# TDDD55- Compilers and Interpreters
# Lesson 2

Zeinab Ganjei (zeinab.ganjei@liu.se)

Department of Computer and Information Science

Linköping University

# 1. Grammars and Top-Down Parsing

- Some grammar rules are given
- Your task:
  - Rewrite the grammar
  - Add attributes and attribute rules to the grammar
  - Implement your grammar in a C++ class named **Parser**. The **Parser** class should contain a method named **Parse** that returns the value of a single statement in the language.

# 2. Scanner Specification

- Finish a scanner specification given in a *scanner.l flex* file, by adding rules for C and C++ style comments, identifiers, integers, and floating point numbers.
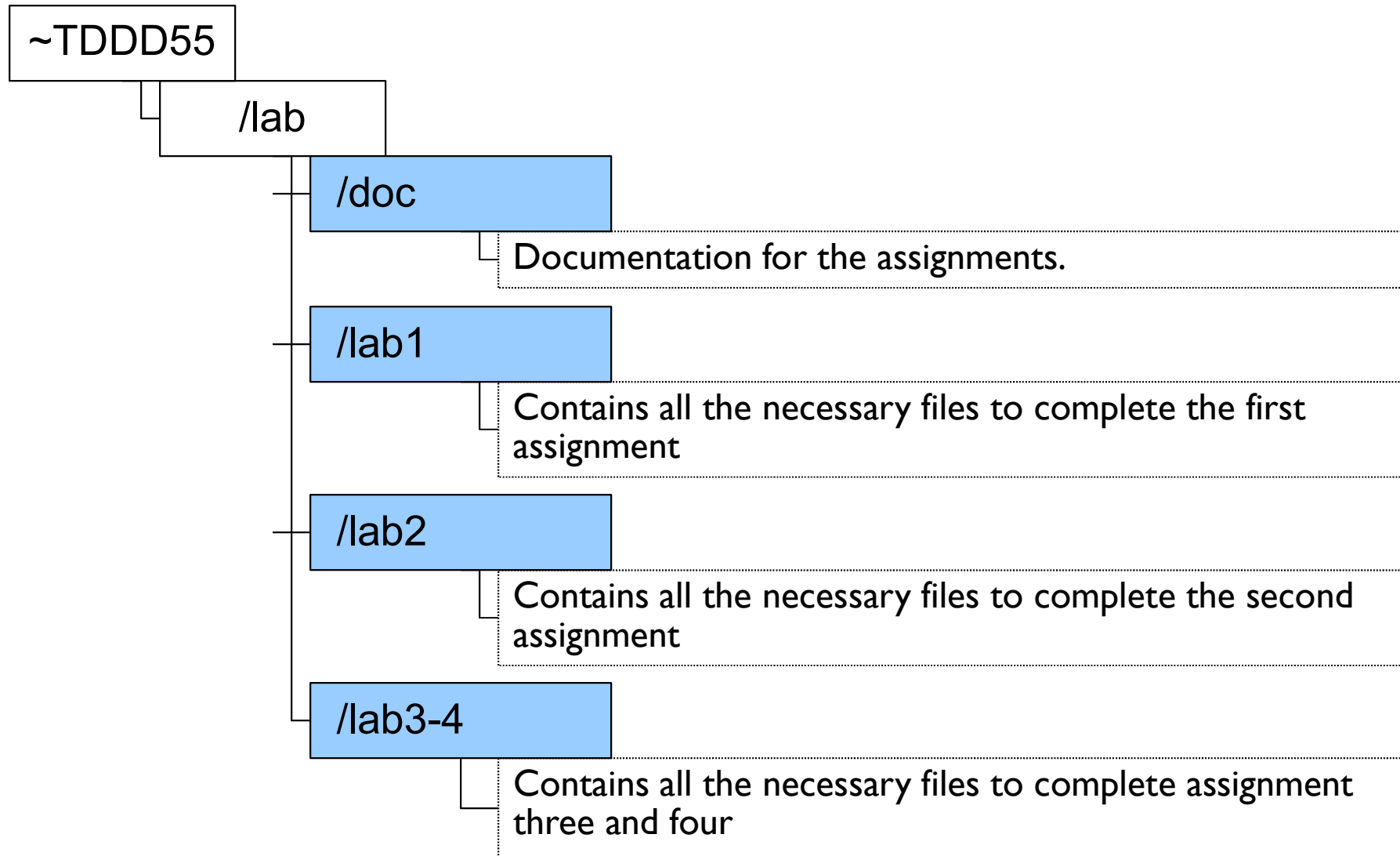
# 3. Parser Generators

- Finish a parser specification given in a *parser.y* <span style="color:red">bison</span> file, by adding rules for expressions, conditions and function definitions, .... You also need to augment the grammar with error productions.

- More on bison in lesson 3.

# 4. Intermediate Code Generation

- The purpose of this assignment is to learn about how abstract syntax trees can be translated into intermediate code.

- You are to finish a generator for intermediate code by adding rules for some language statements.

- More in lesson 3.

# Laboratory Skeleton

~TDDD55
/lab
/doc
Documentation for the assignments.
/lab1
Contains all the necessary files to complete the first assignment
/lab2
Contains all the necessary files to complete the second assignment
/lab3-4
Contains all the necessary files to complete assignment three and four

# Installation

- Take the following steps in order to install the lab skeleton on your system:
  - Copy the source files from the course directory onto your local account:

    ```
    mkdir TDDD55
    cp -r ~TDDD55/lab TDDD55
    ```

  - You might also have to load some modules (more information in the laboratory instructions).

# Today

- Introduction to the flex scanner generator tool.
- Introduction to laboratory assignment 2.
- Exercises in formal languages and automata theory.

# Flex

# Scanners

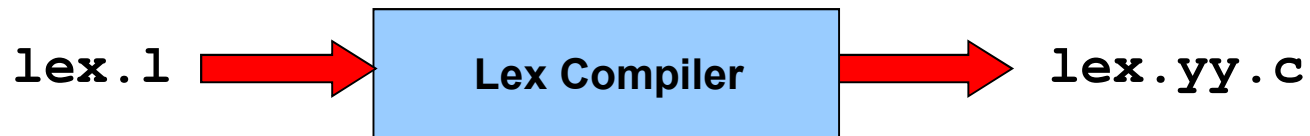**Scanners are programs that recognize lexical patterns in text**

- Its **input** is text written in some language.

- Its **output** is a sequence of tokens from that text. The tokens are chosen according to the language.

- Building a scanner manually is tedious.

- Mapping the regular expressions to finite state machine/automata is straightforward, so why not automate the process?

- Then we just have to type in regular expressions and actions and get the code for a scanner back.

# Scanner Generators

- Automation is exactly what **flex** does!

- **flex** is a fast lexical analyser generator, a tool for generating programs that perform **pattern matching** on text

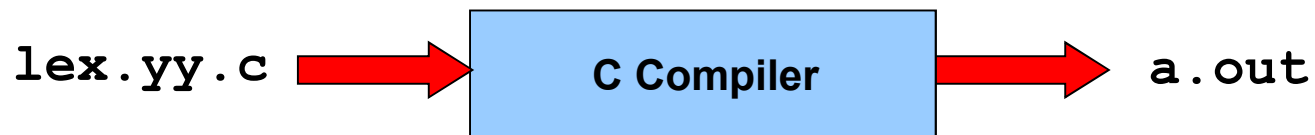- **flex** is a free implementation of the well-known **lex** program

# How it works

**flex** generates at output a **C** source file lex.yy.c
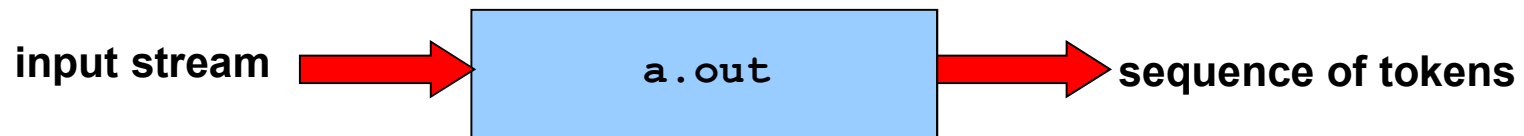which defines a routine yylex()

`lex.l` ➡️ **Lex Compiler** ➡️ `lex.yy.c`

`>> flex lex.l`

# How it works

lex.yy.c is compiled and linked with the -lfl library to produce an executable, which is the scanner

`lex.yy.c` ➡️ [ **C Compiler** ] ➡️ `a.out`

```
>> g++ lex.yy.c -lfl
```

**input stream** ➡️ [ `a.out` ] ➡️ **sequence of tokens**

```
>> a.out < input.txt
```

# Flex Specifications

## Lex programs are divided into three components

```
/* Definitions – name definitions
 *            – variables defined
 *            – include files specified
 *            – etc
 */
%%


/* Translation rules – regular expressions together with actions in C/C++ */


%%


/* User code – support routines for the above C/C++ code */
```

# 1. Name Definitions

- *Definitions* are intended to simplify the scanner specification and have the form:

  **name** **definition**

- Subsequently the definition can be referred to by {**name**}, which then will expand to the **definition**.

- Example:

  **DIGIT** **[0-9]**
  **{DIGIT}+".".{DIGIT}***

  is identical/will be expanded to:

  **([0-9])+".".([0-9])***

# Flex Specifications

## Lex programs are divided into three components

```
/* Definitions – name definitions
 *           – variables defined
 *           – include files specified
 *           – etc
 */
%%


/* Translation rules – regular expressions together with actions in C/C++ */


%%


/* User code – support routines for the above C/C++ code */
```

# 2. Pattern Actions

- The *translation rules* section of the **flex** input file, contains a series of rules of the form:

  **pattern  action**

- Example:

  **[0-9]*   { printf ("%s is a number", yytext);     }**

# Flex Matching

- Match as much as possible.

- If more than one rule can be applied, then the **<span style="color:red">first</span> appearing** in the flex specification file is preferred.

# Simple Patterns

- Match only one specific character

**"x"**   The character '**x**'

**.**    Any character except newline

# Character Class Patterns

- Match any character within the class

**[xyz]**   The pattern matches either '**x**', '**y**',   or '**z**'

**[abj-o]**  This pattern spans over a range of

characters and matches '**a**', '**b**', or

any letter ranging from '**j**' to '**o**'

# Negated Patterns

Match any character not in the class

**[^z]**      This pattern matches any character
            EXCEPT **z**


**[^A-Z]**    This pattern matches any character
            EXCEPT an uppercase letter


**[^A-Z\n]**  This pattern matches any character
            EXCEPT an uppercase letter or a
            newline

# Some Useful Patterns

**r\***  Zero or more '**r**', '**r**' is any regular expr.

**\\0**   **NULL** character (ASCII code 0)

**\123** Character with octal value **123**

**\x2a** Character with hexadecimal value **2a**

**p|s**   Either '**p**' or '**s**'

**p/s**   '**p**' but only if it is followed by an '**s**',
        which is not part of the matched text

**^p**    '**p**' at the beginning of a line

**p$**    '**p**' at the end of a line, equivalent to '**p/\n**'

```
/* Definitions – name definitions
 *          – variables defined
 *          – include files specified
 *          – etc
 */
%%

/* Translation rules – regular expressions together with actions in C/C++ */

%%

/* User code – support routines for the above C/C++ code */
```

# 3.Flex User Code

- The presence of this user code is optional.

- Finally, the _user code_ section is simply copied to lex.yy.c verbatim. It is used for companion routines which call, or are called by the scanner.

- If the lex program is to be used on its own, this section will contain a **main** function. If you leave this section empty you will get the default main.

# Flex Program Variables and Counters

**yytext**   Whenever the scanner matches a token, the text of the token is stored in the null terminated string yytext

**yyleng**   The length of the string yytext

**yylex()**   The scanner created by the Lex has the entry point yylex(), which can be called to start or resume scanning. If **lex** action returns a value to a program, the next call to yylex() will continue from the point of that return

# Flex Program Variables and Functions

**yymore()**   Do another match and append its
result to the current yytext (instead of replacing it)

**yyless(int n)**   Push all but the first n characters
back to the input stream (to be
matched next time). yytext will
contain only the first n of the
matched characters.

# yymore() Example

What will be printed by the following scanner if the input string is "hypertext"?

```
%%

"hyper" yymore();

"text" {  printf("Token is %s\n", yytext);  }
```

The output will be "Token is hypertext".

# Flex Examples

# Example: Recognition of Verbs

```
%{
/* includes and defines should be stated in this section */
%}

%%

[ \t]+              /* ignore white space */

"do"|"does"|"did"|"done"|"has" { printf ("%s: is a verb\n", yytext); }
[a-zA-Z]+           { printf ("%s: is not a verb\n",yytext); }
.|\n                { ECHO; /* normal default anyway */ }

%%

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

Mary has a little lamb

# Example: Character Counting

Write a scanner that counts the number of characters and lines in its input

```
int num_lines = 0, num_chars = 0; /* Variables */

%%

\n     { ++num_lines; ++num_chars; } /* Take care of newline */
.      { ++num_chars; }     /* Take care of everything else */

%%
main() { yylex();
  printf("lines: %d, chars: %d\n", num_lines, num_chars );
}
```

# Example: HTML Tags

Write a scanner that counts the number of characters and lines in its input

```
/*Declarations */
%{
#include <stdio.h>
%}

/*Exclusive start condition, i.e. only rules specific to <html_tag> will match */
%x html_tag
%%

[^<]*                    /* matches any char (zero or more times) except "<" */

"<"                      BEGIN(html_tag);  /* activate start condition <html_tag> */

<html_tag>[^>]*          printf("%s\n", yytext);

<html_tag>">"            BEGIN(INITIAL); /* Return to initial/original state*/
%%
```

- More details about Flex syntax and semantic:

http://www.ida.liu.se/~TDDD55/laboratories/FLEX.ps

# Laboratory Assignment 2

# Laboratory Assignment 2

- Finish a scanner specification given in a *scanner.l* flex file.

- Add regular expressions for floating point numbers, integer numbers, C comments (both /* */ comments and // one line comments), identifiers, empty space, newline.

- Rules for the language keywords are already given in the *scanner.l* file. Add your rules below them.

# Comments

- Skip characters in comments, both single-line C++ comments and multi line C style comments.

- If the scanner sees /* within a C comment, it must print a warning message.

- If end of line is encountered within a C style comment, print an error message and then terminate.

# Comments Example

Rules for comments.

```
"//".*\n                        /* Do nothing */

"/*"                            BEGIN(c_comment)

<c_comment> {

"*/"                            … …
"/*"                            fprintf(stderr, "Warning: Nested comments\n");
… …

}
```

# Integers

- Integers are simply sequences of digits that are not part of identifiers or floating-point numbers.

# Floating Point Numbers

- Floating-point numbers consist of an **integer part** followed by a **decimal point**, **decimal part** and an **exponent part**.
  - e.g. 56.11E-2
- The integer and decimal parts are sequences of digits. The **exponent part** consists of the character *e* or *E* followed by an optional sign *+* or *-* and a sequence of digits.
- Either the integer or the decimal part (or both) must be given.
- The exponent is optional.
- If the integer part and exponent are both given, the decimal point and decimal part are optional.

# Floating Point Numbers Example

- 1.1
- .1
- 1.
- 1.1E2
- 2E-3
- 1E-4


- Use **?** as *optional* pattern. Example: [+-]?

# Identifiers

- Identifiers must start with a letter, followed by any number of digits, letters or underscore characters.

# Questions?