

Implementação de Analisador Sintático para Linguagem *Lang*

Diego Paiva - 201565516C
Thaynara Ferreira - 201565254C

Outubro, 2020

Resumo

O presente documento tem como objetivo esclarecer as principais decisões de projeto tomadas na implementação do analisador sintático para a linguagem fictícia *Lang*, referente ao trabalho prático da disciplina *DCC045 - Teoria dos Compiladores*.

1 Estrutura Léxica

A saber, na implementação deste trabalho, existem dois tipos de *tokens* em *Lang*, aqueles com estrutura fixa (separadores, operadores, etc) e os que seguem uma das formas a seguir:

- ID: Sequência de letras, dígitos e sobrescritos (underscore) que, obrigatoriamente, começa com uma letra. Ex: `var`, `var1`, `fun10`, etc.
- TYPE_NAME: Identificador cuja primeira letra é maiúscula. Ex: `Racional`, `Point`, etc.
- INT: Sequência de um ou mais dígitos. Ex: 0, 1, 10, 128, etc.
- FLOAT: Sequência de zero ou mais dígitos, seguido por um ponto e uma sequência de um ou mais dígitos. Ex: 3.141526535, 1.0, .12345, etc.
- CHAR: Um único caractere delimitado por aspas simples. Os caracteres especiais quebra-de-linha, tabulação, *backspace* e *carriage return* são definidos usando os caracteres de escape `\n`, `\t`, `\b` e `\r`, respectivamente. Para especificar um caractere `\`, é usado `\\` e para aspas simples o `'`. Ex: `'a'`, `'\n'`, `'\t'`, `'\\'` e `'\''`, etc.

2 Estrutura Sintática

A linguagem fictícia *Lang* é uma linguagem com propósitos meramente educacionais. Ela possui construções simples como tipos de dados, declarações, funções, comandos e expressões. Sua estrutura pode ser representada pela gramática a seguir.

```
prog → data* func*
data → data TYPE_NAME '{' decl* '}'
decl → ID '::' type ';'
func → ID '(' params? ')' (':' type '(' type? ')' '{' cmd* '}'
params → ID '::' type (',' ID '::' type)*
type → type '[' ']'
      | btype
btype → Int
      | Char
```

```

| Bool
| Float
| TYPE_NAME
cmd → '{' cmd* '}'
| if '(' exp ')' cmd
| if '(' exp ')' cmd else cmd
| iterate '(' exp ')' cmd
| read lvalue ';'
| print exp ';'
| return exp '(' exp)* ';'
| lvalue '=' exp ';'
| ID '(' exps? ')' ('<' lvalue '(' exp)* '>')? ';'
exp → exp '&&' exp
| rexp
rexp → aexp '<' aexp
| rexp '==' aexp
| rexp '!=' aexp
| aexp
aexp → aexp '+' mexp
| aexp '-' mexp
| mexp
mexp → mexp '*' sexp
| mexp '/' sexp
| mexp '%' sexp
| sexp
sexp → '!' sexp
| '-' sexp
| true
| false
| null
| INT
| FLOAT
| CHAR
| pexp
pexp → lvalue
| '(' exp ')'
| 'new' btype '[' sexp ']'?
| ID '(' exps? ')' '[' exp ']'
lvalue → ID
| lvalue '[' exp ']'
| lvalue '.' ID
exps → exp '(' exp)*

```

Onde a^* e $a^?$ denotam, respectivamente, zero ou mais ocorrências e zero ou uma ocorrência

da regra a . Cadeias entre aspas simples representam *tokens* e palavras em negrito são palavras reservadas da linguagem.

3 Ajuste da Gramática

É possível perceber que a linguagem *Lang*, assim como muitas outras, possui regras que são recursivas à esquerda, e.g., $\text{aexp} \rightarrow \text{aexp} \text{'+' mexp}$. Além disso, também há regras que contém prefixos em comum (fatoração à esquerda), e.g., $\text{mexp} \rightarrow \text{mexp} \text{'*'} \text{sexp} \mid \text{mexp} \text{'/'} \text{sexp} \mid \text{mexp} \text{'\%'} \text{sexp}$, o que acaba gerando ambiguidade na construção da linguagem.

Para tratar estes problemas, é necessário que sejam feitos ajustes nas regras que contém tais inconsistências. Nas Seções a seguir mostraremos como foram feitos estes tratamentos na gramática em questão.

3.1 Remoção da Recursão à Esquerda

Para uma determinada regra $A \rightarrow A\alpha \mid \beta$, é possível remover sua recursividade à esquerda substituindo A pelas regras $A \rightarrow \beta A'$ e $A' \rightarrow \alpha A' \mid \varepsilon$.

3.2 Refatoração à Esquerda

Para uma determinada regra $A \rightarrow \alpha \beta \mid \alpha \gamma$, é possível refatorá-la à esquerda substituindo A pelas regras $A \rightarrow \alpha A'$ e $A' \rightarrow \beta \mid \gamma$.

3.3 Ajustes em *Lang*

Apresentados os métodos para tratamento de recursão à esquerda e de fatoração à esquerda, enumeramos as regras que foram modificadas em *Lang*.

1. $\text{type} \rightarrow \text{btype type}'$
 $\text{type}' \rightarrow \text{'[' ']' type}'$
 $\mid \varepsilon$
2. $\text{cmd} \rightarrow \text{'\{' cmd* '\}'}$
 $\mid \text{if '(\ exp ')} \text{cmd cmd}'$
 $\mid \text{iterate '(\ exp ')} \text{cmd}$
 $\mid \text{read lvalue ';'}$
 $\mid \text{print exp ';'}$
 $\mid \text{return exp '(, exp)* ';'}$
 $\mid \text{lvalue '=' exp ';'}$
 $\mid \text{ID '(\ exps? ')} \text{'(< lvalue '(, lvalue)* '>')? ';'}$
 $\text{cmd}' \rightarrow \text{else cmd}$
 $\mid \varepsilon$
3. $\text{exp} \rightarrow \text{rexpr exp}'$
 $\text{exp}' \rightarrow \text{'\&\&' exp exp}'$
 $\mid \varepsilon$

4. $\text{rexp} \rightarrow \text{aexp} \text{ rexp}' \text{ rexp}'''$
 $\text{rexp}''' \rightarrow \text{rexp}'' \text{ rexp}'''$
 $\quad \mid \varepsilon$
 $\text{rexp}' \rightarrow '<' \text{ aexp}$
 $\quad \mid \varepsilon$
 $\text{rexp}'' \rightarrow '== ' \text{ aexp}$
 $\quad \mid '!= ' \text{ aexp}$
5. $\text{aexp} \rightarrow \text{mexp} \text{ aexp}''$
 $\text{aexp}'' \rightarrow \text{aexp}' \text{ aexp}''$
 $\quad \mid \varepsilon$
 $\text{aexp}' \rightarrow '+' \text{ mexp}$
 $\quad \mid '-' \text{ mexp}$
6. $\text{mexp} \rightarrow \text{sexp} \text{ mexp}''$
 $\text{mexp}'' \rightarrow \text{mexp}' \text{ mexp}''$
 $\quad \mid \varepsilon$
 $\text{mexp}' \rightarrow '*' \text{ sexp}$
 $\quad \mid '/' \text{ sexp}$
 $\quad \mid '\%' \text{ sexp}$
7. $\text{lvalue} \rightarrow \text{ID} \text{ lvalue}''$
 $\text{lvalue}'' \rightarrow \text{lvalue}' \text{ lvalue}''$
 $\quad \mid \varepsilon$
 $\text{lvalue}' \rightarrow '[' \text{ exp } ']$
 $\quad \mid '.' \text{ ID}$

4 Implementação

Para gerar um analisador sintático para a linguagem, utilizou-se a ferramenta **ANTLR**. Segundo [1], "ANTLR é um gerador de analisador poderoso para ler, processar, executar ou traduzir texto estruturado ou arquivos binários. É amplamente utilizado para construir linguagens, ferramentas e estruturas. A partir de uma gramática, ANTLR gera um analisador que pode construir e percorrer árvores de análise."

5 Building

Para realizar o *build* da aplicação, recomenda-se a utilização da ferramenta de gerenciamento **Maven** [2]. Com o Maven instalado, navegue até o diretório do projeto que contém o arquivo `pom.xml`. Após isso, rode o comando:

```
mvn clean compile assembly:single
```

Este comando irá compilar o código, incluir as dependências necessárias (ANTLR) e empacotar em um formato distribuível `jar` dentro da pasta `target`. Para executar o programa principal, rode o seguinte comando:

```
java -jar target/lang-compiler-1.0-jar-with-dependencies.jar -bs
```

Onde `-bs` é a opção para realizar a bateria de testes sintáticos.

Referências

- [1] “Antlr.” <https://www.antlr.org/>. Acesso em: 25/10/2020.
- [2] “Apache maven project.” <http://maven.apache.org/>. Acesso em: 04/10/2020.