

Projeto e Implementação de um analisador léxico para a linguagem C-

Eduardo Kurek¹

¹Departamento Acadêmico de Computação (DACOM)
Universidade Tecnológica Federal do Paraná (UTFPR)

Abstract

Lexical analyzers are the first step in creating a compiler. They split the code into tokens, which are later used in subsequent compilation phases. This work aims to implement a lexical analyzer for the C- language using automata. For this implementation, a Mealy machine was employed. As it is equivalent to a Moore machine, it can also work similarly. After the lexical analysis phase, the output will be a list of tokens that describe the code structure.

Resumo

Analisadores léxicos são a primeira etapa na criação de um compilador. Eles separam o código em marcas, usadas posteriormente nas próximas fases de compilação. Este trabalho tem por finalidade a implementação de um analisador léxico para a linguagem C- usando autômatos. Neste trabalho usei a máquina de *Mealy* para a implementação, como ela é equivalente à máquina de *Moore*, também funcionará. Após a fase de análise léxica, teremos como saída uma lista de marcas, que descreve o código.

1 Introdução

Analisadores léxicos são programas capazes de identificar uma sequência de caracteres (lexemas) e converter em um *token* (marca). Os tokens que serão identificados pela linguagem C- estão descritos na Tabela 1. Além de identificar números positivos e negativos (NUMBER), variáveis que possuem letras e números (ID), e ignorar comentários entre barra e asterisco (*/* Exemplo de comentário */*).

1. Exemplo de código na linguagem C-

```

1  int gcd (int u, int v){
2      if (v == 0) return u;
3      else return gcd(v,u-u/v*v);
4      /* u-u/v*v == u mod v */
5  }
6
7  void main(void){
8      int x; int y;
9      x = input();
10     y = input();
11     output(gcd(x,y));
12 }
```

Lexema	Token
if	IF
else	ELSE
int	INT
float	FLOAT
return	RETURN
void	VOID
while	WHILE
+	PLUS
-	MINUS
*	TIMES
/	DIVIDE
<	LESS
<=	LESS_EQUAL
>	GREATER
>=	GREATER_EQUAL
==	EQUALS
!=	DIFFERENT
(LPAREN
)	RPAREN
[LBRACKETS
]	RBRACKETS
{	LBRACES
}	RBRACES
=	ATtribution
;	SEMICOLON
,	COMMA

Tabela 1: Tokens da linguagem C-.

Neste trabalho, será utilizado o aplicativo JFLAP (Rodger & Finley) para a visualização do autômato, e a biblioteca [automata_python](#) na linguagem python para a implementação do analisador léxico.

A máquina de Moore é um modelo de computação em que se transita de um estado para outro, e quando um estado é atingido, uma determinada sentença é escrita na saída, no nosso caso, a saída padrão do sistema operacional. Já a máquina de Mealy, é semelhante à de Moore, porém a sua saída é na transição, e não no estado.

1.1 Comparação entre máquina de Moore e Mealy

Para a nossa linguagem, devemos identificar que o lexema '<=' resulta no token LESS_EQUAL. Caso haja qualquer caractere diferente de '=' após o '<', então será impresso o token LESS, e também token do próximo caractere. Devido a esses casos, a máquina se torna mais complexa.

Considerando as letras do alfabeto (+,<=) e que não há o lexema '==' na linguagem, veja na Figura 1 a solução para esse problema usando a máquina de Moore, e na Figura 2 usando uma máquina de Mealy.

Podemos notar a diferença de complexidade entre as duas máquinas. Isso acontece pois na máquina de Moore, é necessário criar um estado para cada saída diferente, e consequentemente, cada estado deve ter transição para todas as letras do alfabeto. Considerando um alfabeto de

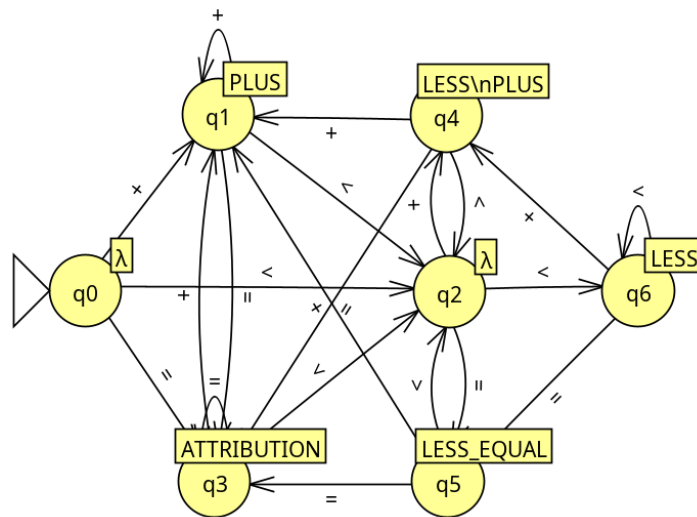


Figura 1: Exemplo da máquina de Moore

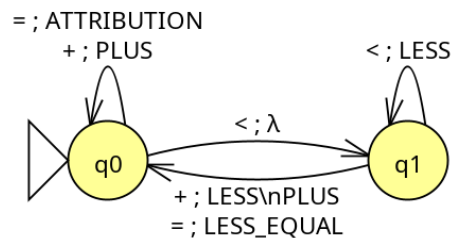


Figura 2: Exemplo da máquina de Mealy

mais de 30 letras, ficaria inviável compreender a máquina de Moore. Quando a saída está na transição, esse processo é simplificado. Logo, ambas as máquinas são equivalentes, porém, é muito mais complexo utilizar a máquina de Moore para o caso em questão. Por esse motivo, estarei utilizando a máquina de Mealy neste projeto, mas ainda assim, seria possível convertê-la para a máquina de Moore.

2 Autômato no JFLAP

Para facilitar a visualização, não criei explicitamente todas as transições. Todas as transições não listadas em um estado, ficam listadas na *label* do estado.

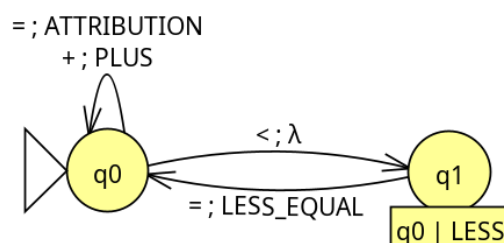


Figura 3: Exemplo de transição implícita em q1: cópia de q0 com o prefixo LESS e sobreposição da transição =

Na figura 3 podemos ver no estado q1 o seguinte texto em sua label: **q0 | LESS**. Posso explicar essa sintaxe em 3 passos:

Passo 1, q1 recebe uma cópia das transições de q0:

- $<; \lambda \rightarrow q1$

- =;ATtribution -> q0
- +;PLUS -> q0

Passo 2, Todas as saídas das transições de q1 recebem o prefixo LESS:

- <;LESS -> q1
- =;LESS Attribution -> q0
- +;LESS PLUS -> q0

Depois, a transição (=;LESS Attribution -> q0) é sobreposta pela (=;LESS_EQUAL -> q0). Assim, as transições do estado q1 são:

- <;LESS -> q1
- =;LESS_EQUAL -> q0
- +;LESS PLUS -> q0

Ou seja, as figuras 2, 1 e 3 são equivalentes. Apenas utilizarei dessa sintaxe para facilitar a visualização e compreensão do autômato.

2.1 Letras únicas

Para identificar os lexemas que eram apenas uma única letra, basta fazer uma transição para o mesmo estado, e que imprime o token. Veja a Figura 4.

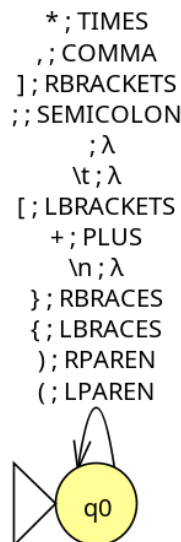


Figura 4: Transições dos lexemas de um único caractere

2.2 Comentários

Para ignorar comentários, é preciso fazer um estado em que o autômato fica preso ao reconhecer '/*', e só é liberado depois de reconhecer '*/'.

Na Figura 5 podemos ver que caso o estado q0 identifique o caractere '/' ele passa ao estado q1. O estado q1 possui uma definição implícita de todas as transições que não sejam '*', isso permite que sejam identificadas sequências como '/+' corretamente. Em q14 o autômato fica preso, e só

passa para o estado q15 se reconhecer '*', e posteriormente retorna ao estado inicial se reconhecer '/' novamente.

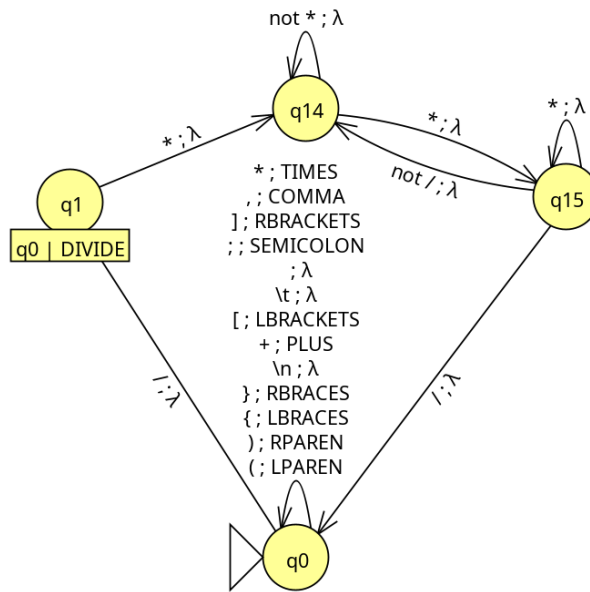


Figura 5: Autômato para ignorar comentários

2.3 Lexemas duplos

Lexemas duplos são aqueles que precisam de dois caracteres para serem reconhecidos, todos seguem o mesmo princípio já visto na Figura 3.

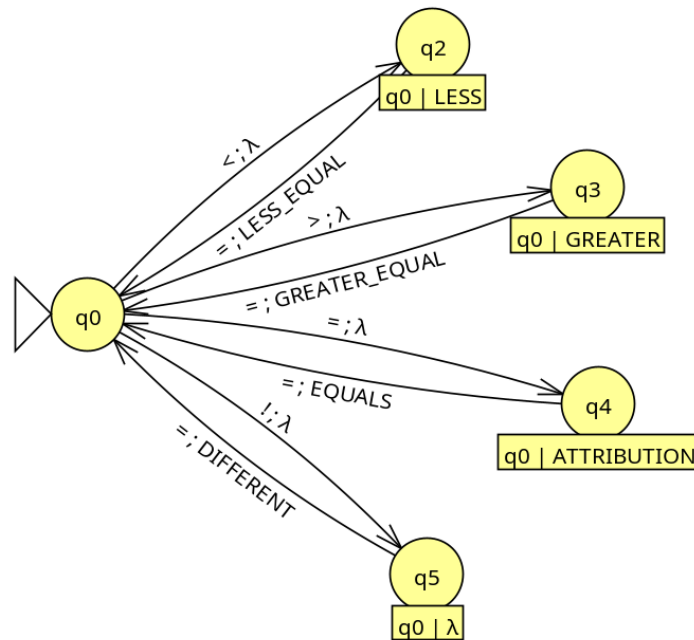


Figura 6: Autômato para identificar os lexemas duplos

2.4 Identificadores

Identificadores são palavras iniciadas por letra, seguidas de repetições de letras ou números. Um simples autômato para reconhecer identificadores está descrito na Figura 7. A abreviação *an* significa todos os alfa-numéricos.

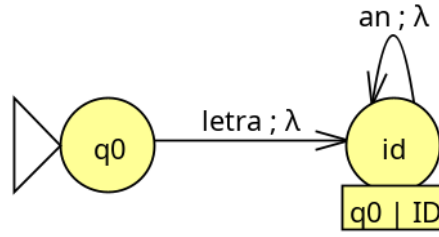


Figura 7: Autômato para reconhecer IDs

O autômato da Figura 7 não está incluso no autômato final dessa forma, há mais restrições. Essa é apenas uma forma genérica para reconhecer identificadores.

2.5 Palavras reservadas

A linguagem C- deve reconhecer as seguintes palavras reservadas: **if**, **int**, **else**, **float**, **return**, **void** e **while**. A Figura 8 exibe o reconhecimento da palavra 'if'.

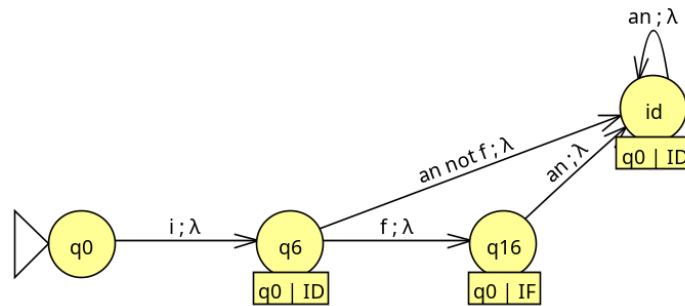


Figura 8: Autômato para reconhecer a palavra reservada if

Ainda na Figura 8, pode-se notar que no estado q6 há uma transição com qualquer alfa-numérico (an) que não seja 'f' para id. Essa transição existe para que cadeias como **ia** ou **i3a** sejam consideradas ID. Assim também, a transição de q16 para id identifica cadeias como **ifa**, **if1**, etc.

As transições implícitas nos estados permitem identificar o encerramento das palavras. Por exemplo, com a transição implícita no estado q6, é possível reconhecer a palavra **i+** como: ID PLUS \rightarrow q0.

As palavras **if** e **int** começam pela mesma letra, logo, deve haver um tratamento especial para esses dois casos. Veja na Figura 9.

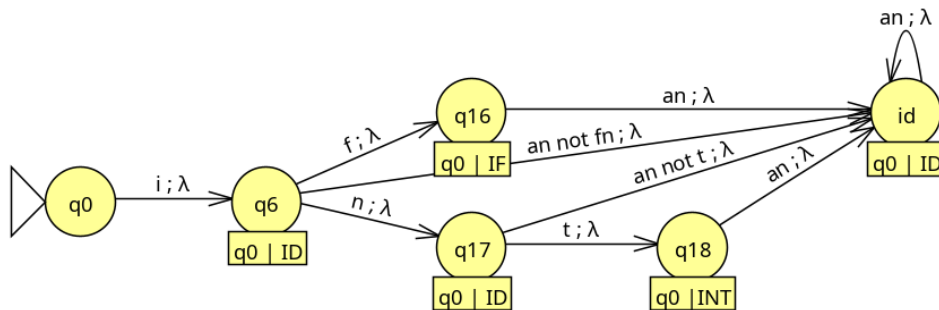


Figura 9: Autômato para reconhecer as palavras reservadas if e int

Todo estado intermediário de uma palavra reservada, deve ter uma transição para id que abrange todos os alfa-numéricos restantes daquele estado. q6 tem uma transição para id com todos os alfa-numéricos diferentes de 'f' e 'n', pois estes caracteres já possuem transições para os estados q16 e q17 respectivamente.

Todo estado final de uma palavra reservada deve ter uma transição para o estado `id` com todos os alfa-numéricos. Como exemplo, temos os estados `q16` e `q18`.

Esse processo de construção do autômato se repete para todas as palavras reservadas restantes. Para melhor visualização, irei omitir todas as transições para `id`. As Figuras 10 e 11 destacam essas diferenças para a palavra reservada `else`.

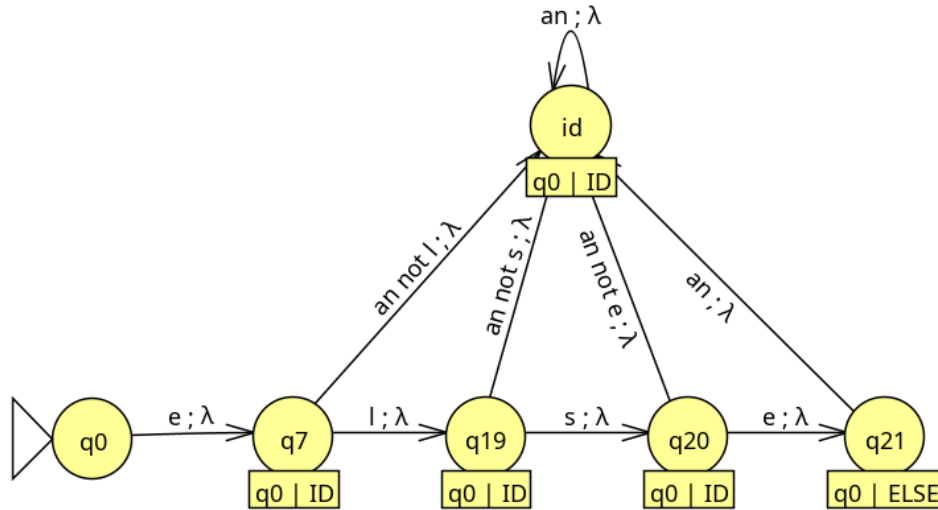


Figura 10: Autômato para reconhecer a palavra reservada `else` com todas as transições para `id`

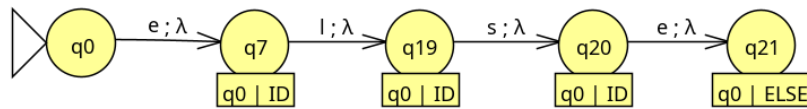


Figura 11: Autômato para reconhecer a palavra reservada `else` sem as transições para `id`

2.6 Números

Os números podem ser negativos, por isso, o símbolo '-' pode estar antes de uma sequência numérica e deve ser reconhecido corretamente. Caso haja qualquer caractere diferente de número após o sinal de menos, será reconhecido o token `MINUS`.

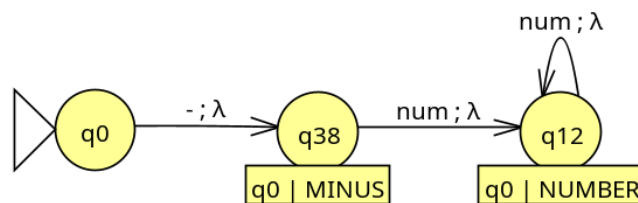


Figura 12: Autômato para reconhecer números

2.7 Resultado do autômato

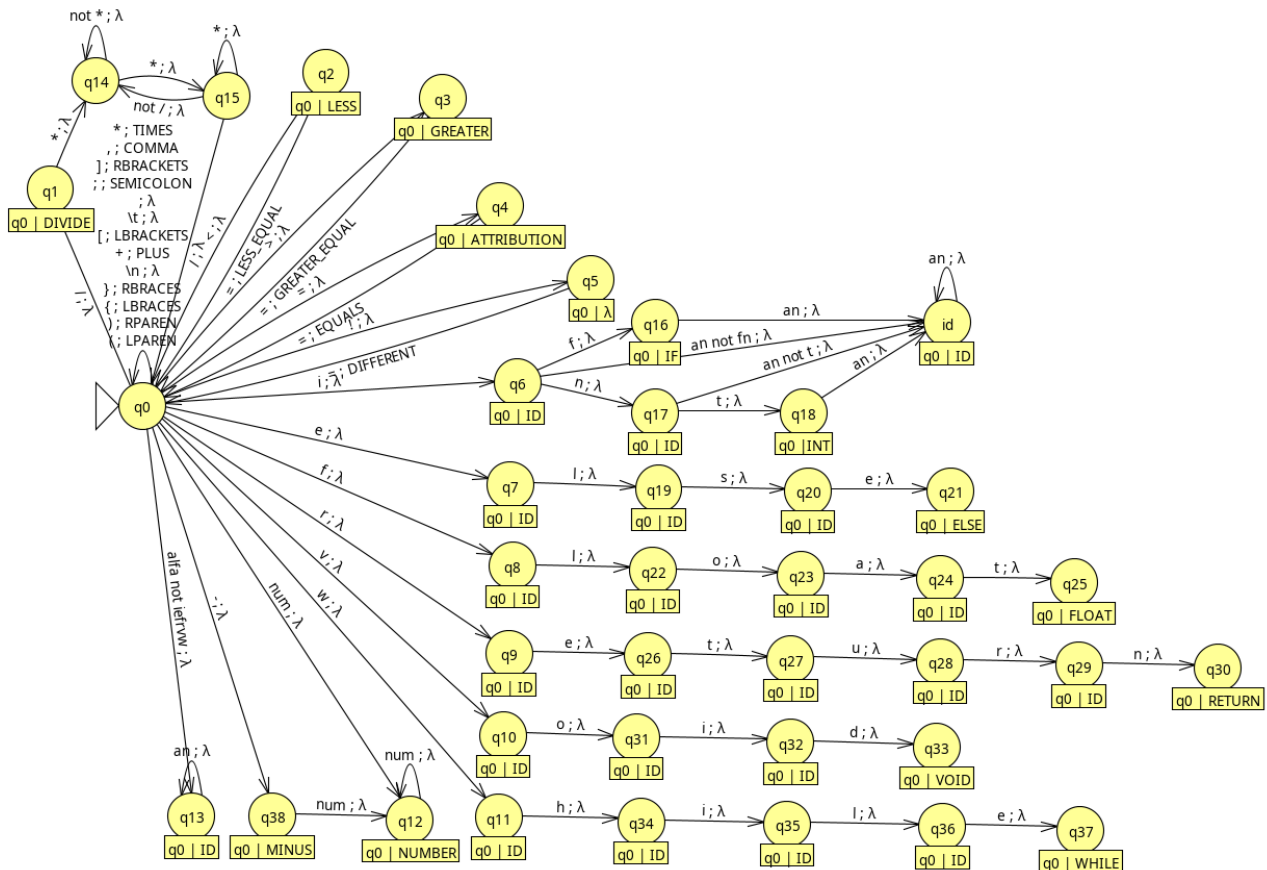


Figura 13: Resultado final do autômato para o analisador léxico

Na Figura 13, o estado q13 não é necessário, pois poderia ser substituído pelo estado id. No entanto, criei outro para melhorar a visualização.

3 Implementação em Python

Para a implementação do autômato na linguagem de programação Python, foi escolhida a biblioteca [automata_python](#). Para criar uma máquina de Mealy com essa biblioteca é necessário informar os seguintes dados:

2. Interface para criar uma máquina Mealy

```
1 mealy = Mealy(
2     states,
3     input_alphabet,
4     output_alphabet,
5     transitions,
6     initial_state
7 )
```

Como as transições do estado inicial seriam copiadas para todos os outros estados, declarei uma variável global que continha essas transições.

3. Transições do estado inicial (q0)

```
1 initial = {
2     # Transições para o estado q0
3     '\n': ('q0', ''),
4     '\t': ('q0', ''),
```



```

5      '\n': ('q0', ''),
6      '+': ('q0', 'PLUS\n'),
7      '*': ('q0', 'TIMES\n'),
8      '(': ('q0', 'LPAREN\n'),
9      ')': ('q0', 'RPAREN\n'),
10     '[': ('q0', 'LBRACKETS\n'),
11     ']': ('q0', 'RBRACKETS\n'),
12     '{': ('q0', 'LBRACES\n'),
13     '}': ('q0', 'RBRACES\n'),
14     ';': ('q0', 'SEMICOLON\n'),
15     ',': ('q0', 'COMMA\n'),
16
17     # Transições para outros estados
18     '-': ('q38', ''),
19     '/': ('q1', ''),
20     '<': ('q2', ''),
21     '>': ('q3', ''),
22     '=': ('q4', ''),
23     '!': ('q5', ''),
24     'i': ('q6', ''),
25     'e': ('q7', ''),
26     'f': ('q8', ''),
27     'r': ('q9', ''),
28     'v': ('q10', ''),
29     'w': ('q11', ''),
30
31     # Números
32     '0': ('q12', ''), '1': ('q12', ''),
33     '2': ('q12', ''), '3': ('q12', ''),
34     '4': ('q12', ''), '5': ('q12', ''),
35     '6': ('q12', ''), '7': ('q12', ''),
36     '8': ('q12', ''), '9': ('q12', ''),
37
38     # Restante das letras
39     'a': ('id', ''), 'b': ('id', ''), 'c': ('id', ''), 'd': ('id', ''), 'g': ('id', ''),
40     'h': ('id', ''), 'j': ('id', ''), 'k': ('id', ''), 'l': ('id', ''), 'm': ('id', ''),
41     'n': ('id', ''), 'o': ('id', ''), 'p': ('id', ''), 'q': ('id', ''), 's': ('id', ''),
42     't': ('id', ''), 'u': ('id', ''), 'x': ('id', ''), 'y': ('id', ''), 'z': ('id', ''),
43 }

```

Criei também uma classe para facilitar a manipulação das transições. Essa classe inicia com as transições contidas na variável `initial`, e possui métodos para adicionar prefixo, sobrescrever transições, e `forEachs` que servem como funções auxiliares.

4. Classe Transition

```

1 class Transition:
2     def __init__(self):
3         self.transition = initial.copy()
4
5     def prefix(self, pre):
6         for i, j in self.transition.items():
7             self.transition[i] = (j[0], pre + j[1])
8         return self
9
10    def overwrite(self, values):
11        self.transition.update(values)
12        return self
13
14    def forEach(self, func):
15        for i, j in self.transition.items():

```

```

16         self.transition[i] = func(i, j)
17     return self
18
19     def forEachAlfaNum(self, func):
20         for i, j in self.transition.items():
21             if i in alfanum:
22                 self.transition[i] = func(i, j)
23         return self
24
25     def forEachAlfa(self, func):
26         for i, j in self.transition.items():
27             if i in alfa:
28                 self.transition[i] = func(i, j)
29         return self
30
31     def forEachNum(self, func):
32         for i, j in self.transition.items():
33             if i in num:
34                 self.transition[i] = func(i, j)
35         return self
36
37     def __getitem__(self, key):
38         return self.transition[key]

```

Para auxiliar a criação dos estados das palavras reservadas, defini duas funções que realizam os processos vistos na seção 2.5.

5. Funções auxiliares para criar estados das palavras reservadas

```

1 def LastReservedWordState(word):
2     transitions = Transition()
3     transitions.prefix(word + '\n')
4     transitions.forEachAlfaNum(lambda _, __: ('id', ''))
5     return transitions
6
7 def IntermediateReservedWordState(input, nextState):
8     transitions = Transition()
9     transitions.prefix('ID\n')
10    transitions.forEachAlfaNum(lambda _, __: ('id', ''))
11    for i in range(len(input)):
12        transitions.override({ input[i]: (nextState[i], '') })
13    return transitions

```

Por fim, criei a transição de cada estado da máquina.

6. Transições do analisador léxico

```

1 transitions = {
2     'id': Transition()
3         .prefix("ID\n")
4         .forEachAlfaNum(lambda _, __: ('id', '')),
5
6     'q0': Transition(),
7
8     'q1': Transition()
9         .prefix("DIVIDE\n")
10        .override({ '*': ('q14', '') }),
11
12    'q14': Transition()
13        .forEach(lambda _, __: ('q14', ''))
14        .override({ '*': ('q15', '') }),
15

```

```

16      'q15': Transition()
17          .forEach(lambda _, __: ('q14', ''))
18          .overwrite({ '/': ('q0', '') }),
19
20      'q2': Transition()
21          .prefix("LESS\n")
22          .overwrite({ '=': ('q0', 'LESS_EQUAL\n') }),
23
24      'q3': Transition()
25          .prefix("GREATER\n")
26          .overwrite({ '=': ('q0', 'GREATER_EQUAL\n') }),
27
28      'q4': Transition()
29          .prefix("ATtribution\n")
30          .overwrite({ '=': ('q0', 'EQUALS\n') }),
31
32      'q5': Transition()
33          .overwrite({ '=': ('q0', 'DIFFERENT\n') }),
34
35
36      # IF e INT
37      'q6': IntermediateReservedWordState(['f', 'n'], ['q16', 'q17']),
38      'q16': LastReservedWordState("IF"),
39      'q17': IntermediateReservedWordState(['t'], ['q18']),
40      'q18': LastReservedWordState("INT"),
41
42      # ELSE
43      'q7': IntermediateReservedWordState(['l'], ['q19']),
44      'q19': IntermediateReservedWordState(['s'], ['q20']),
45      'q20': IntermediateReservedWordState(['e'], ['q21']),
46      'q21': LastReservedWordState("ELSE"),
47
48      # FLOAT
49      'q8': IntermediateReservedWordState(['l'], ['q22']),
50      'q22': IntermediateReservedWordState(['o'], ['q23']),
51      'q23': IntermediateReservedWordState(['a'], ['q24']),
52      'q24': IntermediateReservedWordState(['t'], ['q25']),
53      'q25': LastReservedWordState("FLOAT"),
54
55      # RETURN
56      'q9': IntermediateReservedWordState(['e'], ['q26']),
57      'q26': IntermediateReservedWordState(['t'], ['q27']),
58      'q27': IntermediateReservedWordState(['u'], ['q28']),
59      'q28': IntermediateReservedWordState(['r'], ['q29']),
60      'q29': IntermediateReservedWordState(['n'], ['q30']),
61      'q30': LastReservedWordState("RETURN"),
62
63      # VOID
64      'q10': IntermediateReservedWordState(['o'], ['q31']),
65      'q31': IntermediateReservedWordState(['i'], ['q32']),
66      'q32': IntermediateReservedWordState(['d'], ['q33']),
67      'q33': LastReservedWordState("VOID"),
68
69      # WHILE
70      'q11': IntermediateReservedWordState(['h'], ['q34']),
71      'q34': IntermediateReservedWordState(['i'], ['q35']),
72      'q35': IntermediateReservedWordState(['l'], ['q36']),
73      'q36': IntermediateReservedWordState(['e'], ['q37']),
74      'q37': LastReservedWordState("WHILE"),
75

```

```

76     # Números
77     'q12': Transition()
78         .prefix("NUMBER\n")
79         .forEachNum(lambda _, __: ('q12', '')),
80
81     'q38': Transition()
82         .prefix("MINUS\n")
83         .forEachNum(lambda _, __: ('q12', '')),
84 }

```

4 Resultados

Para executar o analisador léxico, podemos informar o arquivo a ser analisado como parâmetro, e então a saída será os tokens daquele código. É possível adicionar a *flag* -k para mostrar apenas as chaves dos tokens e erros, desse jeito, não serão mostradas informações adicionais.

7. Código para teste do analisador léxico

```

1  int a;
2  float b;
3
4  int main(void){
5      int x;
6      float y;
7
8      return(0);
9  }

```

8. Resultado do analisador léxico

```

1  INT
2  ID
3  SEMICOLON
4  FLOAT
5  ID
6  SEMICOLON
7  INT
8  ID
9  LPAREN
10 VOID
11 RPAREN
12 LBRACES
13 INT
14 ID
15 SEMICOLON
16 FLOAT
17 ID
18 SEMICOLON
19 RETURN
20 LPAREN
21 NUMBER
22 RPAREN
23 SEMICOLON
24 RBRACES

```

Referências

Rodger, Susan H. & Thomas W. Finley. 2006. *Jflap: An interactive formal languages and automata package*. Sudbury, MA: Jones & Bartlett Learning.