vdash c \xRightarrow{\textbf{while}\ e\ \textbf{do}\ s} c'} \qquad \left[\text{WhileFalse}\right]$$

$$\frac{\begin{array}{c} \Phi \vdash c_{j-1} \xRightarrow{e_j}_{\mathscr{E}} c_j = \langle \sigma_j, i_j, o_j, v_j \rangle \\ \Phi\, f = \textbf{fun}\ f\,(\overline{a})\ \textbf{local}\ \overline{l}\ \{s\} \\ \textbf{skip}, \Phi \vdash \langle \textbf{enter}\ \sigma_k\ (\overline{a}@\overline{l})\ [\overline{a_j \leftarrow v_j}], i_k, o_k, - \rangle \xRightarrow{s} \langle \sigma', i', o', - \rangle \\ \textbf{skip}, \Phi \vdash \langle \textbf{leave}\ \sigma'\ \sigma_0, i', o', - \rangle \xRightarrow{K} \langle \sigma'', i'', o'', n \rangle \end{array}}{K, \Phi \vdash c_0 = \langle \sigma_0, \_, \_, \_ \rangle \xRightarrow{f(\overline{e_k})} \langle \sigma'', i'', o'', n \rangle} \qquad \left[\text{Call}\right]$$

$$K, \Phi \vdash c \xRightarrow{\textbf{return}} c \qquad \left[\text{ReturnEmpty}\right]$$

$$\frac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} c'}{K, \Phi \vdash c \xRightarrow{\textbf{return}\ e} c'} \qquad \left[\text{Return}\right]$$

Figure 7: Continuation-passing Style Semantics for Statements

$$\dfrac{\begin{array}{cc} \Phi \vdash c \stackrel{e}{\Longrightarrow}_{\mathscr{E}} c' & \Phi \vdash c' \stackrel{j}{\Longrightarrow}_{\mathscr{E}} c'' \\ l = \textbf{val}\ c' & j = \textbf{val}\ c'' \\ l \in \mathscr{L} & j \in \mathbb{N} \\ (n, f) = \textbf{mem}\ l & j < n \end{array}}{\Phi \vdash c \stackrel{e[j]}{\Longrightarrow}_{\mathscr{E}} \textbf{ret}\ c''(f\ j)} \quad \left[\textsc{ArrayElement}\right]$$

$$\dfrac{\begin{array}{c} \Phi \vdash c_j \stackrel{e_j}{\Longrightarrow}_{\mathscr{E}} c_{j+1}, j \in [0..k] \\ \langle s, \mu, l_m, i, o, \_\rangle = c_{k+1} \end{array}}{\Phi \vdash c_0 \stackrel{[e_0,e_1,...,e_k]}{\Longrightarrow}_{\mathscr{E}} \langle s, \mu[l_m \leftarrow (k+1, \lambda n.\textbf{val}\ c_n)], l_{m+1}, i, o, l_m\rangle} \quad \left[\textsc{Array}\right]$$

$$\dfrac{\begin{array}{c} \Phi \vdash c \stackrel{e}{\Longrightarrow}_{\mathscr{E}} c' \\ l = \textbf{val}\ c' \\ l \in \mathscr{L} \\ (n, f) = (\textbf{mem}\ c')\ l \end{array}}{\Phi \vdash c \stackrel{e.\texttt{length}}{\Longrightarrow}_{\mathscr{E}} \textbf{return}\ c'\ n} \quad \left[\textsc{ArrayLength}\right]$$

Figure 8: Big-step Operational Semantics for Array Expressions

## 8.1 Adding arrays on expression level

On expression level, abstractly/concretely:

$$\begin{array}{llll} \mathscr{E}+= & \mathscr{E}[\mathscr{E}] & (a[e]) & \text{taking an element} \\ & | \quad [\mathscr{E}^*] & ([e_1,e_2,..,e_k]) & \text{creating an array} \\ & | \quad \mathscr{E}.\texttt{length} & (\texttt{e.length}) & \text{taking the length} \end{array}$$

The semantics of enriched expressions is modified as follows. First, we add two additional premises to the rule for binary operators:

$$\dfrac{\begin{array}{cc} \Phi \vdash c \stackrel{A}{\Longrightarrow}_{\mathscr{E}} c' & \Phi \vdash c' \stackrel{B}{\Longrightarrow}_{\mathscr{E}} c'' \\ \textbf{val}\ c' \in \mathbb{Z} & \textbf{val}\ c'' \in \mathbb{Z} \end{array}}{\Phi \vdash c \stackrel{A \otimes B}{\Longrightarrow}_{\mathscr{E}} \textbf{ret}\ c''\ (\textbf{val}\ c' \oplus \textbf{val}\ c'')} \quad \left[\textsc{Binop}\right]$$

These two premises ensure that both operand expressions are evaluated into integer values. Second, we have to add the rules for new kinds of expressions (see Figure 8).

## 8.2 Adding arrays on statement level

On statement level, we add the single construct:

$$\mathscr{S} += \mathscr{E}[\mathscr{E}] := \mathscr{E}$$

This construct is interpreted as an assignment to an element of an array. The semantics of this construct is described by the following rule:

$$\cfrac{\begin{array}{ccc} \Phi \vdash c \xRightarrow{e}_{\mathscr{E}} c' & \Phi \vdash c' \xRightarrow{j}_{\mathscr{E}} c'' & \Phi \vdash c'' \xRightarrow{g}_{\mathscr{E}} \langle s, \mu, l_m, i, o, v \rangle \\ l = \mathbf{val}\, c' & i = \mathbf{val}\, c'' & \\ l \in \mathscr{L} & i \in \mathbb{N} & \end{array} \\ \begin{array}{c} (n, f) = \mu\, l \\ i < n \end{array} \\ \mathtt{skip}, \Phi \vdash \langle s, \mu[l \leftarrow (n, f[i \leftarrow x])], l_m, i, o, - \rangle \xRightarrow{K} \widetilde{c}}{K, \Phi \vdash c \xRightarrow{e[j] := g} \widetilde{c}} \;\; \left[\textsc{ArrayAssign}\right]$$

## 8.3  Strings

With arrays in our hands, we can easily add strings as arrays of characters. In fact, on the source language the strings can be introduced as a syntactic extension:

1. we add a character constants `'c'` as a shortcut for their integer codes;

2. we add a string literals `"abcd ..."` as a shortcut for arrays `['a', 'b', 'c', 'd', ...]`.

Nothing else has to be done — now we have mutable reference-representable strings.

# 9  S-expressions and Pattern Matching

S-expressions can be considered as arrays with additionally attached *tags* (of constructors). We denote the set of all constructors as $\mathscr{C}$:

$$\mathscr{C} = \{C_1, C_2, \dots\}$$

Thus, the set of all S-expressions can be described as

$$\mathscr{S}_{exp} = \mathscr{C} \times \mathscr{A}(\mathscr{V})$$

In the concrete syntax constructors are represented by identifiers started with capital letters (A-Z). As being represented as augmented arrays, the values of S-expressions can be manipulated by array primitives (in particular, their elements can be accessed and assigned using square brackets); like arrays they are stored in abstract memory and accessed via locations.

## 9.1 S-expressions on expression level

At the expression level S-expressions are represented by a single additional syntactic form:

$$\mathscr{E}+=\mathscr{C}\,\mathscr{E}^*$$

In the concrete syntax this form is represented by expressions "$C\,(e_1,\ e_1,\ ...,\ e_k)$" or just "$C$" if no constructor arguments are specified. The semantics of this form of expression is straightforward: all arguments are sequentially evauated left-to right and their values are packed into an array with appropriate tag attached:

$$\frac{\begin{array}{c}\Phi\vdash c_j\xLongrightarrow{e_j}_\mathscr{E}c_{j+1},\,j\in[0..k]\\ \langle s,\mu,l_m,i,o,\_\rangle=c_{k+1}\end{array}}{\Phi\vdash c_0\xLongrightarrow{C(e_0,e_1,...,e_k)}_\mathscr{E}\langle s,\mu[l_m\leftarrow(C,(k+1,\lambda n.\mathbf{val}\,c_n))],l_{m+1},i,o,l_m\rangle}\quad[\mathrm{S}_{exp}]$$

## 9.2 Patterns Matching

Pattern-matching is a new control construct which allows to discriminate between different constructors of S-expressions. Pattern-matching is added at the statement level, its abstract syntax is as follows:

$$\mathscr{S}+=\mathbf{case}\ \mathscr{E}\,(\mathscr{P}\times\mathscr{S})^*$$

Here $\mathscr{P}$ stands for a new syntactic category — *patterns* (see below). In the concrete syntax:

**case** $e$ **of** $p_1\ \to s_1\ \mid\ ...\ \mid\ p_k\ \to s_k$ **esac**

where $e$ — expression (*scrutinee*), $p_i$ — patterns, $s_i$ — statements.
The syntactic category of patterns is defined as follows:

$$\mathscr{P}=\_\mid\mathscr{X}\mid C\mathscr{P}^*$$

In the concrete syntax: "$\_$", "$x$", "$C\ (p_1,p_2,...,p_k)$" or just "$C$", where $x$ — variable, $C$ — constructor, $p_i$ — patterns. Informally speaking, each pattern describes a set of S-expressions and a set of bindings of variables to sub sub-expression of each of these expressions. To specify this precisely we introduce the following semantics for patterns: for a pattern $p$ and abstract memory function $\mu$ we define

$$[\![p]\!]^\mu:\mathscr{V}\to\bot\uplus(2^\mathscr{X}\times(\mathscr{X}\to\mathscr{V}))$$

$[\![p]\!]^\mu$ defines a procedure to inspect a value $v$; as a result either $\bot$ (which can be interpreted as a witness that $v$ does not belong to the set of values defined by $p$) or a set of variables in the patterns with their bindings is returned. The definition is as follows:

$$
\begin{aligned}
\llbracket \_ \rrbracket^\mu \, v &= \langle \varnothing, \text{empty state} \rangle \\
\llbracket x \rrbracket^\mu \, v &= \langle \{x\}, [x \leftarrow v] \rangle \\
\llbracket C(p_1, p_2, \ldots, p_k) \rrbracket^\mu \, l &= \bigoplus_i (\llbracket p_i \rrbracket^\mu \, v_i) \qquad , \qquad \mu(l) = C(v_1, v_2, \ldots, v_k) \\
\llbracket p \rrbracket^\mu \, v &= \bot \qquad\qquad\quad\;\; , \qquad \text{otherwise}
\end{aligned}
$$

where $\bigoplus$ unions the bindings:

$$
\begin{aligned}
\bot \oplus \_ &= \bot \\
\_ \oplus \bot &= \bot \\
\langle S_1, \sigma_1 \rangle \oplus \langle S_2, \sigma_2 \rangle &= \left\langle S_1 \cup S_2, \lambda x. \left\{ \begin{array}{lll} \sigma_2 x & , & x \in S_2 \\ \sigma_1 x & , & x \in S_1 \setminus S_2 \end{array} \right. \right\rangle
\end{aligned}
$$

As pattern matching may introduce new bindings, we need also to modify the definition of state:

$$
\Sigma = (\mathscr{X} \to \mathscr{V}) \times (2^{\mathscr{X}} \times (\mathscr{X} \to \mathscr{V}))^*
$$

Now we have one global state and a list of local scopes with sets of variables and their bindings. The redefined state primitives look like the follows. Evaluating in the state:

$$
\begin{aligned}
\langle \sigma_g, \varepsilon \rangle \, x &= \sigma_g \, x \\
\langle \sigma_g, \langle S, \sigma_l \rangle :: t \rangle \, x &= \left\{ \begin{array}{lll} \sigma_l \, x & , & x \in S \\ \langle \sigma_g, t \rangle \, x & , & x \notin S \end{array} \right.
\end{aligned}
$$

Updating the state:

$$
\langle \sigma_g, s \rangle \, [x \leftarrow z] = \mathbf{update} \, \langle \sigma_g, s \rangle \, \varepsilon \, x \, z
$$

where **update** is defined as follows:

$$
\begin{aligned}
\mathbf{update} \, \langle \sigma_g, \varepsilon \rangle \, s \, x \, z &= \langle \sigma_g[x \leftarrow z], s \rangle \\
\mathbf{update} \, \langle \sigma_g, \langle S, \sigma_l \rangle :: t \rangle \, s \, x \, z &= \left\{ \begin{array}{lll} \langle \sigma_g, s @ \langle S, \sigma_l[x \leftarrow z] \rangle t \rangle & , & x \in S \\ \mathbf{update} \, \langle \sigma_g, t \rangle \, s @ \langle S, \sigma_l \rangle \, x \, z & , & x \notin S \end{array} \right.
\end{aligned}
$$

Empty state:

$$
\Lambda = \langle \Lambda, \varepsilon \rangle
$$

Primitives **enter** /**leave** :

$$
\begin{aligned}
\mathbf{enter} \, \langle \sigma_g, \_ \rangle \, S &= \langle \sigma_g, \langle S, \Lambda \rangle \rangle \\
\mathbf{leave} \, \langle \sigma_g, \_ \rangle \, \langle \_, t \rangle &= \langle \sigma_g, t \rangle
\end{aligned}
$$

We will also need two additional primitives to push/drop new scopes:

$$
\begin{aligned}
\mathbf{push} \, \langle \sigma_g, t \rangle \, s &= \langle \sigma_g, s :: t \rangle \\
\mathbf{drop} \, \langle \sigma_g, \_ :: t \rangle &= \langle \sigma_g, t \rangle
\end{aligned}
$$

$$\dfrac{\mathtt{skip}, \Phi \vdash \mathbf{drop}\ c \xRightarrow{K}_{\mathscr{E}} c'}{K, \Phi \vdash c \xRightarrow{\mathbf{leave}}_{\mathscr{E}} c'} \qquad \left[\text{Leave}\right]$$

$$\dfrac{\Phi \vdash c \xRightarrow{e}_{\mathscr{E}} c' \qquad K, \Phi \vdash c' \xRightarrow{\langle \mathbf{val}\ c',\ ps \rangle}_{\mathscr{P}} c''}{K, \Phi \vdash c \xRightarrow{\mathbf{case}\ e\ ps}_{\mathscr{E}} c''} \qquad \left[\text{Case}\right]$$

$$\dfrac{[\![p]\!]^{\mathbf{mem}\ c}\ v = r\ (\neq \bot) \qquad \mathbf{leave} \diamond K, \Phi \vdash \mathbf{push}\ c\ r \xRightarrow{s}_{\mathscr{P}} c'}{K, \Phi \vdash c \xRightarrow{\langle v, \langle p, s \rangle :: t \rangle}_{\mathscr{P}} c'} \qquad \left[\text{PatternMatched}\right]$$

$$\dfrac{[\![p]\!]^{\mathbf{mem}\ c}\ v = \bot \qquad K, \Phi \vdash c \xRightarrow{\langle v, t \rangle}_{\mathscr{P}} c'}{K, \Phi \vdash c \xRightarrow{\langle v, \langle p, s \rangle :: t \rangle}_{\mathscr{P}} c'} \qquad \left[\text{PatternNotMatched}\right]$$

Figure 9: Continuation-passing Style Semantics for Pattern Matching

We will also use these primitives for the whole configurations, meaning that they apply only to their state parts and leave all other components intact.

To specify big-step semantics for pattern-matching we need a yet another "intrinsic" statement **leave** (remember, the first one was "$\diamond$"). The rules are shown on Fig. 9; note, we need an additional transition relation "$\Rightarrow_{\mathscr{P}}$" to reflect the top-down pattern-matching discipline.