

Analizador Léxico

René Nolio
Ricardo Boing

Ferramenta utilizada: ANTLR4

O ANTLR (ANother Tool for Language Recognition), pronunciado *antler*, é um poderoso gerador de analisador léxico e sintático LL(k). É amplamente usado para construir linguagens, ferramentas e frameworks. A partir de uma gramática, o ANTLR gera um analisador sintático que pode construir e analisar árvores sintáticas. A última versão estável foi lançada em 2017. É o sucessor do PCCTS (Purdue Compiler Construction Tool Set), desenvolvido em 1989. Foi escrito em java e é *cross-plataform*.

Foi desenvolvida por Terence Parr, professor de ciências da computação da Universidade de São Francisco.



Arquivo.g4

- CLASS : 'class';
- EXTENDS : 'extends';
- ABRCHAVE : '{';
- PTVIR : '.';
- FECHCHAVE : '}';
- INT : 'int';
-

Arquivo.g4

- fragment LOWERCASE: [a-z];
- fragment UPPERCASE: [A-Z];
- fragment DIGIT: [0-9];
-

Arquivo.g4

- fragment LOWERCASE: [a-z];
 - fragment UPPERCASE: [A-Z];
 - fragment DIGIT: [0-9];
 -
-
- INTCONSTANT : (DIGIT)+;
 - STRINGCONSTANT : ["](LOWERCASE | UPPERCASE | DIGIT | ESPECIAL)*["];

Arquivo.g4

- fragment EXEMPLO:
- (
• [\u0060]
• | [\u005C]
• | [\u00C0-\u00C5]
•);

ANTLR

- `XccLexer lexer = new XccLexer(input);`
 - `CommonTokenStream tokens = new CommonTokenStream(lexer);`
 - `Vocabulary vocabularioAntlr = lexer.getVocabulary();`
-
- `List<Token> listaDeTokens = tokens.getTokens();`
 - `Token token = listaDeTokens.get(c);`
 - `nomeDoLexema = token.getText();`
 - `linhaDaOcorrencia = tokenAntlr.getLine();`
 - `colunaDaOcorrencia = tokenAntlr.getCharPositionInLine();`
 - `nomeDoToken = vocabularioAntlr.getSymbolicName(token.getType());`

Tarefa AS1: Questão 1

A gramática $X++$ é recursiva à esquerda?

- $P' = \varnothing$, conjunto de não terminais (V_n) sem recursão à esquerda
- $\forall A$, se $\forall P$ são da forma $A \rightarrow \beta\alpha$, $\alpha \in V^*$, $\beta \in V_t$, $A \in V_n$ então $P' = P' \cup \{A\}$
- Demais produções: se todas as produções de A são do tipo $A \rightarrow C\alpha$, onde $C \in P'$, então $P' = P' \cup \{A\}$
- Demais produções: se $\nexists A \rightarrow B\alpha$, então $P' = P' \cup \{B\}$ (pode $\exists B \rightarrow \varepsilon$, seria necessário um passo adicional)
- Neste ponto, $P' = P$, então não existe recursão à esquerda em nenhuma produção.

Tarefa AS1: Questão 2

A gramática X++ está fatorada à esquerda?

- Sejam f_1, f_2, \dots, f_n os conjuntos FIRST das produções de A. G está fatorada se $\forall A, f_1 \cap f_2 \cap \dots \cap f_n = \varnothing, A \in V_n$
- Computamos os conjuntos first para todas as produções:
- $\text{FIRST}(\text{statement}) = \{\text{INT}, \text{STRING}, \text{IDENT}\} \cup \{\text{IDENT}\} \cup \dots$

$\text{statement} \rightarrow (\text{INT} \mid \text{STRING} \mid \text{IDENT}) \text{IDENT} (\text{ABRCOL FECHCOL})^* (\text{VIRG IDENT} (\text{ABRCOL FECHCOL})^*)^* \text{PTVIR} \mid \text{IDENT} (\text{ABRCOL expression FECHCOL} \mid \text{PONTO IDENT} (\text{ABRPARG arglist FECHPAR } ?)^* \text{ATR} (\text{expression} \mid \text{alocexpression}) \text{PTVIR})$

Fatoramos:

$\text{statement} \rightarrow (\text{INT} \mid \text{STRING}) \text{IDENT} (\text{ABRCOL FECHCOL})^* (\text{VIRG IDENT} (\text{ABRCOL FECHCOL})^*)^* \text{PTVIR} \mid \text{IDENT statementaux}$

$\text{statementaux} \rightarrow \text{IDENT} (\text{ABRCOL FECHCOL})^* (\text{VIRG IDENT} (\text{ABRCOL FECHCOL})^*)^* \text{PTVIR} \mid (\text{ABRCOL expression FECHCOL} \mid \text{PONTO IDENT} (\text{ABRPCOL arglist FECHCOL} \mid \text{ATR (expression} \mid \text{alocexpression)}) \text{PTVIR}$

- As demais produções estão fatoradas

Tarefa AS1: Questão 3

A gramática X++ está em LL(3). Por que o 3?

- G está fatorada se qualquer palavra puder ser derivada de forma determinística, isto é, $\exists !$ produção $A \rightarrow \beta\alpha$, $\forall A \in V_n$. Ou seja, consultamos um símbolo à frente. LL(1)
- Para dois símbolos se qualquer palavra puder ser derivada de forma determinística, mas consultando dois símbolos à frente, isto é, $\exists !$ produção $A \rightarrow \beta\delta\alpha$, $\forall A \in V_n$; mas podem existir outras produções $A \rightarrow \beta\theta\alpha$, $\theta \neq \delta$; $\beta, \delta, \theta \in V_t$. LL(2)
- Não é LL(1): vide questão 2. Podemos demonstrar que não é LL(2). Portanto, é LL(3)

- Como G é $LL(3)$, é necessário verificar três símbolos à frente para decidir deterministicamente qual produção utilizar para derivar uma palavra.
- $LL(k)$ significa que faz o análise da esquerda pra direita (L = Left to right) e constroi a derivação mais à esquerda (L = leftmost derivation), o k indica quantos símbolos à frente o analisador irá verificar.

Referências bibliográficas

- <https://wwwantlr.org/>
- <https://unicode-table.com/pt/#control-character>
- <https://tomassetti.me/antlr-mega-tutorial/>
- DELAMARO, Márcio Eduardo. Como Construir um compilador. São Paulo, Novatec, 2004.