

**MINISTRY OF EDUCATION AND RESEARCH OF REPUBLIC OF MOLDOVA**  
**TECHNICAL UNIVERSITY OF MOLDOVA**  
**FACULTY OF COMPUTERS, INFORMATICS AND MICROELECTRONICS**  
**DEPARTMENT OF SOFTWARE ENGINEERING AND AUTOMATICS**

**Laboratory work 4:**  
**Regular expressions**

**Course: Formal Languages & Finite Automata**

**Author: Bostan Victor, FAF-222**

**Chişinău, 2024**

# THEORY

Regular expressions (regex) are a powerful tool used in computing for matching patterns within text. They originate from formal language theory and automata theory in computer science and are commonly used for searching, editing, and manipulating text. A regular expression is essentially a sequence of characters that forms a search pattern. It can be used for everything from finding and replacing text in a word processor to validating the format of email addresses or phone numbers in a database.

The syntax of regular expressions allows them to be highly versatile and expressive. Patterns can range from simple, such as finding a single word, to complex, such as identifying email addresses or validating user input. In programming, regular expressions are used in string searching algorithms for "find" or "find and replace" operations, as well as for input validation. They are implemented in various programming languages, including Python, Java, and JavaScript, and are utilized in many contexts, from software development and web development to data analysis and network security.

The use of regular expressions enables programmers and data scientists to work more efficiently by providing a concise and flexible means for matching strings of text. For example, they can quickly sift through large datasets to find relevant information or validate user inputs against expected patterns. Despite their power, regular expressions can be complex and may require a learning curve to understand and use effectively. However, mastering them can greatly enhance one's ability to process and analyze text, making them an invaluable tool in the arsenal of modern computing.

# OBJECTIVES

1. Write and cover what regular expressions are, what they are used for;
2. Below you will find 3 complex regular expressions per each variant. Take a variant depending on your number in the list of students and do the following:
  - a. Write a code that will generate valid combinations of symbols conform given regular expressions (examples will be shown).
  - b. In case you have an example, where symbol may be written undefined number of times, take a limit of 5 times (to evade generation of extremely long combinations);
  - c. Bonus point: write a function that will show sequence of processing regular expression (like, what you do first, second and so on)

## IMPLEMENTATION DESCRIPTION

For this laboratory we are given 3 regular expressions for which we have to generate valid strings. Here is my variant, Variant 1:

Variant 1:

$(a|b)(c|d)E^+G?$

$P(Q|R|S)T(UV|W|X)^*Z^+$

$1(0|1)^*2(3|4)^536$

Once again, I am going to complete the tasks in Python. The first thing I did was to define the regular expressions in string form (will be needed for the explanation in the bonus point) and in token form in a way which will help me generate the strings.

```
110 # Define example regular expressions and their tokenized forms
111 reg_expressions = [
112     '(a|b)(c|d)E+G?',
113     'P(Q|R|S)T(UV|W|X)*Z+',
114     '1(0,1)*2(3|4)^536'
115 ]
116 reg_expressions_tokens = [ # Regular expressions tokens
117     [(1, ['a', 'b']), (1, ['c', 'd']), ('+', ['E']), ('?', ['G'])],
118     [(1, ['P']), (1, ['Q', 'R', 'S']), (1, ['T']), (*, ['UV', 'W', 'X']), (+, ['Z'])],
119     [(1, ['1']), (*, ['0', '1']), (1, ['2']), (5, ['3', '4']), (1, ['36'])]
120 ]
121 limit = 5 # Limit for symbols written an undefined number of times
122 n = 5 # Number of strings to generate
```

The idea is the following. Let's take the first expression:  $(a|b)(c|d)E^+G?$ . I split this expression into tokens as following:  $(a|b)$ ,  $(c|d)$ ,  $E^+$ ,  $G?$ . Each token represents a step in generating the string and after iterating through every step a valid word is generated. In my code, in the list of the regular expression tokens, I represent the tokens as tuples where the first element of the tuple is the number of repetitions a symbol is allowed to have (1, ? - 0 or 1, 5, + - 1 and more, \* - 0 and more), and the second element is the token itself. The next step is implementing the actual function which will iterate through the tokens and generate valid strings for any regular expression given as input:

```

4  # Function to generate strings based on provided regular expression patterns
5  def generate_strings(reg_exp, n, limit):
6      strings = [] # List to store the generated strings
7
8      for i in range(n):
9          string = '' # Initialize the string to be built
10         # Iterate through each token in the regular expression
11         for tup in reg_exp:
12             if tup[0] == '1': # '1' denotes a mandatory symbol
13                 string += random.choice(tup[1])
14             elif tup[0] == '?': # '?' denotes an optional symbol
15                 nr = random.randint(0, 1)
16                 if nr:
17                     string += random.choice(tup[1])
18             elif tup[0] == '5': # '5' denotes a symbol repeated exactly 5 times
19                 char = random.choice(tup[1])
20                 for j in range(5):
21                     string += char
22             elif tup[0] == '+': # '+' denotes one or more repetitions of a symbol
23                 char = random.choice(tup[1])
24                 nr_of_chars = random.randint(1, limit)
25                 for j in range(nr_of_chars):
26                     string += char
27             else: # For other types, assume '*', denoting zero or more repetitions
28                 char = random.choice(tup[1])
29                 nr_of_chars = random.randint(0, limit)
30                 for j in range(nr_of_chars):
31                     string += char
32
33         strings.append(string)
34
35     return strings

```

As you can see in the image, the *generate\_strings* function first initializes an empty list *strings* for the output and generates 5 strings. A string is generated by first iterating through the tuples and for each one checking what is the number of repetitions allowed. If the allowed repetitions is 1 then we just append to the string a random symbol (or the only one if the list has length 1) from the token list. If the allowed repetitions is ?, we first have to generate a random number between 0 and 1 to check if we will append a symbol or not. If a 1 is generated then we append a random symbol from the token list. If the number of repetitions is 5, then we append exactly 5 random symbols. If the allowed repetitions is +, then we first generate the number of times the symbol is going to be appended, from 1 to limit (as it was in the condition) and append a random symbol from the token that number of times. If the allowed repetitions is \*, then it's the same process as with +, but we generate a number between 0 and the limit. Here are some examples on how the function generates the strings (I chose  $n = 5$  for simplicity):

```
5 random strings for Regular Expression (a|b)(c|d)E+G?:  
['acEEEG', 'bdEE', 'bdEEG', 'bdEEE', 'bcEEEEEE']  
  
5 random strings for Regular Expression P(Q|R|S)T(UV|W|X)*Z+:  
['PRTXXXXXZ', 'PQTUVUVUVVZZZZ', 'PSTXXXXXZ', 'PSTXXXXXZ', 'PSTXXXXZZZZ']  
  
5 random strings for Regular Expression 1(0,1)*2(3|4)^36:  
['10024444436', '1000023333336', '10000024444436', '100023333336', '10024444436']
```

```
5 random strings for Regular Expression (a|b)(c|d)E+G?:  
['bdEEEEEE', 'bdEEEEE', 'acEEE', 'adEEEEE', 'bcEEEEE']  
  
5 random strings for Regular Expression P(Q|R|S)T(UV|W|X)*Z+:  
['PRTXXXXZZZ', 'PRTUVUVUVUVVZZ', 'PSTUVUVUVZZZZ', 'PQTXZZZZ', 'PRTUVZ']  
  
5 random strings for Regular Expression 1(0,1)*2(3|4)^36:  
['1124444436', '1124444436', '111124444436', '10000024444436', '111123333336']
```

```
5 random strings for Regular Expression (a|b)(c|d)E+G?:  
['bdEEEEEE', 'bcEE', 'adEEEEE', 'bdEEE', 'bcEEEEE']  
  
5 random strings for Regular Expression P(Q|R|S)T(UV|W|X)*Z+:  
['PRTXXXZZZ', 'PSTWWWZZZ', 'PRTWWWZZZZZ', 'PQTUVUVVZZZ', 'PRTWWWZZ']  
  
5 random strings for Regular Expression 1(0,1)*2(3|4)^36:  
['1124444436', '123333336', '111123333336', '1124444436', '1244444436']
```

It is easy to check that every string that was generated belongs to their respective regular expressions. Now, all that remains to be done is the bonus task. For the bonus task we had to create a function that provides the sequence of the processing of the string generation. So, I created a function which generates a random string based on the regular expression given as input, and also provides a step-by-step explanation on how it is done:

```

38 # Function to generate a string with a step-by-step explanation
39 def generate_string_with_explanation(exp, reg_exp, n, limit):
40     string = ''
41     step = 1 # Step counter for explanation
42     print(f'\nGenerating a string for Regular Expression {exp}:\n')
43     print(f'Step-by-step explanation:')
44
45     for tup in reg_exp:
46         l = len(tup[1]) # Get the number of options for the current component
47         if tup[0] == '1':
48             string += random.choice(tup[1])
49             if l == 1:
50                 print(f'Step {step}. Append 1 instance of the symbol "{tup[1][0]}". String = {string}')
51                 step += 1
52             else:
53                 print(f'Step {step}. Append 1 instance of one of these symbols - {"", ".join(tup[1])}. String = {string}')
54                 step += 1
55         elif tup[0] == '?':
56             nr = random.randint(0, 1)
57             print(f'Step {step}. Generate a random number between 0 and 1 to determine if a symbol is going to be appended. Generated number = {nr}')
58             step += 1
59             if nr:
60                 string += random.choice(tup[1])
61                 if l == 1:
62                     print(f'Step {step}. Append 1 instance of the symbol "{tup[1][0]}". String = {string}')
63                     step += 1
64                 else:
65                     print(f'Step {step}. Append 1 instance of one of these symbols - {"", ".join(tup[1])}. String = {string}')
66                     step += 1
67             else:
68                 print(f'Step {step}. Nothing is being appended. String = {string}')
69                 step += 1

```

```

70         elif tup[0] == '5':
71             char = random.choice(tup[1])
72             for j in range(5):
73                 string += char
74             if l == 1:
75                 print(f'Step {step}. Append 5 instances of the symbol "{tup[1][0]}". String = {string}')
76                 step += 1
77             else:
78                 print(f'Step {step}. Append 5 instances of one of these symbols - {"", ".join(tup[1])}. String = {string}')
79                 step += 1
80         elif tup[0] == '+':
81             char = random.choice(tup[1])
82             nr_of_chars = random.randint(1, limit)
83             print(f'Step {step}. Generate how many symbols will be appended between 1 and the limit, which is {limit}. In this case we append {nr_of_chars} symbols')
84             step += 1
85             for j in range(nr_of_chars):
86                 string += char
87             if l == 1:
88                 print(f'Step {step}. Append {nr_of_chars} instances of the symbol "{tup[1][0]}". String = {string}')
89                 step += 1
90             else:
91                 print(f'Step {step}. Append {nr_of_chars} instances of one of these symbols - {"", ".join(tup[1])}. String = {string}')
92                 step += 1
93         else:
94             char = random.choice(tup[1])
95             nr_of_chars = random.randint(0, limit)
96             print(f'Step {step}. Generate how many symbols will be appended between 0 and the limit, which is {limit}. In this case we append {nr_of_chars} symbols')
97             step += 1
98             for j in range(nr_of_chars):
99                 string += char
100             if l == 1:
101                 print(f'Step {step}. Append {nr_of_chars} instances of the symbol "{tup[1][0]}". String = {string}')
102                 step += 1
103             else:
104                 print(f'Step {step}. Append {nr_of_chars} instances of one of these symbols - {"", ".join(tup[1])}. String = {string}')

```

As you can see, I used the same method of generating strings as in the previous function. I just added *print* functions along the way to showcase each step of doing it. Here is an example of output:

```
Generating a string for Regular Expression P(Q|R|S)T(UV|W|X)*Z+:
```

```
Step-by-step explanation:
```

```
Step 1. Append 1 instance of the symbol "P". String = P
```

```
Step 2. Append 1 instance of one of these symbols - Q, R, S. String = PR
```

```
Step 3. Append 1 instance of the symbol "T". String = PRT
```

```
Step 4. Generate how many symbols will be appended between 0 and the limit, which is 5. In this case we append 4 symbols
```

```
Step 5. Append 4 instances of one of these symbols - UV, W, X. String = PRTXXXX
```

```
Step 6. Generate how many symbols will be appended between 1 and the limit, which is 5. In this case we append 2 symbols
```

```
Step 7. Append 2 instances of the symbol "Z". String = PRTXXXXZZ
```

```
The resulting string is: PRTXXXXZZ
```

## CONCLUSION

In conclusion, this laboratory work provided a practical understanding of regular expressions and their application in string generation and processing. By implementing functions to generate strings based on predefined patterns, I explored the versatility and power of regular expressions in automating text analysis and manipulation tasks. Through the step-by-step generation and explanation process, I deepened my comprehension of how different components of regular expressions affect string outcomes, enhancing my ability to construct and interpret regex patterns effectively. This hands-on experience reinforces the importance of regular expressions as a fundamental tool in software development and data processing. Now, as it comes to the difficulties met in this lab, there actually weren't any, as it is a pretty straight forward task.