

Yacc: A Syntactic Analysers Generator

Compiler-Construction Tools

The compiler writer uses specialised tools (in addition to those normally used for software development) that produce components that can easily be integrated in the compiler and help implement various phases of a compiler.

- Scanner generators
- Parser generators
- Syntax-directed translation
- Code-generator generators
- Data-flow analysis: key part of code optimisation

Relationship between Parser and Scanner

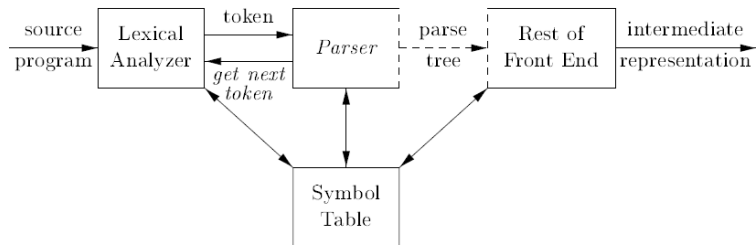


Figure 4.1: Position of parser in compiler model

Parser guides scanner by asking tokens one by one. Interaction with the symbol table is crucial.

Lex and Yacc

There exist various tools for generating parsers.

We will discuss **Yacc**, the companion of Lex: the programs they generate can share variables and procedures and therefore can be compiled jointly.

Lex & YACC are known to GNU/Linux users as **Flex & Bison**, where Flex is a Lex implementation by Vern Paxson and Bison the GNU version of YACC.

The program `lex.yy.c` produced by Lex is called by Yacc via the function `yylex()` defined in Yacc.

Typically, the Lex actions are not visible on the standard output; rather Lex return values and control to the caller (i.e. the parser). Function `yylex()` returns the token names, while the token values are shared by means of global variables such as `yylval` of type `int`.

Yacc

LALR parser generator whose first version is due to S.C. Johnson in the early 1970s.

To construct a translator, Yacc operates as follows:

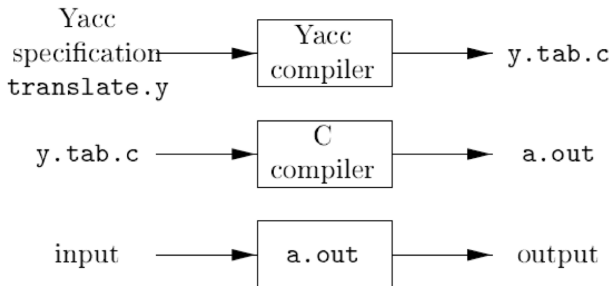


Figure 4.57: Creating an input/output translator with Yacc

Yacc Specifications

Structure of a Yacc program:

```
%{ C declarations %}  
declarations  
%%  
translation rules  
%%  
supporting C routines
```

A production of the form

$$\textit{nonterm} \rightarrow \textit{corpo}_1 \mid \cdots \mid \textit{corpo}_k$$

is expressed in Yacc by the **rules**:

```
nonterm : body1 {semantic action1}  
          ...  
          | bodyk {semantic actionk}  
          ;
```

Declarations

The declaration part of a Yacc program has two (both optional) sections:

- ordinary C declarations delimited by `%{` and `%}` (e.g. standard header files)
- declarations of grammar tokens of the form

```
%token DIGIT
```

These declarations are made available to the analyzer generated by Lex when `lex.yy.c` is compiled together with the Yacc output.

Supporting C-Routines

The third part of a Yacc specification contains procedures in C:

- The **lexical analyzer** is provided as `yylex()`. Usually this is produced by Lex.
- Other procedures such as **error recovery routines**.

If the pair **(token-name, attribute-value)** is produced by `yylex()`, then `token-name` must be declared in the first section of the Yacc specification, while `attribute-value` is communicated to the parser through the variable `yylval` (defined by Yacc).

Example

Construct a simple desk calculator that reads an arithmetic expression, evaluates it and then prints its numeric value.

Start with the following grammar for arithmetic expressions:

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid \mathbf{digit}$$

The token **digit** is a single digit between 0 and 9.

Example ctd.

```
%{
#include <ctype.h>
%}

%token DIGIT

%%
line   : expr '\n'          { printf("%d\n", $1); }
      ;
expr   : expr '+' term     { $$ = $1 + $3; }
      | term
      ;
term   : term '*' factor   { $$ = $1 * $3; }
      | factor
      ;
factor : '(' expr ')'      { $$ = $2; }
      | DIGIT
      ;

%%
yylex() {
    int c;
    c = getchar();
    if (isdigit(c)) {
        yylval = c-'0';
        return DIGIT;
    }
    return c;
}
```

Creating Yacc analyzers

The Lex library [11](#) provides the driver program `yylex()` to Yacc. When using Lex and Yacc together, we can replace the routine `yylex()` in the third part of the Yacc specification by the statement

```
#include "lex.yy.c"
```

Since the Lex output file is compiled as part of the Yacc output file `y.tab.c`, the lexical analyzer can then have access to Yacc's names for tokens.