

Gerador de Código Intermediário

Lucas Mafra Chagas - 12/0126443

Prof.^a Dra. Cláudia Nalon

Resumo Esse projeto tem como intuito construir um analisador sintático para um subconjunto da linguagem C, mostrando sua nova primitiva de linguagem, sua descrição, instrução e gramática.

Keywords: Tradutores · Compiladores · Mini-C.

1 Objetivo

Tendo como base a metodologia de ensino baseada em projeto, determinou-se algumas diretrizes que o discente da Ciência da Computação necessita seguir para alcançar um melhor entendimento do conteúdo apresentado. Dessa forma, as diretrizes definidas consistem na implementação de um analisador de um subconjunto da linguagem C. Por fim, este projeto é constituído por quatro etapas: Implementação do analisador léxico; Implementação do analisador sintático; Implementação do analisador semântico; e Implementação do gerador de código intermediário.

2 Motivação

Um compilador é um programa que lê um programa escrito numa linguagem, conhecida como linguagem fonte, e o traduz num programa equivalente numa outra linguagem, conhecida como linguagem alvo. Essa ação é feita observando possíveis erros presentes no programa fonte [1].

Sendo assim, a linguagem projetada para o desenvolvimento do trabalho final segue as tratativas básicas da linguagem C e a adição de uma nova primitiva de dados, que é declarado como *set*. A declaração de uma variável desse tipo segue o molde de C, sem ter um tipo associado aos seus elementos, onde o *set* tem como papel ser composto por *elem*, onde *elem* é polimórfico, podendo assumir o papel de *int*, *float* ou de *set* [5].

Além disso, o novo tipo de dado da linguagem, o *set* possui operações específicas, como verificação de tipagem, *is_set*, a inclusão de elemento em um conjunto, *add*, a remoção de elemento de um conjunto, *remove*, a seleção dentro desse conjunto, *exists*, e a iteração dentro desse conjunto, *forall* [5].

3 Descrição da Análise Léxica

A análise léxica é a primeira fase de um projeto de compilador. O analisador léxico é responsável por converter uma sequência de caracteres em uma sequência de lexemas, removendo qualquer comentário ou espaço em branco. O analisador léxico é responsável por varrer o código fonte do programa e identificar cada token, um por um [6]. Durante essa fase, para uma melhor compreensão da linguagem desenvolvida, é possível construir uma gramática livre de contexto que determina o seu funcionamento. Além da gramática, a construção de uma tabela de símbolos é responsável por classificar os nomes das variáveis, as funções, as classes e as constantes também auxilia no entendimento [4].

No projeto desenvolvido, além das regras e expressões regulares, o analisador léxico necessitou da estruturação de algumas funcionalidades para a preparação do algoritmo. Essas novas estruturas servem de base para a construção da tabela de símbolos na próxima etapa. Para o tratamento de erros, foi construída uma função que destaca um token inválido ao ser encontrado, mostrando a linha e a coluna onde esse erro aconteceu.

4 Descrição da Análise Sintática

A análise sintática é a segunda fase de um projeto de compilador. O analisador sintático é responsável por garantir que os lexemas identificados pelo analisador léxico sejam descritos de forma adequada aos seus programas [2]. Os tokens, que vem agrupados nos lexemas, possuem dois componentes: os seus nomes e seus valores de atributos. Os nomes dos tokens aparecem como símbolos terminais da gramática da linguagem e seus atributos direcionam para a tabela de símbolos que possuem mais informações sobre os tokens [1].

Nessa segunda etapa, além da revisão do analisador léxico, foi necessário construir uma árvore sintática abstrata, que representa a estrutura sintática hierárquica do programa; e a estrutura da tabela de símbolos, que são estruturas de dados utilizados pelo compilador para obter informações sobre as construções do programa-fonte.

Para a simplificação da construção da árvore sintática abstrata, foi utilizado o conceito de árvore binária para facilitar o desenvolvimento desta etapa, onde a derivação de um símbolo terminal pode crescer para a esquerda ou direita. Nesta estrutura de dados, os nodos intermediários representam os símbolos não terminais, os nodos folhas representam os tokens presentes no código fonte e a raiz representa o programa analisado.

Para os nodos intermediários, foram armazenados apenas as informações consideradas necessárias para a construção da árvore, mostrando apenas as derivações que chegam num símbolo terminal. Nos nodos folhas, foram armazenadas as suas informações, mostrando o valor atribuído aquele símbolo terminal.

No projeto, foi utilizado uma tabela *hash* para a construção das Tabelas de Símbolos. Dentro desse dicionário, foi armazenada o símbolo terminal, com o seu valor, sua tipagem e seu escopo. O intuito desse armazenamento é garantir

com que seja possível fazer a análise semântica das variáveis presentes no código, passando as informações de declarações para o uso do analisador semântico.

5 Descrição da Análise Semântica

A análise semântica é a terceira fase de um projeto de compilador. Isto porque o analisador semântico é responsável por coletar as informações dos tokens e colocar na tabela de símbolos, sintática verificação de tipo, verificação de etiqueta e controle de fluxo. O analisador semântico produz algum tipo de representação do programa, seja o código do objeto ou uma representação intermediária do programa [3] [1].

Nessa terceira etapa, além da revisão do analisador sintático, foi necessário otimizar a tabela de símbolos, para fazer a verificação de tipo, de etiqueta e um controle melhor do fluxo. Além disso, o programa foi otimizado para um compilador de duas etapas. Para isso, o código foi construído para construção da árvore abstrata apenas ao fim das passagens de leitura do código.

Para garantir o pleno funcionamento da análise semântica, é necessário verificar se os seguintes pontos estão presentes no trabalho: a conversão implícita de tipos, a verificação da chamada da função ou variável após a sua declaração, as regras de escopo, as regras para passagem de parâmetros e a existência de uma função main.

A conversão implícita de tipos na linguagem desenvolvida segue os mesmos parâmetros do C, onde em expressões envolvendo operadores com operandos de tipos diferentes, os valores dos operandos são avaliados e convertidos para o mesmo tipo antes da operação ser executada. Nesse caso, os tipos mais simples, como int, são convertidos para os tipos mais complexos, como o float.

No projeto, para a realização dessa etapa, foi necessário armazenar dentro da árvore sintática anotada a tipagem do valor correspondente.

A verificação da chamada de uma função ou variável após a sua declaração tem o papel de verificar a checagem de tipos e a unicidade do identificador. Nesse caso, é observado se a função ou variável foi declarada apenas uma vez, se foi declarada ou definida antes do seu primeiro uso, se é declarada mas nunca utilizada, a que essas declarações se refere, se os tipos de uma expressão são compatíveis e se as dimensões apresentam uma consistência com o que foi declarado.

No projeto, para a verificação da chamada de uma função ou variável, foi utilizado uma estrutura que tem os valores da tabela de símbolos armazenados para verificar a sua redeclaração.

Por se tratar de uma linguagem que faz parte de um subconjunto da linguagem C, as suas regras de escopo também respeitam as regras de escopo do C, sendo um escopo estático. Na execução do projeto, foi necessário verificar se as variáveis são locais ou globais, e como elas se comportam em cada ambiente. Caso uma variável seja declarada globalmente e localmente, dentro de funções as variáveis locais irão sobrepor o valor da variável global.

A função main é considerada o local de início da execução de um programa C. ela é única no código fonte e difere dos outros identificadores. Para isso, é

necessário verificar se ela está presente no programa e se ela não é declarada mais de uma vez.

6 Gerador de Código Intermediário

Como parte final do *front end* de um compilador, o gerador de código intermediário é responsável por receber a entrada de sua fase predecessora, o analisador semântico, na forma de uma árvore de sintaxe anotada e tabela de símbolos. A tabela de símbolos é responsável por mostrar quais as variáveis presentes no programa e a árvore sintática deve ser convertida numa representação linear, sendo a sua geração de código.

Para a geração de código intermediário, foi separada as declarações e as operações em um operador e dois argumentos. Ao realizar uma ação, o operador recebe o seu tipo, como uma operação aritmética, uma operação lógica ou mesmo uma declaração condicional, e comprime as expressões dentro dessa ação em duas variáveis, os dois argumentos.

A abordagem utilizada para o compilador foi a de apenas uma passagem. Embora não seja tão eficiente quanto os compiladores de duas passagens, o compilador de uma passagem é menor e mais rápido, condizente para a linguagem proposta e para a simplificação do projeto.

Para a realização da tradução, foi utilizada uma estrutura simples, onde foi armazenado o símbolo declarado na tabela de símbolos, o valor do seu símbolo, se utiliza algum registrador e o número do seu registrador. Tudo isso armazenado de tal forma a mostrar no final o código gerado.

Para a implementação do tipo set, as variáveis foram armazenadas em um arranjo, onde foi verificado se não acontecia a repetição de seus valores.

7 Descrição dos Arquivos de Teste

Os arquivos de teste buscam apresentar códigos sem erros e códigos que possuem algum erro léxico, sintático ou semântico. Os códigos inválidos buscam mostrar erros observados nas três etapas do projeto, seja com a introdução de um caractere inexistente, uma sequência de comandos que não respeitam a gramática do projeto ou até mesmo um erro semântico no programa fonte, como declaração de mais de uma variável, passagem de parâmetros erradas ou a não inclusão de uma função main.

No código `incorreto1.mf`, o código apresenta erros léxicos, sintáticos e semânticos. Os erros são os seguintes:

- Redecaração da variável `j`;
- Ausência de ponto e vírgula na linha 16, coluna 15;
- Símbolo não aceito na linha 22, coluna 12;

No código `incorreto2.mf`, o código apresenta erros léxicos, sintáticos e semânticos. Os erros são os seguintes:

- Declaração de **EMPTY** logo após um identificador sem tipagem na linha 3, coluna 15;
- Função *forall* sem a operação de conjuntos na linha 5, coluna 16;
- Não declaração da variável *ans*;
- Não declaração da função *copySet*;
- Ausência da função *main*;

8 Instruções para Compilação

O projeto está rodando sobre os seguintes requerimentos:

- gcc (Ubuntu 9.3.0-17ubuntu1 20.04) 9.3.0
- flex 2.6.4
- bison (GNU Bison) 3.7

Para compilar, é necessário executar os seguintes comandos:

Para os códigos corretos.

1. \$ make
2. \$./mafralang tests/corretoX.mf, onde X representa o número do exemplo, para rodar os códigos corretos;
3. tac mafralang.tac

Para os códigos incorretos.

1. \$ make
2. \$./mafralang tests/incorrectoX.mf, onde X representa o número do exemplo, para rodar os códigos incorretos.

Depois de compilar o código, ele gerará no próprio terminal a árvore sintática abstrata e a tabela de símbolos. Os erros apareceram no terminal, informando a linha e a coluna onde ocorreu, além da tabela de símbolos.

Referências

1. Aho, A.V., Lam, M.S., Sethi, R., Ullman, J.D.: Compilers: Principles, Techniques, and Tools (2nd Edition). Addison-Wesley Longman Publishing Co., Inc., USA (2006)
2. Grune, D., Jacobs, C.J.H.: Parsing Techniques: A Practical Guide. Ellis Horwood, USA (1990)
3. Knuth, D.E.: Semantics of context-free languages, vol. 2. Springer (1968)
4. Loudon, K.C.: Compiler construction. Cengage Learning (1997)
5. Nalon, C.: Trabalho prático - descrição da linguagem (2021), <https://aprender3.unb.br/mod/page/view.php?id=294131>, visitou em 26-02-2021
6. Ritchie, D.M., Kernighan, B.W., Lesk, M.E.: The C programming language. Prentice Hall Englewood Cliffs (1988)

A Gramática

A gramática foi feita a partir das adaptações da gramática apresentada por Brian W. Kernighan e Dennis M. Ritchie no livro K&R [6]. Os terminais são apresentados entre aspas simples e em caixa alta, enquanto as variáveis são apresentadas entre parênteses.

$$\begin{aligned}
 \langle program \rangle & ::= \langle translation_unit \rangle \\
 \langle translation_unit \rangle & ::= \langle external_declaration \rangle \\
 & \quad | \langle translation_unit \rangle \langle external_declaration \rangle \\
 \langle external_declaration \rangle & ::= \langle function_definition \rangle \\
 & \quad | \langle declaration \rangle \\
 \langle function_definition \rangle & ::= \langle declaration_specifiers \rangle \langle declarator \rangle \langle declaration_list \rangle \\
 & \quad \langle compound_statement \rangle \\
 & \quad | \langle declaration_specifiers \rangle \langle declarator \rangle \langle compound_statement \rangle \\
 \langle declaration_list \rangle & ::= \langle declaration \rangle \\
 & \quad | \langle declaration_list \rangle \langle declaration \rangle \\
 \langle declaration \rangle & ::= \langle declaration_specifiers \rangle ';' \\
 & \quad | \langle declaration_specifiers \rangle \langle init_declarator_list \rangle ';' \\
 \langle declaration_specifiers \rangle & ::= 'TYPE' \\
 & \quad | 'TYPE' \langle declaration_specifiers \rangle \\
 \langle init_declarator_list \rangle & ::= \langle declarator \rangle \\
 & \quad | \langle init_declarator_list \rangle ',' \langle declarator \rangle \\
 \langle declarator \rangle & ::= \langle direct_declarator \rangle \\
 \langle direct_declarator \rangle & ::= 'ID' \\
 & \quad | 'MAIN' \\
 & \quad | '(' \langle declarator \rangle ')'' \\
 & \quad | \langle direct_declarator \rangle '(' ')'' \\
 & \quad | \langle direct_declarator \rangle '(' \langle identifier_list \rangle ')'' \\
 & \quad | \langle direct_declarator \rangle '(' \langle parameter_list \rangle ')''
 \end{aligned}$$

$$\begin{aligned}
\langle identifier_list \rangle &::= 'ID' \\
&| \langle identifier_list \rangle ', ' ID' \\
\\
\langle parameter_list \rangle &::= \langle parameter_declaration \rangle \\
&| \langle parameter_list \rangle ', ' \langle parameter_declaration \rangle \\
\\
\langle parameter_declaration \rangle &::= \langle declaration_specifiers \rangle \langle declarator \rangle \\
&| \langle declaration_specifiers \rangle \langle abstract_declarator \rangle \\
\\
\langle abstract_declarator \rangle &::= \langle direct_abstract_declarator \rangle \\
\\
\langle direct_abstract_declarator \rangle &::= '(' \langle abstract_declarator \rangle ') ' \\
&| '(', ' ' \\
&| \langle direct_abstract_declarator \rangle '(' ' ' \\
&| \langle direct_abstract_declarator \rangle '(' \langle parameter_list \rangle ') ' \\
\\
\langle compound_statement \rangle &::= '{ ' ' ' \\
&| '{ ' \langle block_item_list \rangle ' ' ' \\
\\
\langle block_item_list \rangle &::= \langle block_item \rangle \\
&| \langle block_item_list \rangle \langle block_item \rangle \\
\\
\langle block_item \rangle &::= \langle declaration \rangle \\
&| \langle statement \rangle \\
\\
\langle statement \rangle &::= \langle expression_statement \rangle \\
&| \langle compound_statement \rangle \\
&| \langle conditional_statement \rangle \\
&| \langle iteration_statement \rangle \\
&| \langle input_statement \rangle \\
&| \langle output_statement \rangle \\
&| \langle return_statement \rangle \\
\\
\langle expression_statement \rangle &::= ';' \\
&| \langle expression \rangle ';' \\
\\
\langle conditional_statement \rangle &::= 'IF' '(' \langle expression \rangle ') ' \langle statement \rangle \\
&| 'IF' '(' \langle expression \rangle ') ' \langle statement \rangle 'else' \langle statement \rangle
\end{aligned}$$

$$\begin{aligned}
\langle \text{iteration_statement} \rangle &::= \text{'FOR' ' ('} \langle \text{expression_statement} \rangle \langle \text{expression_statement} \rangle \\
&\quad \langle \text{expression} \rangle \text{' ' } \langle \text{statement} \rangle \\
&\quad | \text{'FOR' ' ('} \langle \text{expression_statement} \rangle \langle \text{expression_statement} \rangle \text{' ' } \\
&\quad \langle \text{statement} \rangle \\
&\quad | \text{'FOR' ' ('} \langle \text{declaration} \rangle \langle \text{expression_statement} \rangle \langle \text{expression} \rangle \text{' ' } \\
&\quad \langle \text{statement} \rangle \\
&\quad | \text{'FOR' ' ('} \langle \text{declaration} \rangle \langle \text{expression_statement} \rangle \text{' ' } \langle \text{statement} \rangle \\
&\quad | \text{'FORALL' ' ('} \langle \text{expression} \rangle \text{' ' } \langle \text{statement} \rangle
\end{aligned}$$

$$\langle \text{input_statement} \rangle ::= \text{'READ' ' ('} \langle \text{expression} \rangle \text{' '}$$

$$\begin{aligned}
\langle \text{output_statement} \rangle &::= \text{'WRITE' ' ('} \langle \text{expression} \rangle \text{' ' } \\
&\quad | \text{'WRITELN' ' ('} \langle \text{expression} \rangle \text{' ' }
\end{aligned}$$

$$\langle \text{return_statement} \rangle ::= \text{'RETURN' } \langle \text{expression} \rangle \text{' ;'}$$

$$\begin{aligned}
\langle \text{set_expression_list} \rangle &::= \langle \text{is_set_expression} \rangle \\
&\quad | \langle \text{remove_expression} \rangle \\
&\quad | \langle \text{add_expression} \rangle \\
&\quad | \langle \text{exists_expression} \rangle
\end{aligned}$$

$$\langle \text{is_set_statement} \rangle ::= \text{'IS_SET' ' ('} \text{'ID' ' '}$$

$$\langle \text{remove_statement} \rangle ::= \text{'REMOVE' ' ('} \langle \text{expression} \rangle \text{' '}$$

$$\langle \text{add_statement} \rangle ::= \text{'ADD' ' ('} \langle \text{expression} \rangle \text{' '}$$

$$\langle \text{exists_statement} \rangle ::= \text{'EXISTS' ' ('} \langle \text{expression} \rangle \text{' '}$$

$$\begin{aligned}
\langle \text{expression} \rangle &::= \langle \text{expression} \rangle \text{' , ' } \langle \text{assignment_expression} \rangle \\
&\quad | \langle \text{assign_expression} \rangle
\end{aligned}$$

$$\begin{aligned}
\langle \text{assign_expression} \rangle &::= \langle \text{unary_expression} \rangle \text{' = ' } \langle \text{assignment_expression} \rangle \\
&\quad | \langle \text{relational_expression} \rangle
\end{aligned}$$

$$\begin{aligned}
\langle \text{relational_expression} \rangle &::= \langle \text{arithmetic_expression} \rangle \\
&\quad | \langle \text{relational_expression} \rangle \text{' = ' } \langle \text{arithmetic_expression} \rangle \\
&\quad | \langle \text{relational_expression} \rangle \text{' < ' } \langle \text{arithmetic_expression} \rangle \\
&\quad | \langle \text{relational_expression} \rangle \text{' \leq ' } \langle \text{arithmetic_expression} \rangle
\end{aligned}$$

$\langle \text{relational_expression} \rangle \text{'>' } \langle \text{arithmetic_expression} \rangle$
 $\langle \text{relational_expression} \rangle \text{'\geq'} \langle \text{arithmetic_expression} \rangle$
 $\langle \text{relational_expression} \rangle \text{'!=' } \langle \text{arithmetic_expression} \rangle$

$\langle \text{arithmetic_expression} \rangle ::= \langle \text{logical_expression} \rangle$
 $\quad \mid \langle \text{arithmetic_expression} \rangle \text{'+' } \langle \text{logical_expression} \rangle$
 $\quad \mid \langle \text{arithmetic_expression} \rangle \text{'-' } \langle \text{logical_expression} \rangle$
 $\quad \mid \langle \text{arithmetic_expression} \rangle \text{'*'} \langle \text{logical_expression} \rangle$
 $\quad \mid \langle \text{arithmetic_expression} \rangle \text{'/' } \langle \text{logical_expression} \rangle$

$\langle \text{logical_expression} \rangle ::= \langle \text{set_expression} \rangle$
 $\quad \mid \text{'!' } \langle \text{set_expression} \rangle$
 $\quad \mid \langle \text{logical_expression} \rangle \text{'\&\&' } \langle \text{set_expression} \rangle$
 $\quad \mid \langle \text{logical_expression} \rangle \text{'||' } \langle \text{set_expression} \rangle$

$\langle \text{set_expression} \rangle ::= \langle \text{cast_expression} \rangle$
 $\quad \mid \langle \text{set_expression} \rangle \text{'IN' } \langle \text{cast_expression} \rangle$

$\langle \text{cast_expression} \rangle ::= \langle \text{unary_expression} \rangle$
 $\quad \mid \text{'(' } \langle \text{type_name} \rangle \text{'}'$

$\langle \text{type_name} \rangle ::= \langle \text{specifier_qualifier_list} \rangle$
 $\quad \mid \langle \text{specifier_qualifier_list} \rangle \langle \text{abstract_declarator} \rangle$
 $\quad \mid \langle \text{specifier_qualifier_list} \rangle \langle \text{declarator} \rangle$

$\langle \text{specifier_qualifier_list} \rangle :: \text{'TYPE' } \langle \text{specifier_qualifier_list} \rangle$
 $\quad \mid \text{'TYPE'}$

$\langle \text{unary_expression} \rangle ::= \langle \text{postfix_expression} \rangle$
 $\quad \mid \langle \text{set_expression_list} \rangle$
 $\quad \mid \langle \text{function_expression} \rangle$
 $\quad \mid \text{'-' } \langle \text{cast_expression} \rangle$

$\langle \text{function_expression} \rangle ::= \text{'ID' '(' } \langle \text{initializer_list} \rangle \text{'}'$

$\langle \text{postfix_expression} \rangle ::= \langle \text{primary_expression} \rangle$
 $\quad \mid \text{'(' } \langle \text{type_name} \rangle \text{'}' ' \{ ' } \langle \text{initializer_list} \rangle \text{'}'$
 $\quad \mid \text{'(' } \langle \text{type_name} \rangle \text{'}' ' \{ ' } \langle \text{initializer_list} \rangle \text{'}' ' \}'$

$$\begin{aligned} \langle \textit{initializer_list} \rangle &::= \langle \textit{initializer} \rangle \\ &\quad | \quad \langle \textit{initializer_list} \rangle \text{ ',' } \langle \textit{initializer} \rangle \end{aligned}$$

$$\begin{aligned} \langle \textit{initializer} \rangle &::= \langle \textit{assignment_expression} \rangle \\ &\quad | \quad \text{'{' } \langle \textit{initializer_list} \rangle \text{'}} \\ &\quad | \quad \text{'{' } \langle \textit{initializer_list} \rangle \text{' , ' '}} \end{aligned}$$

$$\begin{aligned} \langle \textit{primary_expression} \rangle &::= \text{'ID'} \\ &\quad | \quad \text{'INTEGER'} \\ &\quad | \quad \text{'FLOAT'} \\ &\quad | \quad \text{'STRING'} \\ &\quad | \quad \text{'EMPTY'} \\ &\quad | \quad \text{'(' } \langle \textit{expression} \rangle \text{')'} \end{aligned}$$