

A brief lex & Yacc tutorial

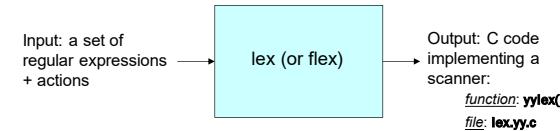
Slides originally created by:
Saumya Debray (The Univ. of Arizona)
Modifications by:
Amey Karkare (IIT Kanpur)



flex (and lex): Overview

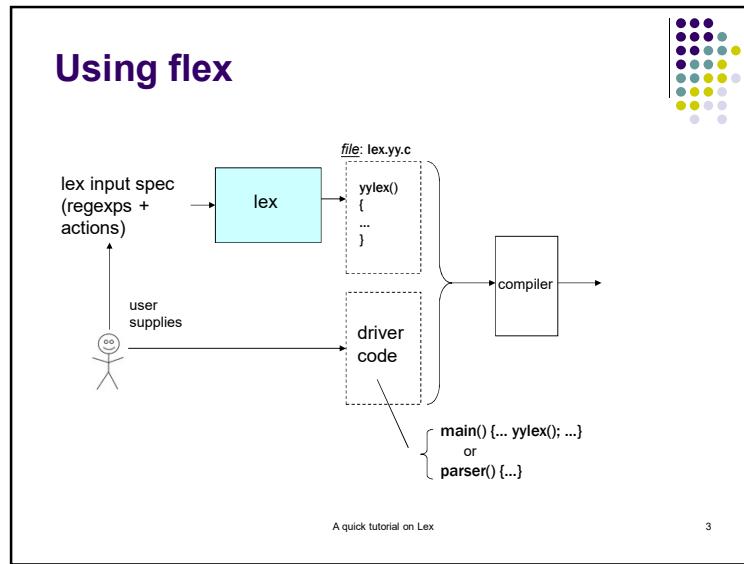
Scanner generators:

- Helps write programs whose control flow is directed by instances of regular expressions in the input stream.



A quick tutorial on Lex

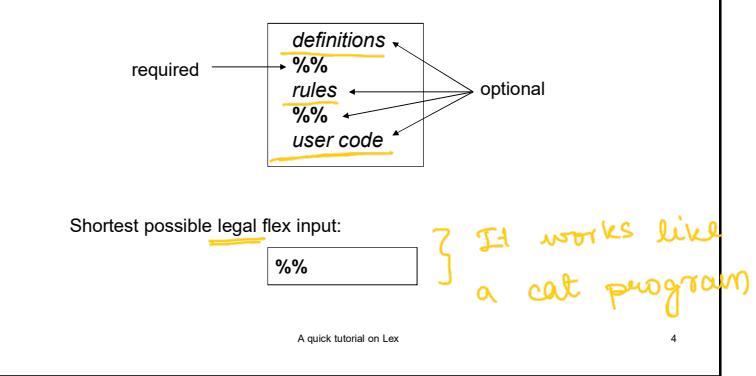
2



3

flex: input format

An input file has the following structure:



4

Definitions

- A series of:
 - ▶ *name definitions*, each of the form
`name definition`

e.g.:

```
DIGIT      [0-9]
CommentStart  /*"
ID         [a-zA-Z][a-zA-Z0-9]*
```

- ▶ *start conditions*
- ▶ stuff to be copied verbatim into the flex output (e.g., declarations, `#includes`):
- enclosed in `%{ ... %}`

A quick tutorial on Lex

5

Rules

- The *rules* portion of the input contains a sequence of rules.

- Each rule has the form

`pattern action`

where:

- ▶ *pattern* describes a pattern to be matched on the input
- ▶ *pattern* must be un-indented
- ▶ *action* must begin on the same line.(version dependent), for multi lined action : use `{ }`

A quick tutorial on Lex

6

Example

A flex program to read a file of (positive) integers and compute the average:

```
%{
#include <stdio.h>
#include <stdlib.h>
%
dgt [0-9]    ]> definition
%{
(dgt)+ return atoi(yytext);
%
void main()
{
    int val, total = 0, n = 0;
    while ( (val = yylex()) > 0 ) {
        total += val;
        n++;
    }
    if (n > 0) printf("ave = %d\n", total/n);
}
```

rules
action

A quick tutorial on Lex

7

Example

A flex program to read a file of (positive) integers and compute the average:

```
%{
#include <stdio.h>
#include <stdlib.h>
%
dgt [0-9]    ]> definition
%{
(dgt)+ return atoi(yytext);
%
void main()
{
    int val, total = 0, n = 0;
    while ( (val = yylex()) > 0 ) {
        total += val;
        n++;
    }
    if (n > 0) printf("ave = %d\n", total/n);
}
```

* 12
lexeme → "12"
yytext → convert to integer .

A quick tutorial on Lex

8

Example

A flex program to read a file of (positive) integers and compute the average:

```
%{  
#include <stdio.h>  
#include <stdlib.h>  
%}  
{ dg[1-9]  
% %  
dg) return atoi(yytext);  
% %  
void main()  
{  
    int val, total = 0, n = 0;  
    while ((val = yylex()) > 0 ) {  
        total += val;  
        n++;  
    }  
    if (n > 0) printf("ave = %d\n", total/n);  
}
```

A quick tutorial on Lex

9

Example

A flex program to read a file of (positive) integers and compute the average:

```
{%  
#include <stdio.h>  
#include <stdlib.h>  
%}  
{digit} [0-9]  
%%  
{digit} return atoi(yytext);  
%%  
void main()  
{  
    int val, total = 0, n = 0;  
    while ((val = yylex()) > 0) {  
        total += val;  
        n++;  
    }  
    if (n > 0) printf("ave = %d\n", total/n);  
}
```

A quick tutorial on Lex

11

- * Here (≤ 0) can't be handled by this simple calculator.
- * It can't handle " $\backslash n$ ", thus it's printed again.

Example

A flex program to read a file of (positive) integers and compute the average:

```

definitions { %{
    #include <stdio.h>
    #include <stdlib.h>
    %}
    (digit) [0-9]
    %%}

rules { (digit) return atoi(yytext); %%}

void main()
{
    int val, total = 0, n = 0;
    while ( (val = yylex()) > 0 ) {
        total += val;
        n++;
    }
    if (n > 0) printf("ave = %d\n", total/n);
}

```

defining and using a name

char * yytext;
a buffer that holds the input
characters that actually match the
pattern

A quick tutorial on Lex

1

LHS: Non-terminal
RHS: - " — + Terminal

Hands On Example

BNF format
yacc syntax

program : slist
| /* empty */
; → yacc ;
slist : assign ';' → source language
| exit_command ';' ←
| print_exp ';' ←
| slist assign ';' ←
| slist print_exp ';' ←
| slist exit_command ';' ←
;
assign : identifier '=' exp

```
exp : term
      | exp '+' exp
      | exp '-' exp
      | exp '*' exp
      | exp '/' exp
      | '(' exp ')'
;
```

term : number → tokens
| identifier → terminals

In YACC, we define our tokens separately

13

Lex script

```
%{
#include <iostream>
#include "y.tab.h"
#include <cstring>
void yyerror (char *s);
int yylex();
%} ↳ returns token no.s,
digit [0-9]
alpha [a-zA-Z_]
alphanum ({alpha}|{digit})
ws [\t\n] ↳ whitespace
%%
```

definitions

A quick tutorial on Lex

13



Matching the Input

- When more than one pattern can match the input, the scanner behaves as follows:
 - the longest match is chosen;
 - if multiple rules match, the rule listed first in the flex input file is chosen;
 - if no rule matches, the default is to copy the next character to **stdout**. (**ECHO**)

```
(cs) {printf("Department");}
(cs)[0-9]{3} {printf("Course");}
[a-zA-Z]+[0-9]+ {printf("AnythingElse");}
Input: cs335
```

{3} → 3 repetition

*② Smaller & special
↳ first*
*③ longer & general
↳ later*

A quick tutorial on Lex

15



Lex script (contd...)

```
"print" { return print; } this needs to be defined somewhere.
"calculate" { return print; }
"exit" { return exit_command; }

{digit}+ { yyval.num = atoi(yytext); return number; } ↳ returning the "number" class not the value.
{alpha}{alphanum}* { yyval.id = strdup(yytext); ③ The no. itself is stored in yyval. return identifier; }

{ws} /* nothing */

[-+\(\)=/*\n; { return yytext[0]; } → we get the corresponding ASCII value as their token number.
④ ECHO, yyerror("::"); matches every character
```

↳ matches every character

↳ the char won't be printed on stdout.

*⑤ ECHO the character back → no error.
→ Everything matches
→ cat program*



④ The no. itself is stored in yyval.

14

Control flow of lexer

```
yylex() {
    /* scan the file pointed to by yyin (default stdin) */
    1. Repeated call of input() to get the next character from the input .....
    2. Occasional calls of unput() .....
    3. Try matching with the regular expression and when matched do the action part.....

    /* got EOF */
    int status = yywrap(); /* default behaviour - return 1 */
    if(1 == status)
        exit();
    else
        yylex();
}

/* Redefine yywrap to handle multiple files */
int yywrap() {
    ....
    if(exists other files to process) { yyin = nextFilePtr; return 1; }
    else { return 0; }
}
```

A quick tutorial on Lex

16

0

1

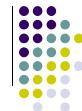
1



Command Line Parsing

Lookahead

Start Conditions

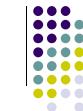


- Used to activate rules conditionally.
 - Any rule prefixed with <S> will be activated only when the scanner is in start condition S.
- This will only match when we're in MAGIC state.*
- ```
%s MAGIC ←——Inclusive start condition
%%
<MAGIC>.+ {BEGIN INITIAL; printf("Magic: "); ECHO;}
magic {BEGIN MAGIC}
Input: magic two three → go into state MAGIC
```
- Warning:** A rule without an explicit start state will match regardless of what start state is active.
- This will match inside & outside magic state*
- ```
%s MAGIC    ←——Inclusive start condition
%%
magic    {BEGIN MAGIC}
ECHO;
<MAGIC>.+ {BEGIN INITIAL; printf("Magic: "); ECHO;}
```
- This becomes useless.*

A quick tutorial on Lex

17

Start Conditions (cont'd)



WayOut:

↑ Exclusive start state

- Use of explicit start state using %x MAGIC
- For versions that lacks %x

```
%s NORMAL MAGIC
%%
%
{
  BEGIN NORMAL;
}
<NORMAL>magic    {BEGIN MAGIC}
<NORMAL>.+     ECHO;
<MAGIC>.+   {BEGIN NORMAL; printf("Magic: "); ECHO;}
```

A quick tutorial on Lex

18

Putting it all together

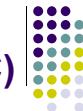


- Scanner implemented as a function `int yylex();`
 - return value indicates type of token found (encoded as a +ve integer);
 - the actual string matched is available in `ytext`.
- When compiling, link in the flex library using “-l” (or equivalent flag)
- You need token type encodings → `y.tab.h`
 - Without Yacc, manually define const/enum/#define names in a file `y.tab.h`
 - use “#include `y.tab.h`” in the definitions section of the flex input file.
- Later, it can be replaced by Yacc generated `y.tab.h`

A quick tutorial on Lex

19

Syntax Analyzer Generator (YACC) “Yet Another Compiler Compiler”



bison



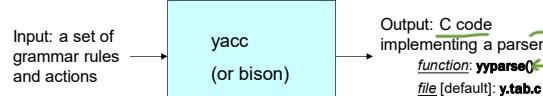
A quick tutorial on Lex

20

Yacc: Overview

Parser generator:

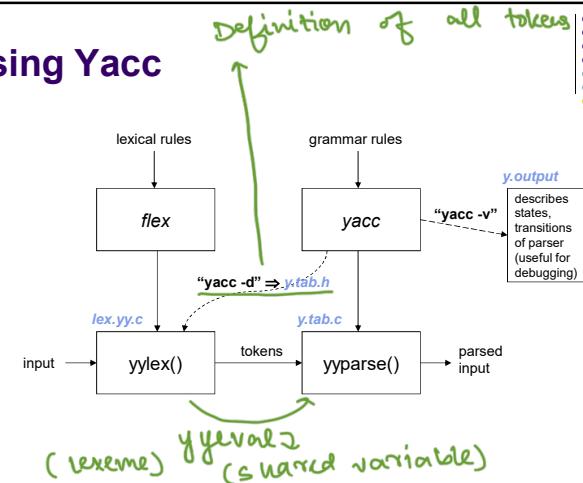
- ▶ Takes a specification for a context-free grammar.
- ▶ Produces code for a parser.



A quick tutorial on yacc

21

Using Yacc



A quick tutorial on yacc

22

int yyparse()

↳ supposed to take whole file

- Called once from main() [user-supplied]
- Repeatedly calls yylex() until done:
 - ▶ On syntax error, calls yyerror() [user-supplied]
 - ▶ Returns 0 if all of the input was processed;
 - ▶ Returns 1 if aborting due to syntax error.

Example:

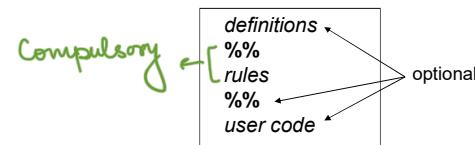
```
int main() { return yyparse(); }
```

A quick tutorial on yacc

23

yacc: input format

A yacc input file has the following structure:



A quick tutorial on yacc

24

Parser

```
%{
    #include<...>
    Declarations .... Copied verbatim into the output yyparse()
%}

%token NAME NUMBER ??.

statement : NAME '=' expression
            | expression [printf("%d\n", $1);]
            | expression '+' NUMBER { $$ = $1 + $3; }
            | expression '*' NUMBER { $$ = $1 * $3; }
            | NUMBER { $$ = $1; }

main() { yyparse(); }

Lexer
```

1st non-terminal

Lexer and Yacc have to agree on terminal names.

No error handling

Parser doesn't have any rule for # → Parser produces error.



: Terminals.

→ optional [] → Action

25

Yacc script (demo) – v1

- Calculator. Try with Inputs:

- 3 + 5
- 3 + 5 + 6
- 2 * 7
- 3 + 5 * 2
- val = 3 + 3

error.

Name handling not done properly.

A quick tutorial on Lex

26

Declaring Return Value Types

User has to do this

union allows these to use the same memory (never used together).

- Default type for nonterminal return values is **int**.
- Need to declare return value types if nonterminal return values can be of other types:

- Declare the union of the different types that may be returned:


```
%union {
    struct symtab *st_ptr;
    double dvalue;
}
```
- Specify which union member a particular grammar symbol will return:


```
%token <st_ptr> NAME;
%token <dvalue> NUMBER;
%type <dvalue> expression;
```

For non-terminals as well we have to specify the type.

A quick tutorial on yacc

27

Yacc script (demo) – v2

- Calculator. Try with Inputs:

- 3 + 5
- 3 + 5 + 6
- 2 * 7
- 3 + 5 * 2
- val = 3 + 3

A quick tutorial on Lex

28

Specifying Operator Properties

- Binary operators: `%left`, `%right`, `%nonassoc`:

```
%left '+' '-'
%left '*' '/'
%right '^'
```

Operators in the same group
have the same precedence

- Unary operators: `%prec`

- Changes the precedence of a rule to be that of the token specified. E.g.:

```
%left '+' '-'
%left '*' '/'
Expr: expr '+' expr
| '-' expr %prec '*'
| ...
```

A quick tutorial on yacc



29

Specifying Operator Properties

- Binary operators: `%left`, `%right`, `%nonassoc`:

```
%left '+' '-'
%left '*' '/'
%right '^'
```

Operators in the same group
have the same precedence

- Unary operators: `%prec`

- Changes the precedence of a rule to be that of the token specified. E.g.:

```
%left '+' '-'
%left '*' '/'
Expr: expr '+' expr
| '-' expr %prec '*'
| ...
```

A quick tutorial on yacc



30

Specifying Operator Properties

- Binary operators: `%left`, `%right`, `%nonassoc`:

```
%left '+' '-'
%left '*' '/'
%right '^'
```

Operators in the same group
have the same precedence

Across groups, precedence
increases going down.

- Unary operators: `%prec`

- Changes the precedence of a rule to be that of the token specified. E.g.:

```
%left '+' '-'
%left '*' '/'
Expr: expr '+' expr
| '-' expr %prec '*'
| ...
```

A quick tutorial on yacc



31

Yacc script (demo) – v3



A quick tutorial on Lex

32

* when yacc encounters an error, it skips it and looks for other errors. This can give us more than one error in one time.

* Better turn-around time.

; \n →
alignment
tokens.

Error Handling

- The “token” ‘**error**’ is reserved for error handling:
 - can be used in rules;
 - suggests places where errors might be detected and recovery can occur.

Example:

```

stmt : IF '(' expr ')' stmt
      | IF '(' error ')' stmt
      | FOR ...
      | ...
  
```

Intended to recover from errors in 'expr'

A quick tutorial on yacc 33

Parser Behavior on Errors

When an error occurs, the parser:

- pops its stack until it enters a state where the token ‘error’ is legal;
- then behaves as if it saw the token ‘error’
 - performs the action encountered;
 - resets the lookahead token to the token that caused the error.
- If no ‘error’ rules specified, processing halts.
program exits.

A quick tutorial on yacc

34

Error Messages

- On finding an error, the parser calls a function


```
void yyerror(char *s) /* s points to an error msg */
```

 - user-supplied, prints out error message.
- More informative error messages:
 - int ychar:** token no. of token causing the error.
 - user program keeps track of line numbers, as well as any additional info desired.

A quick tutorial on yacc 35

Conflicts

- A conflict occurs when the parser has multiple possible actions in some state for a given next token.
- Two kinds of conflicts:
 - reduce-reduce conflict:
 - Start: x B C
 - | y B D;
 - x: A;
 - y: A;
 - y.output** generated using : yacc -v
 - 1: Reduce/reduce conflict (reduce 3, reduce 4) on B.
 - state 1
 - x: A_ (3)
 - y: A_ (4)
 - reduce 3

$S \rightarrow X \mid Y$
 $X \rightarrow a$
 $Y \rightarrow a$
 $S \rightarrow a \rightarrow a \mid b \rightarrow b$ which over?

default : one listed first

A quick tutorial on yacc

36

Conflicts

- shift-reduce conflict:

- Start: x
- | yR;
- x: A R;
- y: A;

- y.output generated using : yacc -v
 - 4: shift/reduce conflict (shift 6, reduce 4) on R.
 - state 4
 - x: A_R
 - y: A_ (4)
 - R shift 6

$$\begin{aligned} S &\rightarrow x \rightarrow A R \\ S &\rightarrow y R \rightarrow A R \end{aligned}$$

* Should we take entire matching part (yR) or only a part (x)?

A quick tutorial on yacc

37

Example of a conflict

Grammar rules:

```
expr : TERMINAL
      | expr '-' expr
```

y.output:

```
4: Shift/reduce conflict (shift 3, reduce 2) on -
State 4:
expr : expr__ - expr
expr : expr - expr_
```

Way Out : define precedence and associativity of the operators.



38

Example of a conflict

Grammar rules:

```
rule : command optional_keyword '(' id_list ')'
;
Optional_keyword:
| '(' keyword ')'
;
Shift / reduce conflict on first parenthesis ( in input stream.
```

Way Out : flattening

```
rule : command '(' id_list ')'
| command '(' keyword ')' '(' id_list ')';
;
```



39

Example of a conflict

Grammar rules:

```
rule : grp_A
      | grp_B
;
grp_A: A_1
      | COMM
;
grp_B: B_1
      | COMM
;
```

Way Out : flattening
Add COMM in a separate rule.
grp_comm : COMM;



40

Putting it all together



- Scanner implemented as a function

```
int yylex();
```

 - return value indicates type of token found (encoded as a +ve integer);
 - the actual string matched is available in `ytext`.
- When compiling, link in the flex library using “-l” (or equivalent flag)
- Scanner and parser need to agree on token type encodings
 - let yacc generate the token type encodings
 - yacc places these in a file `y.tab.h`
 - use “`#include y.tab.h`” in the definitions section of the flex input file.

A quick tutorial on Lex

41