

The Compiler Generator Coco/R

User Manual

Hanspeter Mössenböck
Johannes Kepler University Linz
Institute of System Software

Coco/R¹ is a compiler generator, which takes an attributed grammar of a source language and generates a scanner and a parser for this language. The scanner works as a deterministic finite automaton. The parser uses recursive descent. LL(1) conflicts can be resolved by a multi-symbol lookahead or by semantic checks. Thus the class of accepted grammars is $LL(k)$ for an arbitrary k .

There are versions of Coco/R for C#, Java, C++, Delphi, Modula-2, Oberon and other languages. This manual describes the versions for C#, Java and C++ from the University of Linz.

Download from: <http://ssw.jku.at/Coco/>

Compiler Generator Coco/R,
Copyright © 1990, 2010 Hanspeter Mössenböck, University of Linz

This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 59 Temple Place - Suite 330, Boston, MA 02111-1307, USA.

As an exception, it is allowed to write an extension of Coco/R that is used as a plugin in non-free software.

If not otherwise stated, any source code generated by Coco/R (other than Coco/R itself) does not fall under the GNU General Public License.

¹ Coco/R stands for *compiler compiler* generating *recursive descent* parsers.

Contents

1. Overview	3
1.1 Sample Production.....	3
1.2 Sample Parsing Method	4
1.3 Summary of Features	4
2. Input Language.....	5
2.1 Vocabulary.....	5
2.2 Overall Structure.....	6
2.3 Scanner Specification	7
2.3.1 Character sets	7
2.3.2 Tokens.....	8
2.3.3 Pragmas.....	9
2.3.4 Comments	10
2.3.5 White space.....	10
2.3.6 Case sensitivity	11
2.4 Parser Specification	11
2.4.1 Productions	11
2.4.2 Semantic Actions	12
2.4.3 Attributes	12
2.4.4 The Symbol ANY	14
2.4.5 LL(1) Conflicts	14
2.4.6 LL(1) Conflict Resolvers	17
2.4.7 Syntax Error Handling	20
2.4.8 Frame Files.....	23
3. User Guide.....	23
3.1 Installation	23
3.2 Options.....	23
3.3 Invocation	24
3.4 Interfaces of the Generated Classes	25
3.4.1 Scanner.....	25
3.4.2 Token	25
3.4.3 Buffer	25
3.4.4 Parser	26
3.4.5 Errors	26
3.5 Main Class of the Compiler.....	27
3.6 Grammar Tests.....	28
4. A Sample Compiler	30
5. Applications of Coco/R	32
6. Acknowledgements	33
A. Syntax of Cocol/R	34
B. Sources of the Sample Compiler.....	35
B.1 Taste.ATG	35
B.2 SymTab.cs (symbol table).....	38
B.3 CodeGen.cs (code generator).....	40
B.4 Taste.cs (main program).....	42

1. Overview

Coco/R is a compiler generator, which takes an attributed grammar of a source language and generates a scanner and a recursive descent parser for this language. The user has to supply a main class that calls the parser as well as semantic classes (e.g. a symbol table handler or a code generator) that are used by semantic actions in the parser. This is shown in Figure 1.

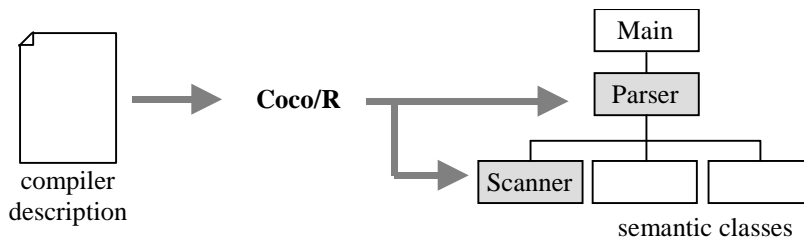


Figure 1 Input and output of Coco/R

1.1 Sample Production

In order to give you an idea of how attributed grammars look like in Coco/R, let us look at a sample production for variable declarations in a Pascal-like language:

```

VarDeclaration<ref int adr>  (. string name; TypeDesc type; .)
= Ident<out name>           (. Obj x = symTab.Enter(name);
                             int n = 1; .)
    { ',' Ident<out name>    (. Obj y = symTab.Enter(name);
                             x.next = y; x = y;
                             n++; .)
    }
    ':' Type<out type>       (. adr += n * typ.size;
                             for (int a = adr; x != null; x = x.next) {
                                 a -= type.size;
                                 x.adr = a;
                             } .)
    ';' .
  
```

The core of this specification is the *EBNF production*

```
VarDeclaration = Ident {',' Ident} ':' Type ';' .
```

It is augmented with attributes and semantic actions. The *attributes* (e.g. <out name>) specify the parameters of the symbols. There are input attributes (e.g. <x, y>) and output attributes (e.g. <out z> or <ref z>). A *semantic action* is a piece of code that is written in the target language of Coco/R (e.g. in C#, Java or C++) and is executed by the generated parser at its position in the production.

1.2 Sample Parsing Method

Every production is translated into a parsing method. The method for `VarDeclaration`, for example, looks like this in C# (code parts originating from attributes or semantic actions are shown in gray):

```
void VarDeclaration(ref int adr) {
    string name; TypeDesc type;
    Ident(out name);
    Obj x = symTab.Enter(name);
    int n = 1;
    while (la.kind == comma) {
        Get();
        Ident(out name);
        Obj y = symTab.Enter(name);
        x.next = y; x = y;
        n++;
    }
    Expect(colon);
    Type(out type);
    adr += n * type.size;
    for (int a = adr; x != null; x = x.next) {
        a -= type.size;
        x.adr = a;
    }
    Expect(semicolon);
}
```

Coco/R also generates a scanner that reads the input stream and returns a stream of tokens to the parser.

1.3 Summary of Features

Scanner

- The scanner is specified by a list of token declarations. Literals (e.g. "if" or "while") do not have to be declared as tokens but can be used directly in the productions of the grammar.
- The scanner is implemented as a deterministic finite automaton (DFA). Therefore the terminal symbols (or tokens) have to be described by a regular EBNF grammar.
- Comments may be nested. One can specify multiple kinds of comments for a language.
- The scanner supports Unicode characters encoded in UTF-8.
- The scanner can be made case-sensitive or case-insensitive.
- The scanner can recognize tokens depending on their context in the input stream.
- The scanner can read from any input stream (not just from a file). However, all input must come from a single stream (no includes).
- The scanner can handle so-called *pragmas*, which are tokens that are not part of the syntax but can occur anywhere in the input stream (e.g. compiler directives or end-of-line characters).
- The user can suppress the generation of a scanner and can provide a hand-written scanner instead.

Parser

- The parser is specified by a set of EBNF productions with attributes and semantic actions. The productions allow for alternatives, repetition and optional parts. Coco/R translates the productions into an efficient recursive descent parser. The parser is reentrant, so multiple instances of it can be active at the same time.
- Nonterminal symbols can have any number of input and output attributes (the Java version allows just one output attribute, which may, however, be an object of a suitable composite class). Terminal symbols do not have explicit attributes, but the tokens returned by the scanner contain information that can be viewed as attributes. All attributes are evaluated during parsing (i.e. the grammar is processed as an L-attributed grammar).
- Semantic actions can be placed anywhere in the grammar (not just at the end of productions). They may contain arbitrary statements or declarations written in the language of the generated parser (e.g. C#, Java or C++).
- The special symbol ANY can be used to denote a set of complementary tokens.
- In principle, the grammar must be LL(1). However, Coco/R can also handle non-LL(1) grammars by using so-called *resolvers* that make a parsing decision based on a multi-symbol lookahead or on semantic information.
- Every production can have its own local variables. In addition to these, one can declare global variables or methods, which are translated into fields and methods of the parser. Semantic actions can also access other objects or methods from user-written classes or from library classes.
- Coco/R checks the grammar for completeness, consistency and non-redundancy. It also reports LL(1) conflicts.
- The error messages printed by the generated parser can be configured to conform to a user-specific format.
- The generated parser and scanner can be specified to belong to a certain namespace (or package).

2. Input Language

This section specifies the compiler description language *Cocol/R* that is used as the input language for Coco/R. A compiler description consists of a set of grammar rules that describe the lexical and syntactical structure of a language as well as its translation to a target language.

2.1 Vocabulary

The basic elements of Cocol/R are identifiers, numbers, strings and character constants, which are defined as follows:

```
ident  = letter {letter | digit}.
number = digit {digit}.
string = '"' {anyButQuote} '"'.
char   = '\' anyButApostrophe '\'
```

Upper case letters are distinct from lower case letters. Strings must not extend across multiple lines. Both strings and character constants may contain the following escape sequences:

<code>\\</code>	backslash	<code>\r</code>	carriage return	<code>\f</code>	form feed
<code>\'</code>	apostrophe	<code>\n</code>	new line	<code>\a</code>	bell
<code>\"</code>	quote	<code>\t</code>	horizontal tab	<code>\b</code>	backspace
<code>\0</code>	null character	<code>\v</code>	vertical tab	<code>\uxxxx</code>	hex char value

The following identifiers are reserved keywords (in the C# version of Cocol/R the identifier `using` is also a keyword, in the Java version the identifier `import`):

ANY	CONTEXT	IGNORE	PRAGMAS	TOKENS
CHARACTERS	END	IGNORECASE	PRODUCTIONS	WEAK
COMMENTS	FROM	NESTED	SYNC	
COMPILER	IF	out	TO	

Comments are enclosed in `/*` and `*/` and may be nested. Alternatively they can start with `//` and go to the end of the line.

EBNF

All syntax descriptions in Cocol/R are written in Extended Backus-Naur Form (EBNF) [Wirth77]. By convention, identifiers starting with a lower case letter denote terminal symbols, identifiers starting with an upper case letter denote nonterminal symbols. Strings denote themselves. The following meta-characters are used:

symbol	meaning	example
<code>=</code>	separates the sides of a production	<code>A = a b c .</code>
<code>.</code>	terminates a production	<code>A = a b c .</code>
<code> </code>	separates alternatives	<code>a b c d e</code> means a b or c or d e
<code>()</code>	groups alternatives	<code>(a b) c</code> means a c or b c
<code>[]</code>	option	<code>[a] b</code> means a b or b
<code>{ }</code>	iteration (0 or more times)	<code>{a} b</code> means b or a b or a a b or ...

Attributes are written between `<` and `>`. Semantic actions are enclosed in `(. and .)`. The operators `+` and `-` are used to form character sets.

2.2 Overall Structure

A Cocol/R compiler description has the following structure:

```
Cocol =
  [Imports]
  "COMPILER" ident
  [GlobalFieldsAndMethods]
  ScannerSpecification
  ParserSpecification
  "END" ident '.'
```

The name after the keyword `COMPILER` is the *grammar name* and must match the name after the keyword `END`. The grammar name also denotes the topmost nonterminal symbol (the *start symbol*). The parser specification must contain a production for this symbol.

Imports. In front of the keyword `COMPILER` one can import namespaces (in C#) or packages (in Java) or include header files (in C++), for example:

```
using System;
using System.Collections;
```

GlobalFieldsAndMethods. After the grammar name one may declare arbitrary fields and methods of the generated parser, for example:

```
int sum;

void Add(int x) {
    sum = sum + x;
}
```

These declarations are written in the language of the generated parser (i.e. in C#, Java or C++) and are not checked by Coco/R. They can be used in the semantic actions of the parser specification. In the C++ version of Coco/R global fields and methods are copied to the header file of the generated parser.

The remaining parts of the compiler description specify the scanner and the parser that are to be generated. They are now described in more detail.

2.3 Scanner Specification

A scanner has to read source text, skip meaningless characters, recognize tokens and pass them to the parser. This is described in a scanner specification, which consists of five optional parts:

```
ScannerSpecification =
  [ "IGNORECASE" ]
  [ "CHARACTERS" {SetDecl} ]
  [ "TOKENS" {TokenDecl} ]
  [ "PRAGMAS" {PragmaDecl} ]
  {CommentDecl}
  {WhiteSpaceDecl}.
```

2.3.1 Character sets

This section allows the user to declare character sets such as letters or digits. Their names can then be used in the other sections of the scanner specification. Coco/R supports the Unicode character set (UTF-8-encoded).

```
SetDecl = ident '=' Set '.'.
Set = BasicSet {('+'|'-') BasicSet}.
BasicSet = string | ident | char [".." char] | "ANY".
```

SetDecl associates a name with a character set. Basic character sets are denoted as:

string	a set consisting of all the characters in the string
ident	a previously declared character set with this name
char	a set containing the character char
char1..char2	the set of all characters from char1 to char2
ANY	the set of all characters in the range 0 .. 65535

Character sets may be formed from basic sets using the operators

+	set union
-	set difference

Examples

```
digit    = "0123456789".      /* the set of all digits */
hexDigit = digit + "ABCDEF".  /* the set of all hexadecimal digits */
letter   = 'A' .. 'Z'.        /* the set of all upper case letters */
eol      = '\r'.              /* the end-of-line character */
noDigit  = ANY - digit.        /* any character that is not a digit */
```

2.3.2 Tokens

This is the main section of the scanner specification, in which the tokens (or terminal symbols) of the language are declared. Tokens may be divided into literals and token classes.

- *Literals* (such as `while` or `>=`) have a fixed representation in the source language. In the grammar they are written as strings (e.g. `"while"` or `">="`) and denote themselves. They don't have to be declared in the tokens section but are implicitly declared at their first use in the productions of the grammar.
- *Token classes* (such as identifiers or numbers) have a certain structure that must be explicitly declared by a regular expression in EBNF. There are usually many instances of a token class (e.g. many different identifiers), which have the same token code, but different lexeme values.

The syntax of token declarations is as follows:

```
TokenDecl  = Symbol ['=' TokenExpr '.'].
TokenExpr  = TokenTerm {'|' TokenTerm}.
TokenTerm  = TokenFactor {TokenFactor} ["CONTEXT" '(' TokenExpr ')'].
TokenFactor = Symbol
            | '(' TokenExpr ')'
            | '[' TokenExpr ']'
            | '{' TokenExpr '}'.
Symbol     = ident | string | char.
```

A token declaration defines the syntax of a terminal symbol by a regular EBNF expression. This expression may contain strings or character constants denoting themselves (e.g. `">="` or `';'`) as well as names of character sets (e.g. `letter`) denoting an arbitrary character from this set. It must not contain other token names, which implies that EBNF expressions in token declarations cannot be recursive.

Examples

```
ident = letter {letter | digit | '_'}.
number = digit {digit}
        | "0x" hexDigit hexDigit hexDigit hexDigit.
float = digit {digit} '.' {digit} ['E' ['+'|-'] digit {digit}].
```

The token declarations need not be LL(1) as can be seen in the declaration of `number`, where both alternatives can start with a `'0'`. Coco/R automatically resolves any ambiguities and generates a deterministic finite scanner automaton.

Tokens may be declared in any order. However, if a token is declared as a literal that matches an instance of a more general token, the literal has to be declared *after* the more general token.

Example

```
ident = letter {letter | digit}.
while = "while".
```

Since the string `"while"` matches both the tokens `while` and `ident`, the declaration of `while` must come after the declaration of `ident`. In principle, literal tokens don't have to be declared in the token declarations at all, but can simply be introduced directly in the productions of the grammar. In some situations, however, it makes sense to declare them explicitly, for example, in order to get a token name for them that can be used in resolver methods (see Section 2.4.6).

Context-dependent tokens. The `CONTEXT` phrase in a `TokenTerm` means that the term is only recognized if its context (i.e. the characters that follow the term in the input stream) matches the `TokenExpr` specified in brackets. Note that the `TokenExpr` is *not* part of the token.

Example

```
number = digit {digit}
        | digit {digit} CONTEXT ("..").
float   = digit {digit} '.' {digit} ['E' ['+'|'-'] digit {digit}].
```

The `CONTEXT` phrase in this example allows the scanner to distinguish between `float` tokens (e.g. 1.23) and integer ranges (e.g. 1..2) that could otherwise not be scanned with a single character lookahead. This works as follows: after having read "1." the scanner still works on both tokens. If the next character is a '.' the characters ".." are pushed back to the input stream and a `number` token with the value 1 is returned to the parser. If the next character is not a '.' the scanner continues with the recognition of a `float` token.

Hand-written scanners. If the right-hand sides of the token declarations are missing no scanner is generated. This gives the user the chance to provide a hand-written scanner, which must conform to the interface described in Section 3.4.1.

Example

```
TOKENS
  ident
  number
  "if"
  "while"
  ...
```

Tokens are assigned numbers in the order of their declaration. The first token gets the number 1, the second the number 2, and so on. The number 0 is reserved for the end-of-file token. The hand-written scanner must return the token numbers according to these conventions. In particular, it must return an end-of-file token if no more input is available.

It is hardly ever necessary to supply a hand-written scanner, because the scanner generated by Coco/R is highly optimized. A user-supplied scanner would be needed, for example, if the scanner were required to process include directives.

2.3.3 Pragmas

Pragmas are tokens that may occur anywhere in the input stream (for example, end-of-line symbols or compiler directives). It would be too tedious to handle all their possible occurrences in the grammar. Therefore they are excluded from the token stream that is passed to the parser. Pragmas are declared like tokens, but they may have a semantic action associated with them that is executed whenever they are recognized by the scanner.

```
PragmaDecl = TokenDecl [SemAction].
SemAction  = "(." ArbitraryStatements ".)".
```

Example

```
PRAGMAS
  option = '$' {letter}.    (. foreach (char ch in la.val)
                             if (ch == 'A') ...
                             else if (ch == 'B') ...
                             ... .)
```

This pragma defines a compiler option that can be written, for example, as \$A. Whenever it occurs in the input stream it is *not* forwarded to the parser but immediately processed by executing its associated semantic action. Note that `la.val` accesses the value of the lookahead token `la`, which is in this case the pragma that was just read (see Section 3.4.4).

2.3.4 Comments

Comments are difficult to specify with regular expressions; nested comments are even impossible to specify that way. This makes it necessary to have a special construct to define their structure.

Comments are declared by specifying their opening and closing brackets. The keyword `NESTED` denotes that they can be nested.

```
CommentDecl = "COMMENTS" "FROM" TokenExpr "TO" TokenExpr ["NESTED"].
```

Comment delimiters must be sequences of 1 or 2 characters, which can be specified as literals or as single-element character sets. They must not be structured (for example with alternatives). It is possible to declare multiple kinds of comments.

Example

```
COMMENTS FROM "/*" TO "*/" NESTED
COMMENTS FROM "//" TO eol
```

Alternatively, if comments cannot be nested one can define them as pragmas, e.g.:

```
CHARACTERS
  other = ANY - '/' - '*'.
PRAGMAS
  comment = "/*" {'/' | other | '*' {'*'} other} '*' {'*'} '/'.
```

This has the advantage that such comments can be processed semantically, for example, by counting them or by processing compiler options within them.

2.3.5 White space

Characters such as blanks, tabulators or end-of-line symbols are usually considered as white space that should be ignored by the scanner. Blanks are ignored by default. If other characters should be ignored as well the user has to specify them in the following way:

```
WhiteSpaceDecl = "IGNORE" Set.
```

Example

```
IGNORE '\t' + '\r' + '\n'
```

2.3.6 Case sensitivity

Some languages such as Pascal are case insensitive. In Pascal, for example, one can write the keyword `while` also as `While` or `WHILE`. By default, Coco/R generates scanners that are case sensitive. If this is not desired, one has to write `IGNORECASE` at the beginning of the scanner specification.

The effect of `IGNORECASE` is that all input to the scanner is treated in a case-insensitive way. The production

```
WhileStatement = "while" '(' Expr ')' Statement.
```

will therefore also recognize while statements that start with `While` or `WHILE`. Similarly, the declaration:

```
TOKENS
float = digit {digit} '.' ['E' ('+'|'-') digit {digit}].
```

will cause the scanner to recognize not only `1.2E2` but also `1.2e2` as a `float` token. However, the original casing of tokens is preserved in the `val` field of every token (see Section 3.4.2) so that the lexical value of tokens such as identifiers and strings is delivered exactly as it was written in the input text.

2.4 Parser Specification

The parser specification is the main part of a compiler description. It contains the productions of an attributed grammar, which specify the syntax of the language to be parsed as well as its translation.

```
ParserSpecification = "PRODUCTIONS" {Production}.
Production = ident [FormalAttributes] [LocalDecl] '=' Expression '.'.
Expression = Term {'|' Term}.
Term = [[Resolver] Factor {Factor}].
Factor = ["WEAK"] Symbol [ActualAttributes]
        | '(' Expression ')'
        | '[' Expression ']'
        | '{' Expression '}'
        | "ANY"
        | "SYNC"
        | SemAction.
Symbol = ident | string | char.
SemAction = "(." ArbitraryStatements ".)".
LocalDecl = SemAction.
FormalAttributes = '<' ArbitraryText '>'.
ActualAttributes = '<' ArbitraryText '>'.
Resolver = "IF" '(' {ANY} ')'.

```

2.4.1 Productions

A production specifies the syntactical structure of a nonterminal symbol. It consists of a left-hand side and a right-hand side which are separated by an equal sign. The left-hand side specifies the name of the nonterminal together with its formal attributes and the local variables of the production. The right-hand side consists of an EBNF expression that specifies the structure of the nonterminal as well as its translation in form of attributes and semantic actions.

The productions may be given in any order. References to as yet undeclared nonterminals are allowed. For every nonterminal there must be exactly one production. In particular, there must be a production for the grammar name, which is the start symbol of the grammar.

2.4.2 Semantic Actions

A semantic action is a piece of code written in the target language of Coco/R (i.e. in C#, Java or C++). It is executed by the generated parser at the position where it has been specified in the grammar. Semantic actions are simply copied to the generated parser without being checked by Coco/R.

A semantic action can also contain the declarations of local variables. Every production has its own set of local variables, which are retained in recursive productions. The optional semantic action on the left-hand side of a production (`LocalDecl`) is intended for such declarations, but variables can also be declared in any other semantic action.

Here is an example that counts the number of identifiers in an identifier list:

```
IdentList =
    ident          (. int n = 1; .)
    {',' ident     (. n++; .)
    }              (. Console.WriteLine("n = " + n); .)
    .
```

As a matter of style, it is good practice to write all syntax parts on the left side and all semantic actions on the right side of a page. This makes a production better readable because the syntax is separated from its processing.

Semantic actions cannot only access local variables but also fields and methods declared at the beginning of the attributed grammar (see Section 2.2) as well as fields and methods of imported classes.

2.4.3 Attributes

Productions are considered as (and are actually translated to) parsing methods. The occurrence of a nonterminal on the right-hand side of a production can be viewed as a call of that nonterminal's parsing method.

Nonterminals may have attributes, which correspond to parameters of the nonterminal's parsing method. There are *input attributes*, which are used to pass values to the production of a nonterminal, and *output attributes*, which are used to return values from the production of a nonterminal to its caller (i.e. to the place where this nonterminal occurs in some other production).

As with parameters, we distinguish between *formal attributes*, which are specified at the nonterminal's declaration on the left-hand side of a production, and *actual attributes*, which are specified at the nonterminal's occurrence on the right-hand side of a production.

Attributes are enclosed in angle brackets (e.g., `< ... >`). If attributes contain the operators '`<`' or '`>`' or generic types like `List<T>` the attribute brackets must be written as `<.>`.

Coco/R checks that nonterminals with attributes are always used with attributes and that nonterminals without attributes are always used without attributes. However, it does not check the correspondence between formal and actual attributes, which is left to the compiler of the target language.

Attributes in C#. A formal attribute looks like a parameter declaration. In C#, output attributes must be preceded by the keyword `out` or `ref`. The following example declares a nonterminal `s` with an input attribute `x` and two output attributes `y` and `z`:

```
S <int x, out int y, ref string z> = ... .
```

An actual attribute looks like an actual parameter. Actual input attributes may be expressions, which are evaluated and assigned to the corresponding formal attributes. In C#, actual output attributes must be preceded by the keywords `out` or `ref`. They are passed by reference like output parameters in C#. Here is an example (`a` and `b` are assumed to be of type `int`, `c` is assumed to be of type `string`):

```
... S <3*a + 1, out b, ref c> ...
```

The production of the nonterminal `s` is translated to the following parsing method:

```
void S(int x, out int y, ref string z) {
    ...
}
```

Attributes in Java. Since Java does not support output parameters, the Java version of Coco/R allows only a single output attribute which is passed to the caller as a return value. However, the return value can be an object of a class that contains multiple values.

If a nonterminal has an output attribute it must be the first attribute. It is denoted by the keyword `out` both in its declaration and in its use. The following example shows a nonterminal `s` with an output attribute `x` and two input attributes `y` and `z` (for compatibility with older versions of Coco/R the symbol `^` can be substituted for the keyword `out`):

```
S<out int x, char y, int z> = ... .
```

This nonterminal is used as follows:

```
... S<out a, 'b', c+3> ...
```

The production of the nonterminal `τ` is translated to the following parsing method:

```
int S(char y, int z) {
    int x;
    ...
    return x;
}
```

Attributes in C++. In the C++ version of Coco/R, input attributes are translated to value parameters and output attributes to reference parameters. The following example declares a nonterminal `s` with an input attribute `x` and an output attribute `y`:

```
S<int x, int &y> = ... .
```

Actual attributes are written like actual parameters in C++, i.e., there is no distinction between value parameters and reference parameters:

```
... S<a+3, b> ...
```

Attributes of terminal symbols. Terminal symbols do not have attributes in Coco/R. For every token, however, the scanner returns the token value (i.e. the token's string representation) as well as the line and column number of the token (see Section 3.4.4). This information can be viewed as output attributes of that token. If users want to access this data they can wrap a token into a nonterminal with the desired attributes, for example:

```

Ident <out string name> =
    ident    (. name = t.val; .) .

Number <out int value> =
    number   (. value = Convert.ToInt32(t.val); .) .

```

The variable `t` is the most recently recognized token. Its field `t.val` holds the textual representation of the token (see Section 3.4.4).

2.4.4 The Symbol ANY

In the productions of the grammar the symbol `ANY` denotes any token that is not an alternative to that `ANY` symbol in the current production. It can be used to conveniently parse structures that contain arbitrary text. The following production, for example, processes an attribute list in Cocol/R and returns the number of characters between the angle brackets:

```

Attributes < out int len> =
    '<'      (. int beg = t.pos + 1; .)
    {ANY}
    '>'      (. len = t.pos - beg; .) .

```

In this example the token `'>'` is an implicit alternative of the `ANY` symbol in curly braces. The meaning is that this `ANY` matches any token except `'>'`. `t.pos` is the source text position of the most recently recognized token (see Section 3.4.4).

Here is another example that counts the number of statements in a block:

```

Block <out int stmts> = (. int n; .)
    '{'                (. stmts = 0; .)
    { ';'              (. stmts++; .)
    | Block<out n>      (. stmts += n; .)
    | ANY
    }
    '}' .

```

In this example the `ANY` matches any token except `';'`, `'{'` and `'}'` which are alternatives of it (`'{'` is a terminal start symbol of `Block`).

2.4.5 LL(1) Conflicts

Recursive descent parsing requires that the grammar of the parsed language is LL(1) (i.e. parsable from **L**eft to **r**ight with **L**eft-canonical derivations and **1** lookahead symbol). This means that at any point in the grammar the parser must be able to decide on the basis of a single lookahead symbol which of several possible alternatives have to be selected. The following production, for example, is not LL(1):

```

Statement = ident '=' Expression ';'
           | ident '(' [ActualParameters] ')' ';'
           | ... .

```

Both alternatives start with the symbol `ident`. When the parser comes to the beginning of a `Statement` and `ident` is the next input token, it cannot distinguish between the two alternatives. However, this production can easily be transformed to

```

Statement = ident ( '=' Expression ';'
                  | '(' [ActualParameters] ')' ';'
                  )
           | ... .

```

where all alternatives start with distinct symbols and the LL(1) conflict has disappeared.

LL(1) conflicts can arise not only from explicit alternatives like those in the example above but also from implicit alternatives that are hidden in optional or iterative EBNF expressions. The following list shows how to check for LL(1) conflicts in these situations (Greek symbols denote arbitrary EBNF expressions such as $a[b]c$; $first(\alpha)$ denotes the set of terminal start symbols of the EBNF expression α ; $follow(A)$ denotes the set of terminal symbols that can follow the nonterminal A in any other production):

▪ **Explicit alternatives**

$A = \alpha | \beta | \gamma$. check that $first(\alpha) \cap first(\beta) = \{\} \wedge first(\alpha) \cap first(\gamma) = \{\} \wedge first(\beta) \cap first(\gamma) = \{\}$.
 $A = (\alpha |) \beta$. check that $first(\alpha) \cap first(\beta) = \{\}$
 $A = (\alpha |)$. check that $first(\alpha) \cap follow(A) = \{\}$

▪ **Options**

$A = [\alpha] \beta$. check that $first(\alpha) \cap first(\beta) = \{\}$
 $A = [\alpha]$. check that $first(\alpha) \cap follow(A) = \{\}$

▪ **Iterations**

$A = \{\alpha\} \beta$. check that $first(\alpha) \cap first(\beta) = \{\}$
 $A = \{\alpha\}$. check that $first(\alpha) \cap follow(A) = \{\}$

It would be very tedious and error-prone to check all these conditions manually for a grammar of a realistic size. Fortunately, Coco/R does that automatically. For example, the grammar

```
A = ( a | B C d ).
B = [ b ] a .
C = c { d } .
```

will result in the following LL(1) warnings:

```
LL1 warning in A: a is start of several alternatives
LL1 warning in C: d is start & successor of deletable structure
```

The first conflict arises because B can start with an a . The second conflict comes from the fact that c may be followed by a d , and so the parser does not know whether it should do another iteration of $\{d\}$ in c or terminate c and continue with the a outside.

Another situation that leads to a conflict is when an expression in curly or square brackets is deletable, e.g.:

```
A = [ B ] a .
B = { b } .
```

If the parser tries to recognize A and sees an a it cannot decide whether to enter the deletable symbol B or to skip $[B]$. Therefore Coco/R prints the warning:

```
LL1 warning in A: contents of [...] or {...} must not be deletable
```

Note that Coco/R reports LL(1) conflicts as warnings, not as errors. Whenever the parser sees two or more alternatives that can start with the same token it always chooses the first one. If this is what the user intends then everything is fine, like in the well-known example of the *dangling else* that occurs in many programming languages:

```
Statement = "if" '(' Expression ')' Statement ["else" Statement]
           | ... .
```

Input for this grammar like

```
if (a > b) if (a > c) max = a; else max = b;
```

is ambiguous: does the "else" belongs to the inner or to the outer if statement? The LL(1) conflict arises because

$$\text{first}(\text{"else" Statement}) \cap \text{follow}(\text{Statement}) = \{\text{"else"}\}$$

However, this is not a big problem, because the parser chooses the first matching alternative, which is the "else" of the inner if statement. This is exactly what we want.

Resolving LL(1) conflicts by grammar transformations

If Coco/R reports an LL(1) conflict the user should try to eliminate it by transforming the grammar as it is shown in the following examples.

Factorization. Most LL(1) conflicts can be resolved by factorization, i.e. by extracting the common parts of conflicting alternatives and moving them to the front. For example, the production

`A = a b c | a b d.`

can be transformed to

`A = a b (c | d).`

Left recursion. Left recursion always represents an LL(1) conflict. In the production

`A = A b | c.`

both alternatives start with `c` (because $\text{first}(A) = \{c\}$). However, left recursion can always be transformed into an iteration, e.g. the previous production becomes

`A = c {b}.`

Hard conflicts. Some LL(1) conflicts cannot be resolved by grammar transformations. Consider the following (simplified) productions from the C# grammar:

```
Expr  = Factor {'+' Factor}.
Factor = '(' ident ')' Factor /* type cast */
      | '(' Expr ')' /* nested expression */
      | ident | number.
```

The conflict arises, because two alternatives of `Factor` start with `'('`. Even worse, `Expr` can also be derived to an `ident`. There is no way to get rid of this conflict by transforming the grammar. The only way to resolve it is to look at the `ident` following the `'('`: if it denotes a type the parser has to select the first alternative otherwise the second one. We will deal with this kind of conflict resolution in Section 2.4.6.

Readability issues. Some grammar transformations can degrade the readability of the grammar. Consider the following example (again taken from a simplified form of the C# grammar):

```
UsingClause = "using" [ident '=' ] Qualident ';'.
Qualident   = ident {'.' ident}.
```

The conflict is in `UsingClause` where both `[ident '=']` and `Qualident` start with `ident`. Although this conflict could be eliminated by transforming the production to

```
UsingClause = "using" ident ( {'.' ident}
                             | '=' Qualident
                             ) ';'.

```

the readability would clearly deteriorate. It is better to resolve this conflict as shown in Section 2.4.6.

Semantic issues. Finally, factorization is sometimes inhibited by the fact that the semantic processing of conflicting alternatives differs, e.g.:

```
A = ident (. x = 1; .) {',' ident (. x++; .) } ':'
    | ident (. Foo(); .) {',' ident (. Bar(); .) } ';'.
```

The common parts of these two alternatives cannot be factored out, because each alternative has its own way to be processed semantically. Again this problem can be solved with the technique explained in Section 2.4.6.

2.4.6 LL(1) Conflict Resolvers

A conflict resolver is a boolean expression that is inserted into the grammar at the beginning of the first of two conflicting alternatives and decides, using a multi-symbol lookahead or a semantic check, whether this alternative matches the actual input. If the resolver yields `true`, the alternative prefixed by the resolver is selected, otherwise the next alternative will be checked. A conflict resolver is written as

```
Resolver = "IF" '(' ... any expression ... ')' .
```

where `any` boolean expression can be written between the parentheses. In most cases this will be a function call that returns `true` or `false`.

Thus we can resolve the LL(1) conflict from Section 2.4.5 in the following way:

```
UsingClause = "using" [IF(IsAlias()) ident '=' ] Qualident ';'.
```

`IsAlias` is a user-defined method that reads two tokens ahead. It returns `true`, if `ident` is followed by `'='`, otherwise it returns `false`.

Conflict resolution by a multi-symbol lookahead

The generated parser remembers the most recently recognized token as well as the current lookahead token in two global variables (see also Section 3.4.4):

```
Token t; // most recently recognized token
Token la; // lookahead token
```

The generated scanner offers a method `Peek()` that can be used to read ahead beyond the lookahead token without removing any tokens from the input stream. When normal parsing resumes the scanner will return these tokens again.

With `Peek()` we can implement `IsAlias()` in the following way:

```
bool IsAlias() {
    Token next = scanner.Peek();
    return la.kind == _ident && next.kind == _eq1;
}
```

The conflict mentioned at the end of Section 2.4.5 can be resolved by the production

```
A = IF(FollowedByColon())
    ident (. x = 1; .) {',' ident (. x++; .) } ':'
    | ident (. Foo(); .) {',' ident (. Bar(); .) } ';'.
```

and the following implementation of the function `FollowedByColon()`:

```
bool FollowedByColon() {
    Token x = la;
    while (x.kind == _comma || x.kind == _ident)
        x = scanner.Peek();
    return x.kind == _colon;
}
```

Token names. For peeking it is convenient to be able to refer to the token numbers by names such as `_ident` or `_comma`. Coco/R generates such names for all tokens declared in the `TOKENS` section of the scanner specification. For example, if the tokens are declared like this:

```
TOKENS
  ident  = letter {letter | digit}.
  number = digit {digit}.
  eql    = '=';
  comma  = ',';
  colon  = ':'.
```

Coco/R will generate the following constant declarations in the parser:

```
const int _EOF = 0;
const int _ident = 1;
const int _number = 2;
const int _eql = 3;
const int _comma = 4;
const int _colon = 5;
```

The token names are preceded by an underscore in order to avoid conflicts with reserved keywords and other identifiers.

Normally the `TOKENS` section will only contain declarations for token classes like `ident` or `number`. However, if the name of a literal token is needed for peeking, it has to be declared there as well. In the productions of the grammar this token can then be referred to either by its name (e.g. `_comma`) or by its literal value (e.g. `' , '`).

Resetting the peek position. The scanner makes sure that a sequence of `Peek()` calls will return the tokens following the lookahead token `la`. In rare situations, however, the user has to reset the peek position manually. Consider the following grammar:

```
A = ( IF (IsFirstAlternative()) ...
      | IF (IsSecondAlternative()) ...
      | ...
    ).
```

Assume that the function `IsFirstAlternative()` starts peeking and finds out that the input does not match the first alternative. So it returns `false` and the parser checks the second alternative. The function `IsSecondAlternative()` starts peeking again, but before that, it should reset the peek position to the first symbol after the lookahead token `la`. This can be done by calling `scanner.ResetPeek()`.

```
bool IsSecondAlternative() {
    scanner.ResetPeek();
    Token x = scanner.Peek(); // returns the first token after the
    ...                      // lookahead token again
}
```

The peek position is reset automatically every time a regular token is recognized by `scanner.Scan()` (see Section 3.4.1).

Translation of conflict resolvers. Coco/R treats resolvers like semantic actions and simply copies them into the generated parser at the position where they appear in the grammar. For example, the production

```
UsingClause = "using" [IF(IsAlias()) ident '='] Qualident ';';
```

is translated into the following parsing method:

```

void UsingClause() {
    Expect(_using);
    if (IsAlias()) {
        Expect(_ident);
        Expect(_eql);
    }
    Qualident();
    Expect(_semicolon);
}

```

Conflict resolution by exploiting semantic information

A conflict resolver can base its decision not only on lookahead tokens but also on any other information. For example it could access a symbol table to find out semantic properties about a token. Consider the following LL(1) conflict between type casts and nested expressions, which can be found in many programming languages:

```

Expr  = Factor {'+' Factor}.
Factor = '(' ident ')' Factor /* type cast */
       | '(' Expr ')'        /* nested expression */
       | ident | number.

```

Since `Expr` can start with an `ident` as well the conflict can be resolved by checking whether this `ident` denotes a type or some other object:

```

Factor = IF (IsCast())
         '(' ident ')' Factor /* type cast */
       | '(' Expr ')'        /* nested expression */
       | ident | number.

```

`IsCast()` looks up `ident` in the symbol table and returns `true`, if it is a type name:

```

bool IsCast() {
    Token x = scanner.Peek();
    if (la.kind == _lpar && x.kind == _ident) {
        object obj = symTab.Find(x.val);
        return obj != null && obj.kind == Type;
    } else return false;
}

```

Placing resolvers correctly

Coco/R checks if resolvers are placed correctly. The following rules must be obeyed:

1. If two alternatives start with the same token, the resolver must be placed in front of the first one. Otherwise it would never be executed because the parser would always choose the first matching alternative. More precisely, a resolver must be placed at the earliest possible point where an LL(1) conflict arises.
2. A resolver may only be placed in front of an alternative that is in conflict with some other alternative. Otherwise it would be illegal.

Here is an example of incorrectly placed resolvers:

```

A =
( a (IF (... ) b) c // misplaced resolver. No LL(1) conflict.
| IF (...) a b      // resolver not evaluated. Place it at first alt.
| IF (...) b        // misplaced resolver. No LL(1) conflict
).

```

Here is how the resolvers should have been placed in this example:

```
A =
( IF (...) a b      // resolves conflict betw. the first two alternatives
| a c
| b
).
```

The following example is also interesting:

```
A =
{ a
| IF (...) b c      // resolver placed incorrectly.
} b.
```

Although the `b` in the second alternative constitutes an LL(1) conflict with the `b` after the iteration, the resolver is placed incorrectly. It should rather be placed at the beginning of the iteration like this:

```
A =
{ IF (AnotherIteration())
( a
| b c
)
} b.
```

The function `AnotherIteration()` could then be implemented as follows:

```
bool AnotherIteration() {
    Token next = scanner.Peek();
    return la.kind == _a ||
           la.kind == _b && next.kind == _c;
}
```

The reason why this resolver is placed incorrectly is that it should be called only once in the parser (namely in the header of the while loop):

```
void A() {
    while (AnotherIteration()) {
        if (la.kind == _a)
            Expect(_a);
        else if (la.kind == _b) {
            Expect(_b); Expect(_c);
        }
    }
    Expect(_b);
}
```

and not both in the while header and at the beginning of the second alternative. Remember, that the resolver must be placed at the earliest possible point where the LL(1) conflict arises.

2.4.7 Syntax Error Handling

If a syntax error is detected during parsing the generated parser reports the error and tries to recover by synchronizing the erroneous input with the grammar. While error messages are generated automatically, the user has to give certain hints in the grammar in order to enable the parser to recover from errors.

Invalid terminal symbols. If a certain terminal symbol was expected but not found in the input the parser just reports that this symbol was expected. For example, if we had a production

```
A = a b c.
```

for which the input was

```
a x c
```

the parser reports

```
-- line ... col ...: b expected
```

Invalid alternative lists. If the lookahead symbol does not match any alternative from a list of expected alternatives in a nonterminal *A* the parser just reports that *A* was invalid. For example, if we had a production

```
A = a (b|c|d) e.
```

for which the input was

```
a x e
```

the parser reports

```
-- line ... col ...: invalid A
```

Obviously, this error message can be improved if we turn the alternative list into a separate nonterminal symbol, i.e.:

```
A = a B e.
B = b|c|d.
```

In this case the error message would be

```
-- line ... col ...: invalid B
```

which is more precise.

Synchronization. After an error was reported the parser continues until it gets to a so-called *synchronization point* where it tries to synchronize the input with the grammar again. Synchronization points have to be specified by the keyword `SYNC`. They are points in the grammar where particularly *safe* tokens are expected, i.e. tokens that hardly occur anywhere else and are unlikely to be mistyped. When the parser reaches a synchronization point it skips all input until a token occurs that is expected at this point.

In many languages good candidates for synchronization points are the beginning of a statement (where keywords like `if`, `while` or `for` are expected) or the beginning of a declaration sequence (where keywords like `public`, `private` or `void` are expected). A semicolon is also a good synchronization point in a statement sequence.

The following production, for example, specifies the beginning of a statement as well as the semicolon after an assignment as synchronization points:

```
Statement =
SYNC
( Designator '=' Expression SYNC ';'
| "if" '(' Expression ')' Statement ["else" Statement]
| "while" '(' Expression ')' Statement
| '{' {Statement} '}'
| ...
).
```

In the generated parser, these synchronization points look as follows (written in pseudo code here):

```

void Statement() {
    while (la.kind ∉ {_EOF, _ident, _if, _while, _lbrace, ...}) {
        Report an error;
        Get next token;
    }
    if (la.kind == _ident) {
        Designator(); Expect(_eq1); Expression();
        while (la.kind ∉ {_EOF, _semicolon}) {
            Report an error;
            Get next token;
        }
    } else if (la.kind == _if) { ...
    } ...
}

```

Note that the end-of-file symbol is always included in the set of synchronization symbols. This guarantees that the synchronization loop terminates at least at the end of the input.

In order to avoid a proliferation of error messages during synchronization, an error is only reported if at least two tokens have been recognized correctly since the last error.

Normally there are only a handful of synchronization points in a grammar for a real programming language. This makes error recovery cheap in Coco/R and does not slow down error-free parsing.

Weak tokens. Error recovery can further be improved by specifying tokens that are "weak" in a certain context. A weak token is a symbol that is often mistyped or missing such as a comma in a parameter list, which is often mistyped as a semicolon. A weak token is preceded by the keyword `WEAK`. When the parser expects a weak token but does not find it in the input stream it adjusts the input to the next token that is either a legal successor of the weak token or a token expected at any synchronization point (symbols expected at synchronization points are considered to be particularly "strong" so that it makes sense to never skip them).

Weak tokens are often separator symbols that occur at the beginning of an iteration. For example, if we have the productions

```

ParameterList = '(' Parameter {WEAK ',' Parameter} ')'.
Parameter = ["ref"|"out"] Type ident.

```

and the parser does not find a `','` or a `')` after the first parameter it reports an error and skips the input until it finds either a legal successor of the weak token (i.e., a legal start of `Parameter`), or a successor of the iteration (i.e. `')`), or any symbol expected at a synchronization point (including the end-of-file symbol). The effect is that the parsing of the parameter list would not be terminated prematurely but would get a chance to synchronize with the start of the next parameter after a possibly mistyped separator symbol.

In order to get good error recovery the user of Coco/R should perform some experiments with erroneous inputs and place `SYNC` and `WEAK` keywords appropriately to recover from the most likely errors.

2.4.8 Frame Files

The scanner and the parser are generated from template files with the names `Scanner.frame` and `Parser.frame`. Those files contain fixed code parts as well as textual markers that denote positions at which grammar-specific parts are inserted by Coco/R. In rare situations advanced users may want to modify the fixed parts of the frame files by which they can influence the behavior of the scanner and the parser to a certain degree. Optionally, a file named `Copyright.frame` can be provided, which will be included at the top of the generated scanner and parser.

3. User Guide

3.1 Installation

Coco/R can be downloaded from <http://ssw.jku.at/Coco/>.

C# and C++ version. Copy the following files to a new directory:

<code>Coco.exe</code>	the executable
<code>Scanner.frame</code>	the frame file from which the scanner is generated
<code>Parser.frame</code>	the frame file from which the parser is generated

Java version. Copy the following files to a new directory:

<code>Coco.jar</code>	an archive containing all classes of Coco/R
<code>Scanner.frame</code>	the frame file from which the scanner is generated
<code>Parser frame</code>	the frame file from which the parser is generated

3.2 Options

Coco/R supports several options that can be provided as command line arguments (see Section 3.3); some of them can also be provided as directives at the beginning of the attributed grammar. If an option is provided both as a command line argument and as a directive in the attributed grammar the command line argument takes precedence.

namespace. The user can specify the namespace (in Java: the package) to which the generated scanner and parser should belong (e.g. `at.jku.ssw.Coco`). If no namespace is specified the generated classes belong to the default namespace. The namespace can be provided as a command line argument or as a directive in the attributed grammar, in which case it has to have the form:

```
$namespace=namespaceName (in Java: $package=packageName)
```

frames. The command line option `frames` can be used to specify the directory that contains the frame files `Scanner.frame`, `Parser.frame` and optionally `Copyright.frame` (see Section 2.4.8). If this option is missing Coco/R expects the frame files to be in the same directory as the attributed grammar.

output directory. The command line option `o` specifies the output directory for the generated scanner and parser. By default, the output directory is the one that contains the attributed grammar.

checkEOF. With the option `checkEOF` the user can specify whether the generated parser should check if the entire input has been consumed after parsing, i.e., if the token after the start symbol of the grammar is an end-of-file token. The user can enable or disable this check by the following directive in the attributed grammar:

```
$checkEOF=true    // enable the end of file check (default)
$checkEOF=false   // disable the end of file check
```

trace. The option `trace` allows the user to specify a string of switches (e.g. `ASX`) that cause internal data structures of Coco/R to be dumped to the file `trace.txt`. The switches are denoted by the following characters:

- A print the states of the scanner automaton
- F print the *first* sets and *follow* sets of all nonterminals
- G print the syntax graph of all productions
- I trace the computation of *first* sets
- J list the ANY and SYNC sets used in error recovery
- P print statistics about the run of Coco/R
- S print the symbol table and the list of declared literals
- X print a cross reference list of all terminals and nonterminals

These switches can be set on in the command line or by a directive in the attributed grammar, which has the form:

```
${letter}
```

For example, the option `$ASX` will cause the states of the automaton, the symbol table and a cross reference list to be printed to the file `trace.txt`.

3.3 Invocation

Coco/R can be invoked from the command line as follows:

```
C# or C++:  Coco fileName [Options]
Java:       java -jar Coco.jar fileName [Options]
```

`fileName` is the name of the file containing the Cocol/R compiler description. As a convention, compiler descriptions have the extension `.ATG` (for attributed grammar).

Options. The following options can be specified:

```
Options =
{
  "-namespace" namespaceName // in Java: "-package" packageName
  "-frames" framesDirectory
  "-trace" traceString
  "-o" outputDirectory
}.
```

A detailed description of these options can be found in Section 3.2.

Output files. Coco/R translates an attributed grammar into the following files:

- `Scanner.cs` (in Java: `Scanner.java`; in C++: `Scanner.h` and `Scanner.cpp`) containing the classes `Scanner`, `Token` and `Buffer`.
- `Parser.cs` (in Java: `Parser.java`; in C++: `Parser.h` and `Parser.cpp`) containing the classes `Parser` and `Errors`.
- `trace.txt` containing trace output (if any).

By default, all files are generated in the directory containing the attributed grammar.

3.4 Interfaces of the Generated Classes

This section specifies the interfaces for the C# version of Coco/R. For Java and C++ the interfaces differ slightly (see the frame files `Scanner.frame` and `Parser.frame`).

3.4.1 Scanner

The generated scanner has the following interface:

```
public class Scanner {
    public Buffer buffer;

    public      Scanner(string sourceFile);
    public      Scanner(Stream s);

    public Token Scan();
    public Token Peek();
    public void  ResetPeek();
}
```

The main class of the compiler (see Section 3.5) has to create a scanner object and pass it either an input stream or the name of a file from where the tokens should be read. The scanner's input buffer is exported in the field `buffer`. It can be used to access the input text at random addresses (see Section 3.4.3).

The method `Scan()` is the actual scanner. The parser calls it whenever it needs the next token. Once the input is exhausted `Scan()` returns the end-of-file token, which has the token number 0. For invalid tokens (caused by illegal token syntax or by invalid characters) `Scan()` returns a special token kind, which normally causes the parser to report an error.

`Peek()` can be used to read one or several tokens ahead without removing them from the input stream. With every call of `Scan()` (i.e. every time a token has been recognized) the peek position is set to the scan position so that the first `Peek()` after a `Scan()` returns the first yet unscanned token. The method `ResetPeek()` can be used to reset the peek position to the scan position after several calls of `Peek()`.

3.4.2 Token

Every token returned by the scanner is an object of the following class:

```
public class Token {
    public int    kind;      // token code (EOF has the code 0)
    public string val;      // token value
    public int    pos;      // token position in the source text
                          // (in bytes starting at 0)
    public int    charPos;  // token position in the source text
                          // (in characters starting at 0)
    public int    line;     // line number (starting at 1)
    public int    col;      // column number (starting at 1)
}
```

3.4.3 Buffer

This is an auxiliary class that is used by the scanner (and possibly by other classes) to read the source stream into a buffer and retrieve portions of it:

```

public class Buffer {
    public const int EOF = char.MaxValue + 1;

    public          Buffer(Stream s);

    public int      Read();
    public int      Peek();
    public int      Pos {get; set;}
    public string    GetString(int beg, int end);
}

```

A buffer is initialized with the source stream. `Read()` returns the next character or 65536 if the input is exhausted. `Peek()` allows the scanner to read characters ahead without consuming them. `Pos` allows the scanner to get or set the reading position, which is initially 0. `GetString(beg, end)` can be used to retrieve the text interval `[beg..end[` from the input stream, where `beg` and `end` are byte positions.

3.4.4 Parser

The generated parser has the following interface:

```

public class Parser {
    public Scanner scanner; // the scanner of this parser
    public Errors  errors;  // the error message stream
    public Token   t;       // most recently recognized token
    public Token   la;      // lookahead token

    public          Parser(Scanner scanner);

    public void      Parse();
    public void      SemErr(string msg);
}

```

The field `t` holds the most recently recognized token. It can be used in semantic actions to access the token value or the token position. The field `la` holds the lookahead token, i.e. the first token after `t`, which has not yet been parsed.

After creating a scanner, the main class of the compiler (see Section 3.5) has to create a parser object and call its method `Parse` in order to start parsing.

The method `SemErr(msg)` can be used to report semantic errors. It calls `errors.SemErr` (see Section 3.4.5) and suppresses error messages that are too close to the position of the previous error, thus avoiding spurious error messages (see Section 2.4.7).

3.4.5 Errors

This class is used to print error messages. Coco/R distinguishes four kinds of errors: syntax errors, semantic errors, warnings and fatal errors. Here is the interface of `Errors`:

```

class Errors {
    public int      count = 0;
    public string    errorStream = Console.Out;
    public string    errMsgFormat = "-- line {0} col {1}: {2}";

    public void      SynErr(int line, int col, int n);
    public void      SemErr(int line, int col, string msg);
    public void      SemErr(string msg);
    public void      Warning(int line, int col, string msg);
    public void      Warning(string msg);
}

```

The field `count` holds the number of errors reported by `SynErr` and `SemErr`. The field `errorStream` denotes the output stream to which error messages are written. By default, this is the console, but the error stream can also be set to any other stream.

Syntax errors are automatically reported by the generated parser, which calls the method `SynErr`. Semantic errors should be reported by calling `Parser.SemErr` which in turn calls `Errors.SemErr`. Warnings can be reported by calling the method `Warning`. Warnings do not increase the error counter.

If `SynErr` and `SemErr` are called with line and column numbers the error message is printed in the format specified by the string `errMsgFormat`, which can be changed by the user to obtain a custom format. The placeholder `{0}` is replaced by the line number, `{1}` is replaced by the column number, and `{2}` is replaced by the error message.

The user can modify the methods `SynErr`, `SemErr` and `Warning` in the file `Parser.frame`. This can be used, for example, to collect all error messages in a data structure instead of writing them to the output stream.

In case of a fatal error from which the compiler cannot recover the user should throw a `FatalError` exception.

```
public class FatalError: Exception { // in Java derived from
    public FatalError(string msg);    // RuntimeException (i.e. unchecked)
}
```

In Coco/R, for example, a `FatalError` is thrown if the frame files cannot be found or are corrupt. The user can catch a `FatalError` in the main method of the compiler and can terminate the compilation.

3.5 Main Class of the Compiler

The main class of a compiler generated with Coco/R has to be provided by the user. It has to create a scanner and a parser object, initiate parsing and possibly report the number of errors detected. In its simplest form it can look like this:

```
public class Compiler {
    public static void Main(string[] arg) {
        Scanner scanner = new Scanner(arg[0]);
        Parser parser = new Parser(scanner);
        parser.Parse();
        Console.WriteLine(parser.errors.count + " errors detected");
    }
}
```

3.6 Grammar Tests

Coco/R checks if the grammar in the compiler specification is well-formed. This includes the following tests:

- **Completeness**

For every nonterminal symbol there must be a production. If a nonterminal x does not have a production Coco/R prints the message

```
No production for X
```

- **Lack of redundancy**

If the grammar contains productions for a nonterminal x that does not occur in any other productions derived from the start symbol Coco/R prints the message

```
X cannot be reached
```

- **Derivability**

If the grammar contains nonterminals that cannot be derived into a sequence of terminals, such as in

```
X = Y ';' .
Y = '(' X ')' .
```

Coco/R prints the messages

```
X cannot be derived to terminals
Y cannot be derived to terminals
```

- **Lack of circularity**

If the grammar contains circular productions, i.e. if nonterminals can be derived into themselves (directly or indirectly) such as in

```
A = [a] B .
B = (C | b) .
C = A {c} .
```

Coco/R prints the messages

```
A --> B
B --> C
C --> A
```

- **Lack of ambiguity**

If two or more tokens are declared so that they can have the same structure and thus cannot be distinguished by the scanner, as in the following example where the input 123 could either be recognized as an integer or as a float:

```
TOKENS
  integer = digit {digit} .
  float   = digit {digit} ['.' {digit}] .
```

Coco/R prints the message

```
Tokens integer and float cannot be distinguished
```

In all these cases the compiler specification is erroneous and no scanner and parser is generated.

Warnings

There are also situations in grammars that—although legal—might lead to problems. In such cases Coco/R prints a warning but nevertheless generates a scanner and a parser. The user should carefully check if these situations are acceptable and, if not, repair the grammar.

▪ Deletable symbols

Sometimes, nonterminals can be derived into the empty string such as in the following grammar:

```
A = B [a].
B = {b}.
```

In such cases Coco/R prints the warnings

```
A deletable
B deletable
```

▪ LL(1) conflicts

If two or more alternatives start with the same token such as in

```
Statement = ident '=' Expression ';'
           | ident '(' Parameters ')' ';'.
```

Coco/R prints the warning

```
LL(1) warning in Statement: ident is start of several alternatives
```

If the start symbols and the successors of a deletable EBNF expression {...} or [...] are not disjoint such as in

```
QualId = [id '.'] id.
IdList = id {' id} [' ,'].
```

Coco/R prints the warnings

```
LL1 warning in QualId: id is start & successor of deletable structure
LL1 warning in IdList: ',' is start & successor of deletable structure
```

The resolution of LL(1) conflicts is discussed in Section 2.4.5.

4. A Sample Compiler

This section shows how to use Coco/R for building a compiler for a tiny programming language called *Taste*. Taste bears some similarities with C# or Java. It has variables of type `int` and `bool` as well as functions without parameters. It allows assignments, procedure calls, `if` and `while` statements. Integers may be read from a file and written to the console, each of them in a single line. It has arithmetic expressions (+,-,*,/) and relational expressions (==,<,>). Here is an example of a Taste program:

```

program Test {
    int i; // global variable
    // compute the sum of 1..i
    void SumUp() {
        int sum;
        sum = 0;
        while (i > 0) { sum = sum + i; i = i - 1; }
        write sum;
    }

    // the program starts here
    void Main() {
        read i;
        while (i > 0) {
            SumUp();
            read i;
        }
    }
}

```

Of course Taste is too restrictive to be used as a real programming language. Its purpose is just to give you a taste of how to write a compiler with Coco/R.

The Taste compiler is a compile-and-go compiler, which means that it reads a source program and translates it into a target program which is executed (i.e. interpreted) immediately after the compilation. In order to run it type

```
Taste Test.TAS
```

The file `Test.TAS` holds the sample program shown above. This file is now compiled and immediately executed. If a program requires input (like `Test.TAS` does) the input file is always `Taste.IN`. For our sample program `Taste.IN` looks like this:

```
3 5 10 0
```

Classes

Figure 2 shows the classes of the compiler.

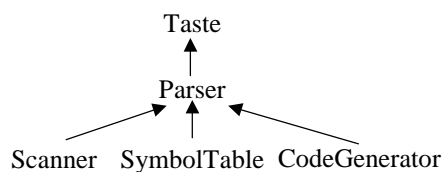


Figure 2 Classes of the Taste compiler

`Taste` is the main class. It creates the scanner and the parser and then calls the parser and the interpreter. The *symbol table* has methods to handle scopes and to store and retrieve object information. The *code generator* has methods to emit instructions. It also contains the interpreter and its data structures. The source code of all classes as well as the attributed grammar `Taste.ATG` can be found in Appendix B.

Target Code

We define an abstract stack machine for the interpretation of `Taste` programs. The compiler translates a source program into instructions of that machine, which are then interpreted. The machine uses the following data structures:

```
char[] code;    // object code (filled by the compiler)
int[]  globals; // data area for global variables
int[]  stack;   // stack with frames for local variables
int    top;     // stack pointer (points to next free stack slot)
int    pc;      // program counter
int    bp;      // base pointer of current frame
```

The architecture of the `Taste` VM is shown in Figure 3.

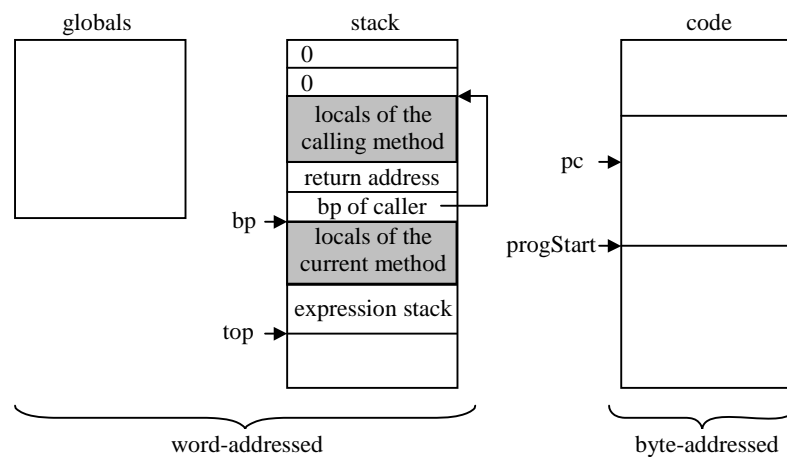


Figure 3: Data structures of the `Taste` VM

Global variables are stored in the word-addressed array `globals` at fixed addresses. *Local variables* are stored in stack frames that are linked with the stack frame of their caller. They are addressed with a word offset relative to the base pointer (`bp`) of the frame. At the end of the topmost stack frame there is the *expression stack* that is used for expression evaluation. After every statement the expression stack is empty.

The *machine code* is stored in the byte-addressed array `code`. The program counter `pc` points to the currently executed instruction. `progStart` is the address of the `Main` method. This is the point where the execution of the program starts.

The machine instructions are described by the following table (the initial values are: `stack[0] = 0; top = 1; bp = 0;`):

CONST	n	Load constant	<code>Push(n);</code>
LOAD	a	Load local variable	<code>Push(stack[bp+a]);</code>
LOADG	a	Load global variable	<code>Push(globals[a]);</code>
STO	a	Store local variable	<code>stack[bp+a]=Pop();</code>
STOG	a	Store global variable	<code>globals[a]=Pop();</code>
ADD		Add	<code>Push(Pop()+Pop());</code>
SUB		Subtract	<code>Push(-Pop()+Pop());</code>
DIV		Divide	<code>x=Pop(); Push(Pop()/x);</code>
MUL		Multiply	<code>Push(Pop()*Pop());</code>

NEG	Negate	Push(-Pop());
EQU	Compare if equal	if (Pop()==Pop()) Push(1); else Push(0);
LSS	Compare if less	if (Pop()>Pop()) Push(1); else Push(0);
GTR	Compare if greater	if (Pop()<Pop()) Push(1); else Push(0);
JMP a	Jump	pc = a;
FJMP a	Jump if false	if (Pop()==0) pc=adr;
READ	Read integer	Push(ReadInt());
WRITE	Write integer	WriteLine(Pop());
CALL a	Call method	Push(pc+2); pc=a;
RET	Return from method	pc = Pop(); if (pc==0) return;
ENTER n	Enter method	Push(bp); bp=top; top+=n;
LEAVE	Leave method	top=bp; bp=Pop();

For example, the method

```
void Foo() {
    int a, b, max;
    read a; read b;
    if (a > b) max = a; else max = b;
    write max;
}
```

is translated into the following code

```
1: ENTER 3
4: READ
5: STO 0
8: READ
9: STO 1
12: LOAD 0
15: LOAD 1
18: GTR
19: FJMP 31
22: LOAD 0
25: STO 2
28: JMP 37
31: LOAD 1
34: STO 2
37: LOAD 2
40: WRITE
41: LEAVE
42: RET
```

Appendix B contains the source code of the following files, which can also be downloaded from <http://ssw.jku.at/Coco/>:

Taste.ATG	the attributed grammar
Taste.cs	the main program
SymTab.cs	the symbol table
CodeGen.cs	the code generator and interpreter

5. Applications of Coco/R

Coco/R can be used not only to write proper compilers, but also to build many kinds of tools that process structured input data. Various people have used Coco/R for the following applications:

- An analyzer for the static complexity of programs. The analyzer evaluates the kind of operators and statements, the nesting of statements and expressions as well as the

use of local and global variables to obtain a measure of the program complexity and an indication if the program is well structured.

- A cross reference generator which lists all occurrences of the objects in a program according to their scope together with information where the objects have been assigned a value and where they have been referenced.
- An pretty printer which uses the structure and the length of statements for proper indentation.
- A program which generates an index for books and reports. The index is generated from a little language that describes page numbers and the keywords occurring on those pages.
- The front end of a syntax oriented editor. A program is translated into a tree representation which is the internal data structure of the editor.
- A program that builds a repository of symbols and their relations in a program. The repository is accessed by a case tool.
- A profiler that inserts counters and timers into the source code of a program and evaluates them after the program has been run.
- A white-box test tool that inserts counters into the source code of a program to find out which paths of the programs have been executed.
- Various compilers for special-purpose scripting languages.
- A log file analyzer that reads machine-generated information and evaluates it.

6. Acknowledgements

The author gratefully acknowledges the help of the following people, who contributed ideas and improvements to Coco/R or ported it to other programming languages:

Pat Terry, Markus Löberbauer, Albrecht Wöß, Csaba Balazs, Frankie Arzu, Peter Rechenberg, Josef Templ and John Gough.

References

- [Möss90] Mössenböck, H.: A Generator for Production Quality Compilers. 3rd Intl. Workshop on Compiler Compilers (CC'90), Schwerin, LNCS 477, Springer-Verlag 1990
- [Terry04] Terry, P.: Compiling with C# and Java. Pearson, 2004.
- [Terry97] Terry, P.: Compilers and Compiler Generators – An Introduction Using C++. International Thomson Computer Press, 1997.
- [Wirth77] Wirth, N.: What Can We Do about the Unnecessary Diversity of Notation for Syntactic Definitions? Communications of the ACM, November 1977
- [WLM03] Wöß A., Löberbauer M., Mössenböck H.: LL(1) Conflict Resolution in a Recursive Descent Compiler Generator, Joint Modular Languages Conference (JMLC'03), Klagenfurt, 2003

A. Syntax of Cocol/R

```

Cocol =
  {ANY}           // using clauses in C#, import clauses in Java,
                  // #include clauses in C++

  "COMPILER" ident
  {ANY}           // global fields and methods
  ScannerSpecification
  Parserspecification
  "END" ident '.'.

```

```

ScannerSpecification =
  [ "IGNORECASE" ]
  [ "CHARACTERS" {SetDecl} ]
  [ "TOKENS" {TokenDecl} ]
  [ "PRAGMAS" {PragmaDecl} ]
  {CommentDecl}
  {WhiteSpaceDecl}.

SetDecl = ident '=' Set.
Set = BasicSet { '(' '+' | '-' ')' BasicSet }.
BasicSet = string | ident | char [ ".." char ] | "ANY".
TokenDecl = Symbol [ '=' TokenExpr '.' ].
TokenExpr = TokenTerm { '|' TokenTerm }.
TokenTerm = TokenFactor { TokenFactor } [ "CONTEXT" '(' TokenExpr ')' ].
TokenFactor = Symbol
  | '(' TokenExpr ')'
  | '[' TokenExpr ']'
  | '{' TokenExpr '}'.

Symbol = ident | string | char.
PragmaDecl = TokenDecl [ SemAction ].
CommentDecl = "COMMENTS" "FROM" TokenExpr "TO" TokenExpr [ "NESTED" ].
WhiteSpaceDecl = "IGNORE" (Set | "CASE").

```

```

ParserSpecification = "PRODUCTIONS" {Production}.
Production = ident [Attributes] [SemAction] '=' Expression '.'.
Expression = Term { '|' Term }.
Term = [[Resolver] Factor {Factor}].
Factor = ["WEAK"] Symbol [Attributes]
  | '(' Expression ')'
  | '[' Expression ']'
  | '{' Expression '}'
  | "ANY"
  | "SYNC"
  | SemAction.
Attributes = '<' {ANY} '>' | "<." {ANY} ">.".
SemAction = "(." {ANY} ".)".
Resolver = "IF" '(' {ANY} ')' .

```

B. Sources of the Sample Compiler

B.1 Taste.ATG

COMPILER Taste

```
const int // types
    undef = 0, integer = 1, boolean = 2;
```

```
const int // object kinds
    var = 0, proc = 1;
```

```
public SymbolTable tab;
public CodeGenerator gen;
```

CHARACTERS

```
letter = 'A'..'Z' + 'a'..'z'.
digit = '0'..'9'.
```

TOKENS

```
ident = letter {letter | digit}.
number = digit {digit}.
```

COMMENTS FROM "/*" TO "*/" NESTED

COMMENTS FROM "//" TO '\n'

IGNORE '\r' + '\n' + '\t'

PRODUCTIONS

```
AddOp<out Op op>
=
    ( '+'
    | '-'
    ).
/*-----*/
Expr<out int type>      (. int type1; Op op; .)
= SimExpr<out type>
    [ RelOp<out op>
        SimExpr<out type1> (. if (type != type1) SemErr("incompatible types");
                             gen.Emit(op); type = boolean; .)
    ].
/*-----*/
Factor<out int type>      (. int n; Obj obj; string name; .)
=
    ( Ident<out name>      (. obj = tab.Find(name); type = obj.type;
                             if (obj.kind == var) {
                                 if (obj.level == 0) gen.Emit(Op.LOADG, obj.adr);
                                 else gen.Emit(Op.LOAD, obj.adr);
                             } else SemErr("variable expected"); .)
    | number               (. n = Convert.ToInt32(t.val);
                             gen.Emit(Op.CONST, n); type = integer; .)
    | '-'
        Factor<out type>    (. if (type != integer) {
                                 SemErr("integer type expected"); type = integer;
                             }
                             gen.Emit(Op.NEG); .)
    | "true"               (. gen.Emit(Op.CONST, 1); type = boolean; .)
    | "false"              (. gen.Emit(Op.CONST, 0); type = boolean; .)
    ).
/*-----*/
Ident<out string name>
= ident      (. name = t.val; .).
/*-----*/
```

```

MulOp<out Op op>
=
    ( '*'
    | '/'
    ).
    (. op = Op.MUL; .)
    (. op = Op.DIV; .)
    ).
/*-----*/
ProcDecl
= "void"
    Ident<out name>
        (. obj = tab.NewObj(name, proc, undef); obj.adr = gen.pc;
        if (name == "Main") gen.progStart = gen.pc;
        tab.OpenScope(); .)

    '(' ' ' ')'
    '{'
    { VarDecl | Stat }
    '}'
        (. gen.Emit(Op.ENTER, 0); adr = gen.pc - 2; .)
        (. gen.Emit(Op.LEAVE); gen.Emit(Op.RET);
        gen.Patch(adr, tab.topScope.nextAdr);
        tab.CloseScope(); .).
/*-----*/
RelOp<out Op op>
=
    ( "=="
    | '<'
    | '>'
    ).
    (. op = Op.EQU; .)
    (. op = Op.LSS; .)
    (. op = Op.GTR; .)
    ).
/*-----*/
SimExpr<out int type>
= Term<out type>
    { AddOp<out op>
        Term<out type>
        (. if (type != integer || type1 != integer)
        SemErr("integer type expected");
        gen.Emit(op); .)
    }.
/*-----*/
Stat
= Ident<out name>
    ( '='
        Expr<out type> ';'
        (. int type; string name; Obj obj;
        int adr, adr2, loopstart; .)
        (. obj = tab.Find(name); .)
        (. if (obj.kind != var) SemErr("cannot assign to procedure"); .)
        (. if (type != obj.type) SemErr("incompatible types");
        if (obj.level == 0) gen.Emit(Op.STOG, obj.adr);
        else gen.Emit(Op.STO, obj.adr); .)
        (. if (obj.kind != proc) SemErr("object is not a procedure");
        gen.Emit(Op.CALL, obj.adr); .)
        )

    | "if"
        '(' Expr<out type> ')'
        (. if (type != boolean) SemErr("boolean type expected");
        gen.Emit(Op.FJMP, 0); adr = gen.pc - 2; .)

        Stat
        [ "else"
            (. gen.Emit(Op.JMP, 0); adr2 = gen.pc - 2;
            gen.Patch(adr, gen.pc);
            adr = adr2; .)

            Stat
        ]
        (. gen.Patch(adr, gen.pc); .)

    | "while"
        '(' Expr<out type> ')'
        (. loopstart = gen.pc; .)
        (. if (type != boolean) SemErr("boolean type expected");
        gen.Emit(Op.FJMP, 0); adr = gen.pc - 2; .)

        Stat
        (. gen.Emit(Op.JMP, loopstart); gen.Patch(adr, gen.pc); .)

    | "read"
        Ident<out name> ';'
        (. obj = tab.Find(name);
        if (obj.type != integer) SemErr("integer type expected");
        gen.Emit(Op.READ);
        if (obj.level == 0) gen.Emit(Op.STOG, obj.adr);
        else gen.Emit(Op.STO, obj.adr); .)

```

```

| "write"
  Expr<out type> ';'      (. if (type != integer) SemErr("integer type expected");
                           gen.Emit(Op.WRITE); .)

| '{' { Stat | VarDecl } '}' .
/*-----*/
Taste                (. string name; .)
= "program"           (. gen.Init(); tab.Init(); .)
  Ident<out name>      (. tab.OpenScope(); .)
  '{'
  { VarDecl | ProcDecl }
  '}'                  (. tab.CloseScope();
                        if (gen.progStart == -1) SemErr("main function never defined");
                        .).
/*-----*/
Term<out int type>    (. int type1; Op op; .)
= Factor<out type>
  { MulOp<out op>
    Factor<out type1>   (. if (type != integer || type1 != integer)
                        SemErr("integer type expected");
                        gen.Emit(op); .)
  }.
/*-----*/
Type<out int type>
=
( "int"               (. type = integer; .)
| "bool"              (. type = boolean; .)
).
/*-----*/
VarDecl              (. string name; int type; .)
= Type<out type>
  Ident<out name>      (. tab.NewObj(name, var, type); .)
  { ',' Ident<out name> (. tab.NewObj(name, var, type); .)
  } ' '.
END Taste.

```

B.2 SymTab.cs (symbol table)

```
using System;

namespace Taste {

public class Obj { // object describing a declared name
    public string name; // name of the object
    public int type; // type of the object (undef for procs)
    public Obj next; // to next object in same scope
    public int kind; // var, proc, scope
    public int adr; // address in memory or start of proc
    public int level; // nesting level; 0=global, 1=local
    public Obj locals; // scopes: to locally declared objects
    public int nextAdr; // scopes: next free address in this scope
}

public class SymbolTable {

    const int // types
        undef = 0, integer = 1, boolean = 2;

    const int // object kinds
        var = 0, proc = 1, scope = 2;

    public int curLevel; // nesting level of current scope
    public Obj undefObj; // object node for erroneous symbols
    public Obj topScope; // topmost procedure scope

    Parser parser;

    // open a new scope and make it the current scope (topScope)
    public void OpenScope () {
        Obj scop = new Obj();
        scop.name = ""; scop.kind = scope;
        scop.locals = null; scop.nextAdr = 0;
        scop.next = topScope; topScope = scop;
        curLevel++;
    }

    // close the current scope
    public void CloseScope () {
        topScope = topScope.next; curLevel--;
    }

    // create a new object node in the current scope
    public Obj NewObj (string name, int kind, int type) {
        Obj p, last, obj = new Obj();
        obj.name = name; obj.kind = kind; obj.type = type;
        obj.level = curLevel;
        p = topScope.locals; last = null;
        while (p != null) {
            if (p.name == name) parser.SemErr("name declared twice");
            last = p; p = p.next;
        }
        if (last == null) topScope.locals = obj; else last.next = obj;
        if (kind == var) obj.adr = topScope.nextAdr++;
        return obj;
    }
}
```

```

// search the name in all open scopes and return its object node
public Obj Find (string name) {
    Obj obj, scope;
    scope = topScope;
    while (scope != null) { // for all scopes
        obj = scope.locals;
        while (obj != null) { // for all objects in this scope
            if (obj.name == name) return obj;
            obj = obj.next;
        }
        scope = scope.next;
    }
    parser.SemErr(name + " is undeclared");
    return undefObj;
}

public SymbolTable (Parser parser) {
    this.parser = parser;
    topScope = null;
    curLevel = -1;
    undefObj = new Obj();
    undefObj.name = "undef"; undefObj.type = undef; undefObj.kind = var;
    undefObj.adr = 0; undefObj.level = 0; undefObj.next = null;
}

} // end SymbolTable

} // end namespace

```

B.3 CodeGen.cs (code generator)

```

using System;
using System.IO;

namespace Taste {

public enum Op { // opcodes
    ADD, SUB, MUL, DIV, EQU, LSS, GIR, NEG,
    LOAD, LOADG, STO, STOG, CONST,
    CALL, RET, ENTER, LEAVE, JMP, FJMP, READ, WRITE
}

public class CodeGenerator {

    string[] opcode =
        {"ADD ", "SUB ", "MUL ", "DIV ", "EQU ", "LSS ", "GIR ", "NEG ",
         "LOAD ", "LOADG", "STO ", "STOG ", "CONST", "CALL ", "RET ", "ENTER",
         "LEAVE", "JMP ", "FJMP ", "READ ", "WRITE"};

    public int progStart; // address of first instruction of main program
    public int pc; // program counter
    byte[] code = new byte[3000];

    // data for Interpret
    int[] globals = new int[100];
    int[] stack = new int[100];
    int top; // top of stack
    int bp; // base pointer

    //----- code generation methods -----

    public void Put(int x) { code[pc++] = (byte)x; }

    public void Emit (Op op) { Put((int)op); }

    public void Emit (Op op, int val) { Emit(op); Put(val>>8); Put(val); }

    public void Patch (int adr, int val) {
        code[adr] = (byte)(val>>8); code[adr+1] = (byte)val;
    }

    public void Decode() {
        int maxPc = pc; pc = 1;
        while (pc < maxPc) {
            Op code = (Op)Next();
            Console.WriteLine("{0,3}: {1} ", pc-1, opcode[(int)code]);
            switch(code) {
                case Op.LOAD: case Op.LOADG: case Op.CONST: case Op.STO: case Op.STOG:
                case Op.CALL: case Op.ENTER: case Op.JMP: case Op.FJMP:
                    Console.WriteLine(Next2()); break;
                case Op.ADD: case Op.SUB: case Op.MUL: case Op.DIV: case Op.NEG:
                case Op.EQU: case Op.LSS: case Op.GIR: case Op.RET: case Op.LEAVE:
                case Op.READ: case Op.WRITE:
                    Console.WriteLine(); break;
            }
        }
    }

    //----- interpreter methods -----

    int Next () {
        return code[pc++];
    }

    int Next2 () {
        int x, y;
        x = (sbyte)code[pc++]; y = code[pc++];
        return (x << 8) + y;
    }
}

```



```

int Int (bool b) {
    if (b) return 1; else return 0;
}

void Push (int val) {
    stack[top++] = val;
}

int Pop() {
    return stack[--top];
}

int ReadInt(FileStream s) {
    int ch, sign, n = 0;
    do {ch = s.ReadByte();} while (!(ch >= '0' && ch <= '9' || ch == '-'));
    if (ch == '-') {sign = -1; ch = s.ReadByte();} else sign = 1;
    while (ch >= '0' && ch <= '9') {
        n = 10 * n + (ch - '0');
        ch = s.ReadByte();
    }
    return n * sign;
}

public void Interpret (string data) {
    int val;
    try {
        FileStream s = new FileStream(data, FileMode.Open);
        Console.WriteLine();
        pc = progStart; stack[0] = 0; top = 1; bp = 0;
        for (;;) {
            switch ((Op)Next()) {
                case Op.CONST: Push(Next2()); break;
                case Op.LOAD:  Push(stack[bp+Next2()]); break;
                case Op.LOADG: Push(globals[Next2()]); break;
                case Op.STO:   stack[bp+Next2()] = Pop(); break;
                case Op.STOG:  globals[Next2()] = Pop(); break;
                case Op.ADD:   Push(Pop()+Pop()); break;
                case Op.SUB:   Push(-Pop()+Pop()); break;
                case Op.DIV:   val = Pop(); Push(Pop()/val); break;
                case Op.MUL:   Push(Pop()*Pop()); break;
                case Op.NEG:   Push(-Pop()); break;
                case Op.EQU:   Push(Int(Pop()==Pop())); break;
                case Op.LSS:   Push(Int(Pop()>Pop())); break;
                case Op.GTR:   Push(Int(Pop()<Pop())); break;
                case Op.JMP:   pc = Next2(); break;
                case Op.FJMP:  val = Next2(); if (Pop()==0) pc = val; break;
                case Op.READ:  val = ReadInt(s); Push(val); break;
                case Op.WRITE: Console.WriteLine(Pop()); break;
                case Op.CALL:  Push(pc+2); pc = Next2(); break;
                case Op.RET:   pc = Pop(); if (pc == 0) return; break;
                case Op.ENTER: Push(bp); bp = top; top = top + Next2(); break;
                case Op.LEAVE: top = bp; bp = Pop(); break;
                default:      throw new Exception("illegal opcode");
            }
        }
    } catch (IOException) {
        Console.WriteLine("--- Error accessing file {0}", data);
        System.Environment.Exit(0);
    }
}

public CodeGenerator () { pc = 1; progStart = -1; }
} // end CodeGen
} // end namespace

```

B.4 Taste.cs (main program)

```
using System;

namespace Taste {

class Taste {

    public static void Main (string[] arg) {
        if (arg.Length > 0) {
            Scanner scanner = new Scanner(arg[0]);
            Parser parser = new Parser(scanner);
            parser.tab = new SymbolTable(parser);
            parser.gen = new CodeGenerator();
            parser.Parse();
            if (parser.errors.count == 0) {
                parser.gen.Decode();
                parser.gen.Interpret("Taste.IN");
            }
        } else {
            Console.WriteLine("-- No source file specified");
        }
    }

}

} // end namespace
```