

Implementação do Analisador Léxico

Gabriel Fonseca*

Fevereiro 2021

Resumo

Este relatório descreve a implementação de um analisador léxico parte de projeto maior de um tradutor. O tradutor está sendo construído para um subconjunto da linguagem C adicionado de conjuntos bem como suas operações usais. Estão descritos nesse documento a estratégia de implementação desse analisador léxico bem como exemplos para teste.

Palavras-chave: tradutores. analisador léxico. flex.

1 Introdução

No escopo de um tradutor, o analisador léxico é responsável por receber como entrada uma cadeia de caracteres e agrupa-las em uma sequência de lexemas. É função também do analisador léxico interagir com a tabela de símbolos, encontrando lexemas do tipo Indentificador e inserindo-os na tabela. Cada lexema é mapeado para um valor correspondente o qual é utilizado para gerar-se tokens. Sendo esses últimos a saída de um analisador léxico e também a entrada para a próxima etapa de tradução: Análise Sintática.

O analisador descrito neste relatório é limitado à receber a entrada, imprimir em tela os tokens válidos obtidos e listar todas as cadeias de caracteres que não puderam ser agrupadas em lexemas. Além disso, será descrita a gramática contruída para o projeto maior do tradutor bem como exemplos para execução do analisador sintático.

2 A linguagem à ser traduzida

A linguagem à ser traduzida trata-se de um subconjunto da linguagem C adicionado dos tipos *set* e *elem*. Os novos tipos emulam os conjuntos matemáticos e devem implementar operações de pertinência, adição e remoção de elementos, seleção de elementos do conjunto e iteração dos mesmos.

*Departamento de Ciência da Computação, Universidade de Brasília, Brasília, DF, Brasil.
nunesgrf@gmail.com

3 Analisador Léxico

O analisador léxico foi gerado com auxílio da ferramenta Flex. Esta ferramenta gera código-fontes para o reconhecimento de padrões léxicos em uma cadeia de caracteres com base em regras fornecidas ao mesmo. As regras definidas para este analisador podem ser encontradas no arquivo `src/main.lex`.

Descritas todas as regras e expressões regulares, basta delegar ao analisador gerado pelo flex a leitura de caractere à caractere em busca de casa-los no maior padrão definido pelas regras. Isto é feito por meio da função `yylex()`, que deve ser invocada enquanto houver caracteres à serem lidos. No momento que um lexema é encontrado, então este é impresso em tela para ciência do usuário.

Além disso, foi implementado dentre as regras os comentários. Então basta utilizar os lexemas `/*` seguido de `*/` ou até mesmo `//` que tais caracteres do escopo serão ignorados pelo *analisador sintático* assim como são ignorados em tradutores para as mais usuais linguagens de programação da modernidade.

3.1 Tratamento de Erros

Uma vez que o analisador léxico não consegue casar padrão de parte da cadeia de entrada com nenhuma regra, o tratamento de erro é invocado. Uma vez que ainda não existe um erro, um objeto é criado para agrupar todos os caracteres relacionados à aquele erro. Após isso, este objeto é salvo em uma lista de erros.

Ao fim, esta lista de erros é iterada e cada objeto é passa por uma formatação para ser exibido em tela. O usuário então pode comparar a saída de erros com a cadeia de entrada.

3.2 Testes

Para testes, estão disponíveis quatro arquivos-exemplos. Dos quais, dois não devem disparar qualquer tratamento de erro dentro do analisador e dois que devem disparar alguns erros.

O arquivo `samples/correct_sample_1.sample` possui apenas declarações de variáveis e operações simples. As variáveis devem ser exibidas como identificadores e os demais lexemas devem ser exibidos conforme as regras definidas na construção do analisador.

O arquivo `samples/correct_sample_2.sample` possui declarações de variáveis dos tipos *set* e *elem*, parte gramática da linguagem. Bem como operações utilizando essas variáveis.

O arquivo `samples/incorrect_sample_1.sample` possui padrões que não podem ser construídos pela gramática linguagem. Todos os padrões devem ser tratados e exibidos em tela.

O arquivo `sample/incorrect_sample_2.sample` possui padrões sintáticos incorretos que não devem ser tratados como erros pela análise léxica e também possui padrões léxicos incorretos que devem sim ser reportados para o usuário.

4 Conclusão

A construção do analisador léxico é simplificada pela ferramenta Flex, o trabalho fica limitado à definir as regras e dá-las como entrada para a ferramenta. Dado isso, toda a parte de identificação dos lexemas fica abstraída.

Adicionalmente à isso, o tratamento de erros também teve que ser implementado, mas com auxílio dos constructos do Flex, muito foi abstraído. Dessa forma, o resultado final é um analisador léxico simplificado que pode ser usado para gerar a entrada para o próximo passo da tradução.

5 Gramática da linguagem

$\langle \text{program} \rangle$	$::= \langle \text{declaration-list} \rangle$
$\langle \text{declaration-list} \rangle$	$::= \langle \text{declaration-list} \rangle \langle \text{declaration} \rangle \mid \langle \text{declaration} \rangle$
$\langle \text{statement} \rangle$	$::= \langle \text{scope} \rangle$ $\mid \langle \text{var-declaration} \rangle$ $\mid \langle \text{assignment} \rangle$ $\mid \langle \text{print} \rangle$ $\mid \langle \text{scan} \rangle$ $\mid \langle \text{expression} \rangle \text{ ‘;’}$ $\mid \langle \text{condition} \rangle$ $\mid \langle \text{iteration} \rangle$ $\mid \langle \text{return} \rangle$
$\langle \text{declaration} \rangle$	$::= \langle \text{var-declaration} \rangle \mid \langle \text{function-declaration} \rangle$
$\langle \text{var-declaration} \rangle$	$::= \langle \text{type} \rangle \langle \text{identifier} \rangle \text{ ‘;’}$ $\mid \langle \text{type} \rangle \langle \text{identifier} \rangle \text{ ‘[’ } \langle \text{integer} \rangle \text{ ‘]’ ‘;’}$ $\mid \langle \text{type} \rangle \langle \text{identifier} \rangle \text{ ‘{’ ‘}’ ‘;’}$
$\langle \text{function-declaration} \rangle$	$::= \langle \text{type} \rangle \langle \text{identifier} \rangle \text{ ‘(’ } \langle \text{parameters} \rangle \text{ ‘)’ } \langle \text{scope} \rangle$
$\langle \text{parameters} \rangle$	$::= \langle \text{parameters} \rangle \text{ ‘,’ } \langle \text{parameter} \rangle \mid \langle \text{parameter} \rangle \mid \varepsilon$
$\langle \text{parameter} \rangle$	$::= \langle \text{type} \rangle \langle \text{identifier} \rangle$
$\langle \text{scope} \rangle$	$::= \text{ ‘{’ } } \langle \text{statement-list} \rangle \text{ ‘}’$
$\langle \text{statement-list} \rangle$	$::= \langle \text{statement-list} \rangle \langle \text{statement} \rangle \mid \varepsilon$
$\langle \text{print} \rangle$	$::= \langle \text{print-word} \rangle \langle \text{expression} \rangle \text{ ‘;’}$
$\langle \text{scan} \rangle$	$::= \langle \text{scan-word} \rangle \langle \text{identifier} \rangle \text{ ‘;’}$
$\langle \text{print-word} \rangle$	$::= \text{ ‘print’}$
$\langle \text{scan-word} \rangle$	$::= \text{ ‘scan’}$
$\langle \text{condition} \rangle$	$::= \text{ ‘if’ ‘(’ } \langle \text{expression} \rangle \text{ ‘)’ } \langle \text{statement} \rangle \langle \text{condition-mid} \rangle$
$\langle \text{condition-mid} \rangle$	$::= \text{ ‘else if’ ‘(’ } \langle \text{expression} \rangle \text{ ‘)’ } \langle \text{statement} \rangle \langle \text{condition-mid} \rangle$ $\mid \langle \text{condition-end} \rangle$
$\langle \text{condition-end} \rangle$	$::= \text{ ‘else’ } \langle \text{statement} \rangle \mid \varepsilon$
$\langle \text{iteration} \rangle$	$::= \text{ ‘while’ ‘(’ } \langle \text{expression} \rangle \text{ ‘)’ } \langle \text{statement} \rangle$ $\mid \text{ ‘for’ ‘(’ } \langle \text{expression} \rangle \text{ ‘?’ ‘;’ } \langle \text{expression} \rangle \text{ ‘?’ ‘;’ } \langle \text{expression} \rangle \text{ ‘)’ } \langle \text{statement} \rangle$
$\langle \text{return} \rangle$	$::= \text{ ‘return’ } \langle \text{expression} \rangle \text{ ‘?’ ‘;’}$

$\langle \text{assignment} \rangle$	$::= \langle \text{identifier} \rangle \langle \text{assignment-op} \rangle \langle \text{expression} \rangle \text{' ; '}$
$\langle \text{expression} \rangle$	$::= \langle \text{and-expression} \rangle$
$\langle \text{and-expression} \rangle$	$::= \langle \text{or-expression} \rangle$ $\langle \text{and-expression} \rangle \text{' \&\& ' } \langle \text{or-expression} \rangle$
$\langle \text{or-expression} \rangle$	$::= \langle \text{bw-or-expression} \rangle$ $\langle \text{or-expression} \rangle \text{' '}$
$\langle \text{eq-expression} \rangle$	$::= \langle \text{relational-expression} \rangle$ $\langle \text{eq-expression} \rangle \text{' == ' } \langle \text{rel-expression} \rangle$ $\langle \text{eq-expression} \rangle \text{' != ' } \langle \text{rel-expression} \rangle$
$\langle \text{rel-expression} \rangle$	$::= \langle \text{rel-expression} \rangle \langle \text{rel-op} \rangle$
$\langle \text{add-expression} \rangle$	$::= \langle \text{mult-expression} \rangle$ $\langle \text{add-expression} \rangle \text{' + ' } \langle \text{mult-expression} \rangle$ $\langle \text{add-expression} \rangle \text{' - ' } \langle \text{mult-expression} \rangle$
$\langle \text{mult-expression} \rangle$	$::= \langle \text{cast-expression} \rangle$ $\langle \text{mult-expression} \rangle \langle \text{mul-op} \rangle \langle \text{cast-expression} \rangle$
$\langle \text{cast-expression} \rangle$	$::= \langle \text{unary-expression} \rangle$ $\text{' (' } \langle \text{type} \rangle \text{') ' } \langle \text{cast-expression} \rangle$
$\langle \text{unary-expression} \rangle$	$::= \langle \text{postfix-expression} \rangle$ $\langle \text{unary-op} \rangle \langle \text{cast-expression} \rangle$ $\text{sizeof } \langle \text{cast-expression} \rangle$
$\langle \text{postfix-expression} \rangle$	$::= \langle \text{primary-expression} \rangle$ $\langle \text{postfix-expression} \rangle [\langle \text{expression} \rangle]$ $\langle \text{postfix-expression} \rangle (\langle \text{param-values} \rangle)$
$\langle \text{param-values} \rangle$	$::= \langle \text{param-values} \rangle \text{' , ' } \langle \text{expression} \rangle \mid \langle \text{expression} \rangle \mid \varepsilon$
$\langle \text{primary-expression} \rangle$	$::= \langle \text{identifier} \rangle$ $\langle \text{constant} \rangle$ $\langle \text{string} \rangle$ $\text{' (' } \langle \text{expression} \rangle \text{') '}$
$\langle \text{constant} \rangle$	$::= \langle \text{integer} \rangle$ $\text{' ' } \langle \text{symbol} \rangle \text{' '}$ $\langle \text{integer} \rangle \text{' . ' } \langle \text{integer} \rangle$
$\langle \text{integer} \rangle$	$::= \langle \text{digit} \rangle +$
$\langle \text{identifier} \rangle$	$::= \langle \text{letter} \rangle \{ \langle \text{letter} \rangle \mid \langle \text{digit} \rangle \mid \text{' _ ' } \}^*$
$\langle \text{type} \rangle$	$::= \text{' int ' } \mid \text{' float ' } \mid \text{' elem ' } \mid \text{' set '}$
$\langle \text{symbol} \rangle$	$::= \text{Printables ASCII chars.}$
$\langle \text{letter} \rangle$	$::= \text{' a ' } \mid \text{' b ' } \mid \dots \mid \text{' z ' } \mid \text{' A ' } \mid \dots \mid \text{' Z '}$
$\langle \text{digit} \rangle$	$::= \text{' 0 ' } \mid \text{' 1 ' } \mid \text{' 2 ' } \mid \dots \mid \text{' 9 '}$
$\langle \text{assignment-op} \rangle$	$::= \text{' = ' } \mid \text{' += ' } \mid \text{' -= ' } \mid \text{' *= ' } \mid \text{' /= '}$
$\langle \text{unary-op} \rangle$	$::= \text{' + ' } \mid \text{' - ' } \mid \text{' ! '}$
$\langle \text{rel-op} \rangle$	$::= \text{' < ' } \mid \text{' > ' } \mid \text{' <= ' } \mid \text{' >= ' } \mid \text{' in '}$
$\langle \text{mul-op} \rangle$	$::= \text{' * ' } \mid \text{' / '}$