



Universidade Federal de Ouro Preto - UFOP
Instituto de Ciências Exatas e Biológicas - ICEB
Departamento de Computação - DECOM
Disciplina: COMPILADORES – BCC 328
Aluno: Thiago Oliveira de Santana
Matrícula: 15.1.4313



TRABALHO PRÁTICO I – ANALISADOR LÉXICO

Sabemos que o processo da análise léxica consiste em ler os caracteres de entrada e agrupá-los em conjuntos que são chamados de tokens, sendo estes, os identificadores e símbolos especiais de uma linguagem de computação. O compilador tem no seu analisador léxico o primeiro componente a entrar em contato com o código fonte, onde o seu papel basicamente é reconhecer os tokens existentes no código e passar essa informação para o parser.

Tendo isso em mente, foi implementado um analisador léxico para a linguagem de programação apresentada por Torben em seu livro de compiladores, no capítulo 4 (gramática 4.1). Para a implementação utilizou-se um gerador de analisador léxico escrito em Java, conhecido como **JFlex**.

1.0 - Linguagem

As palavras pertencentes à linguagem descritas pela gramática são: literais, identificadores, sinais de pontuação, operadores e palavras-chaves. A linguagem possui também características como comentários em linha e bloco (podendo ser aninhado) bem como delimitadores de espaço.

1.1 – Literais

É sabido que os literais podem ser classificados como sequência de caracteres que representa uma constante, como um número inteiro, um número em ponto flutuante, um caractere, uma string, um valor verdade (verdadeiro ou falso), etc. Desta forma, os literais da linguagem em questão pertencem ao tipo inteiro, podendo ser classificadas pelo analisador léxico como um token **LITERAL_INT**. Em seguida, contemplaremos alguns exemplos desse grupo.

Exemplos: 1 - “26342 ”;
2 - “0”;
3 - “100000”.

1.2 – Identificadores

Podemos compreender que os identificadores são palavras utilizadas para nomear entidades do programa, como funções, métodos, classes, módulos, etc. Desta maneira, o analisador léxico implementado classifica todas as palavras que pertencem ao conjunto **[a-zA-Z][a-zA-Z0-9_]*** como um token **ID**. É possível notar que as

palavras pertencentes a esse conjunto são todas aquelas que iniciam com qualquer letra do alfabeto, podendo ser maiúsculas ou minúsculas, seguidas de qualquer letra ou número. A seguir, contemplaremos alguns exemplos dessa categoria.

Exemplos: 1 - “Nome”;
2 - “João09”;
3 - “T_Oliveira”.

1.3 – Sinais de pontuação

É notório que os sinais de pontuação são sequências de caracteres que auxiliam na construção das estruturas do programa, como por exemplo , servindo de separador de expressões em uma lista de expressões. Sendo assim, os sinais de pontuação pertencentes a linguagem são “ (” , “) ” e “ , ”, na qual são classificados pelo analisador como sendo tokens **T_PARENTESEL**, **T_PARENTESER**, **T_VIRGULA**. Em seguida, veremos alguns exemplos dessa classe.

Exemplos: 1 - “ if(num == 0) ”;
2 - “ num1, num2”;
3 - “ (x, y, z)”.

1.4 – Operadores

Os operadores que fazem parte da linguagem são da forma “+” e “=” e que ao serem analisados, são classificados como sendo tokens **T MAIS** e **T_IGUAL**, respectivamente. Em seguida, podemos ver alguns exemplos desse grupo.

Exemplos: 1 - “ num1 = num2 ”;
2 - “ soma = num1+num2”;
3 - “ +=”.

1.5 – Palavras-chaves

As palavras-chaves são usadas para expressar estruturas da linguagem, como comandos condicionais, comandos de repetição, etc. Geralmente são reservadas, não podendo ser utilizadas como identificadores. Portanto, as palavras reservadas desta linguagem são: “**bool**”, “**int**”, “**if**”, “**then**”, “**else**”, “**let**” e “**in**”. Ao serem examinadas pelo analisador léxico, respectivamente, serão gerados tokens **BOOL**, **INT** , **IF**, **THEN**, **ELSE**, **LET** e **IN**. Em seguida, contemplaremos alguns exemplos dessa categoria.

Exemplos: 1 - “ int numero ”;
2 - “ bool verificador”;
3 - “let id = Exp in Exp”; .
4 - “if Exp then Exp else Exp”;

1.6 – Comentários e delimitadores de espaço

Os comentários de linha, bloco e delimitadores de espaço ao serem examinados pelo analisador léxico são desconsiderados pelo fato destes, assim como nas linguagens de programação conhecidas, não serem importantes para o processo de decodificação realizada pelo compilador. O comentário de linha pode ser realizado utilizando % e o comentário de bloco utilizando {% %}. Já os delimitadores de espaço podem ser efetuados com \t\r\n. A seguir, veremos exemplos dessa categoria.

Exemplos: 1 - “ Print\n”;

2 - “ % Comentário em linha “;

3 - “ {% Comentário em bloco %}”.

2.0 – Testes

Para cada categoria descrita anteriormente, foram realizados testes automatizados a fim de averiguar a efetividade do analisador léxico. A seguir, contemplaremos alguns exemplos de testes.

// Comentário de Linha

```
trun("% a line comment\n", "2:1-2:1 EOF");
```

// Comentário de Bloco

```
trun("{% a block comment %}", "1:22-1:22 EOF");
```

```
trun("{% outer {% inner %} outer %}", "1:30-1:30 EOF");
```

// Pontuação

```
trun("=", "1:1-1:2 T_IGUAL", "1:2-1:2 EOF");
```

```
trun(",", "1:1-1:2 T_VIRGULA", "1:2-1:2 EOF");
```

// Operador

```
trun("+", "1:1-1:2 T_MAIS", "1:2-1:2 EOF");
```

// Literais inteiros

```
trun("26342", "1:1-1:6 LITERAL_INT(26342)", "1:6-1:6 EOF");
```

```
trun("0", "1:1-1:2 LITERAL_INT(0)", "1:2-1:2 EOF");
```

// Palavras Reservadas

```
trun("bool", "1:1-1:5 BOOL", "1:5-1:5 EOF");
```

```
trun("int", "1:1-1:4 INT", "1:4-1:4 EOF");
```

```
trun("if", "1:1-1:3 IF", "1:3-1:3 EOF");
```

// Identificadores

```
trun("nome", "1:1-1:5 ID(nome)", "1:5-1:5 EOF");
```

```
trun("camelCase", "1:1-1:10 ID(camelCase)", "1:10-1:10 EOF");
```