

# A Gentle Introduction to *Compiler Generation* Using HACS\*

Kristoffer H. Rose

August 18, 2014

## Abstract

Higher-order Attribute Contraction Schemes—or HACS—is a language for programming compilers. With HACS it is possible to create a fully functional compiler from a single source file. This document explains how to get HACS up and running, and walk through the code of a simple example with each of the main stages of a compiler in HACS: lexical analysis, syntax analysis, semantic analysis, intermediate code generation and optimization, and code generation.

**Contents:** 1. Introduction (1), 2. Lexical Analysis (2), 3. Syntax Analysis (4), 4. Semantic Sorts and Schemes (6), 5. Semantic Analysis I (8), 6. Semantic Analysis II (10), 7. Intermediate and Code Generation (13), 8. Main (13), A. Tables (17), B. Bonus Examples (18), C. Common Errors (20).

## 1 Introduction

HACS abbreviates *Higher-order Attribute Contraction Schemes*, which is a formal system for symbolic rewriting extended with programming idioms commonly used in compiler rewriting.

—or HACS—is a language for programming compilers.

In this section we make sure you have a functional HACS installation on your computer.

**1.1 Requirements.** To run the HACS examples here you need a \*nix system (including a shell and the usual utilities) with these common programs: a Java development environment (at least Java 1.6 SE SDK, with `java` and `javac` commands); a C99 compiler with the *icu* libraries; the *flex* lexical generator; the `wget` internet retrieval and `unzip` unpacking utility; and the GNU `make` command. In addition, the setup process needs internet access to retrieve the JavaCC parser generator [3].

**1.2 Commands.** Retrieve the *hacs.zip* archive, extract it to a new directory, and initialize the Makefile, for example with the following commands:<sup>1</sup>

```
$ mkdir myfirst
$ cd myfirst
$ wget http://krisrose.net/hacs.zip
$ unzip hacs.zip
$ cp Makefile-hacs-template Makefile
```

---

\*UNFINISHED DRAFT—Feedback Appreciated.

<sup>1</sup>User input is blue.

Now check that all the commands defined in the SETUP section near the top of you freshly copied Makefile are correct; if any need changing then uncomment the line and change it to reflect the required commands on your system. Then check that it works with these command, which will need Internet access the first time:<sup>2</sup>

```
$ make first.run
...
$ ./first.run --action=Compile \
               --term="{initial := 1; rate := 1.0; position := initial + rate * 60;}"
LDF T, #1
STF name, T_51
LDF T_84, #1.0
STF name_40, T_56
LDF T_98, name_43
LDF T_70, name_3
LDF T_90, #60
MULF T_96, T_62, T_86
ADDF T_50, T_82, T_177
STF name_23, T_66
```

Congratulations—you just built your first compiler! (It may have taken a few minutes, as HACS also performed the initial setup, which is only needed on the first run.)

### 1.3 Example (module wrapper).

**1.4 Manual** (grammar structure). A HACS compiler is specified as a single *.hx* module file with the following structure:

```
module modulename
{
    // Lexical Analysis
    // Syntax Analysis
    // Semantic Analysis
    // Code Generator
    // Main
}
```

the *modulename* should be a string with a Java style fully qualified class name, where the last component is capitalized, like "net.sf.crsx.samples.gentle.First" in the example above. The individual sections specify the compiler, and the possible contents is documented in the manual blocks throughout this document.

**1.5 Plan.** In the remainder of this document we introduce the most important features of the HACS language by explaining the corresponding parts of the *first.hx* specification used above, which is adapted from [1, Fig. 1.7], as well as several other minor examples. We explain lexical analysis in Sec. 2, syntax analysis in Sec. 3, basic semantic sorts, schemes, and rules in Sec. 4, semantic analysis in Sec. 5 and Sec. 6, Code Generation in Sec. 7, and the invocation conventions in Sec. 8. Appendices include tables of Unicode characters and current limitations (A), some bonus examples (B) and explanations for a few cryptic error messages (C).

## 2 Lexical Analysis

Lexical analysis is the process of splitting the input text into tokens. HACS uses a usual variation of *regular expressions* for this. Unique tokens such as keywords can be declared along with the syntax productions where they are used.

---

<sup>2</sup>Note that the spacing of the output may vary—this is not a stable capability yet.

**2.1 Example** (tokens and white space). The first part of *first.hx* is the lexical analyzer. It uses regular expression notation to define the tokens used as “terminal symbols” by the compiler. Here is the relevant fragment for setting up the concrete syntax of integers, basic floating point numbers, identifiers, and white space:

```

1  /* 1. LEXICAL ANALYSIS. */
2  space [ \t\n] ;
3  token Int      | <Digit>+ ;
4  token Float    | <Digit>* "." <Digit>+ ;
5  token Id       | <Lower>+ ('_'? <Int>)? ;
6  token fragment Digit | [0-9] ;
7  token fragment Lower | [a-z] ;

```

The example illustrates the following particulars of HACS lexical expressions:

- Declarations generally start with a keyword or two and are terminated by a ; (semicolon).
- token declarations in particular have the token name followed by a regular expression between a | (vertical bar) and a ; (semicolon).
- A regular expressions is a sequence of units, corresponding to the concatenation of sequences of characters that match the units. Each unit can be a *character class* such as [a-z], which matches a single character in the indicated range, a string such as ' . ', or a reference to a token or fragment such as <Lower>.<sup>3</sup>
- A token fragment declaration means that the defined token can only be used in other token declarations, and not in syntax productions defined below.
- Every regular expression component can be followed by a repetition marker ?, +, or \*.
- The regular expression for white space is setup by space followed by the regular expression of what to skip – here spaces, tabs, and newlines, where HACS uses backslash for escaping in character classes with usual C-style language escapes.
- HACS supports C/Java-style comments.

**2.2 Commands** (lexical analysis). The generated command, *first.run*, from above, can be used as a lexical analyzer with two arguments: a *token sort* and a *token term*.<sup>4</sup> Execution then proceeds by parsing the string following the syntax of the token. We can, for example, check the lexical analysis of a number:

```

$ ./first.run --sort=Float --term=34.56
34.56

```

If there is an error, the lexical analyzer will inform us of this:

```

$ ./first.run --sort=Int --term=34.56
Exception in thread "main" java.lang.RuntimeException: net.sf.crsx.CRSEException:
net.sf.crsx.parser.ParseException:

```

<sup>3</sup>The used Unicode characters are summarized in Appendix A.1.

<sup>4</sup>The command has more options that we shall introduce as we need them.

```

Parse error in embedded firstInt term at line 1, column 27: [ 34.56 ] ...
Encountered " <T_T_M_FLOAT> "34.56 "" at line 1, column 29.
Was expecting one of:
    <T_meta1_T_M_INT> ...
    "\u27e9" ...
    <T_T_M_INT> ...

```

(where the trail of Java exceptions has been truncated: the important information is in the first few lines).

**2.3 Manual** (lexical declarations). A token is declared with the keyword `token` followed by the token (sort) name, a `|` (vertical bar), and a *regular expression*, which has one of the following forms (with increasing order of precedence):

1. Several alternative regular expressions can be combined with further `|` characters.
2. Concatenation denotes the regular expression recognizing concatenations of what matches the subexpressions.
3. A regular expression (of the forms following this one) can be followed by a *repetition marker*: `?` for zero or one, `+` for one or more, and `*` for zero or more.
4. A simple word without special characters stands for itself.
5. A string in single or double quotes stands for the contents of the string except that `\` introduces an *escape code* that stands for the encoded character in the string.
6. A stand-alone `\` followed by an *escape code* stands for that character: escape codes include the usual C and Java escapes: `\n`, `\r`, `\a`, `\f`, `\t`, octal escapes like `\177`, special character escapes like `\\`, `\'`, `\"`, and Unicode hexadecimal escapes like `\u27e9`.
7. A *character class* is given in `[]`, with these rules:
  - (a) if the first character is `~` then the character class is negated;
  - (b) if the first (after `~`) character is `]` then that character is (not) permitted;
  - (c) if a `\` followed by an *escape code* is encountered then it stands for the encoded character;
  - (d) if two characters are connected with a `-` (dash) then all characters in the indicated *range* (inclusive) are permitted (or excluded).

Note that a character class cannot be empty.

8. The `.` (period) character stands for the character class `[^\n]`.
9. A nested regular expression can be given in `()`.
10. An entire other token `T` can be included (by literal substitution, so recursion is not allowed) by writing `<T>` (the angle brackets are unicode characters U+27E8 and U+27E9). As a special convenience, tokens declared with `token fragment` can *only* be used this way.
11. The special declaration `space` defines what constitutes white space for the generated grammar. (Note that this does not influence what is considered space in the specification itself.) A spacing declaration permits the special alternative `nested` declaration for nested comments, the following defines usual C/Java style spacing with comments, for example:

```
space [ \t\f\r\n] | nested "/*" "*/" | "//" .* ;
```

Notice that spacing is not significant in regular expressions, except (1) in character classes, (2) in literal strings, (3) if escaped (as in `\`).

### 3 Syntax Analysis

Once we have tokens, we can use HACS to program a complete syntax analysis with a grammar that specifies how the input text is decomposed according to a *concrete syntax* and how the desired *abstract syntax tree* (AST) is constructed from that. Notice that HACS does not provide a “parse tree” in the traditional sense, *i.e.*, a tree that represents the full concrete syntax parse: only the AST is built. Grammars are structured following the *sorts* of AST nodes, with concrete syntax details managed through annotations and “syntactic sugar” declarations.

**3.1 Example.** The second part of our *first.hx* example is the syntax analyzer. Our small example source language merely has blocks, statements, and a few forms of expression, like so:

```

1  /* 2. SYNTAX ANALYSIS. */

2  sort Stat | [ [ <Name> := <Exp> ; ] | [ { <Stat*> } ] ] ;

3  sort Exp  | [ [ <Exp@1> + <Exp@2> ]@1
4             | [ [ <Exp@2> * <Exp@3> ]@2
5             | [ [ <Int> ]@3
6             | [ [ <Float> ]@3
7             | [ [ <Name> ]@3
8             | sugar [ ( <Exp#> ) ]@3 → # ;

9  sort Name | symbol [ [ <Id> ] ] ;

```

The grammar structures the input as three sorts: *Stat* for statements, *Exp* for expressions, and *Name* for names (which we shall need later for symbol tables). In addition a “derived” sort, *Stat\**, is used. HACS grammars follow these conventions:<sup>5</sup>

- Each sort is defined by a sort declaration followed by a number of *productions*, each introduced by a | (bar). (The first | corresponds to what is usually written “::=” or “→” in grammars.)
- Concrete syntax is enclosed in [ ... ] (“double” or “white” brackets). Everything inside double brackets should be seen as literal syntax, even \ (backslash), *except* for HACS white space (corresponding to [ \t\n\r ]), which is ignored, and fragments in <...> (angle brackets), which are special.
- References to *nonterminals* (other productions) are wrapped in <...> (angle brackets).
- *Precedence* is indicated with @*n*, where higher numbers *n* designate higher (tighter) precedence. Any reference in as well as the alternative itself may have a precedence in this way; in the example we establish that \* binds tighter than +, and that both operators are left recursive. Note that we specify the precedence of both the entire expression, after the [ ]s, and of each component, inside the < >. (In fact where the precedence is omitted we could have written @0.)
- The special sugar declaration expresses that the concrete syntax can use parentheses to raise the precedence of the enclosed expression to 3: it is the first example of a *rewrite rule* with a → that we see, where we remark that the expression is marked # so we can use the # to indicate that it is extracted as the abstract result of concrete syntax with parenthesis. (In fact the general rule is that when an → is used then all sort specifiers must be “disambiguated” with distinct markers like # or #5 in this way.)
- The *Name* sort is defined as a *symbol*, which is only allowed because the underlying *Id* token permits a trailing *\_n* (underscore and count); this permits the use as binders and automatic symbol generation.

**3.2 Commands.** We can parse an expression from the command line:

<sup>5</sup>As before the details of the Unicode characters are in the appendix.

```
$ ./first.run --sort=Exp --term="(2+(3*(4+5)))"
2 + 3 * ( 4 + 5 )
```

Notice that the printout differs slightly from the input term as it has been “resugared” from the AST with minimal insertion of parentheses.

**3.3 Manual** (syntactic sorts). Formally, HACS uses the following notations for specifying the syntax to use for terms.

1. HACS *production names* are capitalized words, so we can for example use `Exp` for the production of expressions. The name of a production also serves as the name of its *sort*, *i.e.*, the semantic category that is used internally for abstract syntax trees with that root production. If particular instances of a sort need to be referenced later they can be *disambiguated* with a `#i` suffix, *e.g.*, `Exp#2`, where *i* is an optional number or other simple word.
2. A sort is declared by one or more `sort` declarations of the name optionally followed by a number of *abstract syntax production* alternatives, each starting with a `|`. A sort declaration sets the *current sort* for subsequent declarations and in particular any stand-alone production alternatives. All sort declarations for a sort are cumulative.
3. Double square brackets `[[...]]` (unicode U+27E6 and U+27E7) are used for *concrete syntax* but can contain nested angle brackets `<...>` (unicode U+27E8 and U+27E9) with *production references* like `<Exp>` for an expression (as well as several other things that we will come to later). We for example write `[[<Exp>+<Exp>]]` to describe the form where two expressions are separated by a `+` sign.
4. Concrete syntax specification can include `¶` characters to indicate where *newlines* should be inserted in the printed output. (The system can also control indentation but that is not enabled yet.)
5. A trailing `@p` for some precedence integer *p* indicates that either the subexpression or the entire alternative (as appropriate) should be considered to have the indicated precedence, with higher numbers indicating higher precedence, *i.e.*, tighter association. (For details on the limitations of how the precedence and left recursion mechanisms are implemented, see Appendix A.2.)
6. `sugar [[...]]→...` alternatives specify equivalent forms for existing syntax: anything matching the left alternative will be interpreted the same as the right one (which must have been previously defined); references must be disambiguated.
7. A simple sort which contains only a reference to a token, where furthermore the token is defined such that it can end with `_n` (an underscore followed by a count), then the sort can be qualified as a *symbol* sort, which can be used for variables and binders.

## 4 Semantic Sorts and Schemes

When analyzing the AST, which we cover in the next section, we shall need some utility sorts and functions. This section details how these are defined.

**4.1 Example.** The third part of the *first.hx* example has the semantic sorts and operations used. For our toy language that just means the notion of *type* used and the way that types are unified.

```
1  /* 3. SEMANTIC SORTS & SCHEMES. */
2
3  sort Type | Int | Float ;
4
5  | scheme Unif(Type,Type) ;
6  Unif(Int, Int) → Int;
7  Unif(#t1, Float) → Float;
8  Unif(Float, #t2) → Float;
```

The code declares a new sort, `Type`, which is a *semantic* sort because it does not include any syntactic cases: all the possible values (as usual listed after leading `|`s) are simple term structures written without any `[]`s. Structures are written with a leading “constructor,” which should be a capitalized word (the same as sort names), optionally followed by some “arguments” in `()`s, where the declaration gives the sort for each argument.

The semantic sort also includes a *scheme* declaration for the `Unif` constructor followed by an argument list with two `Type` arguments. The scheme declaration is followed by *rules* of the form “pattern→replacement,” which must specify for each possible shape of `Unif`-construction how it should be simplified by the scheme. Rules may include “meta-variables” starting with # (hash), like #t1, to designate “function arguments” that should be copied from the pattern to the replacement.

**4.2 Manual** (raw terms, schemes, and rules). “Raw” declarations consist of the following elements:

1. A *constructor* is a capitalized word (similar to a sort name but in a separate name space).
2. A *variable* is a lower case word (subject to scoping, described below).
3. A sort can be given a *semantic production* as a | (bar) followed by a *form*, which consists of a constructor name, optionally followed by a ()ed , -separated list of *scope forms*, which each consist of a *sort* optionally preceded by a *binder form*, which is a list of sorts followed by a . (dot). Thus in the most general case, a semantic production has the form

$$| C ( S_{11} \cdots S_{1n_1} \cdot S_1 , \dots , S_{m1} \cdots S_{mn_m} \cdot S_m )$$

with  $C$  a constructor name and all  $S_i$  and  $S_{ij}$  sort names. The  $S_i$  declares the *argument sort* for the  $i$ th argument of the construction term, and the  $S_{ij}$  is the *binder sort* of the  $j$ th binder for the  $i$ th argument;  $m$  is the *arity* of the construction and  $n_i$  the *rank* of the  $i$ th argument.

4. A semantic production can be qualified as a *scheme*, which marks the declared construction as a candidate for rewrite rules (defined below).
5. A *raw term* is either a *construction*, a *variable use*, or a *meta-application*, as follows
  - (a) A *construction* term is a constructor name followed by an optional ()ed , -separated list of *scope arguments*, which each consist of a term optionally preceded by a *binder list*, which is a list of variables followed by a . (dot). So in the most general case, a term looks like this:

$$C ( x_{11} \cdots x_{1n_1} \cdot t_1 , \dots , x_{m1} \cdots x_{mn_m} \cdot t_m )$$

The “ $C$ -construction” is said to have the *subterms*  $t_1, \dots, t_m$ , and the arity  $m$  and ranks  $n_1 \dots n_m$  must correspond to a semantic production. If present, the binder prefix of each introduces the specified variables *only* for the appropriate subterm modulo usual renaming, i.e., writing  $A(x\ y.x, x\ y.y)$  and  $A(a\ b.a, a\ b.b)$  and even  $A(s\ t.s, t\ s.s)$  all denote the same term following the conventions of  $\alpha$ -equivalence. In a scope argument  $x.t$  we say that occurrences of  $x$  in  $t$  are *bound* by the binder.

- (b) A *variable use* term is a variable, subject to the usual lexical scoping rules.
  - (c) A *meta-application* term is a *meta-variable*, consisting of a # (hash) followed by a number or word and optionally by a meta-argument list of , -separated terms enclosed in []. Examples include #t1 (with no arguments), #[a,b,c], and #1[OK,#].
6. A term can have a *sort prefix*. So the term `Type Unif(Type #t1, Type Float)` is the same as `Unif(#t1,Float)` provided `Unif` was declared with the raw production `|Unif(Type,Type)`.
  7. A *rewrite rule* is a pair of terms separated by  $\rightarrow$  (arrow, U+2192), with a few additional constraints: in the rule  $p \rightarrow t$ ,  $p$  must be a *pattern*, which means it must be a construction term that has been declared as a *scheme* (syntactic or raw) and with the restriction that all contained arguments to meta-applications must be bound variables, and all meta-applications in  $t$  must have meta-variables that also occur in  $p$  with the same number of meta-arguments.  
Rule declarations must either occur with the appropriate current sort or have a pattern with a sort prefix.
  8. One rule per scheme can be prefixed with the qualifier `default`. If so then the pattern can have no structure: all subterms of the pattern scheme construction must be plain meta-applications. Such a default rule is applied *after* it has been ensured that all other rules fail for the scheme.
  9. Finally, a rule can be prefixed with the word `rule` for clarity, which can be followed by a raw term and a : (TODO: document possible name/option choices).

Rules are used for *rewriting*, a definition of which is beyond the scope of this document; please refer to the literature on higher order rewriting for details [2].

## 5 Semantic Analysis I

Once the AST is loaded we can analyze it. We demonstrate this by fully developing a simple type checking analysis. In this section we deal with analysis of expressions.

**5.1 Example.** In Ex. 3.1 we presented the abstract syntax of the small language processed by *first.hx*. A type analysis of the expressions of language might look as follows as a standard SDT (syntax directed translation), where we use  $E$  for the Exp non-terminal, and these attributes:

- $E.e$  inherited “environment” symbol table mapping names to types.
- $E.t$  synthesized type of expression.

In the notations of [1], the SDT can be specified something like this:

PRODUCTION	SEMANTIC RULES
$E \rightarrow E_1 + E_2$	$E_1.e = E.e; E_2.e = E.e; E.t = \text{Unif}(E_1.t, E_2.t)$ (E1)
$  E_1 * E_2$	$E_1.e = E.e; E_2.e = E.e; E.t = \text{Unif}(E_1.t, E_2.t)$ (E2)
$  \text{int}$	$E.t = \text{Int}$ (E3)
$  \text{float}$	$E.t = \text{Float}$ (E4)
$  \text{name}$	$E.t = \text{if defined}(E.e, \text{name.sym}) \text{ then lookup}(E.e, \text{name.sym})$ (E5) else TypeError

The SDT has one locations where the attribute dependencies are interesting: (E5) defines the synthesized attribute  $E.t$  in terms of the inherited attribute  $E.e$  in a “down then up” configuration. To translate this into HACS, we create a recursive scheme to propagate inherited attributes as needed.

The following code contains the fourth part of *first.hx*, which implements the SDT in HACS.

```

1  /* 4. SEMANTIC EXPRESSION ANALYSIS. */

2  attribute ↑t(Type);      // synthesized expression type
3  attribute ↓e{Name:Type}; // inherited type environment

4  // TA, expression case.

5  sort Exp | scheme [ TA <Exp> ] ↓e ;

6  [ TA <Int#> ] → [ <Int#> ];
7  [ TA <Float#> ] → [ <Float#> ];
8  [ TA (<Exp#1> + <Exp#2>) ] → [ (TA <Exp#1>) + (TA <Exp#2>) ];
9  [ TA (<Exp#1> * <Exp#2>) ] → [ (TA <Exp#1>) * (TA <Exp#2>) ];
10 [ TA <Name id> ] ↓e{id : #t} → [ id ] ↑t(#t);
11 [ TA <Name id> ] ↓e{¬id} → error[Undefined identifier <id>];

12 sort Exp; ↑t;
13 [ <Exp#1 ↑t(#t1)> + <Exp#2 ↑t(#t2)> ] ↑t(Unif(#t1,#t2));
14 [ <Exp#1 ↑t(#t1)> * <Exp#2 ↑t(#t2)> ] ↑t(Unif(#t1,#t2));
15 [ <Int#> ] ↑t(Int);
16 [ <Float#> ] ↑t(Float);

```



The first block sets up the synthesized  $E.t$  and inherited  $E.e$  attributes, using  $\uparrow$  and  $\downarrow$  for synthetic/inherited, respectively, and distinguishing between a value (Type) and symbol mapping  $\{\text{Name}:\text{Type}\}$ .

This is followed by a definition of the recursive TA scheme for *distributing* all types in the symbol table to actual symbol occurrences. Specifically,  $\text{scheme}[\llbracket \text{TA}(\text{Exp}) \rrbracket \downarrow e]$  means that (the sort  $\text{Exp}$ ) has a scheme with the specified syntax that distributes the inherited attribute  $e$ . Note that the scheme uses *concrete syntax* for applications of the scheme: this is a convenience that makes it possible to invoke specific schemes directly from expressions entered in the command line. It does mean that we have to be careful with using the syntax correctly: we have not, for example, given a precedence to TA, so in the recursive rules we have expressions like

$$\llbracket (\text{TA}(\text{Exp}\#1)) + (\text{TA}(\text{Exp}\#2)) \rrbracket$$

which uses the  $(\dots)$  syntactic sugar to properly nest the expression. We could also have written

$$\llbracket \langle \text{Exp} \llbracket \text{TA}(\text{Exp}\#1) \rrbracket \rangle + \langle \text{Exp} \llbracket \text{TA}(\text{Exp}\#2) \rrbracket \rangle \rrbracket$$

using nested  $\langle \dots \rangle$  constructs instead, which does not depend on the syntactic sugar.

Otherwise, the scheme just exploits the syntax to express *where* the recursion goes, which then automatically means that the inherited attribute will follow, as HACS inherited attributes always follow a scheme. The first two “TA-rules” have the cases for looking up a symbol in the environment: if it is there, *i.e.*, if the  $\text{id}$  specified can make the environment pattern  $\downarrow\{\text{id}:\#t\}$  match, instantiating  $\#t$  to the type of  $\text{id}$ , then the first rule is used and the result is the  $\text{id}$  with the added synthesized attribute value  $\uparrow t(\#t)$ , *i.e.*, the  $t$  attribute gets the value we value we extracted from the environment. The rule with the negated pattern  $\downarrow e\{\neg \text{id}\}$  will match if the  $\text{id}$  symbol is *not* defined in the environment, then printing a simple error message. In a sense the two Name cases correspond to  $(E5)$  in the SDT.

Finally, we associate the synthetic attribute  $t$  to the  $\text{Exp}$  sort—all synthetic attributes are associated with one or more sorts—and we give a special *synthesis rule* that shows for each form of  $\text{Exp}$ , where  $t$  is not populated by TA, how synthetic attributes matches on fragments propagate to the entire expression, *e.g.*, the rule for the SDT’s (1) is

$$\llbracket \langle \text{Exp}\#1 \uparrow t(\#t1) \rangle + \langle \text{Exp}\#2 \uparrow t(\#t2) \rangle \rrbracket \uparrow t(\text{Unif}(\#t1, \#t2));$$

and should be read “When considering an  $\text{Exp}$  (the current sort) which has the shape  $\llbracket \langle \langle \text{Exp} \rangle + \langle \text{Exp} \rangle \rrbracket$  where furthermore the first expression has a value matching  $\#t1$  for the synthesized attribute  $t$ , and the second expression has a value matching  $\#t2$  for the synthesized attribute  $t$ , then the entire expression has the value  $\text{Unif}(\#t1, \#t2)$  for the synthesized attribute  $t$ .” Notice that the attribute patterns have the same shape as the attribute declaration—a recurring feature of HACS. The following three lines correspond in the same way to  $(E2, E3, E4)$ .

## 5.2 Manual (attributes and synthesis rules).

1. Attributes are declared by attribute declarations followed by an *attribute form* of one of the following shapes:
  - (a)  $\uparrow \text{Name}(\text{ValueSort})$  defines that the synthesized attribute  $\text{Name}$  has  $\text{ValueSort}$  values;
  - (b)  $\downarrow \text{Name}(\text{ValueSort})$  similarly for a simple inherited attribute;
  - (c)  $\downarrow \text{Name}\{\text{SymbolSort}:\text{ValueSort}\}$  defines the inherited symbol table attribute  $\text{Name}$  which for each constant or variable of  $\text{SymbolSort}$  has a distinct  $\text{ValueSort}$  value.
2. One can add a simple *synthesized attributes* after a raw data term as  $\uparrow \text{name}(\text{value})$ , where the  $\text{name}$  is an attribute name and the  $\text{value}$  can be any term.
3. Simple *inherited attributes* are added similarly after a raw scheme term as  $\downarrow \text{name}(\text{value})$ .

4. An *inherited symbol table attribute extension* is added to a raw scheme term as  $\downarrow name\{symbol:value\}$ , where the *symbol* is either a variable or a constant (of the appropriate sort).
5. A *synthesized attribute reference* has the simple form  $\uparrow name$ ; and declares that the current sort synthesizes *name* attributes.
6. A scheme declaration can include *inherited attribute references* of the form  $\downarrow name$ , which declares that the scheme inherits the *name* attributes.
7. A *synthesis rule* is a special rule of the form  $t \uparrow name(t')$ , where the term  $t$  may contain subterms with attribute constraints. The rule specifies how terms of the current sort and shape  $t$  synthesize *name* attributes.

Inherited attributes are managed with regular rules (for schemes) with inherited attribute constraints and extensions.

**5.3 Manual** (parsed terms). Man. 4.2 documented terms without concrete syntax. The full term model combines this with *parsed terms*.

1. Double square brackets  $\llbracket \dots \rrbracket$  (unicode U+27E6 and U+27E7) can be used for *concrete terms*, provided the *sort* is clear, either
  - (a) by immediately prefixing with the sort (as in  $\text{Exp}\llbracket 1+2 \rrbracket$ ), or
  - (b) by using as the argument of a defined constructor (as  $\text{IsType}(\llbracket \text{mytype} \rrbracket)$ ), or
  - (c) by using as an attribute value, or
  - (d) by using as a top level rule pattern or replacement term with a defined current sort.
2. Concrete terms can contain nested raw terms in  $\langle \dots \rangle$  (unicode U+27E8 and U+27E9). Such nested raw terms *must* have an explicit sort prefix.
3. The special term  $\text{error}\llbracket \dots \rrbracket$  will print the error message embedded in  $\llbracket \dots \rrbracket$ , where one is permitted to embed *symbol*-declared variables in  $\langle \dots \rangle$ .

## 6 Semantic Analysis II

We continue the analysis by explaining how the environment is initialized and maintained across statements, where new names are introduced.

**6.1 Example.** Our type analysis from Ex. 5.1 is extended to  $S, S^*$  for the non-terminals for the sorts *Stat* and *Stat\**, respectively, with these attributes:

- $S.e, S^*.e$  inherited “environment” symbol table mapping names to types.
- $S.sym$  synthesized “defines” attribute with name introduced by  $S$ , or  $\epsilon$  if none.
- $S.t$  synthesized type of  $S.sym$ , if any.

with these rules:

PRODUCTION	SEMANTIC RULES
$S \rightarrow \mathbf{name} := E$	$E.e = S.e; S.sym = \mathbf{name}.sym; S.t = E.t$ (1)
$  \{ S^* \}$	$S^*.e = S.e; S.sym = \epsilon; S.t = \epsilon$ (2)
$S^* \rightarrow S_1 S_2^*$	$S_1.e = S^*.e$ (3)
$  \epsilon$	$S_2^*.e = \text{if } S_1.sym \neq \epsilon \text{ then extend}(S^*.e, S_1.sym, S_1.t) \text{ else } S^*.e$ (4)

The SDT has one locations where the attribute dependencies are interesting: (3) copies the synthesized  $S_1.sym$  and  $S_1.t$  “left to right” into the inherited attribute  $S_2^*.e$ . To translate this into HACS, we create a specific scheme for each “left to right” dependency. Also note that in this particular SDT, the  $S.sym$  and  $S.t$  attributes are used in a trivial way to just allow the containing  $S^*$  sequence to extract information from assignments  $\mathbf{name} := E$  and avoid nested blocks.

Here is the code with the fifth part of *first.hx*, which implements the SDT above in HACS.

```

1  /* 5. SEMANTIC ANALYSIS. */
2  // TA, Statement case.
3  sort Stat | scheme [ TA <Stat> ] ↓e ;
4  [ TA <Name id> := <Exp#2>; ] → [ <Name id> := TA <Exp#2>; ];
5  [ TA {} ] → [ {} ];
6  [ TA { <Name id> := <Exp#2>; <Stat*#3> } ] → [ TA2 { <Name id> := TA <Exp#2>; <Stat*#3> } ];
7  {
8    | scheme [ TA2 <Stat> ] ↓e;
9    [ TA2 { <Name id> := <Exp#2> ↑t(#t2); <Stat*#3> } ] →
10     [ { <Name id> := <Exp#2>; <Stat [TA {<Stat*#3>}] ↓e{id:#t2}> } ];
11 }
12 [ TA { {<Stat*#1>} <Stat*#2> } ] → [ { TA {<Stat*#1>} TA {<Stat*#2>} } ];

```

For (2,3,4) a slight refactoring is in order. Imagine that we could have written the following pseudo-SDT instead, inlining the leading statement of a sequence and unfolding the rather silly “one-step”  $S.sym$  and  $S.t$  attributes and their “absent”  $\epsilon$  value:

$S \rightarrow \mathbf{name} := E$	$E.e = S.e$	(1)
$  \{ S^* \}$	$S^*.e = S.e$	(2)
$S^* \rightarrow \epsilon$		(3)
$  \mathbf{name} := E; S_2^*$	$S_2^*.e = \text{extend}(S^*.e, \mathbf{name}.sym, E.t)$	(4.1)
$  \{ S_1^* \} S_2^*$	$S_2^*.e = S_1^*.e$	(4.2)

This is not a usual SDT, because the cases do not follow the productions strictly, however, it is easy to see that these rules cover the full language just in a slightly different way. This trick, which we can call “inlining the sequence head,” is frequently very useful. For (1, 2, 3, 4.2) these modified rules are straightforward propagations, implemented directly in the code. (The code in fact goes even further than this: in order to keep the TA scheme on just plain Stat they implement the cases for  $S^*$  by reconstructing a single statement  $\{S^*\}$  in each case.)

Next, we handle (4.1), the “left-to-right” dependency, where the inherited attribute  $S_2^*.e$  depends on the synthesized attribute  $E.t$ . In HACS this means that a *wrapper scheme* is needed, and we implement (4.1) as follows (with Stat the current sort):

```

[ TA { <Name id> := <Exp#2>; <Stat*#3> } ] → [ TA2 { <Name id> := TA <Exp#2>; <Stat*#3> } ];
{
  | scheme [ TA2 <Stat> ] ↓e;
  [ TA2 { <Name id> := <Exp#2> ↑t(#t2); <Stat*#3> } ] →
    [ { <Name id> := <Exp#2>; <Stat [TA {<Stat*#3>}] ↓e{id:#t2}> } ];
}

```

The first rule recurses the TA scheme into the expression to associate types with any names. A new, second scheme, TA2, then takes over once the synthesized attribute has been instantiated for the expression, passing the environment extension to the subsequent statement. (The  $\{\}$  wrappers serve as a grouping for the “current sort” status, which here is not in fact changed.)

Here is the breakdown of the TA2 rule. The pattern is an instance of the TA2 syntactic scheme with these fragments:

- TA2—the prefix marking this as an instance of the scheme.
- $\{ \dots \}$  wrapper making the contained  $\text{Stat}^*$  sequence into a single  $\text{Stat}$  as required.
- $\langle \text{Name id} \rangle :=$  an identifier being assigned as part of the first assignment statement; matching will remember the identifier symbol as  $\text{id}$ .
- $\langle \text{Exp\#2 } \uparrow t(\#t2) \rangle$  is the assigned expression, which includes a constraint that the rule can only fire if there is in fact a type annotation that can match  $\#t2$ , corresponding to the idea that it has already been analyzed.
- $;$  the end of the assignment statement.
- $\langle \text{Stat}\#3 \rangle$  is the following statement sequence that we shall annotate next.

Then comes the arrow ( $\rightarrow$ ) followed by the replacement term:

- $\{ \dots \}$  wrapper making the contained  $\text{Stat}^*$  sequence into a single result  $\text{Stat}$  as required.
- $\langle \text{Name\#1} \rangle := \langle \text{Exp\#2} \rangle;$  is the assignment that we matched, unchanged (omitting the attribute means it is not modified).
- $\langle \text{Stat } \dots \rangle$  denotes the replacement follow-on (single) statement, where the “ $\dots$ ” in this case is an embedded computation, namely—
- $\llbracket \text{TA } \{ \langle \text{Stat}\#3 \rangle \} \rrbracket \downarrow e\{\text{id}:\#t2\}$  is the embedded computation to perform to annotate the following statements properly. It is a recursive call of TA on a single statement containing the remainder of the statements, annotated with the inherited type environment extended with the binding from the symbol to its type.

Notice how the environment extension uses nested  $\llbracket \dots \rrbracket$ -brackets because the invocation of TA on the following statement needs the inherited attribute argument, which is not permitted in the language syntax.

## 6.2 Manual (attributes and synthesis rules).

1. Attributes are declared by *attribute declarations* followed by an *attribute form* of one of the following shapes:
  - (a)  $\uparrow \text{Name}(\text{ValueSort})$  defines that the synthesized attribute  $\text{Name}$  has  $\text{ValueSort}$  values;
  - (b)  $\downarrow \text{Name}(\text{ValueSort})$  similarly for a simple inherited attribute;
  - (c)  $\downarrow \text{Name}\{\text{SymbolSort}:\text{ValueSort}\}$  defines the inherited symbol table attribute  $\text{Name}$  which for each constant or variable of  $\text{SymbolSort}$  has a distinct  $\text{ValueSort}$  value.
2. One can add a simple *synthesized attributes* after a raw data term as  $\uparrow \text{name}(\text{value})$ , where the  $\text{name}$  is an attribute name and the  $\text{value}$  can be any term.
3. Simple *inherited attributes* are added similarly after a raw scheme term as  $\downarrow \text{name}(\text{value})$ .
4. An *inherited symbol table attribute extension* is added to a raw scheme term as  $\downarrow \text{name}\{\text{symbol}:\text{value}\}$ , where the  $\text{symbol}$  is either a variable or a constant (of the appropriate sort).
5. A *synthesized attribute reference* has the simple form  $\uparrow \text{name}$ ; and declares that the current sort synthesizes  $\text{name}$  attributes.
6. A scheme declaration can include *inherited attribute references* of the form  $\downarrow \text{name}$ , which declares that the scheme inherits the  $\text{name}$  attributes.
7. A *synthesis rule* is a special rule of the form  $t \uparrow \text{name}(t')$ , where the term  $t$  may contain subterms with attribute constraints. The rule specifies how terms of the current sort and shape  $t$  synthesize  $\text{name}$  attributes.

Inherited attributes are managed with regular rules (for schemes) with inherited attribute constraints and extensions.

**6.3 Manual** (special replacement terms). The special term  $\text{error}\llbracket \dots \rrbracket$  will print the error message embedded in  $\llbracket \dots \rrbracket$ , where one is permitted to embed *symbol*-declared variables in  $\langle \dots \rangle$ .

## 7 Intermediate and Code Generation

After the analysis we are ready for generating code.

**7.1 Example.** The sixth part of *first.hx* is the translation from abstract syntax to the intermediate representation, shown in Fig. 1. The fragment contains the usual components: a syntax specification, rewrite schemes, and rewrite rules for the ICG scheme.

The code only uses two new features: ¶ markers in the syntax to indicate newlines, and rules that introduce *fresh* variables (of Tmp sort): when the replacement of a rule uses a symbol, which was not in the pattern, then this corresponds to *generating* a new globally unique symbol. So each time the rule

```
[[ ICG id := <Exp#2 ↑HasType(#t2)>; ]]  
→ [[ { <I_Progr [[ICGExp T <Exp#2>]] ↓TmpType{T:#t2}> } id = T; ] ] ;
```

is used, T denotes a new so-called “fresh” symbol. When printed, the various incarnations of T will be named T\_1, T\_86, etc.

**7.2 Example.** The seventh part of *first.hx* is the final translation CG from the intermediate representation to assembly code. This uses no new features, and is shown in Fig. 2, however, it is still worth a sanity check, walking through the CG scheme and checking that all syntactic cases are covered.

**7.3 Remark** (concrete vs. raw syntax). In the presentation we have chosen to use *concrete syntax* even for semantic operations. This has the advantage of allowing direct invocation of even complex structured calculations from the command line but it does “pollute” the syntax of the defined language. (Production versions of HACS (not yet released) will have the option of generating parsers that ignore concrete syntax of schemes when running the compiler.) It is sometimes practical to define “bridge schemes” that make schemes available both in syntax and raw; we give an example of this in the following section.

## 8 Main

Finally, we put everything together.

**8.1 Example.** The eighth and last fragment of *first.hx* is the main compilation function, which pulls together the three stages defined in previous sections.

```
1 /* 8. MAIN. */  
  
2 sort A_Progr | scheme [[ Compile <Stat> ] ] ;  
3 [[ Compile <Stat#1> ] ] → [[ CG ICG TA <Stat#1> ] ] ;
```

It is also possible to define wrapper schemes that can be invoked with an option, this is useful when the input term is read from a file where it is not obvious to include the compiler instructions.

```
4 | scheme Compile(Stat);  
5 Compile(#stat) → [[ Compile <Stat#stat> ] ] ;
```

This is the --action we invoked back in Com. 1.2. Such wrapper raw schemes must have a single argument.

```

1  /* 6. INTERMEDIATE CODE GENERATION. */

2  token T | T ('_' <Int>)? ; // temporary

3  // Concrete syntax & abstract syntax sorts.

4  sort I_Progr | [[<I_Instr> <I_Progr>]] | [[]] ;

5  sort I_Instr | [[<Tmp> = <I_Arg> + <I_Arg>; ¶]]
6                | [[<Tmp> = <I_Arg> * <I_Arg>; ¶]]
7                | [[<Tmp> = <I_Arg>; ¶]]
8                | [[<Name> = <Tmp>; ¶]] ;

9  sort I_Arg | [[<Name>]]
10             | [[<Float>]]
11             | [[<Int>]]
12             | [[<Tmp>]] ;

13 sort Tmp | symbol [[ <T> ]] ;

14 // Translation scheme.

15 attribute ↓TmpType{Tmp:Type} ;

16 sort I_Progr ;

17 | scheme [[ ICG <Stat> ]] ↓TmpType ;
18 [[ ICG id := <Exp#2 ↑HasType(#t2)>; ]]
19   → [[ { <I_Progr [[ICGExp T <Exp#2>]] ↓TmpType{T:#t2}> } id = T; ]] ;
20 [[ ICG { } ]] → [[]] ;
21 [[ ICG { <Stat#s> <Stat*#ss> } ]] → [[ { ICG <Stat#s> } ICG { <Stat*#ss> } ]] ;

22 | scheme [[ ICGExp <Tmp> <Exp> ]] ;

23 [[ ICGExp T <Int#1> ]] → [[ T = <Int#1>; ]] ;
24 [[ ICGExp T <Float#1> ]] → [[ T = <Float#1>; ]] ;
25 [[ ICGExp T id ]] → [[ T = id; ]] ;

26 [[ ICGExp T <Exp#1> + <Exp#2> ]]
27   → [[ {ICGExp T_1 <Exp#1>} {ICGExp T_2 <Exp#2>} T = T_1 + T_2; ]] ;

28 [[ ICGExp T <Exp#1> * <Exp#2> ]]
29   → [[ {ICGExp T_1 <Exp#1>} {ICGExp T_2 <Exp#2>} T = T_1 * T_2; ]] ;

30 // Helper to flatten code sequence.
31 | scheme [[ {<I_Progr>} <I_Progr> ]] ;
32 [[ { } <I_Progr#3> ]] → #3 ;
33 [[ {<I_Instr#1> <I_Progr#2>} <I_Progr#3> ]] → [[ <I_Instr#1> {<I_Progr#2>} <I_Progr#3> ]] ;

```

Figure 1: Intermediate Code Generation.

```

1  /* 7. CODE GENERATOR. */

2  sort A_Progr | [[ ⟨A_Instr⟩ ⟨A_Progr⟩ ]] | [[]] ;

3  sort A_Instr | [[ LDF ⟨Tmp⟩, ⟨A_Arg⟩ ¶ ]]
4                  | [[ STF ⟨Name⟩, ⟨Tmp⟩ ¶ ]]
5                  | [[ ADDF ⟨A_Arg⟩, ⟨A_Arg⟩, ⟨A_Arg⟩ ¶ ]]
6                  | [[ MULF ⟨A_Arg⟩, ⟨A_Arg⟩, ⟨A_Arg⟩ ¶ ]] ;

7  sort A_Arg | [[ #⟨Float⟩ ]] | [[ #⟨Int⟩ ]] | [[ ⟨Name⟩ ]] | [[ ⟨Tmp⟩ ]] ;

8  sort A_Progr | scheme [[ CG ⟨I_Progr⟩ ]] ;

9  [[ CG ]] → [[]] ;

10 [[ CG T = ⟨I_Arg#1⟩ + ⟨I_Arg#2⟩ ; ⟨I_Progr#⟩ ]]
11   → [[ ADDF T, [⟨I_Arg#1⟩], [⟨I_Arg#2⟩] CG ⟨I_Progr#⟩ ]] ;

12 [[ CG T = ⟨I_Arg#1⟩ * ⟨I_Arg#2⟩ ; ⟨I_Progr#⟩ ]]
13   → [[ MULF T, [⟨I_Arg#1⟩], [⟨I_Arg#2⟩] CG ⟨I_Progr#⟩ ]] ;

14 [[ CG T = ⟨I_Arg#1⟩ ; ⟨I_Progr#⟩ ]]
15   → [[ LDF T, [⟨I_Arg#1⟩] CG ⟨I_Progr#⟩ ]] ;

16 [[ CG name = T ; ⟨I_Progr#⟩ ]]
17   → [[ STF name, T CG ⟨I_Progr#⟩ ]] ;

18 sort A_Arg | scheme [[ [⟨I_Arg⟩] ]] ;
19 [[ [T] ]] → [[ T ]] ;
20 [[ [name] ]] → [[ name ]] ;
21 [[ [⟨Float#1⟩] ]] → [[ #⟨Float#1⟩ ]] ;
22 [[ [⟨Int#1⟩] ]] → [[ #⟨Int#1⟩ ]] ;

```

Figure 2: Code Generation.

**8.2 Manual** (building and running). To use HACS you need a copy of the *hacs* directory somewhere. If you are working on a HACS specification, say *mycompiler.hx*, then your working directory should have a *Makefile* containing (at least) the following:

```
1 # Makefile for mycompiler.

2 # HACS configuration.
3 HACS = $(abspath hacs)
4 include $(HACS)/Makefile-hx

5 # Dependencies.
6 mycompiler.run: mycompiler.hx
```

where you have replaced *hacs* with the path to the HACS directory on your system (as written the system expects to find *hacs* as a local subdirectory of your working directory). With this setup, and a suitable *mycompiler.hx*, you have the following options:

1. `make mycompiler.run` will generate the script *mycompiler.run*, which implements the compiler you specify in *mycompiler.hx*. Generation creates a number of support files, specifically
  - (a) *build* subdirectory has runtime resources needed by the script.
  - (b) *src* subdirectory contains auxiliary Java resources used by the build process.
  - (c) *mycompiler.crs-installed* and *mycompiler.pg-installed* record when the generated parser and rewrite system were installed.

Note that the first time, some utility programs are compiled from C to ensure that HACS is fully enabled on your system.

2. `make clean` will remove all temporaries not needed for running the compiler script.
3. `make realclean` will remove all traces of the generated compiler.
4. `make distclean` will remove all traces of the generated compiler as well as all generated HACS tooling files. Do not do this unless you really mean it: recovering a useable HACS system requires getting a fresh copy of *hacs.zip* or half an hour of CPU time on nontrivial hardware.

The generated script refers absolutely to files under the *hacs* directory, so the generated *mycompiler.run* script itself can be moved but the *hacs* directory cannot.

The script accepts a number of options:

1. `--sort=Sort` sets the expected sort (and thus parser productions) for the input to *Sort*. The input is read, normalized, and printed.
2. `--action=Constructor` sets the computation for the compiler to *Constructor*, which must be a unary raw scheme; the argument sort of *Constructor* defines the parser productions to use. The input is read, wrapped in the action, normalized, and printed.
3. `--term=term` use the *term* as the input.
4. `--input=file` reads the input from *file*.
5. `--output=file` sends the input to *file* (the default is the standard output).
6. `--verbose=n` sets the verbosity of the underlying CRSX rewrite engine to *n*. The default is 0 (quiet) but 1–3 are useful (above 3 you get a lot of low level diagnostic output).
7. `--parse-verbose` activates verbose output from JavaCC of the parsing.

You must provide one of `--sort` or `--action`, and one of `--term` and `--input`.

HACS will eventually contain a convention for defining “main” sorts and schemes such that defaults can be provided for the configuration options.



## A Tables

### A.1 Notation (used unicode characters).

<i>Glyph</i>	<i>Code Point</i>	<i>Character</i>
¬	U+00AC	logical negation sign
¶	U+00B6	paragraph sign
↑	U+2191	upwards arrow
→	U+2192	rightwards arrow
↓	U+2193	downwards arrow
⌈	U+27E6	mathematical left white square bracket
⌋	U+27E7	mathematical right white square bracket
⟨	U+27E8	mathematical left angle bracket
⟩	U+27E9	mathematical right angle bracket

If you are using a Linux computer (or any where the display is controlled by the X Window System), then you can create a file called `.XCompose`, which has your special characters. A minimal one looks something like this (adjust for your locale, etc.):

```
include "/usr/share/X11/locale/en_US.UTF-8/Compose"
<Multi_key> <P> <P>          : "¶" U00B6 # PARAGRAPH SIGN
<Multi_key> <u> <greater>      : "↑" U2191 # UPWARDS ARROW
<Multi_key> <minus> <greater>   : "→" U2192 # RIGHTWARDS ARROW
<Multi_key> <d> <greater>       : "↓" U2193 # DOWNWARDS ARROW
<Multi_key> <bracketleft> <bracketleft> : "⌈" U27E6 # MATHEMATICAL LEFT WHITE SQUARE BRACKET
<Multi_key> <bracketright> <bracketright> : "⌋" U27E7 # MATHEMATICAL RIGHT WHITE SQUARE BRACKET
<Multi_key> <less> <period>      : "⟨" U27E8 # MATHEMATICAL LEFT ANGLE BRACKET
<Multi_key> <greater> <period>    : "⟩" U27E9 # MATHEMATICAL RIGHT ANGLE BRACKET
```

(We have included a slightly more extensive one in the *gentle* directory.) You should then map the *MultiKey* key to some key on your keyboard in the keyboard options (this is done in the keyboard settings), and then you are ready to enter the fancy characters as the indicated three-key combinations. In some programs (like Eclipse) you need to select the “Input Method” to be the “X Input Method” for this to work.

### A.2 Manual (limitations).

- At most one nested declaration per token.
- Precedence can only be used on self references, *i.e.*,  $\langle E@2 \rangle$  can only occur inside productions for the sort *E*.
- It is not possible to use binders and left recursion in the same production with the same precedence.
- Only *direct* left recursion is currently supported, *i.e.*, the left recursion should be within a single production.
- Productions can share a prefix but only within productions for the same sort, and the prefix has to be literally identical unit by unit, *i.e.*,

```
sort S | [ [ ⟨A⟩ then ⟨B⟩ then C ]
        | [ [ ⟨A⟩ then ⟨B⟩ or else D ] ] ;
```

is fine but

```
sort S | [ [ ⟨A⟩ then ⟨B⟩ then C ]
        | [ [ ⟨A⟩ ⟨ThenB⟩ or else D ] ] ;
sort ThenB | [ [ then ⟨B⟩ ] ] ;
```

is not.

- It is not possible to left-factor a binder (so multiple binding constructs cannot have the same binder prefix).
- Binders must occur to the left of all their occurrences.
- Repeated production references cannot be used in sorts that also use precedence, *i.e.*, you cannot mix  $\langle E^* \rangle$  and  $\langle E@2 \rangle$  inside the productions for the sort *E*.
- Variables embedded in `error[...]` instructions must start with a lower case letter.

- When using the `symbol` qualifier on a reference to a token then the token *must* allow ending in `_n` for  $n$  any natural number.
- When using the same name for a symbol inside of `[[...]]` and the corresponding raw variable outside of the `[[...]]`, then the common symbol and variable name must be a plain word starting with a lower case letter.
- Special terms like `error[[...]]` cannot be used as raw subterms.
- The default rule qualifier is rather fragile and does not yet always work.

## B Bonus Examples

Some additional examples. *Note: still need to be worked through.*

### B.1 Example (`crsx/samples/gentle/bool.hx`).

```
module "net.sf.crsx.samples.gentle.Bool" {

  // Boolean sort.
  sort B | [[ t ]@4 | [f]@4                // constants
           | sugar [( <B#> ) ]@4 → [[<B#>]] // parenthesis
           ;

  // Disjunction.
  | scheme [[<B@2> ∨ <B@1> ]@1 ;
  [[t ∨ <B#>]] → [[t] ;
  [[f ∨ <B#>]] → B# ;

  // Conjunction.
  | scheme [[<B@3> ∧ <B@2> ]@2 ;
  [[t ∧ <B#>]] → B# ;
  [[f ∧ <B#>]] → [f] ;

  // Negation.
  | scheme [[¬ <B@3> ]@3 ;
  [[¬t] → [f] ;
  [[¬f] → [t] ;

}
```

### B.2 Example (`crsx/samples/gentle/deriv.hx`).

```
module "net.sf.crsx.samples.gentle.Deriv" {

  import "bool.hx"(B);

  //// SYNTAX.

  // Arithmetic and basic functions.

  sort Exp | [[ <Exp@1> + <Exp@2> ]@1
            | [[ <Exp@1> - <Exp@2> ]@1
```

```

| [[ <Exp@2> * <Exp@3> ]]@2
| [[ <Exp@2> / <Exp@3> ]]@2
| [[ <Fun> <Exp@4> ]]@3
| [[ <Int> ]]@4

| sugar [[ ( <Exp#> ) ]]@4 → Exp#
| sugar [[ + <Exp#1@2> ]]@1 → [[ 0 + <Exp#1> ]]
| sugar [[ - <Exp#1@2> ]]@1 → [[ 0 - <Exp#1> ]]
;

sort Fun | [[sin]]@2 | [[cos]]@2 | [[ln]]@2 | [[exp]]@2 ;

token Int | [0-9]+ ;

// Functionals.

sort Exp | symbol [[<Var> ]]@4 ;
token Var | [a-z] [A-Za-z0-9]* ;
sort Fun | [[ [ <Var#1> ] ↦ <Exp[Var#1:Exp]> ] ]@2 ;

//// SCHEMES.

sort Fun | scheme [[d<Fun@1>]]@1 ;

[[d sin]] → [[cos]] ;
[[d cos]] → [[a ↦ - sin a]] ;
[[d ln]] → [[z ↦ 1/z]] ;
[[d exp]] → [[exp]] ;

[[d[x ↦ <Exp#1[x]>]]] → [[y ↦ D(y)[z ↦ <Exp#1[z]>]]] ;

sort Exp | scheme [[D <Exp> [ <[Var#1]> ↦ <Exp[Var#1:Exp]> ] ]@3 ;

[[ D<Exp#1>[x ↦ <Int#2>] ]] → [[0]] ;

[[ D<Exp#1>[x ↦ x] ]] → [[1]] ;
[[ D<Exp#1>[x ↦ y] ]] → [[0]] ;

[[ D<Exp#0>[x ↦ <Exp#1[x]> + <Exp#2[x]> ] ]] → [[ D<Exp#0>[y ↦ <Exp#1[y]> ] + D<Exp#0>[z ↦ <Exp#2[z]>] ]] ;
[[ D<Exp#0>[x ↦ <Exp#1[x]> - <Exp#2[x]> ] ]] → [[ D<Exp#0>[y ↦ <Exp#1[y]> ] - D<Exp#0>[z ↦ <Exp#2[z]>] ]] ;

[[ D<Exp#>[x ↦ <Exp#1[x]> * <Exp#2[x]> ] ]]
→ [[ D<Exp#>[x ↦ <Exp#1[x]> ] * <Exp#2[#> ] + <Exp#1[#> ] * D<Exp#>[x ↦ <Exp#2[x]> ] ]] ;

[[ D<Exp#>[x ↦ <Exp#1[x]> / <Exp#2[x]> ] ]]
→ [[ ( D<Exp#>[x ↦ <Exp#1[x]> ] * <Exp#2[#> ] - <Exp#1[#> ] * D<Exp#>[x ↦ <Exp#2[x]> ] )
/ ( <Exp#2[#> ] * <Exp#2[#> ] ) ]] ;

```

```
[[ D⟨Exp#⟩[x→⟨Fun#f⟩⟨Exp#2[x]⟩] ]→[[d⟨Fun#f⟩ ⟨Exp#2[#]⟩ * D⟨Exp#⟩[x→⟨Exp#2[x]⟩] ]];
```

```
////// STATIC REDUCTIONS.
```

```
sort Exp;
```

```
[[ 0 + ⟨Exp#⟩ ]→# ;
[[ ⟨Exp#⟩ + 0 ]→# ;
[[ ⟨Exp#⟩ - 0 ]→# ;
[[ 1 * ⟨Exp#⟩ ]→# ;
[[ ⟨Exp#⟩ * 1 ]→# ;
```

```
[[ 0 * ⟨Exp#⟩ ]→[[0] ;
[[ ⟨Exp#⟩ * 0 ]→[[0] ;
```

```
[[ ⟨Exp#1⟩ * (1 / ⟨Exp#2⟩) ]→[[⟨Exp#1⟩ / ⟨Exp#2⟩ ]];
```

```
[[ [x ↦ ⟨Exp#[x]⟩] ⟨Exp#2⟩ ]→#[#2] ;
```

```
////// MISCELLANEOUS
```

```
space [ \t\n\r ] | "/" [ ^\n\r]* | nested "/" "*" "*" / ;
```

```
}
```

## C Common Errors

### C.1 Error (HACS syntax).

```
Exception in thread "main" java.lang.RuntimeException: net.sf.crsx.CRSEException:
Encountered " "." " " at line 35, column 6.
Was expecting one of:
    <MT_Repeat> ...
    "%Repeat" ...
    <MT_Attributes> ...
```

Indicates a simple syntax errors in the *.hx* file.

### C.2 Error (user syntax).

```
Exception in thread "main" java.lang.RuntimeException:
net.sf.crsx.CRSEException: net.sf.crsx.parser.ParseException:
mycompiler.crs: Parse error in embedded myDecSome term at line 867, column 42:
[[ $TA_Let2b ⟨Dec (#d)⟩{ ⟨DecSome (#ds)⟩} ] at line 867, column 42
Encountered " "\u27e9" "\u27e8Dec (#d)\u27e9 "" at line 867, column 53
...
```

This indicates a concrete syntax error in some parsed syntax—inside  $\llbracket \dots \rrbracket$ —in the *.hx* file. The offending fragment is given in double angles in the message. Check that it is correctly entered in the HACS specification in a way that corresponds to a syntax production. Note that the line/column numbers refer to the generated *build/...Rules.crs* file, which is not immediately helpful (this is a known bug). In error messages a sort is typically referenced as a lower case prefix followed by the sort name—here *myDecSome* indicates that the problem is with parsing the *DecSome* sort of the *My* parser.

### C.3 Error (JavaCC noise).

```
Java Compiler Compiler Version ??._? (Parser Generator)
(type "javacc" with no arguments for help)
Reading from file net/sf/crsx/samples/gentle/FirstParser.jj . . .
Warning: Line 769, Column 51: Non-ASCII characters used in regular expression.
Please make sure you use the correct Reader when you create the parser,
one that can handle your character set.
File "TokenMgrError.java" does not exist. Will create one.
File "ParseException.java" does not exist. Will create one.
File "Token.java" does not exist. Will create one.
File "SimpleCharStream.java" does not exist. Will create one.
Parser generated with 0 errors and 1 warnings.
Note: net/sf/crsx/samples/gentle/FirstParser.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
```

These are “normal” messages from JavaCC.

### C.4 Error (missing library).

```
gcc -std=c99 -g -c -o crsx_scan.o crsx_scan.c
crsx.c:11:30: fatal error: unicode/umachine.h: No such file or directory
```

The HACS tools only use one library in C: ICU. You should get the *libicu-dev* package (or similar) for your system.

### C.5 Error (meta-variable mistake).

```
Error in rule Tiger-Ty2222222222111_9148-1: contractum uses undefined meta-variable (#es)
Errors prevent normalization.
make: *** [pr3.crs-installed] Error 1
```

A rule uses the metavariable `#es` in the replacement without defining it in the corresponding pattern.

### C.6 Error.

```
/home/krisrose/Desktop/teaching/.../hacs/cookmain PG pr3.hxt > pr3.pg
cookmain: crsx.c:528: bufferEnd: Assertion
'(((childTerm)->descriptor == ((void *)0)) ? 0 :
(childTerm)->descriptor->arity) == bufferTop(buffer)->index' failed.
/bin/sh: line 1: 14278 Aborted
(core dumped) /home/krisrose/Desktop/teaching/.../hacs/cookmain PG pr3.hxt > pr3.pg
```

This indicates an arity error: a raw term in the *.hx* file does not have the right number of arguments.

## References

- [1] Alfred V. Aho, Monica S. Lam, Ravi Sethi, , and Jeffrey D. Ullman. *Compilers: Principles, Techniques and Tools*. Pearson Education, Inc, 2006.
- [2] Jan Willem Klop, Vincent van Oostrom, and Femke van Raamsdonk. Combinatory reduction systems: Introduction and survey. *Theor. Computer Science*, 121:279–308, 1993. doi:10.1016/0304-3975(93)90091-7.
- [3] Sreeni Viswanadha, Sriram Sankar, et al. *Java Compiler Compiler (JavaCC) - The Java Parser Generator*. Sun, 4.0 edition, January 2006. URL: <https://javacc.dev.java.net/>.