

[Trabalho] 1^a Fase - Analisador Léxico

Equipe: Charles Lelis Braga (202035015), Gabriella Carvalho (202165047AC)

2024

Resumo

Este trabalho descreve a implementação de um analisador léxico Java para a linguagem Lang utilizando Jflex. Neste projeto será possível captura e printar em tela todos os tokens obtidos, além de processar seus tipos e valores.

1 Estratégias de Implementação

1.1 Modelagem do Trabalho

O projeto do analisador léxico é estruturado em três diretórios: Source, Lib e Samples.

A Lib armazena todos os pacotes necessários, nesse caso, apenas o JFlex, que é usado para gerar o analisador léxico (`Lexer.java`) a partir da especificação contida no arquivo `java.flex`.

A pasta Source contém os arquivos e o código-fonte do projeto, onde a Main é o ponto de entrada, coordenando as chamadas para o Lexer e outros componentes. O subcomponente Lexer contém o arquivo `java.flex`, que é um arquivo de definição para o JFlex. A partir deste arquivo, são gerados dois elementos principais: o `TOKEN_TYPE`, que define os tipos de tokens que o lexer irá identificar, e o `Lexer.java`, que é o código-fonte do lexer gerado automaticamente e responsável por reconhecer e categorizar os tokens do código de entrada.

O Utils inclui o `SampleFileManager`, que é um gerenciador de arquivos responsável por lidar com as operações de leitura e escrita necessárias para o processamento dos arquivos de entrada e saída no projeto.

Já o diretório Samples armazena exemplos utilizados para validação sintática e semântica. Dentro deste componente, o submódulo Sintatic contém amostras para validação sintática, com resultados esperados categorizados como `True` ou `False`, dependendo se o exemplo atende ou não às regras sintáticas definidas. Da mesma forma, o submódulo Semantic armazena amostras para validação semântica, também com resultados `True` ou `False`, verificando se as amostras são semanticamente corretas de acordo com as regras especificadas. Obviamente, nesse ponto do trabalho, só valida-se a captura de tokens.

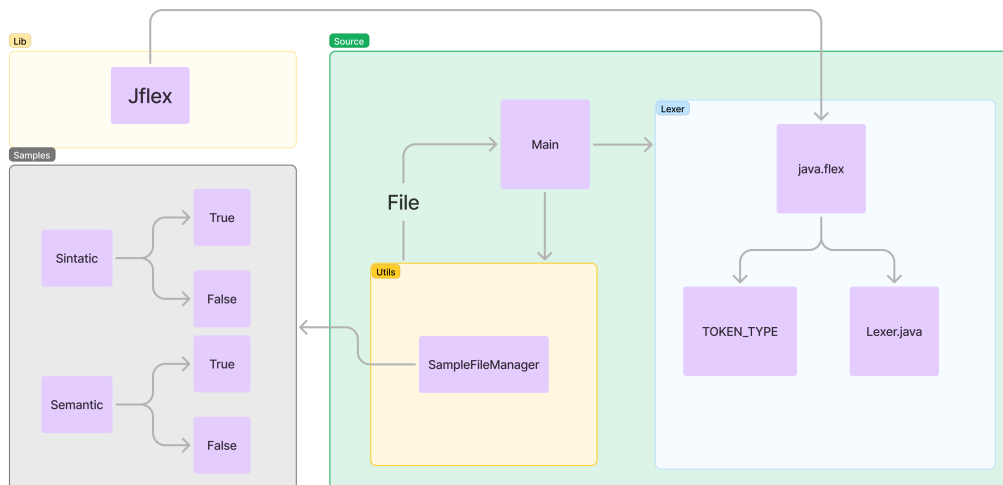


Figura 1 – Arquitetura do projeto

1.2 Organização de Pastas

Desta forma, segue os diretórios ficaram da seguinte forma:

```
1 samples/
2     semantic/
3         true/
4         false/
5     sintatic/
6         true/
7         false/
8 lib/
9     jflex-full-1.8.2.jar
10 src/
11     core/
12         Lexer/
13             lang.flex
14             Lexer.java
15             LexerProcessor.java
16             TOKEN_TYPE.java
17             Token.java
18     utils/
19         SampleFileManager
20         App.java
21 build.sh
```

1.3 Inicialização do processamento

A função `main` da classe `App` é o ponto de entrada do programa. Ela utiliza o método `getSampleFile` da classe `SampleFileManager` para obter um arquivo de entrada, que pode ser escolhido pelo usuário ou pré-definido. Esse arquivo é então passado para o método `process` da classe `Lexer`, que realiza a análise léxica, identificando e categorizando os tokens do conteúdo do arquivo.

1.4 Tokens

O uso de tokens no analisador léxico é uma etapa fundamental no processo de compilação e interpretação de linguagens de programação.

No código, o `Token` é definido pela classe `Token`, que possui um tipo (`t`) e o objeto (`Object info`). A classificação de tipos possui três categorias que facilitam o entendimento deles:

Os tokens de palavras reservadas (`reservedWordsTypes`) incluem elementos como `if`, `else`, `while` e outras palavras-chave que controlam o fluxo do programa. Esses tokens são retornados diretamente como strings representando suas palavras originais.

Os tipos de valores (`valuableTypes`) abrangem tokens que representam dados, como `int`, `float`, `char`, e identificadores (`ID`). Para esses tokens, o método retorna o tipo do token junto com a informação adicional associada (`info`), que pode incluir o valor ou nome específico do token.

Por fim, os operadores (`operatorsTypes`) incluem símbolos como `+`, `-`, `*`, `==`, `&&`, `||`, parênteses e outros operadores da linguagem. Para esses tokens, o método retorna apenas a `info`, que descreve o operador específico encontrado no código. Essa estrutura facilita a compreensão de como cada token é utilizado na análise léxica, categorizando-os com base no papel que desempenham no código.

1.5 Lexer

O analisador é configurado com várias opções do JFlex que controlam seu comportamento:

- `%unicode`: Garante o suporte a caracteres Unicode.
- `%line` e `%column`: Incluem informações de linha e coluna para melhor rastreamento de erros.
- `%public`: Define a classe gerada como pública.
- `%class LexerProcessor`: Nomeia a classe do analisador léxico.
- `%function nextToken`: Define a função que retorna o próximo token.
- `%type Token`: Define o tipo dos tokens gerados.

1.6 Inicialização e Utilitários

O bloco de inicialização e os métodos utilitários fornecem suporte para a contagem de tokens e a criação de instâncias de tokens:

- **Variável `ntk`**: Conta o número de tokens processados.
- **Método `readedTokens`**: Retorna a quantidade total de tokens lidos.
- **Método `symbol`**: Cria um token com um tipo específico e um valor associado. Existem duas sobrecargas deste método, uma para tokens simples e outra para tokens com valores adicionais.

```
1 %{\n2     private int ntk;\n3\n4     public int readedTokens(){\n5         return ntk;\n6     }\n7\n8     private Token symbol(TOKEN_TYPE t) {\n9         ntk++;\n10        return new Token(t, yytext());\n11    }\n12\n13    private Token symbol(TOKEN_TYPE t, Object value) {\n14        ntk++;
```

```

15     return new Token(t, value);
16 }
17 %}
18
19 %init{
20     ntk = 0;
21 %init}

```

Definições de Tokens

As definições de tokens são fundamentais para reconhecer e classificar diferentes elementos no código-fonte:

- **Alfabeto e Identificadores:**
 - ALPHA, ALPHA_UPPERCASE, ALPHA_LOWERCASE definem conjuntos de caracteres alfabéticos.
 - IDENT_UPPERCASE e IDENT_LOWERCASE definem padrões para identificadores, considerando letras, dígitos e sublinhados.
- **Números e Tipos de Dados:**
 - INT e FLOAT definem padrões para números inteiros e de ponto flutuante.
- **Caracteres Especiais:**
 - Definem os caracteres que são utilizados para formatação e controle, como novas linhas, espaços em branco e caracteres de escape.

```

1 ALPHA = [A-Za-z]
2 ALPHA_UPPERCASE = [A-Z]
3 ALPHA_LOWERCASE = [a-z]
4
5 INT = [:digit:] [:digit:]*
6 FLOAT = [-+]?([:digit:]+ \. [:digit:]* )
7 IDENT_UPPERCASE = {ALPHA_UPPERCASE}({ALPHA}|:digit:|_)*
8 IDENT_LOWERCASE = {ALPHA_LOWERCASE}({ALPHA}|:digit:|_)*
9
10 NEW_LINE = \r|\n|\r\n
11 WHITE_SPACE_CHAR = [\n\r\ \t\b]
12
13 CHAR_NEWLINE = \\n
14 CHAR_TAB = \\t
15 CHAR_BACKSPACE = \\b
16 CHAR_CARRIAGE = \\r
17 CHAR_BACKSLASH = \\ \\
18 CHAR_QUOTE = \\ \"

```

1.7 Regras de Token

As regras de token são agrupadas em diferentes estados, cada um lidando com aspectos específicos da linguagem:

Estado Inicial (YYINITIAL)

O estado inicial é responsável por reconhecer e classificar palavras reservadas, identificadores, tipos de dados e operadores. Ele também lida com comentários e caracteres especiais, como parênteses e operadores.

```
1 <YYINITIAL>{
2     /* RESERVED_WORDS */
3     "if" { return symbol(TOKEN_TYPE.IF); }
4     "else" { return symbol(TOKEN_TYPE.ELSE); }
5     "while" { return symbol(TOKEN_TYPE.WHILE); }
6     "for" { return symbol(TOKEN_TYPE.FOR); }
7     "return" { return symbol(TOKEN_TYPE.RETURN); }
8     "break" { return symbol(TOKEN_TYPE.BREAK); }
9     "continue" { return symbol(TOKEN_TYPE.CONTINUE); }
10    "new" { return symbol(TOKEN_TYPE.NEW); }
11    "void" { return symbol(TOKEN_TYPE.VOID); }
12    "struct" { return symbol(TOKEN_TYPE.STRUCT); }
13    "typedef" { return symbol(TOKEN_TYPE.TYPDEF); }
14    "switch" { return symbol(TOKEN_TYPE.SWITCH); }
15    "case" { return symbol(TOKEN_TYPE.CASE); }
16    "default" { return symbol(TOKEN_TYPE.DEFAULT); }
17    "null" { return symbol(TOKEN_TYPE.NULL); }
18    "true" { return symbol(TOKEN_TYPE.BOOL); }
19    "false" { return symbol(TOKEN_TYPE.BOOL); }
20    "print" { return symbol(TOKEN_TYPE.PRINT); }
21    "scan" { return symbol(TOKEN_TYPE.SCAN); }
22
23    "Int" { return symbol(TOKEN_TYPE.BTYPE); }
24    "Float" { return symbol(TOKEN_TYPE.BTYPE); }
25    "Char" { return symbol(TOKEN_TYPE.BTYPE); }
26    "Bool" { return symbol(TOKEN_TYPE.BTYPE); }
27
28    {IDENT_LOWERCASE} { return symbol(TOKEN_TYPE.ID, yytext()); }
29    {IDENT_UPPERCASE} { return symbol(TOKEN_TYPE.ID, yytext()); }
30    {INT} { return symbol(TOKEN_TYPE.INT, yytext()); }
31    {FLOAT} { return symbol(TOKEN_TYPE.FLOAT, yytext()); }
32
33    "--" { yybegin(LINE_COMMENT); }
34    "{-" { yybegin(COMMENT); }
35
36    "=" { return symbol(TOKEN_TYPE.ASSIGNMENT); }
37    "==" { return symbol(TOKEN_TYPE.EQ); }
38    "!=" { return symbol(TOKEN_TYPE.NOT_EQ); }
39    ";" { return symbol(TOKEN_TYPE.SEMI); }
40    "*" { return symbol(TOKEN_TYPE.TIMES); }
```

```

41  "+" { return symbol(TOKEN_TYPE.PLUS); }
42  "%" { return symbol(TOKEN_TYPE.MOD); }
43  "," { return symbol(TOKEN_TYPE.COMMA); }
44  "::" { return symbol(TOKEN_TYPE.DOUBLE_COLON); }
45  ":" { return symbol(TOKEN_TYPE.COLON); }
46
47  "'" { yybegin(CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.SINGLE_QUOTE); }
48
49  "(" { return symbol(TOKEN_TYPE.LEFT_PAREN); }
50  ")" { return symbol(TOKEN_TYPE.RIGHT_PAREN); }
51  "[" { return symbol(TOKEN_TYPE.LEFT_BRACKET); }
52  "]" { return symbol(TOKEN_TYPE.RIGHT_BRACKET); }
53  "{" { return symbol(TOKEN_TYPE.LEFT_BRACE); }
54  "}" { return symbol(TOKEN_TYPE.RIGHT_BRACE); }
55  "." { return symbol(TOKEN_TYPE.DOT); }
56  "-" { return symbol(TOKEN_TYPE.MINUS); }
57  "/" { return symbol(TOKEN_TYPE.DIVIDE); }
58  "<>" { return symbol(TOKEN_TYPE.NOT_EQUAL); }
59  "<=" { return symbol(TOKEN_TYPE.LESS_THAN_OR_EQUAL); }
60  "<" { return symbol(TOKEN_TYPE.LESS_THAN); }
61  ">=" { return symbol(TOKEN_TYPE.GREATER_THAN_OR_EQUAL); }
62  ">" { return symbol(TOKEN_TYPE.GREATER_THAN); }
63  "&&" { return symbol(TOKEN_TYPE.DOUBLE_AMPERSAND); }
64  "&" { return symbol(TOKEN_TYPE.AMPERSAND); }
65  "!" { return symbol(TOKEN_TYPE.EXCLAMATION_MARK); }
66  "||" { return symbol(TOKEN_TYPE.DOUBLE_PIPE); }
67  "|" { return symbol(TOKEN_TYPE.PIPE); }
68
69  {WHITE_SPACE_CHAR} { }
70 }

```

Estado de Citação de Caractere Único (CHAR_SINGLE_QUOTE)

Este estado lida com a leitura de caracteres dentro de aspas simples. Se um caractere de nova linha ou um caractere especial é encontrado, o analisador muda para o estado END_CHAR_SINGLE_QUOTE.

```

1  <CHAR_SINGLE_QUOTE> {
2      {CHAR_NEWLINE} { yybegin(END_CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.CHAR); }
3      {CHAR_TAB} { yybegin(END_CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.CHAR); }
4      {CHAR_BACKSPACE} { yybegin(END_CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.CHAR); }
5      {CHAR_CARRIAGE} { yybegin(END_CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.CHAR); }
6      {CHAR_BACKSLASH} { yybegin(END_CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.CHAR); }
7      {CHAR_QUOTE} { yybegin(END_CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.CHAR); }
8      {ALPHA} { yybegin(END_CHAR_SINGLE_QUOTE); return symbol(TOKEN_TYPE.CHAR); }
9  }

```

Estado de Fim de Citação de Caractere Único (END_CHAR_SINGLE_QUOTE)

Este estado lida com o fechamento da citação de caractere único e retorna o token apropriado ao mudar de volta para o estado YYINITIAL.

```
1 <END_CHAR_SINGLE_QUOTE> {  
2     "\"'" { yybegin(YYINITIAL); return symbol(TOKEN_TYPE.SINGLE_QUOTE); }  
3 }
```

Estado de Comentário Multilinha (COMMENT)

O estado de comentário multilinha lida com blocos de comentários que começam com {- e terminam com -}.

```
1 <COMMENT>{  
2     "-}" { yybegin(YYINITIAL); }  
3     [^"-}" { }  
4 }
```

Estado de Comentário de Linha (LINE_COMMENT)

O estado de comentário de linha lida com comentários que se estendem até o final da linha.

```
1 <LINE_COMMENT>{  
2     {NEW_LINE} { yybegin(YYINITIAL); }  
3     [^\n\r] { }  
4 }
```

Tratamento de Erros

Qualquer caractere não reconhecido é tratado como um erro, lançando uma exceção com uma mensagem descritiva.

```
1 [^] { throw new RuntimeException("Illegal character <"+yytext()+">"); }
```

2 Instruções de Execução

É necessário ter o Java instalado na máquina. O projeto foi desenvolvido com Java 21.0.2 e testado tanto nessa versão quanto na 17. Não é utilizado nenhuma funcionalidade específica de uma versão, portanto é provável que seja possível rodar em qualquer versão recente.

2.1 Linux

No Linux, rode o comando:

```
./build.sh
```


3 Conclusão

O desenvolvimento do analisador léxico foi uma etapa crucial na construção do compilador para a linguagem fictícia. Com a implementação bem-sucedida do analisador léxico, o sistema agora é capaz de identificar e categorizar tokens a partir dos arquivos de entrada, facilitando a análise inicial do código-fonte.

É importante destacar que esta fase representa apenas a primeira parte do projeto. As etapas restantes são essenciais para completar o sistema e garantir sua funcionalidade completa.