

Tradutores - Analisador Sintático

Thiago Veras Machado^[160146682]

Universidade de Brasília cic@unb.com

1 Introdução

O projeto da disciplina Tradutores tem como principal objetivo estudar os aspectos teóricos relacionados à implementação de tradutores quanto à prática de sua implementação.

Nesse trabalho será implementado um tradutor para a linguagem C adaptada, no qual utilizaremos como base o livro da disciplina [ALSU07].

2 Motivação

O trabalho possui um desafio interessante pois será integrado uma nova estrutura de dados *set* e de uma primitiva *elem*. Juntamente com essa primitiva, novos métodos serão acrescentados (*add*, *remove*, *in*, *is_set*, *exists* e *forall*) nos quais executam tais operações em um *set* já declarado.

A primitiva *elem* seria um novo tipo de dado, no qual pode ser do tipo *int*, *float* e *set*, trazendo uma versatilidade maior para o trabalho.

A introdução da nova primitiva *set* possui uma certa importância pois irá complementar uma falta dessa estrutura de dados em C e também a falta de tipos de variáveis genéricas, que podemos ver presente na linguagem *C++*, com isso estamos melhorando a versão da linguagem *C*.

3 Descrição da análise léxica

C é uma linguagem de propósito geral, compilada, de alto nível, com sintaxe estruturada, imperativa e possui tipagem estática. Para o trabalho a linguagem conterá as seguintes estruturas básicas:

- Estrutura condicional: **if** e **else**.
- Estrutura de repetição: **for**.
- Tipos de dados: **int** (números inteiros), **float** (números no formato de ponto flutuante) **set** (estrutura de dados sobre conjuntos) e **elem** (tipo de dado polimórfico).
- Definição e chamada de subrotinas com passagem de parâmetros.
- Operações aritméticas (+, -, *, /), lógicas básicas (!, &&, !) e relacionais (==, >, <, !=, >=, <=).

Além da utilização do livro-texto [ALSU07], também foi utilizado o flex [Est], que consiste em uma ferramenta geradora de programas que reconhecem padrões léxicos em textos. A estrutura do código flex foi pensada e implementada com o intuito de facilitar a criação da tabela de símbolos futuramente.

4 Análise sintática

A análise sintática foi realizada utilizando o Bison [CS21], já para as regras, foi utilizada a gramática contida no relatório (final da página).

O bison utiliza as regras definidas para gerar um analisador sintático LR(1), como ele é da forma bottom-up, então a árvore é gerada do terminal até a raiz.

```
typedef struct Symbol {
    int line , colum;
    char *classType , *type , *body;
} Symbol;
```

Listing 1.1. Definição de Symbol (table.h)

A Estrutura 1.1 foi utilizada para a criação dos símbolos, que será utilizada tanto para a criação da árvore, quanto da tabela. Os campos de sua estrutura foram escolhidos com a ideia de armazenar em qual linha e coluna o token se encontra, qual sua classe (variável ou função), seu tipo (inteiro, float, set, elem), o seu corpo (o nome da variável ou função em si).

4.1 Criação tabela de símbolos

A ideia por trás da criação da tabela é simplesmente a utilização de um vetor de símbolos no qual é inserido cada token que aparece na árvore (limitado somente em tokens de variáveis e funções), com a finalidade de ser utilizado na análise semântica posteriormente.

LINE	COLUMN	CLASS	TYPE	SCOPE	BODY
1	5	function	INT	0	main
2	2	variable	INT	1	x
3	3	variable	set	1	y

4.2 Criação árvore sintática

A criação da árvore sintática foi feita com uma estrutura de dados dinâmica (tree), no qual temos um nó da árvore, que é definido da seguinte forma:

```
typedef struct TreeNode {
    struct TreeNode *children , *nxt;
    char* rule;
    Symbol* symbol;
} TreeNode;
```

Listing 1.2. Estrutura da árvore (tree.h)

A ideia se baseia na premissa de que árvore é um conjunto de blocos, no qual cada nó da árvore pode ser ligado em um filho, e com isso, a profundidade é incrementada. Cada "vizinho" é representado pelo campo *nxt*, que é um ponteiro para o próximo bloco da árvore. Um nó, por sua vez, possui um campo *rule*, no

qual armazena qual o *head* dessa regra e o *body* da mesma pode conter algum símbolo, com isso optei por armazená-lo a fim de mostrar posteriormente na estrutura da árvore.

Vale se atentar que foram omitidas algumas regras na hora de imprimir a estrutura, com a finalidade de se comportar como uma árvore abstrata.

```

├─ start
│  └─ function_definition
│     │  └─ function_declaration - INT  ─ [1:5] variable : main
│     │  └─ statements
│     │     │  └─ variables_declaration - INT  ─ [2:9] variable : x
│     │     │  └─ statements
│     │     │     │  └─ variables_declaration - SET  ─ [3:9] variable : y
│     │     │     │  └─ statements
│     │     │     │     │  └─ set_assignment_expression ─ [4:16] variable : x
│     │     │     │     │  └─ statements
│     │     │     │     │     │  └─ io_statement - READ  ─ [5:10] variable : x
│     │     │     │     │     │  └─ io_statement - WRITELN ─ [6:13] string : "Hello"

```

5 Análise semântica

A análise semântica foi realizada utilizando o Bison [CS21] e todas as estruturas que foram produzidas durante a análise sintática, como tabela de símbolos, pilha de escopos e árvore sintática.

Semanticamente, um trecho de código é analisado com base na lógica implementada, podemos dizer que é uma análise para ver se “faz sentido”. Um exemplo dessa análise seria a quantidade de parâmetros de uma função, variáveis duplicadas, cast de tipos, variáveis sendo utilizadas como função, necessidade de uma função main e etc.

5.1 Escopo

A *pilha de escopo* serve para saber quais escopos podem ser acessados com base no escopo atual (topo da pilha), com isso, ela pode ser utilizada para diversas verificações, como variável / função duplicada ou não declarada.

A regra para alteração de escopo é definida quando se encontra um token '{', no qual indica que um novo escopo será criado, dando push no novo id do escopo na pilha. Quando o analisador encontra um token '}', ele dá pop na pilha, pois agora esse escopo fechou e não pode ser acessado diretamente no contexto atual, como por exemplo, podem ser declarado novas variáveis com o mesmo nome das variáveis do escopo recém fechado.

5.2 Checagem de Variável / função

Para analisar se uma função / variável está duplicada, utilizo 2 estruturas de dados criada na etapa sintática: *Tabela de símbolos* e *pilha de escopos*. A ideia é realizar uma travessia pela pilha e verificar se existe uma variável com aquele nome naquele escopo específico.

Após a travessia pela pilha, é possível tanto determinar que a variável / função não foi declarada (pois só podem ser utilizadas após a declaração) quanto

Vale atentar que quando estamos atribuindo o valor a uma variável, precisamos *"forçar"* o tipo a ser o mesmo da variável, pois como dou prioridade para *elem* e *float*, para deixar a operação mais genérica e fácil de resolver os conflitos internos das expressões, quando atribuímos o tipo na variável pode-se resultar em um tipo diferente, como por exemplo, atribuir um *float* para um *int*, neste caso, precisamos forçar o *float* a virar *int*. Com isso, teremos *cast* de *cast*.

```

| return
|   | expression_additive -- [4:25] additive operator : (int)(ELEM + (elem)INT)
|   |   | expression_additive -- [4:20] additive operator : (elem)FLOAT + ELEM
|   |   |   | value -- [4:18] variable : FLOAT b
|   |   |   | value -- [4:22] variable : ELEM e1
|   |   | const -- [4:27] INT : 10

```

6 Descrição dos arquivos de teste

Os arquivos de testes se encontram na pasta ***Input*** no qual possuem o prefixo *error_* representam os testes que possuem erros a partir de um análise léxica, análogo para os arquivos que possuem o prefixo *success_*.

Arquivos de erro:

```

# bash run.sh
# ./bison Input/error_1.c

Input/error_1.c:9:5: error: called object 'globalSet' is not a function or function pointer
Input/error_1.c:1:5: note: declared here

Input/error_1.c:11:9: error: incompatible types for argument of 'funcSuave'
Expected: funcSuave(SET, INT, INT)
Got: funcSuave(INT, SET, INT)

Input/error_1.c:?:?: undefined reference to 'main'
Program analysis failed!
Analysis terminated with 3 error(s)

# bash run.sh
# ./bison Input/error_2.c

Input/error_2.c:3:33: error: redeclaration of 'h' with no linkage
Input/error_2.c:3:24: note: previous definition of 'h' was here

Input/error_2.c:5:11: error: miss type expression
Expected: FLOAT == FLOAT
Got: FLOAT == SET

Input/error_2.c:6:21: error: miss type return of 'funcRepetido'
Expected: return (SET);
Got: return (FLOAT);

Input/error_2.c:7:17: error: invalid operands to binary == (have 'SET' and 'SET')
Input/error_2.c:12:5: error: redeclaration of 'x' with no linkage
Input/error_2.c:1:5: note: previous definition of 'x' was here

[LEXIC] ERROR line: 15 columns: 18 unidentified char: #

[SYNTATIC] [16,12] ERROR syntax error, unexpected INT_VALUE, expecting ID

Program analysis failed!
Analysis terminated with 7 error(s)

```

7 Compilação e execução do programa.

Para facilitar a compilação e execução do programa, criei um script em *bash* que executa todas as etapas automaticamente para você:

```
bash run.sh
```

Para executar individualmente cada teste com o intuito de ver a análise de cada token e geração de árvore / tabela de símbolos, basta executar:

```
./bison Input/NOME_DO_ARQUIVO.C
```

Ambiente utilizado para a criação do trabalho foi um Ubuntu 20, com terminal *bash* e versão do *bison* 3.7.6.

Em caso de não rodar na máquina do usuário, providenciei uma imagem *docker* que possui um ambiente *ubuntu* no qual é possível realizar os testes (vem com *bison*, *flex* e *valgrind*).

```
docker build -t 160146682_semantico .
```

```
docker run --rm -it --entrypoint="sh" 160146682_semantico
```

Referências

- [ALSU07] A.V. Aho, M.S. Lam, R. Sethi, and J.D. Ullman. *Compilers: Principles, Techniques, Tools*. Pearson/Addison Wesley, 2nd edition, 2007.
- [CS21] Robert Corbett and Richard Stallman. Bison. <https://www.gnu.org/software/bison/manual/bison.pdf>, Online; acessado 18 de Março de 2021.
- [Est] W. Estes. Flex: Fast lexical analyser generator. <https://github.com/westes/flex>. Online; acessado 24 de Fevereiro de 2021.
- [Gup21] Ajay Gupta. The syntax of c in backus-naur form. <https://tinyurl.com/max5eep>, Online; acessado 24 de Abril de 2021.

A Gramática

A gramática abaixo teve que ser alterada comparada com a anterior e descreve a linguagem para qual o compilador será implementado [Gup21].

$\langle \text{start} \rangle$	$::= \langle \text{program} \rangle$
$\langle \text{program} \rangle$	$::= \langle \text{function-definition} \rangle$ $ \langle \text{function-definition} \rangle \langle \text{program} \rangle$ $ \langle \text{variables-declaration} \rangle \langle \text{program} \rangle$
$\langle \text{function-definition} \rangle$	$::= \langle \text{function_declaration} \rangle ' (' \langle \text{parameters} \rangle ') ' \langle \text{function_body} \rangle$ $ \langle \text{function_declaration} \rangle ' (' ' ' \langle \text{function_body} \rangle$
$\langle \text{function-declaration} \rangle$	$::= \langle \text{type-identifier} \rangle \langle \text{ID} \rangle$
$\langle \text{function-body} \rangle$	$::= \{ ' \langle \text{statements} \rangle ' \}$ $ \{ ' ' ' \}$
$\langle \text{parameters} \rangle$	$::= \langle \text{parameters-list} \rangle$
$\langle \text{parameters-list} \rangle$	$::= \langle \text{parameters_list} \rangle ' , ' \langle \text{parameter} \rangle$ $ \langle \text{parameter} \rangle$
$\langle \text{parameter} \rangle$	$::= \langle \text{type-identifier} \rangle \langle \text{ID} \rangle$
$\langle \text{type-identifier} \rangle$	$::= \text{INT}$ $ \text{FLOAT}$ $ \text{ELEM}$ $ \text{SET}$
$\langle \text{statements} \rangle$	$::= \langle \text{statement} \rangle \langle \text{statements} \rangle$ $ \langle \text{statement} \rangle$
$\langle \text{statements_braced} \rangle$	$::= \{ ' \langle \text{statements} \rangle ' \}$ $ \{ ' ' ' \}$
$\langle \text{statement} \rangle$	$::= \langle \text{variables-declaration} \rangle$ $ \langle \text{statements_braced} \rangle$ $ \langle \text{return} \rangle$ $ \langle \text{conditional} \rangle$ $ \langle \text{for} \rangle$ $ \langle \text{expression_statement} \rangle$ $ \langle \text{io_statement} \rangle$ $ \langle \text{set_pre_statement} \rangle$
$\langle \text{set_pre_statement} \rangle$	$::= \langle \text{set_statement_for_all} \rangle$
$\langle \text{set_statement_add_remove} \rangle$	$::= \text{ADD} ' (' \langle \text{set_boolean_expression} \rangle ') '$ $ \text{REMOVE} ' (' \langle \text{set_boolean_expression} \rangle ') '$
$\langle \text{set_statement_for_all} \rangle$	$::= \text{FOR_ALL} ' (' \langle \text{set_assignment_expression} \rangle ') ' \langle \text{statements} \rangle$
$\langle \text{set_statement_exists} \rangle$	$::= \text{EXISTS} ' (' \langle \text{set_assignment_expression} \rangle ') ' \langle \text{statements} \rangle$
$\langle \text{set_boolean_expression} \rangle$	$::= \langle \text{expression} \rangle \text{IN} \langle \text{set_statement_add_remove} \rangle$ $ \langle \text{expression} \rangle \text{IN ID}$
$\langle \text{set_assignment_expression} \rangle$	$::= \text{ID IN} \langle \text{set_statement_add_remove} \rangle$ $ \text{ID IN ID}$

$$\begin{aligned}
\langle \text{expression_statement} \rangle &::= \langle \text{expression} \rangle \text{ ';' } \\
\langle \text{expression_or_empty} \rangle &::= \langle \text{expression} \rangle \\
&\quad | \quad \langle \%empty \rangle \\
\langle \text{expression} \rangle &::= \langle \text{expression_assignment} \rangle \\
\langle \text{expression_assignment} \rangle &::= \langle \text{expression_logical} \rangle \\
&\quad | \quad \text{ID '=' } \langle \text{expression} \rangle \\
\langle \text{expression_logical} \rangle &::= \langle \text{expression_relational} \rangle \\
&\quad | \quad \langle \text{set_boolean_expression} \rangle \\
&\quad | \quad \langle \text{expression_logical} \rangle \text{ AND_OP } \langle \text{expression_relational} \rangle \\
&\quad | \quad \langle \text{expression_logical} \rangle \text{ OR_OP } \langle \text{expression_relational} \rangle \\
\langle \text{expression_relational} \rangle &::= \langle \text{expression_additive} \rangle \\
&\quad | \quad \langle \text{expression_relational} \rangle \text{ RELATIONAL_OP } \langle \text{expression_additive} \rangle \\
\langle \text{expression_additive} \rangle &::= \langle \text{expression_multiplicative} \rangle \\
&\quad | \quad \langle \text{expression_additive} \rangle \text{ ADDITIVE_OP } \langle \text{expression_multiplicative} \rangle \\
\langle \text{expression_multiplicative} \rangle &::= \langle \text{expression_value} \rangle \\
&\quad | \quad \langle \text{expression_multiplicative} \rangle \text{ MULTIPLICATIVE_OP } \\
&\quad \quad \langle \text{expression_value} \rangle \\
\langle \text{expression_value} \rangle &::= \text{'(' } \langle \text{expression} \rangle \text{' '}' \\
&\quad | \quad \text{'!' '(' } \langle \text{expression} \rangle \text{' '}' \\
&\quad | \quad \text{ADDITIVE_OP '(' } \langle \text{expression} \rangle \text{' '}' \\
&\quad | \quad \langle \text{value} \rangle \\
&\quad | \quad \text{'!' } \langle \text{value} \rangle \\
&\quad | \quad \text{ADDITIVE_OP } \langle \text{value} \rangle \\
&\quad | \quad \langle \text{is_set_expression} \rangle \\
&\quad | \quad \langle \text{set_statement_exists} \rangle \\
&\quad | \quad \langle \text{set_statement_add_remove} \rangle \\
\langle \text{is_set_expression} \rangle &::= \text{IS_SET '(' expression '}' \\
&\quad | \quad \text{'!' IS_SET '(' expression '}' \\
\langle \text{for} \rangle &::= \text{FOR '(' } \langle \text{for_expression} \rangle \text{' '}' \langle \text{statements} \rangle \\
\langle \text{for_expression} \rangle &::= \langle \text{expression_or_empty} \rangle \text{ ';' } \langle \text{expression_or_empty} \rangle \text{ ';' } \\
&\quad \langle \text{expression_or_empty} \rangle \\
\langle \text{io_statement} \rangle &::= \text{READ '(' ID ') '}' \\
&\quad | \quad \text{WRITE '(' STRING ') '}' \\
&\quad | \quad \text{WRITE '(' CHAR ') '}' \\
&\quad | \quad \text{WRITE '(' expression ') '}' \\
&\quad | \quad \text{WRITELN '(' STRING ') '}' \\
&\quad | \quad \text{WRITELN '(' CHAR ') '}' \\
&\quad | \quad \text{WRITELN '(' expression ') '}' \\
\langle \text{arguments_list} \rangle &::= \langle \text{arguments_list} \rangle \text{ ',' } \langle \text{expression} \rangle \\
&\quad | \quad \langle \text{expression} \rangle \\
\langle \text{conditional} \rangle &::= \text{IF '(' } \langle \text{expression} \rangle \text{' '}' \langle \text{statement} \rangle \text{ prec THEN} \\
&\quad | \quad \text{IF '(' } \langle \text{expression} \rangle \text{' '}' \langle \text{statement} \rangle \text{ ELSE } \langle \text{statement} \rangle \\
\langle \text{return} \rangle &::= \text{RETURN } \langle \text{expression} \rangle \text{ ';' } \\
&\quad | \quad \text{RETURN ';' }
\end{aligned}$$


```
 $\langle value \rangle$  ::= ID  
          |  $\langle const \rangle$   
          |  $\langle function\_call \rangle$   
 $\langle function\_call\_statement \rangle$  ::=  $\langle function\_call \rangle$  ;  
 $\langle function\_call \rangle$  ::= ID '('  $\langle arguments\_list \rangle$  ')'  
                  | ID '(' ')'  
 $\langle variables\_declaration \rangle$  ::=  $\langle type\_identifier \rangle$  ID ';' ;  
 $\langle const \rangle$  ::= INT_VALUE  
            | FLOAT_VALUE  
            | EMPTY
```