

## **Algebraic semantics of an imperative programming language as a compiler abstract model\***

A. V. Zamulin

**Abstract.** It is shown in the paper that state-based algebraic semantics of an imperative program can be regarded as a compiler abstract model and can serve as a good assistant of the compiler designer.

### **1. Introduction**

A new mechanism of definition of the semantics of an imperative programming language is proposed in [4]. The mechanism is based on the construction of a mathematical model of a program, so that model components are used as denotations of language constructions. This approach follows the style of algebraic specifications [1] where first a signature, then a set of its models, then construction and interpretation of the terms of this signature, and finally what is a specification and what is the set of its models are defined. According to this, it is defined in [4]:

- how the signature of a Pascal-like program containing the symbols to be used in expressions and statements is constructed;
- what is the set of models (arbitrary object programs) of this signature;
- how terms (expressions and statements of this language) are constructed and interpreted;
- what is a program specification and what is the set of its models (valid object programs).

The program state is represented in this model by a many-sorted algebra whose components directly match the components of the state of a real program, and program functions and procedures are represented by mathematical functions analyzing and/or changing the program state. In this way, the abstract program model consists only of those components that have direct counterparts in a real program.

The semantics of expressions and statements of a programming language is defined in terms of the components of the abstract program model. Since

---

\*The work is supported in part by Russian Foundation for Basic Research under Grant N 04-01-00272.

the components are high-level entities (mainly functions), the semantics is natural and easy to understand.

In this paper, we want to show that this kind of semantics of a programming language is also an abstract model of a compiler, and it can serve as an effective assistant of compiler designers.

## 2. Program signature and compiler symbol table

### 2.1. Basic signature

We illustrate the approach basing on a simple type system with the following grammar:

$$\mathcal{T} ::= \text{BASE} \mid \text{array}(\text{BASE}, \text{BASE}) \mid \text{loc}(\mathcal{T}), \quad (1)$$

where **BASE** is a set of primitive types (for instance, **Boolean**, **integer**, **real**, **from\_one\_to\_ten** and **char**), **array** is a Pascal-like array type constructor, and **loc** is a location (pointer) type constructor. A special type **void** not belonging to  $\mathcal{T}$  will be used as a result type of a function producing nothing except a new program state.

Thus, our type system contains the part **BASE** with the basic signature  $\Sigma_{dat} = (S_{dat}, F_{dat})$ , where  $S_{dat}$  is the set of primitive data types and  $F_{dat}$  is the set of their operation symbols (names or operators) indexed by their functionalities (i.e., by the types of their arguments and results; an operation without arguments is called a *constant*).

Given a programming language, the basic signature is the basic part of the compiler symbol table in the process of compilation of any program of this language. In our case, for instance, the symbol table will contain the names **Boolean**, **integer**, **real**, **from\_one\_to\_ten** and **char** as the names of primitive types, the names **array** and **loc** as the names of type constructors, and the symbols of all operations of primitive data types with indication of their functionalities.

### 2.2. Program schema

The program schema defines the symbols that may be used in the program in addition to the symbols of the basic signature. It is constructed on the base of variable, constant and function declarations. Formally, the *program schema*  $\mathcal{S}$  over the type system  $\mathcal{T}$  consists of:

- 1) a finite set of array types  $AT$ ,
- 2) a finite set of location types  $LT$ ,
- 3) three finite mappings  $var : \text{IDENT} \rightarrow \mathcal{T}$ ,  $const : \text{IDENT} \rightarrow \text{BASE}$ , and  $func : \text{IDENT} \times \mathcal{T}^* \rightarrow \mathcal{T}^v$ , where **IDENT** is a nonempty set of *identifiers*,  $\mathcal{T}^*$  the set of strings of elements each being either  $t$  or  $ref\ t$ , ( $t$  is a type from the set  $\mathcal{T}$ ), and  $\mathcal{T}^v$  the set  $\mathcal{T} \cup \{\text{void}\}$ . Note that constants are not provided with

values and functions are not provided with bodies in the program schema, they are part of the *program specification* discussed in the end of the paper.

In the sequel, the elements of the mappings  $var(x)$  and  $const(x)$  are denoted by pairs  $(x, t)$ , and the elements of the mapping  $func$  are denoted by triples  $(f, r, t)$ , where  $r$  is a sequence of parameter types. Functions with the result type `void` are called *procedures* (in a real program, the result type `void` in the procedure declaration may be omitted).

**Example.** The schema of the following fragment of a Pascal program:

```
const c: integer, d: char;
var x: real, q: ↑ integer,
    m: array from_one_to_ten of Boolean;
func f(integer, var Boolean): char;
proc p(real);
```

is as follows:

```
AT = {array(from_one_to_ten, Boolean)};
LT = {loc(integer)};
const = {(c, integer), (d, char)};
var = {(x, real), (q, loc(integer)),
      (m, array(from_one_to_ten, Boolean))};
func = {(f, (integer, ref Boolean), char), (p, (real), void)}.
```

### 2.3. Program and block signature

The program schema  $\mathcal{S}$  naturally defines the signature  $\Sigma = (S, F)$  extending the signature  $\Sigma_{dat} = (S_{dat}, F_{dat})$  as follows:

- $S = S_{dat} \cup AT \cup LT \cup \{loc(t) \mid \exists (x, t) \in var\}$ ;
- $F = F_{dat} \uplus var \uplus const \uplus (func \cup F_{acc})$ , where  $F_{acc}$  is a set containing the *access function* symbol  $(- \uparrow, loc(t), t)$  for each  $loc(t) \in S$ .

Thus, a location type is created for each variable in the schema, and the access function symbol is created for each location type.

**Example.** Denote the type `array(from_one_to_ten, Boolean)` by `AR`. Then, ignoring the operation symbols of the primitive types, we can compose the following signature for our program schema:

```
S = {Boolean, integer, real, from_one_to_ten,
     char, AR, loc(integer), loc(real),
     loc(loc(integer)), loc(AR)};
F = {var'(x, real), var'(q, loc(integer)), var'(m, AR),
     const'(c, integer), const'(d, char),
     func'(f, (integer, ref Boolean), char),
     func'(p, real, void), func'(- ↑, loc(integer), integer),
```

```

func'(- ↑, loc(real), real),
func'(- ↑, loc(loc(integer)), loc(integer)),
func'(- ↑, loc(AR), AR)}.

```

It is easy to notice that the program signature is a “template” of the compiler symbol table (the term “template” means here that constants are not yet provided with values, variables with addresses, etc.). Thus, the set  $S$  indicates which types are used in the program, and the set  $F$  indicates which symbols may be used in its expressions and statements. All symbols in  $F$  are partitioned into the groups corresponding to variable, constant, and function declarations. Each symbol is accompanied by its *functionality*. This is the type of a value for a constant, the type of a location for a variable, and the type of a function for a function symbol.

A set of symbols can also be declared in a block. Formally, this set is modelled by the notion of signature increment defined in the simplest case as follows. Let  $t_1, \dots, t_n$  be types and  $X = \{x_1, \dots, x_n\}$  a set of variable identifiers, then the set

$$\Delta = \{(x_1, t_1), \dots, (x_n, t_n)\}$$

is a *signature increment*.

Using a source signature  $\Sigma = (S, F)$  and the signature increment  $\Delta$ , we construct the block signature  $\Sigma' = (S, F')$  as follows:

1) the signature  $\Sigma$  is reduced to the intermediate signature  $\Sigma'' = (S, F'')$  by defining  $F'' = F \setminus \{(x, t) \mid \exists t'. (x, t') \in \Delta\}$  (i.e., by deleting the names redeclared in the block from the source signature);

2) the signature  $\Sigma' = (S, F')$  is extended to the signature  $\Sigma' = (S, F')$  by defining  $F' = F'' \cup \Delta$  (i.e., by adding the names declared in the block to the intermediate signature  $\Sigma''$ ).

In this way, the *block signature*  $\Sigma'$  is constructed as an extension of the reduced source signature by new symbols. Note that, as a result, the block signature does not possess the symbols redeclared in the block, which corresponds to the principle of symbol hiding in the traditional block-structured programming languages.

**Example.** Let the following variables be declared in block  $Bl$ :

```
var x: integer, q: real.
```

This gives us the following signature increment:

$$\Delta = \{(x, \text{integer}), (q, \text{real})\}$$

which extends our signature  $\Sigma$  to the following block signature:

```

S = {Boolean, integer, real, from_one_to_ten,
     char, AR, loc(integer), loc(real),
     loc(loc(integer)), loc(AR)};
F' = {var'(x, integer), var'(q, real), var'(m, AR),

```

```

const'(c, integer), const'(d, char),
func'(f, ⟨integer, ref Boolean⟩, char),
func'(p, real, void), func'(- ↑, loc(integer), integer),
func'(- ↑, loc(real), real),
func'(- ↑, loc(loc(integer)), loc(integer)),
func'(- ↑, loc(AR), AR)}.

```

When we exit the block, we just return to the source signature.

To avoid storing a chain of signatures (the chain can be quite long in case of a deep block nesting), the symbol table in the compiler is partitioned into blocks corresponding to the program signature and signature increments, and identifier search starts in the deepest block. The return to the previous signature is done just by deleting the latest signature increment.

### 3. Program algebra and program state

The program state is formally represented by a many-sorted algebra whose components are the sets of data type values, primitive type operations, access functions, constant values, variable locations, etc. Two main parts can be distinguished in this algebra: the *static part* (basic algebra) mainly consisting of values and operations of the primitive types, and *dynamic part* mainly consisting of access functions and dynamically created memory locations.

#### 3.1. Basic algebra

The *basic algebra* defines the sets and operations that are basic for the other program components. It is a formal counterpart of the implementations of the primitive types in a given compiler. A *basic algebra*  $\mathbf{B}$  associates with each type  $t$  a set  $\mathbf{B}_t$ , its *carrier*, with each operation symbol of  $t$  a partial function, its *implementation*, and with each constant an element of the corresponding carrier. In addition to this, the components of the basic algebra are a special set  $\mathbf{B}_{\text{Locs}}$  containing elements of all location types, a special element `nil` not belonging to any carrier, and element  $\perp$  as a single value of type `void`. The carriers are pairwise disjoint.

In a particular compiler, the set  $\mathbf{B}_{\text{Locs}}$  may be a set of addresses to be allocated to variables in the stack or heap, the element `nil` the null address, the set  $\mathbf{B}_{\text{char}}$  a set of bytes, the set  $\mathbf{B}_{\text{real}}$  a set of 8-byte elements, and so on. A function associated with an operation symbol may be one or more machine instructions implementing this operation (it is an inline function as a rule).

### 3.2. Dynamic part of the program state

The dynamic part of the program state extends the basic algebra  $\mathbf{B}$  to a *state algebra*  $\mathbf{A}$  in accordance with the extension of the basic signature  $\Sigma_{dat}$  to the program signature  $\Sigma = (S, F)$ . This extension, called the  $\Sigma$ -*state*, is created as follows:

- $\mathbf{A}_t = \mathbf{B}_t$  for any primitive type  $t$  (i.e., the carriers of all primitive types remain the same);
- a finite set  $\mathbf{A}_{loc(t)} = \mathbf{A}_{loc(t)}^o \cup \{\mathbf{nil}\}$ , where  $\mathbf{A}_{loc(t)}^o \subset \mathbf{B}_{Locs}$ , is associated with each type  $loc(t)$  so that  $\mathbf{A}_{loc(t)}^o \cap \mathbf{A}_{loc(t')}^o = \emptyset$  for two different types  $t$  and  $t'$  (i.e., the sets of elements of each location type in each program state have only one common element,  $\mathbf{nil}$ );
- a set  $\mathbf{A}_{array(ind,t)}$  of finite injective mappings  $\mathbf{A}_{ind} \rightarrow \mathbf{A}_{loc(t)}^o$  is associated with each array type  $array(ind,t)$  (i.e., an array is a mapping of the index type to the element type such that each index value is mapped to a separate location);
- an element,  $c^A \in \mathbf{A}_t$ , is associated with each constant declaration  $(c,t)$  (i.e., a value of the indicated type is associated with each constant identifier);
- an element,  $x^A \in \mathbf{A}_{loc(t)}^o$ , called *location constant*, is associated with each variable declaration  $(x,t)$  (in contrast to the constant, a location of the indicated type is associated with a variable identifier);
- a partial function  $!^A : \mathbf{A}_{loc(t)} \rightarrow \mathbf{A}_t$  is associated with each access function symbol  $!_{loc(t),t}$  (i.e., a variable location may be initialized with a value of the corresponding type).

In this way, each location type is provided with a set of locations (called a *location sort* in the sequel) and a special element  $\mathbf{nil}$ , each constant is provided with a value, and each variable is provided with a location that may be initialized.

A state algebra is well-formed if:

- each location constant is supplied with a different location,
- the range of the finite mapping associated with an array variable does not intersect with both the range of the finite mapping associated with any other array variable and the set of locations assigned to variables,
- each array location is initialized.

Thus, the notion of the well-formedness of the state algebra is a direct indication for the compiler designer that each program variable should be supplied with a separate address, and, in addition, each array variable should

be supplied with the needed number of element type locations (this means that an array itself and each of its elements are provided with a location). Note that the array element locations are not required to be adjacent in a well-formed state algebra. The placement of array elements in adjacent locations may be considered as a natural means of program optimization. It is interesting to note that although in a real object program the address of an array and the address of its first element are the same address, formally they are different addresses (locations) of different types.

A location  $\mathbf{l}$  is said to be *allocated* in a state  $\mathbf{A}$  if  $\mathbf{l} \in \mathbf{A}_{\text{loc}(t)}$  for some  $\text{loc}(t) \in S$ . The set of all allocated locations in  $\mathbf{A}$  is denoted by  $\mathbf{A}_{\text{loc}}$  in the sequel.

If  $\mathbf{l}$  and  $\mathbf{a}$  are elements of respective sorts  $\text{loc}(t)$  and  $t$  such that  $\mathbf{l} \uparrow = \mathbf{a}$ , then  $\mathbf{a}$  is the *content* of the location  $\mathbf{l}$  ( $\mathbf{a}$  is *stored* in  $\mathbf{l}$ ). The application of the function " $\uparrow$ " to its argument  $\mathbf{l}$  is called the *dereferencing* of  $\mathbf{l}$ .

As it is indicated above, when we enter a block with a signature increment  $\Delta$ , we move from a source signature  $\Sigma$  to the block signature  $\Sigma' = (S, F')$  that contains new symbols of location constants (variable identifiers). In this case, we need to transform some  $\Sigma$ -algebra  $\mathbf{A}$  to a  $\Sigma'$ -algebra  $\mathbf{A}'$  by allocating memory for local variables. Formally, this is performed in two steps:

- 1) the algebra  $\mathbf{A}$  is defined to be a  $\Sigma''$ -algebra  $\mathbf{A}''$  (in this way, the components redeclared in the block become invisible);
- 2) the algebra  $\mathbf{A}''$  is extended with a constant  $\mathbf{x}^{\mathbf{A}'}$  of type  $\text{loc}(t)$  for each pair  $(x, t) \in \Delta$  (a finite mapping is also associated with  $\mathbf{x}^{\mathbf{A}'}$  if  $x$  is an array variable).

In this way, we have constructed a block algebra that does not possess the redeclared components and possesses a location for each local variable. The corresponding sets of locations  $\mathbf{A}_{\text{loc}(t)}$  and finite mappings  $\mathbf{A}_{\text{ind}} \rightarrow \mathbf{A}_{\text{loc}(t)}^{\circ}$  are also extended in the new algebra (to take into account the newly allocated locations). The requirement of algebra correctness is also extended in an evident way.

When the block finishes its work, we get some  $\Sigma'$ -algebra  $\mathbf{A}'_1$  in which some locations of the original algebra  $\mathbf{A}$  may have a new content. To fix this new content and get rid of local block locations, we do the following actions:

- 1) the algebra  $\mathbf{A}'_1$  is reduced to a  $\Sigma''$ -algebra  $\mathbf{A}''_1$  by deleting the components associated with the elements of the set  $F' \setminus F''$  (i.e., by deleting the components declared in the block);
- 2) the algebra  $\mathbf{A}''_1$  is extended to a  $\Sigma$ -algebra  $\mathbf{A}_1$  by adding the components of the algebra  $\mathbf{A}$  associated with the elements of the set  $F \setminus F''$  (i.e., the components redeclared in the block are actually restored).

The resulting algebra  $\mathbf{A}_1$  possesses all updates of the original algebra. It is clear that this formalism of entering and exiting the block naturally models the program stack pushed in when the block is entered and popped

out when it is exited. Note that, in contrast to the traditional algebraic semantics (see [3]), we did not need to define an extra type "stack" with corresponding operations; everything is easily modelled with the use of the notions of algebra extension and reduction.

A location sort or access function may be different in different states, i.e., locations may be dynamically allocated and deleted in the process of program execution, and a variable may be assigned different values. However, different  $\Sigma$ -states must have the same basic algebra (i.e., the same data type implementation). We denote the set of all  $\Sigma$ -states with the same base  $B$  by  $\text{state}_B(\Sigma)$  and mean by a  $\Sigma_B$ -state a  $\Sigma$ -state with the basic algebra  $B$ .

### 3.3. State update

One state can be transformed into another by a *state update* which is either a *location update* or *sort update*. The first one serves for an update of a program variable and the other one serves for extending a location sort with a new location.

A location update in a  $\Sigma_B$ -state  $\mathbf{A}$  is a pair  $(\mathbf{l}, \mathbf{a})$ , where  $\mathbf{l}$  is an element of the set  $\mathbf{A}_{\text{loc}(\mathfrak{t})}$  such that  $\mathfrak{t}$  is not an array type and  $\mathbf{a}$  an element of the set  $\mathbf{A}_{\mathfrak{t}}$ . ■

Transformation of a  $\Sigma_B$ -state  $\mathbf{A}$  into a  $\Sigma_B$ -state  $\mathbf{A}\alpha$  by the location update  $\alpha = (\mathbf{l}, \mathbf{a})$  is performed by updating the access function  $\uparrow_{\text{loc}(\mathfrak{t}), \mathfrak{t}}$  in such a way that  $(\mathbf{l} \uparrow)^{\mathbf{A}\alpha} = \mathbf{a}$ . It is evident that the content of only one location is changed by this update which corresponds to the assignment of a new value to a variable or the machine instruction of updating a memory location. Note that an array location may not be updated, i.e., it cannot be assigned a new finite mapping (otherwise, the locations of the array elements would be changed); however, the location of any array element may be updated.

A sort update  $\delta$  in a  $\Sigma_B$ -state  $\mathbf{A}$  is a pair  $(\mathbf{A}_{\text{loc}(\mathfrak{t})}, \mathbf{l})$ , where  $\mathbf{l}$  is a location such that  $\mathbf{l} \in \mathbf{B}_{\text{Locs}}$  and  $\mathbf{l} \notin \mathbf{A}_{\text{loc}}$ . ■

Transformation of a  $\Sigma_B$ -state  $\mathbf{A}$  into a  $\Sigma_B$ -state  $\mathbf{A}\delta$  by the sort update  $\delta = (\mathbf{A}_{\text{loc}(\mathfrak{t})}, \mathbf{l})$  is performed by inserting the location  $\mathbf{l}$  into the set  $\mathbf{A}_{\text{loc}(\mathfrak{t})}$ . In this way, the operation of allocating a new memory location is modelled. The answer to the question of where the new location is stored will be given in the next section.

## 4. Program

The abstract program model is a counterpart of the object program possessing a set of states and a set of functions and procedures manipulating the state. Recall that the notion of function in programming languages is different from the notion of function in mathematics since the program



function may have a side effect (update of the program state). One can however convert the program function into the mathematical one by representing its result as a pair, *state* and *value*. Let  $|\mathbf{P}(\mathbf{B})| \subset \mathbf{state}_{\mathbf{B}}(\Sigma)$  be the set of the states of a program  $\mathbf{P}(\mathbf{B})$  and  $t$  be a type. Then  $|\mathbf{P}(\mathbf{B})|_{\mathbf{t}} = \{\langle \mathbf{A}, \mathbf{v} \rangle \mid \mathbf{A} \in |\mathbf{P}(\mathbf{B})| \ \mathbf{v} \in \mathbf{A}_{\mathbf{t}}\}$  is a set of pairs with a state algebra as the first element and an element of the state algebra as the second element. This set may be used as the codomain of the functions declared in the program. In what follows, for any  $\Sigma_{\mathbf{R}}$ -state  $\mathbf{A}$  and any sequence of types  $r = t_1, \dots, t_n$ , we denote the Cartesian product  $\mathbf{A}_{t_1} \times \dots \times \mathbf{A}_{t_n}$  by  $\mathbf{A}_{\mathbf{r}}$ , which is a singleton set if  $n = 0$ .

Program  $\mathbf{P}(\mathbf{B})$  of signature  $\Sigma$  consists of:

- 1) a subset  $|\mathbf{P}(\mathbf{B})|$  of  $\mathbf{state}_{\mathbf{B}}(\Sigma)$  called the *carrier* of  $\mathbf{P}(\mathbf{B})$  (each  $\mathbf{A} \in |\mathbf{P}(\mathbf{B})|$  is a *program state*), and
- 2) for each  $(f, r, t) \in \mathit{func}$ , where  $r = u_1, \dots, u_n$  and  $u_i$  is either  $t_i$  or *ref*  $t_i$ , and each  $\mathbf{A} \in |\mathbf{P}(\mathbf{B})|$ , a partial function

$$f_{\mathbf{r}}^{\mathbf{A}} : \mathbf{A}_{u_1} \times \dots \times \mathbf{A}_{u_n} \rightarrow |\mathbf{P}(\mathbf{B})|_{\mathbf{t}},$$

where  $\mathbf{A}_{u_i}$  is  $\mathbf{A}_{\mathbf{loc}(t_i)}^{\circ}$  if  $u_i$  is *ref*  $t_i$ , and  $\mathbf{A}_{t_i}$  in the opposite case.

Thus, a program possesses a number of states, and, in each state, a function associated with the corresponding function declaration. Each function uses a location sort  $\mathbf{A}_{\mathbf{loc}(t)}^{\circ}$  for a parameter declared as *ref*  $t$  to provide call by reference substitution. Note that functions associated with the same function declaration may be different in different states (for instance, because they depend on the content of global variables). In other words, the program state is an implicit function parameter. A function always produces a pair, an algebra and value (an element of this algebra).

The carrier  $|\mathbf{P}(\mathbf{B})|$  must satisfy the following program invariant: any constant (including a location constant) and a mapping associated with an array variable must be the same in any  $\mathbf{A} \in |\mathbf{P}(\mathbf{B})|$ . This is a natural requirement stating that a constant may not have different values and neither a variable nor an array element may be provided with different locations in different program states. Note that the location update described in the previous section cannot violate the program invariant.

## 5. Function call

It is shown in the previous section that the set  $\mathbf{A}_{\mathbf{loc}(t)}$  is used as a function domain component in algebra  $\mathbf{A}$  for a function parameter declared as *ref*  $t$ . From this, it follows that the corresponding argument in the function call expression must also be a location. In contrast to this, a value of type  $t$  must be the argument for a function parameter declared as  $t$  (of course

the argument is a location if  $t$  is a location type). At the same time, a programming language usually permits a location to be also an argument in this case. Therefore, the rule of construction of the function call expression is as follows: if  $(f, r, t)$ , where  $r = u_1, \dots, u_n$  and  $u_i$  is either  $t_i$  or  $\text{ref } t_i$ , is a function declaration and

$$e_i \text{ is } \begin{cases} \text{expression of type } \text{loc}(t_i), & \text{if } u_i \text{ is } \text{ref } t_i, \\ \text{expression either of type } t_i \text{ or } \text{loc}(t_i), & \text{otherwise,} \end{cases}$$

then  $f(e_1, \dots, e_n)$  is an expression of type  $t$  called the *function call*.

Thus, only a location type expression may be an argument of the function call for a parameter declared as  $\text{ref } t$ , and both an expression of type  $t$  and an expression of type  $\text{loc}(t)$  may be an argument in the opposite case. However, the argument of the function itself in the second case may be only a value of type  $t$ , i.e., the location has to be dereferenced.

Denote the interpretation of an expression  $e$  in algebra  $\mathbf{A}$  by  $\llbracket e \rrbracket^{\mathbf{A}}$ . Recall also that an argument expression may be a function call, which, as it was shown in the previous section, yields a pair, an algebra and its element. Therefore, the interpretation of an expression also generally produces such a pair (the same algebra  $\mathbf{A}$ , if no expression component updates the state). Assume that the function  $\mathbf{fst}$  applied to the pair produces its first component and the function  $\mathbf{snd}$  produces its second component. The interpretation of the function call is then as follows: let  $\mathbf{A}_0$  be a  $\Sigma$ -algebra,  $\mathbf{A}_1 = \mathbf{fst}(\llbracket e_1 \rrbracket^{\mathbf{A}_0})$ ,  $\dots$ ,  $\mathbf{A}_n = \mathbf{fst}(\llbracket e_n \rrbracket^{\mathbf{A}_{n-1}})$ ,

$$v_i = \begin{cases} \mathbf{snd}(\llbracket e_i \rrbracket^{\mathbf{A}_{i-1}}) & \text{if } u_i \text{ is } \text{ref } t_i \text{ and } e_i \text{ is} \\ & \text{expression of type } \text{loc}(t_i), \\ \mathbf{snd}(\llbracket e_i \rrbracket^{\mathbf{A}_{i-1}}) & \text{if } u_i \text{ is } t_i \text{ and } e_i \text{ is} \\ & \text{expression of type } t_i, \\ \mathbf{snd}(\llbracket e_i \uparrow \rrbracket^{\mathbf{A}_{i-1}}) & \text{if } u_i \text{ is } t_i \text{ and } e_i \text{ is} \\ & \text{expression of type } \text{loc}(t_i), \end{cases}$$

then

$$\llbracket f(e_1, \dots, e_n) \rrbracket^{\mathbf{A}_0} = \mathbf{f}_r^{\mathbf{A}_n}(v_1, \dots, v_n).$$

Thus, the interpretation of the function call in algebra  $\mathbf{A}$  is the invocation of the corresponding function in algebra  $\mathbf{A}_n$  produced by the interpretation of the last argument. The argument is dereferenced when a location instead of a value is passed as argument and used directly in all the other cases. This interpretation rule directly tells the compiler designer how function calls should be compiled.

**Example.** Let the function  $f$  be invoked in block  $Bl$  as follows:

$$f(x, m[1])$$

and let  $\mathbf{A}$  be a program algebra of the signature  $\Sigma$  and  $\mathbf{A}_{B1}$  an algebra of the block  $B1$  extending  $\mathbf{A}$  with the location constant  $x^{A_{B1}}$ . Note that since the function is called in the block  $B1$  containing the declaration of a local variable  $x$ , the location associated with  $x$  in algebra  $\mathbf{A}_{B1}$  will be used in the call; the location associated with  $m$  in algebra  $\mathbf{A}_{B1}$  remains the same as in algebra  $\mathbf{A}$ , i.e.,  $m^{A_{B1}} = m^{\mathbf{A}}$ . Note also that the algebra produced by evaluation of the actual parameters remains the same algebra  $\mathbf{A}_{B1}$  since evaluation of the parameters does not update the state. The interpretation of this function call then causes the following function invocation:

$$f^{\mathbf{A}}(x^{A_{B1}} \uparrow, m^{\mathbf{A}}(1)).$$

## 6. Updating statements

Let us see how state updates can be used as semantics of the assignment statement and dynamic variable creation statement. The semantics of a statement  $St$  interpreted in algebra  $\mathbf{A}$  is denoted by  $\llbracket St \rrbracket^{\mathbf{A}}$ . Like the semantics of expressions it generally produces a pair,  $\langle \mathbf{A}', v \rangle$ , where  $\mathbf{A}'$  is a new algebra and  $v$  is a value. The value may be the label indicated in the goto statement or the value produced by the return statement or the value  $\perp$  if the statement does not produce a value.

The assignment statement normally has the following form:  $e_l := e_r$ , where  $e_l$  is an expression of type  $loc(t)$  and  $e_r$  expression either of type  $t$  or of type  $loc(t)$ .

We first define the semantics of the statement for the case where  $t$  is not an array type assuming that the interpretation of any part of the statement may update the state.

$$\text{Let } \llbracket e_l \rrbracket^{\mathbf{A}} = \langle \mathbf{A}_l, l \rangle, \llbracket e_r \rrbracket^{\mathbf{A}} = \langle \mathbf{A}_r, v_r \rangle,$$

$$\alpha = \begin{cases} (l, v_r), & \text{if } e_r \text{ has type } t, \\ (l, v_r \uparrow), & \text{if } e_r \text{ has type } loc(t), \end{cases}$$

$$\mathbf{A}_a = \mathbf{A}_r \alpha. \text{ Then } \llbracket e_l := e_r \rrbracket^{\mathbf{A}} = \langle \mathbf{A}_a, \perp \rangle.$$

According to this semantics, the following actions should be undertaken when executing the statement:

- 1) evaluation of the left side of the statement in the source state with its possible update and getting the address of the location to be updated;
- 2) evaluation of the right side of the statement in the new state with possible its update and getting the updating value or its location;
- 3) update of the location with getting a new state.

If  $ar1$  and  $ar2$  are array variables with  $n$  elements of type  $t$ , then the interpretation of the assignment  $ar1 := ar2$  produces an update set

$$\Gamma_{\mathbf{a}} = \{(\mathbf{l}_1, \mathbf{v}_1), \dots, (\mathbf{l}_n, \mathbf{v}_n)\},$$

where  $\mathbf{l}_i = ar1 \uparrow (i)$  and  $\mathbf{v}_i = ar2 \uparrow (i) \uparrow$ . Then  $\llbracket ar1 := ar2 \rrbracket^{\mathbf{A}} = \langle \mathbf{A}_{\mathbf{a}}, \perp \rangle$ , where  $\mathbf{A}_{\mathbf{a}} = \mathbf{A}_{\mathbf{r}}\Gamma$ .

In this case, all updates in  $\Gamma$  are fired simultaneously, and each element of the first array is replaced with the corresponding element of the second array. “Simultaneously” means here that the order of the array element update is unimportant, which gives the compiler designer some freedom in choosing the updating algorithm.

The *dynamic variable creation* statement is written in Pascal as  $new(p)$ , where  $p$  is a variable of type (i.e., a pointer). Its formal semantics is as follows.

Let  $\delta = (\mathbf{A}_{loc(t)}, \mathbf{l})$ , where  $\mathbf{l} \in \mathbf{B}_{Locs}$  and  $\mathbf{l} \notin \mathbf{A}_{loc}$  (sort update),  $\mathbf{A}' = \mathbf{A}\delta$ ,  $\alpha = (\mathbf{p}^{\mathbf{A}}, \mathbf{l})$ , and  $\mathbf{A}_{\alpha} = \mathbf{A}'\alpha$ . Then  $\llbracket new(p) \rrbracket^{\mathbf{A}} = \langle \mathbf{A}_{\alpha}, \perp \rangle$ .

Thus, the following actions should be undertaken when executing this statement:

- 1) allocation of a new location for the sort  $\mathbf{A}_{loc(t)}$ ;
- 2) update of the location  $p$  by the address of the new location.

If  $t$  is an array type  $array(ind, t')$ , then formally the location sort  $loc(t')$  is extended with an extra set  $L$  of locations, the injective finite mapping  $\mathbf{B}_{ind} \rightarrow L$  is created, and  $\mathbf{l}$  is initialized with this mapping. In practice, the address of an array and the address of its first element are normally the same address, and the injective finite mapping  $\mathbf{B}_{ind} \rightarrow L$  is used just for compiling the indexed variable.

Note that the formal semantics of the statement does not indicate where the new location should be placed in the object program. However, some simple reasoning can help to answer this question. Indeed, it may happen that the variable  $p$  is declared in an outer block and the dynamic variable creation statement is executed in an inner block. According to the rules of algebra transformation discussed in Section 3.2, a location dynamically created in a block should exist when the block is exited. It is clear that such a location may not be placed in the stack, otherwise it would be deleted when the block is exited. Thus, we automatically come to the conclusion that a new memory plot, normally called a “heap”, should exist in the object program. Note that we did not need to use explicitly this notion in the semantics of the language.

## 7. Function definition

In the previous sections, we assumed *some function* is of its parameters and result match the types of parameters and result in the function declaration. A real program requires a *definite function* to be associated with a function declaration. In the theory of algebraic specifications [1], they distinguish the notions of *function declaration* (or *function signature*) and *axiom* limiting the set of functions that may be associated with a given function declaration. In programming languages, the function declaration and the axiom are normally combined in a so called *function definition*, which can be formally defined as follows.

Let  $\Sigma$  be a program signature and  $(f, r, t)$  a function declaration, where  $r = u_1, \dots, u_n$  and  $u_i$  is either  $t_i$  or  $ref\ t_i$ ,  $\Delta = \{(p_1, t_1), \dots, (p_n, t_n)\}$  is a signature increment, and  $St$  a  $\Sigma'$ -block, where  $\Sigma'$  is a block signature constructed with the use of  $\Sigma$  and  $\Delta$  as described in Section 2.3. Then

$$f(p_1 : u_1, \dots, p_n : u_n) : t; St;$$

is a *function definition*. In terms of algebraic specifications, the function definition is just a combination of the function declaration  $f(u_1, \dots, u_n) : t$  and the axiom  $f(p_1, \dots, p_n) = St$ , where  $p_1, \dots, p_n$  are algebraic variables that may be used in the statement (function body)  $St$  (indeed, the requirement that  $St$  is a  $\Sigma'$ -block with the signature increment  $\Delta$  indicates that these variables may be used in the statement).

Let now  $\mathbf{A}$  be a  $\Sigma$ -algebra,  $|\mathbf{A}|$  its carrier, and  $\xi : \Delta \rightarrow |\mathbf{A}|$  a variable valuation such that  $\xi(p_i) \in \mathbf{A}_t$  if  $u_i$  is  $t_i$ , and  $\xi(p_i) \in \mathbf{A}_{loc(t)}$  if  $u_i$  is  $ref\ t_i$ . In the same manner as it is described in Section 3.2, basing on algebra  $\mathbf{A}$ , we construct the  $\Sigma'$ -algebra  $\mathbf{A}'$ , setting  $\mathbf{p}_i^{\mathbf{A}'} \dagger = \xi(\mathbf{p}_i)$ , if  $u_i$  is  $t_i$ , and  $\mathbf{p}_i^{\mathbf{A}'} = \xi(\mathbf{p}_i)$  if  $u_i$  is  $ref\ t_i$ . Note that  $p_i$  is a local variable initialized with the value of the actual parameter in the first case, and  $p_i$  is a local constant equal to the value of the actual parameter (both  $\mathbf{p}_i^{\mathbf{A}'}$  and  $\xi(p_i)$  are the same location) in the second case.

Let now

$$f^{\mathbf{A}}(\xi(p_1), \dots, \xi(p_n)) = \langle \mathbf{A}_f, \mathbf{v} \rangle \text{ and } \llbracket St \rrbracket^{\mathbf{A}'} = \langle \mathbf{A}_{st}, \mathbf{v}_{st} \rangle.$$

The program  $\mathbf{P}(\mathbf{B})$  satisfies the function definition

$$f(p_1 : u_1, \dots, p_n : u_n) : t; St;$$

iff for any  $\Sigma$ -algebra  $\mathbf{A}$  and any variable valuation  $\xi : \Delta \rightarrow |\mathbf{A}|$  it holds:

$$\mathbf{A}_f = \mathbf{A}_{st} \text{ and } \mathbf{v} = \mathbf{v}_{st}.$$

Thus, the semantics of the function definition is a direct instruction for the compiler designer to associate with a function declaration a function

equivalent to the statement used as the function body. Note that *extensional equivalence* (i.e., producing the same results for the same values of the actual parameters) is required. This means that, depending on the choice of the optimization technique, intensionally different functions may be associated with a function declaration; only the result of their execution is important.

**Fact.** *The semantics of the function definition does not depend on the signature and algebra of the block where the function is called.* Indeed, let  $f(e_1, \dots, e_n)$  be a function call in signature  $\Sigma''$  (it may be the signature of some block), where  $e_1, \dots, e_n$  are  $\Sigma''$ -expressions representing actual parameters of the function call and  $\mathbf{A}''$  a  $\Sigma''$ -algebra where the function call is interpreted. We know that the interpretation of a function call in algebra  $\mathbf{A}''$  causes the invocation of the corresponding function in algebra  $\mathbf{A}_n''$  produced as a result of evaluation of the last argument of the call. Let  $\mathbf{v}_1, \dots, \mathbf{v}_n$  be the results of evaluation of expressions  $e_1, \dots, e_n$ . Since  $\Sigma''$  does not include new data types with respect to  $\Sigma$ , there exists a  $\Sigma$ -algebra  $\mathbf{A}$  containing the elements  $\mathbf{v}_1, \dots, \mathbf{v}_n$  (recall that any set of elements may be associated with a data type) and  $\mathbf{f}^{\mathbf{A}_n''}(\mathbf{v}_1, \dots, \mathbf{v}_n) = \mathbf{f}^{\mathbf{A}}(\mathbf{v}_1, \dots, \mathbf{v}_n)$ . The block *St* is interpreted in the algebra extending  $\mathbf{A}$  and therefore its interpretation does not depend on  $\mathbf{A}_n''$ . The fact is proved<sup>1</sup>.

Thus, the semantics of the function definition clearly indicates that, when compiling the function body, we ignore the local variables of the block invoking the function and thus avoid the danger of misinterpretation of identifiers redeclared in the block.

Note another instruction for the compiler designer dictated by the semantics of the function definition. The statement *St* is interpreted in the block containing the declarations of the function formal parameters. As we already know, local block variables are allocated in the stack. Therefore, the values of the actual parameters are stored in the stack when the function is invoked, and they are automatically deleted from the stack after the execution of the function body. A value of type  $t$  is stored in the stack for a function parameter declared as  $t$ , and an address is stored for a function parameter declared as *ref t*. Passing a location as an argument to the function permits the statement *St* to update the location. Thus, the semantics of the function definition defines the way of compiling the function calls with different kinds of function parameters.

**Example.** Let the definition of our function  $f$  be as follows:

```
func f(p1: integer, p2: ref Boolean: char; begin p1 := p1 + 1;
p2 := false; x := 1.5; return 'a' end;
```

According to the semantics of the function definition, the result of the func-

<sup>1</sup>This fact was not recognized by the author when writing the paper [4], and therefore a more complex semantics of the function definition is given in that paper.

tion invocation with some arguments in the algebra  $\mathbf{A}$  should be the same as the result of the interpretation of the function body in the algebra  $\mathbf{A}$  extended with the values of these arguments. Let us have an invocation with the arguments  $v_1$  (an integer) and  $v_2$  (a location). The extension of  $\mathbf{A}$  then produces the algebra  $\mathbf{A}'$  with extra elements  $p1^{\mathbf{A}'}$  and  $p2^{\mathbf{A}'}$ , where  $p1^{\mathbf{A}'}$  is a location initialized by the value  $v_1$  and  $p2^{\mathbf{A}'} = v_2$ . The signature of this algebra enriches the signature of the algebra  $\mathbf{A}$  by the elements  $(p1, integer)$  and  $(p2, Boolean)$  that are used in the process of the function body interpretation. The interpretation will add 1 to the content of the local variable  $p1$ , store `false` in the location  $p2^{\mathbf{A}'}$  (recall that this is the location of the array element  $m[1]$ ), and update the content of the global variable  $x$  (recall again that the local variable  $x$  is forgotten because of the reduct). Following the rules of block interpretation, when the function body is exited, its algebra is reduced to a new algebra of signature  $\Sigma$  and therefore the locations of the formal parameters  $p1^{\mathbf{A}'}$  and  $p2^{\mathbf{A}'}$  are deleted from the stack. As a result, the modification of the parameter  $p1$  will not influence the program state while the modification of the parameter  $p2$  will result in modification of the location  $v_2$ . In contrast to  $\mathbf{A}$ , in the new  $\Sigma$ -algebra,  $\mathbf{A}$ , we have:  $v_2 \uparrow = \text{false}$  and  $x \uparrow = 1.5$ . Since the interpretation of the function body produces the pair  $\langle \mathbf{A}_1, a' \rangle$ , the function  $f$  must also produce the same pair (although an optimizer may simplify the code by deleting the unneeded first statement).

## 8. Program specification

We can also introduce a notion of constant definition in addition to function definition. So, if  $(c, t)$  is a constant declaration and  $e$  a  $\Sigma_{dat}$ -expression (constant expression) of type  $t$ , then  $c : t = e$  is a *constant definition*. A program  $P(\mathbf{B})$  satisfies a constant definition  $c : t = e$  iff for any  $\Sigma$ -algebra  $\mathbf{A}$  it holds:  $c^{\mathbf{A}} = e^{\mathbf{B}}$ .

A *program specification*  $\mathcal{PS}$  of program schema  $\mathcal{S}$  consists of the definitions of all functions and constants declared in  $\mathcal{S}$ . A program  $P(\mathbf{B})$  satisfies a  $\mathcal{PS}$  iff it satisfies the definitions of each function and each constant in  $\mathcal{PS}$ .

In terms of programming languages, a program specification is just a program text. In Pascal-like languages the program specification is accompanied by a statement, which may be regarded as a term of the program signature that should be evaluated in an algebra satisfying the program specification. In the process of compilation of a program specification, function codes are produced, the symbol table is augmented with function and variable addresses, memory location in the stack and heap are reserved, i.e., all actions needed for the construction of an object program as a model of the specification are performed.

## 9. Conclusion

It is shown in the paper that the formal semantics of an imperative programming language based on a formal model of the program can serve as an abstract model of some compiler components and algorithms. In particular, it is shown that the program signature is a template of the compiler symbol table, program algebra is a mathematical counterpart of the object program, semantics of function definition and call indicates the mechanisms of their compilation, semantics of blocks requires the stack memory organization, and semantics of dynamic variable creation requires availability of the heap. Thus, the proposed method of formal definition of the semantics of an imperative programming language both provides a way of unambiguous understanding of the language and serves as a good adviser of the compiler designer.

An interesting direction of further research is the use of this kind of semantics for model-based testing of compilers [2]. We hope to proceed with research in this direction.

## References

- [1] Asteziano E., Kreowski H.-J., Krieg-Brückner B. (Eds.). Algebraic Foundations of System Specification. — Springer, 1999.
- [2] El-Far I., Whittaker J. A. Model-based software testing // Encyclopedia of Software Engineering. — 2001. — Vol. 1. — P. 825–837.
- [3] Goguen J. A., Malcolm G. Algebraic Semantics of Imperative Programs. — The MIT Press, 1996.
- [4] Zamulin A.V. A State-based Semantics of a Pascal-like Language // Proc. Intern. Workshop on Program Understanding. — Novosibirsk, 2003. — P. 47–58.