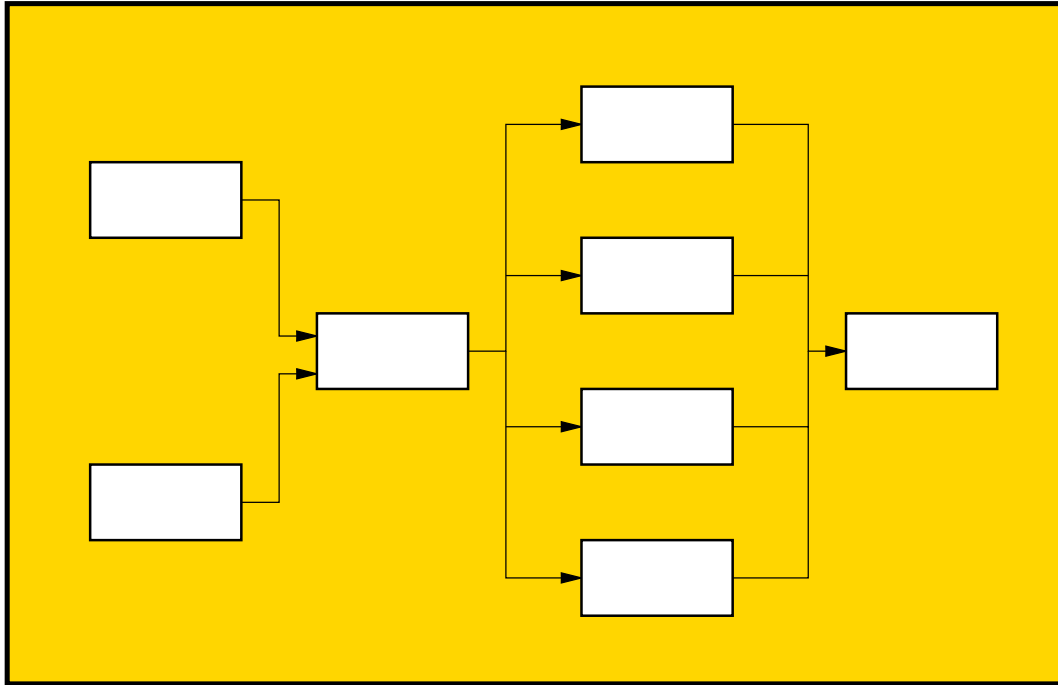


*University of Amsterdam*  
*Programming Research Group*

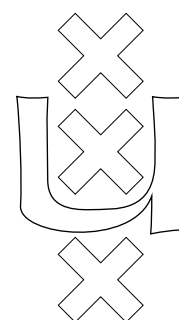


## A Family of Syntax Definition Formalisms

Eelco Visser







University of Amsterdam  
Department of Computer Science  
Programming Research Group

---

## A family of syntax definition formalisms

Eelco Visser

E. Visser

Programming Research Group  
Department of Computer Science  
University of Amsterdam

Kruislaan 403  
NL-1098 SJ Amsterdam  
The Netherlands

tel. +31 20 525 7590  
e-mail: visser@wins.uva.nl

An earlier version of this report appeared in M. G. J. van den Brand et al., *ASF+SDF'95. A Workshop on Generating Tools from Algebraic Specifications*, Technical Report P9504, Programming Research Group, University of Amsterdam. pp 89–126, May 1995.

**Acknowledgements** The author thanks Arie van Deursen, Jan Heering, Tobias Kuipers, Paul Klint, Mark van den Brand, Merijn de Jonge and Alex Sellink for useful suggestions and comments on previous versions of this report.

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

---

# Contents

---

<b>Contents</b>	<b>i</b>
<b>1 A Family of Syntax Definition Formalisms</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 An Overview of SDF2 . . . . .	4
1.3 Design . . . . .	8
1.4 Organization . . . . .	11
<b>2 Context-Free Grammars</b>	<b>13</b>
2.1 Symbols . . . . .	13
2.2 Grammars . . . . .	15
2.3 Context-Free Grammars (Kernel) . . . . .	17
2.4 Basic Symbols . . . . .	21
2.5 Parse Trees . . . . .	27
<b>3 Disambiguation and Abbreviation</b>	<b>39</b>
3.1 Priorities . . . . .	39
3.2 Regular Expressions . . . . .	47
3.3 Lexical and Context-Free Syntax . . . . .	54
3.4 Restrictions . . . . .	61
<b>4 Renaming and Modularization</b>	<b>65</b>
4.1 Renamings . . . . .	65
4.2 Aliases . . . . .	71
4.3 Modules . . . . .	74
<b>5 The Syntax Definition Formalism SDF2</b>	<b>83</b>
5.1 SDF2 . . . . .	83
5.2 Comparison to SDF . . . . .	87
5.3 Discussion and Concluding Remarks . . . . .	89
<b>A Auxiliary Modules for the Specification of SDF2</b>	<b>93</b>
A.1 Literals . . . . .	93
A.2 ATerms . . . . .	93
A.3 Renamings . . . . .	98
A.4 SDF2 . . . . .	100
<b>B Bibliography</b>	<b>103</b>



---

# A Family of Syntax Definition Formalisms

---

In the next chapters we present the design and specification of a family of syntax definition formalisms. The kernel of this family of formalisms is formed by context-free grammars. A number of orthogonal extensions to the kernel is defined. Many of these extensions are defined in terms of the primitives of the kernel by means of normalization functions. This provides a framework for constructing new formalisms by adapting and extending previous ones.

Included in the family are the following extensions of context-free grammars: uniform definition of lexical and context-free syntax, variables, disambiguation by priorities, follow restrictions and reject productions, a rich set of regular expressions defined in terms of context-free productions, character classes, aliases, parameterized modules with hidden imports and renamings. The accumulation of these extensions is the syntax definition formalism SDF2.

This chapter provides an introduction to SDF2 and gives an overview of the design and specification of the family of formalisms.

---

## 1.1 Introduction

New programming, specification and special purpose languages are being developed continuously. Syntax definition formalisms play a crucial role in the design and implementation of new languages. Syntax definition formalisms also play a role embedded in other languages: regular expressions in edit operations, macro definitions for macro preprocessors, user definable infix or distfix operators in programming languages, grammars as signatures in algebraic specification formalisms, and documents that contain a description of their own syntax.

The core of many syntax definition formalisms is formed by context-free grammars, which are widely used in computer science since their introduction by Chomsky (1956). A context-free grammar is a set of string rewrite rules of the form  $\alpha \rightarrow \mathcal{A}$  with  $\alpha$  a string  $\mathcal{A}_1 \dots \mathcal{A}_n$  of zero or more symbols and  $\mathcal{A}$  a symbol. A string (a sequence of symbols)  $w$  is a member of the language described by a grammar  $\mathcal{G}$  if it can be rewritten to the start symbol  $S$ , i.e., if there is a sequence  $w = \alpha_0 \rightarrow \alpha_1 \rightarrow \dots \rightarrow \alpha_n = S$  such that each step has the form  $\alpha_i \rightarrow \alpha_{i+1}$  with  $\alpha_i = \beta_1 \beta_2 \beta_3$  and  $\alpha_{i+1} = \beta_1 B \beta_3$  and  $\mathcal{G}$  contains a production  $\beta_2 \rightarrow B$ .

Despite, or maybe due to, the simplicity of this basic structure there has never emerged a standard formalism for syntax definition. The Backus Naur

Form (BNF), originally developed by Backus (1959) and Naur *et al.* (1960) for the definition of the syntax of Algol, is a commonly used notation for context-free grammars, but it does not have the status of a standard and many variants are in use. Several standard notations for syntax definition have been proposed (e.g., Wirth, 1977, Williams, 1982). None of these has been convincing, instead a number of similar or overlapping formalisms exist.

The reason for this divergence is that a practical syntax definition formalism serves not only to define languages, i.e., sets of strings. Syntax definitions are also interpreted as recognizers that decide whether a string belongs to a language, as parsers that map strings to parse trees, as mappings from parse trees to abstract syntax trees and as syntax directed editors. Plain context-free grammars are not adequate for this purpose. To support the compact definition of languages, formalisms can provide a variety of features as extensions to the basic structure: character classes, regular expressions, disambiguation by associativity and priority declarations, reuse by modularization, parameterization of language definitions, interfacing between the formalism and its environment, e.g., mapping to abstract syntax.

Various extensions of context-free grammars have been developed for attaching semantics to grammars: Attribute grammars (Knuth, 1968) attach attribute evaluation rules to productions. The computation of the semantics of a parse tree consists in computing the values of all attributes. This computation is orthogonal to parsing. Affix grammars (Koster, 1971) and extended affix grammars (Watt, 1977) are similar to attribute grammars, but predicates on affix values can play a role in disambiguation during parsing. Definite clause grammars (Pereira and Warren, 1980) are based on the Horn clauses of logic programming. Parsing is performed by the SLD resolution evaluation mechanism. Semantic values are represented by means of terms and passed around using unification and logic variables. Other approaches including algebraic specification use a separate formalism to define the semantics.

Traditionally, compiler construction is the main application area for syntax definition formalisms. The most well-known is the pair LEX/YACC. The formalism LEX (Lesk and Schmidt, 1986) is used to define the lexical syntax of a language using regular expressions. According to the regular expressions a string is analyzed and divided into tokens. In case more than one regular expression can be matched, a number of disambiguation rules such as prefer longest match and prefer regular expressions appearing earlier in the file. The ‘compiler compiler’ YACC (Johnson, 1975) is used for the definition of the context-free syntax of a language. An LALR(1) parser generator translates grammars to C programs if the grammar is LALR(1). Some conflicts in the parse table caused by ambiguous expression grammars can be solved by means of binary and unary precedences based on the ideas of Aho *et al.* (1975). Furthermore, the order of productions in the grammar is used to solve conflicts. Trees for a string can be constructed by calling C functions from the grammar productions.

Recent formalisms are generally based on the same deterministic parsing techniques, but extend the expressivity and declarativeness of syntax definition by providing mechanisms for building trees and coupling to other phases of compilation. Some examples are: The Cocktail compiler generator (Grosch, 1990) provides a BNF-like formalism with an LALR parser generator. The Eli system (Gray *et al.*, 1992) is a collection of tools for developing all aspects of compilers. The syntax definition formalism is based on context-free grammars. Attribute

rules are added to define semantics computations. The tree transformation language TXL (Cordy and Carmichael, 1993) is a programming language for source to source transformations by means of transformation rules on parse trees. The syntax definition formalism of TXL is based on context-free grammars extended with some regular operators. Lexical syntax is defined by means of predefined lexical notions and by means of regular expressions over character classes. PC-CTS (Parr and Quong, 1994) is a formalism based on top-down LL(1) parsing. The problems of unbounded lookahead are dealt with by means of backtracking and syntactic predicates that can be used to try out a variant before deciding which production to predict.

An application domain derived from compiler construction is the area of programming environments. A programming environment is a collection of tools for interactively developing and testing programs in some programming language. These tools are usually centered around an interactive syntax directed editor. A syntax editor has knowledge of the language of the programs being edited and provides support for checking the syntax of programs and for presenting and manipulating the structure of programs. In order to rapidly process changes to a program, incremental parsing and incremental compilation are used. Syntax definition formalisms developed for derivation of programming environments include the grammar formalism of the Synthesizer Generator (Reps and Teitelbaum, 1989), PSG (Bahlke and Snelting, 1986), METAL (Kahn *et al.*, 1983) and SDF (Heering *et al.*, 1989). The ASF+SDF Meta-Environment (Klint, 1993) is a programming environment for developing and generating programming environments from algebraic specifications. To speed up the development cycle for syntax definitions, incremental parser generation is used to only regenerate those parts of the parser that have been affected by a change.

Syntax definition in algebraic specification takes the form of grammars as algebraic signatures. The motivation here is to provide flexible notation for functions and constructors in abstract data type specifications and less the description of real programming languages. Therefore, the requirements on lexical syntax are not so strong. The correspondence of context-free grammars and many-sorted signatures was first described by Rus (see Hatcher and Rus, 1976). Goguen *et al.* (1977) showed that this correspondence could be used to define the semantics of programming languages. The correspondence was exploited in a number of algebraic specification formalisms to provide flexible, user-definable notation for functions and constructors. The first formalism to incorporate this was OBJ (Futatsugi *et al.*, 1985). Others are Cigale (Voisin, 1986), ASF+SDF (Heering *et al.*, 1989, Bergstra *et al.*, 1989a), the Meta Notation, used in action semantics (Mosses, 1992), and Elan (Vittekk, 1994).

The combination of features that a formalism provides is, necessarily, rather arbitrary and strongly influenced by the expected application of definitions and the environment in which generated tools have to operate. Although it is not desirable to include all conceivable features in a formalism—some features can not be combined with others and too many features results in an unmanageable formalism—the similarities between different formalisms can be exploited by reusing parts of the design and implementation of old formalisms. However, formalisms are conventionally designed in a monolithic way, containing an intertwined mix of features, resulting in a formalism with a lack of orthogonality and uniformity that is difficult to implement, extend and use for other applications than the originally intended ones. Syntax definition formalisms form no

exception to this rule.

Here we set out to design syntax definition formalisms in a modular way, as a family of formalisms each extending a small kernel with some feature for syntax definition. This approach should result in more orthogonal and uniform formalisms and should make it easier to (a) construct formalisms that use some subset of a set of known features, (b) adapt formalisms for use in other application areas, (c) implement tools for such formalisms and (d) design new formalisms that combine new features with existing ones.

As a first step to accomplish this goal we design a concrete formalism with a set of features that is useful in many application areas, but in particular in the application of grammars as signatures for algebraic specifications of programming languages. The result is the syntax definition formalism SDF2 that is a generalization of SDF. It incorporates several concepts and techniques introduced by Heering *et al.* (1989) in a more orthogonal and uniform way and adds several new features.

We use the algebraic specification formalism ASF+SDF to formally specify the family of syntax definition formalisms. For an introduction to ASF+SDF see Van Deursen *et al.* (1996).

In this chapter we outline the main features of SDF2 and examine the structure and design principles of the specification.

## 1.2 An Overview of SDF2

SDF2 is a syntax definition formalism based on context-free grammars, extended with character classes, sorts, literals, priorities, regular expressions, renamings, aliases and modules and combines the definition of lexical and context-free syntax into one formalism. The syntax definition in Figure 1.1 on page 7, taken from Visser (1997c), presents (the syntax of) a small untyped, first-order functional programming language, the data type environments and the evaluation function that interprets such functional programs using an instantiation of the environments data type. A program in this language might contain the following definition of a function `map` that applies a function `F` to all elements of a list `L`:

```
function map(F, L) is
  if(is-nil(L), nil(),
    cons(call(F, head(L)), map(F, tail(L))))
```

We sketch the main features of SDF2 and use the syntax definition in Figure 1.1 as running example.

### 1.2.1 Context-free Productions

The basis of the formalism is formed by *context-free productions*. These are rules of the form  $\alpha \rightarrow \mathcal{A}$ , where  $\alpha$  is a list of symbols  $\mathcal{A}_1 \dots \mathcal{A}_n$  ( $n \geq 0$ ) and  $\mathcal{A}$  a symbol. A production declares that a string of category  $\mathcal{A}$  can be constructed by concatenating strings of the categories  $\mathcal{A}_i$ . For instance, the production

```
Fun "(" Terms ")" -> Term
```

defines that a term can be constructed by means of a function symbol followed by a list of terms separated by commas between parentheses. Conventionally,



context-free productions are written as  $\mathcal{A} \rightarrow \alpha$  or as  $\mathcal{A} ::= \alpha$ . In SDF2 productions are written the other way around to make the similarity to function declarations more apparent. This is useful because SDF2 definitions are used as signatures in algebraic specifications such that productions correspond to algebraic operators. For example, in a conventional signature one would declare the evaluation function that computes the value of a term with respect to a program and some environment by means of the function declaration

```
eval : Program # Term # Env -> Term
```

The production

```
eval "[[" Program "]" "(" Term ")" "_" Env -> Term
```

not only defines a function with the same input types, but also the syntax for its applications, i.e., the program argument should be enclosed in double square brackets and the term argument should be enclosed in parentheses.

### 1.2.2 Character Classes

Syntax definitions describe languages consisting of strings of *characters*, where the set of all characters can be encoded by a finite set of consecutive natural numbers. *Character Classes* are compact descriptions of sets of characters and are typically used in the definition of lexical categories such as layout, identifiers and numbers. The example contains the following character classes: the characters space, tab and newline `[\ \t\n]`, all characters except newline `~[\n]`, all uppercase letters `[A-Z]`, all lowercase letters `[a-z]`, all letters and digits and the hyphen character `[a-zA-Z0-9\ -]`.

### 1.2.3 Literals

*Literals* are strings of characters between double quotes that stand for exactly that string of characters. These are used to represent keywords—such as "function" and "program"—and operators and other literal symbols—such as "[", "|->", "|>", and "\*". The definition of the function `eval` does not use quotes for the function name. This is an exception to the general rule: identifiers starting with a lowercase letter can also be used as literals.

### 1.2.4 Sorts

The basic nonterminal symbols used in productions are *sorts*, which are written as identifiers starting with a capital letter. Sorts should be declared in a sorts section. The example defines the sorts `Var`, `Fun`, `Term`, etc.

### 1.2.5 Regular Expressions

More complex nonterminal symbols can be formed by means of *regular expressions* that provide abbreviations for tupling `(_)`, iteration `_*` and `_+`, optional constructs `_?`, and alternatives `_|_`. For example, `[a-zA-Z0-9\ -]*` denotes a list of zero or more characters from the set of letters, digits and hyphens, `{Term " , "}*` declares a list of terms separated by commas and the expression `(Key "|->" Value)*` denotes lists of zero or more tuples consisting of a key, the symbol `"|->"` and a value.

### 1.2.6 Aliases

Since such regular expressions can become quite tedious to type, it can be useful to introduce a shorter name for such symbols. This can be done by introducing a *symbol alias*. For example the declaration

```
aliases
  {Term ","}* -> Terms
```

introduces `Terms` as an alias for the regular expression `{Term ","}*.`

### 1.2.7 Priorities

Some productions that have a sensible type are syntactically ambiguous. For instance, the two productions for destructive and consistent environment update

```
Env "|>" Env -> Env
Env "*" Env -> Env
```

are ambiguous with respect to themselves and to each other, e.g., the environment expression `Env |> Env * Env` can be constructed as `(Env |> Env) * Env` or as `Env |> (Env * Env)`. *Associativity* and *priority* declarations are a way to resolve most ambiguities of this type. In the example, the ambiguity in the expressions above is resolved by means of the priority declaration

```
Env "|>" Env -> Env > Env "*" Env -> Env
```

that declares `"|>"` to have higher priority than `"*"`, which entails the `(Env |> Env) * Env` interpretation. The `left` attribute of a production declares that the operator is left-associative.

### 1.2.8 Lexical and Context-free Syntax

The phrases making up a string over a language are usually divided into lexical tokens—the words of a sentence—and context-free phrases. The distinction between tokens and phrases is that the tokens making up a phrase can be separated by layout (whitespace and comments) while the characters comprising a token cannot. In definitions this distinction is indicated by means of *lexical* and *context-free* productions. For example, the lexical definition

```
[a-z][a-zA-Z0-9\-*]* -> Fun
```

indicates that function symbols consist of a number of adjacent characters starting with a lowercase letter, followed by zero or more letters, digits or hyphens. While the tokens in the term `succ ( zero() )` can be separated by spaces, the characters in the token `succ` cannot. The layout that can occur between tokens should also be specified. The symbol `LAYOUT` is reserved for this purpose. In the example, layout is declared as

```
[\ \t\n] -> LAYOUT
"%" ~[\n]* -> LAYOUT
```

meaning that spaces, tabs and newlines (also called whitespace) are layout and that any suffix of a line starting with two percent signs is comment.

```

module Functional-Programs
exports
  sorts Var Fun Term FunDef Program
  aliases
    {Term ","}* -> Terms
    {Var  ","}* -> Vars
  lexical syntax
    [\ \t\n]          -> LAYOUT
    "%%" ~[\n]*       -> LAYOUT
    [A-Z][a-zA-Z0-9\-*] -> Var
    [a-z][a-zA-Z0-9\-*] -> Fun
  context-free syntax
    Var          -> Term
    Fun "(" Terms ")" -> Term

    "function" Fun "(" Vars ")" "is" Term -> FunDef
    "program" FunDef*                    -> Program

module Environments
exports
  sorts Key Value Env
  context-free syntax
    "[" (Key "|->" Value)* "]" -> Env
    Env "(" Key ")"           -> Value
    Env ">" Env                -> Env {left}
    Env "*" Env               -> Env {left}
    "(" Env ")"               -> Env {bracket}
  context-free priorities
    Env ">" Env -> Env > Env "*" Env -> Env

module Function-Eval
imports Functional-Programs
      Environments [Key => Var Value => Term]
exports
  context-free syntax
    eval "[" Program "]" "(" Terms ")" "_" Env -> Terms

```

Figure 1.1: SDF2 definition of the *syntax* of a small functional programming language and its main evaluation function.

### 1.2.9 Modules

Grammars can be divided in a number of *modules* such that parts of a grammar can be reused in various language definitions. Modules consist of a list of exports and hiddens sections. An import of a module  $M$  into a module  $N$  denotes the inclusion of the exported grammar of  $M$  into  $N$ . Thus the import of module **Terms** in module **Functions** means that the syntax of terms is included in the syntax of programs. To prevent name clashes or to instantiate generic modules,

*renamings* of symbols and productions can be applied to imported modules. For example, the module `Function-Eval` specifying the evaluation function, imports the generic module defining environments by means of

```
imports Environments [Key => Var Value => Term]
```

renaming the sort `Key` to `Var` and the sort `Value` to `Term`, thus instantiating it for use with the terms of the functional programming language.

Modules can also be parameterized with a list of parameter symbols that can be instantiated on import. For instance, module `Environments` might also have been declared as

```
module Environments[Key Value]
```

declaring `Key` and `Value` as parameters. The import

```
imports Environments[Var Term]
```

would then perform the instantiation.

A complete syntax definition consists of a list of modules and a designated top module. The language defined by such a definition is the one defined by the grammar associated to the top module. Of course, in a programming environment for SDF2 this list does not have to reside in a single file. More likely, each module will be defined in one file with the module name as the file name.

### 1.3 Design

The next chapters give a formal algebraic specification of the syntax and semantics of SDF2. The semantics of a syntax definition is characterized by the well-formed trees it generates. A tree is associated with a sentence—its *yield*. The language associated to a definition is the set of sentences that are yields of trees generated by the definition. A parser is a function that given a sentence, produces the tree (or set of trees) that have that sentence as yield. We do not describe parsing as part of the specification of SDF2, but specify the output required of a parser and allow any implementation that does so. Parsing for SDF2 grammars is described in Visser (1997e).

#### 1.3.1 Modularization

The formalism SDF2 is not designed monolithically, but modularized, as a family of formalisms. The kernel of this family is formed by context-free grammars. All features are defined as independent extensions of the kernel. The combination of the features forms SDF2. This setup makes it easier to construct a variant of the formalism by adding, removing or modifying features. Figure 1.2 depicts the structure of the family by means of (an abstraction of) the import graph of the specification.

Furthermore, the specification of SDF2 covers several aspects. The *syntax* of the formalism consists of the definition of the form of all its constructs. *Projection* functions on these constructs are defined in order to extract information from them. *Normalization* functions transform a syntax definition in order to simplify it. The specification of *parse trees* consists of several parts. A generic format for the representation of structured data called ATerms (Van den Brand

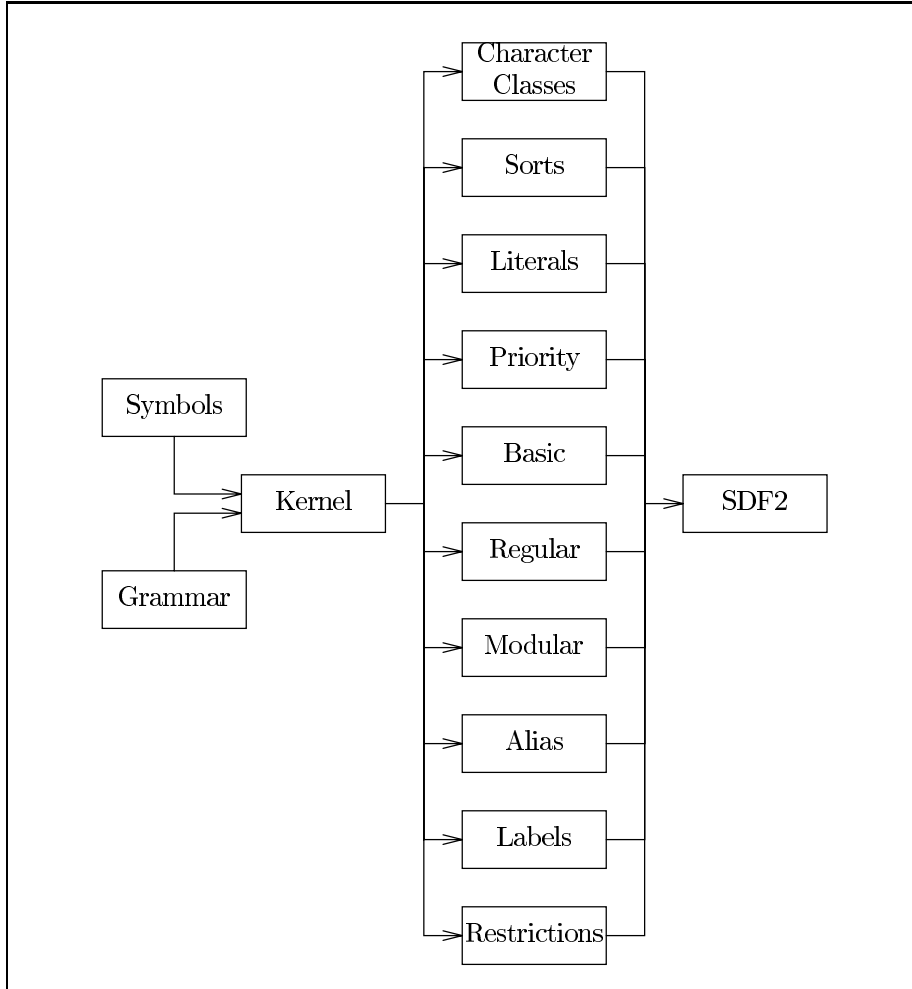


Figure 1.2: Import graph for the definition of SDF2.

*et al.*, 1997) is used to represent parse trees. In order to use this format for a specific purpose, *constructor* names have to be defined. To represent grammar information in parse trees, several constructs of the formalism have to be *encoded as ATerms*. Given this framework, the well-formedness of a tree with respect to a grammar can be defined. Furthermore, the *yield* of trees and the *equality* of trees are defined.

For each feature a number of modules are defined that each define an aspect of the formalism for that feature. The result is the matrix of modules listed in Table 1.1. The rows of the matrix contain the modules for one feature. The columns of the matrix contain all modules for one aspect. Each module in the matrix has a name consisting of the name of the feature and the name of the aspect separated by -Sdf-. For instance module Kernel-Sdf-Syntax specifies the syntax of the constructs in the kernel. So for each feature X we have modules X-Sdf-Syntax, X-Sdf-Projection, X-Sdf-Normalization, X-Sdf-Renamings, X-Sdf-Constructors, X-Sdf-ATerms, X-Sdf-Trees and X-Sdf-Equality. With the

	<i>Syntax</i>	<i>Projection</i>	<i>Normalization</i>	<i>Renaming</i>	<i>Constructors</i>	<i>ATerms</i>	<i>Trees</i>	<i>Equality</i>
Symbols	2.1					2.5.3		
Grammar	2.2		2.3.3		A.2.1	A.2.2		
Kernel	2.3.1	2.3.2	2.3.3	4.1.3	2.5.2	2.5.4	2.5.5	2.5.8
Sorts	2.4.1	2.4.1	2.4.1	A.3	A.2.1	A.2.2		
CC	2.4.2		2.4.2	A.3	A.2.1	A.2.2	2.5.6	
Literals	2.4.3		2.4.3	A.3	A.2.1	A.2.2		
Priority	3.1.1	3.1.2	3.1.3	A.3		A.2.2	3.1.4	
Regular	3.2.1		3.2.2	A.3	A.2.1	A.2.2		
Basic	3.3.1		3.3.2	A.3	A.2.1	A.2.2		3.3.3
Restrictions	3.4.1	3.4.2	3.4.3	A.3				
Renaming	4.1.1	4.1.2		4.1.3			4.1.4	
Alias	4.2.1	4.2.2	4.2.3	4.2.3				
Modular	4.3.1	4.3.2	4.3.3	4.3.4				
Label								
Sdf2	5.1.1	A.4	5.1.2	A.4	A.4	A.4	A.4	A.4

Table 1.1: Modules of the family of syntax definition formalisms. The last row contains the collecting modules for SDF2. There are no collection modules for the rows. The numbers refer to the sections presenting the modules.

exception that if some feature does not change some aspect, the module is omitted.

### 1.3.2 Normalization

An important role in the design of SDF2 is played by the normalization function. In general, a normalization function defines a transformation on an expression that yields an expression in the same language, which uses less features. The normalized expression has the same meaning as the original one. Thus, a normalization is a mapping from the language onto (a subset of) the same language. Ideally, a normalization function should be idempotent, i.e., yield the same result when applied twice. An implementation for such a language only has to consider the simplified expressions, while users have a more expressive language at their disposal.

The requirement that normalization produces an expression in the language itself entails that all constructs used for encodings should also be present in the original language, i.e., the language should be closed under normalization. For example, one of the normalizations in this definition renames a symbol  $\mathcal{A}$  into  $\langle \mathcal{A}\text{-LEX} \rangle$  if it occurs in the lexical syntax. Therefore, the constructor  $\langle \_ \text{-LEX} \rangle$  introduced for the purpose of normalization also becomes a construct of the language before normalization.

A consideration in the definition of a normalization is whether two different expressions that are equivalent with respect to the semantics have the same

normal form. This can be useful when expressions have to be compared. This is for example the case in the normalization of character classes. In Visser (1997b) a normalization of character classes to a unique normal forms is defined such that two character classes that represent the same set of characters are normalized to the same character class expression. In general, however, the normalizations in this paper will not have this property. For instance, all permutations of a list of productions are equivalent. Although such lists could be ordered by imposing an ordering on productions, this is not done here, since comparisons of lists of productions are not needed. In such cases definitions can not use syntactic equality to determine equivalence.

Using this approach SDF2 is an expressive formalism that depends on a small set of features, i.e., we have:

SDF2	=	context-free grammars
	+	priorities
	+	character-classes
	+	reject productions
	+	follow restrictions

Features that are provided in the formalism, but that are eliminated, i.e., expressed using the features above are: literals, regular expressions, lexical and context-free syntax, variables, modules, renamings, and aliases. Furthermore, character classes, priority declarations and grammar composition are simplified considerably.

The normalization of SDF2 is defined as a pipeline of normalizations, as is illustrated in Figure 1.3. This modularization of the definition of normalization makes it easy to define an extension and express it in existing features using a new normalization function. The overall normalization is extended by adding the new function to the normalization pipeline.

## 1.4 Organization

The next chapters discuss the specification of the family of syntax definition formalisms that is the basis of SDF2. Chapter 2 defines context-free grammars, the basic symbols sorts, character classes and literals and defines the well-formed parse trees characterized by a grammar. Chapter 3 defines disambiguation by means of priorities, regular expressions, lexical and context-free syntax and restrictions for lexical disambiguation. Chapter 4 introduces renamings, aliases and modules. Chapter 5 these extensions are combined in the formalism SDF2. The formalism is compared to SDF and a discussion of possible improvements and extensions is given. Appendix A gives some auxiliary modules for the specification of SDF2.

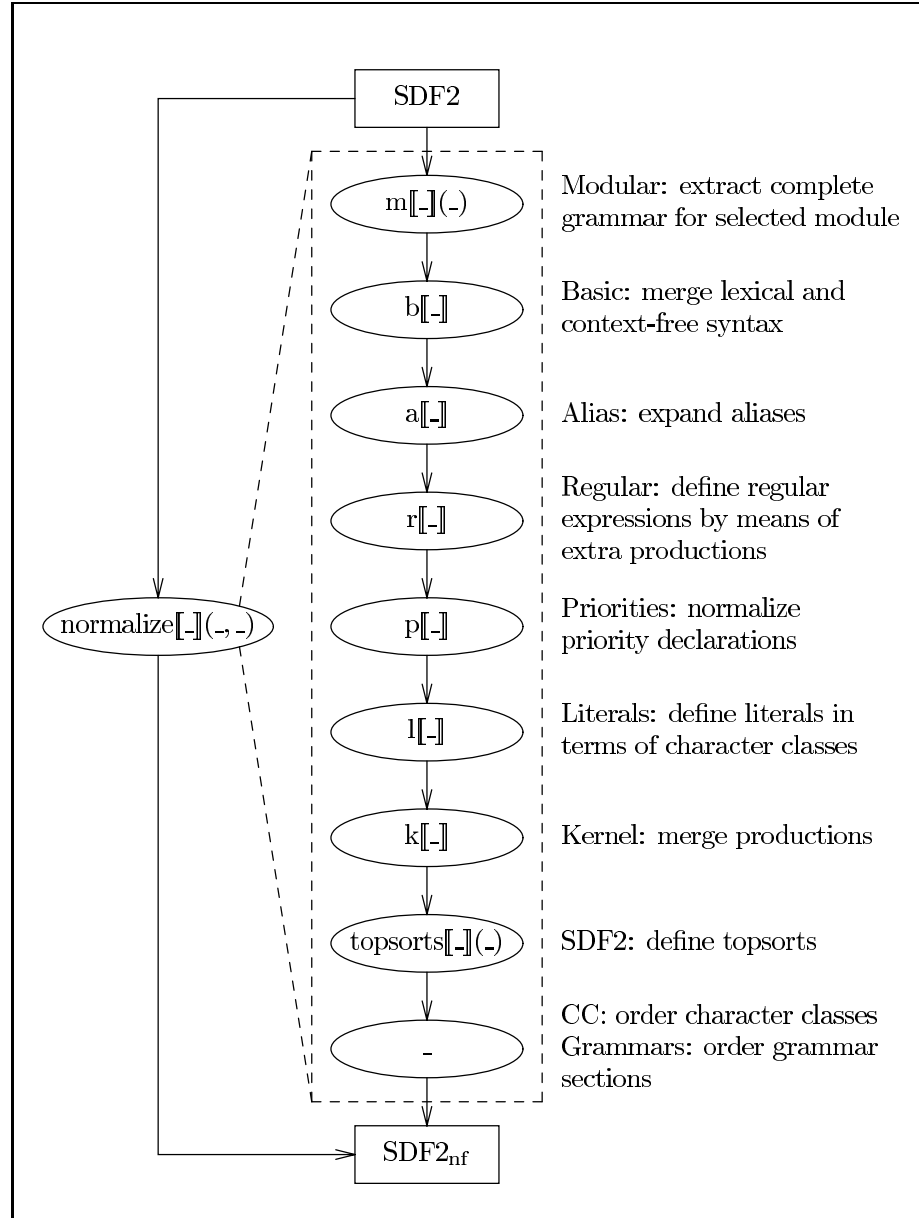


Figure 1.3: The normalization of SDF2 definitions consists of a series of independent transformations. The last step is not performed by a transformation function, but by rewrite rules acting on the constructors themselves.



---

## 2

# Context-Free Grammars

---

In this chapter a context-free grammar formalism is defined. First an abstract framework of symbols and grammars is introduced. In this framework a grammar is interpreted by means of several predicates and functions that characterize the trees and strings of symbols generated by a grammar. A well-formedness predicate on parse trees characterizes the trees over a grammar. From the parse trees over a grammar the strings of the language defined by the grammar are derived. Parse trees are represented in the annotated term format ATerms. An instance of this framework is set up with as kernel context-free productions. We introduce three kinds of basic symbols to be used in productions: sorts, character classes and literals.

---

### 2.1 Symbols

Syntax definitions define languages, i.e., sets of strings of *symbols*. A string of symbols is a list of zero or more symbols. The sort Symbol is declared without actually specifying any constructors for it. This entails that the sort is empty at this point, but can be extended later on with constructors. We do not make a distinction between terminal and nonterminal symbols. Whether a symbol is a terminal or nonterminal symbol is determined by the interpretation and is not fixed syntactically. A symbol that plays the role of a terminal in one view can be a nonterminal in another view. An example is a literal that can be considered as a terminal token or as a nonterminal that is defined in terms of characters.

```
module Symbols
imports Layout
exports
  sorts Symbol Symbols
  context-free syntax
    "(" Symbol ")" → Symbol {bracket}
    Symbol* → Symbols
  variables
    [ABC][0-9']* → Symbol
    [αβγ][0-9']* → Symbol*
    [αβγ]"+"[0-9']* → Symbol+
```

### 2.1.1 Projection

The function  $++$  concatenates strings and  $| \cdot |$  gives the length of a string. The predicate  $\in$  decides list membership.

**module** Symbols-Projection  
**imports** Symbols<sup>2.1</sup> Booleans Integers  
**exports**

**context-free syntax**

Symbols “++” Symbols  $\rightarrow$  Symbols {**assoc**}  
 “|” Symbols “|”  $\rightarrow$  Int  
 Symbol “ $\in$ ” Symbols  $\rightarrow$  Bool

**equations**

Concatenation, length and membership of symbol lists.

- |     |  |
|-----|--|
| [1] | $\alpha ++ \beta = \alpha \beta$   |
| [2] | $  \cdot   = 0$  |
| [3] | $  \mathcal{A} \alpha   =   \alpha   + 1$  |
| [4] | $\mathcal{A} \in = \perp$  |
| [5] | $\mathcal{A} \in \mathcal{A} \alpha = \top$  |
| [6] | $\mathcal{A} \in \mathcal{B} \alpha = \mathcal{A} \in \alpha$ <b>when</b> $\mathcal{A} \neq \mathcal{B}$ |

The concatenation function on sorts such as Symbols is needed because the concatenation of the built-in associative lists (e.g., Symbol\*) of ASF+SDF is not inherited through the injection  $\text{Symbol*} \rightarrow \text{Symbol}$ . The injection is needed because list sorts cannot be output sorts of functions in ASF+SDF.

### 2.1.2 Sets

From lists of symbols we construct sets of symbols by means of the operation  $\{ \cdot \}$ . Although this constructor does not remove double elements from the list, it signifies that the number of occurrences in the list does not matter. Operations on sets are union ( $\cup$ ), difference ( $/$ ) and membership ( $\in$ ). The union  $v_1 \cup v_2$  adds only those elements of  $v_1$  to  $v_2$  that do not already occur in  $v_2$ . If a set is constructed by means of union from singleton sets, the resulting set does not contain double elements. Although this is not strictly necessary it is useful when later on something has to be done once for each symbol in some set.

**module** Symbol-Sets  
**imports** Symbols-Projection<sup>2.1.1</sup> Booleans  
**exports**

**sorts** SymbolSet

**context-free syntax**

“{” Symbols “}”  $\rightarrow$  SymbolSet  
 SymbolSet “ $\cup$ ” SymbolSet  $\rightarrow$  SymbolSet {**right**}  
 SymbolSet “/” SymbolSet  $\rightarrow$  SymbolSet {**left**}  
 “(” SymbolSet “)”  $\rightarrow$  SymbolSet {**bracket**}  
 Symbol “ $\in$ ” SymbolSet  $\rightarrow$  Bool

**priorities**

SymbolSet “/”SymbolSet  $\rightarrow$  SymbolSet  $>$

SymbolSet “ $\cup$ ”SymbolSet  $\rightarrow$  SymbolSet

**variables**

“ $v$ ” $[\theta\text{-}\theta']^* \rightarrow$  SymbolSet

**equations**

Membership

$$[1] \quad \mathcal{A} \in \{\alpha\} = \mathcal{A} \in \alpha$$

Union

$$[2] \quad \{\} \cup v = v$$

$$[3] \quad v \cup \{\} = v$$

$$[4] \quad \{\alpha^+ \beta^+\} \cup v = \{\alpha^+\} \cup \{\beta^+\} \cup v$$

$$[5] \quad \{\mathcal{A}\} \cup v = v \quad \textbf{when } \mathcal{A} \in v = \top$$

$$[6] \quad \{\mathcal{A}\} \cup \{\alpha\} = \{\mathcal{A} \alpha\} \quad \textbf{otherwise}$$

Difference

$$[7] \quad \{\} / v = \{\}$$

$$[8] \quad v / \{\} = v$$

$$[9] \quad \{\alpha^+ \beta^+\} / v = \{\alpha^+\} / v \cup \{\beta^+\} / v$$

$$[10] \quad \{\mathcal{A}\} / v = \{\} \quad \textbf{when } \mathcal{A} \in v = \top$$

$$[11] \quad \{\mathcal{A}\} / v = \{\mathcal{A}\} \quad \textbf{otherwise}$$

## 2.2 Grammars

A syntax definition consists of a grammar. The only generic operations on grammars that we define at this point are an associative composition operation that is used to combine grammars and the constant  $\emptyset$  representing the empty grammar.

**module** Grammar-Syntax

**imports** Layout

**exports**

**sorts** Grammar

**context-free syntax**

“ $\emptyset$ ”  $\rightarrow$  Grammar

Grammar Grammar  $\rightarrow$  Grammar {**assoc**}

“(” Grammar “)”  $\rightarrow$  Grammar {**bracket**}

**variables**

“ $\mathcal{G}$ ” $[\theta\text{-}\theta']^* \rightarrow$  Grammar

### 2.2.1 Interpretation

A grammar defines a set of strings of symbols, a language. We specify the language derived by a grammar indirectly, via the trees it generates. Parse trees will be represented by means of ATerms, a term format for the representation

and exchange of structured data (Van den Brand *et al.*, 1997). The format will be introduced in §2.5.1. In fact we will not just define parse trees, but parse forests. A parse forest is a compact encoding of a collection of parse trees in which contexts are shared. A parse forest is used to represent all parse trees for an ambiguous sentence.

The interpretation of a grammar is now given by two predicates and a function. The predicate  $\mathcal{G} \vdash T$  characterizes the terms  $T$  that are *well-formed parse forests* over grammar  $\mathcal{G}$ . The function  $\text{yield}[\![\mathcal{G}]\!](T)$  maps a parse tree  $T$  to a string of symbols. The predicate  $T \in T'$  determines membership of a tree  $T$  in a forest  $T'$ . The set of trees generated by a grammar is denoted by  $\mathcal{T}[\![\mathcal{G}]\!]$ .

Given these functions we can derive the notion of the language generated by a grammar and the notion of a parser for a grammar. The *language* generated by a grammar corresponds to the set of yields of the parse trees it generates. In other words, a string of symbols  $\alpha$  is an element of the language defined by a grammar  $\mathcal{G}$ , if there exists a well-formed tree  $T$  over the grammar with  $\alpha$  as its yield. A *parser*  $\Pi[\![\mathcal{G}]\!](\alpha)$  is a function that maps a string  $\alpha$  in  $\mathcal{L}[\![\mathcal{G}]\!]$  to a parse forest  $T$  containing all well-formed trees such that their yield is  $\alpha$ . The module below summarizes these definitions. The equations define the set of parse trees  $\mathcal{T}[\![\mathcal{G}]\!]$ , the language  $\mathcal{L}[\![\mathcal{G}]\!]$  and the parser  $\Pi[\![\mathcal{G}]\!](\alpha)$  in terms of well-formedness, yield and parse forest membership.

**module** Grammar-Interpretation

**imports** Grammar-Syntax<sup>2.2</sup> Symbols<sup>2.1</sup> Booleans ATerms<sup>2.5.1</sup> Symbols-Sets  
ATerm-Sets

**exports**

**context-free syntax**

Grammar “ $\vdash$ ” ATerm  $\rightarrow$  Bool  
 ATerm “ $\in$ ” ATerm  $\rightarrow$  Bool  
 yield “[” Grammar “]” “(” ATerm “)”  $\rightarrow$  Symbols  
  
 “T” “[” Grammar “]”  $\rightarrow$  ATermList  
 “L” “[” Grammar “]”  $\rightarrow$  SymbolsSet  
 “ $\Pi$ ” “[” Grammar “]” “(” Symbols “)”  $\rightarrow$  ATerm

**equations**

$$\begin{array}{ll}
 [1] & \frac{\mathcal{G} \vdash T = \top}{T \in \mathcal{T}[\![\mathcal{G}]\!] = \top} \\
 [2] & \frac{\mathcal{G} \vdash T = \top, \text{yield}[\![\mathcal{G}]\!](T) = \alpha}{\alpha \in \mathcal{L}[\![\mathcal{G}]\!] = \top} \\
 [3] & \frac{\mathcal{G} \vdash T = \top, \text{yield}[\![\mathcal{G}]\!](T) = \alpha}{T \in \Pi[\![\mathcal{G}]\!](\alpha) = \top}
 \end{array}$$

Note that these equations are non-constructive, i.e., do not provide decision procedures, but are merely a specification of the required behaviour. (Module Symbols-Sets defines sets of strings of symbols in a similar way as sets of symbols are defined in module Symbol-sets.)

### 2.2.2 Overview

In the rest of this chapter we will provide the specifications of well-formedness, yield and forest membership for a context-free grammar formalism. In §2.3 we define the syntax and normalization of context-free productions. In §2.4 we define basic symbols to be used in grammars: sorts, character classes and literals. In §2.5 we define the well-formed parse trees generated by a context-free grammar. In the next chapters this formalism is extended with a number of features.

An implementation of a parser is not specified, but can be chosen such as to comply with this specification. One possible implementation is discussed in Visser (1997e).

## 2.3 Context-Free Grammars (Kernel)

### 2.3.1 Syntax

The kernel of SDF2 is formed by context-free grammars. A context-free production is a structure  $\alpha \rightarrow \mathcal{A}$ , where  $\alpha$  is a list of symbols and  $\mathcal{A}$  a symbol. A context-free grammar is formed by a list of productions preceded by the keyword `syntax`.

Conventionally, since their introduction by Chomsky (1956), context-free productions are written as  $\mathcal{A} \rightarrow \alpha$  to emphasize the generative view of grammars. A grammar *generates* a string from a symbol, by repeatedly replacing some symbol in a string by the symbols on the right-hand side of a production. There exist many variants of this ‘standard’ notation, e.g.,  $\mathcal{A} ::= \alpha$  in BNF (Backus, 1959) and  $\mathcal{A} : \alpha_1 | \dots | \alpha_n$ ; in YACC (Johnson, 1975).

The unconventional  $\alpha \rightarrow \mathcal{A}$  notation for productions introduced by Heering *et al.* (1989) emphasizes the *functional* view of productions when used in the context of algebraic specification. A production coincides with the declaration of the name and type of a function. This notation is a unification of the definition of context-free productions with the declaration of mixfix functions in algebraic specification formalisms. For example, the declaration of the infix addition operator on natural numbers that is declared as `Nat ::= Nat "+" Nat` in BNF, is declared as `op _ + _ : nat nat -> nat .` in OBJ (Futatsugi *et al.*, 1985) and as `@ + @ : (nat nat) nat` in Elan (Vittekk, 1994). In SDF this becomes `Nat "+" Nat -> Nat`.

All these notations are equivalent in expressive power and could be used instead of the current one. We could effortlessly define a version of SDF that uses the  $\mathcal{A} ::= \alpha$  notation of BNF and define its meaning by translation to the notation used here. Note, however, that this does not mean that other aspects of these formalisms have the same expressive power nor that the parsing techniques coupled to these formalisms all have the same power.

Optionally, productions can have a list of attributes. An *attribute* is an annotation of a production that gives some extra syntactic or semantic information about the production. An example of an attribute that will be introduced in §3.1 is `left` that indicates left associativity of the production. Productions can have any number of attributes. The kernel does not provide any attributes, but to be able to introduce attributes later on without having to introduce an extra constructor for productions, the attribution of a production is defined here.

```

module Kernel-Sdf-Syntax
imports Symbols2.1 Grammar-Syntax2.2
exports
  sorts Attribute Attributes Production Productions
  context-free syntax
    "{" {Attribute ","}* "}"      → Attributes
    "}" {Attribute ","}* "}"      → Attributes
    Symbols "→" Symbol Attributes → Production
    Production*                    → Productions
    "syntax" Productions           → Grammar
  variables
    "attr"[0-9']* → Attribute
    "attr"*"[0-9']* → {Attribute ","}*
    "attr"+"[0-9']* → {Attribute ","}+
    "$"[0-9']* → Attributes
    [p][0-9']* → Production
    [p]*"[0-9']* → Production*
    [p]+"[0-9']* → Production+

```

### 2.3.2 Projection

We define concatenation functions for lists of productions and lists of attributes. The concatenation function for attributes removes duplicates. A production with an empty list of attributes is equal to a production without attributes. The projection function ‘P’ gives the productions of a grammar, the function ‘ $\bar{P}$ ’ gives the non-production parts of a grammar (to be defined later), and the function ‘ $P_A$ ’ gives all productions defining a symbol  $A$ . The function ‘symbols’ gives the set of all symbols in a grammar. The function ‘reachable’ gives all productions reachable from some set of symbols, i.e., used in the definition of those symbols.

```

module Kernel-Sdf-Projection
imports Kernel-Sdf-Syntax2.3.1 Symbol-Sets2.1.2
exports
  context-free syntax
    Productions "++" Productions      → Productions {right}
    Attributes "++" Attributes         → Attributes {right}
    Production "∈" Productions        → Bool
    Productions "⊆" Productions       → Bool
    Production "≅" Production         → Bool

    "P"(Grammar)                      → Productions
    " $\bar{P}$ " "(" Grammar ")"              → Grammar
    "P" "-" Symbol "(" Productions ")" → Productions

    symbols(Productions)              → SymbolSet
    symbols(Grammar)                  → SymbolSet

    reachable(SymbolSet, SymbolSet, Productions) → Productions
    reachable(SymbolSet, Grammar)         → Grammar
  equations

```

Concatenation of lists of productions, membership and subset of a list of productions.

$$\begin{aligned}
[1] \quad & p_1^* ++ p_2^* = p_1^* p_2^* \\
[2] \quad & p \in p_1^* p p_2^* = \top \\
[3] \quad & p \in p^* = \perp \quad \textbf{otherwise} \\
[4] \quad & \subseteq p^* = \top \\
[5] \quad & p p_1^* \subseteq p_2^* = p \in p_2^* \wedge p_1^* \subseteq p_2^*
\end{aligned}$$

Concatenation of attribute lists. Attributes occurring in both lists are added only once.

$$\begin{aligned}
[6] \quad & \{attr_1^+, attr_2^+\} ++ \{attr^*\} = \{attr_1^+\} ++ \{attr_2^+\} ++ \{attr^*\} \\
[7] \quad & \{attr\} ++ \{attr_1^*, attr, attr_2^*\} = \{attr_1^*, attr, attr_2^*\} \\
[8] \quad & \{attr\} ++ \{attr^*\} = \{attr, attr^*\} \quad \textbf{otherwise} \\
[9] \quad & \$ ++ = \$ \\
[10] \quad & ++ \$ = \$ \\
[11] \quad & \{\} =
\end{aligned}$$

The last equation states that an empty list of attributes  $\{\}$  is equal to no attributes.

Two productions are similar if they are the same up to their attributes

$$\begin{aligned}
[12] \quad & \alpha \rightarrow \mathcal{A} \$ _1 \cong \alpha \rightarrow \mathcal{A} \$ _2 = \top \\
[13] \quad & p_1 \cong p_2 = \perp \quad \textbf{otherwise}
\end{aligned}$$

Function ‘P’ gives all productions of a grammar and function ‘ $\bar{P}$ ’ gives all non-syntax parts of a grammar. The function  $P_{\mathcal{A}}$  gives all productions defining the symbol  $\mathcal{A}$ .

$$\begin{aligned}
[14] \quad & P(\text{syntax } p^*) = p^* \\
[15] \quad & P(\mathcal{G}_1 \mathcal{G}_2) = P(\mathcal{G}_1) ++ P(\mathcal{G}_2) \\
[16] \quad & P(\mathcal{G}) = \quad \textbf{otherwise} \\
[17] \quad & \bar{P}(\text{syntax } p^*) = \emptyset \\
[18] \quad & \bar{P}(\mathcal{G}_1 \mathcal{G}_2) = \bar{P}(\mathcal{G}_1) \bar{P}(\mathcal{G}_2) \\
[19] \quad & \bar{P}(\mathcal{G}) = \mathcal{G} \quad \textbf{otherwise} \\
[20] \quad & P_{\mathcal{A}}() = \\
[21] \quad & P_{\mathcal{A}}(\alpha \rightarrow \mathcal{A} \$ p^*) = \alpha \rightarrow \mathcal{A} \$ ++ P_{\mathcal{A}}(p^*) \\
[22] \quad & P_{\mathcal{A}}(p p^*) = P_{\mathcal{A}}(p^*) \quad \textbf{otherwise}
\end{aligned}$$

The function ‘symbols’ gives the set of symbols of a list of productions or a grammar.

$$\begin{aligned}
[23] \quad & \text{symbols}() = \{\} \\
[24] \quad & \text{symbols}(\alpha \rightarrow \mathcal{A} \$ p^*) = \{\alpha\} \cup \{\mathcal{A}\} \cup \text{symbols}(p^*) \\
[25] \quad & \text{symbols}(\text{syntax } p^*) = \text{symbols}(p^*) \\
[26] \quad & \text{symbols}(\mathcal{G}_1 \mathcal{G}_2) = \text{symbols}(\mathcal{G}_1) \cup \text{symbols}(\mathcal{G}_2) \\
[27] \quad & \text{symbols}(\mathcal{G}) = \{\} \quad \textbf{otherwise}
\end{aligned}$$

The function ‘reachable’ gives the subgrammar with those productions reachable from some set of symbols  $v$ . It is defined by applying the auxiliary ‘reachable’

function to the productions of the grammar. Observe how the resulting grammar is a composition of the reachable productions and the non-production parts of a grammar. The auxiliary function selects for each symbol in the original set the productions for that symbol from the original grammar. This is applied recursively to the symbols used in the left-hand sides of those productions. The first set argument of the auxiliary function represents the symbols already handled. The second set contains the symbols for which the productions still have to be looked up.

$$[28] \quad \text{reachable}(v, \mathcal{G}) = \text{syntax reachable}(\{\}, v, P(\mathcal{G})) \bar{P}(\mathcal{G})$$

$$[29] \quad \text{reachable}(v, \{\}, p^*) =$$

$$[30] \quad \frac{\mathcal{A} \in v_1 = \perp, P_{\mathcal{A}}(p_1^*) = p_2^*, \text{symbols}(p_2^*) / v_1 = v_2}{\text{reachable}(v_1, \{\mathcal{A} \alpha\}, p_1^*) = p_2^* \mathbin{++} \text{reachable}(v_1 \cup \{\mathcal{A}\}, v_2 \cup \{\alpha\}, p_1^*)}$$

$$[31] \quad \text{reachable}(v_1, \{\mathcal{A} \alpha\}, p^*) = \text{reachable}(v_1, \{\alpha\}, p^*) \\ \text{otherwise}$$

### 2.3.3 Normalization

*Grammar Normalization* Composition of grammars is commutative and the empty grammar is a unit for grammar composition. Since commutativity cannot be expressed by means of a terminating rewrite system, the following module normalizes grammar compositions as a right associative list, where the grammars are ordered or merged as specified by the operation  $\diamond$ . If  $\mathcal{G}_1 \diamond \mathcal{G}_2$  yields a pair  $\langle \mathcal{G}_3, \mathcal{G}_4 \rangle$  this means that the composition  $\mathcal{G}_1 \mathcal{G}_2$  should be replaced by  $\mathcal{G}_3 \mathcal{G}_4$ . The definition of  $\diamond$  can either merge the two grammars into one, yielding the pair  $\langle \mathcal{G}_3, \emptyset \rangle$ , or exchange the grammars yielding  $\langle \mathcal{G}_2, \mathcal{G}_1 \rangle$ . The termination of this normalization depends on the property of  $\diamond$  that a swap cannot be undone. The definition of  $\diamond$  should be extended for each new grammar constructor. For example, the merging of the productions of two adjacent syntax sections is expressed in the next module. The ordering could also be defined directly on the grammar composition operator, but that would entail that two equations would have to be written for each pair of constructors that have to be merged or swapped, corresponding to the last two equations below.

**module** Grammar-Normalization

**imports** Grammar-Syntax<sup>2.2</sup>

**exports**

**sorts** Grammar-Grammar

**context-free syntax**

Grammar “ $\diamond$ ” Grammar  $\rightarrow$  Grammar-Grammar

“ $<$ ” Grammar “ $,$ ” Grammar “ $>$ ”  $\rightarrow$  Grammar-Grammar

**equations**

The empty grammar  $\emptyset$  is a unit for composition and composition is associative.

$$[1] \quad \emptyset \mathcal{G} = \mathcal{G}$$

$$[2] \quad \mathcal{G} \emptyset = \mathcal{G}$$

$$[3] \quad \mathcal{G}_1 (\mathcal{G}_2 \mathcal{G}_3) = \mathcal{G}_1 \mathcal{G}_2 \mathcal{G}_3$$



Subgrammars can be swapped as specified by the function  $\diamond$ .

$$\begin{aligned} [4] \quad & \mathcal{G}_1 \mathcal{G}_2 = \mathcal{G}'_1 \mathcal{G}'_2 \quad \textbf{when} \quad \mathcal{G}_1 \diamond \mathcal{G}_2 = \langle \mathcal{G}'_1, \mathcal{G}'_2 \rangle \\ [5] \quad & \mathcal{G}_1 \mathcal{G}_2 \mathcal{G}_3 = \mathcal{G}_1 \mathcal{G}'_2 \mathcal{G}'_3 \quad \textbf{when} \quad \mathcal{G}_2 \diamond \mathcal{G}_3 = \langle \mathcal{G}'_2, \mathcal{G}'_3 \rangle \end{aligned}$$

*Context-free Grammar Normalization* The normalization function  $k[\cdot]$  for the kernel, merges productions with the same arguments and result symbols. If such productions have different attributes, these are joined. This normalization entails that two occurrences of the same production are identified and do therefore not cause an ambiguity. Consequently, other normalization functions can generate a production more than once, without changing the meaning of the grammar. This strategy will be relevant later on when we introduce modularization of grammars. The identification of productions means that a production that is declared in two or more different modules is identified when these modules are imported in the same module.

**module** Kernel-Sdf-Normalization

**imports** Kernel-Sdf-Projection<sup>2.3.2</sup> Grammar-Normalization<sup>2.3.3</sup>

**exports**

**context-free syntax**

“k” “ $\llbracket$ ” Grammar “ $\rrbracket$ ”  $\rightarrow$  Grammar

merge(Productions)  $\rightarrow$  Productions

**equations**

An empty list of productions is equivalent to an empty grammar and multiple syntax sections are merged into one.

$$\begin{aligned} [1] \quad & \text{syntax} = \emptyset \\ [2] \quad & \text{syntax } p_1^* \diamond \text{syntax } p_2^* = \langle \text{syntax } p_1^* p_2^*, \emptyset \rangle \end{aligned}$$

The normalization function ‘k’ merges productions with the same arguments and result, using the auxiliary function ‘merge’.

$$\begin{aligned} [3] \quad & k[\mathcal{G}] = \text{syntax merge}(P(\mathcal{G})) \bar{P}(\mathcal{G}) \\ [4] \quad & \frac{p_1 = \alpha \rightarrow \mathcal{A} \$1, \quad p_2 = \alpha \rightarrow \mathcal{A} \$2, \quad \$1 \mathbin{++} \$2 = \$3, \quad p_3 = \alpha \rightarrow \mathcal{A} \$3}{\text{merge}(p_1^* p_1 p_2^* p_2 p_3^*) = \text{merge}(p_1^* p_3 p_2^* p_3^*)} \\ [5] \quad & \text{merge}(p^*) = p^* \\ & \quad \textbf{otherwise} \end{aligned}$$

## 2.4 Basic Symbols

The kernel formalism presented in the previous section is a complete definition of context-free grammars, except for the notation of symbols. In this section we present three extensions of the kernel that provide notation for basic symbols needed in syntax definition. Sorts represent the non-terminals of grammars, the categories or domains that the grammar introduces. Character classes are used to represent the terminals of grammars, the characters from which strings are

built. Literals are convenient abbreviations for fixed strings of characters. With these extensions we will have a complete notation for context-free grammars. The extensions in later sections will provide features to make this formalism more expressive.

### 2.4.1 Sorts

*Syntax* Sorts are the symbols that represent the basic domains or categories of a syntax definition. A sort identifier is a word starting with an uppercase letter followed by zero or more letters or digits. Hyphens can be used between the first and last character. Sorts used in the productions of a grammar should be declared in a separate sorts section that consists of the keyword ‘sort’ and a list of symbols.

```
module Sorts-Sdf-Syntax
imports Kernel-Sdf-Syntax2.3.1
exports
  sorts Sort
  lexical syntax
    [A-Z] → Sort
    [A-Z][A-Za-z0-9\-*][A-Za-z0-9] → Sort
  context-free syntax
    Sort → Symbol
    “sorts” Symbols → Grammar
  variables
    “S”[0-9']* → Sort
```

*Normalization* Ordering of sorts and syntax sections (sorts are placed before syntax sections) and merging of sorts sections.

```
module Sorts-Sdf-Normalization
imports Sorts-Sdf-Syntax2.4.1 Kernel-Sdf-Normalization2.3.3
equations
```

```
[1]          sorts = ∅
[2]      sorts α ◇ sorts β = ⟨sorts α β, ∅⟩
[3]      syntax p* ◇ sorts α = ⟨sorts α, syntax p*⟩
```

*Projection* The projection function ‘S’ gives the list of sorts of a grammar.

```
module Sorts-Sdf-Projection
imports Kernel-Sdf-Projection2.3.2 Sorts-Sdf-Syntax2.4.1
exports
  context-free syntax
    “S”(Grammar) → Symbols
equations
```

The declared sorts of a grammar.

```
[1]      S(sort α) = α
[2]      S(G1 G2) = S(G1) ++ S(G2)
[3]      S(G) = otherwise
```

### 2.4.2 Character Classes

A character class is an expression such as, for example, `[a-z\']` that denotes a set of characters, in this case the set of all lower case letters and a prime. For example, the following definition defines identifiers as lists of characters starting with a lowercase letter followed by zero or more lowercase letters or digits.

```

sorts Id
syntax
  [a-z]          -> Id
  Id [a-z0-9]    -> Id

```

The meaning of character classes could be defined in terms of productions and characters, effectively eliminating them from the formalism. For instance, the character class `[a-z]` is completely defined by 26 productions of the form:

```
[a] -> [a-z]  [b] -> [a-z]  ...  [z] -> [a-z]
```

However, this would cause an enormous increase in the number of productions. Therefore, the interpretation of character classes is not defined by translating character classes out of the language. This means that interpretation functions should be extended to character classes.

We do not give the complete specification of character classes and character class arithmetic. A full specification of character classes can be found in Visser (1997b). The normalization defined there ensures that two classes that contain the same elements have the same normal form.

*Characters* A character is a constant of the form  $\backslash d_1 \dots d_n$ , where the  $d_i$  are decimal digits, denoting the  $d_1 \dots d_n$ -th member of some finite, linearly ordered universe of characters. Since specifying characters by their index in some encoding scheme is difficult, we provide easier syntax for specification of characters. Alphanumeric characters (letters and digits) can be specified as themselves. Other visible characters in the ASCII set can be specified by escaping them using a backslash, e.g., `\(` for left parenthesis, `\-` for a hyphen and `\`  (a backslash followed by a space) for space. The characters `\t` and `\n` represent tabs and newlines. Finally, there are two special characters, `\EOF` and `\TOP`. `\EOF` is the character used to indicate the end of a file. `\TOP` is used to represent the largest character in the character universe.

**module** Character-Syntax

**imports** Layout

**exports**

**sorts** Character NumChar ShortChar

**lexical syntax**

```

  [0-9]+          -> NumChar
  [a-zA-Z0-9]     -> ShortChar
  [~\000-\037A-Za-mo-su-z0-9] -> ShortChar

```

**context-free syntax**

```

  NumChar -> Character
  ShortChar -> Character
  "\TOP"   -> Character
  "\EOF"   -> Character

```

**variables**

`"c"[0-9]*`  $\rightarrow$  Character

*Character Classes* A set of characters—a character class—is represented by a list of characters and character ranges between square brackets [ and ]. A list is constructed by an injection of characters into lists and by a right associative binary concatenation operator on lists. Operations on character classes are difference (/), intersection (^), union (v) and complement with respect to the complete character set, i.e., the characters in the range \0-\TOP, (~).

**module** Character-Class-Syntax

**imports** Character-Syntax<sup>2.4.2</sup>

**exports**

**sorts** CharRange CharRanges OptCharRanges CharClass

**context-free syntax**

Character  $\rightarrow$  CharRange

Character "—" Character  $\rightarrow$  CharRange

CharRange  $\rightarrow$  CharRanges

CharRanges CharRanges  $\rightarrow$  CharRanges **{right}**

"(" CharRanges ")"  $\rightarrow$  CharRanges **{bracket}**

$\rightarrow$  OptCharRanges

CharRanges  $\rightarrow$  OptCharRanges

"[" OptCharRanges "]"  $\rightarrow$  CharClass

"~" CharClass  $\rightarrow$  CharClass

CharClass "/" CharClass  $\rightarrow$  CharClass **{left}**

CharClass "^" CharClass  $\rightarrow$  CharClass **{left}**

CharClass "v" CharClass  $\rightarrow$  CharClass **{left}**

"(" CharClass ")"  $\rightarrow$  CharClass **{bracket}**

**priorities**

"~" CharClass  $\rightarrow$  CharClass > CharClass "/" CharClass  $\rightarrow$  CharClass

> CharClass "^" CharClass  $\rightarrow$  CharClass >

CharClass "v" CharClass  $\rightarrow$  CharClass

**variables**

"cr"[0-9]\*  $\rightarrow$  CharRange

"cr"\*[0-9]\*  $\rightarrow$  OptCharRanges

"cr"+"[0-9]\*  $\rightarrow$  CharRanges

"cc"[0-9]\*  $\rightarrow$  CharClass

*Syntax* The kernel formalism is extended by adding character classes as symbols.

**module** CC-Sdf-Syntax

**imports** Character-Class-Syntax<sup>2.4.2</sup> Kernel-Sdf-Syntax<sup>2.3.1</sup>

**exports**

**context-free syntax**

CharClass  $\rightarrow$  Symbol

*Normalization* Character classes can be normalized to a unique normal form by ordering the ranges such that all characters are translated to their numeric equivalent and such that smaller characters are before larger characters and by fusing adjacent or overlapping ranges. For example, the class [A-Z0-9\%] has

normal form `[\37\48-\57\65-\90]`, because `\37` is the numerical representation of `\%`, `\48-\57` of `0-9`, `\65-\90` of `A-Z` and these do not overlap and are ordered. This normalization is specified in module `Character-Class-Normalization` that can be found in Visser (1997b).

```
module CC-Sdf-Normalization
imports CC-Sdf-Syntax2.4.2 Character-Class-Normalization
        Kernel-Sdf-Normalization2.3.3
```

### 2.4.3 Literals

Literals are abbreviations for fixed lists of characters. For example, the following production uses literals to define the keywords of a conditional statement.

```
"if" Exp "then" Stat "else" Stat -> Stat
```

The meaning of literals is expressed by means of a production that specifies the sequence of characters that makes up the literal. For instance, the meaning of the literals above is expressed by the productions

```
[\105] [\102]           -> "if"
[\116] [\104] [\101] [\110] -> "then"
[\101] [\108] [\115] [\101] -> "else"
```

Literals that are identifiers starting with a lowercase letter can be specified without the double quotes.

Another useful abbreviation in this category is the definition of the syntax of prefix functions in the form

```
add(Nat, Nat) -> Nat
```

as an abbreviation of

```
"add" "(" Nat "," Nat ")" -> Nat
```

*Syntax* Literals consist of a list of characters between double quotes. For the complete syntax of literals see §A.1. Literals that start with a lowercase letter can be written without quotes, hence the name ‘unquoted literals’. Prefix functions can be declared by means of a special form of productions, where the double quotes for the parentheses and commas can be omitted.

```
module Literals-Sdf-Syntax
imports Kernel-Sdf-Syntax2.3.1 Literals4.1
exports
  sorts UQLiteral
  lexical syntax
    [a-z]                -> UQLiteral
    [a-z][A-Za-z0-9\-*][A-Za-z0-9] -> UQLiteral
  context-free syntax
    UQLiteral             -> Literal
    Literal               -> Symbol
    Literal "(" {Symbol ","}* ")" "→" Symbol Attributes -> Production
```

*Normalization* The normalization function ‘ $\mathcal{I}$ ’ generates a defining production for each literal that is used as a symbol in one of the productions of the grammar.

```

module Literals-Sdf-Normalization
imports Literals-Sdf-Syntax2.4.3 CC-Sdf-Normalization2.4.2
exports
  context-free syntax
    “[” Grammar “]”      → Grammar
    literals(SymbolSet)   → Productions
    chars(Literal)        → Symbols
    symbols({Symbol “,”}*) → Symbols
  variables
    “c”[0-9']* → CHAR
    “c”“+”[0-9']* → CHAR+
    “c”“*”[0-9']* → CHAR*
    “L”        → Literal
    “γ”[0-9']* → {Symbol “,”}*
    “γ”“+”[0-9']* → {Symbol “,”}+
  equations
  Unquoted literals are translated to quoted literals.

```

$$[1] \quad \text{uqliteral}(c^+) = \text{literal}(\text{""} c^+ \text{""})$$

The function  $\mathcal{I}[\_]$  generates a production for each literal symbol in the grammar. The production generated for a literal  $L$  has the form  $\alpha \rightarrow L$ , where  $\alpha$  is a list of singleton character classes representing the characters of  $L$ . This list is produced by the function ‘chars’.

$$\begin{aligned}
 [2] \quad & \mathcal{I}[\mathcal{G}] = \mathcal{G} \text{ syntax literals}(\text{symbols}(\mathcal{G})) \\
 [3] \quad & \text{literals}(\{L\}) = \text{chars}(L) \rightarrow L \\
 [4] \quad & \text{literals}(\{\mathcal{A}\}) = \text{otherwise} \\
 [5] \quad & \text{literals}(\{\}) = \\
 [6] \quad & \text{literals}(\{\alpha^+ \beta^+\}) = \text{literals}(\{\alpha^+\}) ++ \text{literals}(\{\beta^+\})
 \end{aligned}$$

The function ‘chars’ scans the characters in the literal string, translating them to short characters. These are then normalized to numeric character codes by character normalization. The third equation tries if the first character of the string is a short-character by normalizing it and then testing whether it has reduced to a numeric character. This works for letters and digits. If this fails, the fourth equation translates the character to an escaped short-character, which succeeds for all other characters. Characters that are already escaped are handled by the second equation.

$$\begin{aligned}
 [7] \quad & \text{chars}(\text{""}) = \\
 [8] \quad & \frac{[\text{shortchar}(\text{"\" } c)] = cc}{\text{chars}(\text{literal}(\text{""} \text{"\" } c c^* \text{""})) = cc ++ \text{chars}(\text{literal}(\text{""} c^* \text{""}))} \\
 [9] \quad & \frac{[\text{shortchar}(c)] = cc, \quad cc = [\text{numchar}(c^+)]}{\text{chars}(\text{literal}(\text{""} c c^* \text{""})) = cc ++ \text{chars}(\text{literal}(\text{""} c^* \text{""}))}
 \end{aligned}$$

$$\begin{array}{l}
[10] \quad \frac{\begin{array}{l} [\text{shortchar}("\backslash" c)] = cc, \\ \alpha = \text{chars}(\text{literal}(""" c^* """)) \end{array}}{\text{chars}(\text{literal}(""" c c^* """)) = cc ++ \alpha} \quad \text{otherwise}
\end{array}$$

Prefix function productions are translated to normal productions by enclosing the parentheses and commas in double quotes.

$$\begin{array}{ll}
[11] & L(\gamma) \rightarrow \mathcal{A} \$ = L "(" ++ \text{symbols}(\gamma) ++ ")" \rightarrow \mathcal{A} \$ \\
[12] & \text{symbols}() = \\
[13] & \text{symbols}(\mathcal{A}) = \mathcal{A} \\
[14] & \text{symbols}(\gamma_1^+, \gamma_2^+) = \text{symbols}(\gamma_1^+) ++ ", " ++ \text{symbols}(\gamma_2^+)
\end{array}$$

## 2.5 Parse Trees

Now we can define the interpretation of grammars, that is, the well-formed trees characterized by a grammar and the yield of those trees. The general idea is that a context-free production  $p = \mathcal{A}_1 \dots \mathcal{A}_n \rightarrow \mathcal{A}_0$  constructs trees of type  $\mathcal{A}_0$  labeled with the production  $p$  and with a list of direct descendants of type  $\mathcal{A}_1 \dots \mathcal{A}_n$ . Such trees are represented by means of terms. The constructor ‘appl’ builds an application of a production  $p$  to a list of trees, i.e., if  $T_1 \dots T_n$  are trees of type  $\mathcal{A}_1 \dots \mathcal{A}_n$  then  $\text{appl}(p, [T_1, \dots, T_n])$  is a tree of type  $\mathcal{A}_0$ . Parse forests are constructed by representing choice nodes or *ambiguity nodes* by means of the constructor ‘amb’. If  $T_1 \dots T_n$  are all trees of the same type  $\mathcal{A}$ , then  $\text{amb}([T_1, \dots, T_n])$  is an ambiguity node of type  $\mathcal{A}$ .

To formally define this notion of trees we introduce the notion of terms. We first present the generic term format that is used to encode parse trees and the encoding of symbols and grammars in that format. With those tools in place, we define the well-formedness rules in §2.5.5.

### 2.5.1 Term Format

Van den Brand *et al.* (1997) introduce the generic, annotated term format ATerms for the representation and exchange of structured data. The format is designed such that all kinds of data can be represented in a single, fixed format, with the purpose of exchanging such data between tools and providing generic operations on these data. The definition of the format comes with an extensive library of (higher-order) functions. We will use the ATerm format to represent parse trees.

The syntax of *ATerms* is defined in module ATerms below. Terms are constructed by means of four constructors, i.e., an ATerm is one of the following:

- A *constant* (ACon), which is either an integer constant or a real number constant.
- A *list* of terms (ATermList), which is either empty [], or a list of one or more terms separated by commas between square brackets  $[T_1, \dots, T_n]$ . The sort ATerms represents lists of one or more terms separated by commas.
- A *function symbol* (AFun).

- An *application* of a function symbol to a list of one or more terms separated by commas.

Furthermore, each of these constructors can be *annotated* by a list of one or more terms between { and } (Ann). Literals are strings of characters between double quotes. Integer constants are lists of digits and real constants are floating point numbers with an optional exponent. For the syntax of literals, integers and reals see Van den Brand *et al.* (1997).

```

module ATerms
imports LiteralsA.1 IntCon RealCon
exports
  sorts ATerms ATermList ACon AFun ATerm Ann
  context-free syntax
    ATerm                → ATerms
    ATerm “,” ATerms      → ATerms
    “[” “[”              → ATermList
    “[” ATerms “[”       → ATermList
    IntCon                → ACon
    RealCon               → ACon
    Literal               → AFun
    ACon                  → ATerm
    ATermList             → ATerm
    AFun                  → ATerm
    AFun “(” ATerms “)”   → ATerm
    “{” ATerms “}”       → Ann
    ACon Ann              → ATerm
    ATermList Ann         → ATerm
    AFun Ann              → ATerm
    AFun “(” ATerms “)” Ann → ATerm
  variables
    “Ts”[0-9']*          → ATerms
    “Tl”[0-9']*          → ATermList
    “ACon”[0-9']*        → ACon
    “AFun”[0-9']*        → AFun
    “T”[0-9']*           → ATerm
    “Ann”[0-9']*         → Ann

```

### 2.5.2 Constructors for Parse Trees

Function symbols can be literals—strings of characters between double quotes—or identifiers. Specification of the identifiers is not included in the ATerm format. For each application of ATerms, an appropriate set of AFuns should be declared, with the requirement that they are restricted to names of the form [a-z][a-zA-Z0-9\-\*]. For the representation of grammars, symbols, productions and trees we define the following function symbols.

```

module Kernel-Sdf-Tree-Constructors
imports Grammar-Tree-ConstructorsA.2.1
exports
  context-free syntax
    “prod” → AFun

```



```

“no-attrs” → AFun
“attrs”    → AFun
“atr”      → AFun
“syntax”   → AFun

“appl”     → AFun
“amb”      → AFun

```

The function symbols ‘appl’ and ‘amb’ will be used to represent parse trees. The others will be used in the encoding of grammar structures. Each extension of the kernel that adds new constructors for symbols or grammars should also add the corresponding ATerm function symbols. These are included in §A.2.

### 2.5.3 ATerm Encoding

Now we can encode symbols, grammars and productions as terms. For each sort  $\mathcal{S}$  a function  $\text{aterm}(\mathcal{S}) \rightarrow \text{ATerm}$  is defined that encodes  $\mathcal{S}$ -expressions as ATerms. The encoding is injective. For each sort  $\mathcal{S}$  a decoding function  $\text{s}(\text{ATerm}) \rightarrow \mathcal{S}$  is defined such that  $\text{s}(\text{aterm}(s)) = s$ . Figure 2.1 illustrates the encoding of symbols and productions as ATerms.

The following module defines the encoding of lists of symbols. The encoding of constructors for symbols is defined in a module for each extension of the kernel; see §A.2.

```

module Symbols-ATerms
imports Symbols-Projection2.1.1 ATerm-Lists
exports
  context-free syntax
    aterm(Symbol)      → ATerm
    atermList(Symbols) → ATermList
    symbol(ATerm)      → Symbol
    symbols(ATermList) → Symbols

```

#### equations

Encoding lists of symbols.

- [1]  $\text{atermList}(\mathcal{A}) = [\text{aterm}(\mathcal{A})]$
- [2]  $\text{atermList}(\alpha) = []$  **when**  $\alpha =$
- [3]  $\text{atermList}(\alpha^+ \beta^+) = \text{atermList}(\alpha^+) ++ \text{atermList}(\beta^+)$

Decoding lists of symbols.

- [4]  $\text{symbols}([]) =$
- [5]  $\text{symbols}([T]) = \text{symbol}(T)$
- [6]  $\text{symbols}([T, Ts]) = \text{symbol}(T) ++ \text{symbols}([Ts])$

where we have the requirement that

- [7]  $\text{symbol}(\text{aterm}(\mathcal{A})) = \mathcal{A}$

The last equation requires of each future definition that it should be such that the decoding of an encoded symbol gives the original symbol.

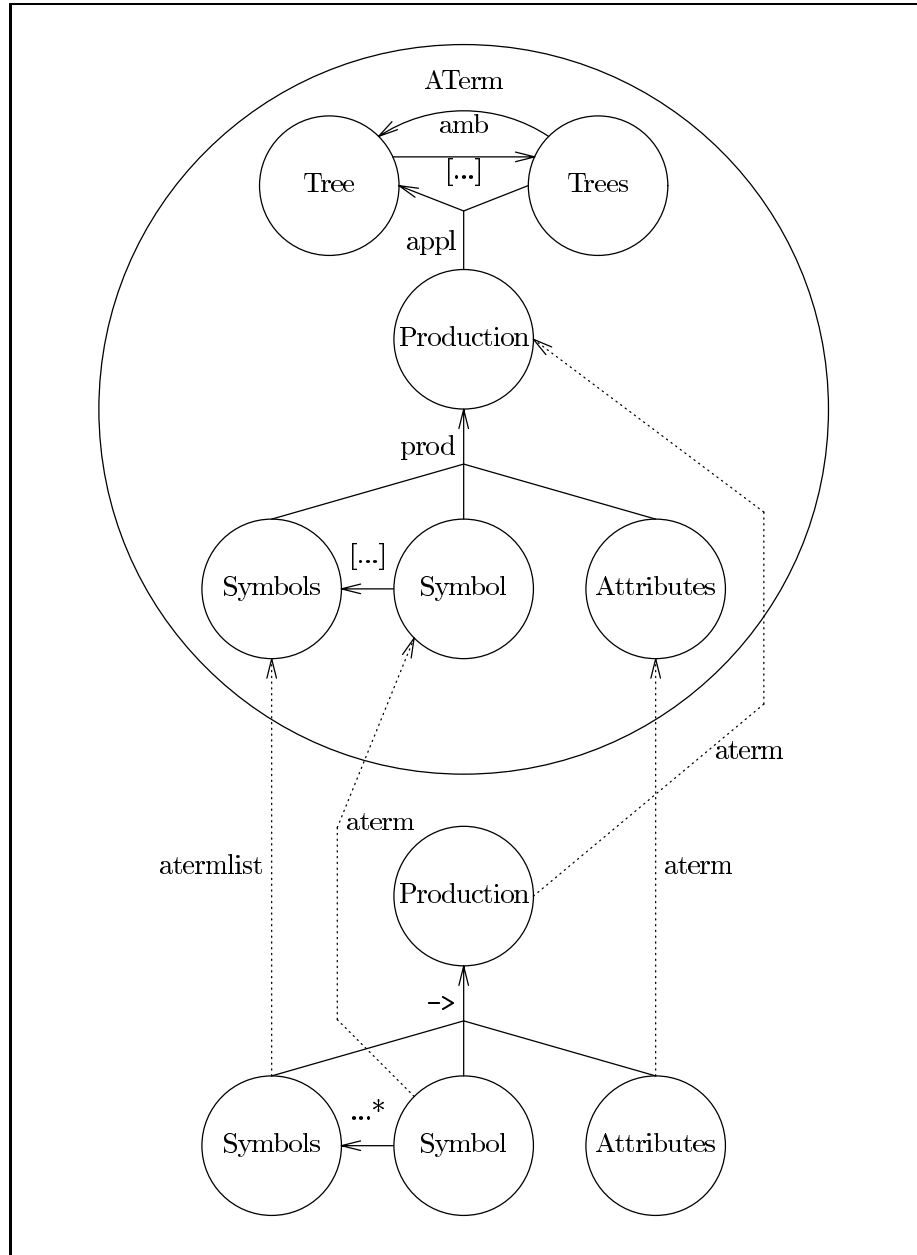


Figure 2.1: Encoding symbols, productions and parse trees in a fixed term format. Grammar domains such as Symbol and Production are mapped (the dotted arrows) onto subsets of the set of ATerms. Parse Trees are another subset of ATerms formed by the constructor ‘appl’ from a production and a list of trees, or by the constructor ‘amb’ from a list of trees.

### 2.5.4 Encoding Productions

The following module defines the encoding of productions. As an example, consider the production

Exp "+" Exp  $\rightarrow$  Exp {left}

which is encoded as

```
prod([sort("Exp"),lit("+"),sort("Exp")],
     sort("Exp"),attrs([atr("left")]))
```

A production is represented by the function symbol ‘prod’ and the attributes of the production are represented by ‘attrs’. Note that this makes use of the encoding of sort and literal symbols that is defined in §A.2.

**module** Kernel-Sdf-ATerms

**imports** Kernel-Sdf-Projection<sup>2.3.2</sup> Kernel-Sdf-Tree-Constructors<sup>2.5.2</sup>

Symbols-ATerms<sup>2.5.3</sup> Grammar-ATerms<sup>4.2.2</sup> ATerm-Lists

**exports**

**context-free syntax**

aterm(Production)	$\rightarrow$ ATerm
aterm(Attributes)	$\rightarrow$ ATerm
atermlist({Attribute “,”}*)	$\rightarrow$ ATermList
aterm(Attribute)	$\rightarrow$ ATerm
atermlist(Productions)	$\rightarrow$ ATermList
production(ATerm)	$\rightarrow$ Production
attributes(ATerm)	$\rightarrow$ Attributes
attribute(ATerm)	$\rightarrow$ Attribute
productions(ATermList)	$\rightarrow$ Productions

**equations**

Encoding productions and attributes.

- [1]  $\text{aterm}(\alpha \rightarrow \mathcal{A} \$) = \text{prod}(\text{atermlist}(\alpha), \text{aterm}(\mathcal{A}), \text{aterm}(\$))$
- [2]  $\text{aterm}() = \text{no-attrs}$
- [3]  $\text{aterm}(\{attr^*\}) = \text{attrs}(\text{atermlist}(attr^*))$
- [4]  $\text{atermlist}(attr^*) = [] \text{ when } \{attr^*\} = \{\}$
- [5]  $\text{atermlist}(attr) = [\text{aterm}(attr)]$
- [6]  $\text{atermlist}(attr_1^+, attr_2^+) = \text{atermlist}(attr_1^+) ++ \text{atermlist}(attr_2^+)$

Decoding productions and attributes.

- [7]  $\text{production}(\text{prod}(Tl, T, T')) = \text{symbols}(Tl) \rightarrow \text{symbol}(T) \text{ attributes}(T')$
- [8]  $\text{attributes}(\text{no-attrs}) =$
- [9]  $\text{attributes}(\text{attrs}(Tl)) = \text{attributes}(Tl)$
- [10]  $\text{attributes}([]) = \{\}$
- [11]  $\text{attributes}([T]) = \{\text{attribute}(T)\}$
- [12]  $\text{attributes}([T, Ts]) = \{\text{attribute}(T)\} ++ \text{attributes}([Ts])$

where we have the requirement that

- [13]  $\text{attribute}(\text{aterm}(attr)) = attr$

Encoding grammars and lists of productions.

- [14]  $\text{aterm}(\text{syntax } p^*) = \text{syntax}(\text{atermlist}(p^*))$
- [15]  $\text{atermlist}(p) = [\text{aterm}(p)]$
- [16]  $\text{atermlist}(p^*) = [] \text{ when } p^* =$
- [17]  $\text{atermlist}(p_1^+ p_2^+) = \text{atermlist}(p_1^+) ++ \text{atermlist}(p_2^+)$

Decoding of grammars and lists of productions.

- [18]  $\text{grammar}(\text{syntax}(Tl)) = \text{syntax productions}(Tl)$
- [19]  $\text{productions}([]) =$
- [20]  $\text{productions}([T]) = \text{production}(T)$
- [21]  $\text{productions}([T, Ts]) = \text{production}(T) ++ \text{productions}([Ts])$

### 2.5.5 Well-formed Parse Trees

Now we have prepared all equipment for the characterization of the terms that represent the well-formed parse trees over a grammar. The predicate  $\mathcal{G} \vdash T$  determines whether a tree  $T$  is well-formed with respect to grammar  $\mathcal{G}$ . This is defined in terms of  $\vdash T : T'$ , which checks whether  $T$  is a tree of type  $T'$ . Since a tree contains all type information explicitly in the form of productions, this can be checked without reference to the grammar. The well-formedness with respect to the grammar is then defined by checking that all productions in a tree are actually productions of the grammar.

The main constructor for trees is the function ‘`appl`’ that creates an *application* of a context-free production to a list of trees such that the types of the argument trees correspond to the symbols in the left-hand side of the production. As an example, consider the grammar

```

sorts E
syntax
  [a-z]    -> E
  [\+]    -> "+"
  E "+" E -> E {left}

```

The following ATerm is a well-formed parse tree over this grammar for the sentence `a+b`.

```

appl(prod([sort("E"),lit("+"),sort("E")],
          sort("E"), attrs([atr("left")])),
      [appl(prod([char-class([range(97,122)])],
                sort("E"),no-attrs),
            [97]),
        appl(prod([char-class([43])],
                  lit("+"),no-attrs),
            [43]),
        appl(prod([char-class([range(97,122)])],
                  sort("E"),no-attrs),
            [98])])])

```

Observe that the main `appl` has the `aterm` encoding of the production `E "+" E -> E {left}` as first argument and as second argument a list of three trees with types that correspond to the arguments of that production. The numbers at the leafs of the trees denote the ASCII values of the characters, i.e., 97 denotes `a`, 43 denotes `+` and 98 denotes `b`.

Context-free grammars can be ambiguous, i.e., generate more than one tree for a single sentence. The constructor ‘`amb`’ is introduced to represent parse forests, i.e., compact representations of sets of parse trees. A term `amb([ $T_1, \dots, T_n$ ])` represents the set of parse trees containing the terms  $T_1, \dots, T_n$  (which can again contain ‘`amb`’ nodes). For example, the string `E "*" E "+" E` is ambiguous with respect to the grammar

```
sorts E
syntax
  E "+" E -> E
  E "*" E -> E
```

and therefore has the following parse forest that represents the two possible parses, left-associative (`E "*" E`) `"+" E` and right associative `E "*" (E "+" E)`.

```
amb([appl(prod([sort("E"),lit("+"),sort("E")],
               sort("E"),no-attrs),
  [appl(prod([sort("E"),lit("*"),sort("E")],
             sort("E"),no-attrs),
        [sort("E"),
         lit("*"),
         sort("E")]),
    lit("+"),
    sort("E")]),
  appl(prod([sort("E"),lit("*"),sort("E")],
            sort("E"),no-attrs),
    [sort("E"),
     lit("*"),
     appl(prod([sort("E"),lit("+"),sort("E")],
               sort("E"),no-attrs),
        [sort("E"),
         lit("+"),
         sort("E")])])])])
```

Note that in order to reduce the size of the term, the subtrees for `E` and `"*"` and `"+"` are symbols. The tree represents the parse tree for a sentential form.

The yield of a tree is the concatenation of the characters at the leafs of the tree. For instance, the yield of the first tree above is `[97] [43] [98]`, i.e., the list of characters `a+b`.

```
module Kernel-Sdf-Trees
imports Kernel-Sdf-ATerms2.5.4 Kernel-Sdf-Projection2.3.2
exports
  context-free syntax
    Grammar "⊢" ATerm → Bool
    "⊢" ATerm ":" ATerm → Bool

  yield "[[" Grammar "]" "(" ATerm ")" → Symbols
```

args(ATerm)	→ ATermList
type(ATerm)	→ ATerm
prods(ATerm)	→ Productions

**variables**

"Prod"[0-9']*	→ ATerm
"Res"[0-9']*	→ ATerm
"Attrs"[0-9']*	→ ATerm
"Args"[0-9']*	→ ATermList

**equations**

A term  $T$  is a well-formed parse tree *over grammar*  $\mathcal{G}$ , if it is a well-formed tree and if all its productions are productions of  $\mathcal{G}$ .

$$[1] \quad \frac{\vdash T : \text{type}(T) \wedge \text{prods}(T) \subseteq P(\mathcal{G}) = \top}{\mathcal{G} \vdash T = \top}$$

Otherwise, the tree is not well-formed.

$$[2] \quad \mathcal{G} \vdash T = \perp \quad \textbf{otherwise}$$

The function ‘prods’ gives the productions of a term, i.e., the list of all productions used in its applications.

$$\begin{aligned}
[3] \quad & \text{prods}(\text{appl}(\text{Prod}, \text{Args})) = \text{production}(\text{Prod}) \uplus \text{prods}(\text{Args}) \\
[4] \quad & \text{prods}(\text{amb}(\text{Args})) = \text{prods}(\text{Args}) \\
[5] \quad & \text{prods}([\ ] ) = \\
[6] \quad & \text{prods}([T]) = \text{prods}(T) \\
[7] \quad & \text{prods}([T, Ts]) = \text{prods}(T) \uplus \text{prods}([Ts]) \\
[8] \quad & \text{prods}(T) = \textbf{otherwise}
\end{aligned}$$

For the definition of  $\mathcal{G} \vdash T$  we need several auxiliary functions on terms. The function ‘args’ gives the arguments of a production. The ‘type’ of a production is its result type. The ‘type’ of an application is the result type of its production.

$$\begin{aligned}
[9] \quad & \text{args}(\text{prod}(Tl, Res, Attrs)) = Tl \\
[10] \quad & \text{args}(\text{appl}(\text{Prod}, \text{Args})) = \text{Args} \\
[11] \quad & \text{type}(\text{prod}(Tl, Res, Attrs)) = Res \\
[12] \quad & \text{type}(\text{appl}(\text{Prod}, \text{Args})) = \text{type}(\text{Prod}) \\
[13] \quad & \text{type}(\text{amb}([Ts])) = \text{type}(\text{first}([Ts]))
\end{aligned}$$

An application is a well-formed term of type  $T$  if the type of its production is  $T$  and if its arguments are well-formed terms with types that correspond to the argument types of the production of the application.

$$[14] \quad \frac{\text{type}(\text{Prod}) = T, \vdash \text{Args} : \text{args}(\text{Prod}) = \top}{\vdash \text{appl}(\text{Prod}, \text{Args}) : T = \top}$$

A list of trees or symbols is well-formed if each element is well-formed.

$$\begin{aligned}
[15] \quad & \vdash [\ ] : [\ ] = \top \\
[16] \quad & \vdash [T] : [T'] = \vdash T : T' \\
[17] \quad & \vdash [T, Ts] : [T', Ts'] = \vdash T : T' \wedge \vdash [Ts] : [Ts']
\end{aligned}$$

An ambiguous node is well-formed if all possibilities have the same type.

$$\begin{array}{l}
 [18] \quad \vdash \text{amb}([T]) : T' = \vdash T : T' \\
 [19] \quad \vdash \text{amb}([T, Ts]) : T' = \vdash T : T' \wedge \vdash \text{amb}([Ts]) : T'
 \end{array}$$

A symbol term  $T$  is a well-formed tree of type  $T$ . This is used to represent trees for sentential forms.

$$[20] \quad \frac{\text{symbol}(T) = \text{symbol}(T')}{\vdash T : T' = \top}$$

The yield of a term is the concatenation of all its leaf symbols.

$$\begin{array}{l}
 [21] \quad \text{yield}[\mathcal{G}](\text{appl}(\text{Prod}, Tl)) = \text{yield}[\mathcal{G}](Tl) \\
 [22] \quad \text{yield}[\mathcal{G}](\text{amb}(Tl)) = \text{yield}[\mathcal{G}](\text{first}(Tl)) \\
 [23] \quad \text{yield}[\mathcal{G}]([\ ] ) = \\
 [24] \quad \text{yield}[\mathcal{G}]([T]) = \text{yield}[\mathcal{G}](T) \\
 [25] \quad \text{yield}[\mathcal{G}]([T, Ts]) = \text{yield}[\mathcal{G}](T) ++ \text{yield}[\mathcal{G}]([Ts]) \\
 [26] \quad \text{yield}[\mathcal{G}](T) = \text{symbol}(T) \quad \textbf{otherwise}
 \end{array}$$

### 2.5.6 Trees with Characters and Literals

In the previous section we defined character classes and literals as symbols in grammars. Since the meaning of literals is defined in terms of character classes by means of context-free productions, the definition of well-formedness of parse trees does not need to be extended for literals. However, for character classes we have to extend the definition such that a character is a tree with as type any character class that contains it, i.e., a parse tree is a well-formed tree of a character class type if it is a character (represented by a natural number) that is an element of the character class. For example, the term

```

appl(prod([sort("Id"),char-class([range(97,122)])]),
      sort("Id"),no-attrs),
[appl(prod([char-class([range(97,122)])]),
        sort("Id"),no-attrs),
 [97]),
 98])

```

is a well-formed parse tree for the identifier **ab**. The definition of well-formedness and yield is extended as follows:

```

module CC-Sdf-Trees
imports CC-Sdf-ATermsA.2.2 Kernel-Sdf-Trees2.5.5
        Character-Class-Normalization
equations

```

A character is represented by an integer. The characteristic functions for trees are extended for this new tree constructor. The type of a character is the character code itself. The yield of a character is a character class containing the single character. A character does not contain any productions.

$$[1] \quad \text{type}(n) = n$$

$$\begin{array}{l}
[2] \quad \text{yield}[\mathcal{G}](n) = [\text{character}(n)] \\
[3] \quad \text{prods}(n) =
\end{array}$$

A character code  $n$  is a well-formed tree of type  $T$  if  $T$  represents a character class that contains the character corresponding to  $n$ .

$$[4] \quad \frac{\text{symbol}(T) = cc, \text{character}(n) \in cc = \top}{\vdash n : T = \top}$$

### 2.5.7 Cyclic Parse Forests

Some grammars generate infinitely many parse trees for a single string. For instance, the grammar

```

syntax
  -> S
  [a] -> S
  S S -> S

```

generates infinitely many trees for the string **a**. The collection of parse trees for strings over such grammars can be finitely represented by means of a cyclic parse forest (Billot and Lang, 1989, Rekers, 1992). However, in the term format defined here we have no provisions for cyclic forests. The parser for SDF2 in Visser (1997e) *does* generate a cyclic parse forest.

In ASF+SDF cyclic structures cannot be expressed in a natural way. This can be simulated by explicitly representing the pointer structure by means of a table of references and tags to represent these references. But this would complicate the entire specification. Since the application of SDF2 will be mainly to non-cyclic grammars we have not gone into the trouble of defining cyclic forests in ASF+SDF.

### 2.5.8 Equality of Trees

We define an equality predicate  $\doteq$  on parse trees. At this point this comes down to syntactic equality. Later on we will extend the predicate such that trees that are not syntactically equal can be equal. This will be useful to abstract from certain details in parse trees. For instance, in §3.3 we will introduce parse trees containing layout. This is useful for applications that are aware of layout. But most applications will want to abstract from the specific layout in a tree and consider two trees equal up to layout. Another application is the equality of trees with associative operators such as list concatenation. The definition of the equality predicate below is intended to specify the details of such equality considerations. Furthermore, we define membership of a tree in a parse forest.

```

module Kernel-Sdf-Equality
imports Kernel-Sdf-Trees2.5.5
exports
  context-free syntax
    ATerm “ $\doteq$ ” ATerm  $\rightarrow$  Bool
    ATerm “ $\in$ ” ATerm  $\rightarrow$  Bool

```



**equations**

Equality of applications. If the productions are the same two applications are equal if the argument lists are.

$$[1] \quad \frac{Args \doteq Args' = \top}{\text{appl}(Prod, Args) \doteq \text{appl}(Prod, Args') = \top}$$

Argument lists are equal if the elements are pairwise equal.

$$\begin{aligned} [2] \quad & [] \doteq [] = \top \\ [3] \quad & [T] \doteq [T'] = T \doteq T' \\ [4] \quad & [T, Ts] \doteq [T', Ts'] = T \doteq T' \wedge [Ts] \doteq [Ts'] \end{aligned}$$

An ambiguity node is equal to a tree if it all its possibilities are contained in the tree and vice versa.

$$[5] \quad \frac{\text{amb}([Ts]) \dot{\in} T \wedge T \dot{\in} \text{amb}([Ts]) = \top}{\text{amb}([Ts]) \doteq T = \top}$$

$$[6] \quad \frac{\text{amb}([Ts]) \dot{\in} T \wedge T \dot{\in} \text{amb}([Ts]) = \top}{T \doteq \text{amb}([Ts]) = \top}$$

If none of the cases above apply the terms are not equal.

$$[7] \quad T_1 \doteq T_2 = \perp \quad \textbf{otherwise}$$

A tree is member of a parse forest (tree containing ambiguities) if it is contained in one of the possibilities of an ambiguity.

$$[8] \quad \frac{T_1 \dot{\in} T_2 = \top}{T_1 \dot{\in} \text{amb}([T_2]) = \top}$$

$$[9] \quad \frac{T_1 \dot{\in} T_2 \vee T_1 \dot{\in} \text{amb}([Ts]) = \top}{T_1 \dot{\in} \text{amb}([T_2, Ts]) = \top}$$

An ambiguity is contained in a forest if *all* its possibilities are contained in the forest.

$$[10] \quad \frac{T_1 \dot{\in} T_2 = \top}{\text{amb}([T_1]) \dot{\in} T_2 = \top}$$

$$[11] \quad \frac{T_1 \dot{\in} T_2 \wedge \text{amb}([Ts]) \dot{\in} T_2 = \top}{\text{amb}([T_1, Ts]) \dot{\in} T_2 = \top}$$

An application is contained in an application, if the arguments of the first are contained in the arguments of the second.

$$[12] \quad \frac{Args_1 \dot{\in} Args_2 = \top}{\text{appl}(Prod, Args_1) \dot{\in} \text{appl}(Prod, Args_2) = \top}$$

Lists

$$[13] \quad [] \dot{\in} [] = \top$$

$$\begin{aligned}
 [14] \quad & [T] \dot{\in} [T'] = T \dot{\in} T' \\
 [15] \quad & [T, Ts] \dot{\in} [T', Ts'] = T \dot{\in} T' \wedge [Ts] \dot{\in} [Ts']
 \end{aligned}$$

If none of the cases above apply membership does not hold

$$[16] \quad T_1 \dot{\in} T_2 = \perp \quad \textbf{otherwise}$$

---

# 3

## Disambiguation and Abbreviation

---

In this chapter we present features for disambiguation of ambiguous grammars and abbreviation of common patterns. Priorities are used to disambiguate ambiguous expression syntax, providing support for compact abstract syntax. Priorities are defined by means of an extension of the well-formedness predicate on parse trees. Regular expressions abbreviate common patterns of productions such as lists, optional constructs, alternatives, etc. Regular expressions are defined by generating the defining productions for each expression in the grammar. Lexical and context-free syntax sections separate the definition of tokens and phrases. These are integrated into a single context-free grammar by normalization such that no interference between the two levels is created. Furthermore, the definition of the placement of layout between tokens is handled by this normalization. Follow restrictions and reject productions are provided to express lexical disambiguation rules such as ‘prefer longest match’ and ‘prefer literals’.

---

### 3.1 Priorities

Context-free grammars can be ambiguous. There are many methods for the disambiguation of context-free grammars. Most programming language oriented formalisms provide some kind of precedence based method. Here we adopt the method of disambiguation by associativity and priority as used in SDF. New with respect to the design of priorities in Heering *et al.* (1989) is (a) disambiguation of lexical syntax by lexical priorities, (b) a more uniform notation for priority declarations, and (c) derivation of productions from priority declarations, which provides a more compact notation by avoiding multiple declarations of productions. A feature not present in SDF2 is the abbreviation of productions in priority declarations by the list of literals of the left-hand side. (For example, `"*" > "+"` as an abbreviation of `E "*" E -> E > E "+" E -> E`.) The reason for this omission is the unclear semantics in combination with modularization. Also `<` priority-chains are not included in SDF2 because these can also be expressed using `>` chains.

We first define syntax, projection functions and normalization of priority declarations. In §3.1.4 we describe an extension of the well-formedness predicate on parse trees that characterizes parse trees without priority conflicts.

**Example 1.1** The following grammar defines priority and associativity relations over the syntax of expressions with unary negation and binary operators

for exponentiation, multiplication, addition and subtraction. Note that, unlike in SDF, the syntax for arithmetic operators can be defined *inside* the priorities section.

```

syntax
  "(" E ")" -> E {bracket}
priorities
  "-" E -> E
> E "^" E -> E {right}
> E "*" E -> E {left}
> {left:
  E "+" E -> E {assoc}
  E "-" E -> E {left}}

```

This grammar declares that unary  $-$  has *higher priority* than  $\wedge$ , which has higher priority than binary  $*$ , which has higher priority than  $+$  and binary  $-$ . The latter two are mutually left associative as declared by the *group associativity*. The bracket production declares that parentheses can also be used to disambiguate expressions. With this grammar the expression  $E--E*E+E-E^E$  should be interpreted as  $((E-((-E)*E))+E)-(E^E)$ .  $\square$

### 3.1.1 Syntax

The priorities section of a grammar defines the priority relation  $>$  on productions and the associativity relations ‘left’, ‘right’, ‘assoc’, and ‘non-assoc’. A priority declaration is either a  $>$  chain or an associativity declaration. The objects of these declarations are single productions or groups of productions. A group can have an X-associativity, which declares the productions in the group to be mutually X-associative.

The ‘bracket’ attribute declares a production of the form  $l \mathcal{A} r \rightarrow \mathcal{A} \{\text{bracket}\}$ , with ‘l’ (‘r’) denoting the syntax for the left- (right-) bracket, to be the identity function on  $\mathcal{A}$ . Such productions can be used to explicitly disambiguate some text or to indicate a different disambiguation than the one given by priority rules.

```

module Priority-Sdf-Syntax
imports Kernel-Sdf-Syntax2.3.1

```

```

exports

```

```

  sorts Associativity Group Priority Priorities

```

```

  context-free syntax

```

“left”	→ Associativity
“right”	→ Associativity
“non-assoc”	→ Associativity
“assoc”	→ Associativity
“bracket”	→ Attribute
Associativity	→ Attribute
Production	→ Group
“{” Productions “}”	→ Group
“{” Associativity “:” Productions “}”	→ Group
{Group “>”}+	→ Priority
Group Associativity Group	→ Priority
{Priority “,”}*	→ Priorities

“priorities” Priorities  $\rightarrow$  Grammar  
**variables**  
 “g”[0-9']\*  $\rightarrow$  Group  
 “gg\*”[0-9']\*  $\rightarrow$  {Group “>”}\*  
 “gg+”[0-9']\*  $\rightarrow$  {Group “>”}+  
 “pr”[0-9']\*  $\rightarrow$  Priority  
 “pr” “\*”[0-9']\*  $\rightarrow$  {Priority “,”}\*  
 “pr” “+”[0-9']\*  $\rightarrow$  {Priority “,”}+  
 “as”[0-9']\*  $\rightarrow$  Associativity

### 3.1.2 Projection

The projection function ‘Pr’ yields the list of all priority declarations of a grammar. The projection function ‘ $\overline{\text{Pr}}$ ’ yields the grammar without its priority declarations.

**module** Priority-Sdf-Projection  
**imports** Priority-Sdf-Syntax<sup>3.1.1</sup> Kernel-Sdf-Projection<sup>2.3.2</sup>  
**exports**  
**context-free syntax**  
 Priorities “++” Priorities  $\rightarrow$  Priorities {**assoc**}  
 “Pr”(Grammar)  $\rightarrow$  Priorities  
 “ $\overline{\text{Pr}}$ ”(Grammar)  $\rightarrow$  Grammar  
 Priority “ $\in$ ” Priorities  $\rightarrow$  Bool

#### equations

Concatenation of priority declarations.

$$[1] \quad pr_1^* ++ pr_2^* = pr_1^*, pr_2^*$$

The priorities and non-priorities of a grammar.

$$\begin{aligned}
 [2] \quad & \text{Pr}(\text{priorities } pr^*) = pr^* \\
 [3] \quad & \text{Pr}(\mathcal{G}_1 \mathcal{G}_2) = \text{Pr}(\mathcal{G}_1) ++ \text{Pr}(\mathcal{G}_2) \\
 [4] \quad & \text{Pr}(\mathcal{G}) = \text{otherwise} \\
 [5] \quad & \overline{\text{Pr}}(\text{priorities } pr^*) = \emptyset \\
 [6] \quad & \overline{\text{Pr}}(\mathcal{G}_1 \mathcal{G}_2) = \overline{\text{Pr}}(\mathcal{G}_1) \overline{\text{Pr}}(\mathcal{G}_2) \\
 [7] \quad & \overline{\text{Pr}}(\mathcal{G}) = \mathcal{G} \text{ otherwise}
 \end{aligned}$$

Membership of a priority declaration. A pair is member of a declaration if the declaration contains a pair with similar productions. Recall from §2.3.2 that two productions are similar if they are the same except for their attributes, which may be different.

$$[8] \quad \frac{p_1 \cong p_3 = \top, p_2 \cong p_4 = \top}{p_1 > p_2 \in pr_1^*, p_3 > p_4, pr_2^* = \top}$$

$$[9] \quad \frac{p_1 \cong p_3 \wedge p_2 \cong p_4 \vee p_1 \cong p_4 \wedge p_2 \cong p_3 = \top}{p_1 \text{ as } p_2 \in pr_1^*, p_3 \text{ as } p_4, pr_2^* = \top}$$

$$[10] \quad pr \in pr^* = \perp \text{ otherwise}$$

### 3.1.3 Normalization

The complex syntax for priority declarations can be expressed by means of only binary declarations for the  $>$  relation and the associativity relations as follows: (1) Priority chains of the form  $p_1 > \dots > p_n$  are normalized to lists of simple priorities of the form  $p_i > p_{i+1}$ . This relation is closed transitively. (2) Associativity declarations in attributes and group associativities are expressed as binary associativity declarations of the form  $p_1 \text{ as } p_2$ . (3) The productions that are mentioned in priorities sections are added to the syntax section of the grammar.

**Example 1.2** The normalization of the grammar in Example 1.1 on page 39 is:

```

syntax
  "-" E    -> E
  E "^" E  -> E {right}
  E "*" E  -> E {left}
  E "+" E  -> E {assoc}
  E "-" E  -> E {left}
priorities
  "-" E -> E          >      E "^" E -> E {right},
  E "^" E -> E {right} right E "^" E -> E {right},
  E "^" E -> E {right} >      E "*" E -> E {left},
  E "*" E -> E {left} >      E "+" E -> E {assoc},
  E "*" E -> E {left} >      E "-" E -> E {left},
  E "*" E -> E {left} left  E "*" E -> E {left},
  E "+" E -> E {assoc} left  E "-" E -> E {left},
  E "+" E -> E {assoc} assoc E "+" E -> E {assoc},
  E "-" E -> E {left} left  E "-" E -> E {left}

```

Observe that all productions mentioned in the priorities declaration are also declared as productions in the ‘syntax’ part. Furthermore, the  $>$  chain is made into a binary relation, which is transitively closed. All associativity attributes are expressed by means of binary declarations.  $\square$

**module** Priority-Sdf-Normalization

**imports** Priority-Sdf-Syntax<sup>3.1.1</sup> Booleans Kernel-Sdf-Normalization<sup>2.3.3</sup>

Priority-Sdf-Projection<sup>3.1.2</sup>

**exports**

**context-free syntax**

```

  "p[" Grammar "]"      -> Grammar
  "assoc" "[" Grammar "]" -> Priorities
  "assoc" "[" Productions "]" -> Priorities
  "syntax" "[" Priorities "]" -> Grammar
  "norm" "[" Priorities "]" -> Priorities
  "trans" "[" Priorities "]" -> Priorities

```

**equations**

The normalization function  $p[\_]$  extracts syntax information from priorities and priority information from syntax, normalizes the priorities declarations and takes the transitive closure.

$$[1] \quad \frac{\text{Pr}(\mathcal{G}) = pr_1^*, \text{norm}[pr_1^*] = pr_2^*, \overline{\text{Pr}}(\mathcal{G}) \text{ syntax}[pr_2^*] = \mathcal{G}'}{p[\mathcal{G}] = \mathcal{G}' \text{ priorities trans}[pr_2^* \text{ ++ assoc}[\mathcal{G}']]}$$

The function  $\text{norm}[\_]$  normalizes a priority declaration to a list of pairs of the form  $p > p'$  or  $p \text{ as } p'$  by eliminating  $>$ -chains and  $\{\_ \}$  groups.

$$\begin{aligned} [2] \quad & \text{norm}[p > p'] = p > p' \\ [3] \quad & \text{norm}[p \text{ as } p'] = p \text{ as } p', p' \text{ as } p \end{aligned}$$

Each of the priority declarations in the list is normalized.

$$\begin{aligned} [4] \quad & \text{norm}[\_] = \\ [5] \quad & \text{norm}[p] = \\ [6] \quad & \text{norm}[pr_1^+, pr_2^+] = \text{norm}[pr_1^+] \text{ ++ } \text{norm}[pr_2^+] \end{aligned}$$

A  $>$  chain is broken into binary  $>$  declarations. The transitive closure defined below ensures that  $p_1 > p_3$  if  $p_1 > p_2 > p_3$  was specified.

$$[7] \quad \text{norm}[gg_1^+ > gg_2^+ > gg_3^+] = \text{norm}[gg_1^+ > gg_2^+, gg_2^+ > gg_3^+]$$

Groups and priority. A group is an abbreviation for a pointwise extension of the declared relation to the members of the group.

$$\begin{aligned} [8] \quad & \text{norm}[\{p\}] = \\ [9] \quad & \text{norm}[\{ \} > g] = \\ [10] \quad & \text{norm}[g > \{ \}] = \\ [11] \quad & \text{norm}[\{p_1^+ p_2^+\} > g] = \text{norm}[\{p_1^+\} > g, \{p_2^+\} > g] \\ [12] \quad & \text{norm}[g > \{p_1^+ p_2^+\}] = \text{norm}[g > \{p_1^+\}, g > \{p_2^+\}] \\ [13] \quad & \text{norm}[\{p\} > g] = \text{norm}[p > g] \\ [14] \quad & \text{norm}[g > \{p\}] = \text{norm}[g > p] \\ [15] \quad & \text{norm}[\{as : p^*\} > g] = \text{norm}[\{as : p^*\}, \{p^*\} > g] \\ [16] \quad & \text{norm}[g > \{as : p^*\}] = \text{norm}[\{as : p^*\}, g > \{p^*\}] \end{aligned}$$

Groups and associativity.

$$\begin{aligned} [17] \quad & \text{norm}[\{ \} \text{ as } g] = \text{norm}[g] \\ [18] \quad & \text{norm}[\{p p^*\} \text{ as } g] = \text{norm}[p \text{ as } g, \{p^*\} \text{ as } g] \\ [19] \quad & \text{norm}[g \text{ as } \{ \}] = \text{norm}[g] \\ [20] \quad & \text{norm}[g \text{ as } \{p p^*\}] = \text{norm}[g \text{ as } p, g \text{ as } \{p^*\}] \\ [21] \quad & \text{norm}[\{as_1 : p^*\} \text{ as}_2 g] = \text{norm}[\{as_1 : p^*\}, \{p^*\} \text{ as}_2 g] \\ [22] \quad & \text{norm}[g \text{ as}_2 \{as_1 : p^*\}] = \text{norm}[\{as_1 : p^*\}, g \text{ as}_2 \{p^*\}] \end{aligned}$$

Associativity groups are abbreviations. The members of an associativity group are mutually associative with respect to the declared relation. If the group contains a single production, it is taken to define the associativity for that production. Otherwise, the associativities are defined only *between* the members of the group and are not defined reflexively. This makes it possible, for instance, to have a production that is left-associative with respect to itself, but right-associative with respect to a group of other productions.

$$\begin{aligned} [23] \quad & \text{norm}[\{as : p\}] = p \text{ as } p \\ [24] \quad & \text{norm}[\{as : p_1 p_2\}] = p_1 \text{ as } p_2 \\ [25] \quad & \text{norm}[\{as : p_1 p_2 p^+\}] = \text{norm}[p_1 \text{ as } p_2, \{as : p_1 p^+\}, \{as : p_2 p^+\}] \end{aligned}$$

The function  $\text{trans}[\_]$  takes the transitive closure of the  $>$  relation.

$$[26] \quad \frac{pr_1^*, p_1 > p_2, pr_2^* = pr^*, pr^* = pr_3^*, p_2 > p_3, pr_4^*, p_1 > p_3 \in pr^* \neq \top}{\text{trans}[pr^*] = \text{trans}[p_1 > p_3, pr^*]}$$

$$\begin{aligned}
[27] \quad & \text{trans}[pr_1^*, pr, pr_2^*, pr, pr_3^*] = \text{trans}[pr_1^*, pr, pr_2^*, pr_3^*] \\
[28] \quad & \text{trans}[pr^*] = pr^* \\
& \quad \text{otherwise}
\end{aligned}$$

The function `assoc[.]` derives associativity declarations from the productions of a grammar. Productions that have an attribute declaring them as left, right, or non-associative produce a declaration of that associativity in the priorities declaration.

$$\begin{aligned}
[29] \quad & \text{assoc}[\mathcal{G}] = \text{assoc}[\mathcal{P}(\mathcal{G})] \\
[30] \quad & \text{assoc}[] = \\
[31] \quad & \text{assoc}[p_1^+ p_2^+] = \text{assoc}[p_1^+] ++ \text{assoc}[p_2^+] \\
[32] \quad & \text{assoc}[p] = p \text{ as } p \quad \textbf{when} \quad p = \alpha \rightarrow \mathcal{A} \{attr_1^*, as, attr_2^*\} \\
[33] \quad & \text{assoc}[p] = \text{otherwise}
\end{aligned}$$

The function ‘syntax’ derives from a priorities declaration the list of all productions referred to in that declaration.

$$\begin{aligned}
[34] \quad & \text{syntax}[] = \emptyset \\
[35] \quad & \text{syntax}[pr_1^+, pr_2^+] = \text{syntax}[pr_1^+] \text{ syntax}[pr_2^+] \\
[36] \quad & \text{syntax}[p_1 > p_2] = \text{syntax } p_1 \ p_2 \\
[37] \quad & \text{syntax}[p_1 \text{ as } p_2] = \text{syntax } p_1 \ p_2
\end{aligned}$$

Merging and ordering of grammars.

$$\begin{aligned}
[38] \quad & \text{priorities} = \emptyset \\
[39] \quad & \text{priorities } pr_1^* \diamond \text{priorities } pr_2^* = \langle \text{priorities } pr_1^*, pr_2^*, \emptyset \rangle \\
[40] \quad & \text{priorities } pr^* \diamond \text{syntax } p^* = \langle \text{syntax } p^*, \text{priorities } pr^* \rangle
\end{aligned}$$

### 3.1.4 Parse Trees with Priority Conflicts

We extend the notion of well-formedness of parse trees to well-formedness over a grammar with priorities. A tree is well-formed if it is a well-formed context-free tree and if, moreover, it does not contain priority conflicts.

**module** Priority-Sdf-Trees

**imports** Kernel-Sdf-Trees<sup>2.5.5</sup> Priority-Sdf-Projection<sup>3.1.2</sup>

**exports**

**context-free syntax**

$$\begin{aligned}
& \text{Grammar } \text{“} \vdash_{\text{prio}} \text{” ATerm} && \rightarrow \text{Bool} \\
& \text{conf } \text{“} [ \text{” Priorities } \text{“} ] \text{ “} ( \text{” ATerm } \text{“} ) && \rightarrow \text{Bool} \\
& \text{rootconf } \text{“} [ \text{” Priorities } \text{“} ] \text{ “} ( \text{” Production } \text{“}, \text{” ATermList } \text{“} ) && \rightarrow \text{Bool} \\
& \text{left } \text{“} [ \text{” Priorities } \text{“} ] \text{ “} ( \text{” Production } \text{“}, \text{” ATerm } \text{“} ) && \rightarrow \text{Bool} \\
& \text{middle } \text{“} [ \text{” Priorities } \text{“} ] \text{ “} ( \text{” Production } \text{“}, \text{” ATermList } \text{“} ) && \rightarrow \text{Bool} \\
& \text{right } \text{“} [ \text{” Priorities } \text{“} ] \text{ “} ( \text{” Production } \text{“}, \text{” ATerm } \text{“} ) && \rightarrow \text{Bool}
\end{aligned}$$

**equations**

We introduce an extension of the notion of well-formedness. A tree is well-formed with respect to the priorities in a grammar, if it is well-formed with



respect to the grammar and does not contain a priority conflict.

$$[1] \quad \frac{\mathcal{G} \vdash T = \top, \text{conf}[\![\text{Pr}(\mathcal{G})]\!](T) = \perp}{\mathcal{G} \vdash_{\text{prio}} T = \top}$$

An application has a conflict if it has a root conflict or if any of its descendants has a conflict.

$$[2] \quad \frac{\text{rootconf}[\![pr^*]\!](\text{production}(Prod), Args) \vee \text{conf}[\![pr^*]\!](Args) = Bool}{\text{conf}[\![pr^*]\!](\text{appl}(Prod, Args)) = Bool}$$

For the other constructors, a tree has a conflict if any descendant has.

$$\begin{aligned} [3] \quad & \text{conf}[\![pr^*]\!](\text{amb}([T])) = \text{conf}[\![pr^*]\!](T) \\ [4] \quad & \text{conf}[\![pr^*]\!](\text{amb}([T, Ts])) = \text{conf}[\![pr^*]\!](T) \vee \text{conf}[\![pr^*]\!](\text{amb}([Ts])) \\ [5] \quad & \text{conf}[\![pr^*]\!]([]) = \perp \\ [6] \quad & \text{conf}[\![pr^*]\!](\langle T \rangle) = \text{conf}[\![pr^*]\!](T) \\ [7] \quad & \text{conf}[\![pr^*]\!](\langle T, Ts \rangle) = \text{conf}[\![pr^*]\!](T) \vee \text{conf}[\![pr^*]\!](\langle Ts \rangle) \end{aligned}$$

An application with no descendants does not have a conflict.

$$[8] \quad \text{rootconf}[\![pr^*]\!](p, []) = \perp$$

An application with more than one descendant has a root conflict if it has a left conflict, a middle conflict or a right conflict.

$$[9] \quad \text{rootconf}[\![pr^*]\!](p, [T, Ts]) = \text{left}[\![pr^*]\!](p, T) \vee \text{middle}[\![pr^*]\!](p, [Ts])$$

An injection, i.e., an application with only one child has a root conflict if its production has higher priority than its child's production.

$$[10] \quad \frac{\text{production}(Prod) = p_2, p_1 > p_2 \in pr^* = Bool}{\text{rootconf}[\![pr^*]\!](p_1, [\text{appl}(Prod, Args)]) = Bool}$$

A tree has a left conflict if the productions of root and left-most child are mutually right-associative or non-associative, or if the root production has higher priority than the child production.

$$[11] \quad \frac{\text{production}(Prod) = p_2, p_1 \text{ right } p_2 \in pr^* \vee p_1 \text{ non-assoc } p_2 \in pr^* \vee p_1 > p_2 \in pr^* = Bool}{\text{left}[\![pr^*]\!](p_1, \text{appl}(Prod, Args)) = Bool}$$

If the left-most child is an ambiguity node, the tree has a conflict if there is a left conflict with any of the possibilities of the ambiguity.

$$\begin{aligned} [12] \quad & \text{left}[\![pr^*]\!](p_1, \text{amb}([T])) = \text{left}[\![pr^*]\!](p_1, T) \\ [13] \quad & \text{left}[\![pr^*]\!](p_1, \text{amb}([T, Ts])) = \text{left}[\![pr^*]\!](p_1, T) \vee \text{left}[\![pr^*]\!](p_1, \text{amb}([Ts])) \end{aligned}$$

A tree has a middle conflict if the root production has higher priority than any of the middle child productions.

$$[14] \quad \frac{\text{production}(Prod) = p_2, p_1 > p_2 \in pr^* \vee \text{middle}[\![pr^*]\!](p_1, [Ts]) = Bool}{\text{middle}[\![pr^*]\!](p_1, [\text{appl}(Prod, Args), Ts]) = Bool}$$

$$[15] \quad \text{middle}[\![pr^*]\!](p, [T]) = \text{right}[\![pr^*]\!](p, T)$$

$$[16] \quad \text{middle}[\![pr^*]\!](p_1, [\text{amb}([T]), Ts]) = \text{middle}[\![pr^*]\!](p_1, [T, Ts])$$

$$[17] \quad \frac{\text{middle}[\![pr^*]\!](p_1, [T, Ts]) \vee \text{middle}[\![pr^*]\!](p_1, [\text{amb}([Ts]), Ts]) = \text{Bool}}{\text{middle}[\![pr^*]\!](p_1, [\text{amb}([T, Ts]), Ts]) = \text{Bool}}$$

A tree has a right conflict if the productions of root and right-most child are mutually left-associative, non-associative or associative (a synonym for left), or if the root production has higher priority than the child production.

$$[18] \quad \frac{\begin{array}{l} \text{production}(Prod) = p_2, \\ p_1 \text{ left } p_2 \in pr^* \vee p_1 \text{ assoc } p_2 \in pr^* \\ \vee p_1 \text{ non-assoc } p_2 \in pr^* \vee p_1 > p_2 \in pr^* = \text{Bool} \end{array}}{\text{right}[\![pr^*]\!](p_1, \text{appl}(Prod, Args)) = \text{Bool}}$$

The case of an ambiguity as right-most child.

$$[19] \quad \text{right}[\![pr^*]\!](p_1, \text{amb}([T])) = \text{right}[\![pr^*]\!](p_1, T)$$

$$[20] \quad \frac{\text{right}[\![pr^*]\!](p_1, T) \vee \text{right}[\![pr^*]\!](p_1, \text{amb}([Ts])) = \text{Bool}}{\text{right}[\![pr^*]\!](p_1, \text{amb}([T, Ts])) = \text{Bool}}$$

### 3.1.5 Discussion

Here we have described the requirements on parse trees that a parser should produce, i.e., not containing priority conflicts. There are various ways to implement this requirement. One possible scheme that is further discussed in Klint and Visser (1994) is to interpret the priority rules as a filter on parse forests that prunes the subtrees with conflicts. This scheme is used in the parser in the current ASF+SDF Meta-Environment (Heering *et al.*, 1989, Klint, 1993). An advantage of this approach is that disambiguation is decoupled from parsing and that other disambiguation filters could be added. The drawback of the approach is that the parse forest can become very large, which hampers efficiency. Therefore, applying the priority rules as early as possible in the parsing process will increase efficiency. A parser-generation time interpretation of priorities is described in Visser (1995). There the priorities are completely expressed in the parse tables produced by the parser generator. An implementation of this method is discussed in Visser (1997e).

*Other Disambiguation Methods* Disambiguation by priority conflicts is similar to the methods using precedences of Earley (1975) and Aho *et al.* (1975). The latter also describe a method for interpreting these rules in the parser generation process, which is less general than the one in Visser (1995). Disambiguation by priorities as defined in this section is based on the definition of priorities in Heering *et al.* (1989). In that definition a second interpretation of priorities is defined. Parse trees are interpreted as a multi-set of productions and the priorities are interpreted as an ordering of such multi-sets. This ordering is

used to make a further selection of trees if the filtering by priority conflicts does not solve all ambiguities.

Subtree exclusion is a disambiguation method introduced by Thorup (1994) that works by specifying a finite set of partial parse trees that are forbidden as subtrees of parse trees yielded by the parser. This method allows a more fine tuned disambiguation than is achievable by the priority scheme. Examples are disambiguation of generic operators and internal arguments. Some problems can not be solved appropriately. The if-then-else ambiguity is solved in the same way as with priorities, which is not correct. In Klint and Visser (1994) these and several other disambiguation methods are studied in the framework of filters on parse forests.

*Brackets* Unparsing is complicated in the presence of priorities. When a parse tree is created by a semantics processor, a rewriter for instance, it might create a well-formed tree that does not satisfy the  $\vdash_{\text{prio}}$  predicate, i.e., contains a priority conflict. Such trees are semantically meaningful, but problematic when their yield is considered. Naively translating an abstract syntax tree to a string as described before might lead to a string that, when parsed, does not represent the same tree because it would contain conflicts. To force equivalence of tree and string, brackets should be introduced. In Van den Brand and Visser (1996) the rules for priority conflicts are used to place brackets when unparsing an abstract syntax tree.

### 3.2 Regular Expressions

Certain patterns of context-free productions occur again and again. Examples of such patterns are lists, lists with separators, optional constructs and alternative. For example, a list of one or more identifiers can be specified by the grammar

```
syntax
  Id          -> Id-List
  Id-List Id-List -> Id-List {left}
```

Here a list is defined in terms of two constructors, one for singleton lists and one for concatenation of lists.

Many formalisms provide shortcuts for such patterns by extending the language of context-free grammars with some collection of regular operators on symbols. For instance, BNF provides an alternative at the level of productions, i.e., a production has the form  $A := A_0 | \dots | A_n$ , where the symbol  $|$  has the meaning of *or*. Extended BNF (EBNF) is the canonical extension of BNF with regular operators. In one formulation, Wirth (1977) adds the operators  $\{A\}$  for iteration and  $[A]$  for optionality. Variations on this notation appear in Lee (1972) and Williams (1982). SDF provides iteration  $A^*$  and  $A^+$  and  $\{A L\} \oplus$  iteration for abbreviation of lists of  $A$ s separated by a literal  $L$ .

In this section we give an extension of context-free productions by a set of regular operators on symbols. In all the approaches mentioned above regular operators are given a special treatment. New in our formulation is the treatment of regular operators as first class citizens. They are nothing but constructors of new symbols that spare the specifier the burden of having to invent new names. As a consequence, a regular expression can occur at all positions where a normal symbol can occur, in particular in the right-hand side of a production.

This approach is motivated by the following considerations: (1) It enables us to express the meaning of regular expressions by means of a normalization of the grammar that adds defining productions for each expression. (2) Our grammars function as signatures for algebraic specifications, where each production represents a function. If regular symbols can not be the result of functions, as is the case in SDF, we still have to define an auxiliary symbol to define a function that yields such a result. For example, suppose that we want to define a function `add` that adds an integer to each integer in a list of integers. In the syntax below we can write this as

`add(Int, Int*) -> Int*`

whereas in SDF we should introduce an auxiliary sort `IntList` to represent the result sort of this function.

### 3.2.1 Syntax

We consider the following operators:

- Empty: The symbol  $()$  represents the empty string
- Concatenation: The symbol  $(\mathcal{A}_1 \dots \mathcal{A}_n)$  with  $n > 2$  denotes the concatenation  $a_1 \dots a_n$  of expressions  $a_i$  of type  $\mathcal{A}_i$ .
- Alternative: The symbol  $\mathcal{A}_1 | \dots | \mathcal{A}_n$ , with  $n \geq 2$ , denotes an expression  $a$  of one of the types  $\mathcal{A}_i$ .
- Optional: The symbol  $\mathcal{A}?$  is an optional  $\mathcal{A}$
- Iteration: The symbol  $\mathcal{A}^* (\mathcal{A}^+)$  denotes a list  $a_1 \dots a_n$  with  $n \geq 0$  ( $n \geq 1$ ) of expressions  $a_i$  of type  $\mathcal{A}$ .
- Iteration with separator: The symbol  $\{\mathcal{A} \mathcal{B}\}^* (\{\mathcal{A} \mathcal{B}\}^+)$  denotes a list  $a_1 b_1 \dots b_{n-1} a_n$  with  $n \geq 0$  ( $n \geq 1$ ) of expressions  $a_i$  of type  $\mathcal{A}$  separated by expressions  $b_i$  of type  $\mathcal{B}$ . Observe that, unlike in SDF, in SDF2 there is no limitation of the symbols that can be used as separators. For example, `{Stat [\;] | [\n]}*` denotes lists of statements separated by semicolons *or* newlines.
- Constrained iteration: The symbol  $\{\mathcal{A}\}n^+$  with  $n \geq 2$  denotes a list  $a_1 \dots a_n$  of  $n$  or more expressions  $a_i$  of type  $\mathcal{A}$ . Similarly for  $\{\mathcal{A} \mathcal{B}\}n^+$  with separator  $\mathcal{B}$ .
- Set expressions: The symbol `Set[ $\mathcal{A}$ ]` represents the syntax of set expressions of the form  $\{a_1, \dots, a_n\}$  with the  $a_i$  expressions of type  $\mathcal{A}$
- Product: The symbol  $\mathcal{A}_1 \# \dots \# \mathcal{A}_n$ , with  $n \geq 2$ , denotes tuples  $\langle a_1, \dots, a_n \rangle$  of expressions  $a_i$  of type  $\mathcal{A}_i$ .
- Functions: The symbol  $(\mathcal{A}_1 \dots \mathcal{A}_n \Rightarrow \mathcal{B})$ , with  $n \geq 0$ , denotes function expressions  $f$  that can be used in expressions  $f(a_1 \dots a_n)$  of type  $\mathcal{B}$  with the  $a_i$  expressions of type  $\mathcal{A}_i$ .
- Permutation: The symbol  $\ll \mathcal{A}_1 \dots \mathcal{A}_n \gg$  denotes expressions of the form  $a_1 \dots a_n$  such that for each  $\mathcal{A}_i$  exactly one of the  $a_j$  has type  $\mathcal{A}_i$ .

The syntax of these operators is defined in the following module. Observe that the empty symbol  $()$  and sequences  $(\mathcal{A}_1 \dots \mathcal{A}_n)$  are not defined using a single production  $"(" \text{Symbol} * ")" \rightarrow \text{Symbol}$  because parentheses around a single symbol are already used as brackets; see §2.1.

```

module Regular-Sdf-Syntax
imports Kernel-Sdf-Syntax2.3.1 IntCon
exports
  context-free syntax
    "(" ")" → Symbol
    "(" Symbol Symbol+ ")" → Symbol
    Symbol "?" → Symbol
    Symbol "+" → Symbol
    Symbol "*" → Symbol
    "{" Symbol Symbol "}" "+" → Symbol
    "{" Symbol Symbol "}" "*" → Symbol
    "{" Symbol "}" NatCon "+" → Symbol
    "{" Symbol Symbol "}" NatCon "+" → Symbol
    "Set" "[" Symbol "]" → Symbol
    Symbol "#" Symbol → Symbol {right}
    "(" Symbols "⇒" Symbol ")" → Symbol
    Symbol "|" Symbol → Symbol {right}
    "⟨⟨" Symbols "⟩⟩" → Symbol
  priorities
    {Symbol "?" → Symbol, Symbol "*" → Symbol, Symbol "+" → Symbol,
     Symbol NatCon "+" → Symbol} > Symbol "#" Symbol → Symbol >
    Symbol "|" Symbol → Symbol

```

### 3.2.2 Normalization

We define a normalization function  $r[\_]$  that for each regular expression that is used in the grammar introduces one or more productions that define its meaning. In this interpretation regular expressions form a shorthand for defining extra symbols and productions.

**Example 2.1** The following production defines a single production describing the structure of a block in a While program consisting of an optional declaration followed by a list of statements.

```

syntax
  "begin" (Decl ";")? {Stat ";"}+ "end" -> Stat

```

The normalization of this grammar is:

```

syntax
  "begin" (Decl ";")? {Stat ";"}+ "end" -> Stat
  (Decl ";") → (Decl ";")?
  Decl ";" → (Decl ";")
  Stat → {Stat ";"}+
  {Stat ";"}+ ";" {Stat ";"}+ → {Stat ";"}+ {left}

```

<code>{Stat ";" + ";" {Stat ";"}*</code>	<code>-&gt; {Stat ";" + ";"}</code>
<code>{Stat ";"}* ";" {Stat ";"} +</code>	<code>-&gt; {Stat ";" + ";"}</code>
<code>{Stat ";"}* ";" {Stat ";"}*</code>	<code>-&gt; {Stat ";"}* {left}</code>
	<code>-&gt; {Stat ";"}*</code>
<code>{Stat ";"} +</code>	<code>-&gt; {Stat ";"}*</code>
<code>priorities</code>	
<code>{left :</code>	
<code>  {Stat ";" + ";" {Stat ";"} +</code>	<code>-&gt; {Stat ";" + ";" {left}</code>
<code>  {Stat ";" + ";" {Stat ";"}*</code>	<code>-&gt; {Stat ";" + ";"}</code>
<code>  {Stat ";"}* ";" {Stat ";"} +</code>	<code>-&gt; {Stat ";" + ";"}</code>
<code>  {Stat ";"}* ";" {Stat ";"}*</code>	<code>-&gt; {Stat ";"}* {left}</code>
<code>&gt;}</code>	
<code>{Stat ";"} +</code>	<code>-&gt; {Stat ";"}*</code>

We see that the meaning of the operators is expressed by means of extra productions. Observe how regular expressions are used as target symbols of productions.  $\square$

#### **module** Regular-Sdf-Normalization

**imports** Regular-Sdf-Syntax<sup>3.2.1</sup> Priority-Sdf-Syntax<sup>3.1.1</sup>

Literals-Sdf-Syntax<sup>2.4.3</sup> Kernel-Sdf-Normalization<sup>2.3.3</sup>

#### **exports**

##### **context-free syntax**

<code>"r["</code> Grammar <code>""]</code>	<code>→ Grammar</code>
<code>"r["</code> Symbols <code>""]</code>	<code>→ Grammar</code>
<code>alt(Symbol, Symbol)</code>	<code>→ Grammar</code>
<code>tup(Symbol)</code>	<code>→ Symbols</code>
<code>perm(Symbols)</code>	<code>→ Productions</code>
<code>perm3(Symbols, Symbols)</code>	<code>→ Productions</code>

##### **equations**

The function `r[.]` adds defining productions for each regular expression occurring in one of the productions of the grammar. Existing productions are not affected.

$$[1] \quad r[\mathcal{G}] = \mathcal{G} \text{ r}[\alpha] \quad \textbf{when} \quad \{\alpha\} = \text{symbols}(\mathcal{G})$$

Recall that the function ‘symbols’, defined in §2.3.2, gives the set of all symbols in a grammar. The function `r[.]` generates a grammar for each of the regular expressions in the list of symbols.

$$[2] \quad r[\ ] = \emptyset$$

$$[3] \quad r[\alpha^+ \beta^+] = r[\alpha^+] \text{ r}[\beta^+]$$

*Concatenation* The regular expression  $(\alpha)$  is a *symbol* that abbreviates the *concatenation* of the symbols  $\alpha$ .

$$[4] \quad r[(\ )] = \text{syntax} \rightarrow ( \ )$$

$$[5] \quad \frac{p = \mathcal{A} \alpha^+ \rightarrow (\mathcal{A} \alpha^+)}{r[(\mathcal{A} \alpha^+)] = \text{syntax } p \text{ r}[\mathcal{A} \alpha^+]}$$

Note that `r[ $\mathcal{A} \alpha^+$ ]` recursively produces the productions for regular expressions in the list of symbols  $\mathcal{A} \alpha^+$ .

*Alternative* The *alternative*  $\mathcal{A}|\mathcal{B}$  denotes either  $\mathcal{A}$  or  $\mathcal{B}$ . We could thus define  $r[\mathcal{A}|\mathcal{B}]$  to yield the productions  $\mathcal{A} \rightarrow \mathcal{A}|\mathcal{B}$  and  $\mathcal{B} \rightarrow \mathcal{A}|\mathcal{B}$ . However, if one of the alternatives is again an alternative, an unnecessary chain  $\mathcal{A} \rightarrow \mathcal{A}|\mathcal{B}$  and  $\mathcal{A}|\mathcal{B} \rightarrow \mathcal{A}|\mathcal{B}|\mathcal{C}$  is created. We would rather have  $\mathcal{A} \rightarrow \mathcal{A}|\mathcal{B}|\mathcal{C}$ . Therefore, we define

$$[6] \quad r[\mathcal{A} | \mathcal{B}] = \text{alt}(\mathcal{A} | \mathcal{B}, \mathcal{A} | \mathcal{B})$$

where the function ‘alt’ unpacks the alternative until a symbol is reached that is not an alternative.

$$[7] \quad \text{alt}(\mathcal{B}_1 | \mathcal{B}_2, \mathcal{A}) = \text{alt}(\mathcal{B}_1, \mathcal{A}) \text{ alt}(\mathcal{B}_2, \mathcal{A})$$

$$[8] \quad \text{alt}(\mathcal{B}, \mathcal{A}) = \text{syntax } \mathcal{B} \rightarrow \mathcal{A} \ r[\mathcal{B}] \quad \textbf{otherwise}$$

*Optional* The *optional construct*  $\mathcal{A}?$  is either empty or  $\mathcal{A}$ .

$$[9] \quad \frac{\begin{array}{l} p_1 = \rightarrow \mathcal{A}?, \\ p_2 = \mathcal{A} \rightarrow \mathcal{A}? \end{array}}{r[\mathcal{A}?] = \text{syntax } p_1 \ p_2 \ r[\mathcal{A}]}$$

*Iteration* The *iteration* operator  $\mathcal{A}+$  denotes lists of *one or more*  $\mathcal{A}$ ’s, i.e., either  $\mathcal{A}$  or  $\mathcal{A} \mathcal{A}$  or  $\mathcal{A} \mathcal{A} \mathcal{A}$  or  $\dots$ . The iteration  $\mathcal{A}^*$  denotes a list of *zero or more*  $\mathcal{A}$ ’s, i.e.,  $\epsilon$  (empty) or  $\mathcal{A}$  or  $\mathcal{A} \mathcal{A}$  or  $\mathcal{A} \mathcal{A} \mathcal{A}$  or  $\dots$ . There are several ways to define such lists with productions. It is not sufficient to define a list by means of the productions

$$\begin{array}{ll} \rightarrow \mathcal{A}^* & \mathcal{A} \rightarrow \mathcal{A} + \\ \mathcal{A} + \rightarrow \mathcal{A}^* & \mathcal{A} + \mathcal{A} \rightarrow \mathcal{A} + \end{array}$$

The symbols  $\mathcal{A}^*$  and  $\mathcal{A}+$  can be the right-hand side of any production, i.e., lists can be the result of arbitrary functions. Therefore, an  $\mathcal{A}^*$  expression can also contain function calls and variables. For instance, if a grammar contains the production

`yield(Tree) -> Symbol*`

then `yield(T1) yield(T2)` should also be an expression of type `Symbol*` (with `T1` and `T2` expression of type `Tree`). We have the following rules for the composition of list expressions.

0. A single  $\mathcal{A}$  is an  $\mathcal{A}+$ .
1. An  $\mathcal{A}+$  followed by an  $\mathcal{A}+$  is an  $\mathcal{A}+$ .
2. An  $\mathcal{A}+$  followed by an  $\mathcal{A}^*$  is an  $\mathcal{A}+$ .
3. An  $\mathcal{A}^*$  followed by an  $\mathcal{A}+$  is an  $\mathcal{A}+$ .
4. An  $\mathcal{A}^*$  followed by an  $\mathcal{A}^*$  is an  $\mathcal{A}^*$ .
5. An  $\mathcal{A}^*$  can be empty.
6. An  $\mathcal{A}+$  is an  $\mathcal{A}^*$ .

Productions expressing these rules are generated by the following equation. The priorities section declares the concatenation operators to be mutually left-associative. The priority prevents that the empty production and the injection are used vacuously.

$$\begin{aligned}
 p_0 &= \mathcal{A} \rightarrow \mathcal{A}+ , \\
 p_1 &= \mathcal{A}+ \mathcal{A}+ \rightarrow \mathcal{A}+ \{\text{left}\}, \\
 p_2 &= \mathcal{A}+ \mathcal{A}* \rightarrow \mathcal{A}+ , \\
 p_3 &= \mathcal{A}* \mathcal{A}+ \rightarrow \mathcal{A}+ , \\
 p_4 &= \mathcal{A}* \mathcal{A}* \rightarrow \mathcal{A}* \{\text{left}\}, \\
 p_5 &= \rightarrow \mathcal{A}* , \\
 p_6 &= \mathcal{A}+ \rightarrow \mathcal{A}* \\
 [10] \quad &\frac{}{\text{r}[\mathcal{A}*] = \text{syntax } p_0 \ p_1 \ p_2 \ p_3 \ p_4 \ p_5 \ p_6} \\
 &\text{priorities } \{\text{left} : p_1 \ p_2 \ p_3 \ p_4\} > \{p_5 \ p_6\} \text{r}[\mathcal{A}]
 \end{aligned}$$

$$[11] \quad \text{r}[\mathcal{A}+] = \text{r}[\mathcal{A}*]$$

*Iteration with Separator* The *iteration with separator* operators  $\{\mathcal{A} \ \mathcal{B}\}+$  and  $\{\mathcal{A} \ \mathcal{B}\}*$  denote iteration of  $\mathcal{A}$ 's separated by  $\mathcal{B}$ 's. Their meaning is defined analogously to  $\mathcal{A}+$  and  $\mathcal{A}*$ .

$$\begin{aligned}
 p_0 &= \mathcal{A} \rightarrow \{\mathcal{A} \ \mathcal{B}\}+ , \\
 p_1 &= \{\mathcal{A} \ \mathcal{B}\}+ \mathcal{B} \{\mathcal{A} \ \mathcal{B}\}+ \rightarrow \{\mathcal{A} \ \mathcal{B}\}+ \{\text{left}\}, \\
 p_2 &= \{\mathcal{A} \ \mathcal{B}\}+ \mathcal{B} \{\mathcal{A} \ \mathcal{B}\}* \rightarrow \{\mathcal{A} \ \mathcal{B}\}+ , \\
 p_3 &= \{\mathcal{A} \ \mathcal{B}\}* \mathcal{B} \{\mathcal{A} \ \mathcal{B}\}+ \rightarrow \{\mathcal{A} \ \mathcal{B}\}+ , \\
 p_4 &= \{\mathcal{A} \ \mathcal{B}\}* \mathcal{B} \{\mathcal{A} \ \mathcal{B}\}* \rightarrow \{\mathcal{A} \ \mathcal{B}\}* \{\text{left}\}, \\
 p_5 &= \rightarrow \{\mathcal{A} \ \mathcal{B}\}* , \\
 p_6 &= \{\mathcal{A} \ \mathcal{B}\}+ \rightarrow \{\mathcal{A} \ \mathcal{B}\}* \\
 [12] \quad &\frac{}{\text{r}[\{\mathcal{A} \ \mathcal{B}\}*] = \text{syntax } p_0 \ p_1 \ p_2 \ p_3 \ p_4 \ p_5 \ p_6} \\
 &\text{priorities } \{\text{left} : p_1 \ p_2 \ p_3 \ p_4\} > p_6 \text{r}[\mathcal{A} \ \mathcal{B}]
 \end{aligned}$$

$$[13] \quad \text{r}[\{\mathcal{A} \ \mathcal{B}\}+] = \text{r}[\{\mathcal{A} \ \mathcal{B}\}*]$$

*Constrained Iteration* The iteration operator  $\{\mathcal{A}\}n+$  denotes the iteration of *at least*  $n$   $\mathcal{A}$ s. First of all we define that zero or more  $\mathcal{A}$ s corresponds to  $*$  iteration and that one or more  $\mathcal{A}$ s corresponds to  $+$  iteration. For integers  $n \geq 2$  we define  $\{\mathcal{A}\}n+$  in terms of  $\{\mathcal{A}\}(n-1)+$ , and eventually  $\mathcal{A}+$ , by productions of the form  $\mathcal{A} \ \{\mathcal{A}\}(n-1)+ \rightarrow \{\mathcal{A}\}n+$ .

$$[14] \quad \{\mathcal{A}\} \ 0 \ + = \mathcal{A}*$$

$$[15] \quad \{\mathcal{A}\} \ 1 \ + = \mathcal{A}+$$

$$[16] \quad \frac{n \geq 2 = \top, \ n-1 = n', \ p = \mathcal{A} \ \{\mathcal{A}\} \ n' \ + \rightarrow \{\mathcal{A}\} \ n \ +}{\text{r}[\{\mathcal{A}\} \ n \ +] = \text{syntax } p \ \text{r}[\{\mathcal{A}\} \ n' \ +]}$$

Constrained iteration is defined similarly for lists with separators.

$$[17] \quad \{\mathcal{A} \ \mathcal{B}\} \ 0 \ + = \{\mathcal{A} \ \mathcal{B}\}*$$

$$[18] \quad \{\mathcal{A} \ \mathcal{B}\} \ 1 \ + = \{\mathcal{A} \ \mathcal{B}\}+$$

$$[19] \quad \frac{n \geq 2 = \top, \ n-1 = n', \ p = \mathcal{A} \ \mathcal{B} \ \{\mathcal{A} \ \mathcal{B}\} \ n' \ + \rightarrow \{\mathcal{A} \ \mathcal{B}\} \ n \ +}{\text{r}[\{\mathcal{A} \ \mathcal{B}\} \ n \ +] = \text{syntax } p \ \text{r}[\{\mathcal{A} \ \mathcal{B}\} \ n' \ +]}$$



*Tuples* For the definition of functions that return a tuple of values, new sorts have to be invented. To give sensible types to tuples the notation  $\mathcal{A} \# \mathcal{B}$  is introduced. A symbol  $\mathcal{A}_1 \# \dots \# \mathcal{A}_n$  denotes a tuple of  $\mathcal{A}_1 \dots \mathcal{A}_n$  expressions. A tuple is written as  $\langle T_1, \dots, T_n \rangle$ , where the  $T_i$  are expressions of type  $\mathcal{A}_i$ .

$$[20] \quad \frac{"<" \mathrel{++} \text{tup}(\mathcal{A} \# \mathcal{B}) \mathrel{++} ">" = \alpha, \quad p = \alpha \rightarrow \mathcal{A} \# \mathcal{B}}{\mathbf{r}[\mathcal{A} \# \mathcal{B}] = \text{syntax } p \mathbf{r}[\alpha]}$$

The auxiliary function  $\text{tup}[\_]$  derives the syntax of the body of the tuple by separating the symbols by commas.

$$[21] \quad \text{tup}(\mathcal{A} \# \mathcal{B}) = \text{tup}(\mathcal{A}) \mathrel{++} ", " \mathrel{++} \text{tup}(\mathcal{B})$$

$$[22] \quad \text{tup}(\mathcal{A}) = \mathcal{A} \quad \textbf{otherwise}$$

*Sets* The conventional notation for sets is a list of items between  $\{$  and  $\}$ . The operator  $\text{Set}[\mathcal{A}]$  generates this notation such that if  $T_1 \dots T_n$  are expressions of type  $\mathcal{A}$ , then  $\{T_1, \dots, T_n\}$  is an expression of type  $\text{Set}[\mathcal{A}]$ .

$$[23] \quad \frac{\alpha = "\{" \{ \mathcal{A} ", " \}^* "\}" , \quad p = \alpha \rightarrow \text{Set}[\mathcal{A}]}{\mathbf{r}[\text{Set}[\mathcal{A}]] = \mathbf{r}[\alpha] \text{ syntax } p}$$

*Functions* Sometimes it is convenient to pass functions around as data. The operator  $(\alpha \Rightarrow \mathcal{B})$  can be used to give a type to functions. It denotes the sort of functions from  $\alpha$  to  $\mathcal{B}$ . The operator generates syntax for the prefix application of a function to an argument.

$$[24] \quad \frac{p = (\alpha \Rightarrow \mathcal{B}) \text{ "(" } \alpha \text{ ")" } \rightarrow \mathcal{B}}{\mathbf{r}[(\alpha \Rightarrow \mathcal{B})] = \text{syntax } p \mathbf{r}[\alpha \mathcal{B}]}$$

*Permutation* The permutation symbol  $\ll \alpha \gg$  denotes any concatenation of the symbols in  $\alpha$ , i.e.,  $\beta \rightarrow \ll \alpha \gg$  if  $\beta$  is a permutation of  $\alpha$ .

$$[25] \quad \mathbf{r}[\ll \alpha \gg] = \text{syntax perm}(\alpha) \mathbf{r}[\alpha]$$

The function ‘perm’ generates the productions for all permutations of a set of symbols. In case the permutation consists of two elements it generates the two productions directly. In case of more elements the function ‘perm3’ is used to generate permutations.

$$[26] \quad \text{perm}() = \rightarrow \ll \gg$$

$$[27] \quad \text{perm}(\mathcal{A}) = \mathcal{A} \rightarrow \ll \mathcal{A} \gg$$

$$[28] \quad \text{perm}(\mathcal{A} \mathcal{B}) = \mathcal{A} \mathcal{B} \rightarrow \ll \mathcal{A} \mathcal{B} \gg \quad \mathcal{B} \mathcal{A} \rightarrow \ll \mathcal{A} \mathcal{B} \gg$$

$$[29] \quad \text{perm}(\alpha) = \text{perm3}(\alpha) \quad \textbf{when } |\alpha| > 2 = \top$$

For each symbol in the list a production is generated with that symbol first and a permutation of the other symbols following it.

$$[30] \quad \text{perm3}(\alpha, ) =$$

$$[31] \quad \frac{p = \mathcal{A} \ll \alpha \mathcal{B} \gg \rightarrow \ll \alpha \mathcal{A} \mathcal{B} \gg}{\text{perm3}(\alpha, \mathcal{A} \mathcal{B}) = p \mathrel{++} \text{perm3}(\alpha \mathcal{A}, \mathcal{B}) \mathrel{++} \text{perm}(\alpha \mathcal{B})}$$

It should be observed that this is not a very efficient way to implement permutation constructs. It should be adequate for permutations of 2 or 3 elements, though. What is needed in addition to the generation of these productions, is the normalization of the parse trees over these productions to a form that lists the elements in a fixed order such that semantic functions do not also have to deal with all permutations. Cameron (1993) describes an extension of LL(1) parsing for permutation operators. An alternative approach suggested by Cameron (1993) is the introduction of an intermediate symbol representing the union of the symbols in the permutation and a check after parsing that each symbol in the permutation is represented exactly once.

*Discussion* We could have handled several of these regular expressions by translating them to other regular expressions. For instance, optionality can be expressed by means of empty and alternative via the equation  $\mathcal{A}^? = ()|\mathcal{A}$ . In the specification above we have chosen not follow this route. Except for a few cases involving constrained iteration.

### 3.2.3 Equality of Parse Trees with Lists

Since all new constructs are expressed by means of existing constructs—all regular expressions are expressed by means of generated context-free productions—there is no need to extend the definition of well-formedness of parse trees.

We do extend the definition of the equality of trees. This definition makes lists equal modulo associativity of the concatenation operators. It is the basis for matching modulo associativity. We give the equations that should be considered in matching, where a variable  $a^+$  ( $a^*$ ) ranges over all constructs of type  $\mathcal{A}^+$  ( $\mathcal{A}^*$ ) and  $\epsilon_A$  denotes the tree constructed with  $\rightarrow \mathcal{A}^*$ . Empty sublists are units for concatenation and can be removed.

$$\begin{array}{ll} \epsilon_A a^+ = a^+ & a^+ \epsilon_A = a^+ \\ \epsilon_A a^* = a^* & a^* \epsilon_A = a^* \end{array}$$

Injectons from  $\mathcal{A}^+$  into  $\mathcal{A}^*$  can be removed or lifted over concatenations.

$$\begin{aligned} [a_1^+ \rightarrow \mathcal{A}^*] [a_2^+ \rightarrow \mathcal{A}^*] &= [(a_1^+ a_2^+) \rightarrow \mathcal{A}^*] \\ a_1^+ [a_2^+ \rightarrow \mathcal{A}^*] &= a_1^+ a_2^+ \\ [a_1^+ \rightarrow \mathcal{A}^*] a_2^+ &= a_1^+ a_2^+ \end{aligned}$$

Right-associative concatenations are equal to left-associative ones. Each of these expressions involves different concatenation operators.

$$\begin{array}{ll} a_1^+ (a_2^+ a_3^+) = (a_1^+ a_2^+) a_3^+ & a_1^* (a_2^+ a_3^+) = (a_1^* a_2^+) a_3^+ \\ a_1^* (a_2^* a_3^+) = (a_1^* a_2^*) a_3^+ & a_1^* (a_2^+ a_3^*) = (a_1^* a_2^+) a_3^* \\ a_1^+ (a_2^* a_3^+) = (a_1^+ a_2^*) a_3^+ & a_1^+ (a_2^* a_3^*) = (a_1^+ a_2^*) a_3^* \\ a_1^+ (a_2^+ a_3^*) = (a_1^+ a_2^+) a_3^* & a_1^* (a_2^* a_3^*) = (a_1^* a_2^*) a_3^* \end{array}$$

## 3.3 Lexical and Context-Free Syntax

The syntax of a programming language is usually divided into two levels: lexical syntax and context-free syntax. Lexical syntax is the syntax of the tokens,

the words of the language, e.g., identifiers, numbers and keywords. Context-free syntax is the syntax of the sentences of a language, e.g., expressions, statements, type declarations and function definitions. The division affects both language definition and implementation. Conventionally lexical analysis is restricted to grammars that can be recognized by finite automata, whereas context-free analysis is implemented with push-down automata. Indeed, it is sometimes not clear whether the division is motivated by the implementation or by an inherent concept of lexical syntax.

In many formalisms the separation is even physical; lexical and context-free syntax are defined with completely different formalisms that are written in separate files. For instance, YACC and METAL use LEX to define lexical syntax. This means lexical definitions in the form of a number of regular expressions are defined in a separate file. Context-free and lexical definitions share a declaration of token symbols that constitutes the interface between the lexical and context-free level. The syntax definition formalism of PCCTS uses a lexical syntax similar to LEX, but provides a mechanism to include token definitions in the same file as the context-free syntax definition. In SDF lexical and context-free syntax are integrated in one formalism, but still uses different semantics for both. All these approaches have in common that the distinction between lexical and context-free syntax is identified with the distinction between regular and context-free grammars.

In SDF2 the inherent distinction between the two categories is that context-free symbols can be *separated by layout*, while lexical symbols cannot. Beyond that difference there is none. The exact same features should be available for the definition of lexical and context-free syntax.

New in this approach is that we provide a uniform notation for the definition of lexical and context-free syntax by means of context-free productions. Grammars for lexical and context-free syntax are normalized to the context-free grammars of the kernel. The distinction between lexical and context-free syntax is completely expressed in the resulting productions.

By treating lexical and context-free syntax identically, every extension that is defined for one is also applicable to the other. For instance, in §3.1 we defined priorities for disambiguation. In Heering *et al.* (1989) these are only defined for context-free syntax. As result of our approach we can also provide lexical disambiguation through priorities. Similarly the regular operators introduced in §3.2 can be used in the definition of both lexical and context-free syntax.

In addition to lexical syntax we also define variables. Variable schemes are used in the specification of the semantics of a language. We also introduce the notion of lexical variables that range over constructs introduced in lexical syntax grammars.

The extension in this section is called *Basic Sdf* because it covers the basic idea of the original SDF: integration of lexical and context-free syntax in one formalism.

**Example 3.1** The following definition introduces a simple expression language with variables and addition.

```
sorts Id Exp
lexical syntax
  [\ \t\n]    -> LAYOUT
```

```

[a-z]+      -> Id
context-free syntax
Id          -> Exp
Exp "+" Exp -> Exp {left}
variables
[i]         -> Id
[xyz]       -> Exp

```

The lexical syntax section defines the syntax of layout as spaces, tabs and new-lines and identifiers as lists of one or more lowercase letters. The division in lexical and context-free syntax entails that whitespace can occur between expressions, but not between the letters of an identifier.

To illustrate the power of the integration of lexical and context-free syntax we can extend the layout convention above by introducing C-like comments consisting of a string of comment words between `/*` and `*/`, as follows:

```

sorts ComWord Comment
lexical syntax
~[\ \t\n\\\/\*]+ -> ComWord
context-free syntax
"/*" ComWord* "*/" -> Comment
Comment              -> LAYOUT

```

Because the definition of comments is part of the context-free syntax, comment words can be separated by layout, including layout. This means that we have specified nested comments, which is useful when commenting out pieces of code already containing comment.

We can extend the definition of comments further to include syntactically correct expressions between bars as comment words.

```

context-free syntax
"|" Exp "|" -> ComWord

```

For instance, the text

```

a + b /* an expression |x + y| denotes
      the addition of |x| and |y| */ + c

```

is a syntactically correct expression over the grammar above denoting the expression `a+b+c` with some comment after `b`. In the conventional setting of a separate scanner and parser this would require a call to the parser from the scanner. One application of syntactically correct program fragments in comments is in typesetting programs for documentation. The typesetting algorithms applied to the real program text can also be applied in typesetting the expressions in comments and crossreferences to program variables can be extended to variables occurring comments.  $\square$

### 3.3.1 Syntax

The grammar constructors ‘lexical syntax’ and ‘context-free syntax’ introduce the syntax of lexical constructs and context-free constructs, respectively. The grammar constructors ‘variables’ and ‘lexical variables’ introduce the syntax of

variables over context-free symbols and variables over lexical symbols, respectively. The symbol constructors  $\langle\_ -\text{LEX}\rangle$ ,  $\langle\_ -\text{CF}\rangle$  and  $\langle\_ -\text{VAR}\rangle$  are used to indicate lexical symbols, context-free symbols and variable symbols, respectively. The special symbol `LAYOUT` is used to define layout.

```

module Basic-Sdf-Syntax
imports Kernel-Sdf-Syntax2.3.1
exports
  context-free syntax
    "lexical" "syntax" Productions      → Grammar
    "context-free" "syntax" Productions → Grammar
    "variables" Productions             → Grammar
    "lexical" "variables" Productions   → Grammar

    "<" Symbol "-CF" ">"                → Symbol
    "<" Symbol "-LEX" ">"                → Symbol
    "<" Symbol "-VAR" ">"                → Symbol

    "LAYOUT"                            → Symbol

```

### 3.3.2 Normalization

The normalization function defined below expresses the meaning of lexical and context-free syntax by merging them into a single grammar. To avoid interference between the two levels, the symbols in the lexical syntax are renamed into  $\langle\_ -\text{LEX}\rangle$  symbols and the symbols in the context-free syntax are renamed into  $\langle\_ -\text{CF}\rangle$  symbols. These ideas are illustrated in the following example.

**Example 3.2** The grammar in Example 3.1 is mapped to the following grammar in which lexical and context-free syntax have been merged.

```

sorts Id Exp ComWord Comment
syntax
  <[a-z]+-LEX>                -> <Id-LEX>
  <Id-LEX>                     -> <Id-CF>
  <Id-CF>                      -> <Exp-CF>
  <Exp-CF> <LAYOUT?-CF> "+"
    <LAYOUT?-CF> <Exp-CF> -> <Exp-CF> {left}
  [i]                          -> <<Id-CF>-VAR>
  <<Id-CF>-VAR>                 -> <Id-CF>
  [xyz]                        -> <<Exp-CF>-VAR>
  <<Exp-CF>-VAR>                -> <Exp-CF>
  <~[\ \t\n|\\047\*]+-LEX>     -> <ComWord-LEX>
  <ComWord-LEX>                -> <ComWord-CF>
  "|" <LAYOUT?-CF> <Exp-CF>
    <LAYOUT?-CF> "|"          -> <ComWord-CF>
  "/" <LAYOUT?-CF> <ComWord*-CF>
    <LAYOUT?-CF> "*"          -> <Comment-CF>
  [\ \t\n]                     -> <LAYOUT-LEX>
  <LAYOUT-LEX>                  -> <LAYOUT-CF>
  <Comment-CF>                  -> <LAYOUT-CF>
  <LAYOUT-CF> <LAYOUT-CF>      -> <LAYOUT-CF> {left}

```

<LAYOUT-CF>	-> <LAYOUT?-CF>
	-> <LAYOUT?-CF>

The symbols in lexical productions are renamed into  $\langle \_ \text{-LEX} \rangle$  symbols. The symbols in context-free productions are renamed into  $\langle \_ \text{-CF} \rangle$  symbols. The connection between lexical and context-free syntax is made by an injection from each  $\langle \mathcal{A} \text{-LEX} \rangle$  symbol into the corresponding  $\langle \mathcal{A} \text{-CF} \rangle$  symbol.  $\square$

The following module makes these ideas formal by introducing the normalization function  $b[\_]$ .

**module** Basic-Sdf-Normalization

**imports** Basic-Sdf-Syntax<sup>3.3.1</sup> Regular-Sdf-Syntax<sup>3.2.1</sup> Priority-Sdf-Syntax<sup>3.1.1</sup>

Kernel-Sdf-Normalization<sup>2.3.3</sup>

**exports**

**context-free syntax**

“b[” Grammar “]”	→ Grammar
“b <sub>aux</sub> [” Grammar “]”	→ Grammar
“<” Symbols “-LEXs” “>”	→ Symbols
“<” Production “-LEX” “>”	→ Production
“<” Productions “-LEXs” “>”	→ Productions
“<” Grammar “-LEX” “>”	→ Grammar
“<” Symbols “-CFs” “>”	→ Symbols
“<” Production “-CF” “>”	→ Production
“<” Productions “-CFs” “>”	→ Productions
“<” Grammar “-CF” “>”	→ Grammar
“<” Productions “-VARs” “>”	→ Productions
“<” Productions “-LEXVARs” “>”	→ Productions

**equations**

The normalization function  $b[\_]$  integrates lexical and context-free syntax. It applies the auxiliary function  $b_{aux}$  to each subgrammar of a grammar to transform lexical and context-free sections into normal production sections by renaming symbols and separating context-free symbols by  $\langle \text{LAYOUT?-CF} \rangle$ , which entails that two tokens can optionally be separated by  $\langle \text{LAYOUT?-CF} \rangle$ . Context-free layout is a list of lexical layout. Concatenation of layout is defined by the production added by the function  $b[\_]$ .

$$[1] \quad \frac{p = \langle \text{LAYOUT-CF} \rangle \langle \text{LAYOUT-CF} \rangle \rightarrow \langle \text{LAYOUT-CF} \rangle \{left\}}{b[\mathcal{G}] = \text{syntax } p \text{ } b_{aux}[\mathcal{G}]}$$

The default rule declares that unless otherwise stated  $b_{aux}$  does not affect a grammar. Below we deal with the exceptions.

$$\begin{aligned}
 [2] \quad & b_{aux}[\emptyset] = \emptyset \\
 [3] \quad & b_{aux}[\mathcal{G}_1 \mathcal{G}_2] = b_{aux}[\mathcal{G}_1] b_{aux}[\mathcal{G}_2] \\
 [4] \quad & b_{aux}[\mathcal{G}] = \mathcal{G} \quad \textbf{otherwise}
 \end{aligned}$$

*Lexical Syntax* Lexical syntax grammars are translated to normal syntax grammars by encoding the symbols of the grammar to  $\langle \mathcal{A} \text{-LEX} \rangle$  symbols. Furthermore, for each symbol appearing in a lexical syntax section an injection from

the lexical into the context-free symbol is added.

- [5]  $\text{baux}[\llbracket \text{lexical syntax } p^* \rrbracket] = \langle \text{syntax } p^* \text{-LEX} \rangle$
- [6]  $\langle \alpha \text{-LEXs} \rangle = \text{when } \alpha =$
- [7]  $\langle \mathcal{A} \text{-LEXs} \rangle = \langle \mathcal{A} \text{-LEX} \rangle$
- [8]  $\langle \alpha^+ \beta^+ \text{-LEXs} \rangle = \langle \alpha^+ \text{-LEXs} \rangle \uparrow \langle \beta^+ \text{-LEXs} \rangle$
- [9]  $\langle \alpha \rightarrow \mathcal{A} \$ \text{-LEX} \rangle = \langle \alpha \text{-LEXs} \rangle \rightarrow \langle \mathcal{A} \text{-LEX} \rangle \$$
- [10]  $\langle \alpha \rightarrow \mathcal{A} \$ \text{-LEXs} \rangle = \langle \alpha \rightarrow \mathcal{A} \$ \text{-LEX} \rangle \langle \mathcal{A} \text{-LEXs} \rangle \rightarrow \langle \mathcal{A} \text{-CF} \rangle$
- [11]  $\langle p^* \text{-LEXs} \rangle = \text{when } p^* =$
- [12]  $\langle p_1^+ p_2^+ \text{-LEXs} \rangle = \langle p_1^+ \text{-LEXs} \rangle \uparrow \langle p_2^+ \text{-LEXs} \rangle$
- [13]  $\langle \text{syntax } p^* \text{-LEX} \rangle = \text{syntax } \langle p^* \text{-LEXs} \rangle$
- [14]  $\langle \mathcal{G}_1 \mathcal{G}_2 \text{-LEX} \rangle = \langle \mathcal{G}_1 \text{-LEX} \rangle \langle \mathcal{G}_2 \text{-LEX} \rangle$
- [15]  $\langle \emptyset \text{-LEX} \rangle = \emptyset$

*Context-free Syntax* Context-free syntax is treated similarly to lexical syntax. All symbols in the production are mapped to  $\langle \mathcal{A} \text{-CF} \rangle$  symbols. The important difference is that each adjacent pair of symbols in the left-hand side of a production is separated by the symbol  $\langle \text{LAYOUT?} \text{-CF} \rangle$ .

- [16]  $\text{baux}[\llbracket \text{context-free syntax } p^* \rrbracket] = \langle \text{syntax } p^* \text{-CF} \rangle$
- [17]  $\langle \alpha \text{-CFs} \rangle = \text{when } \alpha =$
- [18]  $\langle \mathcal{A} \text{-CFs} \rangle = \langle \mathcal{A} \text{-CF} \rangle$
- [19]  $\langle \alpha^+ \beta^+ \text{-CFs} \rangle = \langle \alpha^+ \text{-CFs} \rangle \uparrow \langle \text{LAYOUT?} \text{-CF} \rangle \uparrow \langle \beta^+ \text{-CFs} \rangle$
- [20]  $\langle \alpha \rightarrow \mathcal{A} \$ \text{-CF} \rangle = \langle \alpha \text{-CFs} \rangle \rightarrow \langle \mathcal{A} \text{-CF} \rangle \$$
- [21]  $\langle p \text{-CFs} \rangle = \langle p \text{-CF} \rangle$
- [22]  $\langle p^* \text{-CFs} \rangle = \text{when } p^* =$
- [23]  $\langle p_1^+ p_2^+ \text{-CFs} \rangle = \langle p_1^+ \text{-CFs} \rangle \uparrow \langle p_2^+ \text{-CFs} \rangle$
- [24]  $\langle \text{syntax } p^* \text{-CF} \rangle = \text{syntax } \langle p^* \text{-CFs} \rangle$
- [25]  $\langle \mathcal{G}_1 \mathcal{G}_2 \text{-CF} \rangle = \langle \mathcal{G}_1 \text{-CF} \rangle \langle \mathcal{G}_2 \text{-CF} \rangle$
- [26]  $\langle \emptyset \text{-CF} \rangle = \emptyset$

*Variables* Variables and lexical variables grammars introduce tokens that have the status of variables. The symbol constructor  $\langle \mathcal{A} \text{-VAR} \rangle$  is used to denote variables over the symbol  $\mathcal{A}$ . The left-hand sides of variable productions are interpreted as lexical syntax. The lexical value produced by such a left-hand side is given the type of a variable over the symbol in the right-hand side of the production. For each production in a variables grammar, two productions are generated. The first interprets the left-hand side of the production as a lexical pattern, i.e., the symbols on the left-hand side are lexical symbols and no layout between symbols can occur. The right-hand side is  $\langle \langle \mathcal{A} \text{-CF} \rangle \text{-VAR} \rangle$  indicating that the pattern is a variable over the context-free symbols  $\mathcal{A}$ . The second production injects  $\langle \mathcal{A} \text{-CF} \rangle$  variables into  $\langle \mathcal{A} \text{-CF} \rangle$  such that a variable can occur wherever an  $\langle \mathcal{A} \text{-CF} \rangle$  can occur.

- [27]  $\text{baux}[\llbracket \text{variables } p^* \rrbracket] = \text{syntax } \langle p^* \text{-VARs} \rangle$
- [28]  $\langle \text{-VARs} \rangle =$
- [29]  $\langle p_1^+ p_2^+ \text{-VARs} \rangle = \langle p_1^+ \text{-VARs} \rangle \uparrow \langle p_2^+ \text{-VARs} \rangle$

$$[30] \quad \langle \alpha \rightarrow \mathcal{A} \$\text{-VARs} \rangle = \langle \alpha\text{-LEXs} \rangle \rightarrow \langle \langle \mathcal{A}\text{-CF} \rangle\text{-VAR} \rangle \$ \\ \langle \langle \mathcal{A}\text{-CF} \rangle\text{-VAR} \rangle \rightarrow \langle \mathcal{A}\text{-CF} \rangle$$

Lexical variables are treated similarly, but their result sort is the corresponding lexical sort.

$$[31] \quad \text{baux}[\text{lexical variables } p^*] = \text{syntax } \langle p^*\text{-LEXVARs} \rangle \\ [32] \quad \langle \text{-LEXVARs} \rangle = \\ [33] \quad \langle p_1^+ p_2^+\text{-LEXVARs} \rangle = \langle p_1^+\text{-LEXVARs} \rangle ++ \langle p_2^+\text{-LEXVARs} \rangle \\ [34] \quad \langle \alpha \rightarrow \mathcal{A} \$\text{-LEXVARs} \rangle = \langle \alpha\text{-LEXs} \rangle \rightarrow \langle \langle \mathcal{A}\text{-LEX} \rangle\text{-VAR} \rangle \$ \\ \langle \langle \mathcal{A}\text{-LEX} \rangle\text{-VAR} \rangle \rightarrow \langle \mathcal{A}\text{-LEX} \rangle$$

*Ordering Grammars* The following equations specify the ordering of grammars, where the following order is obtained: lexical syntax, context-free syntax, lexical variables, and variables. We only show two of the equations, the other cases are similar.

$$[35] \quad \frac{\mathcal{G}_1 = \text{context-free syntax } p_1^*, \mathcal{G}_2 = \text{context-free syntax } p_2^*}{\mathcal{G}_1 \diamond \mathcal{G}_2 = \langle \text{context-free syntax } p_1^* p_2^*, \emptyset \rangle} \\ [36] \quad \frac{\mathcal{G}_1 = \text{context-free syntax } p_1^*, \mathcal{G}_2 = \text{lexical syntax } p_2^*}{\mathcal{G}_1 \diamond \mathcal{G}_2 = \langle \text{lexical syntax } p_2^*, \text{context-free syntax } p_1^* \rangle}$$

### 3.3.3 Parse Trees

Since we have expressed the meaning of lexical syntax and context-free syntax in terms of normal syntax productions, we do not have to extend the definition of parse trees, except for the encoding of symbols and grammars in the ATerm format. See §A.2 for the encoding and decoding of the newly introduced constructs. This entails that trees for lexical and context-free syntax have the same form. In particular, the structure assigned to lexical tokens by the grammar is retained in parse trees for tokens.

We will refine the equality predicate on trees such that layout is ignored. In considering whether two trees are equivalent it is likely that we do not want to consider layout. For this purpose it is not required to first translate a parse tree to an abstract syntax tree. It suffices to define two arbitrary layout trees as equivalent, as is done in the following extension of the equality predicate on trees.

**module** Basic-Sdf-Equality  
**imports** Kernel-Sdf-Equality<sup>2.5.8</sup> Basic-Sdf-ATerms<sup>A.2.2</sup>  
 Regular-Sdf-ATerms<sup>A.2.2</sup>  
**equations**

$$[1] \quad \frac{\text{symbol}(\text{type}(T_1)) = \langle \text{LAYOUT-CF} \rangle, \text{symbol}(\text{type}(T_2)) = \langle \text{LAYOUT-CF} \rangle}{T_1 \doteq T_2 = \top} \\ [2] \quad \frac{\text{symbol}(\text{type}(T_1)) = \langle \text{LAYOUT?-CF} \rangle, \text{symbol}(\text{type}(T_2)) = \langle \text{LAYOUT?-CF} \rangle}{T_1 \doteq T_2 = \top}$$



### 3.3.4 Discussion

*Lexical Layout* In some languages, such as FORTRAN, tokens can contain some kind of layout. In Heering *et al.* (1989) the symbol IGNORE is introduced for this purpose. This can be dealt with by separating the symbols in a lexical production by a *lexical layout* symbol just as this is done with context-free productions. This is not done in the current version because for most languages this is not necessary, but it is straightforward to add this feature to the normalization above.

*Implementation* A conventional implementation of parsers for lexical and context-free syntax is based on a separate scanner and parser. Such an implementation can be achieved for grammars as introduced here by separating productions for  $\langle\_ -\text{LEX}\rangle$  and  $\langle\_ -\text{VAR}\rangle$  symbols from productions for  $\langle\_ -\text{CF}\rangle$  symbols and generating a scanner based on finite automata for the first set of productions and by generating a parser for the second set of productions based on push-down automata. Scanner and parser communicate through some shared buffer-like data-structure. A requirement for this approach is that the lexical productions form a regular grammar. This can be enforced by specifying static constraints on lexical productions.

The parser generator for SDF2 described in Visser (1997e) does not depend on a separate scanner. Instead ‘lexical analysis’, i.e., parsing according to the productions for  $\langle\_ -\text{LEX}\rangle$  symbols, is incorporated in the parser. To cope with ambiguities and lookahead, generalized LR parsing is used. A similar approach is described by Salomon and Cormack (1989, 1995) under the name scannerless parsing using conventional LR techniques.

## 3.4 Restrictions

When a distinction is made between lexical and context-free syntax, lexical ambiguities have to be solved before tokens can be sent to the parser. This is usually done by applying rules such as ‘prefer longest match’, ‘prefer keywords’ and ‘prefer variables’.

By removing this distinction, as we did in the previous section, lexical ambiguities can be dealt with in the same way as context-free ambiguities. For example, in §3.1 we defined disambiguation by priorities, which applies both to lexical and context-free syntax. Furthermore, many lexical ambiguities are solved by considering the context in which tokens occur. For instance, the well-known problem of distinguishing an occurrence of the subrange 1..10 from two consecutive occurrences of the real numbers 1. and .10 in Pascal is solved automatically, because ranges and reals do not occur in the same context in the grammar.

However, not all lexical ambiguities can be solved by context or by means of priorities. Some lexical ambiguities need to be solved by rules such as ‘prefer longest match’ and ‘prefer literals’.

In this section we introduce two extensions of context-free grammars that are aimed at lexical disambiguation: *follow restrictions* and *reject productions*. A follow restriction  $\mathcal{A} \not\text{--} cc$  declares that the symbol  $\mathcal{A}$  can not be followed by any character in the character class  $cc$ . A reject production  $\alpha \rightarrow \mathcal{A} \{\text{reject}\}$  declares that any tree of type  $\mathcal{A}$  should be rejected if there exists a tree with the same

yield that has this reject production as root production. These constructs suffice for expressing most lexical disambiguation rules.

**Example 4.1** The definition of a simple expression language with nested comments in Example 3.1 contains two lexical ambiguities. First, the definition of lists of comment-words `ComWord*` is ambiguous. The string `abc` can be one comment word (a list of characters from the class `~[\ \t\n\\\/]`), but it can also be considered as a list of two comment-words `ab` and `c` or as `a` and `bc` or as `a` and `b` and `c`. We want to express that the longest possible comment-word should be selected. Second, the definition of identifiers and variables for identifiers and expressions overlap, i.e., `x` can be either an identifier or an expression of sort `Exp`. Here we want to express the rule ‘prefer variables’ that selects a variable over a lexical. These ambiguities are solved by the following rules:

```
lexical restrictions
  ComWord -/- ~[\ \t\n\\\/]
syntax
  <<Id-CF>-VAR> -> <Id-CF> {reject}
  <<Exp-CF>-VAR> -> <Id-CF> {reject}
```

The first rule states that a comment-word should not be followed by any of the characters in `~[\ \t\n\\\/]`. This solves the problem because it rules out all parses, except the one in which `abc` is one word. The last two rules state that variables should be preferred over identifiers.  $\square$

### 3.4.1 Syntax

A follow restriction has the form  $\mathcal{A}/cc$ . Follow restrictions are declared in a grammar starting with the keyword ‘restrictions’ followed by a list of restrictions. A reject production is a normal production attributed with the attribute ‘reject’.

```
module Restrictions-Sdf-Syntax
imports CC-Sdf-Syntax2.4.2
exports
  sorts Restriction Restrictions
  context-free syntax
    Symbols “/” CharClass → Restriction
    Restriction* → Restrictions
    “restrictions” Restrictions → Grammar
    “reject” → Attribute
  variables
    “restr”[0-9']* → Restriction
    “restr*”[0-9']* → Restriction*
    “restr+”[0-9']* → Restriction+
```

### 3.4.2 Projection

The function ‘R’ gives the restrictions of a grammar. The function  $\pi_{\mathcal{A}}$  looks up the restrictions for some symbol.

```
module Restrictions-Sdf-Projection
```

**imports** Restrictions-Sdf-Syntax<sup>3.4.1</sup>

**exports**

**context-free syntax**

Restrictions “+” Restrictions → Restrictions {**right**}  
 “R”(Grammar) → Restrictions  
 $\pi$  “\_” Symbol “(” Restrictions “)” → CharClass

**equations**

Concatenation of restrictions.

$$[1] \quad restr_1^* ++ restr_2^* = restr_1^* restr_2^*$$

The restrictions of a grammar.

$$[2] \quad R(\text{restrictions } restr^*) = restr^*$$

$$[3] \quad R(\mathcal{G}_1 \mathcal{G}_2) = R(\mathcal{G}_1) ++ R(\mathcal{G}_2)$$

$$[4] \quad R(\mathcal{G}) = \text{otherwise}$$

The restrictions for a symbol.

$$[5] \quad \pi_{\mathcal{A}}() = []$$

$$[6] \quad \pi_{\mathcal{A}}(\alpha \not\vdash cc restr^*) = \pi_{\mathcal{A}}(restr^*)$$

$$[7] \quad \pi_{\mathcal{A}}(\mathcal{A} \alpha \not\vdash cc restr^*) = cc \vee \pi_{\mathcal{A}}(\alpha \not\vdash cc restr^*)$$

$$[8] \quad \pi_{\mathcal{A}}(\mathcal{B} \alpha \not\vdash cc restr^*) = \pi_{\mathcal{A}}(\alpha \not\vdash cc restr^*) \text{ otherwise}$$

### 3.4.3 Normalization

No special normalization is needed for restrictions except the normal ordering and merging of grammars.

**module** Restrictions-Sdf-Normalization

**imports** Restrictions-Sdf-Syntax<sup>3.4.1</sup> CC-Sdf-Normalization<sup>2.4.2</sup>

**equations**

Merging and ordering of grammars.

$$[1] \quad \text{restrictions } restr_1^* \diamond \text{restrictions } restr_2^* = \langle \text{restrictions } restr_1^* restr_2^*, \emptyset \rangle$$

$$[2] \quad \text{restrictions } restr^* \diamond \text{syntax } p^* = \langle \text{syntax } p^*, \text{restrictions } restr^* \rangle$$

### 3.4.4 Discussion

The disambiguation rules presented above are derived from similar rules introduced by Salomon and Cormack (1989). The adjacency restriction of Salomon and Cormack (1989) is more general. It has the form  $\mathcal{A} \not\vdash \mathcal{B}$  and declares that symbols  $\mathcal{A}$  and  $\mathcal{B}$  should not be adjacent. Since this may require arbitrary long lookahead, we have chosen for the simpler follow restrictions, which can be implemented by restricting the lookahead of productions. The implementation of reject productions in SGLR parsing described in Visser (1997e) is more general than the implementation based on noncanonical SLR(1) parsing of Salomon and Cormack (1989).

We have not presented the interpretation of follow restrictions and reject productions as disambiguation devices. Follow restrictions can be interpreted as an extension of the well-formedness predicate on parse trees. If a follow restriction applies to a symbol, for any tree with that symbol as type, the character immediately next to the right-most character of its yield should not be contained in the restriction. For a discussion of the semantics of reject productions see Visser (1997e).

In the current situation lexical disambiguation rules have to be invented by the user. In SDF lexical disambiguation is completely taken care of in the scanner by means of a number of heuristics. These heuristics do cause problems in a number of cases. Therefore, it is attractive to have complete control over lexical disambiguation as is provided by restrictions and reject productions introduced here. However, it would be desirable if for most cases the necessary restrictions could be derived automatically from the grammar. Although some schemes have been considered, it is not yet clear how the derivation rules should be defined.

---

# 4

## Renaming and Modularization

---

In this chapter we introduce a module mechanism for reusing parts of syntax definitions. In order to adapt imported modules to specific applications and to avoid name clashes, a renaming mechanism is provided that can be used to rename symbols and productions. The renaming mechanism is also used in the definition of symbol aliases that can be used to define abbreviated names for large regular expressions. Renamings are also used to define symbol parameterization of modules.

---

### 4.1 Renamings

In the previous sections we have presented a number of features that enable more concise definition of syntax than plain context-free grammars. The grammars that can be defined are long monolithic lists of productions. To promote reuse of grammars we will introduce in §4.3 a module layer on top of grammars, such that parts of a language definition can be reused in various other definitions. To make the opportunities for reuse even greater we introduce here a renaming operator on grammars. Renamings enable the adaptation of a generic grammar to specific needs by renaming sorts and productions. A renaming is either a symbol renaming  $\mathcal{A} \Rightarrow \mathcal{B}$  that renames  $\mathcal{A}$  to  $\mathcal{B}$  or a production renaming  $p_1 \Rightarrow p_2$  that renames  $p_1$  to  $p_2$ . For example, the renaming

```
[Key => Var    Value => Term    Table => Subst
  lookup(Table, Key) -> Value
    => Subst "[" Var "]" -> Term ]
```

specifies the renaming of symbols `Key` and `Value` to `Var` and `Term`, respectively, and the renaming of the production `lookup(Table, Key) -> Value` to `Subst "[" Var "]" -> Term`.

Once we have defined renamings on grammars we can apply them in several situations: renaming of imported modules, symbol parameters of modules and symbol aliases. These will be the subject of the next sections.

#### 4.1.1 Syntax

A renaming is a list of symbol renamings of the form  $\mathcal{A} \Rightarrow \mathcal{B}$  and production renamings of the form  $p_1 \Rightarrow p_2$ .

```

module Renaming-Sdf-Syntax
imports Kernel-Sdf-Syntax2.3.1
exports
  sorts Renaming Renamings
  context-free syntax
    "[" Renaming* "]"      → Renamings
    Symbol "⇒" Symbol      → Renaming
    Production "⇒" Production → Renaming
  variables
    "ρ"[0-9']*      → Renamings
    "ρ" "*" [0-9']* → Renaming*
    "ρ" "+" [0-9']* → Renaming+

```

The only requirement on production renamings is that if  $\alpha \rightarrow \mathcal{A} \Rightarrow \beta \rightarrow \mathcal{B}$  is a renaming, then  $\alpha$  and  $\beta$  should be similar, i.e., the non-terminal parts should correspond. This entails that production renamings can only be used to rename literals between the arguments—the ‘syntax’—and not the order of the arguments.

#### 4.1.2 Projection

We define two projection functions for looking up the value of a symbol or a production in a list of renamings.

```

module Renaming-Sdf-Projection
imports Renaming-Sdf-Syntax4.1.1
exports
  context-free syntax
    Renamings "+" Renamings      → Renamings {right}
    "π" "-" Symbol "(" Renamings ")" → Symbol
    "π" "-" Production "(" Renamings ")" → Production
    "(" Symbols "⇒" Symbols ")" → Renamings
  equations
    Concatenation of renamings.

```

$$[1] \quad [\rho_1^*] ++ [\rho_2^*] = [\rho_1^* \rho_2^*]$$

Looking up the renaming of a symbol in a list of renamings.

$$\begin{aligned}
[2] \quad & \pi_{\mathcal{A}}([\mathcal{A} \Rightarrow \mathcal{B} \rho^*]) = \mathcal{B} \\
[3] \quad & \pi_{\mathcal{A}}([\mathcal{A}' \Rightarrow \mathcal{B} \rho^*]) = \pi_{\mathcal{A}}([\rho^*]) \quad \text{otherwise} \\
[4] \quad & \pi_{\mathcal{A}}([p \Rightarrow p' \rho^*]) = \pi_{\mathcal{A}}([\rho^*]) \\
[5] \quad & \pi_{\mathcal{A}}([\ ] ) = \mathcal{A}
\end{aligned}$$

Looking up the renaming of a production in a list of renamings.

$$\begin{aligned}
[6] \quad & \pi_p([\ ] ) = p \\
[7] \quad & \pi_p([p \Rightarrow p' \rho^*]) = p' \\
[8] \quad & \pi_{\alpha \rightarrow \mathcal{A}} \$([\alpha \rightarrow \mathcal{A} \Rightarrow \beta \rightarrow \mathcal{B} \rho^*]) = \beta \rightarrow \mathcal{B} \$ \\
[9] \quad & \pi_p([p' \Rightarrow p'' \rho^*]) = \pi_p([\rho^*]) \quad \text{otherwise} \\
[10] \quad & \pi_p([\mathcal{A} \Rightarrow \mathcal{B} \rho^*]) = \pi_p([\rho^*])
\end{aligned}$$

Abbreviation of a renaming of a list of symbols into another list of symbols.

$$\begin{aligned}
 [11] \quad & (\Rightarrow) = [] \\
 [12] \quad & (\alpha \Rightarrow) = [] \\
 [13] \quad & (\Rightarrow \beta) = [] \\
 [14] \quad & (\mathcal{A} \alpha \Rightarrow \mathcal{B} \beta) = [\mathcal{A} \Rightarrow \mathcal{B}] ++ (\alpha \Rightarrow \beta)
 \end{aligned}$$

This will be used for the instantiation of a list of formal parameters with a list of actual parameters.

### 4.1.3 Normalization

Now we can define the application of a renaming to a grammar. For each sort  $\mathcal{S}$  we define an application function  $(\mathcal{S}) \text{ Renamings} \rightarrow \mathcal{S}$  that applies a renaming to constructs of sort  $\mathcal{S}$ . We start by defining the renaming of symbols and productions. The rest is mainly a distribution of the renamings function over the constructs building a grammar.

**module** Kernel-Sdf-Renaming

**imports** Renaming-Sdf-Projection<sup>4.1.2</sup> Kernel-Sdf-Projection<sup>2.3.2</sup>

**exports**

**context-free syntax**

“(” Symbol “)” Renamings  $\rightarrow$  Symbol  
 “[” Symbol “]” Renamings  $\rightarrow$  Symbol  
 “(” Symbols “)\*” Renamings  $\rightarrow$  Symbols  
 “(” Production “)” Renamings  $\rightarrow$  Production  
 “(” Productions “)\*” Renamings  $\rightarrow$  Productions  
 “(” Grammar “)” Renamings  $\rightarrow$  Grammar

**equations**

Renaming a symbol. If the symbol is defined in the renaming it is replaced by its value in the renaming. Otherwise, the renaming is applied recursively to the subsymbols of the symbol, which is done by the function  $[-]\rho$ .

$$\begin{aligned}
 [1] \quad & (\mathcal{A}) \rho = \mathcal{B} \quad \textbf{when} \quad \pi_{\mathcal{A}}(\rho) = \mathcal{B}, \mathcal{A} \neq \mathcal{B} \\
 [2] \quad & (\mathcal{A}) \rho = [\mathcal{A}] \rho \quad \textbf{otherwise}
 \end{aligned}$$

Renaming a production works similarly. If the production is defined in the renaming it is replaced by its value. Otherwise, its symbols are renamed.

$$\begin{aligned}
 [3] \quad & (p) \rho = p' \quad \textbf{when} \quad \pi_p(\rho) = p', p' \neq p \\
 [4] \quad & (\alpha \rightarrow \mathcal{A} \$) \rho = (\alpha) * \rho \rightarrow (\mathcal{A}) \rho \$ \quad \textbf{otherwise}
 \end{aligned}$$

For all other grammar constructs, renaming is a homomorphism that applies the renamings to symbols and productions contained in the structure.

Renaming lists of symbols.

$$\begin{aligned}
 [5] \quad & (\alpha) * \rho = \quad \textbf{when} \quad \alpha = \\
 [6] \quad & (\alpha^+ \beta^+) * \rho = (\alpha^+) * \rho ++ (\beta^+) * \rho \\
 [7] \quad & (\mathcal{A}) * \rho = (\mathcal{A}) \rho
 \end{aligned}$$

Renaming lists of productions.

$$[8] \quad (p) * \rho = (p) \rho$$

$$\begin{aligned}
[9] \quad & (p_1^*)^* \rho = \text{when } p_1^* = \\
[10] \quad & (p_1^+ \ p_2^+)^* \rho = (p_1^+)^* \rho \mathbin{++} (p_2^+)^* \rho
\end{aligned}$$

Renaming grammars.

$$\begin{aligned}
[11] \quad & (\emptyset) \rho = \emptyset \\
[12] \quad & (\mathcal{G}_1 \ \mathcal{G}_2) \rho = (\mathcal{G}_1) \rho \ (\mathcal{G}_2) \rho \\
[13] \quad & (\text{syntax } p^*) \rho = \text{syntax } (p^*)^* \rho
\end{aligned}$$

The application of a renaming to a renaming denotes the composition of the renamings.

```

module Renaming-Sdf-Renaming
imports Kernel-Sdf-Renaming4.1.3
exports
  context-free syntax
    “(” Renamings “)” Renamings → Renamings
equations

```

A renaming  $\rho_2$  applied to a renaming  $\rho_1$   $(\rho_1)\rho_2$  denotes the composition of the renamings, i.e.,  $(x)(\rho_1)\rho_2 = ((x)\rho_1)\rho_2$ . This can be expressed by means of a single renaming by renaming the targets of  $\rho_1$  with  $\rho_2$  and adding  $\rho_2$  at the end of the list of renamings.

$$\begin{aligned}
[1] \quad & ([\ ] ) \rho = \rho \\
[2] \quad & ([\mathcal{A} \Rightarrow \mathcal{B} \ \rho^*]) \rho = [\mathcal{A} \Rightarrow (\mathcal{B}) \ \rho] \mathbin{++} ([\rho^*]) \rho \\
[3] \quad & ([p \Rightarrow p' \ \rho^*]) \rho = [p \Rightarrow (p') \ \rho] \mathbin{++} ([\rho^*]) \rho
\end{aligned}$$

For each of the extensions of the kernel we have to extend the renaming functions to the new constructors. See §A.3 for the specification of these extensions.

#### 4.1.4 Renaming Trees

If well-formed trees exist over a grammar that is renamed, the trees have to be renamed as well, if they have to be reused in the same context as the renamed grammar. For example, if equations over a grammar are defined, the equations must be renamed as well. Therefore, we extend the definition of renaming to parse trees.

```

module Renaming-Sdf-Trees
imports Kernel-Sdf-Trees2.5.5 CC-Sdf-Trees2.5.6 Basic-Sdf-Trees
  Regular-Sdf-Trees Kernel-Sdf-Renaming4.1.3
  Literals-Sdf-Normalization2.4.3
exports
  context-free syntax
    “(” ATerm “)” Renamings → ATerm
    rnargs(ATermList, ATermList) → ATermList
    mktree(Literal) → ATerm
    chartrees(Symbols) → ATermList
equations

```

Renaming an application. If the production is defined in the renaming, rename



the arguments and then rename the literals in the argument terms according to the new production.

$$[1] \quad \frac{\text{aterm}(\text{production}(\text{Prod})) \rho = \text{Prod}', \text{Prod}' \neq \text{Prod}, \text{Args} \rho = \text{Args}', \text{rnargs}(\text{args}(\text{Prod}'), \text{Args}') = \text{Args}''}{(\text{appl}(\text{Prod}, \text{Args})) \rho = \text{appl}(\text{Prod}', \text{Args}'')}$$

If the production is not defined in the renaming, then only rename the arguments.

$$[2] \quad (\text{appl}(\text{Prod}, \text{Args})) \rho = \text{appl}(\text{Prod}, (\text{Args}) \rho) \quad \textbf{otherwise}$$

Renaming is a homomorphism over the other tree constructors.

$$\begin{aligned} [3] \quad & (n) \rho = n \\ [4] \quad & (\text{amb}(\text{Args})) \rho = \text{amb}((\text{Args}) \rho) \\ [5] \quad & (Tl) \rho = [] \quad \textbf{when} \quad Tl = [] \\ [6] \quad & ([T]) \rho = [(T) \rho] \\ [7] \quad & ([T, Ts]) \rho = [(T) \rho, Ts'] \quad \textbf{when} \quad [Ts'] = ([Ts]) \rho \end{aligned}$$

Renaming the arguments.

$$[8] \quad \text{rnargs}([], []) = []$$

Insert literals of the new pattern.

$$[9] \quad \frac{\text{symbol}(\text{first}([Ts_1])) = L}{\text{rnargs}([Ts_1], Tl_2) = \text{mktree}(L) : \text{rnargs}(\text{rest}([Ts_1]), Tl_2)}$$

Skip literals of the old tree.

$$[10] \quad \frac{\text{symbol}(\text{type}(\text{first}([Ts_2]))) = L}{\text{rnargs}(Tl_1, [Ts_2]) = \text{rnargs}(Tl_1, \text{rest}([Ts_2]))}$$

Copy layout from old tree to new tree if layout is requested in new pattern.

$$[11] \quad \frac{\begin{array}{l} \text{symbol}(\text{first}([Ts_1])) = \langle \text{LAYOUT?}-\text{CF} \rangle, \\ \text{symbol}(\text{type}(\text{first}([Ts_2]))) = \langle \text{LAYOUT?}-\text{CF} \rangle \end{array}}{\text{rnargs}([Ts_1], [Ts_2]) = \text{first}([Ts_2]) : \text{rnargs}(\text{rest}([Ts_1]), \text{rest}([Ts_2]))}$$

Insert empty layout in the new tree.

$$[12] \quad \frac{\begin{array}{l} \text{symbol}(\text{first}([Ts_1])) = \langle \text{LAYOUT?}-\text{CF} \rangle, \\ \text{symbol}(\text{type}(\text{first}([Tl_2]))) \neq \langle \text{LAYOUT?}-\text{CF} \rangle \end{array}}{\text{rnargs}([Ts_1], Tl_2) = \text{appl}(\text{aterm}(\rightarrow \langle \text{LAYOUT?}-\text{CF} \rangle), []) : \text{rnargs}(\text{rest}([Ts_1]), Tl_2)}$$

Skip layout of the old tree.

$$[13] \quad \frac{\begin{array}{l} \text{symbol}(\text{first}(Tl_1)) \neq \langle \text{LAYOUT?}-\text{CF} \rangle, \\ \text{symbol}(\text{type}(\text{first}([Ts_2]))) = \langle \text{LAYOUT?}-\text{CF} \rangle \end{array}}{\text{rnargs}(Tl_1, [Ts_2]) = \text{rnargs}(Tl_1, \text{rest}([Ts_2]))}$$

In the other cases there is no layout or literal in either list. This means that it concerns an argument tree that should be copied from the old tree to the renamed tree.

$$[14] \quad \text{rnargs}([Ts_1], [Ts_2]) = \text{first}([Ts_2]) : \text{rnargs}(\text{rest}([Ts_1]), \text{rest}([Ts_2])) \\ \text{otherwise}$$

The function ‘mktree’ constructs a tree for a literal  $L$ , by constructing the production according to the definition in §2.4.3 and by generating the list of character codes.

$$[15] \quad \frac{\text{chars}(L) = \alpha, \text{aterm}(\alpha \rightarrow L) = \text{Prod}}{\text{mktree}(L) = \text{appl}(\text{Prod}, \text{chartrees}(\alpha))}$$

From a list of singleton character classes generate a term list of integers representing the character codes.

$$[16] \quad \text{chartrees}() = [] \\ [17] \quad \text{chartrees}([c] \alpha) = \text{int}(c) : \text{chartrees}(\alpha)$$

#### 4.1.5 Discussion

It would be desirable that renaming preserves well-formedness, i.e., if a tree  $T$  is well-formed under some grammar  $\mathcal{G}$ , it should also be well-formed when renamed with some renaming  $\rho$ . That is, we want that

$$\mathcal{G} \vdash T \Rightarrow (\mathcal{G})\rho \vdash (T)\rho$$

In fact we would like that renaming preserves the structure defined by a grammar, i.e., renaming the trees generated by a grammar gives the same trees as those generated by the renamed grammar:

$$(\mathcal{T}[\mathcal{G}])\rho = \mathcal{T}[(\mathcal{G})\rho]$$

Unfortunately, this is not the case for all renamings. If the argument sorts of a production are renamed using a production renaming, but the sorts are not renamed independently, the arguments of an application with that production have the wrong type after renaming. For instance, the renaming

[E "+" E -> E => Set "&" Set -> Set]

will change the notation of addition on the sort **E** into a binary operator **&** on the sort **Set**. Other constructs for sort **E** will still have type **E** after renaming, including the arguments of the **&** operator, which will hence not be well-formed. It is sufficient to require that in such cases the corresponding symbol renamings are present as well, i.e., the renaming

[E => Set E "+" E -> E => Set "&" Set -> Set]

does preserve well-formedness.

Also the interaction between regular expressions and renamings spoils the preservation property. For instance, consider the renaming

```
[{Int ",",}* => {Int ";"}*]
```

that is intended to rename lists of integers separated by commas into lists separated by semicolons. This will rename all symbols `{Int ",",}*` , but it will not rename the concatenation operators for this sort. The renaming

```
[{Int ",",}* => {Int ";"}*
  {Int ",",}* "," {Int ",",}* -> {Int ",",}*
=> {Int ";"}* ";" {Int ";"}* -> {Int ";"}*
  {Int ",",}* "+" {Int ",",}* -> {Int ",",}*
=> {Int ";"}* "+" {Int ";"}* -> {Int ";"}*
  {Int ",",}* "||" {Int ",",}* -> {Int ",",}*
=> {Int ";"}* "||" {Int ";"}* -> {Int ";"}*
  {Int ",",}* "+" {Int ",",}* -> {Int ",",}*
=> {Int ";"}* "+" {Int ";"}* -> {Int ";"}*
  {Int ",",}* "+" {Int ",",}* -> {Int ",",}*
=> {Int ";"}* "+" {Int ";"}* -> {Int ";"}*
]
```

is a well-formedness preserving renaming that does have the intended effect.

In all these cases correct renamings can be given that will preserve well-formedness and achieve the intended renaming, but these examples show that care has to be taken when writing down renamings. Ideally we would like to restrict the renamings such that the preservation property holds. It might also be possible to complete a renaming to guarantee well-formedness preservation as in the examples above. This is a matter for further study.

## 4.2 Aliases

The regular expressions introduced in §3.2 provide a way to concisely declare a number of productions without actually having to write them down. A problem with these regular expressions is that they can become rather large. This is a property that might make their use unattractive. Therefore, we introduce symbol aliases. An alias declaration introduces a short name for a complicated regular expression. All occurrences of the alias are replaced by their meaning. For example, the declarations

```
aliases
{Term ",",}*          -> Terms
{Var  ",",}*          -> Vars
Set[(Var "|->" Term)] -> Subst
```

introduce `Terms` and `Vars` as aliases for lists of `Term` and `Var`, respectively, and `Subst` as an alias for sets of pairs of variables and terms. This entails that all operations generated for list constructs also apply to `Terms` and `Vars` and all operations generated for sets apply to `Subst`.

Aliases are defined using the renamings of the previous section. An alias  $\mathcal{A} \rightarrow \mathcal{B}$  induces a grammar renaming  $[\mathcal{B} \Rightarrow \mathcal{A}]$ , which is applied to the entire grammar. Why then introduce this extra feature if we already have renamings? Renamings apply to a fixed grammar. Only the grammar to which the renaming is applied, including all imported grammars, is affected. An alias is a renaming of a symbol that also affects all modules that import the alias.

### 4.2.1 Syntax

An *alias* grammar consists of a list of aliases of the form  $\mathcal{A} \rightarrow \mathcal{B}$  that define the symbol  $\mathcal{B}$  to be an alias of symbol  $\mathcal{A}$ .

```

module Alias-Sdf-Syntax
imports Kernel-Sdf-Syntax2.3.1
exports
  sorts Alias Aliases
  context-free syntax
    “aliases” Aliases      → Grammar
    Symbol “→” Symbol     → Alias
    Alias*                  → Aliases
  variables
    “al”[ $\theta$ - $\theta'$ ]* → Alias
    “al*”[ $\theta$ - $\theta'$ ]* → Alias*
    “al+”[ $\theta$ - $\theta'$ ]* → Alias+

```

### 4.2.2 Projection

Concatenation of alias lists. Projection of the aliases and non-alias parts of a grammar.

```

module Alias-Sdf-Projection
imports Alias-Sdf-Syntax4.2.1
exports
  context-free syntax
    Aliases “++” Aliases → Aliases    {right}
    “Al”(Grammar)        → Aliases
    “Al̄”(Grammar)        → Grammar

```

#### equations

The function ‘Al’ gives all alias declarations of a grammar, ‘Al̄’ the grammar without alias declarations.

- |     |  |
|-----|--|
| [1] | $al_1^* ++ al_2^* = al_1^* al_2^*$   |
| [2] | $Al(\text{aliases } al^*) = al^*$  |
| [3] | $Al(\mathcal{G}_1 \mathcal{G}_2) = Al(\mathcal{G}_1) ++ Al(\mathcal{G}_2)$                               |
| [4] | $Al(\mathcal{G}) = \text{otherwise}$   |
| [5] | $\overline{Al}(\text{aliases } al^*) = \emptyset$  |
| [6] | $\overline{Al}(\mathcal{G}_1 \mathcal{G}_2) = \overline{Al}(\mathcal{G}_1) \overline{Al}(\mathcal{G}_2)$ |
| [7] | $\overline{Al}(\mathcal{G}) = \mathcal{G} \text{ otherwise}$   |

### 4.2.3 Normalization

Aliases are defined by renaming all alias symbols to their defined meaning. The function  $a[\_]$  produces a renaming from the alias declarations in the grammar and applies it to the non-alias parts of the grammar. The alias declarations are then attached to the renamed grammar. This is done in order to keep the following modular property:

$$a[\mathcal{G}_1 \mathcal{G}_2] = a[a[\mathcal{G}_1] a[\mathcal{G}_2]]$$

This entails that aliases can be replaced before flattening a module, after which the aliases are still part of the grammar and keep their forward renaming property.

```
module Alias-Sdf-Normalization
imports Alias-Sdf-Projection4.2.2 Kernel-Sdf-Normalization2.3.3
        Kernel-Sdf-Renaming4.1.3
exports
  context-free syntax
    “a[” Grammar “]” → Grammar
    rn(Aliases)        → Renamings
    symbols(Aliases)   → SymbolSet
```

#### equations

Replace all alias symbols by their definition by applying a renaming derived from the alias declarations to the non-alias parts of the grammar.

$$\begin{aligned} [1] \quad a[\mathcal{G}] &= \text{aliases } al^+ (\overline{AI}(\mathcal{G})) \text{ rn}(al^+) \quad \textbf{when } AI(\mathcal{G}) = al^+ \\ [2] \quad a[\mathcal{G}] &= \mathcal{G} \quad \textbf{otherwise} \end{aligned}$$

Build a renaming from a list of aliases. The target  $\mathcal{B}$  of the alias declaration  $\mathcal{A} \rightarrow \mathcal{B}$  is renamed to the source  $\mathcal{A}$ .

$$\begin{aligned} [3] \quad \text{rn}() &= [] \\ [4] \quad \text{rn}(\mathcal{A} \rightarrow \mathcal{B} \text{ } al^*) &= [\mathcal{B} \Rightarrow \mathcal{A}] \text{ ++ } \text{rn}(al^*) \end{aligned}$$

The symbols occurring in an alias declaration.

$$\begin{aligned} [5] \quad \text{symbols}(\text{aliases } al^*) &= \text{symbols}(al^*) \\ [6] \quad \text{symbols}(al^*) &= \{\} \quad \textbf{when } al^* = \\ [7] \quad \text{symbols}(\mathcal{A} \rightarrow \mathcal{B} \text{ } al^*) &= \{\mathcal{A} \mathcal{B}\} \cup \text{symbols}(al^*) \end{aligned}$$

Merging and ordering of grammars.

$$\begin{aligned} [8] \quad \text{aliases} &= \emptyset \\ [9] \quad \text{aliases } al_1^* \diamond \text{aliases } al_2^* &= \langle \text{aliases } al_1^* \text{ } al_2^*, \emptyset \rangle \\ [10] \quad \text{syntax } p^* \diamond \text{aliases } al^* &= \langle \text{aliases } al^*, \text{syntax } p^* \rangle \end{aligned}$$

Aliases themselves can also be subject to renamings.

```
module Alias-Sdf-Renaming
imports Kernel-Sdf-Renaming4.1.3 Alias-Sdf-Projection4.2.2
exports
  context-free syntax
    “(” Aliases “)”a” Renamings → Aliases
```

#### equations

Renaming of aliases.

$$\begin{aligned} [1] \quad (\text{aliases } al^*) \rho &= \text{aliases } (al^*) \text{a } \rho \\ [2] \quad (al^*) \text{a } \rho &= \quad \textbf{when } al^* = \\ [3] \quad (\mathcal{A} \rightarrow \mathcal{B}) \text{a } \rho &= (\mathcal{A}) \rho \rightarrow (\mathcal{B}) \rho \\ [4] \quad (al_1^+ \text{ } al_2^+) \text{a } \rho &= (al_1^+) \text{a } \rho \text{ ++ } (al_2^+) \text{a } \rho \end{aligned}$$

### 4.3 Modules

In this section we introduce a module framework for grammars to support management and reuse of parts of the grammar of a language. A modular definition consists of a list of named modules. Modules can be reused in other modules by means of imports. The body of a module is a list of exported and hidden grammars. Export and hiding provide a means to control what is visible from a module and what is local to that module. Hidden syntax is useful when the syntax definition formalism is coupled to a semantics formalism for the specification of the semantics of languages. Hidden syntax then plays the role of auxiliary functions. Since imports are abbreviations for grammars, an import can be hidden or exported. Modules can be parameterized by a list of symbols. An import can instantiate these parameters, although this is not required. Parameterization is an abbreviation for a renaming. When a module  $M[\alpha]$  is imported as  $M[\beta]$ , the formal parameters  $[\alpha]$  are renamed into the actual parameters  $[\beta]$ . An import can also be subject to a renaming of symbols and productions.

**Example 3.1 (Aliases and Renaming)** The following module defines the syntax of tables. A table is defined as an alias for a set of mappings from keys to values. The value assigned to a key can be looked up in a table using the access function `lookup`.

```
module Tables
  exports
    sorts Key Value Table
    aliases
      Set[(Key "|->" Value)] -> Table
    context-free syntax
      lookup(Table, Key) -> Value
```

Below we transform tables into mappings from variables to terms, thus obtaining a representation for substitutions. This is achieved by renaming the sorts in module `Tables` such that variables become the keys and terms the values in tables.

```
module Substitutions
  imports Terms
    Tables[Key => Var   Value => Term   Table => Subst
      lookup(Table, Key) -> Value
      => Subst "[" Var "]" -> Term ]
  exports
    context-free syntax
      Subst "(" Term ")" -> Term
```

The additional function applies a substitution to all variables in a term. □

**Example 3.2 (Map)** Using renaming a kind of polymorphic higher-order functions can be expressed. The following module defines a function that maps a function over the elements of a list. The function is defined for a given  $A$  and  $B$  that can be instantiated as needed.

```
module Map[A B]
```

```

exports
  sorts A B
  context-free syntax
    (A => B) "*" "(" A* ")" -> B*

```

The disadvantage of this kind of polymorphism is that for each instance of a polymorphic function, an explicit module import has to be done.  $\square$

**Example 3.3 (Parameterized Modules)** The following module defines the syntax of a list of conditional equations preceded by the keyword ‘equations’. This is the syntax of the equations part of an ASF+SDF module, which is parameterized by the syntax of some language. For each sort, productions defining the syntax of equations over that sort are defined. Note the use of the constrained iteration operator to define the bar (**Implies**) between conditions and conclusion as at least 3 equal signs.

```

module Equations
exports
  sorts Tag TagId CondEquation Equation Implies
    Condition Equations
  lexical syntax
    {[\\=]}3+      -> Implies
    [a-z0-9A-Z\\-]+ -> TagId
  aliases
    {Condition ","}+ -> Conditions
  context-free syntax
    "equations" CondEquation*      -> Equations
    Tag (Conditions Implies)? Equation -> CondEquation
    Tag Equation "when" Conditions -> CondEquation
    "[" TagId? "]"                 -> Tag

```

Next we define generic syntax for sorts Equation and Condition as follows:

```

module X-Equations[X]
exports
  sorts X Equation Condition
  context-free syntax
    X "=" X -> Equation
    X "=" X -> Condition
    X "!=" X -> Condition

```

To define the syntax of the equations part of a module  $M$ , the ASF+SDF Meta-Environment generates a module  $M$ -Equations that defines the syntax of these equations (Klint, 1993). This module imports the language independent syntax of equations and defines equations for the sorts declared in the module. With the parameterized module X-Equations we can express this by a module that contains an import for each declared sort. For instance, for **Boolean-Equations** we get the following module:

```

module Booleans-Equations
imports Booleans
  Equations
  X-Equations[Bool]

```

Observe that the sorts `Condition` and `Equation` are declared in two different modules. This is not problematic when these modules meet, because duplicate definitions are merged.  $\square$

### 4.3.1 Syntax

A modular syntax definition consists of a series of named module declarations. A module declaration consists of a list of sections, which are either exports or hiddens. A module name consists of a module identifier and an optional list of parameters. Module identifiers can contain slashes to enable the use of directory names in module names, e.g., `sdf/kernel/Syntax`. A module can import any number of other modules. An import consists of a module name with optionally a renaming applied to it. An import of a module  $M$  denotes the grammar declared in module  $M$ . An import can be contained in one of the exports or hiddens sections. In the latter case all syntax imported through that module is hidden and thus not exported from the module. Imports can also occur at the start of a module, outside any exports or hiddens section. In this case the imports are exported.

**module** Modular-Sdf-Syntax

**imports** Kernel-Sdf-Syntax<sup>2.3.1</sup> Renaming-Sdf-Syntax<sup>4.1.1</sup>

**exports**

**sorts** ModuleId ModuleName Import Imports Section Sections Module  
Definition ImpSection

**lexical syntax**

[A-Za-z0-9\_\-]+ → ModuleWord  
"/" ModuleWord → ModuleDir  
ModuleWord ModuleDir+ → ModuleId  
ModuleDir+ → ModuleId

**context-free syntax**

Module\* → Definition  
"module" ModuleName ImpSection\* Sections → Module  
"exports" Grammar → Section  
"hiddens" Grammar → Section  
Section\* → Sections  
ModuleId → ModuleName  
ModuleId "[" Symbols "]" → ModuleName  
id(ModuleName) → Attribute  
"imports" Imports → ImpSection  
ImpSection → Grammar  
Import\* → Imports  
ModuleName → Import  
ModuleName Renamings → Import  
 "(" Import ")" → Import {bracket}

**variables**

"Mid"[0-9']\* → ModuleId  
"M"[0-9']\* → ModuleName  
"s"[0-9']\* → Section  
"s" "\*" [0-9']\* → Section\*  
"s" "+" [0-9']\* → Section+  
"m"[0-9']\* → Module



$\text{"}m^*\text{"}[0-g']^* \rightarrow \text{Module}^*$   
 $\text{"}m+\text{"}[0-g']^* \rightarrow \text{Module}^+$   
 $\text{"}d\text{"}[0-g']^* \rightarrow \text{Definition}$   
 $\text{"}i\text{"}[0-g']^* \rightarrow \text{Import}$   
 $\text{"}i\text{"}^*\text{"}[0-g']^* \rightarrow \text{Import}^*$   
 $\text{"}i\text{"}^+\text{"}[0-g']^* \rightarrow \text{Import}^+$   
 $\text{"}is\text{"}[0-g']^* \rightarrow \text{ImpSection}$   
 $\text{"}is^*\text{"}[0-g']^* \rightarrow \text{ImpSection}^*$

### 4.3.2 Projection

Projection functions:  $\pi_M(d)$  yields the body of the module named  $M$ . ‘Exp’ yields the exported part of a module and ‘Hid’ yields the hidden part of a module.

**module** Modular-Sdf-Projection  
**imports** Modular-Sdf-Syntax<sup>4.3.1</sup> Booleans Kernel-Sdf-Projection<sup>2.3.2</sup>  
           Modular-Sdf-Renaming<sup>4.3.4</sup>

**exports**

**context-free syntax**

Import “ $\in$ ” Imports  $\rightarrow$  Bool  
 Sections “ $++$ ” Sections  $\rightarrow$  Sections {assoc}  
 Imports “ $++$ ” Imports  $\rightarrow$  Imports {assoc}  
 “ $\pi$ ” “ $_$ ” ModuleName “(” Definition “)”  $\rightarrow$  Sections  
  
 “Exp”(Sections)  $\rightarrow$  Grammar  
 “Hid”(Sections)  $\rightarrow$  Grammar

**equations**

Membership of a list of imports.

$$\begin{aligned}
 [1] \quad & i \in i_1^* i_2^* = \top \\
 [2] \quad & i \in i^* = \perp \quad \text{otherwise}
 \end{aligned}$$

Concatenation of section and imports lists.

$$\begin{aligned}
 [3] \quad & s_1^* ++ s_2^* = s_1^* s_2^* \\
 [4] \quad & i_1^* ++ i_2^* = i_1^* i_2^*
 \end{aligned}$$

Lookup of a module by its name in a list of modules. If a module name matches the module name searched for, its list of sections is yielded. If a parameterized module is imported without specifying any actual parameters, the parameters are left uninstantiated. If a list of actual parameters is given, these are used to rename the formal parameters into the actual parameters. The function  $(\_ \Rightarrow \_)$  constructs a renaming from the formal parameters to the actual parameters of a parameterized module. If no modules are found the empty list of sections is yielded.

$$\begin{aligned}
 [5] \quad & \pi_M(\text{module } M \ s^* \ m^*) = s^* ++ \pi_M(m^*) \\
 [6] \quad & \pi_{Mid}(\text{module } Mid[\alpha] \ s^* \ m^*) = s^* ++ \pi_{Mid}(m^*) \\
 [7] \quad & \pi_{Mid[\beta]}(\text{module } Mid[\alpha] \ s^* \ m^*) = (s^*) (\alpha \Rightarrow \beta) ++ \pi_{Mid[\beta]}(m^*) \\
 [8] \quad & \pi_M(\text{module } M' \ s^* \ m^*) = \pi_M(m^*) \quad \text{otherwise} \\
 [9] \quad & \pi_M() =
 \end{aligned}$$

Exported grammars from a list of sections.

- [10]  $\text{Exp}() = \emptyset$
- [11]  $\text{Exp}(s_1^+ s_2^+) = \text{Exp}(s_1^+) \text{Exp}(s_2^+)$
- [12]  $\text{Exp}(\text{exports } \mathcal{G}) = \mathcal{G}$
- [13]  $\text{Exp}(\text{hiddens } \mathcal{G}) = \emptyset$

Hidden grammars from a list of sections.

- [14]  $\text{Hid}() = \emptyset$
- [15]  $\text{Hid}(s_1^+ s_2^+) = \text{Hid}(s_1^+) \text{Hid}(s_2^+)$
- [16]  $\text{Hid}(\text{exports } \mathcal{G}) = \emptyset$
- [17]  $\text{Hid}(\text{hiddens } \mathcal{G}) = \mathcal{G}$

### 4.3.3 Normalization

We define the semantics of the modular constructs introduced above by means of a normalization function that yields the flattening of a module in a modular syntax definition by replacing each import by the body of the module it refers to.

Hidden productions are renamed by attaching the name of the hiding module. Since all productions occurring in a hiddens do not occur in another hiddens section (they should have been exported) it can never occur that two such renamed productions are imported into the same module. A consequence of production merging in this case is that an exported function becomes hidden if it is also in the hiddens part of the module.

We define the function  $\text{m}[d](M)$  that yields the grammar corresponding to module  $M$  in the definition  $d$ .

**module** Modular-Sdf-Normalization

**imports** Modular-Sdf-Projection<sup>4.3.2</sup> Modular-Sdf-Renaming<sup>4.3.4</sup>

Kernel-Sdf-Normalization<sup>2.3.3</sup> Grammar-Projection

**exports**

**context-free syntax**

- “m[” Definition “]” “(” ModuleName “)”  $\rightarrow$  Grammar
- hide(ModuleName, Grammar)  $\rightarrow$  Grammar
- hide(ModuleName, Productions)  $\rightarrow$  Productions

**hiddens**

**sorts** IG

**context-free syntax**

- “<” Imports “,” Grammar “>”  $\rightarrow$  IG
- imp “[” Definition “]” “(” Imports “,” Import “)”  $\rightarrow$  IG
- ims “[” Definition “]” “(” Imports “,” Imports “)”  $\rightarrow$  IG
- gra “[” Definition “]” “(” Imports “,” Grammar “)”  $\rightarrow$  IG

**equations**

Normalization of order of grammars.

- [1]  $\text{imports} = \emptyset$
- [2]  $\text{imports } i_1^* \diamond \text{imports } i_2^* = \langle \text{imports } i_1^* i_2^*, \emptyset \rangle$
- [3]  $\mathcal{G} \diamond \text{imports } i^* = \langle \text{imports } i^*, \mathcal{G} \rangle$  **otherwise**

Normalization of module sections. Exports and hiddens sections can be merged.

- [4]  $\text{module } M \text{ is }^* \text{ is } s^* = \text{module } M \text{ is }^* \text{ exports is } s^*$
- [5]  $s_1^* \text{ exports } \mathcal{G}_1 \text{ exports } \mathcal{G}_2 \ s_2^* = s_1^* \text{ exports } \mathcal{G}_1 \ \mathcal{G}_2 \ s_2^*$
- [6]  $s_1^* \text{ hiddens } \mathcal{G}_1 \text{ hiddens } \mathcal{G}_2 \ s_2^* = s_1^* \text{ hiddens } \mathcal{G}_1 \ \mathcal{G}_2 \ s_2^*$
- [7]  $s_1^* \text{ hiddens } \mathcal{G}_1 \text{ exports } \mathcal{G}_2 \ s_2^* = s_1^* \text{ exports } \mathcal{G}_2 \text{ hiddens } \mathcal{G}_1 \ s_2^*$

The semantics of a module named  $M$  in a definition  $d$  is expressed by  $m[d](M)$  and is the composition of the exported and hidden grammars of module  $M$  with all imports replaced by the exported grammars of the modules they refer to.

$$[8] \frac{\pi_M(d) = s^*, \text{ gra}[d](\text{Hid}(s^*)) = \langle i_1^*, \mathcal{G}_1 \rangle, \text{ gra}[d](i_1^*, \text{Exp}(s^*)) = \langle i_2^*, \mathcal{G}_2 \rangle}{m[d](M) = \mathcal{G}_2 \text{ hide}(M, \mathcal{G}_1)}$$

The function ‘hide’ marks all productions in the hiddens part of a module with the module name by attaching the attribute  $\text{id}(M)$  to it.

- [9]  $\text{hide}(M, \emptyset) = \emptyset$
- [10]  $\text{hide}(M, \mathcal{G}_1 \ \mathcal{G}_2) = \text{hide}(M, \mathcal{G}_1) \ \text{hide}(M, \mathcal{G}_2)$
- [11]  $\text{hide}(M, \text{syntax } p^*) = \text{syntax } \text{hide}(M, p^*)$
- [12]  $\text{hide}(M, \mathcal{G}) = \mathcal{G} \quad \text{otherwise}$
- [13]  $\text{hide}(M, \alpha \rightarrow \mathcal{A} \ \$) = \alpha \rightarrow \mathcal{A} \ \$ \ ++ \ \{\text{id}(M)\}$
- [14]  $\text{hide}(M, ) =$
- [15]  $\text{hide}(M, p_1^+ \ p_2^+) = \text{hide}(M, p_1^+) \ ++ \ \text{hide}(M, p_2^+)$

The function ‘gra’ expands all the imports in a grammar. It returns a structure  $\langle i^*, \mathcal{G} \rangle$ , which denotes a flattened grammar with the list of imports  $i^*$  that were expanded to flatten the grammar. This list is passed on to the rest of the flattening process in order to prevent multiple imports of the same module. This is important in particular in the presence of cyclic imports.

$$[16] \frac{\text{gra}[d](i_1^*, \mathcal{G}_1) = \langle i_2^*, \mathcal{G}_1' \rangle, \text{ gra}[d](i_2^*, \mathcal{G}_2) = \langle i_3^*, \mathcal{G}_2' \rangle}{\text{gra}[d](i_1^*, \mathcal{G}_1 \ \mathcal{G}_2) = \langle i_3^*, \mathcal{G}_1' \ \mathcal{G}_2' \rangle}$$

$$[17] \text{ gra}[d](i_1^*, \text{imports } i_2^*) = \text{ims}[d](i_1^*, i_2^*)$$

$$[18] \text{ gra}[d](i^*, \mathcal{G}) = \langle i^*, \mathcal{G} \rangle \quad \text{otherwise}$$

The function ‘ims’ yields the flattened grammars for a list of imports.

$$[19] \text{ ims}[d](i^*, ) = \langle i^*, \emptyset \rangle$$

$$[20] \frac{\text{imp}[d](i_1^*, i) = \langle i_3^*, \mathcal{G}_1 \rangle, \text{ ims}[d](i_3^*, i_2^*) = \langle i_4^*, \mathcal{G}_2 \rangle}{\text{ims}[d](i_1^*, i \ i_2^*) = \langle i_4^*, \mathcal{G}_1 \ \mathcal{G}_2 \rangle}$$

The function ‘imp’ yields the flattened grammar associated with the exported grammar of an import. The first list of imports denotes the imports that are already expanded. If a module was already imported it is not imported again. This is a protection against cyclic imports.

$$[21] \frac{\mathcal{G} = \text{if } M \in i^* \text{ then } \emptyset \text{ else } \text{Exp}(\pi_M(d)) \text{ fi}}{\text{imp}[d](i^*, M) = \text{gra}[d](i^* \ M, \mathcal{G})}$$

$$[22] \quad \frac{\mathcal{G} = \text{if } M \rho \in \mathbf{i}^* \text{ then } \emptyset \text{ else } (\text{Exp}(\pi_M(d))) \rho \text{ fi}}{\text{imp}\llbracket d \rrbracket(\mathbf{i}^*, M \rho) = \text{gra}\llbracket d \rrbracket(\mathbf{i}^* M \rho, \mathcal{G})}$$

As we will see in the next section, the renaming  $\rho$  that is applied to the exported part of the imported module  $M$  in the last equation above is also applied to the imports of that module and hence is applied recursively to all modules imported via  $M$ .

#### 4.3.4 Renaming

We extend the definition of renaming to renaming of module sections and imports. This includes the renaming of imports, and hence the renaming of renamings applied to imported modules.

**module** Modular-Sdf-Renaming

**imports** Renaming-Sdf-Renaming<sup>4.1.3</sup> Modular-Sdf-Syntax<sup>4.3.1</sup>  
Modular-Sdf-Projection<sup>4.3.2</sup>

**exports**

**context-free syntax**

“(” Sections “)” Renamings  $\rightarrow$  Sections

“(” Imports “)” Renamings  $\rightarrow$  Imports

**equations**

Renaming sections.

$$\begin{aligned} [1] \quad & (s^*) \rho = \text{when } s^* = \\ [2] \quad & (\text{exports } \mathcal{G}) \rho = \text{exports } (\mathcal{G}) \rho \\ [3] \quad & (\text{hiddens } \mathcal{G}) \rho = \text{hiddens } (\mathcal{G}) \rho \\ [4] \quad & (s_1^+ s_2^+) \rho = (s_1^+) \rho ++ (s_2^+) \rho \end{aligned}$$

Renaming a list of imports implies applying the renaming to all imported modules, i.e., attaching the renaming to each module name in the list of imports.

$$\begin{aligned} [5] \quad & (\text{imports } \mathbf{i}^*) \rho = \text{imports } (\mathbf{i}^*) \rho \\ [6] \quad & (M) \rho = M \rho \\ [7] \quad & (\mathbf{i}^*) \rho = \text{when } \mathbf{i}^* = \\ [8] \quad & (\mathbf{i}_1^+ \mathbf{i}_2^+) \rho = \mathbf{i}_3^+ \mathbf{i}_4^+ \text{ when } \mathbf{i}_3^+ = (\mathbf{i}_1^+) \rho, \mathbf{i}_4^+ = (\mathbf{i}_2^+) \rho \end{aligned}$$

If the imported module has already a renaming attached to it, the new renaming is applied to the first, yielding the composition of the two renamings.

$$[9] \quad (M \rho_1) \rho_2 = M (\rho_1) \rho_2$$

#### 4.3.5 Discussion

The modularization presented here is an extension with symbol parameters, import renamings and hidden imports of the modularization of ASF+SDF as implemented in the ASF+SDF Meta-Environment (Klint, 1993). The definition here is a pure ‘textual’ inclusion semantics of modularization. Hendriks (1991) describes both a textual normalization semantics and an incremental semantics

for modular constructs without renamings and hidden imports. The incremental implementation of modularization in the Meta-Environment becomes complicated in the presence of renamings, since items are created on the fly and can no longer be associated with a module. We have not addressed the issue of incremental parser generation and modular parser generation in a setting with renamings.

We deviate from the original design of ASF in that we do not incorporate the ‘origin rule’ that forbids identification of names that originate from different modules (Bergstra *et al.*, 1989b). This style forbids to have two modules with partly overlapping signatures, e.g., both introducing the same sort or function, that are imported in the same module, even if the overlap is intentional. The definition here is completely liberal in this respect. Productions that are imported via different routes are identified if they are the same.

In §4.1 we saw that renamings are not guaranteed to preserve well-formedness of trees. A further study of modular properties of grammars in the line of module algebra (Bergstra *et al.*, 1990) should give more insight into properties of good modularization. Some topics for study are: properties of trees and languages under renaming, ambiguity caused by union, interaction of regular expressions and renamings, modular properties of reject productions.



---

# 5

## The Syntax Definition Formalism SDF2

---

This chapter presents the assembly of the syntax definition formalism SDF2 from the features designed in previous chapters. This is mainly a matter of defining collecting modules that import the modules defined earlier. However, some features interfere. In some cases the normalization functions have to be extended to cover constructs introduced for other features. In other cases features have to be extended such that the orthogonality of another feature is maintained. The chapter concludes with a comparison of SDF2 to SDF and a discussion of anomalies and possible improvements to the formalism.

---

### 5.1 SDF2

Now we put the pieces together and define the syntax definition formalism SDF2, which is a generalization of SDF (Heering *et al.*, 1989). It covers all features available in SDF and adds several new ones. Furthermore, up to some small adaptations, SDF is textually (although not structurally) a subset of SDF2. This means that existing SDF definitions can be used almost literally as SDF2 definitions. The differences can be translated automatically by means of a migration tool.

The combination of features described earlier is achieved basically by combining them by means of imports into collecting modules. For each aspect of the definition, such as syntax, projection and normalization, a collecting module is defined. Here we show the collecting modules for the syntax and normalization of SDF2. The other modules can be found in Appendix A.4. Although we have tried to define features orthogonally, some interference between them is unavoidable. For instance, when we extend the syntax of symbols, normalization functions that deal with symbols are affected and have to be extended accordingly.

#### 5.1.1 Syntax

The syntax of SDF2 is simply the collection of the syntax of all features introduced so far. The syntax is extended with *lexical* and *context-free* priorities and restrictions, which arise as a result of the combination of Basic-Sdf with Priority-Sdf and Restriction-Sdf. The constructor ‘definition’ collects a list of modules into a single SDF definition.

The symbols  $\langle \text{Start} \rangle$  and  $\langle \text{START} \rangle$  serve to define grammars with a single

start symbol. In the normalization below productions will be added such that  $\langle \text{START} \rangle$  is the union of all sorts of the grammar. The symbol  $\langle \text{Start} \rangle$  is used to describe files that consist of a string over the language of  $\langle \text{START} \rangle$  followed by the end of file character.

In the Label extension a symbol can be labeled with a literal using the syntax  $L : \mathcal{A}$ . This extension is not further defined here. The priorities section is extended to deal with this extra symbol constructor.

```
module Sdf2-Syntax
imports Kernel-Sdf-Syntax2.3.1 Basic-Sdf-Syntax3.3.1 Modular-Sdf-Syntax4.3.1
        Regular-Sdf-Syntax3.2.1 Priority-Sdf-Syntax3.1.1 CC-Sdf-Syntax2.4.2
        Sorts-Sdf-Syntax2.4.1 Literals-Sdf-Syntax2.4.3 Label-Sdf-Syntax
        Restrictions-Sdf-Syntax3.4.1 Alias-Sdf-Syntax4.2.1
exports
  sorts SDF
  context-free syntax
    "<START>" → Symbol
    "<Start>" → Symbol
    "lexical" "priorities" Priorities → Grammar
    "context-free" "priorities" Priorities → Grammar
    "lexical" "restrictions" Restrictions → Grammar
    "context-free" "restrictions" Restrictions → Grammar
    "definition" Definition → SDF
  priorities
    Symbol "[Symbol" → Symbol > Literal ":" Symbol → Symbol
```

### 5.1.2 Normalization

We define the normalization function that normalizes a syntax definition by applying the normalization functions of the individual features. Here we have to deal with interaction between the normalization functions for the separate features and the constructs added to the formalism in other features.

```
module Sdf2-Normalization
imports Sdf2-Syntax5.1.1 Sdf2-ProjectionA.4 Sdf2-RenamingA.4
        Basic-Sdf-Normalization3.3.2 Modular-Sdf-Normalization4.3.3
        Priority-Sdf-Normalization3.1.3 Regular-Sdf-Normalization3.2.2
        Literals-Sdf-Normalization2.4.3 CC-Sdf-Normalization2.4.2
        Sorts-Sdf-Normalization2.4.1 Sorts-Sdf-Projection2.4.1
        Restrictions-Sdf-Normalization3.4.3 Alias-Sdf-Normalization4.2.3
exports
  context-free syntax
    normalize "[ SDF "]" "(" ModuleName "," Symbol ")" → Grammar
    topsorts "[ Grammar "]" "(" Symbol ")" → Grammar
    topsorts(Symbol, Symbols) → Productions
  equations
```

The normalization of an SDF2 definition is defined by the following equation. The function 'normalize' is parameterized with a module name denoting the top module to be normalized and a sort denoting the topsort of the definition. The definition is normalized by first expanding module  $M$  by means of function 'm'. Then the normalization functions 'b' (Basic), 'a' (Alias), 'r' (Regular), 'p'



(Priorities), ‘l’ (Literals) and ‘k’ (Kernel) are applied to the resulting grammar. The function ‘topsorts’, defined below, is used to add special productions for the top sorts of the definition and to remove productions not reachable from those top sorts.

$$[1] \quad \frac{k[l[p[r[topsorts[a[b[m[d](M)]]](\mathcal{A})]]]] = \mathcal{G}}{\text{normalize}[\text{definition } d](M, \mathcal{A}) = \text{reachable}(\{\langle \text{Start} \rangle\}, \mathcal{G})}$$

The function ‘topsorts’ adds a special production for the symbol  $\langle \text{Start} \rangle$ , which declares that a text over the grammar is a string of sort  $\mathcal{A}$  followed by the character representing the end of file. For each declared sort in the definition a production is added that defines that a text can be a string of that sort which starts and ends with layout.

$$[2] \quad \frac{\mathcal{G}' = \text{syntax } \mathcal{A} [\backslash \text{EOF}] \rightarrow \langle \text{Start} \rangle \text{ ++ topsorts}(\mathcal{A}, S(\mathcal{G}))}{\text{topsorts}[\mathcal{G}](\mathcal{A}) = \mathcal{G}' \mathcal{G}}$$

$$[3] \quad \text{topsorts}(\mathcal{A}, \alpha^+ \beta^+) = \text{topsorts}(\mathcal{A}, \alpha^+) \text{ ++ topsorts}(\mathcal{A}, \beta^+)$$

$$[4] \quad \text{topsorts}(\mathcal{A}, ) =$$

$$[5] \quad \text{topsorts}(\mathcal{A}, \mathcal{B}) = \langle \text{LAYOUT?}-\text{CF} \rangle \langle \mathcal{B}-\text{CF} \rangle \langle \text{LAYOUT?}-\text{CF} \rangle \rightarrow \mathcal{A}$$

### 5.1.3 Interaction

Several of the normalization functions are underdefined, i.e., the full SDF2 formalism contains more constructors than the extension for which they have been defined. Therefore, we must extend these functions accordingly.

#### exports

##### context-free syntax

- “<” Priorities “-LEXp” “>” → Priorities
- “<” Priorities “-CFp” “>” → Priorities
- “<” Restrictions “-LEX” “>” → Restrictions
- “<” Restrictions “-CF” “>” → Restrictions

#### hiddens

##### variables

- “L” → Literal

#### equations

The normalization function for regular expressions must be extended to the symbol constructors added in other extensions. The first equations express that sorts, character classes, literals and the symbols LAYOUT,  $\langle \text{START} \rangle$  and  $\langle \text{Start} \rangle$  do not generate any productions.

$$\begin{array}{ll} [6] & r[\mathcal{S}] = \emptyset \\ [8] & r[cc] = \emptyset \\ [10] & r[L] = \emptyset \end{array} \quad \begin{array}{ll} [7] & r[\text{LAYOUT}] = \emptyset \\ [9] & r[\langle \text{START} \rangle] = \emptyset \\ [11] & r[\langle \text{Start} \rangle] = \emptyset \end{array}$$

The following equations define that the productions generated for some symbol  $\mathcal{A}$  should be transformed into productions for lexical (context-free) productions if a lexical (context-free) version of the symbol occurs. This entails that first the productions for  $\mathcal{A}$  are generated by the recursive call and that these are transformed by the  $\langle \_-\text{LEX} \rangle$  ( $\langle \_-\text{CF} \rangle$ ) function.

$$[12] \quad r[\langle \mathcal{A}-\text{LEX} \rangle] = \langle r[\mathcal{A}]-\text{LEX} \rangle$$

$$\begin{aligned}
[13] \quad & r[\langle A\text{-CF} \rangle] = \langle r[A] \text{-CF} \rangle \\
[14] \quad & r[\langle A\text{-VAR} \rangle] = r[A]
\end{aligned}$$

This is an example of the context-sensitivity of the generation of productions from symbols. The meaning of  $\langle \text{Id} * \text{-CF} \rangle$  is different from that of  $\langle \text{Id} * \text{-LEX} \rangle$ .

*Basic* Literals and character classes do not need the  $\langle \text{-LEX} \rangle$  or  $\langle \text{-CF} \rangle$  constructor, because they are lexical by definition

$$\begin{aligned}
[15] \quad & \langle cc\text{-LEX} \rangle = cc & [16] \quad & \langle cc\text{-CF} \rangle = cc & [17] \quad & \langle cc\text{-VAR} \rangle = cc \\
[18] \quad & \langle L\text{-LEX} \rangle = L & [19] \quad & \langle L\text{-CF} \rangle = L & [20] \quad & \langle L\text{-VAR} \rangle = L
\end{aligned}$$

*Basic + Priorities* Equations for the normalization of lexical and context-free priorities that were added at the level of SDF2.

$$\begin{aligned}
[21] \quad & \frac{\mathcal{G}_1 = \text{context-free priorities } pr^*, \mathcal{G}_2 = \text{context-free syntax } p^*}{\mathcal{G}_1 \diamond \mathcal{G}_2 = \langle \text{context-free syntax } p^*, \text{context-free priorities } pr^* \rangle} \\
[22] \quad & \frac{\mathcal{G}_1 = \text{context-free priorities } pr^*, \mathcal{G}_2 = \text{lexical syntax } p^*}{\mathcal{G}_1 \diamond \mathcal{G}_2 = \langle \text{lexical syntax } p^*, \text{context-free priorities } pr^* \rangle} \\
[23] \quad & \frac{\mathcal{G}_1 = \text{context-free priorities } pr_1^*, \mathcal{G}_2 = \text{context-free priorities } pr_2^*}{\mathcal{G}_1 \diamond \mathcal{G}_2 = \langle \text{context-free priorities } pr_1^*, pr_2^*, \emptyset \rangle}
\end{aligned}$$

Context-free priorities are priority declarations for context-free productions and are abbreviations of normal priorities in the same way that context-free syntax is an abbreviation for a certain style of normal syntax. The productions in the priorities sections are thus treated with the same  $\langle \text{-CF} \rangle$  functions as context-free productions.

$$\begin{aligned}
[24] \quad & \text{baux}[\langle \text{context-free priorities } pr^* \rangle] = \langle \text{priorities } pr^* \text{-CF} \rangle \\
[25] \quad & \langle \text{priorities } pr^* \text{-CF} \rangle = \text{priorities } \langle \text{norm}[\langle pr^* \rangle] \text{-CFp} \rangle \\
[26] \quad & \langle pr^* \text{-CFp} \rangle = \text{when } pr^* = \\
[27] \quad & \langle pr_1^+, pr_2^+ \text{-CFp} \rangle = \langle pr_1^+ \text{-CFp} \rangle ++ \langle pr_2^+ \text{-CFp} \rangle \\
[28] \quad & \langle p_1 > p_2 \text{-CFp} \rangle = \langle p_1 \text{-CF} \rangle > \langle p_2 \text{-CF} \rangle \\
[29] \quad & \langle p_1 \text{ as } p_2 \text{-CFp} \rangle = \langle p_1 \text{-CF} \rangle \text{ as } \langle p_2 \text{-CF} \rangle
\end{aligned}$$

Similarly for lexical priorities.

$$\begin{aligned}
[30] \quad & \text{baux}[\langle \text{lexical priorities } pr^* \rangle] = \langle \text{priorities } pr^* \text{-LEX} \rangle \\
[31] \quad & \langle \text{priorities } pr^* \text{-LEX} \rangle = \text{priorities } \langle \text{norm}[\langle pr^* \rangle] \text{-LEXp} \rangle \\
[32] \quad & \langle pr^* \text{-LEXp} \rangle = \text{when } pr^* = \\
[33] \quad & \langle pr_1^+, pr_2^+ \text{-LEXp} \rangle = \langle pr_1^+ \text{-LEXp} \rangle ++ \langle pr_2^+ \text{-LEXp} \rangle \\
[34] \quad & \langle p_1 > p_2 \text{-LEXp} \rangle = \langle p_1 \text{-LEX} \rangle > \langle p_2 \text{-LEX} \rangle \\
[35] \quad & \langle p_1 \text{ as } p_2 \text{-LEXp} \rangle = \langle p_1 \text{-LEX} \rangle \text{ as } \langle p_2 \text{-LEX} \rangle
\end{aligned}$$

*Basic + Restrictions*

$$\begin{aligned}
[36] \quad & \text{baux}[\langle \text{lexical restrictions } restr^* \rangle] = \text{restrictions } \langle restr^* \text{-LEX} \rangle \\
[37] \quad & \langle restr^* \text{-LEX} \rangle = \text{when } restr^* = \\
[38] \quad & \langle restr_1^+ restr_2^+ \text{-LEX} \rangle = \langle restr_1^+ \text{-LEX} \rangle ++ \langle restr_2^+ \text{-LEX} \rangle
\end{aligned}$$

$$[39] \quad \langle \alpha \not\vdash cc\text{-LEX} \rangle = \langle \alpha\text{-LEXs} \rangle \not\vdash cc$$

$$[40] \quad \text{baux}[\llbracket \text{context-free restrictions } restr^* \rrbracket] = \text{restrictions } \langle restr^*\text{-CF} \rangle$$

$$[41] \quad \langle restr^*\text{-CF} \rangle = \text{when } restr^* =$$

$$[42] \quad \langle restr_1^+ restr_2^+\text{-CF} \rangle = \langle restr_1^+\text{-CF} \rangle ++ \langle restr_2^+\text{-CF} \rangle$$

$$[43] \quad \langle \not\vdash cc\text{-CF} \rangle = \not\vdash cc$$

$$[44] \quad \langle \mathcal{A} \alpha \not\vdash cc\text{-CF} \rangle = \langle \mathcal{A}\text{-CF} \rangle \beta \not\vdash cc \quad \text{when } \langle \alpha \not\vdash cc\text{-CF} \rangle = \beta \not\vdash cc$$

### Labels

$$[45] \quad r[\llbracket L : \mathcal{A} \rrbracket] = r[\llbracket \mathcal{A} \rrbracket]$$

$$[46] \quad \langle L : \mathcal{A}\text{-LEX} \rangle = L : \langle \mathcal{A}\text{-LEX} \rangle$$

$$[47] \quad \langle L : \mathcal{A}\text{-CF} \rangle = L : \langle \mathcal{A}\text{-CF} \rangle$$

$$[48] \quad \langle L : \mathcal{A}\text{-VAR} \rangle = L : \langle \mathcal{A}\text{-VAR} \rangle$$

### Hiding Productions

$$[49] \quad \text{hide}(M, \text{context-free syntax } p^*) = \text{context-free syntax } \text{hide}(M, p^*)$$

$$[50] \quad \text{hide}(M, \text{lexical syntax } p^*) = \text{lexical syntax } \text{hide}(M, p^*)$$

$$[51] \quad \text{hide}(M, \text{variables } p^*) = \text{variables } \text{hide}(M, p^*)$$

$$[52] \quad \text{hide}(M, \text{lexical variables } p^*) = \text{lexical variables } \text{hide}(M, p^*)$$

### Aliases

$$[53] \quad \frac{\mathcal{G}_1 = \text{aliases } al^*, \mathcal{G}_2 = \text{sorts } \alpha}{\mathcal{G}_1 \diamond \mathcal{G}_2 = \langle \text{sorts } \alpha, \text{aliases } al^* \rangle}$$

## 5.2 Comparison to SDF

SDF2 was developed as a generalization of SDF (Heering *et al.*, 1989). We briefly list the differences between the two formalisms.

### 5.2.1 Semantics

SDF defines the semantics of a syntax definition by means of mappings to other formalisms. The lexical syntax is mapped to a regular grammar. The context-free syntax is mapped to a context-free grammar. From the entire definition a first-order many-sorted algebraic signature is derived. A parse tree for a string according to the grammar is translated to a term or abstract syntax tree over the signature. In SDF2 parse trees are defined by means of a well-formedness predicate on ATerms based directly on (the normal form of) a syntax definition. The strings of the language defined by a grammar are obtained via the function *yield*. No external formalism is used to define trees. In this way the notions of grammar and signature that were related via mappings in SDF are completely integrated in SDF2.

### 5.2.2 Lexical and Context-free Syntax

SDF integrates lexical syntax and context-free syntax in one formalism. However, this integration is only at the level of the formalism; on the level of the implementation these are separated. The lexical syntax is mapped to a regular grammar (hence the specification of the lexical syntax should also be regular). The context-free syntax is translated to a context-free grammar. In SDF2 the integration of lexical and context-free syntax is completed. All other features are orthogonal with respect to lexical and context-free syntax. For instance, character classes and regular expressions can be used in exactly the same way in lexical productions and context-free productions.

### 5.2.3 Lexical Disambiguation

SDF has several built-in lexical disambiguation rules that are applied to the token stream before tokens are passed to the parser. SDF2 has no built-in lexical disambiguation rules, but provides reject productions to express the prefer literals rule and follow restrictions to express longest match disambiguation.

### 5.2.4 Character Classes

In SDF the syntax of character classes is defined lexically. In SDF2 character classes are defined by means of context-free constructors. This makes the definition of normalization of character classes much easier. The differences with character classes in SDF are: numeric characters have a decimal interpretation instead of an octal interpretation, there is no syntactic limit to the range of numeric characters, all characters except letters and digits have to be escaped using a slash. In SDF2 character classes have a numeric interpretation, that is, each character class is normalized to an ordered and non-overlapping list of numeric characters and ranges of characters.

### 5.2.5 Lists

SDF only provides list sorts that can be used in the left-hand sides of productions. Furthermore, lists are not orthogonally defined. In the lexical syntax no iteration with separator is provided. In SDF only sorts can be used on the right-hand side of a production. This means that list sorts cannot be the result of functions. In order to define a function with a list as result, a new sort has to be introduced into which the list sort is injected. Furthermore, to concatenate the lists that result from a function a concatenation function should be defined.

SDF2 provides an expressive set of regular expressions that are treated as first-class citizens. Regular expressions can be used where ever any another symbol can be used. In general, all symbols that can be used in the left-hand side of a production can also be used as output symbols.

### 5.2.6 Priorities

The priorities declarations of SDF2 are the same as in SDF with the following exceptions: No abbreviations of productions in priorities are supported because of the problematic semantics in a setting with modules. No <-chains are provided. The implementation does not provide the multi-set filter interpretation

of priorities.

### 5.2.7 Reuse

SDF does not provide renamings, module parameterization, hidden imports, and aliases.

## 5.3 Discussion and Concluding Remarks

We have presented the modular design of a family of syntax definition formalisms. The result is a uniform formalism for syntax definition designed for extensibility. A guiding principle in the design is the orthogonality of the features with respect to one another. As a consequence it is easy to replace a feature by a variant or to add a new feature without affecting the design and implementation of all other features.

### 5.3.1 Parser Generation

The direct motivation for this work was the specification of a parser generator for SDF. Many of the techniques presented in this paper were originally developed for the translation of SDF to intermediate languages like context-free and regular grammars as prescribed by the SDF reference manual (Heering *et al.*, 1989). Gradually it became clear that the difficulty of this project was due to the monolithic design of SDF. The features presented in this paper are combined in the formalism SDF2 that is intended to replace SDF. The specification of a parser generator for SDF2 was easier due to the uniform abstract syntax and elimination of cases by normalization. The tables generated by the generator are interpreted by the generic scannerless generalized-LR parser described in Visser (1997e).

### 5.3.2 Disambiguation

Priorities are interpreted as a well-formedness requirement on parse forests, which could be operationalized as a filter on parse forests as prescribed by Heering *et al.* (1989). This approach can be extended to other disambiguation methods as described in Klint and Visser (1994).

We have provided some features for disambiguation of ambiguous context-free grammars. There remain a large number of ambiguities that can not be solved with these mechanisms. Some more advanced disambiguation methods are described in Klint and Visser (1994). Here we list some ideas for improvements of the current scheme.

The priority relation  $>$  on productions does not allow a distinction between the arguments of the productions to which it applies. In several cases it would be useful to restrict the relation to certain arguments. For instance, the priority declaration

$$T \ T \rightarrow T \ > \ \text{"let"} \ V \ "=" \ T \ \text{"in"} \ T \rightarrow T$$

does correctly forbid the usage of a `let` expression as the first argument of an application. However it also forbids the usage of `let` as the last argument of an application, for which there is no reason. An extension of the notation could be

$$T \ T \rightarrow T \ \{1\} > \text{"let"} \ V \ \text{"="} \ T \ \text{"in"} \ T \rightarrow T$$

to declare the desired disambiguation. There would be no implementation problems with such an extension.

A case for non-standard disambiguation is in ambiguous equations. In §4.3 we gave as an example the specification of the syntax of conditional equations. It can occur that equations are ambiguous due to injections. If a symbol  $\mathcal{A}$  is injected in  $\mathcal{B}$ , then an equation over two  $\mathcal{A}$  expressions can be interpreted both as  $\mathcal{A}$  and  $\mathcal{B}$  equations. A possible interpretation of such an ambiguity is to take both possibilities. This is done in the definition of multi-level algebraic specifications in Visser (1996), where ambiguous equations can occur due to overloading of functions.

In the implementation of SDF in the ASF+SDF Meta-Environment an undocumented disambiguation method is used. As a simplification of the multi-set ordering, trees with fewer injections are preferred over trees with more injections. Such a method is needed to disambiguate conditions of equations. This method has not been implemented as part of the SDF2 tools, but can be added as a post-parse filter without problems.

We have defined follow restrictions and reject productions to express lexical disambiguation rules. We omitted the definition of these methods as an extension of the well-formedness predicate on parse trees. See Visser (1997e) for a discussion of the semantics of these methods and for a discussion of automatic lexical disambiguation.

### 5.3.3 Renaming

Modules associate a name with a grammar. Grammars can be combined by module imports. Export and hiding provide control over visibility of grammars. New with respect to the modularization of SDF are renamings and hidden imports. In the current definition of renaming productions, only the literal skeleton of the production can be changed, but the order of the arguments stays the same. Sometimes it is desirable to change the syntax of a production and also make a permutation of the arguments. A notation for such permutations should be devised by means of some kind of indexing. The problem with such a notation is that the current definition reuses the syntax of productions literally in the definition of renamings. Changing the syntax of productions will thus be applicable everywhere. The label facility (see below) could be used for this purpose. Unfortunately, renamings are not guaranteed to preserve well-formedness of parse trees. Further study is needed to find a set of sufficient requirements on renamings that do guarantee well-formedness.

### 5.3.4 Labels

A feature that has not been discussed are labels. Labels are intended to be used as ‘field names’ of a record. For instance, consider the following production defining the syntax of assignments in an imperative language:

```
var : Var "!=" value : Exp -> Stat {cons(assign)}
```

The two arguments are labeled with `var` and `value`, respectively. From this information we can derive the following syntax for projection functions based on the field names:

```

Stat "." var    -> Var
Stat "." value -> Exp

```

This should be accompanied by the the defining equations for these functions.

### 5.3.5 Derived Syntax

Regular expressions are considered as name constructors that are used to make new names out of existing ones. A normalization function adds canonical productions defining the regular operators. For instance,  $\mathcal{A}?$  denotes an optional  $\mathcal{A}$  and is defined by the productions  $\mathcal{A} \rightarrow \mathcal{A}?$  and  $\rightarrow \mathcal{A}?$ . However, there is no restriction on the use of these name constructors. Other defining productions can be added by the user. In the context of algebraic specification this means for instance that users can specify functions that have lists ( $\mathcal{A}^*$ ) as result.

Regular expressions are an example of derived syntax: Given some symbol or even production in the grammar, other productions are derived. Many other applications of derived syntax could be useful.

Sometimes it is useful to explicitly indicate empty constructs and injections. This could be accommodated by generating syntax for explicitly matching injection functions and  $\epsilon$ -functions, i.e., if  $A \rightarrow B$  is a production then also " $\text{inj}A\text{-}B$ " ( $A \rightarrow B$ ). Similarly for  $\epsilon$ , if  $\rightarrow A$  a production, then also " $\text{empty}A$ "  $\rightarrow A$ . These constructors should of course match with their origins. This can be done by translating these functions internally to the real injection or  $\epsilon$  function. See Dinesh (1995) for some interesting remarks on injections in ASF+SDF.

A structure editor provides facilities to manipulate sentential forms. This requires the specification of the syntax of symbol placeholders. For each symbol  $A$  that is not a literal add a production " $\langle A \rangle$ "  $\rightarrow A$ .

Another case of this kind is the generation of explicit type casts  $S \text{ ":" } "S" \rightarrow S \{ \text{cast} \}$  (like bracket attribute) to constrain the type of an overloaded entity. This would be similar to the `no-operator` attribute in SDF.

It would be even better to make syntax derivation user-definable by providing schemas such as discussed above.

### 5.3.6 Polymorphic Syntax

The definition of regular expressions by introducing new productions is an instance of second order quantification. The generalization of this approach to two-level grammars in Visser (1997d) provides the syntactic counterpart of the two-level specifications in Meinke (1992) and the multi-level specifications in Visser (1996). Generic productions are written as production schemata. The syntax of symbol constructors is described by means of a second level grammar.

### 5.3.7 Dynamic Syntax

Another open problem is the formal description of languages with an extensible syntax. Programs in such languages can contain grammars that define part of the syntax of the program itself. An example of extensible syntax is the syntax of equations in ASF+SDF. Several other instances exist, e.g., Cardelli *et al.* (1994), and Vittek (1994) (Elan). All these approaches treat the meta-language and object language differently. A formal approach to this problem

would specify the syntax of the base language and the grammars it can specify and the lifting of these grammars to meta-level grammars.

### 5.3.8 Design Methodology

We have presented a large specification. We approached this using a rigorous modularization of the specification in a matrix of modules. For each feature the syntax and tools are described in separate modules. In this way it becomes feasible to flexibly include and exclude parts of a language definition. Some parts of the specification such as the ATerm encoding are not very interesting. It would be better if those parts could be generated using a simple rule.

The main technique we have applied is that of definition of features by normalization, i.e., transformation to a subset of the language. The great advantage of normalization is that many features can be provided to enhance the expressiveness of the language while defining the semantics of the formalism on a small set of kernel features in which the other features are expressed. Normalization has also its disadvantages. The semantics of various features is defined indirectly, which makes reasoning about them more troublesome. Furthermore, parse trees over a grammar use the normalized productions, which can look rather different than their origins. It would be desirable to use normalization equations rather than functions in order to be able to reason about equivalence of syntax definitions. The problem with such an approach is the lack of control over normalization. A solution could be the use of strategies such as described in Luttik and Visser (1997).

The modularization of the formalism and hence the modularization of the normalization in separate normalization functions for each feature made the specification of normalization feasible. A normalization function that would in one pass over the grammar normalize it would be a very complex. However, the modularization also hides the interaction between features. When defining a normalization function for an extension of the kernel, only those constructs introduced are normalized. The combination of features prompts the extension of the normalization function to new constructs. Often this can be achieved by innocent distribution equations, but in some cases the interaction between features is more problematic. In particular, the interaction of renamings with other features needs more study.

Language design is a software engineering process. A language definition gets better developed if it is actually used in a prototype implementation. The parts of the specification of SDF2 that are used in the parser generator, i.e., the normalization, were developed on demand. Especially the fragment of SDF2 that corresponds to SDF was developed first, because most syntax definitions fed to the parser generator were converted SDF definitions. Other parts of the specification, such as well-formedness or equality that are not directly used in tools were developed later. But these parts are important because they define the correctness criteria for implementations. The well-formedness checker can be used to validate the output of a parser for a grammar. The equality checker can be used to validate a matching algorithm for terms.

The design approach we have used for SDF2 has led to an infrastructure for further study of syntax definition and experimentation with new features. It is indeed very easy to extend the specification in order to construct subsets or supersets of the formalism, or to replace a feature by a variant.



---

# A

## Auxiliary Modules for the Specification of SDF2

---

In this appendix we include several auxiliary modules used in the specification of SDF2.

---

### A.1 Literals

```
module Literals
imports Layout
exports
  sorts Literal
  lexical syntax
    “\”~□ → EscChar
    “\”[01][0-7][0-7] → EscChar
    ~[\000-\040”\] → L-Char
    [␣\t\n] → L-Char
    EscChar → L-Char
    “\””L-Char*“\”” → Literal
  variables
    “L”[0-9’]* → Literal
```

### A.2 ATerms

#### A.2.1 Constructors

```
module Grammar-Tree-Constructors
imports ATerms2.5.1
exports
  context-free syntax
    “empty-grammar” → AFun
    “conc-grammars” → AFun

module CC-Sdf-Tree-Constructors
imports Kernel-Sdf-Tree-Constructors2.5.2
exports
  context-free syntax
    “char-class” → AFun
```

“range” → AFun

**module** Sorts-Sdf-Tree-Constructors  
**imports** Kernel-Sdf-Tree-Constructors<sup>2.5.2</sup>  
**exports**  
  **context-free syntax**  
  “sort” → AFun

**module** Literals-Sdf-Tree-Constructors  
**imports** Kernel-Sdf-Tree-Constructors<sup>2.5.2</sup>  
**exports**  
  **context-free syntax**  
  “lit” → AFun

**module** Regular-Sdf-Tree-Constructors  
**imports** ATerms<sup>2.5.1</sup>  
**exports**  
  **context-free syntax**

“empty” → AFun  
“seq” → AFun  
“opt” → AFun  
“iter” → AFun  
“iter-star” → AFun  
“iter-sep” → AFun  
“iter-star-sep” → AFun  
“iter-n” → AFun  
“iter-sep-n” → AFun  
“set” → AFun  
“pair” → AFun  
“func” → AFun  
“alt” → AFun  
“perm” → AFun

**module** Basic-Sdf-Tree-Constructors  
**imports** Kernel-Sdf-Tree-Constructors<sup>2.5.2</sup>  
**exports**

**context-free syntax**  
  “lexical-syntax” → AFun  
  “context-free-syntax” → AFun  
  “variables” → AFun  
  “lexical-variables” → AFun  
  “cf” → AFun  
  “lex” → AFun  
  “varsym” → AFun  
  “layout” → AFun

### A.2.2 Encoding and Decoding

**module** Grammar-ATerms  
**imports** Grammar-Tree-Constructors<sup>4.2.1</sup> Grammar-Syntax<sup>2.2</sup>  
**exports**  
  **context-free syntax**

aterm(Grammar) → ATerm  
 grammar(ATerm) → Grammar

**equations**

Encoding of grammars.

- [1]                    aterm( $\emptyset$ ) = empty-grammar
- [2]                    aterm( $\mathcal{G}_1 \ \mathcal{G}_2$ ) = conc-grammars(aterm( $\mathcal{G}_1$ ), aterm( $\mathcal{G}_2$ ))

Decoding of grammars.

- [3]                    grammar(empty-grammar) =  $\emptyset$
- [4]                    grammar(conc-grammars( $T_1, T_2$ )) = grammar( $T_1$ ) grammar( $T_2$ )

**module** CC-Sdf-ATerms

**imports** Kernel-Sdf-ATerms<sup>2.5.4</sup> CC-Sdf-Syntax<sup>2.4.2</sup>

CC-Sdf-Tree-Constructors<sup>A.2.1</sup> Character-Arithmetic

**exports**

**context-free syntax**

atermlist(OptCharRanges) → ATermList  
 aterm(Character)            → NatCon  
 ranges(ATermList)          → CharRanges  
 range(ATerm)               → CharRange  
 character(ATerm)           → Character

**equations**

Encoding character classes.

- [1]                    aterm( $[cr^*]$ ) = char-class(atermlist( $cr^*$ ))
- [2]                    aterm( $cr^*$ ) = [ ] **when**  $cr^* =$
- [3]                    aterm( $cr_1^+ \ cr_2^+$ ) = aterm( $cr_1^+$ ) ++ aterm( $cr_2^+$ )
- [4]                    aterm( $cr$ ) = [aterm( $c$ )] **when**  $cr = c$
- [5]                    aterm( $c_1 - c_2$ ) = [range(aterm( $c_1$ ), aterm( $c_2$ ))]
- [6]                    aterm( $c$ ) = int( $c$ )

Decoding character classes.

- [7]                    symbol(char-class([ ])) = [ ]
- [8]                    symbol(char-class( $[Ts]$ )) = [ranges( $[Ts]$ )]
- [9]                    ranges( $[T]$ ) = range( $T$ )
- [10]                   ranges( $[T, Ts]$ ) = range( $T$ ) ranges( $[Ts]$ )
- [11]                   range( $n$ ) = character( $n$ )
- [12]                   range(range( $n_1, n_2$ )) = character( $n_1$ ) - character( $n_2$ )
- [13]                   character( $n$ ) = char( $n$ )
- [14]                   symbol( $n$ ) = [char( $n$ )]

**module** Sorts-Sdf-ATerms

**imports** Kernel-Sdf-ATerms<sup>2.5.4</sup> Sorts-Sdf-Tree-Constructors<sup>A.2.1</sup>

Sorts-Sdf-Syntax<sup>2.4.1</sup>

**equations**

Encoding and decoding sorts.

- [1]                    aterm(sort( $c^+$ )) = sort(literal("""  $c^+$  """))

$$[2] \quad \text{symbol}(\text{sort}(\text{literal}(\text{"" } c^+ \text{ ""}))) = \text{sort}(c^+)$$

**module** Literals-Sdf-ATerms

**imports** Kernel-Sdf-ATerms<sup>2.5.4</sup> Literals-Sdf-Tree-Constructors<sup>A.2.1</sup>  
 Literals-Sdf-Syntax<sup>2.4.3</sup>

**equations**

Encoding and decoding literals.

$$[1] \quad \text{aterm}(L) = \text{lit}(L)$$

$$[2] \quad \text{symbol}(\text{lit}(L)) = L$$

**module** Priority-Sdf-ATerms

**imports** Kernel-Sdf-ATerms<sup>2.5.4</sup> Priority-Sdf-Syntax<sup>3.1.1</sup>

**equations**

Encoding attributes.

$$[1] \quad \text{aterm}(\text{left}) = \text{atr}(\text{"left"})$$

$$[2] \quad \text{aterm}(\text{right}) = \text{atr}(\text{"right"})$$

$$[3] \quad \text{aterm}(\text{bracket}) = \text{atr}(\text{"bracket"})$$

$$[4] \quad \text{aterm}(\text{assoc}) = \text{atr}(\text{"assoc"})$$

$$[5] \quad \text{aterm}(\text{non-assoc}) = \text{atr}(\text{"non-assoc"})$$

Decoding attributes.

$$[6] \quad \text{attribute}(\text{atr}(\text{"left"})) = \text{left}$$

$$[7] \quad \text{attribute}(\text{atr}(\text{"right"})) = \text{right}$$

$$[8] \quad \text{attribute}(\text{atr}(\text{"bracket"})) = \text{bracket}$$

$$[9] \quad \text{attribute}(\text{atr}(\text{"assoc"})) = \text{assoc}$$

$$[10] \quad \text{attribute}(\text{atr}(\text{"non-assoc"})) = \text{non-assoc}$$

**module** Regular-Sdf-ATerms

**imports** Regular-Sdf-Tree-Constructors<sup>A.2.1</sup> Kernel-Sdf-ATerms<sup>2.5.4</sup>  
 Regular-Sdf-Syntax<sup>3.2.1</sup>

**equations**

Encoding regular expressions.

$$[1] \quad \text{aterm}(( )) = \text{empty}$$

$$[2] \quad \text{aterm}(\mathcal{A} \alpha^+) = \text{seq}(\text{atermlist}(\mathcal{A} \alpha^+))$$

$$[3] \quad \text{aterm}(\mathcal{A}?) = \text{opt}(\text{aterm}(\mathcal{A}))$$

$$[4] \quad \text{aterm}(\mathcal{A}+) = \text{iter}(\text{aterm}(\mathcal{A}))$$

$$[5] \quad \text{aterm}(\mathcal{A}*) = \text{iter-star}(\text{aterm}(\mathcal{A}))$$

$$[6] \quad \text{aterm}(\{\mathcal{A} \mathcal{B}\}+) = \text{iter-sep}(\text{aterm}(\mathcal{A}), \text{aterm}(\mathcal{B}))$$

$$[7] \quad \text{aterm}(\{\mathcal{A} \mathcal{B}\}*) = \text{iter-star-sep}(\text{aterm}(\mathcal{A}), \text{aterm}(\mathcal{B}))$$

$$[8] \quad \text{aterm}(\{\mathcal{A}\} n^+) = \text{iter-n}(\text{aterm}(\mathcal{A}), \text{con}(n))$$

$$[9] \quad \text{aterm}(\{\mathcal{A} \mathcal{B}\} n^+) = \text{iter-sep-n}(\text{aterm}(\mathcal{A}), \text{aterm}(\mathcal{B}), \text{con}(n))$$

$$[10] \quad \text{aterm}(\text{Set}[\mathcal{A}]) = \text{set}(\text{aterm}(\mathcal{A}))$$

$$[11] \quad \text{aterm}(\mathcal{A} \# \mathcal{B}) = \text{pair}(\text{aterm}(\mathcal{A}), \text{aterm}(\mathcal{B}))$$

- [12]  $\text{aterm}((\alpha \Rightarrow \mathcal{B})) = \text{func}(\text{atermlist}(\alpha), \text{aterm}(\mathcal{B}))$
- [13]  $\text{aterm}(\mathcal{A} \mid \mathcal{B}) = \text{alt}(\text{aterm}(\mathcal{A}), \text{aterm}(\mathcal{B}))$
- [14]  $\text{aterm}(\ll \alpha \gg) = \text{perm}(\text{atermlist}(\alpha))$

Decoding regular expressions.

- [15]  $\text{symbol}(\text{empty}) = ()$
- [16]  $\text{symbol}(\text{seq}(Tl)) = (\mathcal{A} \alpha^+ \text{ when symbols}(Tl) = \mathcal{A} \alpha^+$
- [17]  $\text{symbol}(\text{opt}(T_1)) = \text{symbol}(T_1)?$
- [18]  $\text{symbol}(\text{iter}(T_1)) = \text{symbol}(T_1)^+$
- [19]  $\text{symbol}(\text{iter-star}(T_1)) = \text{symbol}(T_1)^*$
- [20]  $\text{symbol}(\text{iter-sep}(T_1, T_2)) = \{\text{symbol}(T_1) \text{ symbol}(T_2)\}^+$
- [21]  $\text{symbol}(\text{iter-star-sep}(T_1, T_2)) = \{\text{symbol}(T_1) \text{ symbol}(T_2)\}^*$
- [22]  $\text{symbol}(\text{iter-n}(T_1, n)) = \{\text{symbol}(T_1)\}^n$
- [23]  $\text{symbol}(\text{iter-sep-n}(T_1, T_2, n)) = \{\text{symbol}(T_1) \text{ symbol}(T_2)\}^n$
- [24]  $\text{symbol}(\text{set}(T_1)) = \text{Set}[\text{symbol}(T_1)]$
- [25]  $\text{symbol}(\text{pair}(T_1, T_2)) = \text{symbol}(T_1) \# \text{symbol}(T_2)$
- [26]  $\text{symbol}(\text{func}(Tl_1, Tl_2)) = (\text{symbols}(Tl_1) \Rightarrow \text{symbol}(Tl_2))$
- [27]  $\text{symbol}(\text{alt}(T_1, T_2)) = \text{symbol}(T_1) \mid \text{symbol}(T_2)$
- [28]  $\text{symbol}(\text{perm}(Tl)) = \ll \text{symbols}(Tl) \gg$

**module** Basic-Sdf-ATerms

**imports** Basic-Sdf-Tree-Constructors<sup>4.2.1</sup> Basic-Sdf-Syntax<sup>3.3.1</sup>

Kernel-Sdf-ATerms<sup>2.5.4</sup>

**equations**

Encoding grammars.

- [1]  $\text{aterm}(\text{lexical syntax } p^*) = \text{lexical-syntax}(\text{atermlist}(p^*))$
- [2]  $\text{aterm}(\text{context-free syntax } p^*) = \text{context-free-syntax}(\text{atermlist}(p^*))$
- [3]  $\text{aterm}(\text{variables } p^*) = \text{variables}(\text{atermlist}(p^*))$
- [4]  $\text{aterm}(\text{lexical variables } p^*) = \text{lexical-variables}(\text{atermlist}(p^*))$

Encoding symbols.

- [5]  $\text{aterm}(\langle \mathcal{A}\text{-LEX} \rangle) = \text{lex}(\text{aterm}(\mathcal{A}))$
- [6]  $\text{aterm}(\langle \mathcal{A}\text{-CF} \rangle) = \text{cf}(\text{aterm}(\mathcal{A}))$
- [7]  $\text{aterm}(\langle \mathcal{A}\text{-VAR} \rangle) = \text{varsym}(\text{aterm}(\mathcal{A}))$
- [8]  $\text{aterm}(\text{LAYOUT}) = \text{layout}$

Decoding grammars.

- [9]  $\text{grammar}(\text{lexical-syntax}(Tl)) = \text{lexical syntax productions}(Tl)$
- [10]  $\text{grammar}(\text{context-free-syntax}(Tl)) = \text{context-free syntax productions}(Tl)$
- [11]  $\text{grammar}(\text{variables}(Tl)) = \text{variables productions}(Tl)$
- [12]  $\text{grammar}(\text{lexical-variables}(Tl)) = \text{lexical variables productions}(Tl)$

Decoding symbols.

- [13]  $\text{symbol}(\text{lex}(T)) = \langle \text{symbol}(T)\text{-LEX} \rangle$
- [14]  $\text{symbol}(\text{cf}(T)) = \langle \text{symbol}(T)\text{-CF} \rangle$
- [15]  $\text{symbol}(\text{varsym}(T)) = \langle \text{symbol}(T)\text{-VAR} \rangle$
- [16]  $\text{symbol}(\text{layout}) = \text{LAYOUT}$

### A.3 Renamings

**module** CC-Sdf-Renaming  
**imports** Kernel-Sdf-Renaming<sup>4.1.3</sup> CC-Sdf-Syntax<sup>2.4.2</sup>  
**equations**  
 Renaming character classes.

$$[1] \quad [cc] \rho = cc$$

**module** Literals-Sdf-Renaming  
**imports** Kernel-Sdf-Renaming<sup>4.1.3</sup> Literals-Sdf-Syntax<sup>2.4.3</sup>  
**hiddens**  
**variables**  
 $"L" \rightarrow \text{Literal}$   
**equations**  
 Renaming literals.

$$[1] \quad [L] \rho = L$$

**module** Sorts-Sdf-Renaming  
**imports** Sorts-Sdf-Syntax<sup>2.4.1</sup> Kernel-Sdf-Renaming<sup>4.1.3</sup>  
**equations**  
 Renaming sorts.

$$[1] \quad [S] \rho = \mathcal{S}$$

$$[2] \quad (\text{sorts } \alpha) \rho = \text{sorts } (\alpha) * \rho$$

**module** Priority-Sdf-Renaming  
**imports** Priority-Sdf-Projection<sup>3.1.2</sup> Basic-Sdf-Renaming<sup>4.3</sup>  
**exports**  
**context-free syntax**  
 $"(" \text{ Priorities } ")" \text{ Renamings } \rightarrow \text{Priorities}$   
 $"(" \text{ Group } ")"_G \text{ Renamings } \rightarrow \text{Group}$   
**equations**  
 Renaming symbols and productions in priority declarations.

$$[1] \quad (\text{priorities } pr^*) \rho = \text{priorities } (pr^*) \rho$$

Lists of priorities.

$$[2] \quad (pr^*) \rho = \text{when } pr^* =$$

$$[3] \quad (pr_1^+, pr_2^+) \rho = (pr_1^+) \rho \mathrel{++} (pr_2^+) \rho$$

Associativity and priority declarations.

$$[4] \quad (g_1 \text{ as } g_2) \rho = (g_1)_G \rho \text{ as } (g_2)_G \rho$$

$$[5] \quad (g_1 > g_2) \rho = (g_1)_G \rho > (g_2)_G \rho$$

$$[6] \quad (g_1 > g_2 > gg^+) \rho = (g_1)_G \rho > g_2' > gg'^+ \text{ when } g_2' > gg'^+ = (g_2 > gg^+) \rho$$

Groups.

$$\begin{aligned}
[7] \quad & (g)_G \rho = p' \quad \text{when } g = p, \ p' = (p) \rho \\
[8] \quad & (\{p^*\})_G \rho = \{(p^*)^* \rho\} \\
[9] \quad & (\{as : p^*\})_G \rho = \{as : (p^*)^* \rho\}
\end{aligned}$$

**module** Regular-Sdf-Renaming

**imports** Kernel-Sdf-Renaming<sup>4.1.3</sup> Regular-Sdf-Syntax<sup>3.2.1</sup>

**equations**

Renaming symbols in regular expressions.

$$\begin{aligned}
[1] \quad & [()] \rho = () \\
[2] \quad & [(\mathcal{A} \alpha^+)] \rho = (\mathcal{B} \beta^+) \quad \text{when } (\mathcal{A} \alpha^+)^* \rho = \mathcal{B} \beta^+ \\
[3] \quad & [\mathcal{A}^?] \rho = (\mathcal{A}) \rho? \\
[4] \quad & [\mathcal{A}^+] \rho = (\mathcal{A}) \rho^+ \\
[5] \quad & [\mathcal{A}^*] \rho = (\mathcal{A}) \rho^* \\
[6] \quad & [(\mathcal{A} \mathcal{B})^+] \rho = \{(\mathcal{A}) \rho (\mathcal{B}) \rho\}^+ \\
[7] \quad & [(\mathcal{A} \mathcal{B})^*] \rho = \{(\mathcal{A}) \rho (\mathcal{B}) \rho\}^* \\
[8] \quad & [\mathcal{A}^n] \rho = (\mathcal{A}) \rho^n \\
[9] \quad & [(\mathcal{A} \mathcal{B})^n] \rho = \{(\mathcal{A}) \rho (\mathcal{B}) \rho\}^n \\
[10] \quad & [\text{Set}[\mathcal{A}]] \rho = \text{Set}[(\mathcal{A}) \rho] \\
[11] \quad & [\mathcal{A} \# \mathcal{B}] \rho = (\mathcal{A}) \rho \# (\mathcal{B}) \rho \\
[12] \quad & [(\alpha \Rightarrow \mathcal{B})] \rho = ((\alpha)^* \rho \Rightarrow (\mathcal{B}) \rho) \\
[13] \quad & [\mathcal{A} \mid \mathcal{B}] \rho = (\mathcal{A}) \rho \mid (\mathcal{B}) \rho
\end{aligned}$$

**module** Basic-Sdf-Renaming

**imports** Basic-Sdf-Normalization<sup>3.3.2</sup> Kernel-Sdf-Renaming<sup>4.1.3</sup>

**exports**

**context-free syntax**

“<” Renamings “-LEX” “>” → Renamings

“<” Renamings “-CF” “>” → Renamings

**equations**

Renaming grammars.

$$\begin{aligned}
[1] \quad & (\text{context-free syntax } p^*) \rho = \text{context-free syntax } (p^*)^* \rho \\
[2] \quad & (\text{lexical syntax } p^*) \rho = \text{lexical syntax } (p^*)^* \rho \\
[3] \quad & (\text{variables } p^*) \rho = \text{variables } (p^*)^* \rho \\
[4] \quad & (\text{lexical variables } p^*) \rho = \text{lexical variables } (p^*)^* \rho
\end{aligned}$$

Renaming symbols.

$$\begin{aligned}
[5] \quad & [\langle \mathcal{A} \text{-LEX} \rangle] \rho = \langle (\mathcal{A}) \rho \text{-LEX} \rangle \\
[6] \quad & [\langle \mathcal{A} \text{-CF} \rangle] \rho = \langle (\mathcal{A}) \rho \text{-CF} \rangle \\
[7] \quad & [\langle \mathcal{A} \text{-VAR} \rangle] \rho = \langle (\mathcal{A}) \rho \text{-VAR} \rangle \\
[8] \quad & [\text{LAYOUT}] \rho = \text{LAYOUT}
\end{aligned}$$

Applying  $\langle \_ \text{-LEX} \rangle$  to a renaming.

$$[9] \quad \langle [\_ \text{-LEX}] \rangle = []$$

$$\begin{aligned}
[10] \quad & \langle [A \Rightarrow B] \text{-LEX} \rangle = [\langle A \text{-LEX} \rangle \Rightarrow \langle B \text{-LEX} \rangle] \\
[11] \quad & \langle [p_1 \Rightarrow p_2] \text{-LEX} \rangle = [\langle p_1 \text{-LEX} \rangle \Rightarrow \langle p_2 \text{-LEX} \rangle] \\
[12] \quad & \langle [\rho_1^+ \ \rho_2^+] \text{-LEX} \rangle = \langle [\rho_1^+] \text{-LEX} \rangle ++ \langle [\rho_2^+] \text{-LEX} \rangle
\end{aligned}$$

Applying  $\langle \_ \text{-CF} \rangle$  to a renaming.

$$\begin{aligned}
[13] \quad & \langle [] \text{-CF} \rangle = [] \\
[14] \quad & \langle [A \Rightarrow B] \text{-CF} \rangle = [\langle A \text{-CF} \rangle \Rightarrow \langle B \text{-CF} \rangle] \\
[15] \quad & \langle [p_1 \Rightarrow p_2] \text{-CF} \rangle = [\langle p_1 \text{-CF} \rangle \Rightarrow \langle p_2 \text{-CF} \rangle] \\
[16] \quad & \langle [\rho_1^+ \ \rho_2^+] \text{-CF} \rangle = \langle [\rho_1^+] \text{-CF} \rangle ++ \langle [\rho_2^+] \text{-CF} \rangle
\end{aligned}$$

**module** Restrictions-Sdf-Renaming  
**imports** Restrictions-Sdf-Syntax<sup>3.4.1</sup> Kernel-Sdf-Renaming<sup>4.1.3</sup>  
**exports**  
**context-free syntax**  
 “(” Restrictions “)” Renamings  $\rightarrow$  Restrictions  
**equations**  
 Renaming restrictions.

$$\begin{aligned}
[1] \quad & (\text{restrictions } restr^*) \ \rho = \text{restrictions } (restr^*) \ \rho \\
[2] \quad & (restr^*) \ \rho = \textbf{when } restr^* = \\
[3] \quad & (\alpha \not\vdash cc \ restr_1^*) \ \rho = (\alpha) * \rho \not\vdash cc \ restr_2^* \quad \textbf{when } restr_2^* = (restr_1^*) \ \rho
\end{aligned}$$

## A.4 SDF2

**module** Sdf2-Projection  
**imports** Kernel-Sdf-Projection<sup>2.3.2</sup> Sorts-Sdf-Projection<sup>2.4.1</sup>  
 Priority-Sdf-Projection<sup>3.1.2</sup> Renaming-Sdf-Projection<sup>4.1.2</sup>  
 Modular-Sdf-Projection<sup>4.3.2</sup> Alias-Sdf-Projection<sup>4.2.2</sup>  
 Restrictions-Sdf-Projection<sup>3.4.2</sup>  
**module** Sdf2-Renaming  
**imports** Sdf2-Syntax<sup>5.1.1</sup> Kernel-Sdf-Renaming<sup>4.1.3</sup> Priority-Sdf-Renaming<sup>4.3</sup>  
 Regular-Sdf-Renaming<sup>4.3</sup> Literals-Sdf-Renaming<sup>4.3</sup>  
 CC-Sdf-Renaming<sup>4.3</sup> Basic-Sdf-Renaming<sup>4.3</sup> Sorts-Sdf-Renaming<sup>4.3</sup>  
 Restrictions-Sdf-Renaming<sup>4.3</sup> Modular-Sdf-Renaming<sup>4.3.4</sup>  
 Alias-Sdf-Renaming<sup>4.2.3</sup>  
**equations**

$$\begin{aligned}
[1] \quad & (\text{context-free priorities } pr^*) \ \rho = \text{context-free priorities } (pr^*) \ \rho \\
[2] \quad & (\text{lexical priorities } pr^*) \ \rho = \text{lexical priorities } (pr^*) \ \rho
\end{aligned}$$

$$\begin{aligned}
[3] \quad & [\langle \text{Start} \rangle] \ \rho = \langle \text{Start} \rangle \\
[4] \quad & [\langle \text{START} \rangle] \ \rho = \langle \text{START} \rangle
\end{aligned}$$



```

module Sdf2-Tree-Constructors
imports Kernel-Sdf-Tree-Constructors2.5.2 Basic-Sdf-Tree-ConstructorsA.2.1
        Modular-Sdf-Tree-Constructors Regular-Sdf-Tree-ConstructorsA.2.1
        Priority-Sdf-Tree-Constructors CC-Sdf-Tree-ConstructorsA.2.1
        Sorts-Sdf-Tree-ConstructorsA.2.1 Literals-Sdf-Tree-ConstructorsA.2.1

module Sdf2-ATerms
imports Sdf2-Tree-ConstructorsA.4 Sdf2-Syntax5.1.1 Kernel-Sdf-ATerms2.5.4
        Basic-Sdf-ATermsA.2.2 Modular-Sdf-ATerms Regular-Sdf-ATermsA.2.2
        Priority-Sdf-ATermsA.2.2 CC-Sdf-ATermsA.2.2 Sorts-Sdf-ATermsA.2.2
        Literals-Sdf-ATermsA.2.2 Restrictions-Sdf-ATerms
equations

```

```

[1]          aterm(⟨START⟩) = sort("<START>")
[2]          aterm(⟨Start⟩) = sort("<Start>")

```

```

module Sdf2-Trees
imports Sdf2-ATermsA.4 Sdf2-Syntax5.1.1 Kernel-Sdf-Trees2.5.5
        Priority-Sdf-Trees3.1.4 CC-Sdf-Trees2.5.6 Renaming-Sdf-Trees4.1.4

module Sdf2-Equality
imports Kernel-Sdf-Equality2.5.8 Regular-Sdf-Equality Basic-Sdf-Equality3.3.3

```



---

# B

## Bibliography

---

Technical reports from the Programming Research Group of the University of Amsterdam can be obtained from <http://www.wins.uva.nl/research/prog/reports/reports.html>. Technical reports from CWI can be found at <http://www.cwi.nl/cwi/publications/reports/reports.html>.

---

- Aho, A. V., Johnson, S. C., and Ullman, J. D. (1975). Deterministic parsing of ambiguous grammars. *Communications of the ACM*, **18**(8), 441–452.
- Backus, J. W. (1959). The syntax and semantics of the proposed international algebraic language of the Zurich ACM-GAMM conference. In *Proceedings of the International Conference on Information Processing*, pages 125–132, Paris. UNESCO.
- Bahlke, R. and Snelting, G. (1986). The PSG system: from formal language definitions to interactive programming environments. *ACM Transactions on Programming Languages and Systems*, **8**(4), 547–576.
- Bergstra, J. A., Heering, J., and Klint, P., editors (1989a). *Algebraic Specification*. ACM Press Frontier Series. The ACM Press in co-operation with Addison-Wesley.
- Bergstra, J. A., Heering, J., and Klint, P. (1989b). The algebraic specification formalism ASF. In J. A. Bergstra, J. Heering, and P. Klint, editors, *Algebraic Specification*, ACM Press Frontier Series, pages 1–66. The ACM Press in co-operation with Addison-Wesley. Chapter 1.
- Bergstra, J. A., Heering, J., and Klint, P. (1990). Module algebra. *Journal of the ACM*, **37**(2), 335–372.
- Billot, S. and Lang, B. (1989). The structure of shared forests in ambiguous parsing. In *Proceedings of the Twenty-Seventh Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Van den Brand, M. G. J. and Visser, E. (1996). Generation of formatters for context-free languages. *ACM Transactions on Software Engineering and Methodology*, **5**(1), 1–41.
- Van den Brand, M. G. J., Klint, P., Olivier, P., and Visser, E. (1997). ATerms: Representing structured data for exchange between heterogeneous tools. Technical report, Programming Research Group, University of Amsterdam.

- Cameron, R. D. (1993). Extending context-free grammars with permutation phrases. *ACM Letters on Programming Languages and Systems*, **2**(1-4), 85–94.
- Cardelli, L., Matthes, F., and Abadi, M. (1994). Extensible syntax with lexical scoping. SRC Research Report 121, Digital Systems Research Center, Palo Alto, California.
- Chomsky, N. (1956). Three models for the description of language. *IRE Transactions on Information Theory*, **2**, 113–124.
- Cordy, J. R. and Carmichael, I. H. (1993). *The TXL Programming Language. Syntax and Informal Semantics. Version 7*. Software Technology Laboratory, Department of Computer and Information Science, Queen’s University at Kingston, Kingston, Canada, 7 edition.
- Van Deursen, A., Heering, J., and Klint, P., editors (1996). *Language Prototyping. An Algebraic Specification Approach*, volume 5 of *AMAST Series in Computing*. World Scientific, Singapore.
- Dinesh, T. B. (1995). Injection misdemeanors. In M. G. J. v. d. Brand *et al.*, editors, *ASF+SDF’95. A Workshop on Generating Tools from Algebraic Specifications*, pages 255–270. Technical Report P9504, Programming Research Group, University of Amsterdam.
- Earley, J. (1975). Ambiguity and precedence in syntax description. *Acta Informatica*, **4**(1), 183–192.
- Futatsugi, K., Goguen, J., Jouannaud, J.-P., and Meseguer, J. (1985). Principles of OBJ2. In B. Reid, editor, *Conference Record of the Twelfth Annual ACM Symposium on Principles of Programming Languages*, pages 52–66. ACM.
- Goguen, J. A., Thatcher, J. W., Wagner, E. G., and Wright, J. B. (1977). Initial algebra semantics and continuous algebras. *Journal of the ACM*, **24**(1), 68–95.
- Gray, R. W., Heuring, V. P., Levi, S. P., Sloane, A. M., and Waite, W. M. (1992). Eli: A complete, flexible compiler construction system. *Communications of the ACM*, **35**, 121–131.
- Grosch, J. (1990). Lalr - a generator for efficient parsers. *Software-Practice & Experience*, **20**, 1115–1135.
- Hatcher, W. S. and Rus, T. (1976). Context-free algebras. *Journal of Cybernetics*, **6**, 65–76.
- Heering, J., Hendriks, P. R. H., Klint, P., and Rekers, J. (1989). The syntax definition formalism SDF – reference manual. *SIGPLAN Notices*, **24**(11), 43–75.
- Hendriks, P. R. H. (1991). *Implementation of Modular Algebraic Specifications*. Ph.D. thesis, University of Amsterdam.
- Johnson, S. C. (1975). YACC—yet another compiler-compiler. Technical Report CS-32, AT & T Bell Laboratories, Murray Hill, N.J.

- Kahn, G., Lang, B., Mèlèse, B., and Morcos, E. (1983). METAL: A formalism to specify formalisms. *Science of Computer Programming*, **3**, 151–188.
- Klint, P. (1993). A meta-environment for generating programming environments. *ACM Transactions on Software Engineering and Methodology*, **2**(2), 176–201.
- Klint, P. and Visser, E. (1994). Using filters for the disambiguation of context-free grammars. In G. Pighizzini and P. San Pietro, editors, *Proc. ASMICS Workshop on Parsing Theory*, pages 1–20, Milano, Italy. Tech. Rep. 126–1994, Dipartimento di Scienze dell’Informazione, Università di Milano.
- Knuth, D. E. (1968). Semantics of context-free languages. *Mathematical Systems Theory*, **2**(2), 127–145. Correction in: *Mathematical Systems Theory* 5(1), pp. 95–96, Springer-Verlag, 1971.
- Koster, C. H. A. (1971). Affix grammars. In J. E. L. Peck, editor, *Algol-68 Implementation*. North-Holland, Amsterdam.
- Lee, J. A. N. (1972). The formal definition of the BASIC language. *Computer Journal*, **15**(1), 37–41.
- Lesk, M. E. and Schmidt, E. (1986). *LEX — A lexical analyzer generator*. Bell Laboratories. UNIX Programmer’s Supplementary Documents, Volume 1 (PS1).
- Luttik, B. and Visser, E. (1997). Specification of rewriting strategies. In A. Sellink, editor, *Second International Conference on the Theory and Practice of Algebraic Specification (ASF+SDF’97)*, Amsterdam, The Netherlands. Programming Research Group, University of Amsterdam.
- Meinke, K. (1992). Equational specification of abstract types and combinators. In E. Boerger, G. Jaeger, H. K. Buening, and M. M. Richter, editors, *Computer Science Logic - CSL’91*, volume 626 of *Lecture Notes in Computer Science*, pages 257–271, Berlin. Springer-Verlag.
- Mosses, P. D. (1992). *Action Semantics*. Cambridge University Press.
- Naur, P. et al. (1960). Report on the algorithmic language ALGOL 60. *Communications of the ACM*, **3**(5), 299–314.
- Parr, T. J. and Quong, R. W. (1994). Adding semantic and syntactic predicates to LL(k): pred-LL(k). In P. A. Fritzson, editor, *Compiler Construction, 5th International Conference, CC’94*, volume 786 of *LNCS*, pages 263–277, Edinburgh, U.K. Springer-Verlag.
- Pereira, F. C. N. and Warren, D. H. D. (1980). Definite Clause Grammars for language analysis—a survey of the formalism and a comparison with augmented transition networks. *Artificial Intelligence*, **13**, 231–278.
- Rekers, J. (1992). *Parser Generation for Interactive Environments*. Ph.D. thesis, University of Amsterdam. <ftp://ftp.cwi.nl/pub/gipe/reports/Rek92.ps.Z>.
- Reps, T. and Teitelbaum, T. (1989). *The Synthesizer Generator: a System for Constructing Language-Based Editors*. Springer-Verlag.

- Salomon, D. J. and Cormack, G. V. (1989). Scannerless NSLR(1) parsing of programming languages. *SIGPLAN Notices*, **24**(7), 170–178.
- Salomon, D. J. and Cormack, G. V. (1995). The disambiguation and scannerless parsing of complete character-level grammars for programming languages. Technical Report 95/06, Department of Computer Science, University of Manitoba, Winnipeg, Canada.
- Thorup, M. (1994). Controlled grammatic ambiguity. *ACM Transactions on Programming Languages and Systems*, **16**(3), 1024–1050.
- Visser, E. (1995). A case study in optimizing parsing schemata by disambiguation filters. In S. Fischer and M. Trautwein, editors, *Proceedings Accolade95*, pages 153–167, Amsterdam. The Dutch Graduate School in Logic.
- Visser, E. (1996). Multi-level specifications. In A. van Deursen, J. Heering, and P. Klint, editors, *Language Prototyping. An Algebraic Specification Approach*, volume 5 of *AMAST Series in Computing*, pages 105–196. World Scientific, Singapore.
- Visser, E. (1997a). A case study in optimizing parsing schemata by disambiguation filters. In A. Nijholt, editor, *International Workshop on Parsing Technology IWPT'97*, Boston, USA. (To appear).
- Visser, E. (1997b). Character classes. Technical Report P9708, Programming Research Group, University of Amsterdam.
- Visser, E. (1997c). Executable specification of programming languages. (Draft).
- Visser, E. (1997d). Polymorphic syntax definition. *Theoretical Computer Science*. (To appear).
- Visser, E. (1997e). Scannerless generalized-LR parsing. Technical Report P9707, Programming Research Group, University of Amsterdam.
- Vitteck, M. (1994). *ELAN: Un cadre logique pour le prototypage de langage de programmation avec contraintes*. Ph.D. thesis, Université Henri Poincaré – Nancy I, Nancy, France.
- Voisin, F. (1986). CIGALE: a tool for interactive grammar construction and expression parsing. *Science of Computer Programming*, **7**, 61–86.
- Watt, D. A. (1977). The parsing problem for affix grammars. *Acta Informatica*, **8**(1), 1–20.
- Williams, M. H. (1982). A flexible notation for syntactic definitions. *ACM Transactions on Programming Languages and Systems*, **4**(1), 113–119.
- Wirth, N. (1977). What can we do about the unnecessary diversity of notation for syntactic definitions. *Communications of the ACM*, **20**(11), 822–823.

## Technical Reports of the Programming Research Group

**Note:** These reports can be obtained using the technical reports overview on our WWW site (URL <http://www.wins.uva.nl/research/prog/reports/>) or using anonymous ftp to [ftp.wins.uva.nl](ftp://ftp.wins.uva.nl), directory `pub/programming-research/reports/`.

- [P9711] L. Moonen. *A Generic Architecture for Data Flow Analysis to Support Reverse Engineering*.
- [P9710] B. Luttik and E. Visser. *Specification of Rewriting Strategies*.
- [P9709] J.A. Bergstra and M.P.A. Sellink. *An Arithmetical Module for Rationals and Reals*.
- [P9708] E. Visser. *Character Classes*.
- [P9707] E. Visser. *Scannerless generalized-LR parsing*.
- [P9706] E. Visser. *A family of syntax definition formalisms*.
- [P9705] M.G.J. van den Brand, M.P.A. Sellink, and C. Verhoef. *Generation of Components for Software Renovation Factories from Context-free Grammars*.
- [P9704] P.A. Olivier. *Debugging Distributed Applications Using a Coordination Architecture*.
- [P9703] H.P. Korver and M.P.A. Sellink. *A Formal Axiomatization for Alphabet Reasoning with Parametrized Processes*.
- [P9702] M.G.J. van den Brand, M.P.A. Sellink, and C. Verhoef. *Reengineering COBOL Software Implies Specification of the Underlying Dialects*.
- [P9701] E. Visser. *Polymorphic Syntax Definition*.
- [P9618] M.G.J. van den Brand, P. Klint, and C. verhoef. *Re-engineering needs Generic Programming Language Technology*.
- [P9617] P.I. Manuel. *ANSI Cobol III in SDF + an ASF Definition of a Y2K Tool*.
- [P9616] P.H. Rodenburg. *A Complete System of Four-valued Logic*.
- [P9615] S.P. Luttik and P.H. Rodenburg. *Transformations of Reduction Systems*.
- [P9614] M.G.J. van den Brand, P. Klint, and C. Verhoef. *Core Technologies for System Renovation*.
- [P9613] L. Moonen. *Data Flow Analysis for Reverse Engineering*.
- [P9612] J.A. Hillebrand. *Transforming an ASF+SDF Specification into a Tool-Bus Application*.
- [P9611] M.P.A. Sellink. *On the conservativity of Leibniz Equality*.

- [P9610] T.B. Dinesh and S.M. Üsküdarlı. *Specifying input and output of visual languages.*
- [P9609] T.B. Dinesh and S.M. Üsküdarlı. *The VAS formalism in VASE.*
- [P9608] J.A. Hillebrand. *A small language for the specification of Grid Protocols.*
- [P9607] J.J. Brunekreef. *A transformation tool for pure Prolog programs: the algebraic specification.*
- [P9606] E. Visser. *Solving type equations in multi-level specifications (preliminary version).*
- [P9605] P.R. D’Argenio and C. Verhoef. *A general conservative extension theorem in process algebras with inequalities.*
- [P9602b] J.A. Bergstra and M.P.A. Sellink. *Sequential data algebra primitives (revised version of P9602).*
- [P9604] E. Visser. *Multi-level specifications.*
- [P9603] M.G.J. van den Brand, P. Klint, and C. Verhoef. *Reverse engineering and system renovation: an annotated bibliography.*
- [P9602] J.A. Bergstra and M.P.A. Sellink. *Sequential data algebra primitives.*
- [P9601] P.A. Olivier. *Embedded system simulation: testdriving the ToolBus.*
- [P9512] J.J. Brunekreef. *TransLog, an interactive tool for transformation of logic programs.*
- [P9511] J.A. Bergstra, J.A. Hillebrand, and A. Ponse. *Grid protocols based on synchronous communication: specification and correctness.*
- [P9510] P.H. Rodenburg. *Termination and confluence in infinitary term rewriting.*
- [P9509] J.A. Bergstra and Gh. Stefanescu. *Network algebra with demonic relation operators.*
- [P9508] J.A. Bergstra, C.A. Middelburg, and Gh. Stefanescu. *Network algebra for synchronous and asynchronous dataflow.*
- [P9507] E. Visser. *A case study in optimizing parsing schemata by disambiguation filters.*
- [P9506] M.G.J. van den Brand and E. Visser. *Generation of formatters for context-free languages.*
- [P9505] J.M.T. Romijn. *Automatic analysis of term rewriting systems: proving properties of term rewriting systems derived from ASF+SDF specifications.*