

Exercise 1 Syntax: Regular Expressions, Grammars and Scanners

PDF Updated on Oct 21, 2019

(Deadline for uploading solutions: Oct 27, 2019, 11:59pm Stuttgart time)

The materials provided for this homework are:

- a pdf file with the text of the homework (this);
- a zip file with the **folder structure and the templates that must be used for the submission.**

The folder structure is shown in Figure 1.

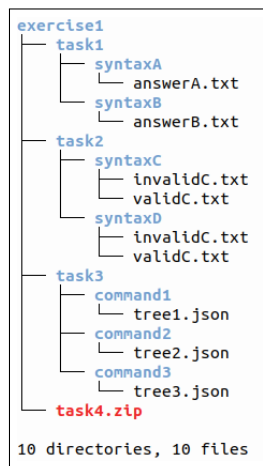


Figure 1: Folder structure to be used for the solution.

The submission must be compressed in a zip file using the given folder structure. The name of folders and files must not be changed or moved, otherwise the homework will not be evaluated.

There are four tasks, which contribute a specific percentage to the overall points for this exercise.

Notes about the symbols used in this exercise:

- *blue tokens* are non-terminals;
- *black tokens* are terminals;
- * is the Kleene star symbol;
- * is a terminal symbol (e.g., multiplication symbol);
- | means "or";

1 Task I (10% of total points of the exercise)

You are given two sets of production rules. For each set, the question is: Do these rules describe a regular expression or a grammar?

1.1 Syntax A

```

start → sentence
sentence → the qualified_noun verb | pronoun verb
qualified_noun → adjective noun
noun → student | programmer | engineer
pronoun → he | she
verb → talks | studies | works
adjective → expert | smart | new

```

Write the following into the file `exercise1/task1/syntaxA/answerA.txt` as show in Figure 2: **GRAMMAR** if the above is a grammar or **REGEX** if it is a regular expression (in capital letters).

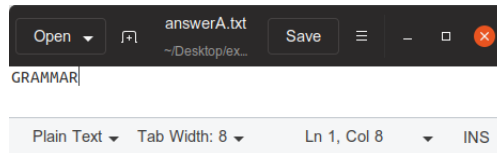


Figure 2: Example of the file containing the answer of task 1.

1.2 Syntax B

```

start → wishlist
wishlist → product $ price | product $ price , wishlist
product → car | computer | smartphone
price → non_zero_digit digit*
non_zero_digit → 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

```

Write the following into the file `exercise1/task1/syntaxB/answerB.txt` as show in Figure 2: **GRAMMAR** if the above is a grammar or **REGEX** if it is a regular expression (in capital letters).

2 Task II (10% of total points of the exercise)

Two grammars are given. Your task is to write three valid strings and three invalid strings for each of them.

2.1 Syntax C

```

start → expression
expression → term | expression - term
term → factor | term * factor
factor → x | y

```

where `*` is a terminal symbol and spaces are allowed.

The three valid strings must be written into the file `exercise1/task2/syntaxC/validC.txt`, one string per line. An example is shown in Figure 3.

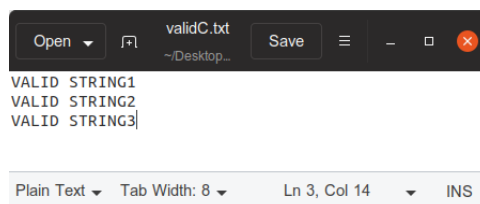


Figure 3: Example of the file containing the three valid strings.

The three invalid strings must be written into the file `exercise1/task2/syntaxC/invalidC.txt`, one string per line.

2.2 Syntax D

```
start → pathfilename
pathfilename → drive : \ pathname* filename fileextension
drive → letter
pathname → id \
filename → id
fileextension → .id
id → letter* | letter* digit*
digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
letter → a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z
```

where * is the Kleene star symbol and spaces are not allowed in all the rules.

The three valid strings must be written into the file `exercise1/task2/syntaxD/validC.txt`, one string per line. (`validC.txt` is the right file here too, but in the folder `syntaxD/`)

The three invalid strings must be written into the file `exercise1/task2/syntaxD/invalidC.txt`, one string per line.

3 Task III (30% of total points of the exercise)

A correct grammar and three correct examples of programs are proposed. The goal of this exercise is to create a parse tree and provide it in a JSON file.

3.1 Syntax E

```
start → command
command → program
program → convert | run | comparison | exit now
exit → shutdown -h
convert → convert argument filename . extension
argument → - RAWtoPNG | - PNGtoJPEG | - JPEGtoPNG
extension → png | jpeg | raw
run → python filename . py
comparison → diff filename .txt filename .txt > filename .txt
filename → letter*(letter|digit)*
digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
letter → a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z
```

where * is the Kleene star symbol.

A example of legal command is: `shutdown -h now`

A parse tree generated from this string is show in Figure 4.



Figure 4: Example of the tree created from the `shutdown -h now` string.

The JSON file that can represent this parser tree is:

```
{
  "id": "start",
  "children": [
    {
      "id": "command",
      "children": [
        {
          "id": "program",
          "children": [
            {
              "id": "exit",
              "children": [
                {
                  "id": "shutdown -h",
                  "children": []
                }
              ]
            }
          ]
        },
        {
          "id": "now",
          "children": []
        }
      ]
    }
  ]
}
```

The goal of this exercise is to create the parse tree of the following three commands. Please use the above JSON format to create the three requested parser trees.

1. `convert -RAWtoPNG newcar2019.raw`
2. `python helloworld.py`
3. `diff oldfile.txt newfile.txt > output.txt`

Write the parse trees into the following JSON files:

- parse tree of command 1 in the file `exercise1/task3/command1/tree1.json`
- parse tree of command 2 in the file `exercise1/task3/command2/tree2.json`
- parse tree of command 3 in the file `exercise1/task3/command3/tree3.json`

4 Task IV (50% of total points of the exercise)

You are given regular expressions that describe the legal tokens of a language:

```
keyword → if | else | for | while | return
punctuator → ( | ) | [ | ] | { | } | . | ; | |
comparison → < | > | == | !=
assign → <=
op → + | - | * | / | && | ||
id → letter*(letter | digit)* (except for keywords)
number → digit*
comment → //(non_newline)* newline
```

where *non_newline* means all characters other than *newline* and *letter* is a lowercase or uppercase letter. The *keywords* are case-sensitive.

The aim of this exercise is to implement a scanner in Java that transforms a string in the language into a sequence of tokens, or reports an error if the string is illegal. You must implement the scanner by hand, i.e., do not use a scanner generator tool.

The *comments* must be skipped, so they don't appear on the output. All white space except for newlines must be skipped. For example, the following two snippets of code correspond to the same sequence of tokens.

```
1 if(x > 0){
2     y <= x;
3 }
4 else{
5     y <= y / 2;
6 }
7
8 //This is a comment.
9 for(i <= 0; i < 10; i++){
10     y <= 2 * y;
11 }
```

```
12 if (x>0){y<=x;}
13 else {y<=y/2;}
14 //This is a comment.
15 for (i <=0;i <10;i++){y<=2*y; }
```

An Eclipse project for the scanner implementation is provided: *exercise1/task4/*. You can import the project template (.zip) into Eclipse as follows: *File-Import-General-Existing project into Workspace-Select archive file-Finish* .

Please, implement your scanner in the method *public static List < String > functionScanner(String input)*, which you find in file *exercise1/task4/src/scanner/TokenScanner.java*.

The input to the scanner is a *String* containing a snippet of code. The output is a *List < String >* with all the tokens found by the a scanner. In case of illegal tokens, the output must be a *List < String >* with a single *String* element: "Illegal".

You find some JUnit tests in folder *exercise1/task4/src/scanner/*. Use them to check whether your scanner works. They can be run using Eclipse. There is a little mistake in the Testcase2: *for(i <= 0; i < 10; i = i + 1)* should be *for(i <= 0; i < 10; i <= i + 1)*. We will use additional tests to evaluate your solution, and you are advised to also add additional tests for your own testing.

For the submission, your Eclipse project must be exported into a .zip archive using Eclipse: *File-Export-General-Archive File*, and then added to the tree structure in Figure 1.