

# Lexical Analysis and Lexical Analyzer Generators [Chapter 3 - Part 2]

Software used for the project  
lex/flex

## Lecture 9

*Edited by Dr. Ismail Hababeh*

*German Jordanian University*

*Adapted from slides by Dr. Robert A. Engelen*

# The Lexical Analyzer Generators

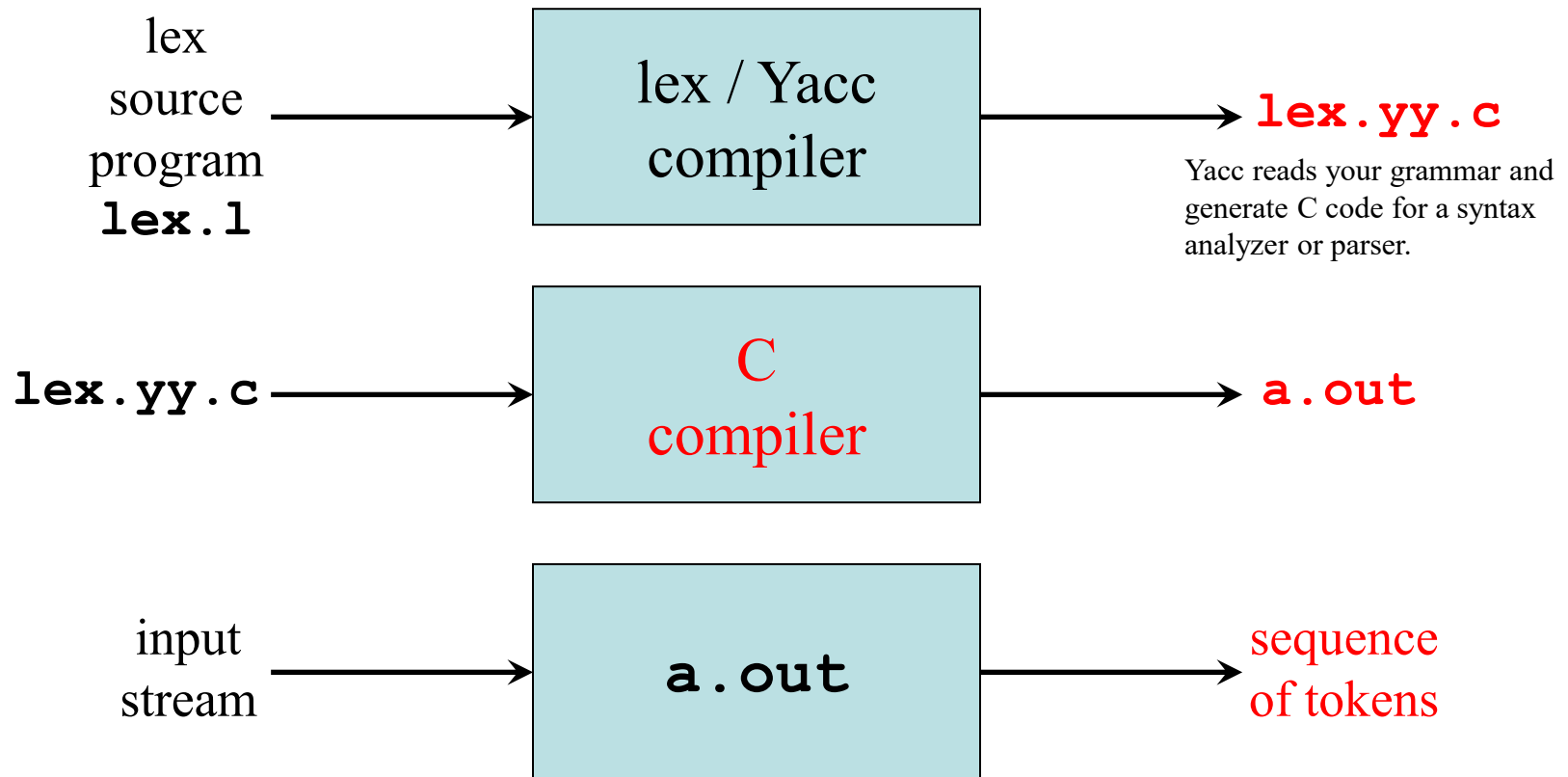
- Lexical Analyzer Generator LAG (**scanner /tokenizer**) is a **program** designed to **recognize lexical patterns in text**.
- LAG is used to **identify specific text strings** in a structured way that converts the source program into **meaningful tokens**.
- LAG tools: Lex and Flex.
  - Lex "lexical analyzer generator".
  - Flex "Fast lexical analyzer generator".
- Lex/Flex translates **regular definitions** into **C source code** for efficient scanning.

# The Parser Generators

- Parser Generator is a program used to produce a **parser for a given grammar**.
- Yacc/Bison is a grammar parser that generates a parse tree from the tokens generated by Lex/Flex.
- Yacc/Bison **reads your own grammar** and **generates C code** for a syntax analyzer or parser.
- The input of Yacc/Bison is the user rules or grammar and the **output is a C program**.

# Creating a Lexical Analyzer with Lex and Yacc

Lex and Yacc are tools used to generate lexical analyzers and parsers



# Lex Program Parts

- Declarations
  - Auxiliary declarations
  - Regular definitions
- Rules
  - Patterns followed by actions to be executed
- Auxiliary Procedures

# Lex Program Specification

- The *lex program parts specifications*:

*% { declarations % }*

*regular definitions*

*% %*

*translation rules*

*% %*

*user-defined auxiliary procedures*

- The *translation rules* are of the form:

Patterns	{	<i>p<sub>1</sub></i>	<i>{ action<sub>1</sub> }</i>	}	actions
		<i>p<sub>2</sub></i>	<i>{ action<sub>2</sub> }</i>		
		<i>...</i>			
		<i>p<sub>n</sub></i>	<i>{ action<sub>n</sub> }</i>		

# Regular Expressions in Lex

7

<b>x</b>	match the character <b>x</b>
<b>\.</b>	match the character <b>.</b>
<b>"string"</b>	match contents of string of characters
<b>.</b>	match any character except newline
<b>^</b>	match beginning of a line
<b>\$</b>	match the end of a line
<b>[xyz]</b>	match one character <b>x</b> , <b>y</b> , or <b>z</b>
<b>[^xyz]</b>	match any character except <b>x</b> , <b>y</b> , and <b>z</b>
<b>[a-z]</b>	match one character <b>a</b> to <b>z</b>
<b>r*</b>	closure (match zero or more occurrences)
<b>r+</b>	positive closure (match one or more occurrences)
<b>r?</b>	optional (match zero or one occurrence)
<b>r<sub>1</sub>r<sub>2</sub></b>	match <b>r<sub>1</sub></b> then <b>r<sub>2</sub></b> (concatenation)
<b>r<sub>1</sub> r<sub>2</sub></b>	match <b>r<sub>1</sub></b> or <b>r<sub>2</sub></b> (union)
<b>( r )</b>	grouping
<b>r<sub>1</sub>\r<sub>2</sub></b>	match <b>r<sub>1</sub></b> when followed by <b>r<sub>2</sub></b>
<b>{d}</b>	match the regular expression defined by <b>d</b>

# Lex Program – Example

// Declarations Part

%{

#include <stdio.h>

int global\_variable;

%}

// Regular definitions

number [0-9]<sup>+</sup>

opr [-|+|\*|/|^|=]

%%

// Rules

%%

// Auxiliary procedures



# Lex Rules - Example

%%

```
{number} { printf("number");}
```

```
{opr}      { printf("operator");}
```

%%

- LEX generates C code for the user rules and places this code into a single function called **yylex()**.

# Lex Auxiliary Procedures

```
int main ()  
{  
yylex();  
}
```

- This part allows to add user code to the *lex.yy.c* file.
- Once the code is written, *lex.yy.c* is generated using the command: `>> lex "filename.l"`
- *lex.yy.c* is compiled using the command:  
`>> gcc lex.yy.c`

# Lex Input/Output – Example 1

Based on the previous example:

Input: 12

Output: number

Input: +

Output: operator

Input: 1+2

Output: number operator number

# Lex Program – Example 2

Translation  
rules



```
%{  
#include <stdio.h>  
%}  
  
%%  
CS    printf("Computer science \n");  
CE    printf("Computer Engineering\n");  
ECE   printf("Electrical & Comm. Eng.\n");  
%%  
...
```

```
..> lex school.1 //generate lex.yy.c  
..> gcc lex.yy.c -ll //generate a.out  
..> ./a.out //generate executable file
```

In gcc, -ll is used to tell the compiler to use the **lex** libraries

# Lex Specification – Example 3

Contains  
the matching  
lexeme

(pointer to the  
beginning of the lexeme)

Translation  
rules

```
%{  
#include <stdio.h>  
%}  
  
%%  
[0-9]+  { printf("%s\n", yytext); }  
.|\\n   { }  
%%  
main()  
{ yylex();  
}
```

Invokes  
the lexical  
analyzer  
(returns the  
token)

# Lex Specification – Example 5

Translation rules

Regular definitions

```

%{
#include <stdio.h>
%}
digit      [0-9]
letter     [A-Za-z]
id         {letter} ({letter}|{digit})*
%%
{digit}+    { printf("number: %s\n", yytext); }
{id}        { printf("ident: %s\n", yytext); }
.           { printf("other: %s\n", yytext); }
%%
main()
{ yylex();
}
  
```

# Lex Specification – Example 6

```
%{ /* definitions of manifest constants */
#define LT (256)
```

```
...
```

```
%}
```

```
delim    [ \t\n]
```

```
ws       {delim}+
```

```
letter   [A-Za-z]
```

```
digit    [0-9]
```

```
id        {letter}({letter}|{digit})*
```

```
number    {digit}+(\.{digit}+)?(E[+\-]?{digit}+)?
```

```
%%
```

```
{ws}     { }
```

```
if        {return IF;}
```

```
then      {return THEN;}
```

```
else      {return ELSE;}
```

```
{id}      {yyval = install_id(); return ID;}
```

```
{number}  {yyval = install_num(); return NUMBER;}
```

```
"<"       {yyval = LT; return RELOP;}
```

```
"<="      {yyval = LE; return RELOP;}
```

```
"="       {yyval = EQ; return RELOP;}
```

```
"<>"      {yyval = NE; return RELOP;}
```

```
">"       {yyval = GT; return RELOP;}
```

```
">="      {yyval = GE; return RELOP;}
```

```
%%
```

```
int install_id(){/* function to install the lexeme whose 1st character
                  is pointed to by yytext and whose length is yyleng, into the
                  symbol table and return a pointer to it*/}
```

```
int install_num(){/* Similar to install_id(), but puts numerical
                   contents into a separate table */}
```

Return  
token to  
parser

Token  
attribute

the token attribute is stored in the shard variable yyval

yvalue stores the attribute value that is returned to the parser