# Universidade Federal de Santa Catarina Relatório T2 Construção de Compiladores

Isac de Souza Campos (17200449) Arthur Philippi Bianco (17203358) Enzo Coelho Albornoz (18100527)

#### Forma BNF

A forma Backus-Naur serve para representação de gramáticas livres de contexto. Devemos eliminar esta notação e converter para forma convencional, a seguir a gramática **antes** de realizar esta conversão:

```
PROGRAM → (STATEMENT|FUNCLIST)?
FUNCLIST → FUNCDEF FUNCLIST|FUNCDEF
FUNCDEF → def ident(PARAMLIST) {STATELIST}
PARAMLIST → ((int|float|string)ident, PARAMLIST|
(int|float|string)ident)?
STATEMENT →
(VARDECL; | ATRIBSTAT; | PRINTSTAT; | READSTAT; | RETURNSTAT; | IFSTAT | FOR
STAT|{STATELIST}|break;|;)
VARDECL → (int|float|string)ident([int_constant])*
ATRIBSTAT → LVALUE=(EXPRESSION|ALLOCEXPRESSION|FUNCCALL)
FUNCCALL → ident(PARAMLISTCALL)
PARAMLISTCALL → (ident, PARAMLISTCALL|ident)?
PRINTSTAT → print EXPRESSION
READSTAT → read LVALUE
RETURNSTAT → return
IFSTAT → if(EXPRESSION)STATEMENT(else STATEMENT)?
FORSTAT → for(ATRIBSTAT; EXPRESSION; ATRIBSTAT) STATEMENT
STATELIST → STATEMENT (STATELIST)?
ALLOCEXPRESSION → new(int|float|string)([NUMEXPRESSION])+
EXPRESSION → NUMEXPRESSION((<|>|<=|>=|==|!=)NUMEXPRESSION)?
NUMEXPRESSION → TERM((+|-)TERM)*
TERM → UNARY EXPR((* |/|%)UNARY EXPR)*
UNARYEXPR \rightarrow ((+|-))?FACTOR
```

```
FACTOR → (int_constant|float_constant|string_constant|null||LV ALUE|(NUMEXPRESSION))

LVALUE → ident( [NUMEXPRESSION] )*
```

O processo de conversão da gramática em notação BNF para notação convencional foi feito de forma a não produzir recursão à esquerda. Vejamos um exemplo:

```
NUMEXPRESSION → TERM ((+|-) TERM)*
```

Esta expressão pode ser convertida para o formato convencional, ingenuamente, para a seguinte expressão:

```
NUMEXPRESSION -> NUMEXPRESSION (+|-) TERM | TERM
```

Esta expressão é equivalente à expressão anterior, porém, o problema deste tipo de conversão é a introdução da recursão direta à esquerda. Outra alternativa seria criar um novo não-terminal:

```
NUMEXPRESSION -> NUMEXPRESSION2 | TERM
NUMEXPRESSION2 -> NUMEXPRESSION ((+|-) TERM)
```

Contudo, esta nova expressão possui também uma recursão, porém indireta à esquerda. Portanto, a melhor solução é a seguinte, que está na forma convencional e sem recursão à esquerda:

```
NUMEXPRESSION -> (TERM(ADDSUBTERM)) | TERM
ADDSUBTERM -> ((+|-)TERM)ADDSUBTERM | ((+|-)TERM)
```

Estando consciente destes problemas de introdução de recursão à esquerda à gramática ao converter para a forma convencional, foram removidas as recursões já durante o processo de conversão. A seguir a gramática resultante:

\*Programa ConvCC-2021-1 já sem recursão à esquerda:

```
PROGRAM -> STATEMENT | FUNCLIST | ε

FUNCLIST -> FUNCDEF FUNCLIST | FUNCDEF
```

```
FUNCDEF -> def ident(PARAMLIST){STATELIST} | def
ident(){STATELIST}
PARAMLIST -> ((int | float | string) ident, PARAMLIST | (int |
float | string ) ident)
STATEMENT->(VARDECL; | ATRIBSTAT; | PRINTSTAT; | READSTAT; |
RETURNSTAT; | IFSTAT | FORSTAT | {STATELIST} | break; | ;)
VARDECL -> ((int|float|string)ident(VARDECLARRAY)) |
((int|float|string)ident)
VARDECLARRAY -> ([int_constant]VARDECLARRAY) | [int_constant]
ATRIBSTAT -> LVALUE=(EXPRESSION|ALLOCEXPRESSION|FUNCCALL)
FUNCCALL -> ident() | ident(PARAMLISTCALL)
PARAMLISTCALL -> (ident, PARAMLISTCALL) | ident
PRINTSTAT -> print EXPRESSION
READSTAT -> read LVALUE
RETURNSTAT -> return
IFSTAT -> if(EXPRESSION)STATEMENT else STATEMENT |
if (EXPRESSION) STATEMENT
FORSTAT -> for (ATRIBSTAT; EXPRESSION; ATRIBSTAT) STATEMENT
STATELIST -> STATEMENT STATELIST | STATEMENT
ALLOCEXPRESSION -> new(int|float|string)NUMEXPRESSIONARRAY
NUMEXPRESSIONARRAY -> ([NUMEXPRESSION]NUMEXPRESSIONARRAY) |
[NUMEXPRESSION]
EXPRESSION -> (NUMEXPRESSION((< | >| <= | >= |
!=)NUMEXPRESSION)) | NUMEXPRESSION
NUMEXPRESSION -> (TERM(ADDSUBTERM)) | TERM
ADDSUBTERM -> ((+|-)TERM)ADDSUBTERM | ((+|-)TERM)
TERM -> UNARYEXPR TERMULTDIVMOD | UNARYEXPR
TERMULTDIVMOD -> (((* | / | %)UNARYEXPR)TERMULTDIVMOD) | ((* | /
| %)UNARYEXPR)
UNARYEXPR -> ((+|-)) FACTOR | FACTOR
FACTOR -> (int_constant|float_constant|string_constant|null||LV
ALUE | (NUMEXPRESSION))
LVALUE -> ident[NUMEXPRESSIONARRAY] | ident
```

### Recursão à esquerda:

A linguagem resultante não possui recursão direta à esquerda, eliminando a necessidade de tratamento desta propriedade. Estes casos são indesejáveis pois podem acabar gerando recursão infinita se o analisador sintático utilizado na compilação for recursivamente descendente, por exemplo, quebrando o algoritmo de compilação. A propriedade é caracterizada pela sua forma direta e indireta. Um exemplo de caso direto é visto a seguir:

•  $S \rightarrow S\alpha \mid \beta$ 

Nota-se que, o símbolo não terminal S, que é cabeça de produção, também é encontrado à esquerda de um dos possíveis caminhos a serem tomados. Porém, apesar da facilidade de se identificar o problema neste primeiro exemplo, em alguns casos a recursão à esquerda pode não ser tão facilmente encontrada, pois não se faz presente diretamente na produção de origem do símbolo não terminal analisado. A seguir vemos um exemplo de recursão indireta:

- S -> Aα | β
- A -> A $\alpha$  | S $\alpha$  |  $\beta$

Percebemos que, apesar de não existir recursão diretamente na produção de S, se o caminho tomado pela gramática seguir por  $A\alpha$  e, então, chegar a  $S\alpha$ , isso pode gerar o mesmo problema de recursão infinita visto no exemplo anterior.

Para provarmos que não há recursão à esquerda na gramática ConvCC-2021-1 obtida, utilizaremos a própria ferramenta *pest* que escolhemos para realizar a parte prática deste projeto. O *pest* é uma biblioteca de parsing de propósito geral que implementa um Parsing Expression Grammar [1]. Esta gramática é bastante parecida com gramáticas livre de contexto (CFG), porém os PEGs não são capazes de reconhecer gramáticas ambíguas, pois eles escolhem a primeira opção de derivação sempre, de modo que para cada string válida existe apenas uma árvore de análise sintática que a-descreva. Mais importante do que isso, PEGs são incapazes de expressar recursão à esquerda, ficando para sempre presos em loop nestas produções recursivas. Portanto, para poder utilizar este tipo de analisador, deve-se remover a recursão à esquerda. Reescrevemos a gramática usada em código para ser idêntica à BNF que produzimos no item anterior, e, caso haja recursão à esquerda, existirão

erros. Como o *pest* conseguiu montar a PEG a partir da descrição BNF, ela necessariamente não possui recursão à esquerda.

#### Fatoração à esquerda

A fatoração à esquerda é importante pois garante o funcionamento eficiente de analisadores sintáticos preditivos que tenham uma gramática com essa característica como entrada. Isso ocorre pois a fatoração garante que nenhuma produção tenha prefixos iguais, eliminando os não-determinismos e permitindo que o analisador não cometa erros ao escolher o caminho a ser tomado.

A gramática ConvCC-2021-1 vista anteriormente possui algumas produções que necessitam de tratamento para se encaixarem neste perfil. A seguir, temos o resultado obtido após a fatoração ser aplicada.

#### ConvCC-2021-1 fatorada à esquerda:

```
PROGRAM -> STATEMENT | FUNCLIST_

FUNCLIST -> FUNCDEF FUNCLIST | &

FUNCDEF -> def ident(FUNCDEF_

FUNCDEF_ -> PARAMLIST) {STATELIST} | ) {STATELIST}

PARAMLIST -> ((int | float | string) ident, PARAMLIST | (int | float | string) ident)

STATEMENT -> VARDECL; | ATRIBSTAT; | PRINTSTAT; | READSTAT; | RETURNSTAT; |

IFSTAT | FORSTAT | {STATELIST} | break; |;

VARDECL -> (int|float|string)ident VARDECL_

VARDECL_ -> VARDECLARRAY | &

VARDECLARRAY -> [int_constant] VARDECL_

ATRIBSTAT -> LVALUE=(EXPRESSION | ALLOCEXPRESSION | FUNCCALL)
```

```
FUNCCALL -> ident(FUNCCALL_
FUNCCALL_ -> ) | PARAMLISTCALL)
PARAMLISTCALL -> ident PARAMLISTCALL_
PARAMLISTCALL -> , PARAMLISTCALL | \epsilon
PRINTSTAT -> print EXPRESSION
READSTAT -> read LVALUE
RETURNSTAT -> return
IFSTAT -> if(EXPRESSION)STATEMENT IFSTAT_
IFSTAT_ -> else STATEMENT | ε
FORSTAT -> for(ATRIBSTAT; EXPRESSION; ATRIBSTAT) STATEMENT
STATELIST -> STATEMENT STATELIST_
STATELIST -> STATELIST | ε
ALLOCEXPRESSION -> new(int | float | string) NUMEXPRESSIONARRAY
NUMEXPRESSIONARRAY -> [NUMEXPRESSION] NUMEXPRESSIONARRAY_
NUMEXPRESSIONARRAY | ε
EXPRESSION -> NUMEXPRESSION EXPRESSION_
EXPRESSION_ -> (< | >| <= | >= | !=) NUMEXPRESSION | \epsilon
NUMEXPRESSION -> TERM(NUMEXPRESSION_)
NUMEXPRESSION_ -> ADDSUBTERM | ε
ADDSUBTERM -> ((+|-)TERM) NUMEXPRESSION_
TERM -> UNARYEXPR TERM_
```

```
TERM_ -> TERMULTDIVMOD | ε

TERMULTDIVMOD -> (* | / | %)UNARYEXPR TERM_

UNARYEXPR -> ((+|-)) FACTOR | FACTOR

FACTOR -> int_constant | float_constant | string_constant | null | LVALUE |
(NUMEXPRESSION)

LVALUE -> ident LVALUE_

LVALUE_ -> [NUMEXPRESSIONARRAY] | ε
```

Algumas alterações além da fatoração foram feitas para fins de otimização da gramática. Após a fatoração, um símbolo não-terminal TERMULTDIVMOD\_ foi criado. Porém, este não-terminal possuía as mesmas configurações do símbolo TERM\_. Por conta disso e para diminuir o número de símbolos da gramática, os TERMULTDIVMOD\_ contidos à direita das produções foram substituídos por TERM\_. Outros dois casos como este foram os de substituição de ADDSUBTERM\_ por NUMEXPRESSION\_ e VARDECLARRAY\_ por VARDECL\_.

Além dessas substituições, uma pequena alteração foi feita na produção de PROGRAM, transformando o trecho "FUNCLIST | E" em FUNCLIST\_, uma vez que a produção de FUNCLIST\_ já possuía como produção o trecho adaptado.

Após obtermos a gramática fatorada, passamos para o próximo passo, que consiste em provar que a gramática está em LL(1).

#### Gramática em LL(1):

O último passo a ser dado no tratamento da gramática deste trabalho é provar que ela está em LL(1). Para isso, devemos primeiro calcular os conjuntos *First* e *Follow* dos símbolos pertencentes à ela. Visando uma maior objetividade do relatório, apenas os conjuntos referentes aos símbolos não-terminais serão demonstrados. O resultado da obtenção de ambos os conjuntos pode ser visto na tabela 1.

FIRST	FOLLOW
int, float, string, ident, print, read, return, if,	
for, {, break, ';'def, &	\$
def	\$
def, E	\$
def	def, \$
int, float, string, )	def, \$
int, float, string	)
int, float, string, ident, print, read, return, if, for, {, break, ';'	\$, int, float, string, ident, print, read, return, if, for, '{', break, ';'
int, float, string	### 7
[, &	"." '
]	"-"
ident	";", )
(, ident	";", )
ident	)
",", ε	)
print	"." '
read	"-"
return	":"
if	\$, int, float, string, ident, print, read, return, if, for, {, break, ';'
else, &	\$, int, float, string, ident, print, read, return, if, for, {, break, ';'
for	\$, int, float, string, ident, print, read, return, if, for, {, break, ';'
int, float, string, ident, print, read, return, if, for, {, break, ';'	}
int, float, string, ident, print, read, return, if, for, {, break, ';', &	}
new	",", )
[	",", ), ]

3, ]	",", ), ]
+', '-', int_constant, float_constant, string_constant, null, ident, '('	"." , , )
<, >, <=, >=, !=, E	";", )
+', '-', int_constant, float_constant, string_constant, null, ident, '('	], <, >, <=, >=, !=, ";", )
"+", "-", E	], <, >, <=, >=, !=, ";", )
"+", "-"	], <, >, <=, >=, !=, ";", )
+', '-', int_constant, float_constant, string_constant, null, ident, '('	+', '-', ']', '<', '>", '<=', '>=', '==', '!=', ';', ')'
"*", "/", "%", E	+', '-', ']', '<', '>", '<=', '>=', '==', '!=', ';', ')'
"*", "/", "%"	+', '-', ']', '<', '>", '<=', '>=', '==', '!=', ';', ')'
+', '-', int_constant, float_constant, string_constant, null, ident, '('	*', '/', '%', +, -, ], <, >, <=, >=, ==, !=, ';', ')'
int_constant, float_constant, string_constant, null, ident, (	*', '/', '%', +, -, ], <, >, <=, >=, !=, ';', ')'
ident	=', '*', '/', '%', +, -, ], <, >, <=, >=, ==, !=, ';', ')'
[, &	=', '*', '/', '%', '+', '-', ']', '<', '>', '<=', '>=', '==', '!=', ';', ')'

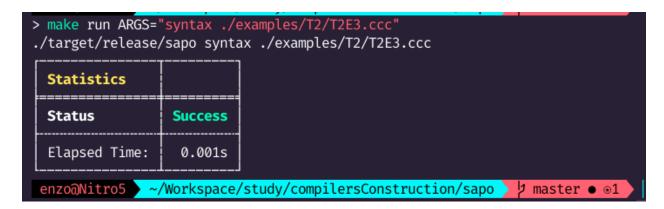
Tabela 1: Conjuntos First e Follow obtidos através da gramática ConvCC-2021-1.

Com base nos conjuntos gerados, uma Tabela de Análise Sintática pôde ser criada para observarmos se a gramática pertence ou não à LL(1). Lembrando que, para que uma gramática pertença à LL(1), cada espaço de relação entre um símbolo terminal e um não-terminal na tabela (que indica o caminho a ser tomado pelo algoritmo) deve contar com apenas uma produção. Isso ocorreu na tabela feita para a gramática ConvCC-2021-1 em estudo. Por conta de seu tamanho extenso, esta tabela não será mostrada no presente documento. Porém, pode-se ter acesso à mesma através do link permanente:

https://docs.google.com/spreadsheets/d/e/2PACX-1vTm8uK7CX\_v2Zc0gl97WdxfHh7iojtpPaPaZfcYF82Gl96x8QhFbRdq3kGymWeDxoewujYBn1LYk7B3/pubhtml?gid=1892824522&single=true

## Programa

Foi adicionado o subcomando **syntax** cuja execução irá fazer uma análise sintática do arquivo fonte dado. Em caso de sucesso, uma mensagem de sucesso irá aparecer.



Por uma questão de como o *parser* é implementado (como ele reconhece mais que uma LL(1), seu reconhecimento não possui o mesmo fluxo que execução, consumindo os tokens de maneira diferente e usando algumas estruturas diferentes) não temos como mostrar a entrada na tabela LL(1). Em outras palavras, não há tabela de análise sintática em LL(1) da qual tirar as informações requisitadas.

Contudo, temos os dados de token lido e tokens esperados. Caso haja um erro na leitura é mostrado o token lido e dito quais tokens seriam os corretos. Vejamos um exemplo abaixo:

Para validar o sistema (tanto para sucesso quanto para erro), executamos o build da aplicação (obtendo um tempo de 32.89 segundos) e verificamos os códigos fonte desenvolvidos para esta entrega. Os testes foram feitos em um I5 8300H (Notebook Acer).

## Referencias

[1] Pest The Elegant Parser. Disponível em: <a href="https://github.com/pest-parser/pest">https://github.com/pest-parser/pest</a>. Acesso em 13 de Agosto de 2021.