

Compilador da Linguagem T++: Relatório de Implementação

Sávio O. Camacam¹

¹Departamento Acadêmico de Computação –
Universidade Tecnológica Federal do Paraná (UTFPR)
Caixa Postal 271 – 87.301-899 – Campo Mourão – PR – Brazil

saviocamacam@alunos.utfpr.edu.br

Abstract. *This work describes the implementation phases of the compilation tool created during the school year in Compilers, aiming to fix the main concepts of the actual compiler construction steps. The complete work comprises 4 phases of implementation: lexical analysis, syntactic analysis, semantic analysis and, finally, code generation. Since the implementation process is incremental, the report is also done in incremental format, with the four main topics related to these steps and their respective internal stages and characteristics.*

Resumo. *Este trabalho descreve as fases de implementação da ferramenta de compilação criada durante o período letivo na disciplina de Compiladores, tendo o objetivo de fixar os principais conceitos das etapas de construção de compilador real. O trabalho completo compreende 4 fases de implementação: análise léxica, análise sintática, análise semântica e por fim, geração de código. Sendo incremental o processo de implementação, o relatório se faz também no formato incremental, com os quatro grandes tópicos relativos a estas etapas e suas respectivas etapas e características internas.*

1. Introdução

Este relatório descreve o processo de construção da ferramenta de compilação para a Linguagem T++, que possui uma extensão de funcionalidades e características, e adaptação para o português, da linguagem-exemplo TINNY, do livro-texto Compiler Construction Principles, de Kenneth Loudon [Louden].

Um programa nessa linguagem tem uma estrutura relativamente simples, mas com características de línguas como C , possuindo sequências de *statements*, possuindo declarações e chamadas de funções, tipos de dados dimensionados e operadores lógicos.

Este trabalho está organizado da forma como segue: na próxima seção é feita uma especificação mais detalhada da linguagem, sobre seus itens e dados e declarações; em seguida são descritas as palavras reservadas da linguagem que são aceitas no Português; adiante são descritos os símbolos e operadores de manipulação da linguagem; logo depois é dada uma especificação formal dos autômatos da linguagem com descrição as expressões regulares usadas para captura dos tokens, bem como os autômatos que os representam; depois é explicado sobre o processo de execução do programa, bem como é exibida uma figura de exemplo de uma saída desse processador e por fim, uma conclusão sobre esta etapa do trabalho.

2. Análise Léxica

O processo de escaneamento, ou análise léxica, é a fase do processo de compilação que a ferramenta tem a tarefa de ler o código fonte a partir de um arquivo [Louden], lendo caractere a caractere e submetendo-os às regras que geram os *Tokens* desta linguagem. Tokens são basicamente palavras em linguagem natural formado pela sequência de caracteres lida do arquivo fonte, e representa uma unidade de informação.

Uma vez que a tarefa de extração dos tokens é executada, é necessário então a aplicação dos padrões de combinação que classificam essas unidades de informação, de acordo com as classes de uma linguagem, tais como: palavras-chave, identificadores ou constantes numéricas.

2.1. Especificação da Linguagem de Programação T++

A linguagem T++ objeto desse estudo é relativamente simples nos termos de suas funções, mas tem a complexidade inerente das questões de implementação de qualquer linguagem. Abaixo estão discriminados os símbolos, descritores, operações e opções suportadas pela linguagem.

2.1.1. Tipos de Dados

Por se tratar uma linguagem experimental para estudo de compiladores, T++ apresenta apenas os tipos de dados primitivos:

- Inteiro;
- Flutuante.

Entretanto, a linguagem permite a criação de Tipos de Vetores e Matrizes a partir desses tipos primitivos, na seguinte sintaxe, além do reconhecimento de números em notação científica:

```
inteiro : var1
flutuante : var1
var1 := 1.54334E-34
inteiro : arr1[indice1]
inteiro : arr2[indice1][indice2]
```

2.1.2. Palavras Reservadas

A linguagem atente ao parâmetro de uma linguagem discorrida no Português, o que significa o suporte à caracteres especiais dessa língua, como acentuação e "ç".

Suas palavras reservadas são: **se, então, senão, fim, repita, flutuante, retorna, até, leia, escreve, e inteiro**. A seguir, um exemplo de código em T++:

```

inteiro: n

inteiro fatorial(inteiro: n)
  inteiro: fat
  se n > 0 então {não calcula se n > 0}
    fat := 1
    repita
      fat := fat * n
      n := n - 1
    até n = 0
    retorna(fat) {retorna o valor do fatorial de n}
senão
  retorna(0)
fim

fim

inteiro principal()
  leia(n)
  escreva(fatorial(n))
  retorna(0)
fim

```

Figura 1. Código de exemplo em T++.

2.1.3. Símbolos e Operadores

Além do conjunto de palavras reservadas, essa linguagem também tem uma série de operadores mapeados para operações específicas:

- “*” : operador de *multiplicação*;
- “-” : operador de *subtração*;
- “+” : operador de *soma*;
- “:=” : operador de *atribuição*;
- “:” : símbolo *dois-pontos*;
- “<” : operador de comparação *menor-quê*;
- “>” : operador de comparação *maior-quê*;
- “=” : operador de *igualdade*;
- “<=” : operador de comparação *menor ou igual a*;
- “>=” : operador de comparação *maior ou igual a*;
- “[” : símbolo abre *colchetes*;
- “]” : símbolo fecha *colchetes*;
- “(” : símbolo abre *parênteses*;
- “)” : símbolo fecha *parênteses*;
- “!” : operador lógico de *negação*;
- “&&” : operador lógico *conjunção*;
- “||” : operador lógico *disjunção*;

2.2. Especificação Formal dos Autômatos

O sistema de varredura foi implementado com base em expressões regulares, descritas como se segue, gerando as marcas para identificação dos lexemas da linguagem:

- Palavra Reservada **t_PR** = (se)|(então)|(senão)|(fim)|(repita)| (flutuante)|(retorna)|(até)|(leia)|(escreva)|(inteiro)

Para cada conjunto de caractere que representa uma palavra-chave é necessário criar uma regra específica que dará forma a um autômato para a aquela palavra, depois se unindo com os demais autômatos de outras regras.

- Identificador **t_ID** = [a-zA-Zá-ñÁ-Ñ][a-zA-Zá-ñÁ-Ñ0-9_]

Regra que combina as definições de números e letras, acentuadas ou não, e underline para identificar marcadores de identificadores da linguagem;
- Inteiro **t_INTEIRO** = ?[0-9][0-9]*

Regra genérica para captura de números inteiros;
- Flutuante **t_FLUTUANTE** = ?[0-9][0-9]*[0-9]+([Ee][-]?[0-9]+)?

Regra genérica para captura de números flutuantes e em notação científica;
- Branco **t_ignore** = [\t]

Regra para captura de símbolos de espaço e outros não imprimíveis que não fazem parte do programa;
- Símbolo **SB** = (:)=|(:)|(>)|(<)|(|)|(|)|(|)(*)|(=)|(-)|(+)|(/)|(=<)|(>=)|(|)|(|)

Da mesma forma como ocorre com as palavras reservadas, um autômato de símbolo é composto pelos estados com transições com os caracteres específicos dos símbolos da linguagem, que se une com as demais regras formando o autômato da linguagem.
- Comentário **t_COMMENT** = {[~]}*

Regra que casa abertura e fechamento de comentários reconhecido pela linguagem;
- Erro de comentário **t_ERRORCT** = [{}]

Regra de exceção para abertura ou fechamento sem sequência de comentários;

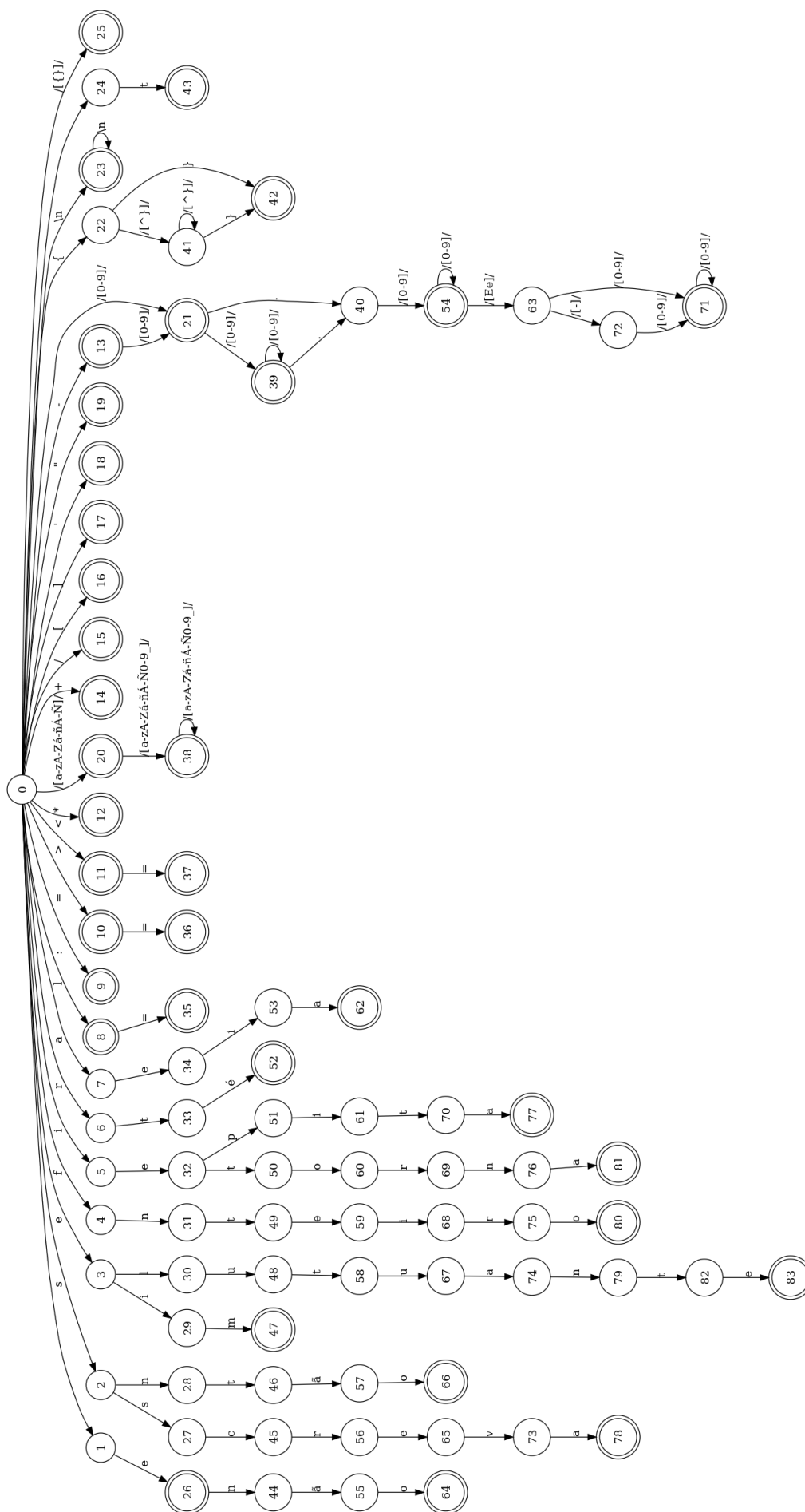


Figura 2. Autômato resultante da união das partes do analisador.

2.3. Detalhes de Implementação da Ferramenta

Considerando experiências passadas com o Java na implementação de ferramentas de compilação, para esse novo projeto foi usado o Python como linguagem-base para implementação com auxílio da biblioteca PLY.

O texto base para desenvolvimento da ferramenta pode ser encontrado em <http://www.dabeaz.com/ply/ply.html> [Beazley], que conta com exemplos básicos, detalhes de formulação de regras e modelo projeto.

Diferentemente da Linguagem Java, para o analisador léxico construído com o PLY é necessária a criação de um único arquivo de programação que representa uma única classe, nesse caso, chamada apenas de *Lexer*, que possui os métodos *Test*.

2.4. Resultados da Varredura

A entrada desse programa é composta de dois arquivos: o primeiro contém a descrição das regras de interpretação da linguagem em *Lex*; o segundo, contém o código da linguagem que deseja-se compilar.

O processo trata da leitura caractere a caractere do arquivo de entrada do programa - código fonte - como um *stream*, onde a partir dele, são buscadas combinações de regras léxicas dada a máquina de estados da linguagem que indica a criação das marcas dessa linguagem.

Como saída, é exibida no console do programa, a combinação <classe, valor> para cada marca válida da linguagem. Em situação de caractere ou sequência inválidos é apontado um erro, com especificação de caractere não identificado ou sequência inválida, como abertura e não fechamento de *chaves* que indica um comentário incompleto.

O projeto foi escrito usando a IDE PyCharm, e pode ser executado em `Run > Run > nome_programa`, sendo necessário também informar por parâmetro o nome do arquivo-fonte do programa a ser analisado.

Assumindo que o programa chegou até o caminho dos arquivos de entrada a saída será exibida dessa forma:

```

lexer.py
1  # -*- coding: utf-8 -*-
2  #
3  # lexer.py
4  # Analisador léxico para a linguagem T++
5  # Autores: Sávio Camacá
6  #
7
8  import ply.lex as lex
9  from ply.ctokens import t_ID
10
11
12 class Lexer:
13     def __init__(self):
14         self.lexer = lex.lex(debug=False, module=self)
15
16     keywords = {
17         u'se': 'SE',
18         u'então': 'ENTÃO',
19         u'senão': 'SENÃO',
20         u'fim': 'FIM',
21         u'repita': 'REPITA',
22         u'flutuante': 'FLUTUANTE',
23         u'retorna': 'RETORNA',
24         u'até': 'ATE',
25         u'leia': 'LEIA',
26         u'escreve': 'ESCREVE',
27         u'inteiro': 'INTEIRO',
28     }
29
30     tokens = ['ASS', 'COLON', 'LPAR', 'RPAR', 'IF',
31              'LEQ', 'GEQ', 'EQ', 'NEQ', 'ID',
32              'ID', 'NOT', 'AND', 'OR'] + list(keywords.values())
33
34     t_ASS = r':='
35     t_COLON = r':'
36     t_LPAR = r'('
37     t_RPAR = r')'
38     t_ADD = r'+'
39     t_SUB = r'-'
40     t_TIMES = r'*'
41     t_DIV = r'/'
42     t_LEQ = r'<='
43     t_GEQ = r'>='
44     t_EQ = r'='

```

```

segundo.tpp
1  {1. (Valor: 2.0) . Dizemos que uma matriz quadrada
2
3  inteiro principal()
4      inteiro: A[tamanho_matriz][tamanho_matriz]
5      inteiro: linha[tamanho_matriz]
6      inteiro: coluna[tamanho_matriz]
7      inteiro: diagonal_principal
8      inteiro: diagonal_secundaria
9      inteiro: i
10     inteiro: j
11     inteiro: contador
12     inteiro: valor
13     inteiro: somador
14     inteiro: somador_temporario
15
16     i := 0
17     j := 0
18     diagonal_principal := 0
19     diagonal_secundaria := 0
20     contador := 0
21     valor := 0
22     somador := 0
23
24     repita
25         repita
26             leia(valor)
27             A[i][j] = valor
28             linha[i] = linha[i] + A[i][j]
29             coluna[j] = coluna[j] + A[i][j]
30
31             se i = j então
32                 diagonal_principal = diagonal_principal + valor
33                 fim
34
35             se (i + j + 1) = tamanho_matriz então
36                 diagonal_secundaria = diagonal_secundaria + valor
37                 fim
38
39             ate j < tamanho_matriz
40             ate i < tamanho_matriz
41
42     repita
43         se coluna[i] = linha[i] então

```

```

lexer (1)
lexer (1)
C:\Users\savio\AppData\Local\Programs
LexToken(INTEIRO,'inteiro',3,311)
LexToken(ID,'principal',3,319)
LexToken(LPAR,'(',3,328)
LexToken(RPAR,')',3,329)
LexToken(INTEIRO,'inteiro',4,332)
LexToken(COLON,':',4,339)
LexToken(ID,'A',4,341)
LexToken(LBR,'[',4,342)
LexToken(ID,'tamanho_matriz',4,343)
LexToken(RBR,']',4,357)
LexToken(LBR,'[',4,358)
LexToken(ID,'tamanho_matriz',4,359)
LexToken(RBR,']',4,373)
LexToken(INTEIRO,'inteiro',5,376)
LexToken(COLON,':',5,383)
LexToken(ID,'linha',5,385)
LexToken(LBR,'[',5,390)
LexToken(ID,'tamanho_matriz',5,391)
LexToken(RBR,']',5,405)
LexToken(INTEIRO,'inteiro',6,408)
LexToken(COLON,':',6,415)
LexToken(ID,'coluna',6,417)
LexToken(LBR,'[',6,423)
LexToken(ID,'tamanho_matriz',6,424)
LexToken(RBR,']',6,438)
LexToken(INTEIRO,'inteiro',7,441)
LexToken(COLON,':',7,448)
LexToken(ID,'diagonal_principal',7,449)
LexToken(INTEIRO,'inteiro',8,470)
LexToken(COLON,':',8,477)
LexToken(ID,'diagonal_secundaria',8,478)
LexToken(INTEIRO,'inteiro',9,500)
LexToken(COLON,':',9,507)
LexToken(ID,'i',9,509)
LexToken(INTEIRO,'inteiro',10,512)
LexToken(COLON,':',10,519)
LexToken(ID,'j',10,521)
LexToken(INTEIRO,'inteiro',11,524)
LexToken(COLON,':',11,531)
LexToken(ID,'contador',11,533)
LexToken(INTEIRO,'inteiro',12,543)
LexToken(COLON,':',12,550)
LexToken(ID,'valor',12,552)

```

Figura 3. Situação de saída do programa.

3. Conclusão

Com a conclusão desta etapa da implementação da ferramenta de compilação para a linguagem T++, os conceitos de design de compiladores continuam a ser fixados e aprimorados com a revisão da estrutura do analisador, agora no uso da biblioteca PLY, o que mostrou ser mais simples, ao contrário do uso do java CUP que envolvia uma quantidade muito maior de classes e arquivos, que complicava o entendimento do processo de geração automática e agravava a migração do projeto para a fase seguinte, de implementação do analisar sintático.

Referências

Beazley, D. M. Ply version: 3.9.

Louden, K. C. *Compiler Construction Principles*.