

Projeto de Implementação de um Compilador para a Linguagem T++ Análise Léxica (Trabalho – 1ª parte)

Rafael Menezes Barboza¹

¹Universidade Tecnológica Federal do Paraná (UTFPR)
Via Rosalina Maria dos Santos, 1233
CEP: 87301-899, Campo Mourão – PR – Brasil

ra29fa@gmail.com

***Resumo.** O seguinte trabalho faz parte da disciplina de Compiladores do curso de ciência da computação, é a primeira parte de total de quatro etapas que ao final resultarão no estudo e implementação de compilador para a linguagem T++. Esta primeira abordagem pretende explorar o funcionamento do Analisador Léxico, um dos componentes principais para separar todo código fonte em pequenos "tokens" devidamente classificados e reconhecidos pela linguagem.*

1. Introdução

Um compilador de forma simplificada nada mais é do que um tradutor de um código fonte em uma linguagem de programação de alto nível para uma linguagem de programação de baixo nível. As linguagens de alto nível são muito utilizadas no desenvolvimento de aplicações, porém a maioria dos computadores processam instruções mais simples e elementares e o compilador tem exatamente este papel de transformar o código para algo mais próximo da linguagem da máquina.

O processo de compilação é dividido em várias etapas bem definidas, cada etapa gera uma saída de dados correspondente a entrada da próxima etapa a ser executada, se assemelhando aos famosos Pipelines. Cada etapa deve ser encarregada de tratar possíveis erros ou notificá-los antes de passar para as próximas fases.

2. Linguagem de Estudo T++

A linguagem de programação escolhida para o trabalho foi a T++ linguagem essa inventada para fins acadêmicos e não é utilizada em meios comerciais. A linguagem T++ se assemelha bastante com as principais linguagem de programação compiladas do mercado por conter uma lógica de programação muito bem estruturada e familiar porém com algumas certas limitações.

Se trata de uma linguagem considerada quase fortemente tipada pois toda e qualquer variável criada deve ser atribuída a algum tipo de dado reconhecido pela linguagem e este deve ser persistente durante toda a execução do código. Por se tratar de uma linguagem quase fortemente tipada nem todos os erros são especificados mas sempre deve ocorrer avisos. O código pode ser estruturado por meio de funções que se relacionam podendo ou não retornar algum valor, no caso de não retornar nada assumimos que a função em questão é do tipo void.

A linguagem conta com dois tipos de dados básicos sendo eles “inteiro” e “flutuante” suporta também estruturas de dados unidimensional (array) e bidimensional (matriz). Com tudo isso é possível realizar uma série de operações aritméticas e lógicas, assim com em outras linguagem como por exemplo C/C++.

A Tabela 1 apresenta os principais símbolos e palavras reservadas suportadas pela linguagem T++, nesta etapa de análise léxica iremos nos referir a cada item da tabela como Tokens. Porém ainda podem ser definidos como tokens, números com um ou mais dígitos que podem ser inteiro, flutuante ou notação científica, identificadores que começam com uma letra e precedem com N letras e números e por fim comentários cercados de chaves {...} podendo conter uma ou mais linhas. Todos estes tokens posteriormente serão devidamente relacionados gerando assim a toda a logica do código fonte.

palavras reservadas	símbolos
se	+ soma
então	- subtração
senão	* multiplicação
fim	/ divisão
repita	= igualdade
flutuante	, vírgula
retorna	:= atribuição
até	< menor
leia	> maior
escreva	<= menor-igual
inteiro	>= maior-igual
	(abre-par
) fecha-par
	: dois-pontos
	[abre-col
] fecha-col
	&& e-logico
	ou-logico
	! negação

Tabela 1. Tokens da linguagem T++

O código fonte da figura 1 em T++ calcula o fatorial de um número passado pelo input do teclado, é possível visualizar o uso de algumas palavras reservadas e símbolos em ação no código. É notório que algumas palavras reservadas como “se, repita, senão ...” estão sempre aninhadas da palavra reservada “fim” que denota o fechamento do bloco, função está realizada pelo abre e fecha chaves na linguagem C/C++.

3. Analise Léxica

Análise léxica é a primeira fase do processo de compilação de um código fonte, nessa etapa o componente léxico realiza a varredura pelo código fonte carácter por carácter

```

1 ▼ inteiro principal()
2     inteiro: digitado
3     inteiro: i
4     i := 1
5 ▼     repita
6         flutuante: f
7         inteiro: int
8         flutuante: resultado
9         f := i/2.0
10        int := i/2
11        resultado := f - int
12
13        se resultado > 0
14            escreva (i)
15        fim
16        i := i+1
17    até i <= digitado
18 fim

```

Figura 1. Fatorial em T++

coletando as palavras e separando em conjuntos de marcas (tokens) que podem ser manipulados de forma mais fácil por um parser leitor de saída.

São reconhecidas as palavras reservadas, identificadores, constantes, símbolos entre outras palavras da linguagem em questão. Também é nessa etapa que é tratado espaços em branco, contagem das linhas e eliminação de comentários simplificando a saída para os próximos componentes do processo.

Cada conjunto de caracteres ou palavra é reconhecida por uma regra, no seguinte trabalho as regras foram elaboradas por meio de expressões regulares bem definidas que realizam a filtragem e separação de cada token com o seu devido valor.

3.1. Expressão Regular e Linguagens Formais

Para esse tipo de tarefa é de extrema importância a percepção da necessidade da utilização de linguagens formais para validar os tokens corretamente. A linguagem natural apesar de muito usual em diversos aspectos ela não é recomendada processos validativos pois pode trazer ambiguidade e uma complexidade ainda maior ao problema. Por isso a escolha de uma linguagem formal que seja simples, concisa, clara e sem ambiguidade é tão importante.

Por meio das expressões regulares é possível determinar uma cadeia de caracteres específicos para ser reconhecido, um exemplo é a expressão regular $r'senão'$ para expressão só será reconhecido palavras com a exata ordem de caracteres determinado. Para reconhecer palavras de forma menos engessada pode ser usado um range de caracteres como por exemplo $r'[A-Z]'$ que reconhece qualquer caractere do alfabeto (maiúsculo) que esteja dentro do range definido. Também é possível definir expressões que aceitam repetição de caracteres, números ou símbolos, um exemplo é a expressão regular $r'[A-Za-z]+[0-9_]*'$ que pode aceitar palavras que comece com pelo menos um carácter maiúsculo ou minúsculo e nenhum ou vários caracteres numéricos de 0 a 9 podendo haver o símbolo underscore ou não, aceitando palavras como $[rafa, RAFAEL, R, a56_, w_, ...]$.

3.2. Autômatos Finitos

Através de autômatos finitos podemos representar o funcionamento de expressões regulares visualizando os estados de aceitação e transições. Nesta sessão serão apresentados os autômatos finitos referentes às expressões regulares implementadas no analisador léxico, é possível veras com mais detalhes na seção 4 deste documento.

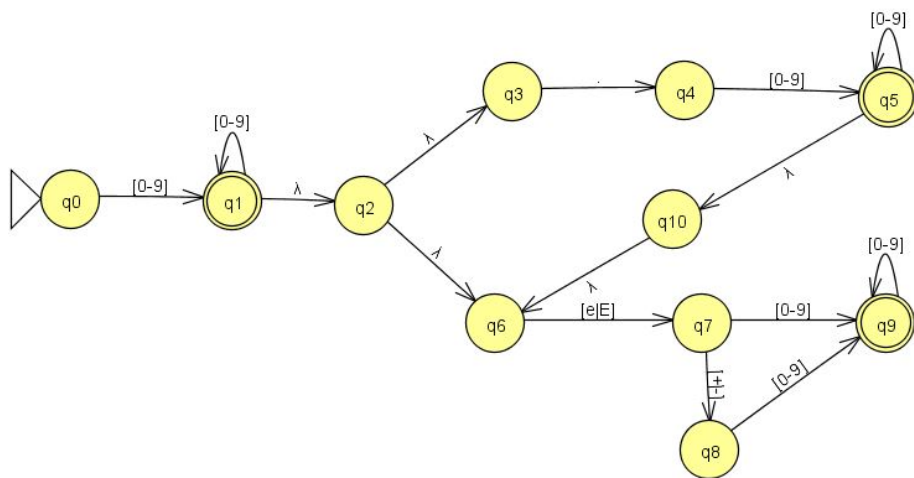


Figura 2. Reconhece números (Inteiros, Ponto Flutuante e Notação Científica)

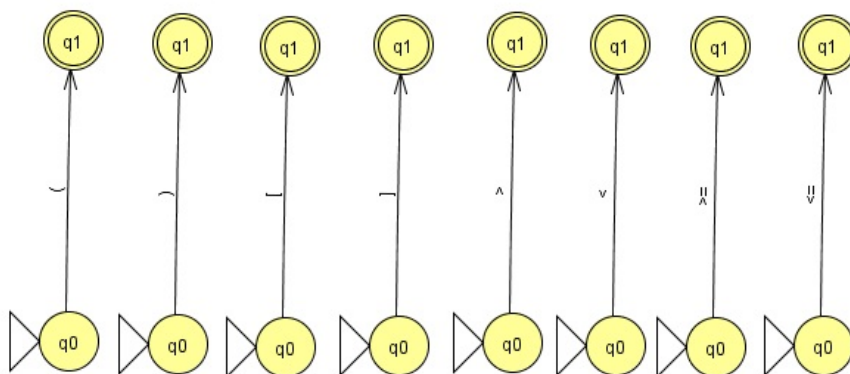


Figura 3. Reconhece operadores aritméticos

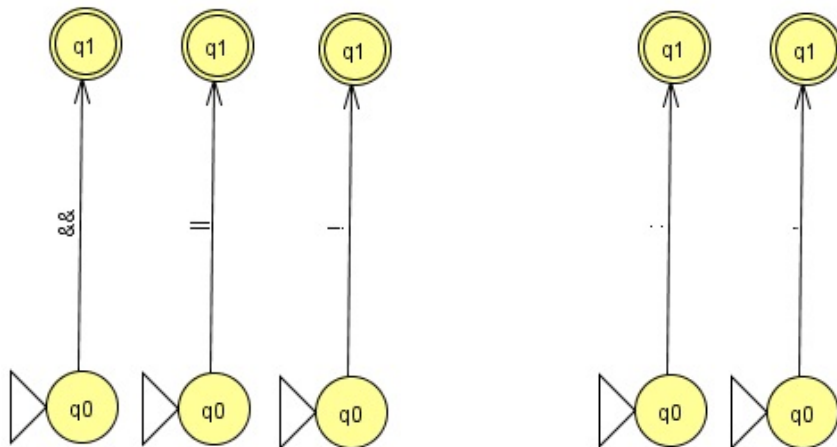


Figura 4. Reconhece símbolos e caracteres especiais da linguagem

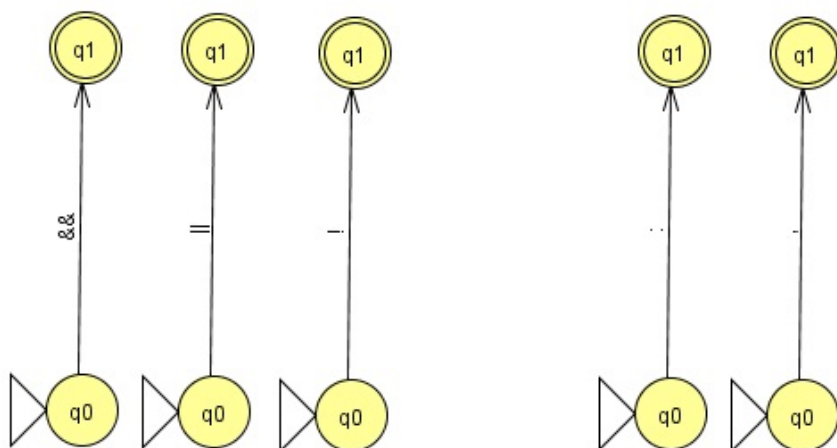


Figura 5. Reconhece Operadores lógicos e pontuações

4. Detalhes da implementação

4.1. Ferramentas auxiliares

Para a implementação e realização deste trabalho algumas ferramentas auxiliares foram utilizadas. As sessões seguintes descrevem um pequeno resumo das ferramentas.

4.1.1. Python

A linguagem de programação usada para a implementação deste projeto foi a linguagem python pois oferece uma gama maior de ferramentas de fácil acesso facilitando a implementação.

4.1.2. PLY

A ferramenta PLY se trata de uma implementação do analisador léxico (lex) e (yacc), com esta ferramenta é possível junto com a linguagem de programação Python definir regras de precedências, gramáticas e etc.

4.2. Implementação

Para a implementação do analisador léxico foi necessário definir tokens e palavras reservadas da linguagem T++. Estes foram definidos em duas listas presentes no código nas linhas [3 a 49], assim que definidos foi necessário definir para cada token uma expressão regular para que fosse reconhecido, podemos verificar isto nas linhas [51 a 87]. Para as palavras reservadas do dicionário RESERVED não foi necessário definir as expressões regulares pois o próprio dicionário de palavras é capaz de mapear isto.

Como o analisador léxico é responsável por notificar possíveis erros e ignorar comentários no código, foi implementado duas funções presentes nas linhas [89 a 98] que realizam esta função respectivamente.

As linhas [100 a 108] apenas realizam a chamada da função lex(), a abertura do arquivo a ser analisado e o processo de impressão dos tokens.

```
1 import ply.lex as lex
2 import sys
3
4 RESERVED = {
5     'se' : 'SE',
6     'senão' : 'SENAO',
7     'então' : 'ENTAO',
8     'repita' : 'REPITA',
9     'escreva' : 'ESCREVA',
10    'leia' : 'LEIA',
11    'até' : 'ATE',
12    'fim' : 'FIM',
13    'inteiro' : 'INTEIRO',
14    'flutuante' : 'FLUTUANTE',
15    'retorna' : 'RETORNA',
16 }
17 tokens = ['ID',
18          'DOIS_PONTOS',
19          'ATRIBUICAO',
20          'MENOR',
21          'MAIOR',
22          'MENOR_IGUAL',
23          'MAIOR_IGUAL',
24          'MULT',
25          'VIRGULA',
26          'MAIS',
27          'MENOS',
28          'DIVIDE',
29          'ABRE_PAREN',
30          'FECHA_PAREN',
31          'ABRE_COUCH',
32          'FECHA_COUCH',
33          'SE',
34          'SENAO',
35          'ENTAO',
36          'REPITA',
37          'ESCREVA',
38          'LEIA',
39          'ATE',
40          'FIM',
41          'INTEIRO',
```

```

42         'FLUTUANTE',
43         'COMPARACAO',
44         'NUMERO',
45         'RETORNA',
46         'E_LOGICO',
47         'OU_LOGICO',
48         'NEGACAO',
49         'QUEBRA_LINHA',
50     ]
51
52     t_ignore = " \t\n"
53
54     ##operadores
55     t MAIS = r'\+'
56     t MENOS = r'\-'
57     t MULT = r'\*'
58     t DIVIDE = r'\/'
59     t COMPARACAO = r'='
60
61     ##símbolos
62     t ABRE_PAREN = r'\('
63     t FECHA_PAREN = r'\)'
64     t ABRE_COUCH = r'\['
65     t FECHA_COUCH = r'\]'
66     t MENOR = r'<'
67     t MAIOR = r'>'
68     t ATRIBUICAO = r':='
69     t MENOR_IGUAL = r'<='
70     t MAIOR_IGUAL = r'>='
71
72     ##operadores lógicos
73     t E_LOGICO = r'&&'
74     t OU_LOGICO = r'\|\|'
75     t NEGACAO = r'!'
76
77     ##pontos
78     t DOIS_PONTOS = r':'
79     t VIRGULA = r','
80
81     ##números
82     t_NUMERO = r'([0-9]+)(\.[0-9]+)?([e|E][+|-]?[0-9]+)?'
83
84     ##identificadores
85     def t_ID(t):
86         r'[A-Za-z_][\w]*'
87         t.type = RESERVED.get(t.value, 'ID')
88         return t
89
90     ##comentários
91     def t_comment(t):
92         r'\{[^\}]*[^\}*\}'
93         pass
94
95     ##tratamento de erros
96     def t_error(t):
97         print("Erro", t.value[0])
98         raise SyntaxError("ERRO", t.value[0])
99         t.lexer.skip(1)
100
101     lex.lex()
102
103     file = open(sys.argv[1], "r", encoding="utf-8")
104     lex.input(file.read())
105     while True:
106         tok = lex.token()
107         if not tok: break
108         print( tok.type, ': ', tok.value )

```

5. Exemplos de saída do Analisador Léxico

Nesta seção serão apresentados alguns exemplos de saída quando o analisador léxico é submetido a códigos testes. A função responsável por imprimir a saída foi manipulada para assumir uma estrutura mais simples e facilitar o entendimento seguindo sempre a lógica de “TOKEN : valor”.

```
1 ▾ inteiro principal()
2     inteiro: digitado
3     inteiro: i
4     i := 1
5 ▾ repita
6     flutuante: f
7     inteiro: int
8     flutuante: resultado
9     f := i/2.0
10    int := i/2
11    resultado := f - int
12
13    se resultado > 0
14    | escreva (i)
15    fim
16    i := i+1
17 até i <= digitado
18 fim
```

Figura 6. Verifica se é numero primo em T++ (ENTRADA)

1 INTEIRO : inteiro	12 ATRIBUICAO ::=	23 ID : resultado	34 ID : resultado	45 ID : i	56 ID : digitado
2 ID : principal	13 NUMERO : 1	24 ID : f	35 ATRIBUICAO ::=	46 FECHA_PAREN :)	57 FIM : fim
3 ABRE_PAREN : (14 REPITA : repita	25 ATRIBUICAO ::=	36 ID : f	47 FIM : fim	
4 FECHA_PAREN :)	15 FLUTUANTE : flutuante	26 ID : i	37 MENOS : -	48 ID : i	
5 INTEIRO : inteiro	16 DOIS_PONTOS ::	27 DIVIDE : /	38 ID : int	49 ATRIBUICAO ::=	
6 DOIS_PONTOS ::	17 ID : f	28 NUMERO : 2.0	39 SE : se	50 ID : i	
7 ID : digitado	18 INTEIRO : inteiro	29 ID : int	40 ID : resultado	51 MAIS : +	
8 INTEIRO : inteiro	19 DOIS_PONTOS ::	30 ATRIBUICAO ::=	41 MAIOR : >	52 NUMERO : 1	
9 DOIS_PONTOS ::	20 ID : int	31 ID : i	42 NUMERO : 0	53 ATE : até	
10 ID : i	21 FLUTUANTE : flutuante	32 DIVIDE : /	43 ESCREVA : escreva	54 ID : i	
11 ID : i	22 DOIS_PONTOS ::	33 NUMERO : 2	44 ABRE_PAREN : (55 MENOR_IGUAL : <=	

Figura 7. Verifica se é numero primo em T++ (SAIDA)

```
1 inteiro: vet[10]
2 inteiro: tam
3
4 tam := 10
5
6 { preenche o vetor no pior caso }
7 preencheVetor()
8     inteiro: i
9     inteiro: j
10    i := 0
11    j := tam
12    repita
13    | vet[i] = j
14    | i := i + 1
15    | j := j - 1
16 até i < tam
17 fim

18 { implementação do bubble sort }
19 ▾ bubble_sort()
20     inteiro: i
21     i := 0
22 ▾ repita
23     inteiro: j
24     j := 0
25 ▾ repita
26 ▾ se vet[i] > v[j] então
27     | inteiro: temp
28     | temp := vet[i]
29     | vet[i] := vet[j]
30     | vet[j] := temp
31     fim
32     j := j + 1
33 até j < i
34     i := i + 1
35 até i < tam
36 fim

37 { programa principal }
38 ▾ inteiro principal()
39     preencheVetor()
40     bubble_sort()
41     retorna(0)
42 fim
```

Figura 8. Algoritmo de ordenação BubbleSort em T++ (ENTRADA)

1 INTEIRO : inteiro	24 NUMERO : 0	47 MENOR : <	70 ID : i	93 ABRE_COUCH : [116 NUMERO : 1
2 DOIS_PONTOS ::	25 ID : j	48 ID : tam	71 FECHA_COUCH :]	94 ID : j	117 ATE : até
3 ID : vet	26 ATRIBUICAO ::=	49 FIM : fim	72 MAIOR : >	95 FECHA_COUCH :]	118 ID : i
4 ABRE_COUCH : [27 ID : tam	50 ID : bubble_sort	73 ID : v	96 ID : vet	119 MENOR : <
5 NUMERO : 10	28 REPITA : repita	51 ABRE_PAREN : (74 ABRE_COUCH : [97 ABRE_COUCH : [120 ID : tam
6 FECHA_COUCH :]	29 ID : vet	52 FECHA_PAREN :)	75 ID : j	98 ID : j	121 FIM : fim
7 INTEIRO : inteiro	30 ABRE_COUCH : [53 INTEIRO : inteiro	76 FECHA_COUCH :]	99 FECHA_COUCH :]	122 INTEIRO : inteiro
8 DOIS_PONTOS ::	31 ID : i	54 DOIS_PONTOS ::	77 ENTAO : então	100 ATRIBUICAO ::=	123 ID : principal
9 ID : tam	32 FECHA_COUCH :]	55 ID : i	78 INTEIRO : inteiro	101 ID : temp	124 ABRE_PAREN : (
10 ID : tam	33 COMPARACAO : =	56 ID : i	79 DOIS_PONTOS ::	102 FIM : fim	125 FECHA_PAREN :)
11 ATRIBUICAO ::=	34 ID : j	57 ATRIBUICAO ::=	80 ID : temp	103 ID : j	126 ID : preencheVetor
12 NUMERO : 10	35 ID : i	58 NUMERO : 0	81 ID : temp	104 ATRIBUICAO ::=	127 ABRE_PAREN : (
13 ID : preencheVetor	36 ATRIBUICAO ::=	59 REPITA : repita	82 ATRIBUICAO ::=	105 ID : j	128 FECHA_PAREN :)
14 ABRE_PAREN : (37 ID : i	60 INTEIRO : inteiro	83 ID : vet	106 MAIS : +	129 ID : bubble_sort
15 FECHA_PAREN :)	38 MAIS : +	61 DOIS_PONTOS ::	84 ABRE_COUCH : [107 NUMERO : 1	130 ABRE_PAREN : (
16 INTEIRO : inteiro	39 NUMERO : 1	62 ID : j	85 ID : i	108 ATE : até	131 FECHA_PAREN :)
17 DOIS_PONTOS ::	40 ID : j	63 ID : j	86 FECHA_COUCH :]	109 ID : j	132 RETORNA : retorna
18 ID : i	41 ATRIBUICAO ::=	64 ATRIBUICAO ::=	87 ID : vet	110 MENOR : <	133 ABRE_PAREN : (
19 INTEIRO : inteiro	42 ID : j	65 NUMERO : 0	88 ABRE_COUCH : [111 ID : i	134 NUMERO : 0
20 DOIS_PONTOS ::	43 MENOS : -	66 REPITA : repita	89 ID : i	112 ID : i	135 FECHA_PAREN :)
21 ID : j	44 NUMERO : 1	67 SE : se	90 FECHA_COUCH :]	113 ATRIBUICAO ::=	136 FIM : fim
22 ID : i	45 ATE : até	68 ID : vet	91 ATRIBUICAO ::=	114 ID : i	
23 ATRIBUICAO ::=	46 ID : i	69 ABRE_COUCH : [92 ID : vet	115 MAIS : +	

Figura 9. Algoritmo de ordenação BubbleSort em T++ (SAIDA)

Referências

- Aho, A., Sethi, R., and Lam, S. (2008). *Compiladores. Princípios, Técnicas e Ferramentas*. LONGMAN DO BRASIL.
- Beazley, D. M. (2018 (Acessado em 01 de Setembro, 2018)). *PLY (Python Lex-Yacc)*. <http://www.dabeaz.com/ply/ply.html>.
- Wikipédia (2018). *Compilador - wikipédia, a enciclopédia livre*. (Acessado em 01 de Setembro, 2018).