

# Polymorphic syntax definition

Eelco Visser \*

*Programming Research Group, Department of Computer Science, University of Amsterdam,  
Kruislaan 403, NL-1098 SJ Amsterdam, Netherlands*

---

## Abstract

Context-free grammars are used in several algebraic specification formalisms instead of first-order signatures for the definition of the structure of algebras, because grammars provide better notation than signatures. The rigidity of these first-order structures enforces a choice between strongly typed structures with little genericity or generic operations over untyped structures. In two-level signatures level 1 defines the algebra of types used at level 0 providing the possibility to define polymorphic abstract data types. Two-level grammars are the grammatical counterpart of two-level signatures. This paper discusses the correspondence between context-free grammars and first-order signatures, the extension of this correspondence to two-level grammars and signatures, examples of the usage of two-level grammars for polymorphic syntax definition, a restriction of the class of two-level grammars for which the parsing problem is decidable, a parsing algorithm that yields a minimal and finite set of most general parse trees for this class of grammars, and a proof of its correctness. © 1998 Published by Elsevier Science B.V. All rights reserved

*Keywords:* Algebraic specification; Grammar formalisms; Polymorphism; Syntax definition; Two-level grammars

---

## 1. Introduction

In the algebraic approach to programming language specification, languages are considered as algebras. A sentence, program or expression in a language is an object of its algebra. The constructs for composition of expressions from smaller expressions and the operations that interpret, translate, transform or analyze expressions are the operations of the algebra. Algebraic specifications describe algebras by means of a finite structure that describes the sorts of the algebra, its operations and the relations between the operations. Any algebra that has the structure prescribed by the specification and that satisfies its relations is a model of the specification. Therefore, a specification always describes a class of algebras instead of precisely the intended algebra. There are many formalisms for algebraic specification. Depending on the expressive power of a

---

\* E-mail: visser@acm.org.

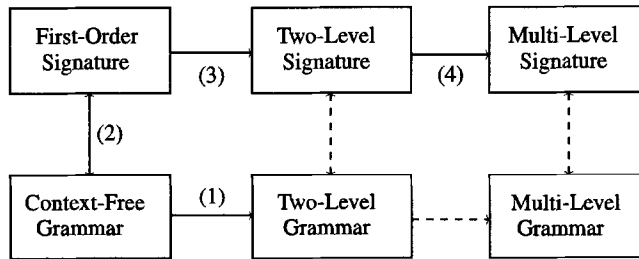
formalism the class of algebras described by a specification can be narrowed down to the intended algebra. First-order algebraic specifications consist of a first-order signature and a set of equations over the terms generated by the signature. A first-order signature consists of a finite set of sorts and a finite number of operations over those sorts.

Grammars describe languages by means of a finite structure that describes the syntactic categories of a language and the sentences of its categories. Context-free grammars and first-order signatures generate the same class of algebras. Parse trees or abstract syntax trees can be considered as terms over a signature and the language of terms over a signature can be described by a context-free grammar [6, 7]. This correspondence is exploited in several algebraic specification formalisms by allowing the use of signatures with mix-fix operators [1, 5] or even arbitrary context-free grammars [10] instead of just prefix function signatures. This provides concrete notation for functions and constructors in data type specifications and it enables definition of operations on programming languages directly in their syntactic constructs.

The rigidity of first-order signatures and context-free grammars makes it difficult to generically describe properties of an algebra. For example, an algebra with lists of integers and lists of strings can be specified with a first-order signature by declaring a sort LI (list of integers) and a sort LS (list of strings) and by defining operations like the empty list, cons, head, tail and concatenation on both sorts. However, if these list sorts have the same properties independent of the contents of the lists for some operations, this cannot be expressed in a first-order specification. Similarly, if for both list sorts an operation exists that applies a function to each element of a list, this cannot be expressed in a generic way in a first-order specification. This lack of genericity makes it difficult to develop libraries with specifications of common data types and generic language constructs.

A higher type algebra [15] is an algebra with an algebraic structure imposed on the set of sorts, i.e., the set of sorts is itself an algebra with operations. These sort operators are interpreted as functions from collections of carrier sets to collections of carrier sets. For instance, the sorts LI and LS above can be seen as sorts constructed from the sorts I (integer) and S (string) by the sort operator L that constructs the sort of sequences of integers and strings, respectively. In such algebras more generic statements about (classes of) objects and operations of the algebra can be made. For example, one can say that, for an arbitrary sort  $x$ , the tail function is a function from  $Lx$  to  $Lx$  that yields the argument sequence without its first element, where we abstract from the fact that  $x$  is equal to I or S. One could say that higher type algebras provide a higher resolution in the sort space of an algebra. Algebraic specifications in higher types [9, 14, 18, 21, 26] describe higher type algebras by means of two (or more) levels of signatures. Each level specifies the sort operations for the next level, i.e., the terms over the signature at level  $i+1$  are the sort expressions of the signature at level  $i$ . Sort expressions with variables are polymorphic sorts that describe all sorts obtained by substituting sorts for the variables. Polymorphic sorts are used to specify polymorphic functions that uniformly apply to many sorts.

In this paper we discuss polymorphic syntax definition by means of context-free and two-level grammars. We argue that the grammatical counterpart of algebraic specifications with two-levels are two-level grammars. This correspondence can be extended to multi-level signatures leading to multi-level grammars. The connections between the various formalisms are summarized by the following diagram:



where we refer to the following literature: (1) van Wijngaarden et al. [29], Pereira and Warren [19], (2) Hatcher and Rus [7], Goguen et al. [6], Futatsugi et al. [5], Heering et al. [10], (3) Poigné [21], Meinke [14, 15], Hearn and Meinke [9], (4) Hearn [8], Visser [26].

The rest of this paper is structured as follows. Section 2 contains a review of first-order signatures, context-free grammars and their correspondence and gives some examples of data type specification with context-free grammars. Section 3 defines two-level grammars and the parsing problem for two-level grammars. Section 4 illustrates how two-level grammars can be used for polymorphic syntax definition. Section 5 discusses several properties of two-level grammars including a characterization of a large class of grammars for which the parsing problem is decidable, although membership of the class is undecidable. Section 6 defines a parsing algorithm, with a correctness proof, for this class of two-level grammars that yields for each string a minimal and finite representation of the set of all parse trees for the string. Section 7 discusses related grammar formalisms and type systems and Section 8 concludes the paper.

## 2. Signatures and grammars

In this section we review many-sorted algebras, context-free grammars, the correspondence between first-order signatures and context-free grammars and the use of context-free grammars in the algebraic specification of languages and data types.

### 2.1. Many-sorted algebra

Many-sorted algebras or  $\Sigma$ -algebras were introduced by Higgins [11] as a generalization of the theory of abstract algebra. Here we give the basic constructs needed in this paper. For a further introduction to the theory of universal algebra see for instance

Meinke and Tucker [16], who also give several example applications. A note on notation: We will frequently use the notion of a family, which is a collection of sets indexed by some, finite or infinite, index set. If  $F$  is a family indexed by  $I$ , we denote by  $F(i)$  the set at index  $i \in I$  and write  $F = (F(i) \mid i \in I)$ . If we want to indicate that  $x$  is an element of some  $F(i)$ , we loosely write  $x \in F$  identifying  $F$  with  $\bigcup_{i \in I} F(i)$ .

**Definition 2.1 (Signature).** A many-sorted signature  $\Sigma$  is a pair  $\langle S, F \rangle$  where  $S = S(\Sigma) \subseteq \mathbf{S}$  is a set of sort names and  $F = F(\Sigma) \subseteq \mathbf{O} \times S(\Sigma)^+$  a set of function declarations (with  $\mathbf{S}$  and  $\mathbf{O}$  some sets of sort names and operation names, respectively). We write  $f : \tau_1 \times \cdots \times \tau_n \rightarrow \tau_0$  if  $\langle f, \tau_1 \dots \tau_n \tau_0 \rangle \in F(\Sigma)$ .  $\Sigma \cup V$  is the extension of a signature  $\Sigma$  with a  $S(\Sigma)$ -indexed family of sets of variables  $V$ . We write  $x : \tau$  if  $x \in V(\tau)$ . The class of all signatures is denoted by SIG.

**Definition 2.2 ( $\Sigma$ -Algebra).** A  $\Sigma$ -algebra  $\mathcal{A}$  is an  $S(\Sigma)$ -indexed family of carrier sets  $A(\tau)$  and an assignment of each  $f : \tau_1 \times \cdots \times \tau_n \rightarrow \tau_0$  in  $F(\Sigma)$  to an  $\mathcal{A}$  function  $f_{\mathcal{A}} : A(\tau_1) \times \cdots \times A(\tau_n) \rightarrow A(\tau_0)$  such that  $f_{\mathcal{A}}(a_1, \dots, a_n) \in A(\tau_0)$  if  $a_i \in A(\tau_i)$  ( $1 \leq i \leq n$ ).  $\text{Alg}(\Sigma)$  denotes the collection of all  $\Sigma$ -algebras.

An equational many-sorted algebraic specification consists of a signature and a set of equations that define the relations between objects of the algebras described by the specification. Note that the theory of universal algebra does not limit algebras to have finitely many operations or sorts, but that an algebraic specification must be a finite structure. The following example illustrates the definitions above. We use the keywords `sorts`, `functions` and `variables` to indicate the declaration of  $S(\Sigma)$ ,  $F(\Sigma)$  and  $V$ , respectively. We write `#` for  $\times$  and `->` for  $\rightarrow$ . Furthermore, we make use of modular specifications consisting of modules that can import other modules, where a module with imports denotes the pointwise union of the imported and importing specification.

**Example 2.3.** The following is an example of a first-order algebraic specification of the algebra of natural numbers.

```
module naturals
  sorts nat;
  functions
  zero : nat;
  succ : nat -> nat;
  add  : nat # nat -> nat;
variables
  I, J : nat;
equations
  add(zero, I)      = I;
  add(succ(I), J) = succ(add(I, J))
```

**Definition 2.4 (Terms).** The  $S(\Sigma)$ -indexed family  $\text{Tree}(\Sigma)$  of well-formed terms (or trees) over signature  $\Sigma$  is defined by the inference rules below such that  $t \in \text{Tree}(\Sigma)(\tau)$

iff  $\Sigma \vdash t : \tau$ .

$$\frac{x \in V(\tau)}{\Sigma \cup V \vdash x : \tau} \quad (\text{Var1})$$

$$\frac{f : \tau_1 \times \cdots \times \tau_n \rightarrow \tau \in F(\Sigma), \Sigma \vdash t_i : \tau_i (1 \leq i \leq n)}{\Sigma \vdash f(t_1, \dots, t_n) : \tau} \quad (\text{App1})$$

**Definition 2.5 (Homomorphism).** A  $\Sigma$ -homomorphism  $h : \mathcal{A} \rightarrow \mathcal{B}$  is an  $S(\Sigma)$ -indexed family of functions  $h_\tau$  such that for any  $f : \tau_1 \times \cdots \times \tau_n \rightarrow \tau \in F(\Sigma)$ ,  $h_\tau(f_{\mathcal{A}}(a_1, \dots, a_n)) = f_{\mathcal{B}}(h_{\tau_1}(a_1), \dots, h_{\tau_n}(a_n))$ . A  $\Sigma$ -algebra  $\mathcal{A}$  is *initial* in  $\text{Alg}(\Sigma)$  if for any  $\mathcal{B} \in \text{Alg}(\Sigma)$  there is a unique homomorphism from  $\mathcal{A}$  to  $\mathcal{B}$ .

Because there is a unique homomorphism  $h_{\mathcal{A}} : \text{Tree}(\Sigma) \rightarrow \mathcal{A}$  for any  $\mathcal{A} \in \text{Alg}(\Sigma)$ , i.e.,  $h_{\mathcal{A}}(f(t_1, \dots, t_n)) = f_{\mathcal{A}}(h_{\mathcal{A}}(t_1), \dots, h_{\mathcal{A}}(t_n))$ , we have

**Proposition 2.6.**  $\text{Tree}(\Sigma)$  is an initial algebra in  $\text{Alg}(\Sigma)$ .

**Definition 2.7 (Substitution).** A *substitution*  $\sigma : V \rightarrow \text{Tree}(\Sigma \cup V)$  is a  $S(\Sigma)$ -indexed function mapping variables to terms. The function  $\bar{\sigma} : \text{Tree}(\Sigma \cup V) \rightarrow \text{Tree}(\Sigma \cup V)$  is the homomorphic extension of a substitution  $\sigma$  that replaces all variables in a term by their  $\sigma$  images. A term  $t$  is an *instance* of term  $t'$  or  $t'$  is *more general* than  $t$ , written as  $t' \geq t$ , if there is some substitution  $\sigma$  such that  $t = \bar{\sigma}(t')$ . A term  $t$  is *strictly more general* than  $t'$ ,  $t > t'$ , if  $t \geq t'$  and not  $t' \geq t$ . In case  $t' \geq t$  we also say that  $t$  *matches*  $t'$  and that  $\sigma$  is the *match*. A substitution  $\sigma$  is a *unifier* for two terms  $t$  and  $t'$  if  $\bar{\sigma}(t) = \bar{\sigma}(t')$ . A unifier  $\sigma$  is a *most general unifier* for  $t$  and  $t'$  if for each unifier  $\sigma'$  we have that  $\bar{\sigma}(t) \geq \bar{\sigma}'(t)$ . A substitution  $\sigma$  is a *renaming* of  $t$  if  $\bar{\sigma}(t) \geq t$ , i.e., if  $\text{range}(\sigma | \text{vars}(t)) \subseteq V$ , with  $\text{vars}(t)$  the set of variables in  $t$ . Two terms  $t$  and  $t'$  are equal up to renaming of variables ( $t \doteq t'$ ) if there is a renaming  $\sigma$  such that  $\sigma(t) = t'$ .

The  $S(\Sigma)$ -indexed family of equations of an algebraic specification with signature  $\Sigma$  is a subfamily of the family  $\text{Eq}(\Sigma \cup V)$  such that  $\text{Eq}(\Sigma \cup V)(\tau) \subseteq \text{Tree}(\Sigma \cup V)(\tau)^2$ . A  $\Sigma$ -algebra  $\mathcal{A}$  satisfies an equation  $t_1 = t_2$ , if for any substitution  $\sigma$ ,  $h \circ \bar{\sigma}(t_1) = h \circ \bar{\sigma}(t_2)$  in  $\mathcal{A}$ , where  $h$  is the unique homomorphism  $h : \text{Tree}(\Sigma) \rightarrow \mathcal{A}$ .

## 2.2. Context-free grammars

Context-free grammars can be used to define languages, i.e., sets of strings and analyses of strings in the form of parse trees. The structure of parse trees corresponds to the structure of terms over a signature as we shall see in the next subsection. However, grammars provide more flexible notation for terms.

**Definition 2.8 (Grammar).** A *context-free grammar*  $\mathcal{G}$  is a triple  $\langle S, L, P \rangle$  with  $S = \text{Ss}(\mathcal{G})$  a finite set of *sort symbols* or *nonterminals*,  $L = \text{Sl}(\mathcal{G})$  a finite set of *literals* or *terminals*, with  $\text{Ss}(\mathcal{G}) \cap \text{Sl}(\mathcal{G}) = \emptyset$  and  $\text{S}(\mathcal{G}) = \text{Ss}(\mathcal{G}) \cup \text{Sl}(\mathcal{G})$  the set of *symbols* of  $\mathcal{G}$ , and  $P = \text{P}(\mathcal{G}) \subseteq \text{S}(\mathcal{G})^* \times \text{Ss}(\mathcal{G})$  a finite set of *productions*. We write  $\alpha \rightarrow \tau$  for a

production  $\langle \alpha, \tau \rangle \in P(\mathcal{G})$ .  $\mathcal{G} \cup V$  is the extension of a grammar with variables. We write  $x \rightarrow \tau$  if  $x \in V(\tau)$ . The class of all context-free grammars is denoted by CFG.

Observe that productions are reversed in order to make them look like function declarations in a signature – conventionally a production  $\alpha \rightarrow \tau$  is written as  $\tau \rightarrow \alpha$  or  $\tau ::= \alpha$ . Also note that in the conventional definition of context-free grammars a single symbol has the role of start symbol from which all sentences of the grammar are generated. In the definition above all sort symbols are start symbols. Rus and Jones [22] make a distinction between context-free grammars that have a single start symbol or *axiom*, *algebraic grammars* that have all nonterminals as start symbol, and *multi-axiom grammars* with a subset of the nonterminals as start symbol. In that terminology our grammars might more appropriately be called algebraic grammars. However, in our definition of language generated by a grammar (below), we distinguish the sets generated by each nonterminal, whereas in the definition of Rus and Jones [22] the language of a grammar is the union of all strings generated by all axioms, weakening the expressive power of the formalism. With Goguen et al. [6] we stick with the familiar ‘context-free grammar’.

As concrete syntax for grammars in examples we adopt the style of the syntax definition formalism SDF [10]. The keywords `sorts`, `syntax` and `variables` indicate the sets of sort symbols, context-free productions and variables declarations, respectively. Strings of characters between double quotes represent the literals of the grammar and identifiers are used as sort symbols. The sort symbols are explicitly declared in the `sorts` section, whereas literals are implicitly declared by their usage in productions. Grammars can be divided in modules and modules can import other modules. A module with imports denotes the pointwise union of the imported and importing grammar.

**Example 2.9.** The following specification uses a context-free grammar as signature in the specification of successor naturals. This specification is similar to the specification in Example 2.3, but in the equations we can use the more natural infix notation familiar from mathematics.

```
module naturals-cfg
  sorts nat;
  syntax
    "0"          -> nat;
    "s" nat      -> nat;
    nat "+" nat -> nat {left};
    "(" nat ")" -> nat {bracket};
  variables
    "I" -> nat; "J" -> nat;
  equations
    0 + I = I;
    s(I) + J = s(I + J)
```

The attributes attached to the productions are meant for disambiguation. The attribute `left` indicate the left associativity of the addition function and the attribute `bracket` indicates that parentheses around a natural number behave as the identity function. Disambiguation will be further discussed below.

**Definition 2.10** (*Parse Trees*). The  $S(\mathcal{G})$ -indexed family  $\text{Tree}(\mathcal{G})$  of *parse trees* over grammar  $\mathcal{G}$  is defined by the inference rules below such that  $t \in \text{Tree}(\mathcal{G})(\tau)$  iff  $\mathcal{G} \vdash t : \tau$ .

$$\begin{array}{c} \frac{L \in \text{Sl}(\mathcal{G})}{\mathcal{G} \vdash L : L} \quad (\text{Lit2}) \\[10pt] \frac{x \in V(\tau)}{\mathcal{G} \cup V \vdash \text{var}(x, \tau) : \tau} \quad (\text{Var2}) \\[10pt] \frac{\tau_1 \dots \tau_n \rightarrow \tau \in P(\mathcal{G}), \mathcal{G} \vdash t_i : \tau_i (1 \leq i \leq n)}{\mathcal{G} \vdash \text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [t_1, \dots, t_n]) : \tau} \quad (\text{App2}) \end{array}$$

**Example 2.11.** As an example of this inference relation consider the following parse tree for the sentence  $0 + I$  over the grammar of Example 2.9.

$\text{app}(\text{nat } "+" \text{ nat} \rightarrow \text{nat}, [\text{app}("0" \rightarrow \text{nat}, ["0"]), "+", \text{var}("I", \text{nat})])$

Rule (App2) defines the construction of *application* tree nodes for productions of a grammar. Observe that the complete production is used as label in such application nodes.

Because the structure of parse trees is different from terms over a signature, we redefine the notion of substitution.

**Definition 2.12** (*Substitution*). A substitution  $\sigma : (V \times S(\mathcal{G})) \rightarrow \text{Tree}(\mathcal{G} \cup V)$  is a  $S(\mathcal{G})$ -indexed family of functions mapping variables to trees. The extension  $\bar{\sigma}$  of  $\sigma$  to trees is defined as

$$\begin{aligned} \bar{\sigma}(L) &= L \\ \bar{\sigma}(\tau)(\text{var}(x, \tau)) &= \sigma(\text{var}(x, \tau)), \\ \bar{\sigma}(\tau)(\text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [t_1, \dots, t_n])) &= \text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [\bar{\sigma}(\tau_1)(t_1), \dots, \bar{\sigma}(\tau_n)(t_n)]). \end{aligned}$$

All other notions defined in Definition 2.7 are defined in the same way for parse trees.

**Definition 2.13** (*Language*). The *language*  $L(\mathcal{G})$  generated by a context-free grammar  $\mathcal{G}$  is the  $S(\mathcal{G})$ -indexed family of strings such that  $L(\mathcal{G})(\tau) = \text{yield}(\text{Tree}(\mathcal{G})(\tau))$ , where the function  $\text{yield} : \text{Tree}(\mathcal{G} \cup V) \rightarrow (S(\mathcal{G}) \cup V)^*$  is defined by

$$\begin{aligned} \text{yield}(L) &= L, \\ \text{yield}(\text{var}(x, \tau)) &= x, \\ \text{yield}(\text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [t_1, \dots, t_n])) &= \text{yield}(t_1) \dots \text{yield}(t_n) \end{aligned}$$

and applied to a set of trees denotes the pointwise extension to sets.

**Definition 2.14 (Parsing).** A parser for a context-free grammar  $\mathcal{G}$  is a function  $\Pi(\mathcal{G}) : S(\mathcal{G})^* \rightarrow \mathcal{P}(\text{Tree}(\mathcal{G}))$  that maps a string of symbols to a subfamily of  $\text{Tree}(\mathcal{G})$  such that

$$\Pi(\mathcal{G})(w)(\tau) = \{t \in \text{Tree}(\mathcal{G})(\tau) \mid \text{yield}(t) = w\}.$$

A recognizer is a predicate  $\in L(\mathcal{G})$  that decides whether a string is in the language generated by  $\mathcal{G}$  or more specifically a predicate  $\in L(\mathcal{G})(\tau)$  that decides whether a string is in the language generated by sort symbol  $\tau$ .

### 2.3. Correspondence of signatures and grammars

There is a correspondence between the trees generated by first-order signatures and context-free grammars such that grammars can be used to describe the structure of algebras [6, 7, 10].

**Proposition 2.15.** *There are mappings  $\text{grm} : \text{SIG} \rightarrow \text{CFG}$  and  $\text{sig} : \text{CFG} \rightarrow \text{SIG}$  such that  $\text{Tree}(\text{grm}(\Sigma)) \cong \text{Tree}(\Sigma)$  and  $\text{Tree}(\text{sig}(\mathcal{G})) \cong \text{Tree}(\mathcal{G})$ .*

**Proof.** Define  $\text{grm}$  such that for a signature  $\Sigma$  a grammar is constructed that expresses the syntax of terms over a signature  $\Sigma$  by taking as nonterminals the sorts of  $\Sigma$  and as literals the operator symbols of  $\Sigma$ , parentheses and commas.

$$\text{Ss}(\text{grm}(\Sigma)) = \text{S}(\Sigma),$$

$$\text{Sl}(\text{grm}(\Sigma)) = \{ "f" \mid f : \tau_1 \times \dots \times \tau_n \rightarrow \tau_0 \in \text{F}(\Sigma) \} \cup \{ "(", ")", ",", " " \},$$

$$\text{P}(\text{grm}(\Sigma)) = \{ "f" \ " (" \tau_1 ", " \dots ", " \tau_n ") " \rightarrow \tau_0 \mid f : \tau_1 \times \dots \times \tau_n \rightarrow \tau_0 \in \text{F}(\Sigma) \}.$$

Now we can translate terms over  $\Sigma$  to parse trees over  $\text{grm}(\Sigma)$  by means of the function  $i_{\text{grm}} : \text{Tree}(\Sigma) \rightarrow \text{Tree}(\text{grm}(\Sigma))$  as follows:

$$i_{\text{grm}}(\tau)(f(t_1, \dots, t_n)) = \text{app}("f" \ " (" \tau_1 ", " \dots ", " \tau_n ") " \rightarrow \tau_0,$$

$$["f" \ " (" i_{\text{grm}}(\tau_1)(t_1) ", " \dots ", " i_{\text{grm}}(\tau_n)(t_n) ") " ])$$

for each  $f : \tau_1 \times \dots \times \tau_n \rightarrow \tau \in \text{P}(\Sigma)$ . Define  $\text{sig}$  such that a grammar is translated to a signature in which the productions of the grammar have the role of function names.

$$\text{S}(\text{sig}(\mathcal{G})) = \text{S}(\mathcal{G}),$$

$$\text{F}(\text{sig}(\mathcal{G})) = \{ " \alpha \rightarrow \tau " : \tau_1 \times \dots \times \tau_n \rightarrow \tau \mid \alpha \rightarrow \tau \in \text{P}(\mathcal{G}), \alpha = \tau_1 \dots \tau_n \}$$

$$\cup \{ "L" : \rightarrow L \mid L \in \text{Sl}(\mathcal{G}) \}.$$



Now we can translate parse trees to terms by means of the function  $i_{\text{sig}} : \text{Tree}(\mathcal{G}) \rightarrow \text{Tree}(\text{sig}(\mathcal{G}))$  as follows:

$$i_{\text{sig}}(L)(L) = L(),$$

$$i_{\text{sig}}(\tau)(\text{var}(x, \tau)) = x,$$

$$i_{\text{sig}}(\tau)(\text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [t_1, \dots, t_n])) = "\tau_1 \dots \tau_n \rightarrow \tau"(i_{\text{sig}}(\tau_1)(t_1), \dots, i_{\text{sig}}(\tau_n)(t_n)).$$

It is clear that  $i_{\text{grm}}$  and  $i_{\text{sig}}$  are isomorphisms. (Note that  $\text{grm}$  and  $\text{sig}$  are not isomorphisms from  $\text{SIG}$  to  $\text{CFG}$  and vice versa:  $\Sigma \neq \text{sig}(\text{grm}(\Sigma))$  and  $\mathcal{G} \neq \text{grm}(\text{sig}(\mathcal{G}))$ .)  $\square$

The following proposition tells us that we can use context-free grammars as many-sorted algebraic signatures, where productions play the role both of function symbol and type declaration. We can thus speak of the class of algebras  $\text{Alg}(\mathcal{G})$  generated by a context-free grammar  $\mathcal{G}$ , where the grammar symbols are interpreted as carrier sets and productions as algebraic operations. It is clear that the family  $\text{Tree}(\mathcal{G})$  of parse trees over  $\mathcal{G}$  is an initial algebra in  $\text{Alg}(\mathcal{G})$ . The language  $L(\mathcal{G})$  is also an element of  $\text{Alg}(\mathcal{G})$ , with yield as the unique homomorphism from  $\text{Tree}(\mathcal{G}) \rightarrow L(\mathcal{G})$ . However,  $L(\mathcal{G})$  is not necessarily initial in  $\text{Alg}(\mathcal{G})$ . A context-free grammar is ambiguous if there is some string  $w \in L(\mathcal{G})$  for which more than one parse tree exists.

**Proposition 2.16.**  $L(\mathcal{G})$  is initial in  $\text{Alg}(\mathcal{G})$  iff  $\mathcal{G}$  is unambiguous.

For if  $\mathcal{G}$  is ambiguous, yield is not injective, hence not an isomorphism. This entails that algebraic properties do not apply to the strings used to denote trees. For example, in a grammar of arithmetic expressions with the production  $e \rightarrow e - e$ , the composition of the strings  $x$ ,  $-$  and  $y - z$  does not correspond with the composition of their trees, i.e.,  $x - (y - z)$ , but with  $(x - y) - z$ , which usually has a different semantic interpretation. We could require the use of unambiguous grammars. However, it is undecidable whether a context-free grammar is ambiguous. There are decidable subclasses of CFG, e.g., the  $\text{LR}(k)$  grammars, that are unambiguous, but these classes are much more restrictive than the class of unambiguous grammars and, moreover, not closed under union of grammars, which is a handicap when developing modular specifications. Furthermore, to disambiguate a grammar it is often necessary to introduce new sort symbols and to restrict the possibility to compose expressions.

In the work of Klint and Visser [12] a method for disambiguation of context-free grammars by means of *disambiguation filters* is proposed. A filter  $\mathcal{F}(\mathcal{G})$  selects a subset from the parse trees for a string, i.e.,  $\mathcal{F}(\mathcal{G})(\Pi(\mathcal{G})(w)) \subseteq \Pi(\mathcal{G})(w)$ . A filter is completely disambiguating if for each string  $w$ ,  $|\mathcal{F}(\mathcal{G})(\Pi(\mathcal{G})(w))| \leq 1$ .

**Proposition 2.17.** If a filter  $\mathcal{F}(\mathcal{G})$  is completely disambiguating, then  $\mathcal{F}(\mathcal{G}) \circ \Pi(\mathcal{G})$  is an injection  $L(\mathcal{G}) \rightarrow \text{Tree}(\mathcal{G})$ .

However, the trees that are not selected by the filter become unreachable with this method, i.e.,  $\mathcal{F}(\mathcal{G}) \circ \Pi(\mathcal{G})$  is not surjective. A solution to this problem is to try to add bracket productions, which are interpreted as identity functions, to the grammar such that all trees become reachable.

**Proposition 2.18.** *If  $\text{Tree}(\mathcal{G}) \cong \text{Tree}(\mathcal{G} \cup \mathcal{G}_{br}) / \equiv_{br}$  and  $\mathcal{F}(\mathcal{G} \cup \mathcal{G}_{br})$  is completely disambiguating, then  $L(\mathcal{G}) \cong \text{Tree}(\Sigma)$ .*

For a further discussion of this topic we refer to the work of Klint and Visser [12]. In the sequel we will assume that we are dealing with such grammars that we can use strings to denote trees. In examples we use a simple method for disambiguation by priority and associativity declarations. For instance, in the grammar of naturals above we used the production attribute `left` to declare the addition operator as left associative. Furthermore, the `bracket` attribute declares the production `"( nat )" -> nat {bracket}` as the identity function on natural numbers and makes all trees in  $\text{Tree}(\text{nat})$  reachable by means of strings.

#### 2.4. Data type specification

By means of grammars as signatures we have a flexible framework for syntax definition in the algebraic specification of data types, for example, the typical stack constructors might be defined as `"[]" -> stack; "push" int "on" stack -> stack`. In algebraic specification of programming languages, context-free grammars can be used for instance to specify the syntax of a programming language as in `var ":"=" exp -> stat` and the syntax of operations on programs such as typecheckers `decl "|-" exp -> bool; decl "|-" stat -> bool` that characterize the well-typed expressions and statements, interpreters `"eval" "[" stat "]" "(" env ")" -> env` that interpret statements as functions from environments to environments and compilers `"trans" "[" stat "]" -> smc` that translate statements to stack machine code.

The disadvantage of first-order signatures and context-free grammars is the rigid monomorphic typing scheme. For instance, we cannot express that for each sort  $\tau$ , the sort  $\tau^*$  of sequences of  $\tau$ 's can be constructed and that for each function  $f: \tau_1 \rightarrow \tau_2 \in F(\Sigma)$  the function  $f^*: \tau_1^* \rightarrow \tau_2^*$  extends  $f$  to sequences such that  $f^*(a_1 \dots a_n) = f(a_1) \dots f(a_n)$ . The consequence is that for each special case of a generic construct such as sequences and for each instance of a generic function such as  $_*$ , a separate definition has to be given.

One solution to overcome this rigidity is to loosen the typing requirements. In the work of Visser [25] terms of typed combinatory logic are encoded as simple untyped applicative terms. In the work of Van den Brand et al. [2] a similar structure is defined for the representation of parse trees and other structured data by means of a generic format for term representation. We study a combination of the ideas from those papers. The following grammar of generic terms (aterms) is defined by Van den Brand et al. [2] to represent parse trees and abstract syntax trees over arbitrary

grammars. A term is a function symbol (afun), an application of a function symbol to a list of arguments  $F(T_2, \dots, T_n)$ , or a list of terms  $[T_1, \dots, T_n]$ . A function symbol is a literal or an identifier. Identifier function symbols have to be defined explicitly. The module `literals` that is imported in module `aterms` defines the syntax of literals, i.e., strings of characters between double quotes.

```
module atersms
  imports literals
  sorts atersms atermlist afun aterm
  syntax
    aterm                                -> atersms;
    aterm "," atersms                    -> atersms;
    "[" "]"                              -> atermlist;
    "[" atersms "]"                      -> atermlist;
    literal                              -> afun;
    afun                                 -> aterm;
    afun "(" atersms ")"                 -> aterm;
    atermlist                            -> aterm;
  variables
    "T" -> aterm; "Ts" -> atersms; "Tl" -> atersms;
```

With this term structure it is possible to define higher-order functions. For instance, the following module defines the function `*` that applies a function  $F$  to each element of a list of terms and the function `:` that adds an element to the front of a list. Functions that are passed as arguments to higher-order functions are also represented as terms. The function `@` defines the application of such symbolically represented functions to their arguments.

```
module listops
  imports atersms;
  syntax
    aterm ":" atermlist -> atermlist {right};
    aterm "*" atermlist -> aterm {right};
    "map"                -> afun;
    aterm "@" aterm       -> aterm {left};
    "(" aterm ")"         -> aterm {bracket};
  variables
    "Fun" -> aterm;
  equations
    T : []          = [T];
    T : [Ts]        = [T, Ts];
    Fun * []        = [];
    Fun * [T]       = [Fun @ T];
    Fun * [T, Ts]   = (Fun @ T) : (F * [Ts]);
    map(Fun) @ T    = Fun * T;
```

Such a definition works well as long as sensible terms are considered. However,  $([] * \text{map})$ , the empty list mapped over the function  $\text{map}$ , is also a syntactically correct term, but does not have a clear interpretation. We would rather forbid this term on the basis of some typing rule without losing the genericity of the term structure.

One application of the generic term structure of  $\text{aterms}$ , is the representation of parse trees. We add the following function symbols

```
module atrees
  imports aterms
  syntax
    "var" -> afun; "app" -> afun; "prod" -> afun; "lit" -> afun;
```

The following proposition shows how this language can be used to represent parse trees over arbitrary grammars. Observe that we use the concrete syntax of  $\text{aterms}$  to represent elements of  $\text{Tree}(\text{atrees})$ .

**Proposition 2.19.** *For any CFG  $\mathcal{G}$ , there is an injection  $\lceil \_ \rceil : \text{Tree}(\mathcal{G}) \rightarrow \text{Tree}(\text{aterms})$  such that  $\text{Tree}(\mathcal{G})$  is isomorphic with its  $\lceil \_ \rceil$  image in  $\text{aterms}$ , i.e.,  $\text{Tree}(\mathcal{G}) \cong \lceil \text{Tree}(\mathcal{G}) \rceil$ .*

**Proof.** Given some CFG  $\mathcal{G}$  first define  $\lceil \_ \rceil : S(\mathcal{G}) \rightarrow \text{Tree}(\text{atrees})$  as

$$\lceil L \rceil = \text{lit}("L")$$

$$\lceil \tau \rceil = "\tau" \text{ if } \tau \in Ss(\mathcal{G})$$

then define  $\lceil \_ \rceil : \text{Tree}(\mathcal{G}) \rightarrow \text{Tree}(\text{atrees})$  as

$$\lceil L \rceil = "L"$$

$$\lceil \text{var}(x, \tau) \rceil = \text{var}("x", \lceil \tau \rceil)$$

$$\begin{aligned} \lceil \text{app}(\tau_1 \dots \tau_n \rightarrow \tau_0, [t_1, \dots, t_n]) \rceil \\ = \text{app}(\text{prod}(\lceil \tau_1 \rceil, \dots, \lceil \tau_n \rceil, \lceil \tau_0 \rceil), \lceil [t_1, \dots, t_n] \rceil) \end{aligned}$$

Now we have  $\text{Tree}(\mathcal{G}) \cong \lceil \text{Tree}(\mathcal{G}) \rceil$ .  $\square$

As a result, any sentence in a context-free language can be represented as a string in the fixed language of  $\text{aterms}$  preserving the structure assigned to it by the context-free grammar describing the language. For example, the parse tree for the string  $s \ 0$  according to the grammar for natural numbers is translated as follows:

$$\begin{aligned} \lceil \text{app}("s \text{ nat} \rightarrow \text{nat}, ["s" \text{ app}("0" \rightarrow \text{nat}, ["0"])]]) \rceil \\ = \text{app}(\text{prod}(\lceil \text{lit}("s") \rceil, \text{"nat"}, \text{"nat"}), \\ \lceil ["s", \text{app}(\text{prod}(\lceil \text{lit}("0") \rceil, \text{"nat"}, \lceil ["0"] \rceil)] \rceil) \end{aligned}$$

The resulting string does not only have a fixed syntax, it is also self-descriptive. The grammar  $\mathcal{G}$  can be derived from the aterm that encodes a parse tree. With this encoding we can define very generic, language independent operations on parse trees like substitution, unification and searching of subtrees. Again, the disadvantage of this scheme is that there are (many) aterms that are not encodings of parse trees, e.g., "abc" ("def") is a syntactically correct aterm but is not an element of  $\text{Tree}(\mathcal{G})$  for any  $\mathcal{G}$ . Therefore, specifications and programs that manipulate aterms encoding parse trees have to type check the terms they receive and have to preserve well-formedness of the terms they process and construct.

### 3. Two-level grammars

Context-free grammars provide either a strongly typed but rigid syntactic structure or a generic but untyped structure. Two-level grammars provide a method for polymorphic syntax definition that supports definition of generic structures with type constraints. Two-level grammars have been defined in several variants after the original formulation for the definition of the syntax of Algol68 in the work of van Wijngaarden et al. [29]. Here we introduce a definition of two-level grammars that is straightforwardly formulated as two levels of context-free grammars, where level 1 defines the syntax of the nonterminals of level 0. The productions at level 0 of a two-level grammar are *production schemata* that uniformly describe sets of context-free productions in the same way that polymorphic functions in a framework like ML [17] describe collections of functions. Given the extension of context-free grammars to two-level grammars, it is straightforward to generalize two-level grammars to multi-level grammars, in the same way as multi-level specifications are defined by Visser [26]. In this paper we will restrict our attention to two-level grammars.

**Definition 3.1** (*Two-level grammar*). A *two-level grammar*  $\Gamma$  is a pair  $\langle \mathcal{G}_1, \mathcal{G}_0 \rangle$  of context-free grammars such that the sort symbols of  $\mathcal{G}_0$  are terms, possibly with variables, over  $\mathcal{G}_1$ , i.e.,  $\text{Ss}(\mathcal{G}_0) \subseteq \text{Tree}(\mathcal{G}_1 \cup V_1)$ .

The following definition gives the meaning of finite two-level grammars in terms of, possibly infinite, context-free grammars.

**Definition 3.2.** A two-level grammar  $\Gamma$  corresponds to a, possibly infinite, context-free grammar  $\llbracket \Gamma \rrbracket$  that is derived from  $\Gamma$  by taking all substitutions of symbols  $S(\llbracket \Gamma \rrbracket) = \{ \bar{\sigma}(\tau) \mid \tau \in S(\mathcal{G}_0), \sigma : V_1 \rightarrow \text{Tree}(\mathcal{G}_1 \cup V_1) \}$  and productions  $P(\llbracket \Gamma \rrbracket) = \{ \bar{\sigma}(\tau_1) \rightarrow \bar{\sigma}(\tau_2) \mid \tau_1 \rightarrow \tau_2 \in P(\mathcal{G}_0), \sigma : V_1 \rightarrow \text{Tree}(\mathcal{G}_1 \cup V_1) \}$ .

Through the translation of a two-level grammar  $\Gamma$  to a CFG  $\llbracket \Gamma \rrbracket$  we immediately have the definitions of the term algebra  $\text{Tree}(\llbracket \Gamma \rrbracket)$  and the language  $L(\llbracket \Gamma \rrbracket)$ . Another

characterization of the trees generated by a two-level grammar is given by means of inference rules in the following definition.

**Definition 3.3.** The  $S(\mathcal{G}_0)$ -indexed family  $\text{Tree}(\Gamma)$  of *parse trees over two-level grammar*  $\Gamma$  is defined by the inference rules below such that  $t \in \text{Tree}(\Gamma)(\tau)$  iff  $\Gamma \vdash t : \tau$ .

$$\frac{L \in S(\mathcal{G}_0)}{\Gamma \vdash L : L} \quad (\text{Lit3})$$

$$\frac{x \in V(\tau'), \tau' \geq \tau}{\Gamma \cup V \vdash \text{var}(x, \tau) : \tau} \quad (\text{Var3})$$

$$\frac{p \in P(\mathcal{G}_0), p \geq \tau_1 \dots \tau_n \rightarrow \tau, \Gamma \vdash t_i : \tau_i \ (1 \leq i \leq n)}{\Gamma \vdash \text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [t_1, \dots, t_n]) : \tau} \quad (\text{App3})$$

Recall from Definition 2.7 that the relation  $p \geq p'$  holds if production  $p'$  is an instance of  $p$ , i.e.,  $p$  is more general than  $p'$ .

We observe that the two ways of defining the terms generated by a two-level grammar are equivalent.

**Proposition 3.4.**  $\Gamma \vdash t : \tau$  iff  $\llbracket \Gamma \rrbracket \vdash t : \tau$

**Proof.**  $(\Rightarrow)$  by induction on  $t$ : (i) if  $t = L$  then  $\llbracket \Gamma \rrbracket \vdash L : L$  by (Lit2) (ii) if  $t = \text{var}(x, \tau)$  then  $\llbracket \Gamma \rrbracket \vdash \text{var}(x, \tau) : \tau$  by (Var2) (iii) if  $t = \text{app}(p', [t_1, \dots, t_n])$ , by induction hypothesis  $\llbracket \Gamma \rrbracket \vdash t_i : \tau_i$ , and by (App3) there is some  $p \in P(\mathcal{G}_0)$  such that  $p \geq p'$ , i.e., there is some  $\sigma$  such that  $\bar{\sigma}(p) = p'$ , but then  $p' \in P(\llbracket \Gamma \rrbracket)$ , therefore, by (App2),  $\llbracket \Gamma \rrbracket \vdash \text{app}(p', [t_1, \dots, t_n])$ .  $(\Leftarrow)$  similarly.  $\square$

**Corollary 3.5.**  $\text{Tree}(\llbracket \Gamma \rrbracket) = \text{Tree}(\Gamma)$  and  $L(\Gamma) = L(\llbracket \Gamma \rrbracket)$

**Definition 3.6 (Substitution).** A two-level substitution  $\varphi$  is a pair  $\langle \sigma_1, \sigma_0 \rangle$  of a type substitution  $\sigma_1 : (V_1 \times S(\mathcal{G}_1)) \rightarrow \text{Tree}(\mathcal{G}_1 \cup V_1)$  and an object substitution  $\sigma_0 : (V_0 \times S(\mathcal{G}_0)) \rightarrow \text{Tree}(\mathcal{G}_0 \cup V_0)$ . The extension  $\bar{\varphi}$  of  $\varphi$  to level 0 trees is defined as

$$\bar{\varphi}(L) = L,$$

$$\bar{\varphi}(\tau)(\text{var}(x, \tau)) = \sigma_0(\text{var}(x, \bar{\sigma}_1(\tau))),$$

$$\bar{\varphi}(\tau)(\text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [t_1, \dots, t_n]))$$

$$= \text{app}(\bar{\sigma}_1(\tau_1) \dots \bar{\sigma}_1(\tau_n) \rightarrow \bar{\sigma}_1(\tau), [\bar{\varphi}(\tau_1)(t_1), \dots, \bar{\varphi}(\tau_n)(t_n)]).$$

All other notions defined in Definition 2.7 are defined in the same way for parse trees. A two-level substitution with  $\sigma_0$  equal to the identity function is also denoted by  $\bar{\sigma}$ , i.e., a function that substitutes type variables throughout a term.

**Definition 3.7 (Parsing).** Given a two-level grammar  $\Gamma$  and a string  $w$  the parsing problem is to find the set of parse trees  $\Pi(\Gamma)(w)$  such that

$$\Pi(\Gamma)(w)(\tau) = \{t \mid \Gamma \vdash t : \tau \wedge \text{yield}(t) = w\}.$$

**Discussion 3.8.** According to the definition above, *trees* over level 1 are used as sort symbols in level 0. However, if we write such grammars, we want to use strings instead of trees, i.e.,  $S(\mathcal{G}_0) \subseteq L(\mathcal{G}_1 \cup V_1) \cup S$  instead of  $S(\mathcal{G}_0) \subseteq \text{Tree}(\mathcal{G}_1 \cup V_1) \cup S$ . This entails that the syntax of two level grammars is not fixed, the syntax of the symbols of level 0 is determined by level 1. To parse a two-level grammar we first have to parse level 1 with a parser for a context-free grammar formalism in order to construct a parser for level 0. Note that we use the same, SDF style, notation for productions and modules at both levels.

## 4. Examples

In this section we discuss several examples of two-level grammars. The syntax of grammars is the adaptation of the syntax of the multi-level specifications of Visser [26] to grammars, i.e., function declarations become productions. It is not our intention to explain every detail of the notation used, but we do want to illustrate the general utility of two-level grammars for specification of data types.

### 4.1. Naturals

Module `nat` defines the syntax of natural number expressions. Level 1 introduces the sort type and the type constant `nat`. The expression `nat` can then be used as sort at level 0. Consider for example the production `"s" nat -> nat` of level 0. The expression `nat` in this production, is the constant `"nat" -> type` defined at level 1.

```

module nat
level 1
  sorts type;
  syntax
    "nat" -> type;
level 0
  syntax
    "0"          -> nat;
    "s" nat      -> nat;
    nat "+" nat -> nat {left};
    "(" nat ")" -> nat {bracket};
  variables
    "I" -> nat; "J" -> nat;
  equations
    0 + I = I;
    s(I) + J = s(I + J);

```

#### 4.2. Booleans and polymorphic conditional

The grammar in module `nat` defines monomorphic syntax for natural numbers. Each production has one instance, i.e., the production itself. The following module defines the data type of Booleans. At level 1 the type constant `bool` is introduced, which is used as sort at level 0. In addition to the ordinary Boolean connectives, the module defines a polymorphic conditional for any type. The type variable `A` in the `if-then-else-fi` production can be instantiated with any type expression. The production actually denotes the set of all instantiations of this production. Furthermore, the module defines a polymorphic bracket function.

```

module bool
imports nat;
level 1
  syntax
    "bool" -> type;
  variables
    "A" -> type; "B" -> type; "C" -> type;
level 0
  syntax
    "true"                -> bool;
    "false"               -> bool;
    "not" bool            -> bool;
    bool "\"/" bool       -> bool {left};
    "if" bool "then" A "else" A "fi" -> A;
    "(" A ")"             -> A {bracket};
  priorities
    "not" bool -> bool > bool "\"/" bool -> bool
  variables
    "B" -> bool; "X" -> A; "X'" -> A;
  equations
    not true  = false;
    not false = true;
    true  \/" B = true;
    false \/" B = B;

    if true  then X else X' fi = X;
    if false then X else X' fi = X';

```

#### 4.3. Polymorphic lists

Most grammar formalisms provide a built-in notion of lists. The next example shows how such notation can be introduced with two-level grammars. Module `list` introduces type operators at level 1 denoting the type of polymorphic lists. The operators `{_ _}+` and `{_ _}*` denote the type of non-empty and possibly-empty lists with separators,



respectively. The operators " $_+$ " and " $_*$ " denote the type of nonempty and possibly-empty lists without separators, respectively. The latter two operators are defined in terms of the former two by means of the equations that define lists without separators as lists with empty separators, where empty is a sep.

At level 0 polymorphic constructor functions for these types are defined. A nonempty list of As separated by Seps is either an A or two lists concatenated by a Sep. The first equation expresses that Sep concatenation associates to the right. An  $\{A \text{ Sep}\}^*$  list is either empty or a nonempty list of As.  $\{A \text{ Sep}\}^*$ -lists can be concatenated by means of the operator  $[_ \text{ Sep } _]$ . Note that " $\wedge$ " is used as a variable to denote separators.

```

module list
imports bool;
level 1
  sorts regtype, sep;
  syntax
    "{" type sep "}" "*" -> regtype; type "*" -> regtype;
    "{" type sep "}" "+" -> regtype; type "+" -> regtype;
    "[" regtype "]" -> type; "empty" -> sep;
  variables
    "Sep" -> sep; "R" -> regtype;
  equations
    A* = {A empty}*; A+ = {A empty}+;
level 0
  syntax
    A -> {A Sep}+;
    {A Sep}+ Sep {A Sep}+ -> {A Sep}+ {right};
    {A Sep}+ -> {A Sep}*;
    {A Sep}* -> {A Sep}*;
    "[" {A Sep}* Sep {A Sep}* "]" -> {A Sep}*;
    -> empty;
    "(" R ")" -> R {bracket};
  variables
    "L" -> {A Sep}*; "Lp" -> {A Sep}+; "\wedge" -> Sep;
  equations
    (Lp1 ^ Lp2) ^ Lp3 = Lp1 ^ (Lp2 ^ Lp3);
    [ ^ L ] = L;
    [L ^ ] = L;
    [Lp1 ^ Lp2] = Lp1 ^ Lp2;

```

Observe again how expressions over the syntax defined at level 1 are used as sorts at level 0. For instance, in the production  $\{A \text{ Sep}\}^+ \rightarrow \{A \text{ Sep}\}^*$ , the syntax of the expression  $\{A \text{ Sep}\}^+$  is defined by the production  $\{" \text{ type sep } "\} \text{ "+" } \rightarrow \text{regtype}$  and by the variables  $\text{"A"} \rightarrow \text{type}$  and  $\text{"Sep"} \rightarrow \text{sep}$ .

We have introduced a new sort `regtype` at level 1 as the sort of list type constructors in order to avoid an infinite chain caused by the injection of arbitrary types in the corresponding list type. If we would have declared the list type constructor as

```
"{" type sep "}" "+" -> type
```

the production  $A \rightarrow \{A \text{ Sep}\}^+$  would give rise to the productions

```
{A Sep}+ -> {{A Sep}+ Sep}+
{{A Sep}+ Sep}+ -> {{{A Sep}+ Sep}+ Sep}+
```

etc., causing each expression to have infinitely many nonunifiable types. By introducing the new sort `regtype`, lists are not automatically embedded in types, i.e.,  $A$  does not unify with  $\{A \text{ Sep}\}^+$  because their sorts are different.

The usage of list types is illustrated in the following grammar of a fragment of an imperative language. A statement is either an assignment, a while-do loop or a list of statements separated by semicolons.

```
module while
imports list, exp;
level 1
  syntax
    "var"  -> type; "exp" -> type;
    "stat" -> type; ";"  -> sep;
level 0
  syntax
    var ":" exp          -> stat;
    "while" exp "do" stat -> stat;
    "begin" {stat ';' }* "end" -> stat;
    ";"      -> ';' ;
```

The expression  $\{\text{stat } ' ; ' \}^*$  is defined by the polymorphic productions in module `list`, which have the following instantiations.

```
stat          -> {stat ' ; ' }+;
{stat ' ; ' }+ ' ; ' {stat ' ; ' }+ -> {stat ' ; ' }+;
                                     -> {stat ' ; ' }*;
{stat ' ; ' }+ -> {stat ' ; ' }*;
```

#### 4.4. Polymorphic operations

Now that we have a polymorphic definition of list construction we can also define polymorphic functions over lists. For instance, the `length` function that computes the number of elements of a list can be polymorphically defined by the following specification:

```

module length
imports list;
level 0
  syntax
    "length" "(" {A Sep}* ")" -> nat;
  equations
    length()          = 0;
    length(X)         = s(0);
    length(Lp1 ^ Lp2) = length(Lp1) + length(Lp2);

```

#### 4.5. Higher-order functions

Another example of a type constructor is the arrow  $\Rightarrow$  of function types. A term of sort  $A \Rightarrow B$ , i.e., a function from  $A$  to  $B$ , can be applied to a term of sort  $A$  yielding a  $B$ .

```

module arrow
imports list;
level 1
  syntax
    type "=>" type -> type {right};
level 0
  syntax
    (A => B) "(" A ")" -> B;
    (A => B) "*" {A Sep}* -> {B Sep}*;
  variables
    "F" -> (A => B);

```

The higher-order function  $*$  (map) takes as arguments a function from  $A$  to  $B$  and a list of  $A$ s and applies the function to each element of the list.

```

equations
  F * ()          = ();
  F * (X)         = F(X);
  F * (Lp1 ^ Lp2) = [(F * Lp1) ^ (F * Lp2)];

```

If we want to pass functions such as `length` and `map` themselves as arguments to some higher-order function we need to define the combinators (curried versions) associated with the functions as follows:

```

syntax
  "if"          -> (bool => A => A => A);
  "flength"     -> ([{A Sep}*] => nat);
  "(*)"         -> (A => B) => ([{A Sep}*] => [{B Sep}*]);

```

```

(A  => [R']) "(" A ")" -> R';
([R] => B)    "(" R ")" -> B;
([R] => [R']) "(" R ")" -> R';

```

equations

```

if(B)(X)(Y) = if B then X else Y fi;
flength(L)  = length(L);
(*) (F)(L)  = F * L;

```

Observe the usage of the operator `[_]` that injects `regtypes` into types in order to reuse the functionality for type expressions. We added extra application operators to apply functions like `flength` to lists.

These examples illustrate how two-level grammars provide user-definable syntax for sort symbols and generic definition of polymorphic mix-fix functions and constructors over data types. More advanced examples of two- and multi-level specifications (with prefix function signatures instead of grammars) can be found in the works of Hearn and Meinke [9], Hearn [8] and Visser [26].

## 5. Properties

We have seen how two-level grammars can be used for polymorphic syntax definition in algebraic specification. To actually use two-level grammars in an executable specification formalism, it is necessary that we can parse strings over the language specified by a grammar. Unfortunately, the parsing problem for two-level grammars is in general undecidable as shown by the following theorem.

**Theorem 5.1** (Sintzoff [23]). *For every semi-Thue system  $T$  we can construct a Van Wijngaarden Grammar  $W$  such that the set  $S(T)$  generated by  $T$  is the set  $S(W)$  generated by  $W$ .*

**Corollary 5.2** (Sintzoff [23]). *Every recursively enumerable set is generated by a Van Wijngaarden grammar.*

**Corollary 5.3** (Sintzoff [23]). *The problem of determining, of a given string, whether or not it is generated by a given Van Wijngaarden grammar, is recursively unsolvable.*

Although the version of two-level grammars defined in this paper is somewhat weaker because it uses trees instead of strings as nonterminals at level 0, these constructs can be translated to our two-level grammars. From these theorems it follows that we cannot construct terminating parsers for arbitrary two-level grammars in a general way. However, for the purpose of polymorphic syntax definition we are interested only in restricted forms of the formalism.

One view on two-level grammars is that they are used to abbreviate frequently occurring patterns in context-free grammars, but that in the end we want only a finite

context-free grammar from a ground subgrammar and the appropriate instantiations of generic productions. For instance, the grammar of the programming language in module `while` gives rise to the instantiation of the list construction functions for `{stat ' ; '}`\* and to the instantiation of the list and map functions for lists of statements. This is the effect that is reached when reuse of functions is obtained by means of parameterized modules for which only finitely many instantiations are requested. Although it is clear by looking at a grammar, which instantiations of productions are needed for the implementation of a certain subgrammar, we have not yet found a syntactic characterization of productions such that such subgrammar operations are possible. A promising approach might be the extension of the layering operations of Hatcher and Rus [7] and Rus and Jones [22] to two-level grammars.

In context-free grammars empty ( $\varepsilon$ ) and chain productions are the cause of infinite ambiguities. In two-level grammars they are the cause of the undecidability of the formalism. In the rest of this section we investigate the restriction of the usage of such productions in order to achieve a subclass of the two-level grammars with a decidable parsing problem that still allows the kind of grammars as shown in Section 4.

**Definition 5.4** ( $\varepsilon$ -elimination). The conventional method for eliminating  $\varepsilon$ -productions from context-free grammars applied to two-level grammars works by adding productions to the level 0 grammar  $\mathcal{G}_0$  according to the rule

$$\frac{\alpha A \beta \rightarrow B \in \mathcal{G}_0, \varepsilon \rightarrow A' \in \mathcal{G}_0, \sigma(A') = \sigma(A)}{\sigma(\alpha \beta) \rightarrow \sigma(B) \in \mathcal{G}_0}$$

where  $\sigma$  is a most general unifier of  $A$  and  $A'$ . After no more productions can be added, all  $\varepsilon$ -productions are removed. Define  $ee(\Gamma)$  to be the result of removing  $\varepsilon$ -productions from two-level grammar  $\Gamma$  by the above procedure.

Note that  $\varepsilon$ -elimination preserves both the language and the trees generated by the grammar (if  $\varepsilon$ -trees are identified).

**Proposition 5.5.**  $L(ee(\Gamma)) = L(\Gamma)$  and  $Tree(\Gamma) \cong Tree(ee(\Gamma))$ .

Deussen [4] shows that this method can turn finite two-level grammars into infinite ones. Consider the following grammar that gives type  $\langle a^n \rangle$  to each sentence  $a^m$  with  $m \leq n$ .

```

level 1
  sorts type, list;
  syntax
    "<" list ">" -> type
    "a"          -> list
    list "a"     -> list
  variables
    "L" -> list

```

level 0

```
syntax
  <L> <a> -> <L a>;
  "a"      -> <a>;
           -> <a>;
```

If we try to eliminate the last production by substituting it in the first production we get the productions

```
<L> -> <L a>;    %% <a> can be empty
      -> <a a>;    %% <L> unifies with the rhs of -> <a>
      -> <a a a>;  %% <L> unifies with the rhs of -> <a a>
```

and all other productions of the form  $\rightarrow \langle a^n \rangle$  for  $n > 0$ . However, for many applications  $\varepsilon$ -productions can be eliminated. For instance, the production  $\rightarrow A^*$  in the list grammar in Section 4 can be eliminated by means of the procedure outlined above, resulting in a finite two-level grammar defining the same language.

In a similar fashion chain productions can be eliminated from grammars.

**Definition 5.6** (*Chain elimination*). To eliminate chain production from a two-level grammar  $\Gamma$ , first take the transitive closure of all chains in the level 0 grammar  $\mathcal{G}_0$ :

$$\frac{\tau_1 \rightarrow \tau_2 \in \mathcal{G}_0, \tau_3 \rightarrow \tau_4 \in \mathcal{G}_0, \sigma(\tau_2) = \sigma(\tau_3)}{\sigma(\tau_1) \rightarrow \sigma(\tau_4) \in \mathcal{G}_0}$$

then use chain productions as substitutions

$$\frac{\tau_1 \rightarrow \tau_2 \in \mathcal{G}_0, \alpha\tau_3\beta \rightarrow \tau_4 \in \mathcal{G}_0, \sigma(\tau_2) = \sigma(\tau_3)}{\sigma(\alpha\tau_1\beta \rightarrow \tau_4) \in \mathcal{G}_0},$$

and finally remove all chain productions from  $\mathcal{G}_0$  resulting in  $ce(\Gamma)$ .

This procedure also preserves the language and trees generated by a grammar.

**Proposition 5.7.**  $L(\Gamma) = L(ce(\Gamma))$  and  $Tree(\Gamma) \cong Tree(ce(\Gamma))$ .

Also chain elimination does not terminate for all grammars. Take for instance the grammar for lists in the previous section. If we redefine the syntax of the list operators as

```
"{" type sep "}" "*" -> type;
"{" type sep "}" "+" -> type
```

then we have that  $A$  unifies with  $\{A \text{ Sep}\}^+$  resulting in infinitely many productions

```
A                -> {A Sep}+;
{A Sep}+         -> {{A Sep}+ Sep}+;
{{A Sep}+ Sep}+ -> {{{A Sep}+ Sep}+ Sep}+;
...
```

An  $A$  is a singleton list of  $A$ s, which is a singleton list of lists of  $A$ s, etc.

So we see that  $\varepsilon$ -elimination and chain elimination will not terminate for arbitrary grammars. However, for the grammars for which it succeeds we have the following corollary from Theorem 6.8 that we will prove in the next section.

**Corollary 5.8.** *If  $\Gamma$  is a finite two-level grammar without  $\varepsilon$ - and chain productions, then the question  $w \in L(\Gamma)$  is decidable.*

The intuition behind this result is that without  $\varepsilon$ - and chain-free grammars at most  $n$  reductions can be done for a string of length  $n$ . Based on the same idea, the next definition defines a characterization of a larger class of two-level grammars for which the parsing problem is decidable.

**Definition 5.9** (*Finite chain property*). A two-level grammar  $\Gamma$  has the *finite chain property* if it is (1)  $\varepsilon$ -free, (2) its chain productions are noncyclic and have a finite transitive closure and (3) it does not contain redundant productions, where a production  $p \in P(\mathcal{G}_0)$  is redundant if there is some  $p' \in P(\mathcal{G}_0)$  such that  $p \neq p'$  and  $p \geq p'$ .

In the next section we will define a parsing algorithm for two-level grammars and prove that it is a decision procedure for membership of languages defined by finite chain two-level grammars.

On the positive side we have a subclass of the two-level grammars with a decidable parsing problem. On the negative side, membership of the class itself is undecidable.

**Proposition 5.10.** *It is undecidable whether a two-level grammar satisfies the finite chain property.*

However, decidability of the finite chain property is not essential for using two-level grammars for language specification. The situation can be compared to ambiguity of context-free grammars. Although it is undecidable whether a context-free grammar is ambiguous, it is a good formalism for defining unambiguous languages. A large class of grammars is evidently nonambiguous and for others ambiguities will turn up when working with the grammars.

The examples presented in Section 4 satisfy the finite chain property, except for the empty production  $\rightarrow \{A \text{ Sep}\}^*$  for lists. As remarked above this production is not a problematic  $\varepsilon$ -production because it can be eliminated from the grammar. In general we can follow the following procedure for determining whether a grammar has the finite chain property: (1) Try to eliminate  $\varepsilon$ -productions by the method of Definition 5.4. (2) Try to eliminate chain rules by means of the method in Definition 5.6. (3) If this terminates we know that the grammar has the finite chain property and that we can parse with it (see next section). (4) If either step (1) or step (2) takes too long, this is a hint that it does not terminate. In such cases we can inspect the list of  $\varepsilon$ -productions or chain productions added by the elimination procedures. These traces will give a clue about the productions that cause the nontermination, because these

will lead to a repetition of similar productions, as we saw in the example above. This information can be used to redesign the grammar such that it satisfies the finite chain property.

## 6. Parsing

In this section we define a parsing algorithm for finite chain two-level grammars. The parsing algorithm below is a parallel bottom-up parsing algorithm that computes all parse trees for a sentence. This procedure is similar to the Hindley–Milner-type assignment procedure used in functional languages, that assigns to each expression a single principal type [3]. The difference is that in two-level grammars strings can have more than one principal type due to ambiguities. It will turn out that for finite chain two-level grammars there are only finitely many principal types for a string. We first define a function that gives the type of a parse tree.

**Definition 6.1.** The type of a parse tree is defined as:

$$\begin{aligned} \text{type}(L) &= L \\ \text{type}(\text{var}(x, \tau)) &= \tau \\ \text{type}(\text{app}(\tau_1 \dots \tau_n \rightarrow \tau, [t_1, \dots, t_n])) &= \tau \end{aligned}$$

Next we define the data structure of parse configurations that is used in parsing.

**Definition 6.2.** A parse configuration  $(\vec{t} \bullet a_1 \dots a_n)_\Phi$  is an element of the set  $\text{Tree}(\Gamma \cup V)^* \times \text{Sl}(\Gamma)^* \times \text{Set}(V_1)$ , i.e., a triple consisting of a list of trees  $\vec{t} = t_1 \dots t_m$  (the stack), a list of literals  $a_1 \dots a_n$  (the remaining input) and a set of sort variables  $\Phi$  (the sort variables over level 1 that are used in  $\vec{t}$ ).

**Algorithm 6.3.** Define the function  $\text{parse}(\Gamma) : \text{S}(\Gamma)^* \rightarrow \text{Set}(\text{Tree}(\Gamma))$  as

$$\text{parse}(\Gamma)(w) = \{t \mid (\varepsilon \bullet w) \Rightarrow_\Gamma^* (t \bullet \varepsilon)\}$$

where  $\Rightarrow_\Gamma^*$  is the transitive closure of the one-step parse relation  $\Rightarrow_\Gamma$  on parse configurations, which is defined by the rules

$$\begin{aligned} (\vec{t} \bullet a_1 a_2 \dots a_n)_\Phi &\Rightarrow_\Gamma (\vec{t} a_1 \bullet a_2 \dots a_n)_\Phi \quad (\text{Shift}) \\ \frac{x \in V(\tau'), \tau = \bar{\rho}(\tau')}{(\vec{t} x \bullet \vec{a})_\Phi &\Rightarrow_\Gamma (\vec{t} \text{var}(x, \tau) \bullet \vec{a})_{\Phi \cup \text{vars}(\tau)}} \quad (\text{Var}) \\ \frac{p \in \text{P}(\mathcal{G}_0), \bar{\rho}(p) = \alpha \rightarrow \tau, |\alpha| = m, \text{mgu}(\alpha; \text{type}(t_1, \dots, t_m)) = \sigma}{(\vec{t} t_1 \dots t_m \bullet \vec{a})_\Phi &\Rightarrow_\Gamma (\vec{t} \bar{\sigma}(\text{app}(p, [t_1, \dots, t_m])) \bullet \vec{a})_{\Phi \cup \text{vars}(\bar{\sigma}(p))}} \quad (\text{Red}) \end{aligned}$$

where  $\rho : V_1 \rightarrow V_1$  is a renaming of sort variables occurring in  $\Phi$  such that  $\rho(\Phi) \cap \Phi = \emptyset$ . We identify configurations that are the same up to renaming of sort variables.



We now prove that the algorithm is a correct implementation of  $\Pi(\Gamma)$  for finite chain two-level grammars. We first show that the trees produced by the parser are correct parse trees.

**Lemma 6.4** (Sound).  $\forall t \in \text{Tree}(\Gamma) : t \in \text{parse}(\Gamma)(w) \Rightarrow \text{yield}(t) = w$ .

**Proof.** We first prove that if  $(\vec{t}_1 \bullet \vec{a}_1) \Rightarrow_{\Gamma} (\vec{t}_2 \bullet \vec{a}_2)$ , then  $\text{yield}(\vec{t}_1)\vec{a}_1 = \text{yield}(\vec{t}_2)\vec{a}_2$ . For (Shift) and (Var) the property clearly holds. In (Red) we see  $\text{yield}(\vec{t} \ t_1 \dots t_m)\vec{a} = \text{yield}(\vec{t} \ \bar{\sigma}(\text{app}(p, [t_1 \dots t_m])))\vec{a}$  by definition of yield and by the fact that type substitutions do not affect the yield of a tree. But then also for  $(\vec{t}_1 \bullet \vec{a}_1) \Rightarrow_{\Gamma}^* (\vec{t}_2 \bullet \vec{a}_2)$  we have  $\text{yield}(\vec{t}_1)\vec{a}_1 = \text{yield}(\vec{t}_2)\vec{a}_2$ . In particular, if  $(\varepsilon \bullet w) \Rightarrow_{\Gamma}^* (t \bullet \varepsilon)$  we have that  $w = \text{yield}(\varepsilon)w = \text{yield}(t)\varepsilon$ .  $\square$

Next we show that the parser is *complete*, in the sense that any parse tree for the sentence can be derived by instantiating one of the parse trees produced by the parser.

**Lemma 6.5** (Complete).  $\forall t \in \text{Tree}(\Gamma) : \text{yield}(t) = w \Rightarrow \exists t' \geq t : t' \in \text{parse}(\Gamma)(w)$ .

**Proof.** By induction on  $t$ : (1) if  $t = L$  then  $\text{parse}(\Gamma)(L) \ni L \geq L$ .

(2) if  $t = \text{var}(x, \tau)$ , there is some  $\tau' \geq \tau$  such that  $x \in V(\tau')$ , but then  $\text{parse}(\Gamma)(x) \ni \text{var}(x, \tau') \geq \text{var}(x, \tau)$ .

(3) If  $t = \text{app}(p, [t_1, \dots, t_n])$  with  $p = \tau_1 \dots \tau_n \rightarrow \tau$  (a) By definition of yield we have  $\text{yield}(t) = \text{yield}(t_1) \dots \text{yield}(t_n) = w_1 \dots w_n$  with  $\text{yield}(t_i) = w_i$  for  $1 \leq i \leq n$ . (b) By induction we have  $t'_i \geq t_i$  for  $1 \leq i \leq n$  – and thus  $\tau'_i \geq \tau_i$  with  $\tau'_i = \text{type}(t'_i)$  – such that  $t'_i \in \text{parse}(\Gamma)(w_i)$ . (c) By (App) there is some  $\tau''_1 \dots \tau''_n \rightarrow \tau'' = p'' \in P(\mathcal{G}_0)$  such that  $p'' \geq p$  (variables of  $p''$  and  $p$  disjunct). By (b) and (c) there is a substitution  $\sigma_0$  such that  $\bar{\sigma}_0(\tau''_1 \dots \tau''_n) = \tau_1 \dots \tau_n = \bar{\sigma}_0(\tau'_1 \dots \tau'_n)$ . Then there is also a most general unifier, say  $\sigma$ . Now take  $t' = \bar{\sigma}(\text{app}(p'', [t'_1 \dots t'_n])) \in \text{Tree}(\Gamma)$ . It is clear that  $t' \geq t$  and that  $(\varepsilon \bullet w_1 \dots w_n) \Rightarrow_{\Gamma}^* (t'_1 \dots t'_n \bullet \varepsilon) \Rightarrow_{\Gamma} (t' \bullet \varepsilon)$ .  $\square$

Next we show that the set of parse trees produced by the algorithm is *minimal* in the sense that it generates only the most general parse trees for a string.

**Lemma 6.6** (Minimal).  $\forall t, t' \in \text{parse}(\Gamma)(w) : t \geq t' \vee t' \geq t \Rightarrow t \doteq t'$ .

**Proof.** Assume that  $t, t' \in \text{parse}(\Gamma)(w)$  and that  $t \succ t'$ . Because both trees are in the set there must be sequences of configurations for their derivation. Because  $t \succ t'$ , the trees have the same structure, i.e., the configuration sequences have the same number of reductions and shifts. But also because  $t \succ t'$ , there must be some point at which the sequences diverge, i.e.,

$$(\varepsilon \bullet w) \Rightarrow_{\Gamma}^* (\vec{t}_1 \vec{t}_2 \bullet \vec{a}) \Rightarrow_{\Gamma} (\vec{t}_1 \ \bar{\sigma}_1(\text{app}(p_1, \vec{t}_2)) \bullet \vec{a}) \Rightarrow_{\Gamma}^* (t \bullet \varepsilon)$$

and

$$(\varepsilon \bullet w) \Rightarrow_{\Gamma}^* (\vec{t}_1 \vec{t}_2 \bullet \vec{a}) \Rightarrow_{\Gamma} (\vec{t}_1 \vec{\sigma}_2(\text{app}(p_2, \vec{t}_2)) \bullet \vec{a}) \Rightarrow_{\Gamma}^* (t' \bullet \varepsilon)$$

for  $t$  and  $t'$ , respectively. Because  $t > t'$  we must have  $\vec{\sigma}_1(p_1) > \vec{\sigma}_2(p_2)$ . Now we have either (1)  $p_1 = p_2$  and  $\sigma_1 > \sigma_2$ , but then  $\sigma_2$  is not a most general unifier and hence  $t' \notin \text{parse}(\Gamma)(w)$  or (2) if  $p_1 > p_2$ , then  $\Gamma$  does not satisfy the finite chain property because it has the redundant production  $p_2$ .  $\square$

Finally we prove that  $\text{parse}$  yields a finite set of parse trees, entailing that  $\text{parse}$  is effectively computable.

**Lemma 6.7** (Finite).  $|\text{parse}(\Gamma)(w)| \in N$ .

**Proof.** (1) For each configuration and each production there is at most one reduction step (Red) because there is at most one most general unifier for  $\alpha$  and  $\tau_1 \dots \tau_m$ . For each configuration there is at most one (Shift) step and one (Var) step. Therefore, the graph of the relation  $\Rightarrow_{\Gamma}$  is finitely branching.

(2) The length of configurations does not increase (no  $\varepsilon$ -productions). For any configuration  $(\vec{t} \bullet \vec{a})$ ,  $|\vec{a}|$  (Shift) steps can be done. A (Red) step with a production  $\alpha \rightarrow \tau$  such that  $|\alpha| > 1$  decreases the length of a configuration, therefore at most  $|\vec{t}|/2$  such reductions can be performed for a configuration  $(\vec{t} \bullet \vec{a})$ . By the finite chain property only finitely many chain reductions can be done, i.e., for each configuration  $(\vec{t} \bullet \vec{a})$  there is a maximal value  $n$  such that  $(\vec{t} \bullet \vec{a}) \Rightarrow_{\Gamma}^n (\vec{t}' \bullet \vec{a})$ . Therefore, the graph of the relation  $\Rightarrow_{\Gamma}$  has no infinite paths.

(3) From any configuration  $(\vec{t} \bullet \vec{a})$  only finitely many configurations are reachable. In particular, for any string  $w$  only finitely many configurations of the form  $(t \bullet \varepsilon)$  are reachable from  $(\varepsilon \bullet w)$ .  $\square$

Finally, we see that Algorithm 6.3 is a correct implementation of a parser for finite chain two-level grammars.

**Theorem 6.8** (Correct). *For any finite chain two-level grammar  $\Gamma$  and any string  $w \in \text{Sl}(\Gamma)^*$ ,  $\text{parse}(\Gamma)(w)$  is a minimal and finite set of parse trees, unique up to renaming of sort variables, that generates  $\Pi(\Gamma)(w)$ .*

**Proof.** By Lemmas 6.4 and 6.5 all and exactly the trees in  $\Pi(\Gamma)(w)$  can be derived from  $\text{parse}(\Gamma)(w)$ . By Lemma 6.7  $\text{parse}(\Gamma)(w)$  is finite and by Lemma 6.6 it is minimal.  $\square$

As a result the recognition problem for finite chain two-level grammars is decidable.

**Corollary 6.9** (Decidable). *For a finite chain two-level grammar  $\Gamma$  it is decidable whether  $w \in L(\Gamma)$  and  $w \in L(\Gamma)(\tau)$ .*

The relation  $\Rightarrow_r$  defines a tree shaped search space. Only the types of trees in the configuration matter for the rest of the process. We would like to identify configurations  $(\vec{t}_1 \bullet \vec{a})$  and  $(\vec{t}_2 \bullet \vec{a})$  for which  $|\vec{t}_1| = |\vec{t}_2|$  and  $\text{type}(\vec{t}_1) = \text{type}(\vec{t}_2)$ . This would lead to a generalization of the graph structured stack and the parse forests of Tomita [24] to parsing for two-level grammars.

## 7. Related formalisms

In the same way that context-free grammars correspond to first-order signatures, two-level grammars correspond to two-level signatures. The type system of the functional programming language ML [17] can be considered as two-level signatures in which the expressions over level 1 are single-sorted expressions of sort type. This system was used to introduce parametric polymorphic functions. Two-level signatures are discussed by Poigné [21], Möller [18] and Meinke [14]. After a two-level signature is expanded, a, possibly infinite, one-level signature results that can again be used as the specification of the sort space of a level 0 signature. In this manner the extension of signatures to two-level signatures can be generalized to signatures with three and more levels. Hearn and Meinke [9] introduce the three-level algebraic specification formalism Atlas, which is generalized by Hearn [8] to a multi-level specification formalism. The complete and formal specification of the related multi-level specification formalism MLS is presented by Visser [26]. MLS supports overloading of function symbols, which entails that a term can have infinitely many types, but only finitely many most general or principal types. This property is not respected by general two-level grammars as discussed in Section 5.

On the grammatical side, many variants of two-level grammars have been proposed in the literature for various purposes. Van Wijngaarden grammars (VWG) [29] were developed to express the syntax and semantics of Algol68. In VWGs strings, instead of trees, over level 1 are used as nonterminals (hypermotions) at level 0. This leads to the problem of grammatical unification – whether two sentential forms over a context-free grammar are unifiable by means of a substitution of nonterminals with strings – which Maluszynski [13] shows to be undecidable. The transparent two-level grammars of Maluszynski [13] are a restriction of VWGs such that grammatical unification comes down to term unification. Another restriction of VWGs are the Extended Affix Grammars (EAG) [28] that restrict the order in which the variables in nonterminals at level 0 can be instantiated.

The observation that two-level grammars are Turing equivalent sparked another development: two-level grammars as logic or functional programming languages. The (context-free) metagrammar (level 1) is used to define the syntax of language and semantic domains. The hypergrammar (level 0) is used to define the operations on the data. See for example Maluszynski [13].

Definite Clause Grammars (DCG) introduced by Pereira and Warren [19] are grammars embedded in Prolog programs. They are equivalent to two-level grammars with

a fixed level 1 equivalent to the following grammar:

```

module dcg
level 1
  sorts fun, term;
  syntax
    [a-z] [A-Za-z0-9]*      -> fun;
    fun      -> term;
    fun      "(" {term " ," }* ")" -> term;
  variables
    [A-Z] [A-Za-z0-9]* -> term;

```

that defines an untyped domain of terms that can be used as grammar symbols in level 0. These terms are then typically used at level 0 in productions such as the following from a tiny natural language grammar:

```

np(N) vp(N)      -> s;
det(N) n(N) rel(N) -> np(N);

```

Parsing of DCGs – parsing as deduction [20] – uses Prolog’s built-in resolution strategy to answer queries like  $w \in L(\mathcal{G})(s)$ . With the normal evaluation strategy of Prolog (SLD resolution) this comes down to top-down backtrack parsing. Problems with this strategy are that it cannot cope with left-recursion and that already computed answers are not reused. The tabulation strategy described by Warren [27] partially overcomes these problems. One of the problems of the latter approach is that unification in Prolog is not many-sorted, disabling solutions like that with `regtype` in Section 4.

## 8. Conclusions

Algebraic specification with first-order signatures or context-free grammars enforce a choice between strongly typed structures with little genericity or generic operations over untyped structures. Polymorphism combines genericity with typedness, making it possible to develop libraries of specifications. In this paper we have discussed how the integration of algebraic specification with user-definable syntax and polymorphism can be materialized. The extension with polymorphism of algebraic specification formalisms that use context-free grammars as signatures, e.g., OBJ or ASF+SDF, leads necessarily to formalisms with two-level grammars as signatures. Likewise, the extension with user-definable syntax of formalisms that have polymorphic signatures, including polymorphic functional and logic programming languages, leads to two-level grammars.

In two-level grammars level 1 defines the syntax of sort symbols used at level 0. Sort terms with variables are interpreted as sort schemata that can have many instantiations. Productions at level 0 with such sorts are production schemata, i.e., declarations of polymorphic functions with mix-fix syntax. Thus two-level grammars combine polymorphism with user-definable syntax, as we illustrated by means of a number of

examples of polymorphic syntax definition in data type and programming language specifications.

Although the parsing problem for context-free grammars and the type-assignment problem for two-level signatures are decidable, the parsing problem for the integration of both formalisms is undecidable if no restrictions are considered. We defined an intuitive restriction of the class of two-level grammars that results in a class of two-level grammars for which the parsing problem is decidable and for which we defined a parsing algorithm that yields a minimal and finite set of most general parse trees for each string.

## Acknowledgements

This research was supported by the Netherlands Computer Science Research Foundation (SION) with financial support from the Netherlands Organisation for Scientific Research (NWO) under grant 612-317-420: Incremental parser generation and context-dependent disambiguation, a multi-disciplinary perspective.

## References

- [1] M. Bidoit, M.-C. Gaudel, A. Mauboussin, How to make algebraic specifications more understandable: an experiment with the PLUSS specification language, *Sci. Comput. Programming* 12 (1989) 1–38.
- [2] M.G.J. Van den Brand, P. Klint, P. Olivier, E. Visser, ATerms: representing structured data for exchange between heterogeneous tools. Technical report, Programming Research Group, University of Amsterdam, 1997.
- [3] L. Damas, R. Milner, Principal type-schemes for functional programs, in: *Proc. Conf. Record 9th Ann. ACM Symp. Principles of Programming Languages*, ACM, New York, 1982, pp. 207–212.
- [4] P. Deussen, A decidability criterion for van Wijngaarden grammars, *Acta Inform.* 5 (1975) 353–375.
- [5] K. Futatsugi, J. Goguen, J.-P. Jouannaud, J. Meseguer, Principles of OBJ2, in: B. Reid (Ed.), *Proc. Conf. Record 12th Ann. ACM Symp. on Principles of Programming Languages*, ACM, New York, 1985, pp. 52–66.
- [6] J.A. Goguen, J.W. Thatcher, E.G. Wagner, J.B. Wright, Initial algebra semantics and continuous algebras, *J. ACM* 24(1) (1977) 68–95.
- [7] W.S. Hatcher, T. Rus, Context-free algebras, *J. Cybernet.* 6 (1976) 65–76.
- [8] B.M. Hearn, The design and implementation of typed languages for algebraic specification, Ph.D. Thesis, University of Wales, Swansea, 1995.
- [9] B.M. Hearn, K. Meinke, ATLAS: a typed language for algebraic specification, in: J. Heering, K. Meinke, B. Möller, T. Nipkow (Eds.), *Proc. 1st Internat. Workshop on Higher-Order, Algebra, Logic and Term Rewriting (HOA '93)*, Lecture Notes in Computer Science, vol. 816, Springer, Berlin, 1994, pp. 146–168.
- [10] J. Heering, P.R.H. Hendriks, P. Klint, J. Rekers, The syntax definition formalism SDF – reference manual, *SIGPLAN Notices* 24(11) (1989) 43–75.
- [11] P.J. Higgins, Algebras with a scheme of operators, *Math. Nachr.* 27 (1963) 115–132.
- [12] P. Klint, E. Visser, Using filters for the disambiguation of context-free grammars, in: G. Pighizzini, P. San Pietro (Eds.), *Proc. ASMICS Workshop on Parsing Theory*, Milano, Italy, pp. 1–20, Tech. Rep. 126–1994, Dipartimento di Scienze dell'Informazione, Università di Milano, 1994.
- [13] J. Maluszynski, Towards a programming language based on the notion of two-level grammar, *Theoret. Comput. Sci.* 28 (1984) 13–43.

- [14] K. Meinke, Equational specification of abstract types and combinators, in: E. Boerger, G. Jaeger, H.K. Buening, M.M. Richter (Eds.), *Computer Science Logic – CSL’91, Lecture Notes in Computer Science*, vol. 626, Springer, Berlin, 1992, pp. 257–271.
- [15] K. Meinke, Universal algebra in higher types, *Theoret. Comput. Sci.* 100 (1992) 385–417.
- [16] K. Meinke, J.V. Tucker, Universal algebra, in: S. Abramsky, D. Gabbay, T.S.E. Maibaum (Eds.), *Handbook of Logic in Computer Science*, vol. I: Mathematical Structures, Oxford University Press, Oxford, 1992, pp. 189–141.
- [17] R. Milner, A theory of type polymorphism in programming, *J. Comput. System Sci.* 17 (3) (1978) 348–375.
- [18] B. Möller, Algebraic specification with higher-order operators, in: L. Meertens (Ed.), *Program Specification and Transformation*, Elsevier, Amsterdam, 1987, pp. 367–398.
- [19] F.C.N. Pereira, D.H.D. Warren, Definite clause grammars for language analysis – a survey of the formalism and a comparison with augmented transition networks. *Artificial Intelligence* 13 (1980) 231–278.
- [20] F.C.N. Pereira, D.H.D. Warren, Parsing as deduction, in: *Proc. 21st Ann. Meeting of the Assoc. Computational Linguistics*, Massachusetts Institute of Technology, Cambridge, MA, 1983.
- [21] A. Poigné, On specifications, theories, and models with higher types, *Inform. and Control* 68 (1986) 1–46.
- [22] T. Rus, J.S. Jones, Multi-layered pipeline parsing from multi-axiom grammars, in: A. Nijholt, G. Scollo, R. Steetskamp (Eds.), *Algebraic Methods in Language Processing AMILP’95, Twente Workshops in Language Technology*, vol. 10, Twente University of Technology, Enschede, The Netherlands, 1995, pp. 65–81.
- [23] M. Sintzoff, Existence of a van Wijngaarden syntax for every recursively enumerable set, *Ann. Soc. Sci. Bruxelles Ser. II* 81 (1967) 115–118.
- [24] M. Tomita, *Efficient Parsing for Natural Languages. A Fast Algorithm for Practical Systems*, Kluwer, Dordrecht, 1985.
- [25] E. Visser, *Combinatory algebraic specification & compilation of list matching*, Master’s thesis, Department of Computer Science, University of Amsterdam, Amsterdam, 1993.
- [26] E. Visser, Multi-level specifications, in: A. van Deursen, J. Heering, P. Klint (Eds.), *Language Prototyping. An Algebraic Specification Approach*, *AMAST Series in Computing*, vol. 5, World Scientific, Singapore, 1996, pp. 105–196.
- [27] D.S. Warren, Memoing for Logic Programs, *Comm. ACM* 35 (3) (1992) 94–111.
- [28] D.A. Watt, The parsing problem for affix grammars, *Acta Informatica* 8(1) (1977) 1–20.
- [29] A. van Wijngaarden, B.J. Mailloux, J.E.L. Peck, C.H.A. Koster, M. Sintzoff, C.H. Lindsey, L.G.L.T. Meertens, R.G. Fisker (Eds.), *Revised Report on the Algorithmic Language Algol 68*, Springer, Berlin, 1976.