

-1- Lex and Yacc

"LCS" copy

Introduction:

Both Lex and yacc programs helps to transform
structured input.

But for a C program we have different units
like variables, constants, strings, operators, punctuation etc.

This division into units (of tokens) is known as

Lexical Analysis: Lex helps us by taking a set of
possible tokens and producing a scanner or scanner
descriptions of tokens which we call a lexical analyzer or scanner. The set of descriptions
that can identify those tokens. The set of descriptions
we give to Lex is called a Lex Specification.

The token descriptions that Lex uses are
known as regular expressions. Lex turns these
regular expressions into a form that the lexer can
use to scan the input text extremely fast.
It is independent of the number of expressions
a program is trying to match.

As the input is divided into tokens, a
program often needs to establish the relationship
among the tokens. This task is known as Parsing
and the set of rules that define the relationships
of the program understands is a grammar.

Yacc (yet another Compiler-Compiler) takes a concise
description of a grammar and produces a C
routine that can parse that grammar, called a
parser (Symbol Analyzer).

The yacc parser automatically detects
whenever a scanner of input tokens matches on
the rules in the grammar and also directly
a syntax error whenever its input doesn't match
any of the rules.

Symbol Table!

2. compiler stores the variable names, (in other words all tokens) in the program in its symbol table. Each naming stand along with information describing the nature of the symbol. In a compiler the information is the type of the symbol, declaration scope, variable type etc.

The symbol table can be defined using any data structure. If we are linked list the each node definition may look as follows:

struct node

{
char *stren_name;

int token-type;

struct node *next;

- whenever we get a token if entered into the symbol table along with the token type.

Grammars

As the input is divided into tokens, a program often needs to establish the relationship among the tokens. A description of such a set of relationships among the tokens and a set of actions is known as a grammar.

Example: One grammar will might be,

1. date : monthname Day ; Year ;

i.e. the rules to read or scan a valid date is monthname followed by space followed by date followed by , followed by year. e.g. October, 10, 2004 will be matched.

The input October 10, 2004 will be matched by the above rule.

2. the other grammar rule might be,

2. date : Monthnumber /' Day /' year ;

The input
to the above rule.

10/10/2004 will be matched by

Yacc communication (Parser-lexer communication)

When we use a lex scanner and a yacc parser together, the parser is higher level routine. It calls the lexer yylex() whenever it needs a token from the input. The lexer then scans through the input recognizing tokens. As soon as it finds a token of interest to the parser, it returns a token's code as the value of yylex().

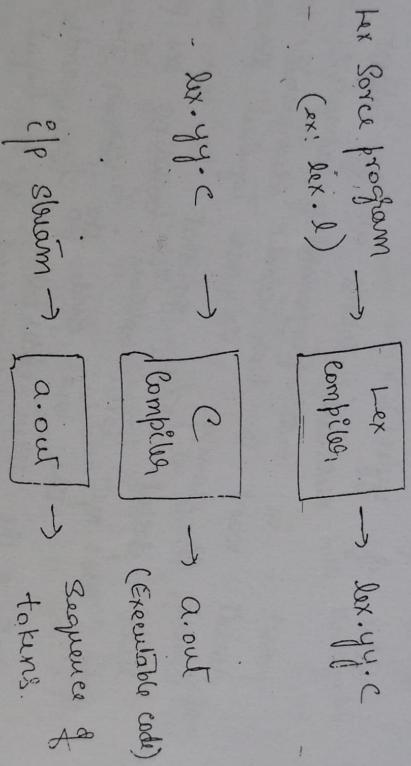
Not all tokens are of interest to the parser - in most programming languages, the parser doesn't care about comments and whitespace, the want to look about for these ignored. Looking for example, so that it can continue on spaces, doesn't return go that token to the parser. Lex doesn't return without knowing the parser. The next token without having to agree the lexer and the parser have to agree is the token code zero what is the logical end of the input. defined for no logical end of the input.

USING LEX

4

Lex is a tool for building lexical analyzer & lexers. A Lex takes an arbitrary input string and tokenizes it, i.e. divides it up into lexical tokens. This tokenized output can then be processed further, usually by yacc or it can be the "end product."

Lex uses the user's expressions and actions into the host-generated purpose language. (usually C language). Creating a lexical analyzer with Lex is shown in following steps.



First, a specification of a lexical analyzer is prepared by writing a program `lex.l` in the Lex language. Then, Lex is run through the Lex compiler to produce a C program `lex.y.c`. The program `lex.y.c` consists of a tabular representation of a transition diagram constructed from the regular expressions of `lex.l`, together with a standard routine that uses the table to recognize lexemes. The actions associated with regular expressions in

Mr. J. do plan of C code and are carried out
from directly to Mr. V.P.C.
Finally, Mr. V.P.C. is run through the
C compiler to produce an object program a.out,
which is the final answer that transforms an
Up stream into a sequence of bytes.

Lex Specification:

- 1. Lex programs consist of 3 parts:
- 1.1 Declarations (Definition section) of
- 1.2 Translation rules
- 1.3 Auxiliary procedures

The definition section include declarations of variables,
and regular definitions.
The background in C code with the special
"%%" and "%". Lex copies the modified C
definitions. "%
and "% directly to the generated C
between "% and "% directly to the generated C code here.
So we may write any aux
file, so we may write any aux file. Each
The first section is the rules section. Each
rule is made up of two parts: a pattern and
rule is made up of two parts: a pattern and
an action, separated by whitespace. The lex will
execute the action when it recognizes the pattern.
The final section is the auxiliary procedure
or user subroutines section which will consist of
any legal C code. Lex copies it to the C
file after the end of the last generated code.

Note: To run a lex program, link the library
1. Use a .L d
2. Use lex.vpc -lQ

(6)

Lex Regular Expressions

A regular expression is a pattern describing 6
using a "meta" language, a language that
use to describe particular patterns of interest.
The characters that form regular expressions are:

- matches any single character except the

new line character \n -

- matches zero or more copies of the preceding expression.

- Matches one or more occurrences of the preceding regular expression.
- Ex: $[0-9]^+$ matches "1", "123", or "123456" but not an empty string.

- [] - A character class which matches any character within the brackets. If the first character is a circumflex ("^") it forces the matching to match out changes the meaning to match out characters except the ones within the brackets.

A dash "-" inside the square brackets indicates a character range.

Ex: $[0-9]$ means $[0123456789]$.

- If we want to include the character "-" in a character class, it should be first or last. Ex: $[-+0-9]$ matches all the digits and "-" and "+".

- matches the end of a line as the last character of a regular expression

8.

Indicates how many times the previous pattern is allowed. To match when containing one or two numbers.

For ex: A{1,3} matches one to three occurrences of the letter A.

? (- optional operator matches zero or one occurrence of the preceding regular expression

Ex: ① ab?c matches either ac abac

② [t]?[o-q]+ matches a + no matching either the preceding regular expression including an optional leading + + so and so matches to the above regular expression.

matches either the preceding regular expression or the following regular expression

Ex: "cow" | "sharp" | "cat" matches any of

the 3 words.

/ - matches the preceding regular expression but only if followed by the following regular expression.

Ex: ① matches "0" in the string "long" "01" but it would not match anything in the strings "0" or "02".

The match is matched by the pattern following the slash is not consumed and remaining to be treated into subsequent tokens. Only one slash is permitted

⑤

Look Diminu Function

()

- Groups a series of regular expression
together into a new regular expression.
for ex: (01) represent the char
sequence 01. Parentheses are used
when building up complex patterns with
*, + and |.

Example of the regular expression:

[-] [0 - 9] + - regular expression for - an integer.
[+] ? [0 - 9] + - "

[t] ? [0 - 9] * \ [0 - 9] + - "

Notice that "\ " before " " do make it -
separate from a wild card character " ".

[a - 2 A - 2] [a - 2 A - 2] * - is a typical regular
expression for identifying -

dynamic finger in programming

Lex Action:

When our expression written as above is
matched, Lex executes the corresponding action.
Note that there is a default action, which consists of
copying the input to the output. If the file pointer
variable is not initialized then it reads from keyboard
else from specified file. If the file pointer
variable is not initialized then it writes to memory
else to specified file.

Usually it leaves the text that matched some regular expression, in an external character array named yylext.

Ex: $[a-zA-Z]^+$ { printf("%s", yytext); }
 will print the string in yylext.
 This action is so common that it can be written as ECHO.

i.e. $[a-zA-Z]^+ \{ ECHO; \}$ is the same as above.

Lex provides a counter variable yylen which is the no. of characters matched.

Ex: $[a-zA-Z]^+$ { words++; chars=char++;
 yylen; }
 which counts both the no. of words and the no. of characters in the words recognized.

Example program:

The following lex program counts the no. of characters, words and newlines in a stream of input characters string.

```
%{  
int c=0, w=0, l=0;
```

```
%}
```

```
WORD [^ \t\n]^+
```

```
EOL [\n]
```

```
%%
```

```
{WORD}
```

```
{EOL}
```

```
}%
```

```
{w++; c=c+yylen; }
```

```
{c++; l++; }
```

```
{c++; }
```

(10) %/
int yywrap()
{
 return (1);
}

main()
{

printf("Enters the string upto ^d\n");

yylex();

printf("No. of chars = %d\n No. of words = %d\n No. of
 new lines = %d\n", c, w, l);

}

The section bracketed by "%{" and "%}" is C code which is copied directly to C code. The next 2 lines are definitions. Lex provides a simple substitution mechanism. To make it easier to define long & complex patterns.

The next section is rules and action part. If there more than one action is there, then the action part should be enclosed between {}'s.

The main program first calls yylex() which is the lexer's entry point. When yylex() reaches the end of its input file, it calls yywrap(), which returns a value of 0 or 1. If the value is 1, the program is done and there is no more input. If the value is 0, on the other hand, the lexer assumes that yywrap() has opened another file for it to read, and continues to read from yyin. The default yywrap() always

returns 1. By providing our own version of yywrap(), we can have our program read all of the files named on the command line, one at a time.

The lex provides the routines input() and unput(). The input() routine handles calls from the lexer to obtain characters. When the current argument is exhausted it moves to the next argument, if there is one, and continues scanning. If there are no more arguments, we treat it as the user's end-of-file condition and return a zero byte. [Note: For example refer variable count ^{to} program]

The unput() routine handles calls from the lexer to "push back" characters into the input stream. It does this by reversing the pointers direction, moving backwards in the string.

start states:

~~start state~~ start-state is a method of capturing context sensitive information within the lexer. Tagging rules with start states tells the lexer only to recognize the rules when the start state is in effect. The start state is defined in the definition section using %s followed by state name. In the rules section we can use Keen state name which is enclosed between < and >. These rules are only

recognized when the lexer is in state ~~STATEMENT~~
stated in the definition section. #12.

We enter a new state with a BEGIN
statement. To change back to the default state,
we use "BEGIN 0". (the default state zero,
is also known as INITIAL.)

Ex:
~~/*~~ To count the no. of comment lines in a
C pgm.

Y{
int com=0;

Y}

/* COMMENT

*/

[\E]* /* { BEGIN COMMENT; }

<COMMENT> In { com++; }

<COMMENT> */ { BEGIN 0; com++; }

<COMMENT> . { ; }

Y,Y

main()

{

yyEx();

printf("no. of comments = %d", com);

}

60 copies

-7-

USING YACC

Roshan Fernandes

(Yet Another Compiler
comp)

Lex divides the input into pieces (tokens) and then Yacc takes these pieces and glues them together logically.

Grammars: Yacc takes a grammar that we specify and writes a parser that recognises "sentences" in that grammar. A grammar is a series of rules that the parser uses to recognise syntactically valid input.

Consider the following grammar;

statement → NAME = expression
expression → NUMBER + NUMBER | NUMBER

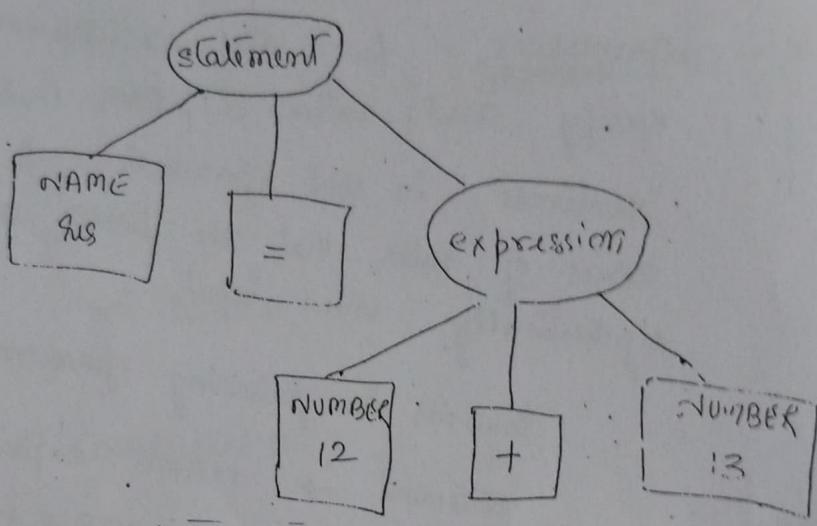
The vertical bar ":" means there are 2 possibilities for the same symbol, i.e. an expression can either be an addition or a subtraction.
The symbol to the left of the " \rightarrow " is known as the Left-Hand side of the rule (LHS) and the symbols to the right are Right-Hand side of rule (RHS).

Symbols that actually appear in the input and are returned by the lexer are terminal symbols or tokens. Symbols that appear on left-hand side of the rules are non-terminal symbols or non-terminals.

Usually Terminal symbols are given in uppercase and Non-Terminal are given in lowercase.

(14) The usual way to represent a parsed sentence is as a tree.

For example, if we parsed the input "res = 12 + 13" with the above grammar, the tree would look like below.



"12 + 13" is an expression, and

"res = expression" is a statement.

Note: A yacc parser doesn't actually create this tree as a data structure.

Every grammar includes a start symbol, the one that has to be at the root of the parse tree.

In the above grammar, Statement is the start symbol.

Recursive rules:

Rules can refer directly or indirectly to themselves. This important ability makes it possible to parse arbitrarily long input sequences.

Example:

-8-

$\text{expression} \rightarrow \text{NUMBER}$

$| \quad \text{expression} + \text{NUMBER}$

$| \quad \text{expression} - \text{NUMBER}$

Now we can parse a sequence like

" $res = 14 + 23 - 11 + 7$ " by applying the expression rules repeatedly.

Shift/Reduce Parsing:

A Yacc parser uses Shift/Reduce Parsing technique.

As the parser reads tokens, each time it reads a token that doesn't complete a rule, it pushes the token on an internal stack and switches to a new state reflecting the token it just read. This action is called Shift. When it has found all the symbols that constitute the right-hand side of a rule, it pops the right-hand side symbols off the stack, pushes the left-hand side symbol onto the stack and switches to a new state reflecting the new symbol on the stack. This action is called Reduction, since it reduces the number of items on the stack. Whenever Yacc reduces a rule, it executes code associated with the rule.

Let us look how it parses the input " $res = 12 + 13$ " using the simple rules,

$\text{statement} \rightarrow \text{NAME} = \text{expression}$

$\text{expression} \rightarrow \text{NUMBER} + \text{NUMBER}$ | $\text{NUMBER} - \text{NUMBER}$

(16)

The parser starts by shifting tokens on to
internal stack one at a time:

 q_{13} $q_{12} =$ $q_{13} = 12$ $q_{13} = 12. +$ $q_{13} = 12 + 13$

At this point it can reduce the rule "expression
 $\rightarrow \text{NUMBER} + \text{NUMBER}$ " so it pops $13,$
 $+$ and 12 from the stack and replaces them
with expression.

i.e. $q_{13} = \text{expression}$.

Now EF reduces the rule " $\text{statement} \rightarrow \text{NAME} =$
 expression ", so it pops $\text{expression}, =$
and q_{13} from the stack and replaces them
with statement . We have reached the end of
the input and the stack has been reduced to
the start symbol, so the input was valid
according to the grammar.

What Yacc cannot Parse:

Yacc cannot handle ambiguous grammars,
ones in which the same input can match more than
one parse tree. [But yacc can deal with a limited
but useful set of ambiguous grammars.]

Yacc also cannot deal with grammars
that need more than one token of lookahead
to tell whether it has matched a rule.

Example:

-9-

Consider the following grammar:

phrase	\rightarrow	cart-animal	(AND)	CART
		walk_animal	, AND	PLOW
cart_animal \rightarrow		HORSE		GOAT
walk_animal \rightarrow		HORSE		OX

Yacc cannot handle it because it requires two symbols of lookahead. In particular, in the input "HORSE AND CART" it cannot tell whether HORSE is a cart-animal or a walk-animal until it sees CART, and Yacc cannot look that far ahead.

If we changed the first rule to this:

phrase	\rightarrow	cart_animal	CART
		walk_animal	PLOW

Yacc would have no problem, since it can look one token ahead to see whether an input of HORSE is followed by CART. In this case the token is a cart-animal; if it is followed by PLOW, then the horse is a walk-animal.

A Yacc Parser!

A Yacc grammar has the following 3-part structure as its specification:

The 3 sections are definition section, auxiliary section and programs | auxiliary production section

(18)

Definition Section:

It can include a declaration part enclosed between `%{` and `%}`.
The tokens & the terminal symbols have to be defined using `%token`.

Ex: `%token NAME NUMBER`

We can use single quoted characters as tokens without declaring them.

Example: `'=' . '+' , etc.`

The Rules Section:

The rules section consists of a list of grammar rules. In ~~most~~ we use a `:` (colon) between the left- and right-hand sides of a rule, and we put a semicolon at the end of each rule:

Ex: `%token NAME NUMBER`

`statement: NAME '=' expression`

`expression`

`expression: NUMBER '+' NUMBER`

`NUMBER '-' NUMBER`

The symbol on the left hand side of the first rule in the grammar is normally in bold symbol.

We can use `%start` declaration to specify starting symbol.

-10- Symbol Values and Actions:

Every symbol in a yacc parser has a value. The value gives additional information about a particular instance of a symbol. If a symbol represents a number, the value would be the particular number. If it represents a literal text string, the value would probably be a pointer to a copy of the string. If it represents a variable in a program, the value would be a pointer to a symbol table entry describing the variable. Some tokens don't have a useful value, ex: a token representing a close parenthesis, since one close parenthesis is the same as another.

Whenever the parser reduces a rule, it executes user code associated with the rule known as the rule's action. The action appears in braces after the end of each rule, before the semicolon or vertical bar.

The action code can refer to the values of the right-hand side symbols as \$1, \$2, ... and can set the value of the left-hand side by setting \$\$.

Example:

ytoken NAME NUMBER;

statement : NAME '=' expression

| expression { printf("ydtu", \$1); }

expression : expression '+' NUMBER, { \$\$ = \$1 + \$3
| expression '-' NUMBER, { \$\$ = \$1 - \$3; }

NUMBER { \$1-\$1; }

the expression building rules. the first and second numbers' values, are \$1 and \$3 respectively. the operator's value would be \$2.

The Lexical

The parser is a higher level routine and calls the lexical yylex() whenever it needs a token from the input. As soon as the lexical finds a token of interest to the parser, it returns to the parser, returning the token code as the value.

Yacc defines the token names in the parser as C preprocessor names in y.tab.h

Example: Here is a simple lexer to provide tokens for our parser:

```
#include "y.tab.h"
extern int yylval;
```

```
[0-9]+ { yylval = atoi(yylext); return NUMBER; }
```

```
[ \t\n\r\f\v ] { } // ignore whitespace
```

```
\n { return 0; } /* End of file */
```

```
.. { return yylext[0]; }
```

```
/* */
```

Whenever the lexer returns a token to the parser, if the token has an associated value,

(21)

the lexer must store the value in `yyval` before returning the token.

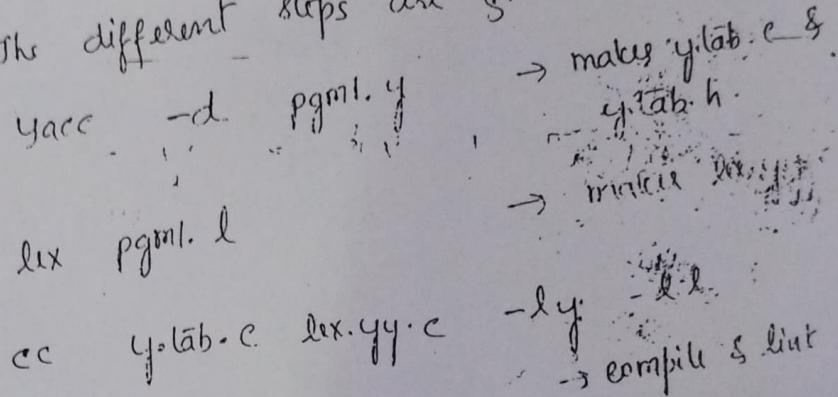
Compiling and running a simple parser:

On a UNIX system, Yacc takes our grammar and creates `y.tab.c`, the C language parser and `y.tab.h`, the include file with the token number definitions. Lex creates `lex.yy.c`, the C language parser. We need only compile them together with the Yacc & Lex libraries.

The main() program calls the parser, if `yyparse()`.

and exits.

The different steps are given below:



/a.out <1>

10+50 <1>

60

/a.out <1>

100 + -50 <1>

syntax error.

It correctly reports a message "syntax error". When we enter something that doesn't conform to the grammar.

We can also call a routine `yyerror`

22)

Ex:

if ($\$2 == 0$)
 yyerror("Divide by zero")

Ambiguity and Conflicts

In any arithmetic expression grammar, operators are grouped into levels of precedence from lowest to highest.
 If there are operators of same precedence then associativity has to be specified.

example:

$a-b-c$ can be $(a-b)-c$ or $a-(b-c)$. So the operator can be grouped as left associative and right associative

left associative example $\rightarrow (a-b)-c$

right associative example $\rightarrow a-(b-c)$

The precedences and associativity are attached to tokens in the declaration section. This is done by a series of lines beginning with a yacc keyword

%left, %right & %nonassoc, followed by a list of

tokens. The keyword %left is used to describe left associative operators, %right is used to describe right associative operators and %nonassoc is used to describe operators, that may not associate with themselves.

Note: All the tokens on the same line are assumed to have the same precedence level and associativity. The lines are listed in order of increasing precedence & binding strength.