



Министерство науки и высшего образования Российской Федерации

Федеральное государственное бюджетное образовательное учреждение высшего образования

«МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ

УНИВЕРСИТЕТ имени Н.Э.БАУМАНА

(национальный исследовательский университет)»

Факультет: Информатика и системы управления

Кафедра: Теоретическая информатика и компьютерные технологии

Лабораторная работа № 3

Раскрутка самоприменимого компилятора

Вариант 2

Работу выполнил

студент группы ИУ9-61

Бакланова А.Д.

Москва, 2020

Цель работы:

Целью данной работы является приобретение навыка разработки простейших лексических анализаторов, работающих на основе поиска в тексте по образцу, заданному регулярным выражением.

Исходные данные:

Стандартная библиотека любого современного языка программирования содержит средства для поиска в тексте образцов, заданных регулярными выражениями. При этом используется расширенный синтаксис записи регулярных выражений, позволяющий по сути выйти за рамки регулярных языков. Механизм поиска по таким регулярным выражениям годится для написания простейших лексических анализаторов. Однако, для этого механизма характерна нелинейная зависимость времени работы от длины распознаваемой лексемы, поэтому в промышленных компиляторах он не используется.

В качестве языка реализации в данной лабораторной работе выберем язык Java, стандартная библиотека которого содержит пакет `java.util.regex`, в котором располагаются классы `Pattern` и `Matcher`, предназначенные для поиска по регулярным выражениям. Документация по этому пакету находится по адресу:

<http://docs.oracle.com/javase/7/docs/api/java/util/regex/package-summary.html>.

Вводную статью по синтаксису регулярных выражений можно прочитать здесь:

<http://www.quizful.net/post/Java-RegExp>.

Идея лексического анализа на основе поиска по регулярным выражениям состоит в использовании групп, представляющих собой фрагменты регулярных выражений, заключённые в круглые скобки, значения которых запоминаются при сопоставлении текста с образцом. Например, на листинге 1 показано, как с использованием групп отличить идентификаторы от числовых литералов.

Индивидуальный вариант:

2	Комментарии: начинаются с «/*», заканчиваются на «*/» и могут пересекать границы строк текста. Идентификаторы: последовательности латинских букв и десятичных цифр, в которых буквы и цифры чередуются. Ключевые слова: «for», «if», «m1».
---	--

Задание:

В лабораторной работе необходимо реализовать на языке Java две первые фазы стадии анализа: чтение входного потока и лексический анализ. Чтение входного потока должно осуществляться из файла (в UTF-8), при этом лексический анализатор должен вычислять текущие координаты в обрабатываемом тексте. В результате работы программы в стандартный поток вывода должны выдаваться описания распознанных лексем в формате

Тег (координаты): значение

Например,

```
IDENT (1, 2): count
ASSIGN (1, 8): :=
NUMBER (1, 11): 100
```

Лексемы во входном файле могут разделяться пробельными символами (пробел, горизонтальная табуляция, маркеры окончания строки), а могут быть записаны слитно (если это не приводит к противоречиям).

Идентификаторы и числовые литералы не могут содержать внутри себя пробельных символов, если в задании явно не указано иного (варианты 4, 14 и 36). Комментарии, строковые и символьные литералы могут содержать внутри себя пробельные символы.

Входной файл может содержать ошибки, при обнаружении которых лексический анализатор должен выдавать сообщение с указанием координаты:

```
syntax error (10,2)
```

После обнаружения ошибки лексический анализатор должен восстанавливаться по следующей схеме: из входного потока пропускаются все подряд идущие символы до нахождения следующей лексемы.

Лексический анализатор должен иметь программный интерфейс для взаимодействия с парсером. Рекомендуется реализовывать его как итератор с методом `nextToken()` для императивных языков или функцию, возвращающую список лексем, для функциональных языков.

В регулярных выражениях рекомендуется использовать классы символов Unicode для обозначения букв, чисел и других подобных множеств. Многие движки регулярных выражений для задания классов используют синтаксис `\r{класс}`. Вместо нумерованных групп рекомендуется использовать именованные (`?<имя>regex`), при использовании нумерованных групп ненумеруемые обозначаются как (`? :regex`).

Варианты языков для лексического анализа приведены в таблицах 1, 2, 3, 4 и 5.

Реализация:

```
1  import java.util.Arrays;
2  import java.util.Scanner;
3  import java.util.function.IntPredicate;
4  import java.util.regex.Matcher;
5  import java.util.regex.Pattern;
6
7  class Position {
8      private String text;
9      private int amountRemoved;
10     private int index, line, col;
11
12     public Position(String text) {
13         this(text, 0, 1, 1);
14     }
15
16     private Position(String text, int index, int line, int col) {
17         this.text = text;
18         this.index = index;
19         this.line = line;
20         this.col = col;
21     }
22
23     public boolean isEmpty() {
24         return text.isEmpty();
25     }
26
27     public int getChar() {
28         return index < text.length() ? text.codePointAt(index) : -1;
29     }
30
31     public boolean satisfies(IntPredicate p) {
32         return p.test(getChar());
33     }
34
35     public Position skip() {
36         int c = getChar();
37         switch (c) {
38             case -1:
39                 return this;
40             case '\n':
41                 return new Position(text.substring(1), index/* + 1*/, line + 1, 1);
42             default:
43                 return new Position(text.substring(1), index/* + (c > 0xFFFF ? 2 : 1)*/, line, col + 1);
44         }
45     }
46
47     public Position skipToken(int range) {
48         Position pos = this;
49         for (int i = 0; i < range; i++) {
50             pos = pos.skip();
51         }
52         pos.text = text.substring(range);
53
54         return pos;
55     }
56
57     public Position skipWhile(IntPredicate p) {
58         Position pos = this;
59         while (pos.satisfies(p)) pos = pos.skip();
60         return pos;
61     }
62
63     public String getText() {
64         return text;
65     }
66
67     public String toString() {
68         return String.format("(%d, %d)", line, col);
69     }
70 }
```

```

70 }
71
72 class SyntaxError extends Exception {
73     public SyntaxError(Position pos, String msg) {
74         super(String.format("Syntax error at %s: %s", pos.toString(), msg));
75     }
76 }
77
78
79 enum Tag {
80     IDENT,
81     comments,
82     sm1,
83     sif,
84     sfor,
85     END_OF_TEXT;
86
87     public String toString() {
88         switch (this) {
89             case IDENT:
90                 return "IDENT";
91             case comments:
92                 return "comments";
93             case sm1:
94                 return "sm1";
95             case sif:
96                 return "sif";
97             case sfor:
98                 return "sfor";
99             case END_OF_TEXT:
100                 return "END OF TEXT";
101         }
102         throw new RuntimeException("unreachable code");
103     }
104 }
105
106 class Token {
107     private Tag tag;
108     private String value;
109     private Position start, follow;
110     private Pattern pattern;
111
112     public Token(String text, Pattern pattern) throws SyntaxError {
113         this(new Position(text), pattern);
114     }
115
116     private Token(Position cur, Pattern pattern) throws SyntaxError {
117         this.pattern = pattern;
118         start = cur.skipWhile(Character::isWhitespace);
119         follow = start.skip();
120
121         Matcher matcher = pattern.matcher(start.getText());
122
123         if (start.isEmpty()) {
124             tag = Tag.END_OF_TEXT;
125             value = "EOF";
126         } else if (!matcher.find()) {
127             throwError("invalid character");
128         } else if (matcher.group("ident") != null) {
129             tag = Tag.IDENT;
130             value = matcher.group("ident");
131             follow = follow.skipToken(matcher.end() - 1);
132         } else if (matcher.group("sm1") != null) {
133             tag = Tag.sm1;
134             value = matcher.group("sm1");
135             follow = follow.skipToken(matcher.end() - 1);
136         } else if (matcher.group("sif") != null) {

```


Тестирование:

```
/Library/Java/JavaVirtualMachines/jdk1.8.0_181.jdk/Contents/Home/bin/java ...  
/* AAAAAAA */ if for m1 a5v6n /*b*/  
^D  
comments (1, 1): /  
comments (1, 2): *  
comments (1, 4): A  
comments (1, 5): A  
comments (1, 6): A  
comments (1, 7): A  
comments (1, 8): A  
comments (1, 9): A  
comments (1, 10): A  
comments (1, 12): *  
comments (1, 13): /  
sif (1, 15): if  
sfor (1, 18): for  
sm1 (1, 22): m1  
IDENT (1, 25): a5v6n  
comments (1, 32): /  
comments (1, 33): *  
comments (1, 34): b  
comments (1, 35): *  
comments (1, 36): /  
  
Process finished with exit code 0
```

Вывод:

В результате выполнения лабораторной работы были получены навыки разработки простейших лексических анализаторов, работающих на основе поиска в тексте по образцу, заданному регулярным выражением, также была реализована программа, считывающая строку и выводящая описания распознанных лексем.