

Analizador léxico

Vinícius T M Sugimoto

CiC-UnB

Resumo Este documento apresenta a construção de um analisador léxico para um subconjunto da linguagem C feito com a ferramenta **flex**.

Keywords: Analisador sintático · Flex · C · Tradutores

1 Motivação

Este projeto foi desenvolvido com o objetivo de construir uma linguagem derivada da linguagem **C** para trabalhar com conjuntos. Assim, os tipos primitivos `set` e `elem` foram adicionados à linguagem.

2 Análise Léxica

A análise léxica é a primeira etapa para a construção de um tradutor para a linguagem proposta. Com a análise léxica, um programa escrito na linguagem proposta é convertido em uma sequência de *tokens* que podem ser usados nas próximas etapas da construção do tradutor, as análises sintática e semântica. Esta análise foi feita usando a ferramenta **flex**, que define um conjunto de “regras” que casam expressões regulares com o programa de entrada e executa comandos específicos para cada regra.

O **flex** foi utilizado para gerar *tokens* de maneira apropriada para cada regra e tentar identificar erros no código na etapa de análise léxica.

3 Testes

Junto aos arquivos deste projeto, que podem ser obtidos seguindo os passos da seção 4 estão quatro arquivos nomeados `subset_sumX.conj`, em que $X = 1, 2, 3, 4$ que server como arquivos de entrada teste para o analisador léxico gerado. Os arquivos com final 1, 2 são arquivos corretos, já os arquivos com final 3, 4 representam arquivos com erros.

4 Reprodução

4.1 Dependências

Este projeto tem as seguintes dependências:

- **gcc** versão 10.2.0
- **flex** versão 2.6.4

4.2 Execução

Os arquivos deste projeto podem ser encontrados no endereço <https://github.com/vinicius-toshiyuki/conjunto.git>. Os arquivos podem ser baixados em formato .zip pelo site ou o repositório pode ser clonado com o programa `git` no terminal. A seguir serão mostrados os passos a serem seguidos para executar o analisador léxico gerado em um terminal Linux.

```
# Clona o repositório
git clone \
https://github.com/vinicius-toshiyuki/conjunto.git
# Muda para o diretório criado
cd conjunto
# Gera o código C do analisador léxico
flex sintaxe.l
# Compila o analisador léxico
gcc lex.yy.c -lfl
# Para executar use o comando cat
cat subset_sum.conj | ./a.out
```

Referências

- UoCa. Berkeley University of California. flex.1. Disponível em http://web.mit.edu/gnu/doc/html/flex_1.html. Acessado em 2021.
- UoCb. Berkeley University of California. Lexical analysis with flex, for flex 2.6.3. Disponível em https://www.cs.virginia.edu/~cr4bd/flex-manual/index.html#SEC_Contents. Acessado em 2021.

5 Anexo

A sintaxe da linguagem é composta de expressões e comandos. A seguir são apresentados os padrões reconhecidos pela gramática desta linguagem. Nos padrões a seguir os caracteres [] são usados para impôr uma precedência maior na expressão regular interna, [[]] são usados para indicar grupos ou intervalos de caracteres em que qualquer um dos caracteres pertencentes pode ser casado, * é usado como a “*Kleene star*”, + é usado para indicar que a expressão regular precedente é obrigatória e pode ser repetida, ? é usado para indicar que a expressão regular precedida é opcional, < > são usados para indicar um grupo de padrões. Vale ressaltar que espaços em branco não estão inclusos nos padrões descritos a seguir, para fins de simplicidade, mas considere-os válidos onde seriam válidos na linguagem C.

5.1 Padrões auxiliares

Os padrões a seguir são utilizados nos outros grupos de padrões e cada um tem sua própria denominação.

1. $-^?[[0 - 9]]^+$
Um número inteiro. É denominado pelo grupo **<INT>**.
2. $-^?[[0 - 9]]^+.[[0 - 9]]^*f^?$ ou $-^?[[0 - 9]]^*.[[0 - 9]]^+f^?$ ou **<INTEGER>**f
Um número real. É denominado pelo grupo **<FLOAT>**.
3. **<INT>** ou **<FLOAT>** ou **EMPTY**
Uma constante. É denominada pelo grupo **<CONST>**.
4. $[[a - zA - Z_]] [[a - zA - Z0 - 9_]]^*$
Um identificador. É denominado pelo grupo **<ID>**.
5. **int** ou **float** ou **set** ou **elem**
Um tipo de dados. É denominado pelo grupo **<TYPE>**.
6. * ou / ou % ou + ou -
Um operador aritmético. É denominado pelo grupo **<OPARIT>**.
7. < ou > ou <= ou >= ou == ou !=
Um operador relacional. É denominado pelo grupo **<OPREL>**.
8. && ou || ou & ou |
Um operador booleano. É denominado pelo grupo **<OPBOOL>**.
9. **<OPARIT>** = ou & = ou | =
Um operador composto. É denominado pelo grupo **<OPCOMP>**.
10. **<OPARIT>** ou **<OPCOMP>** ou **<OPREL>** ou **<OPBOOL>** ou = ou in
Um operador binário. É denominado pelo grupo **<OPBIN>**.
11. ! ou -
Um operador unário. É denominado pelo grupo **<OPUNI>**.
12. ++ ou --
Um operador unário especial. É denominado pelo grupo **<OPSPE>**.

5.2 Expressões

Os padrões a seguir fazem parte do grupo $\langle \mathbf{EXP} \rangle$.

1. $\langle \mathbf{CONST} \rangle$
Uma expressão constante. Seu valor é a própria constante.
2. $\langle \mathbf{ID} \rangle$
Uma expressão identificador. Seu valor é o valor associado ao identificador.
3. $\langle \mathbf{EXP} \rangle \langle \mathbf{OPBIN} \rangle \langle \mathbf{EXP} \rangle$
Uma expressão com operação binária. Seu valor é o valor da operação sobre os valores das duas expressões internas.
4. $\langle \mathbf{OPUNI} \rangle \langle \mathbf{EXP} \rangle$ ou $\langle \mathbf{OPSPE} \rangle \langle \mathbf{EXP} \rangle$
Uma expressão com operação unária. Seu valor é o valor da operação sobre o valor da expressão interna.
5. $\langle \mathbf{EXP} \rangle \langle \mathbf{OPSPE} \rangle$
Uma expressão com operação unária especial. Seu valor é o valor da expressão interna.
6. $\langle \mathbf{ID} \rangle ([\langle \mathbf{EXP} \rangle [, \langle \mathbf{EXP} \rangle]^*]^?)$
Uma expressão chamada de função. Seu valor é o valor de retorno da função computada sobre seus argumentos.

5.3 Comandos

Os padrões a seguir fazem parte do grupo $\langle \mathbf{CMD} \rangle$.

1. $\{ \langle \mathbf{CMD} \rangle^* \}$
Um bloco é um comando. Um bloco é seguido de uma sequência opcional de comandos.
2. $\langle \mathbf{EXP} \rangle^?;$
Um comando vazio. Um comando vazio pode ser precedido de uma expressão.
3. $\mathbf{if}(\langle \mathbf{EXP} \rangle) \langle \mathbf{CMD} \rangle$
Um comando **if**. Um comando **if** é seguido de uma expressão em parênteses e de um comando.
4. $\mathbf{else} \langle \mathbf{CMD} \rangle$
Um comando **else**. Um comando **else** é seguido de um comando. Na etapa de análise léxica não há checagem se um comando **else** é precedido de um **if**.
5. $\mathbf{while}(\langle \mathbf{EXP} \rangle) \langle \mathbf{CMD} \rangle$
Um comando **while**. Um comando **while** é seguido de uma expressão em parênteses e de um comando.
6. $\mathbf{forall}(\langle \mathbf{EXP} \rangle) \langle \mathbf{CMD} \rangle$
Um comando **forall**. Um comando **forall** é seguido de uma expressão em parênteses e de um comando. Na etapa de análise léxica não há checagem se a expressão é uma expressão “**in**”.
7. $\mathbf{for}(\langle \mathbf{EXP} \rangle^?; \langle \mathbf{EXP} \rangle^?; \langle \mathbf{EXP} \rangle^?) \langle \mathbf{CMD} \rangle$
Um comando **for**. Um comando **for** é seguido de três expressões opcionais separadas por ponto-e-vírgula em parênteses e de um comando.

8. $\langle \mathbf{TYPE} \rangle \langle \mathbf{ID} \rangle [, \langle \mathbf{ID} \rangle]^*$;
Uma comando de declaração de variável. Uma declaração é composta de um tipo seguido de um identificador seguido de uma lista de “vírgula e identificador” opcional e de um ponto-e-vírgula.
9. $\langle \mathbf{TYPE} \rangle \langle \mathbf{ID} \rangle ([\langle \mathbf{TYPE} \rangle \langle \mathbf{ID} \rangle [, \langle \mathbf{TYPE} \rangle \langle \mathbf{ID} \rangle]^*]^?) \{ \langle \mathbf{CMD} \rangle \}$
Uma comando de declaração de função. Uma declaração é composta de um tipo seguido de um identificador seguido de uma lista opcional de parâmetros no formato “tipo e identificador” em parênteses e de um comando bloco.
10. **return** $\langle \mathbf{EXP} \rangle$
Um comando de retorno. É composto da palavra reservada **return** seguido de uma expressão e um ponto-e-vírgula.