

Name: Gaurav Kedia

Roll No: 39

Practical No. 2

Theory:

YACC:

Yacc (Yet Another Compiler-Compiler) is a computer program for the Unix operating system developed by Stephen C. Johnson. It is a Look Ahead Left-to-Right (LALR) parser generator, generating a parser, the part of a compiler that tries to make syntactic sense of the source code, specifically a LALR parser, based on an analytic grammar written in a notation similar to Backus–NaurForm (BNF). Yacc is supplied as a standard utility on BSD and AT&T Unix. GNU- based Linux distributions include Bison, a forward-compatible Yacc replacement.

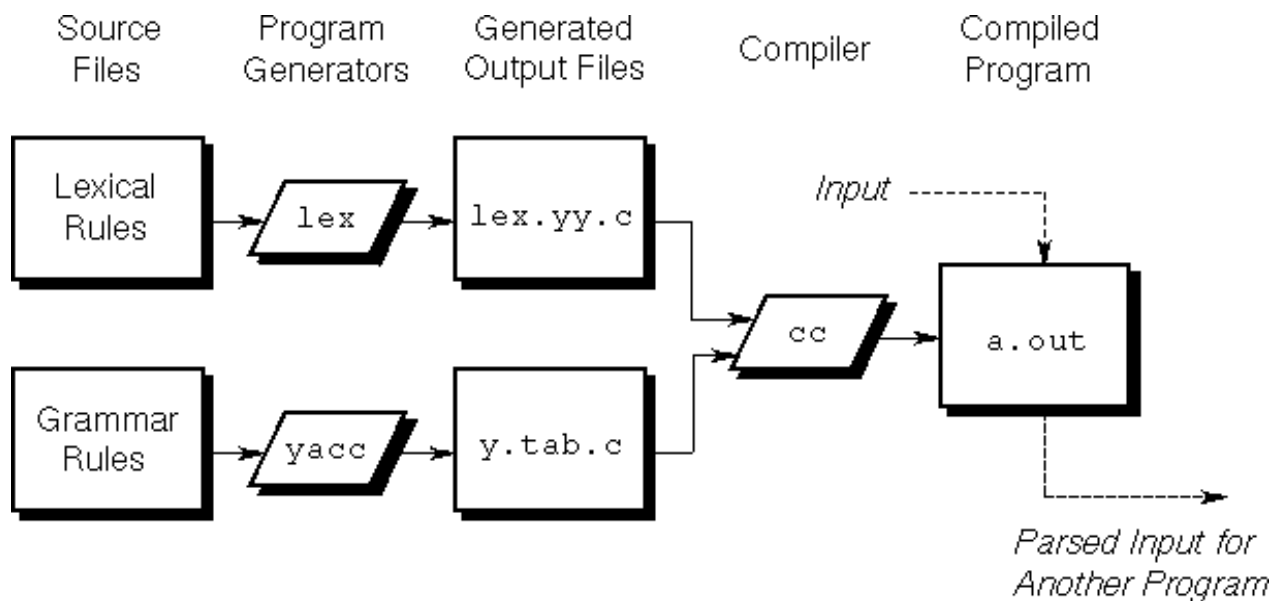
The input to Yacc is a grammar with snippets of C code (called "actions") attached to its rules. Its output is a shift-reduce parser in C that executes the C snippets associated with each rule as soon as the rule is recognized. Typical actions involve the construction of parse trees. Using an example from Johnson, if the call node(label, left, right) constructs a binary parse tree node with the specified label and children, then the rule.

recognizes summation expressions and constructs nodes for them. The special identifiers \$\$, \$1 and \$3 refer to items on the parser's stack.

Yacc produces only a parser (phrase analyzer); for full syntactic analysis this requires an external lexical analyzer to perform the first tokenization stage (word analysis), which is then followed by the parsing stage proper. Lexical analyzer generators, such as Lex or Flex are widely available. The IEEE POSIX P1003.2 standard defines the functionality and requirements for both Lex and Yacc.

Some versions of AT&T Yacc have become open source. For example, source code is available with the standard distributions of Plan 9.

Diagram of YACC



Basic Specifications

Names refer to either tokens or nonterminal symbols. Yacc requires token names to be declared as such. In addition, for reasons discussed in Section 3, it is often desirable to include the lexical analyzer as part of the specification file; it may be useful to include other programs as well. Thus, every specification file consists of three sections: the declarations, (grammar) rules, and programs. The sections are separated by double percent ``%%" marks. (The percent ``%" is generally used in Yacc specifications as an escape character.)

In other words, a full specification file looks like

```
declarations
%%
rules
%%
Programs
```

How the parser works?

Yacc turns the specification file into a C program, which parses the input according to the specification given. The algorithm used to go from the specification to the parser is complex, and will not be discussed here (see the references for more information). The parser itself, however, is relatively simple, and understanding how it works, while not strictly necessary, will nevertheless make treatment of error recovery and ambiguities much more comprehensible.

The parser produced by Yacc consists of a finite state machine with a stack. The parser is also capable of reading and remembering the next input token (called the lookahead token). The current state is always the one on the top of the stack. The states of the finite state machine are given small integer labels; initially, the machine is in state 0, the stack contains only state 0, and no lookahead token has been read.

The machine has only four actions available to it, called shift, reduce, accept, and error. A move of the parser is done as follows:

1. Based on its current state, the parser decides whether it needs a lookahead token to decide what action should be done; if it needs one, and does not have one, it calls yylex to obtain the next token.
2. Using the current state, and the lookahead token if needed, the parser decides on its next action, and carries it out. This may result in states being pushed onto the stack, or popped off of the stack, and in the lookahead token being processed or left alone.

Actions

With each grammar rule, you can associate actions to be performed when the rule is recognized. Actions can return values and can obtain the values returned by previous actions. Moreover, the lexical analyzer can return values for tokens, if desired.

An action is an arbitrary C-language statement and as such can do input and output, call subroutines, and alter arrays and variables. An action is specified by one or more statements enclosed in { and }. For example, the following two examples are grammar rules with actions:

```
A      : '(' B ')'
      {
          hello(1, "abc" );
      }

and

XXX    : YYY ZZZ
      {
          (void) printf("a message\n");
          flag = 25;
      }
```

The \$ symbol is used to facilitate communication between the actions and the parser. The pseudo-variable \$\$ represents the value returned by the complete action.

For example, the action:

```
{ $$ = 1; }
```

returns the value of one; in fact, that's all it does.

To obtain the values returned by previous actions and the lexical analyzer, the action can use the pseudo-variables \$1, \$2, ... \$n. These refer to the values returned by components 1 through n of the right side of a rule, with the components being numbered from left to right. If the rule is

```
A: B C D ;
```

then \$2 has the value returned by C, and \$3 the value returned by D. The following rule provides a common example:

```
expr: '(' expr ')' ;
```

You would expect the value returned by this rule to be the value of the expr within the parentheses. Since the first component of the action is the literal left parenthesis, the desired logical result can be indicated by:

```
expr: '(' expr ')'
{
    $$ = $2 ;
}
```

By default, the value of a rule is the value of the first element in it (\$1). Thus, grammar rules of the following form frequently need not have an explicit action:

```
A : B ;
```

In previous examples, all the actions came at the end of rules. Sometimes, it is desirable to get control before a rule is fully parsed. yacc permits an action to be written in the middle of a rule as well as at the end.

This action is assumed to return a value accessible through the usual \$ mechanism by the actions to the right of it. In turn, it can access the values returned by the symbols to its left. Thus, in the rule below, the effect is to set x to 1 and y to the value returned by C:

```
A      : B
      {
          $$ = 1;
      }
C      : {
      x = $2;
      y = $3;
      }
      ;
```

Actions that do not terminate a rule are handled by yacc by manufacturing a new nonterminal symbol name and a new rule matching this name to the empty string. The interior action is the action triggered by recognizing this added rule

YACC Declaration Summary

<u>Declaration</u>	<u>Description</u>
%start	Specify the grammar's start symbol
%token	Declare a terminal symbol (token type name) with no precedence or associativity specified
%type	Declare the type of semantic values for a nonterminal symbol
%right	Declare a terminal symbol (token type name) that is right-associative
%left	Declare a terminal symbol (token type name) that is left-associative
%nonassoc	Declare a terminal symbol (token type name) that is non-associative (using it in a way that would be associative is a syntax error, <i>Ex: $x \text{ op. } y \text{ op. } z$ is syntax error</i>)

Question:

1. What is the use of yyparse() ?

- You call the function yyparse to cause parsing to occur. This function reads tokens, executes actions, and ultimately returns when it encounters end-of-input or an unrecoverable syntax error. You can also write an action which directs yyparse to return immediately without reading further.

2. What is y.tab.h contains?

- Before writing the LEX program, there must be some way by which the YACC program can tell the LEX program that DIGIT is a valid token that has been declared in the YACC program. The file <y.tab.h> contains definitions for the tokens that the parser program uses. In addition, calc.l contains the rules used to generate the tokens from the input stream.

3. How to declare terminals, nonterminals & start symbol in yacc file ?

- Terminal (or token) names can be declared using the % token declaration, and nonterminal names can be declared using the % type declaration. The % type declaration is not required for nonterminal names. Nonterminal names are defined automatically if they appear on the left side of at least one rule.

4. Justify the need of `yyerror()`. Specify its syntax ?

- The first thing the parser does when it performs the **Error** action is to call a function named **`yyerror()`**. This happens *before* the parser begins going down the state stack in search of a state that can handle the **error** symbol. **`yyerror()`** is a **lex** and **yacc** library function that simply displays a text string argument to **`stderr`** using **`fprintf`**, and returns the integer value received from **`fprintf`**. `yyerror()` is used to cause the parser to initiate error handling
- `void yyerror(parameter){ statements; }`

Practical - 2A

Aim: Write YACC specification to check syntax of a simple expression involving operators +, -, * and /. Also convert the arithmetic expression to postfix.

Program:

```
% {
#include "y.tab.h"
% }
%%
[0-9]+ {yylval=atoi(yytext); return NUMBER;}
[a-zA-Z] {return ID;}
\n {return NL;}
. {return yytext[0];}
%%

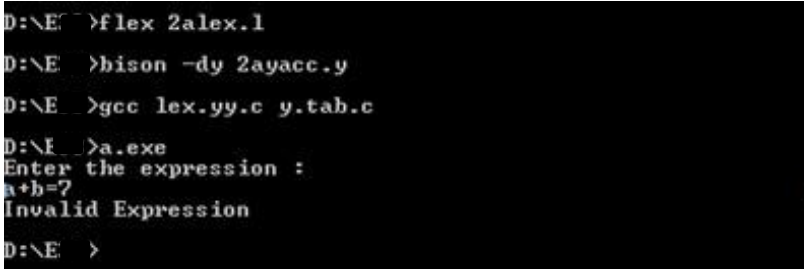
% {
#include<stdio.h>
#include<stdlib.h>
int answer=0;
% }
%token NUMBER ID NL
%left '*' '/'
%left '+' '-'
%%
stmt : exp NL { printf("Answer: %d \n",$1);
               exit(0);}
|
exp1 NL { printf("Valid Expression \nBut, Calculation Can Be Performed On Variables \n");
          exit(0);}
;

exp : exp '+' exp      {$$=$1+$3;}
| exp '-' exp          {$$=$1-$3;}
| exp '*' exp          {$$=$1*$3;}
| exp '/' exp          {$$=$1/$3;}
| '(' exp ')'          {$$=$2;}
| NUMBER               {$$=$1;}

;
exp1 : exp1 '+' exp1
| exp1 '-' exp1
| exp1 '*' exp1
| exp1 '/' exp1
| '(' exp1 ')'
| ID
;
%%
int yyerror(char *msg)
```

```
{printf("Invalid Expression \n");
exit(0);
}
main()
{printf("Enter the expression : \n");
yyparse();}
int yywrap(){return 1;}
Input:
a+b=7
```

Output:



```
D:\E>flex 2alex.l
D:\E>bison -dy 2ayacc.y
D:\E>gcc lex.yy.c y.tab.c
D:\E>./a.exe
Enter the expression :
a+b=7
Invalid Expression
D:\E>
```


Practical - 2B

Aim: To validate syntax of following programming language construct:

Batch E2: while statement

Program:

```
% {
#include "y.tab.h"
% }
chars [A-Za-z]
digit [0-9]
%%
[ \t\n]
while { return WHILE;}
{digit}+ {return NUM;}
{chars}({chars}|{digit})* {return ID;}
"<="|">="|"=="|!=" {return EE;}
"|" {return OR;}
"&&" {return AND;}
. {return yytext[0];}
%%

% {
#include<stdio.h>
#include<stdlib.h>
% }

%token WHILE NUM ID OR EE AND
%right '='
%left AND OR
%left '<' '>' EE
%left '+' '-'
%left '*' '/'
%left '!'
%%
S : ST1 {printf("The given input is valid.\n");
        exit(0);}
    ;

ST1 : WHILE '(' E2 ')' '{' ST '}'

ST : ST ST | E ';'
    ;
```

```

E    : E'+E
      | E'-'E
      | E'*E
      | E'/'E
      | E'<'E
      | E'>'E
      | E EE E
      | E OR E
      | E AND E
      | ID
      | NUM
      ;

```

```

E2   : E'<'E
      | E'>'E
      | E EE E
      | E OR E
      | E AND E
      | ID
      | NUM
      ;

```

```

%%
int yyerror(char *msg)
{ printf("Invalid while Expression \n");
  exit(0);}
main()
{ printf("Enter the while expression: ");
  yyparse();}
int yywrap(){ return 1;}

```

Input:

```
while(a>5;)
```

Output:

```

D:\E_ >flex prac2.1
D:\E_ >bison -dy prac2.y
conflicts: 2 shift/reduce
D:\E_ >gcc lex.yy.c y.tab.c
D:\E_ >a.exe
Enter the while expression: while(a>5;)
Invalid while Expression

```