
Preprocessors

J. Grosch

DR. JOSEF GROSCH

COCOLAB - DATENVERARBEITUNG

GERMANY

Cocktail

Toolbox for Compiler Construction

Preprocessors

Josef Grosch

Aug. 01, 2006

Document No. 24

Copyright © 2006 Dr. Josef Grosch

Dr. Josef Grosch
CoCoLab - Datenverarbeitung
Breslauer Str. 64c
76139 Karlsruhe
Germany

Phone: +49-721-91537544

Fax: +49-721-91537543

Email: grosch@cocolab.com

Preprocessors

1. Introduction

This manual describes the features and the usage of two kinds of preprocessors contained in the Karlsruhe Toolbox for Compiler Construction. The preprocessors *lpp* and *rpp* derive a parser specification and most of a scanner specification from an attribute grammar. The preprocessors *l2r*, *y2l*, and *r2l* convert input for *lex* and *yacc* into specifications for *rex* and *lark* and vice versa. All preprocessors work as filter programs reading input from *standard input* and writing output to *standard output*. Some preprocessors can read from a file specified as argument, too.

2. Specification of Scanner and Parser with an Attribute Grammar

Writing specifications for the scanner generator *rex* [Groa] and the parser generator *lark* [Grob] directly in the tool specific language is a practicable method. However, it has some disadvantages. Most of the tokens are specified twice and their internal representation or code, respectively, has to be selected and kept consistent, manually. Access to attributes using the yacc-style *\$i* construct (see e.g. [Grob]) is less readable and error-prone in case of changes. The following solution avoids these disadvantages.

Instead of using the tool specific language for *rex* and *lark* directly, a language of higher level is used. It replaces the *\$i* construct by named attributes and describes a complete parser and most of a scanner in one document. Two preprocessors transform such a specification into the languages directly understood by *rex* and *lark*. Fig. 1 shows the data flow using arrows between boxes and circles. Boxes represent files or file types and circles represent tools or preprocessors. The intermediate file named *Scanner.rpp* by default contains the part of the scanner specification that can be extracted from the parser specification. Table 1 gives the meaning of the file types used in Fig. 1 and this manual.

Table 1: Meaning of file types

file type	meaning
.prs	scanner and parser specification (including S-attribution)
.scn	rest of scanner specification
.rpp	intermediate data for completion of scanner in .scn
.rex	scanner specification understood by <i>rex</i>
.lrk	parser specification understood by <i>lark</i>
.ell	input for <i>ell</i> (= input for <i>lark</i> with EBNF constructs [GrV])
.lex	input for <i>lex</i>
.yacc	input for <i>yacc</i>
Source	generated module (or linked from library <i>reuse</i>)
Scanner	generated module
Parser	generated module
Errors	generated module (or linked from library <i>reuse</i>)

The formalism used to describe parsers (.prs) is an extension of the input language for the tools *ast* [Groa] and *ag* [Grob] (see section 2.1.). This leads to a rather uniform set of input languages for most of the tools and simplifies their use. The preprocessor *lpp* transforms a parser specification in

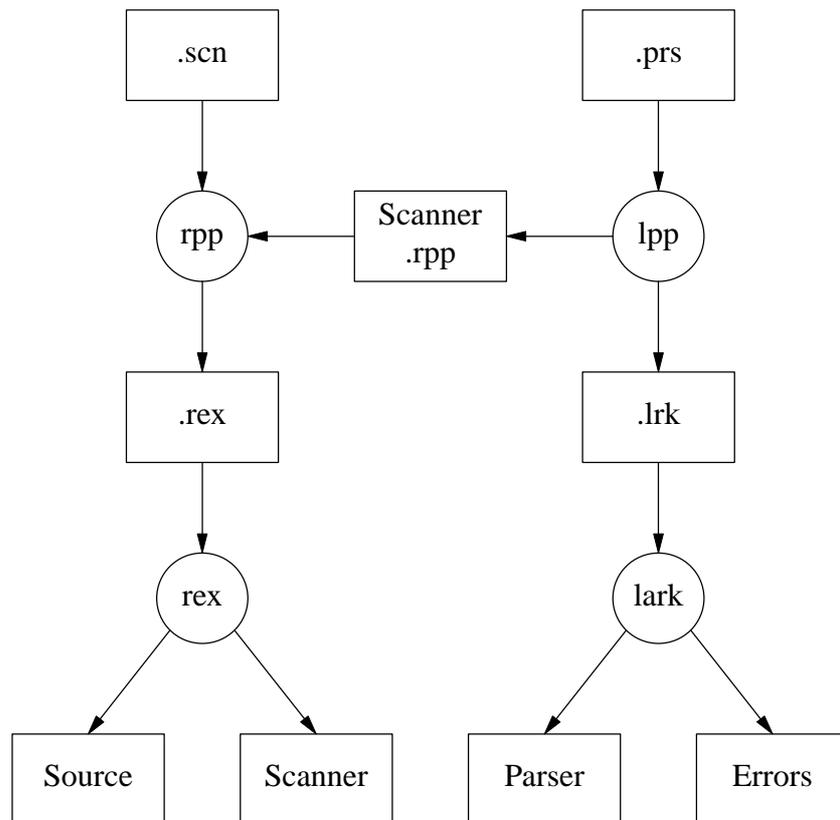


Fig. 1: Data flow during scanner and parser generation

ag notation into one in *lark* notation and extracts most of a scanner specification. The parser specification in *lark* notation is written on a file named `<Parser>.lrk`. `<Parser>` is substituted by the name of the parser module which defaults to *Parser*. The extracted scanner specification is written to a file named `<Scanner>.rpp`. `<Scanner>` is substituted by the name of the scanner module which defaults to *Scanner*. The rest of the scanner specification must be written in the language directly understood by *rex*. It has to contain the part of a scanner specification that can not be derived automatically. This part is usually rather small and comprises the description of user-defined tokens like identifiers and numbers, comments, and the computation of attributes for the tokens. A few insertion points are marked to tell the preprocessor *rpp* where to include the generated parts (see section 2.2.).

2.1. Parser Specification

The input language of *ast* and *ag* can be used to generate a parser. The details of this language can be found in the manuals for these tools [Groa, Grob]. The reader should be familiar with these documents because the current manual describes primarily the extensions necessary for parser generation.

The language can describe concrete as well as abstract syntaxes. Nonterminal, terminal, and abstract symbols or node types are distinguished by the definition characters `'=`', `':'`, and `':='`, respectively, and have to be declared, by default. However, the option `j` of *lpp* allows undeclared terminal symbols and declares them implicitly. In any case, terminal symbols with attributes have to

be declared explicitly.

The following table specifies the keywords of the language. Keywords are not reserved symbols. Keywords can be used as identifiers for grammar symbols in most practical cases. In case of a conflict they should be turned into strings by surrounding them with apostrophes:

BEGIN	CLOSE	DECLARE	DEMAND	END
EVAL	EXPORT	FOR	FUNCTION	GLOBAL
IGNORE	IMPORT	IN	INH	INHERITED
INPUT	LEFT	LOCAL	MODULE	NONE
OUT	OUTPUT	PARSER	PREC	PROPERTY
REMOTE	REV	REVERSE	RIGHT	RULE
SCANNER	SELECT	STACK	START	SUBUNIT
SYN	SYNTHESIZED	THREAD	TREE	VIEW
VIRTUAL	VOID			

The right-hand sides of node types without extensions are interpreted as right-hand sides of grammar rules (see e. g. Assign, Call0, Call, and If in the example below). The children of the right-hand side form a sequence of terminal and nonterminal symbols as usual. The names of those node types serve as rule names. If a symbol occurs several times on one right-hand side, it has to be preceded by different selector names (see e. g. the rule named *If*). Attributes in brackets are not interpreted as grammar symbols but as attribute declarations representing values to be evaluated during parsing.

Not every name of a node type is interpreted as nonterminal or terminal symbol. Only those node types that are used (referenced) on some right-hand and the first node type which is regarded as start symbol are treated as grammar symbols. If a node type is defined as nonterminal then all associated extensions become alternative right-hand sides for this nonterminal symbol. If a node type is defined as terminal it remains a terminal symbol. If a node type is not defined and option *j* is not set an error message is issued. If option *j* is set then undefined node types are implicitly defined as terminals.

The grammar need not be reduced. This means it may contain superfluous terminal and nonterminal symbols. Symbols are superfluous if they are not referenced from any rule. Those symbols are simply ignored and reported as a warning.

Example: input of lpp

```

Stat          = <
  Assign      = Addr ':' '=' Expr .
  Calls       = Ident <
    Call0     = .
    Call      = '(' Actuals ')' .
  > .
  If          = IF Expr THEN Then: Stats ELSE Else: Stats 'END' .
> .
Expr         = <
  Plus       = Lop: Expr '+' Rop: Expr .
  Times      = Lop: Expr '*' Rop: Expr .
  '()'       = '(' Expr ')' .
  Addr       = <
    Name     = Ident .
    Index    = Addr '[' Expr ']' .
  > .
> .

```

Example: output of lpp

```

Stat : Addr ':= ' Expr .
Stat : Ident .
Stat : Ident '(' Actuals ')' .
Stat : IF Expr THEN Stats ELSE Stats 'END' .

Expr : Expr '+' Expr .
Expr : Expr '*' Expr .
Expr : '(' Expr ')' .
Expr : Ident .
Expr : Addr '[' Expr ']' .

Addr : Ident .
Addr : Addr '[' Expr ']' .

```

The rule names and the selector names on the right-hand sides disappear (e. g. If). The extension formalism is expanded (e. g. Calls and Addr) – it is not mapped to chain rules. The expansion includes the inheritance of children and attributes (e. g. Calls). All node types that are used somewhere become nonterminal symbols (e. g. Expr and Addr).

The rule names of non-referenced node types may be omitted. They are necessary for example if modules are used in order to refer from an attribute computation to a grammar rule.

Example without rule names (Modula):

```

Stat = -> [Tree: tTree] <
      = Addr ':= ' Expr      { Tree := mAssign (Addr:Tree, Expr:Tree); } .
      = Ident '(' Actuals ')' { Tree := mCall (Ident:Ident, Actuals:Tree); } .
> .

Ident : [Ident: tIdent] .

```

Example with rule names (Modula):

```

Stat      = <
  Assign = Addr ':= ' Expr .
  Call   = Ident '(' Actuals ')' .
> .

Ident      : [Ident: tIdent] .

MODULE Tree

DECLARE Stat = -> [Tree: tTree] .

Assign      = { Tree := mAssign (Addr:Tree, Expr:Tree); } .
Call        = { Tree := mCall (Ident:Ident, Actuals:Tree); } .

END Tree

```

Attribute declarations and attribute computations are written exactly as for *ast* and *ag*. Attribute computations may be placed everywhere within right-hand sides, not only at the end. These computations are executed from left to right as parsing proceeds. They may only make use of those attributes that have already been computed before or to the left, respectively. The extension or inheritance mechanism for right-hand sides, attribute declarations, and attribute computations is available. The default computations (copy rules) for synthesized attributes are available, too. A specification may be separated into several modules. There could for example be modules for the concrete syntax, for the attribute declarations, and for the attribute computations. It is even possible

to distribute the mentioned kinds of information into several modules.

Attribute declarations and attribute computations with named attributes replace the explicit declaration of the type *tParsAttribute* and the *\$i* construct. The attribute declarations are transformed automatically into a declaration of the type *tParsAttribute*. The attribute computations in the new style are more readable and robust against changes. The given attribute computations are checked for completeness and whether the resulting attribute grammar obeys the SAG property. Only attribute grammars with synthesized attributes can be evaluated during parsing.

Every terminal has two predefined attributes. An attribute named *Position* of type *tPosition* and an attribute named *Scan* of type *tScanAttribute*. The type *tPosition* is a user defined struct (record) type which has to contain at least the members *Line* and *Column*. The attribute *Position* describes the source coordinates of every token and it is computed automatically by *rex* generated scanners. The values of all other attributes of terminals have to be supplied by the scanner by user specified computations. Still, attribute computations for those attributes (except *Position*) have to be specified in the parser specification, too. They are used to generate the procedure *ErrorAttribute*. This procedure is called by the parser whenever tokens are inserted in order to repair syntax errors. The procedure and therefore the mentioned attribute computations have to supply default values for the attributes of tokens.

The attribute *Scan* gives access to the complete record that stores the attributes of the current token. This attribute can be used for tricks such as accessing "undeclared" attributes:

Example:

Parser specification:

```
expr      = <
           = BEGIN { tree := mident (BEGIN:Scan.name.ident); } .
           = name  { tree := mident (name:ident); } .
> .
name      : [ident: tIdent] { ident := Noident; } .
```

Scanner specification:

```
BEGIN    : { Attribute.name.ident = MakeIdent (TokenPtr, TokenLength);
           return BEGIN; }
name     : { Attribute.name.ident = MakeIdent (TokenPtr, TokenLength);
           return name; }
```

In the above example the attribute *ident* is declared for the terminal *name*, only. However, it is used for the terminal *BEGIN* as well. If you are using such tricks you have to know what you are doing: they depend on the characteristics of the target language.

The right-hand side of a node type or a grammar rule, respectively, may contain actions to be executed during parsing. These actions may be placed at the end of the right-hand side or anywhere between the right-hand side elements. In any case these actions are executed from left to right according to the progress of parsing. The syntax of the actions is the one defined for the attribute computations of *ag*. It is assumed that most of the actions will compute attributes. Actions which are not attribute computations are possible but have to be written as some kind of CHECK statement:

```
{ attribute := expression ; }
{ => statement ; }
{ => { statement_sequence } ; }
```

Besides the normal actions, which should be called more precisely conditional actions, there are unconditional actions. See the user's manual of *lark* [Grob] for the exact meaning of these kinds of actions. Unconditional actions are enclosed in a pair of '[' and ']':

```
{[ => statement ; ]}
{[ => { statement_sequence } ; ]}
```

The following extensions of the language of *ast* and *ag* are used for parser generation, only. A grammar may be optionally headed by names for the modules to be generated:

```
SCANNER Name PARSER Name
```

The first identifier specifies the module name of the scanner to be used by the parser. The second identifier specifies a name which is used to derive the names of the parsing module, the parsing routine, the parse tables, etc. The format of the name should be compatible with the target language. For Java, names may include a package name. If the desired name collides with a keyword it may be escaped by quoting, as in

```
SCANNER "TREE"
PARSER "package.VIEW" // Java only
```

If the names are missing they default to *Scanner* and *Parser*. In this document we refer to these names by <Scanner> and <Parser>. The parser name may be followed by a set of target code sections which is copied unchecked and unchanged to the input of the parser generator or the parser module, respectively:

```
SCANNER Name
PARSER Name
IMPORT { target_code }
EXPORT { target_code }
GLOBAL { target_code }
LOCAL { target_code }
BEGIN { target_code }
CLOSE { target_code }
```

The precedence and associativity for operator tokens can be specified after the keyword *PREC* using *LEFT*, *RIGHT*, and *NONE* for left-, right-, and no associativity. Each group headed by one of the latter three keywords introduces a new group of tokens with increasing precedence. The precedence and associativity information is copied unchanged to the parser generator.

Example:

```
PREC LEFT MONOP
      NONE SEQ
      LEFT '+' '-'
      LEFT '*' '/' MOD
```

The precedence and associativity information is usually propagated implicitly to the grammar rules by taking it from the right-most token in a rule. Rules without an operator token can get the precedence and associativity from an operator token by adding a *PREC* clause at the end of a right-hand side.

Example:

```
Expr = <
    = Expr Expr PREC SEQ . /* explicit prec. + assoc. of SEQ */
    = '-' Expr PREC MONOP . /* overwrite prec. + assoc. of '-' by MONOP */
    = Expr '+' Expr . /* implicit prec. + assoc. of '+' */
    = Expr '-' Expr . /* implicit prec. + assoc. of '-' */
> .
```

Start symbols may be listed after the keyword `START`. By default the first node type declared is the start symbol.

Syntactic and semantic predicates can be used in the same way as in the direct input language of the parser generator *lark*. They can be placed at the end of rules or within the right-hand side:

Example:

```
Y = <
    = A ? '+' . /* syntactic terminal predicate */
    = B ? X . /* syntactic nonterminal predicate */
    = C ? - X . /* negated nonterminal predicate */
    = D ? { x == y } . /* semantic predicate at end of rule */
    = E ? { x == y } F . /* semantic predicate within rhs */
> .
```

Tokens or terminal symbols are mapped automatically to integer numbers to be used as internal representation. This mapping can be overwritten by explicitly giving codes for terminals.

Example:

```
Ident : [Ident: tIdent] . /* implicitly coded */
IntConst : 5 [Value: int ] . /* explicitly coded as 5 */
```

Besides an internal representation it is possible to describe an external representation to be used in error messages and a cost value for terminal symbols. Note, an external representation is a string that is always preceded by a comma. See the user's manual of *lark* [Grob] for the exact meaning of these specifications.

Example:

```
Ident : , "noname" [Ident: tIdent] . /* external representation */
IntConst : $ 20 [Value: int ] . /* cost value of 20 */
RealConst: 6 $ 20, "0.0" [Value: float ] . /* all specifications */
```

The attribute declarations for terminals are turned into a declaration of the type *tScanAttribute*. The scanner communicates attribute values of terminals to the parser using a global variable called *Attribute* which is of this type. For a target language other than Java this type is a union type (variant record) with one member (variant) for each terminal with attributes. The names of the terminals are taken for the names of these members (variants). However, this leads to problems if the terminals are named by strings or by keywords of the implementation language. Therefore terminals may have two names. The second one is used as member name in the type *tScanAttribute*. The predefined attribute *Position* mentioned above is always included in this type. Assignments of attribute values in the scanner therefore have to use two selections to describe an attribute:

```
Attribute.<selector name>.<attribute name> = ...
```

If the target language is Java there is a base class `<Scanner>.ScanAttribute` which contains only position information. For each terminal with attributes there is a derived class named by prefixing the terminal with `Xx`. This has a field corresponding to each attribute and a constructor which allows all fields to be set. Assignment of attribute values in the scanner is therefore done by constructing a new attribute value based on the one already constructed by the scanner containing position information:

```
attribute = new Xx<selector name> (attribute, <attribute value>, ...)
```

Example:

definition of terminals including attributes and member selectors in the parser specification

```
Ident          : [Ident: tIdent] .    /* selector name: Ident   */
' := ' sAssign : [Ident: tIdent] .    /* selector name: sAssign */
TRUE sTRUE    : [Ident: tIdent] .    /* selector name: sTRUE   */
'..'         : [Ident: tIdent] .    /* selector name: yy17    */
```

setting of attributes in the scanner specification (C, C++, Modula or Ada)

```
Attribute.Ident.Ident = i
Attribute.sAssign.Ident = i
Attribute.sTRUE.Ident = i
Attribute.yy17.Ident = i
```

setting of attributes in the scanner specification (Java)

```
attribute = new XxIdent (attribute, i)
attribute = new XxsAssign (attribute, i)
attribute = new XxsTRUE (attribute, i)
attribute = new Xx_17 (attribute, i)
```

access of attributes in the parser specification (at node directly)

```
Ident          Position
```

access of attributes in the parser specification (from a child node)

```
Ident:Ident          Ident:Position
' := ':Ident        ' := ':Position
TRUE:Ident           TRUE:Position
' .. ':Ident        ' .. ':Position
```

The preprocessor for the parser specification is implemented by the program *lpp*. It is called by the following command:

```
lpp [ -options ] [ Parser.prs ]
```

The input is read either from the file given as argument or from standard input if the argument is missing. The output is written to a file named `<Parser>.lrk`. The program is controlled by the following options:

- x generate scanner specification for rex
- z generate parser specification for lark
- u generate parser specification for yacc
- v omit actions in the generated parser specification
- j allow undefined node types; define implicitly as terminals
- c generate C source code (default: Modula-2)
- + generate C++ code
- generate Ada code
- J generate Java code
- h print help information
- W suppress warnings
- B allow missing attribute computations in extended node types
- 1 print inserted copy rules
- 2 print inherited attribute computation rules
- 6 generate # line directives
- 7 touch output files only if necessary
- 8 report storage consumption

Appendix 1 of the *ag* manual [Grob] contains the complete formal syntax understood by the program *lpp*. Appendix 2 of this manual shows a parser specification for the example language MiniLAX. It separates context-free grammar and attribute computations in two modules. The attribute computations in the module *Tree* use one attribute also called *Tree* and describe the construction of an abstract syntax tree by calling functions generated by *ast*. The implementation language is C.

2.2. Scanner Specification

The scanner specification has to contain only those parts that can not be extracted automatically from the parser specification. This is as already mentioned above the description of user-defined tokens like identifiers and numbers, comments, and the computation of attributes for the tokens. The formalism to describe this fragmentary scanner specification (.scn) is the input language of rex [Groa]. It may contain three insertion points which instruct the preprocessor *rpp* (rex preprocessor) to include certain generated parts. Moreover, tokens in return statements or expressions can be denoted by the same strings or identifiers as in the parser specification.

```
INSERT tScanAttribute
```

is replaced by the generated declaration of the type *tScanAttribute* and, where the target language requires one, the head or external declaration of the procedure *ErrorAttribute*. This is normally placed in the EXPORT section.

```
INSERT ErrorAttribute
```

is replaced by the body of the generated procedure *ErrorAttribute*. This is usually in the GLOBAL section, but for Java it should be in EXPORT.

The third insertion point lies in the RULE section and has the following syntax (only the brackets [] are used as meta characters and denote optional parts):

```
INSERT RULES [CASE-INSENSITIVE] [[NOT] #<start_states>#] [{ <target_code> }]
```

It is replaced by as many rules as there are tokens extracted automatically from the parser specification. Every rule has the following format:

```
NOT # <start_states> # <token> : { <target_code> return <code>; }
```

The start states including the keyword NOT and the target code are optional and are copied to the generated rule as indicated. If CASE-INSENSITIVE is specified, the regular expressions for the tokens are constructed to match lower as well as upper case letters. Note, only rules for tokens without explicitly declared attributes are constructed automatically.

Within a rule, return (or RETURN) statements are used to report the internal code of a recognized token. The expression of those statements can be any expression of the implementation language or a string or an identifier used in the parser specification. The latter are replaced by their internal representation.

Example:

```
return 5;
return Ident;
return ':'=';
```

The denotation for tokens from the parser specification can be used in arbitrary expressions if it is preceded by the keyword YYCODE. The token, which can be optionally enclosed in parentheses, will be replaced by its internal representation.

Example:

```
return YYCODE Ident;
return YYCODE (':=');
if (condition) Token = YYCODE '+'; else Token = YYCODE '-';
```

The program *rpp* is called by the following command:

```
rpp [ Scanner.rpp ] < Scanner.scn > Scanner.rex
```

The fragmentary scanner specification is read from standard input. The scanner specification extracted from the parser specification is read from a file given as argument. This argument is optional and defaults to *Scanner.rpp*. The scanner specification understood by *rex* is written on standard output. The basename *Scanner* in the command line above is usually substituted by the name of the scanner module.

Appendix 1 contains a scanner specification for the example language MiniLAX. It uses C as implementation language.

3. Conversion of Scanner and Parser Specifications

3.1. Input Languages

The input languages of *rex* and *lark* have been designed to be as readable as possible. Although they contain the same information as inputs for *lex* [Les75] and *yacc* [Joh75] they are syntactically incompatible. Several conversion programs allow the transformation of input for *rex* and *lark* to input for *lex* and *yacc* or vice versa. Fig.2 shows the possible conversions for scanner and parser specifications. Table 2 lists the existing filter programs and the types of their input and output files.

The option '-v' instructs *y2l* and *lpp* to omit the semantic actions in the produced output. The following restrictions or problems are known to exist because they can not be mapped to the target

Table 2: Filters for input conversion

filter	input	output
l2r	.lex	.rex
r2l	.rex	.lex
y2l	.yacc	.lrk
lpp -u	.prs	.yacc
lpp -z	.prs	.lrk
bnf	.ell	.lrk
l2cg	.lrk	.prs

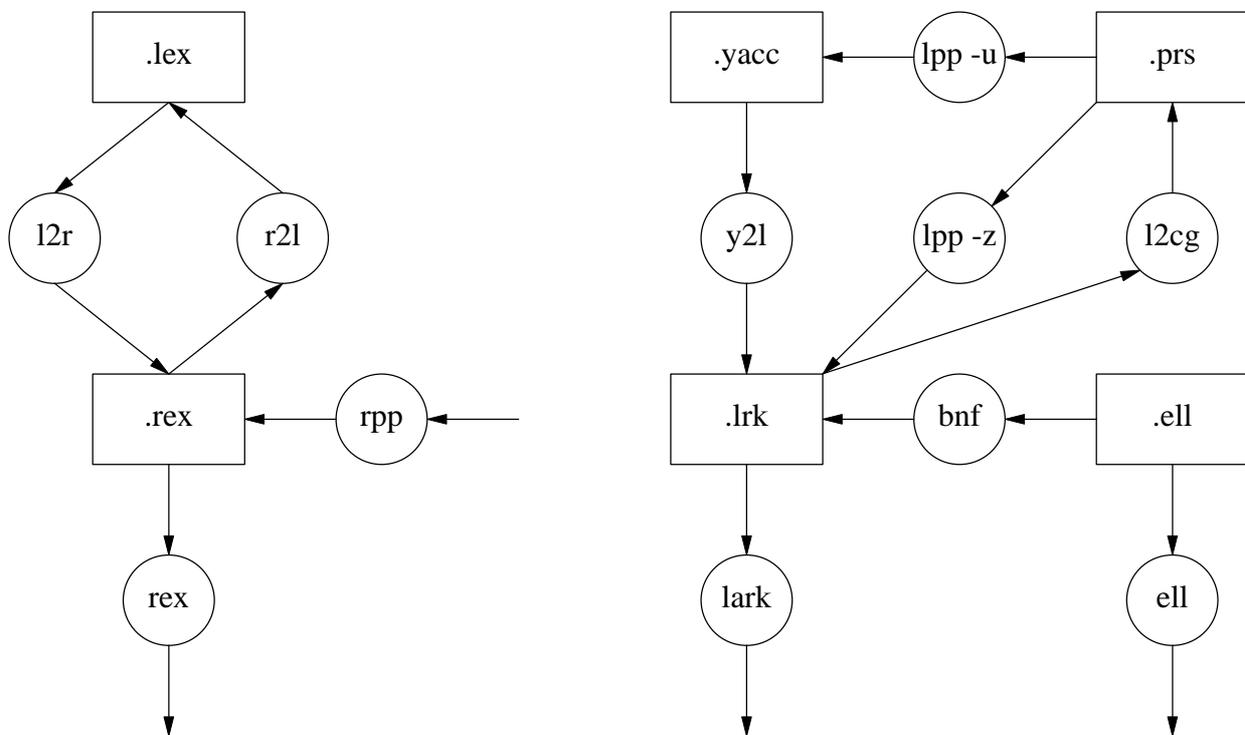


Fig. 2: Conversion programs for scanner and parser specifications

tool:

l2r

- yymore
- REJECT
- yywrap (redirection can be done with *rex*, but differently)
- %T (specification of the character set is not possible)

r2l

- yyPrevious
- EOF (specifies actions to be performed upon reaching end of file)

y2l

- The conversion of token definitions may not be completely automatic.
- If scanning depends on information collected by the parser then parsers generated by *yacc* and *lark* may behave differently because they do not agree upon the order or timing, respectively, of the execution of read (shift) and reduce actions.

bnf

- The attribute computations for *ell* and *lark* are different and are not converted.

3.2. Interfaces

The interfaces of scanners and parsers generated by *lex/yacc* and *rex/lark* are incompatible. The differences are primarily caused by different names for the external (exported) objects. Table 3 lists the interface objects. The interfaces of the scanners and parsers generated by *rex* and *lark* can be switched from the default as listed in Table 3 to an approximation of the *lex* and *yacc* interfaces. This is controlled by *cpp* commands:

```
# define lex_interface
```

in the EXPORT section of a *rex* specification selects a lex-style interface for the scanner.

```
# define lex_interface
```

in the EXPORT section of a *lark* specification tells the parser to use a lex-style interface for the scanner.

```
# define yacc_interface
```

Table 3: Interfaces of generated scanners and parsers

object	yacc/lex	lark/rex
parse routine	int yyparse ();	int Parser ();
stack size	YYMAXDEPTH	yyInitStackSize
attribute type	YYSTYPE	tParsAttribute
global attribute	YYSTYPE yylval;	tScanAttribute Attribute;
position type		typedef struct { short Line, Column; ... } tPosition;
attribute type		typedef struct { tPosition Position; ... } tScanAttribute;
scanner routine	int yylex ();	int GetToken ();
error repair		void ErrorAttribute ();
line number	int yylineno;	<i>member</i> Attribute.Position.Line
token buffer	char yytext [];	
token length	int yyleng;	

in the EXPORT section of a *lark* specification selects a yacc-style interface for the parser.

The output of the preprocessors *l2r* and *y2l* automatically selects lex- and yacc-style interfaces. The following problems are known, currently:

- The output of *l2r* provides the matched string in the array *yytext* to be used in action statements. This is done by calling the procedure *Getword*. However, many actions do not need *yytext*. Deleting superfluous calls of *Getword* will make the scanner significantly faster.
- Access to the line counter *yylineno* has to be replaced by access to

```
Attribute.Position.Line
```

- If both, scanner and parser specification have been converted by *l2r* and *y2l* in order to be fed into *rex* and *lark*, the two preprocessor statements which define *lex_interface* should be deleted in order to select the standard interface. This offers more comfort with respect to the information about the source position.

Acknowledgements

The preprocessor *bnf* has been implemented by Bertram Vielsack. The preprocessors *y2l*, *l2r*, *rpp*, and *lpp* have been implemented by Thomas Müller.

References

- [Groa] J. Grosch, Rex - A Scanner Generator, Cocktail Document No. 5, CoCoLab Germany.
- [Grob] J. Grosch, Lark - An LR(1) Parser Generator With Backtracking, Cocktail Document No. 32, CoCoLab Germany.
- [GrV] J. Grosch and B. Vielsack, The Parser Generators Lalr and Ell, Cocktail Document No. 8, CoCoLab Germany.
- [Groa] J. Grosch, Ast - A Generator for Abstract Syntax Trees, Cocktail Document No. 15, CoCoLab Germany.
- [Grob] J. Grosch, Ag - An Attribute Evaluator Generator, Cocktail Document No. 16, CoCoLab Germany.
- [Joh75] S. C. Johnson, Yacc — Yet Another Compiler-Compiler, Computer Science Technical Report 32, Bell Telephone Laboratories, Murray Hill, NJ, July 1975.
- [Les75] M. E. Lesk, LEX — A Lexical Analyzer Generator, Computing Science Technical Report 39, Bell Telephone Laboratories, Murray Hill, NJ, 1975.

Appendix 1: Scanner Specification for MiniLAX

```

EXPORT {
# include "Idents.h"
# include "Position.h"

INSERT tScanAttribute
}

GLOBAL {
# include <math.h>
# include "Memory.h"
# include "StringM.h"
# include "Idents.h"
# include "Errors.h"

INSERT ErrorAttribute
}

LOCAL { char Word [256]; }

DEFAULT {
MessageI ("illegal character", xxError, Attribute.Position, xxCharacter, TokenPtr);
}

EOF {
    if (yyStartState == Comment)
        Message ("unclosed comment", xxError, Attribute.Position);
}

DEFINE digit = {0-9} .
        letter = {a-z A-Z} .

START Comment

RULE
    "(" :- {yyStart (Comment);}
#Comment# "*" :- {yyStart (STD);}
#Comment# "*" | - {"\t\n"} + :- {}

#STD# digit + : {(void) GetWord (Word);
                Attribute.IntegerConst.Integer = atoi (Word);
                return IntegerConst;}

#STD# digit * "." digit + (E {+|-} ? digit +) ?
        : {(void) GetWord (Word);
            Attribute.RealConst.Real = atof ((char *) Word);
            return RealConst;}

INSERT RULES #STD#

#STD# letter (letter | digit) *
        : {Attribute.Ident.Ident = MakeIdent (TokenPtr, TokenLength);
            return Ident;}

```

Appendix 2: Parser Specification for MiniLAX

PARSER

GLOBAL {# include "Idents.h" }

BEGIN { BeginScanner (); }

```

PREC    LEFT    '<'          /* operator precedence */
        LEFT    '+'
        LEFT    '*'
        LEFT    NOT

```

PROPERTY INPUT

RULE /* concrete syntax */

```

Prog      = PROGRAM Ident ';' 'DECLARE' Decls 'BEGIN' Stats 'END' '.' .
Decl      = <
  Decls1  = Decl .
  Decls2  = Decls ';' Decl .
> .
Decl      = <
  Var     = Ident ':' Type .
  Proc0   = PROCEDURE Ident ';' 'DECLARE' Decls 'BEGIN' Stats 'END' .
  Proc    = PROCEDURE Ident '(' Formals ')' ';'
           'DECLARE' Decls 'BEGIN' Stats 'END' .
> .
Formals   = <
  Formals1 = Formal .
  Formals2 = Formals ';' Formal .
> .
Formal    = <
  Value   = Ident ':' Type .
  Ref     = VAR Ident ':' Type .
> .
Type      = <
  Int     = INTEGER .
  Real    = REAL .
  Bool    = BOOLEAN .
  Array   = ARRAY '[' Lwb: IntegerConst '..' Upb: IntegerConst ']' OF Type .
> .
Stats     = <
  Stats1  = Stat .
  Stats2  = Stats ';' Stat .
> .
Stat      = <
  Assign  = Addr ':=' Expr .
  Call0   = Ident .
  Call    = Ident '(' Actuals ')' .
  If      = IF Expr THEN Then: Stats ELSE Else: Stats 'END' .
  While   = WHILE Expr DO Stats 'END' .
  Read    = READ '(' Addr ')' .
  Write   = WRITE '(' Expr ')' .
> .

```

```

Actuals      = <
  Expr1      = Expr .
  Expr2      = Actuals ',' Expr .
> .
Expr         = <
  Less       = Lop: Expr '<' Rop: Expr .
  Plus       = Lop: Expr '+' Rop: Expr .
  Times      = Lop: Expr '*' Rop: Expr .
  Not        = NOT Expr .
  '()'       = '(' Expr ')' .
  IConst     = IntegerConst .
  RConst     = RealConst .
  False      = FALSE .
  True       = TRUE .
  Addr       = <
    Name     = Ident .
    Index    = Addr '[' Expr ']' .
  > .
> .

                                /* terminals (with attributes) */

Ident        : [Ident: tIdent] { Ident      := NoIdent      ; } .
IntegerConst : [Integer      ] { Integer    := 0              ; } .
RealConst    : [Real : float ] { Real      := 0.0            ; } .

MODULE Tree
                                /* import functions for tree construction */
PARSER GLOBAL {
# include "Tree.h"

tTree nInteger, nReal, nBoolean;
}

BEGIN {
  nInteger = mInteger ();
  nReal    = mReal    ();
  nBoolean = mBoolean ();
}

```

```

                                /* attributes for tree construction          */
DECLARE
  Decls Decl Formals Formal Type Stats Stat Actuals Expr = [Tree: tTree] .

RULE
                                /* tree construction =                      */
                                /* mapping: concrete syntax -> abstract syntax */

Prog   = { => { TreeRoot = mMiniLax (mProc (mNoDecl (), Ident:Ident,
      Ident:Position, mNoFormal (), ReverseTree (Decl:Tree),
      ReverseTree (Stats:Tree)));}; } .
Declsl1 = { Tree := {Decl:Tree->\Decl.Next = mNoDecl (); Tree = Decl:Tree;}; } .
Declsl2 = { Tree := {Decl:Tree->\Decl.Next = Decl:Tree; Tree = Decl:Tree;}; } .
Var     = { Tree := mVar (NoTree, Ident:Ident, Ident:Position, mRef (Type:Tree)); } .
Proc0   = { Tree := mProc (NoTree, Ident:Ident, Ident:Position, mNoFormal (),
      ReverseTree (Decl:Tree), ReverseTree (Stats:Tree)); } .
Proc    = { Tree := mProc (NoTree, Ident:Ident, Ident:Position, ReverseTree
      (Formals:Tree), ReverseTree (Decl:Tree), ReverseTree (Stats:Tree)); } .
Formals1= { Tree := {Formal:Tree->\Formal.Next = mNoFormal ();
      Tree = Formal:Tree;}; } .
Formals2= { Tree := {Formal:Tree->\Formal.Next = Formals:Tree;
      Tree = Formal:Tree;}; } .
Value  = { Tree := mFormal (NoTree, Ident:Ident, Ident:Position,
      mRef (Type:Tree)); } .
Ref    = { Tree := mFormal (NoTree, Ident:Ident, Ident:Position,
      mRef (mRef (Type:Tree))); } .
Int    = { Tree := nInteger; } .
Real   = { Tree := nReal; } .
Bool   = { Tree := nBoolean; } .
Array  = { Tree := mArray (Type:Tree, Lwb:Integer, Upb:Integer, Lwb:Position); } .
Stats1 = { Tree := {Stat:Tree->\Stat.Next = mNoStat (); Tree = Stat:Tree;}; } .
Stats2 = { Tree := {Stat:Tree->\Stat.Next = Stats:Tree; Tree = Stat:Tree;}; } .
Assign = { Tree := mAssign (NoTree, Addr:Tree, Expr:Tree, '=':Position); } .
Call0  = { Tree := mCall (NoTree, mNoActual (Ident:Position), Ident:Ident,
      Ident:Position); } .
Call   = { Tree := mCall (NoTree, ReverseTree (Actuals:Tree), Ident:Ident,
      Ident:Position); } .
If     = { Tree := mIf (NoTree, Expr:Tree, ReverseTree (Then:Tree),
      ReverseTree (Else:Tree)); } .
While  = { Tree := mWhile (NoTree, Expr:Tree, ReverseTree (Stats:Tree)); } .
Read   = { Tree := mRead (NoTree, Addr:Tree); } .
Write  = { Tree := mWrite (NoTree, Expr:Tree); } .
Expr1  = { Tree := mActual (mNoActual (Expr:Tree->\Expr.Pos), Expr:Tree); } .
Expr2  = { Tree := mActual (Actuals:Tree, Expr:Tree); } .
Less   = { Tree := mBinary ('<':Position, Lop:Tree, Rop:Tree, Less); } .
Plus   = { Tree := mBinary ('+':Position, Lop:Tree, Rop:Tree, Plus); } .
Times  = { Tree := mBinary ('*':Position, Lop:Tree, Rop:Tree, Times); } .
Not    = { Tree := mUnary (NOT:Position, Expr:Tree, Not); } .
IConst = { Tree := mIntConst (IntegerConst:Position, IntegerConst:Integer); } .
RConst = { Tree := mRealConst (RealConst:Position, RealConst:Real); } .
False  = { Tree := mBoolConst (FALSE:Position, false); } .
True   = { Tree := mBoolConst (TRUE:Position, true); } .
Name   = { Tree := mIdent (Ident:Position, Ident:Ident); } .
Index  = { Tree := mIndex ('[':Position, Addr:Tree, Expr:Tree); } .

END Tree

```

Contents

1.	Introduction	1
2.	Specification of Scanner and Parser with an Attribute Grammar	1
2.1.	Parser Specification	2
2.2.	Scanner Specification	9
3.	Conversion of Scanner and Parser Specifications	10
3.1.	Input Languages	10
3.2.	Interfaces	12
	References	13
	Appendix 1: Scanner Specification for MiniLAX	14
	Appendix 2: Parser Specification for MiniLAX	15