

Universidade do Minho

Mestrado Integrado em Engenharia Informática

Processamento de linguagens - Trabalho Prático #2 Ano Letivo 2021/2022

Gonçalo Pereira (a93168) Tiago Silva (a93277)

22 de setembro de 2022

1 Introdução

O presente relatório refere-se à realização do trabalho prático 2, cujo o objetivo era de um total de quatro problemas escolher um e proceder à implementação de uma solução. O enunciado escolhido "Tradutor PLY-simple para PLY" que, resumidamente, requer a criação de uma linguagem mais "limpa" e simplificada inspirada no PLY, chamada PLY-simple, deduzir um esquema de tradução e criar um compilador que traduza esta linguagem em PLY. Para isso utilizamos geradores de compiladores baseados em gramáticas tradutoras, concretamente o Yacc, versão PLY do Python, completado pelo gerador de analisadores léxicos Lex, também versão PLY do Python.

2 Sintaxe da linguagem Ply-simple

Começamos por definir a sintaxe a usar na nova linguagem Ply-simple. Embora tenhamos recebido uma sugestão de uma sintaxe possível, esta serviu apenas para inspiração já que a versão final da nossa sintaxe de Ply-simple tem algumas diferenças. Dividindo a sintaxe em dois (Lexer e Yacc) começaremos por mostrar cada uma das partes em mais detalhe.

2.1 Lexer

Exemplo:

```
LEX:
literals = "+-*=()"

ignore = " \t\n "

tokens = ['VAR','NUMBER']

flunc: [a-zA-Z\_][a-zA-ZO-9\_]* : return('VAR', t.value)

flunc: \d+(\.\d+)? : return('NUMBER', float(t.value))

lerror(f"Illegal character '{t.value[0]}', [{t.lexer.lineno}]",

t.lexer.skip(1) )
```

Inicialmente no lexer caso pretendido pode-se, não sendo obrigatório, definir os literals e a regra de ignore,por outro lado, tokens tem que estar sempre definido. Além disso a definição das funções do lexer seguem a seguinte estrutura: *lfunc:* para inciar uma função do lexer, seguido de a expressão regular que pretende usar, seguido de : que marca o fim a expressão regular e inicia o *return* que contêm o nome da função e do valor a retornar.

2.2 Yacc

Exemplo:

```
YACC:
2
     precedend =
3
                   ('left','+','-'),
                   ('left','*','\'),
                   ('right', 'UMINUS'),
6
     ts = \{\}
9
10
     grammar:
11
     stat : VAR '=' exp { ts[p[1]] = p[3] }
12
     stat : exp { print(p[1]) }
13
     exp : exp '+' exp { p[0] = p[1] + p[3] }
14
     exp : exp '-' exp { p[0] = p[1] - p[3] }
15
     exp : exp '*' exp { p[0] = p[1] * p[3] }
16
     exp : exp '\' exp { p[0] = p[1] / p[3] }
     exp : '-' exp %prec UMINUS { p[0] = -p[2] }
18
     exp : '('exp ')' { p[0] = p[2] }
19
     exp : NUMBER \{ p[0] = p[1] \}
20
     exp : VAR { p[0] = getval(p[1]) }
21
22
     yfuncs:
23
     def p_error(t){
24
          print(f"Syntax error at '{t.value}', [{t.lexer.lineno}]")
25
26
     def getval(n){
27
          if n not in ts: print(f"Undefined name '{n}'")
28
          return ts.get(n,0)
29
          }
30
```

Na definição do Yacc é possível definir precedência como pode ser visto acima, de forma semelhante, pode-se iniciar variáveis de estado que podem ser listas, dicionários, tuplos, strings e valores numéricos. Depois chega à parte principal que é a da definição da gramática. A palavrachave grammar: marca o início da definição da gramática. A gramática segue uma estrutura comum, em que cada podução é caracterizadas por Simbolos não terminais: Conjunto de simbolos terminais ou não terminais seguida de $\{Ações semânticas\}$. Por fim usando yfunc: inicia-se a definição das funções do Yacc tal como se mostra no exemplo.

3 Solução implementada

3.1 Analisador léxico

3.1.1 Tokens

Na imagem a seguir podemos ver os tokens que foram definidos e as suas respetivas expressões regulares.

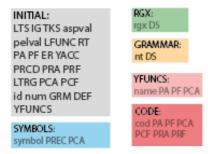
Token -> expressão regular* LEX -> r'LEX:' PA -> r'\(' LTS -> r'literals' $PF \rightarrow r' \downarrow '$ IG -> rignore' PCA -> r'{' TKS -> r'tokens' $PCF \rightarrow r'$ LFUNC -> r'lfunc' PRA -> r'[PRF -> r']' RT -> r'return' ER-> rlemor\(DS -> r':' TVAL -> r't.value' aspval -> r"((\\")|[^0-9\n'])+" TYPE -> r'(float)|(int)|(str) pelval -> r"((\\')|[^0-9\n'])+" cod -> r'((("[^"]+")|(\'[^\']+\'))|([^{}\(\)\[\]]))+' YACC -> rYACC: PRCD -> r'precedent' $id \rightarrow r''[A-Za-z]+"$ LTRG -> r"('left')|('right')" $nt \rightarrow r''[A-Za-z]+'$ GRM -> r'grammar:' symbol -> r"([A-Za-z]+)|(!)|('\")" PREC -> r'%prec' name -> r"[A-Za-z]\w*" YFUNCS -> r'yfuncs:' num -> $r''(\+\-)?\d+(\.\d+)?''$ DEF -> r'def' $rgx \rightarrow r''((\\)[^:])+"$

Tokens

De seguida necessitamos de definir estados, e por isso, segue uma imagem com todos os estados definidos e os todos tokens associados a eles.

States and their defined tokens:

formate python raw string



Estados

Por fim, temos uma imagem que nos ajuda a ver os tokens resultantes de um input exemplo, e que inclusive, ajuda a perceber como a gramática foi construída.

```
PLY-SIMPLE example:
                                                                                                                           RESULT TOKENS:
LEX:
                                                                         LTS = aspval
literals = "+-/"=0"
ignore = ' \t\n'
                                                                         IG = aspval
tokens = [ 'VAR' , 'NUMBER' ]
                                                                         TKS = [pelval, pelval]
func: [a-zA-Z_][a-zA-Z0-9_] : return( 'VAR' , t.value )
                                                                         LFUNC rgx: RT (pelval, TVAL)
                                                                         LFUNC rgx:RT (pelval, TYPE (TVAL))
ER cod {cod [cod]} cod [{cod
lfunc: \d+(\.\d+)?: return( 'NUMBER' , float(t.value) )
lerror(f"Illegal character '{t.value[0]}', [{t.lexer.line
]]",t.lexer.skip(1))
                                                                         }]cod(cod))
YACC:
                                                                         YACC
precedend = [
                                                                         PRCD = [
      ('left' , '+' , '-'),
('left' , '*' , '/'),
                                                                                (LTRG, pelval, pelval),
                                                                                (LTRG, pelval, pelval),
       ('right' , 'UMINUS'),
                                                                                (LTRG, pelval),
                                                                         id = { }
td = [5, "ola", 'adeus']
                                                                         id = [ num , aspval , pelval]
grammar:
                                                                         GRM
stat : VAR'=' exp { ts[t[1]] = t[3] }
                                                                         nt : symbol symbol symbol { cod [ cod ] ] cod [ cod ] }
                                                                         nt : symbol {cod (cod [cod])}
nt : symbol symbol symbol {cod [cod] cod [cod]}
nt : symbol symbol symbol {cod [cod] cod [cod]}
stat : exp { print(t[1]) }
exp: exp'+'exp{t[0] = t[1] + t[3]}
exp: exp'-'exp{t[0] = t[1] - t[3]}
exp : exp '"' exp { t[0] = t[1] * t[3] }
                                                                         nt:symbol symbol symbol { cod [ cod ] cod [ cod ] cod [ cod ] }
exp : exp '/' exp { t[0] = t[1] / t[3] }
exp : '-' exp %prec UMINUS { t[0] = -t[2] }
                                                                         nt : symbol symbol symbol { cod [ cod ] cod [ cod ] cod [ cod ]}
nt : symbol symbol PREC symbol { cod [ cod ] cod [ cod ]}
exp: '(' exp ')' { t[0] = t[2] }
exp: NUMBER { t[0] = t[1] }
                                                                         nt:symbol symbol symbol { cod [ cod ] cod [ cod ]}
                                                                         nt:symbol {cod[cod]cod[cod]}
exp : VAR { t[0] = getval(t[1]) }
                                                                         nt:symbol {cod [cod] cod (cod [cod])}
                                                                         YFUNCS
def p_error(t){
                                                                         DEF name ( name ) {
  print(f"Syntax error at '{t.value}', [{t.lexer.lineno}]")}
                                                                            cod(cod{cod}cod[{cod}])}
def getval(n){
                                                                         DEF name ( name ) {
  if n not in ts: print(f"Undefined name '(n)")
                                                                            cod ( cod { cod } cod )
   return ts.get(n,0)}
                                                                            cod(cod)}
CAPTION (LEXER STATES):
  INITIAL SYMBOLS YFUNC
  CODE REGEX GRAMMAR
```

Ply-simple example

3.2 Analisador semântico

3.2.1 Gramática

Gramática definida:

Ply: Lexer Yc

Lexer: LEX Literals Ignore Tokens L
funcs Lerror Yc: YACC Precedents Declaration Grammar Yfs

Literals: LTS '=' aspval

Literals:

Ignore: IG '=' aspval

Ignore:

Tokens: TKS '=' PRA Tokl PRF

Tokl: Tokl', pelval

Tokl: pelval

Lfuncs: Lfuncs Lfunc

Lfuncs:

Lfunc : LFUNC rgx DS RT PA pelval ',' Tv PF

Tv:TVAL

Tv : TYPE PA TVAL PF Lerror : ER Codes PF

Lerror:

Precedents: PRCD '=' PRA Predlist PRF

Precedents:

Declarations : Declaration

Declarations:

Declaration : id '=' Values Grammar : GRM Productions

Yfs: YFUNCS Funcs

Productions: Production

Productions:

Production: nt DS Symbols PCA Codes PCF

Symbols: Symbols S

Symbols: S:symbol S:PREC symbol

Prcdlist: Prcdlist PA LTRG ',' Pelvals pelval PF ','

Prcdlist:

Pelvals: Pelvals pelval','

Pelvals:

Funcs: Funcs Func

Funcs:

Func: DEF name PA name PF PCA Codes PCF

Codes: Codes Code

Codes: cod

Code: PA Codes PF Code: PRA Codes PRF Code: PCA Codes PCF Values: Values',' Type

Values: Type Type: num Type: aspval Type: pelval

Type: PRA Cont PRF Type: PA Cont PF Type: PCA Cont PCF

Cont: Values

Cont:

4 Verificação de erros

Uma importante característica de um gerador de código é a capacidade de detetar erros e de informar o utilizador dos mesmos. Contudo só implementamos duas verificações, uma que verifica tokens definidos mas não usados e outra que que verifica se estão a ser utilizados símbolos não terminais que não foram definidos. Para a implementação destas verificação são usadas variáveis de estado tt - tabela de terminais, tntDef - tabela de não terminais definidos e tntUsed - tabela de não terminais utilizados. Em tt são guardados todos os tokens e literais definidos, em tntDef são guardados todos os não terminais definidos (lado esquerdo de uma produção) e um campo que contabiliza quantos produções têm um dado não terminal, por fim em tntUsed guarda-se todos os não terminais (usados no lado direito de um produção). A primeira verificação traduz-se num varning para o utilizador enquanto que a segunda traduz-se num varning para o utilizador enquanto que a segunda traduz-se num varning para o utilizador enquanto que a segunda traduz-se num varning para o utilizador enquanto que a segunda traduz-se num varning para o utilizador enquanto que a segunda traduz-se num varning para o utilizador enquanto que a segunda traduz-se num varning para o utilizador enquanto que a segunda traduz-se num varning para o utilizador enquanto que a segunda varning para varning para o utilizador enquanto que a segunda varning para varning para o utilizador enquanto que a segunda varning para varning para

5 Exemplo de Utilização

Segue-se um exemplo de linguagem PLY-Simple usada para definir S expressions e o código PLY gerado.

```
LEX:
2
     literals = "()+*-"
3
     ignore = "\t\n "
     tokens = ['num']
5
     lfunc: \d+ : return('num',int(t.value))
     YACC:
9
10
     grammar:
11
     Z : Sexp {print(p[1])}
12
     Sexp : '(' '+' Lista ')' \{p[0] = somatorio(p[3])\}
13
     Sexp : '(' '*' Lista ')' {p[0] = produtorio(p[3])}
14
     Sexp : num \{p[0] = p[1]\}
15
     Lista: Lista Sexp \{p[0] = p[1] + [p[2]]\}
16
     Lista : Sexp Sexp \{p[0] = [p[1], p[2]] \}
17
18
     yfuncs:
19
      def somatorio(lista){
20
          res = 0
          for n in lista:
22
              res += n
23
          return res
24
     }
25
     def produtorio(lista){
26
          res = 1
27
          for n in lista:
28
              res *= n
29
          return res
30
     }
31
```

```
import ply.lex as lex
     literals = ['(', ')', '+', '*', '-']
     t_ignore = "\t\n "
     tokens = ['num']
     def t_num(t):
         r'\d+ '
         t.value = int(t.value)
         return t
     def t_ANY_error(t):
10
          print(f"Illegal character '{t.value[0]}', [{t.lexer.lineno}]")
11
          t.lexer.skip(1)
13
     lexer = lex.lex()
14
15
      import ply.yacc as yacc
16
17
     def p_Z(p):
18
          "Z : Sexp "
19
         print(p[1])
20
21
      def p_Sexp(p):
22
          "Sexp : '(' '+' Lista ')' "
23
         p[0] = somatorio(p[3])
24
25
      def p_Sexp_1(p):
26
          "Sexp : '(' '*' Lista ')' "
27
         p[0] = produtorio(p[3])
28
     def p_Sexp_2(p):
30
          "Sexp : num "
31
         p[0] = p[1]
32
33
     def p_Lista(p):
34
          "Lista : Lista Sexp "
35
         p[0] = p[1] + [p[2]]
36
37
     def p_Lista_1(p):
38
          "Lista : Sexp Sexp "
39
         p[0] = [p[1], p[2]]
40
41
     def somatorio(lista):
42
          res = 0
43
          for n in lista:
44
              res += n
45
         return res
46
47
     def produtorio(lista):
48
          res = 1
49
          for n in lista:
              res *= n
51
52
         return res
53
     parser = yacc.yacc()
```

6 Conclusão

Com este trabalho sentimos que alcançamos os principais objetivos propostos tais como • aumentar a experiência em engenharia de linguagens e em programação generativa (gramatical), reforçando a capacidade de escrever gramáticas, quer independentes de contexto (GIC), quer tradutoras (GT); • desenvolver processadores de linguagens segundo o método da tradução dirigida pela sintaxe, a partir de uma gramática tradutora; • desenvolver um compilador gerando código para um objetivo específico; • utilizar geradores de compiladores baseados em gramáticas tradutoras, concretamente o Yacc, versão PLY do Python, completado pelo gerador de analisadores léxicos Lex, também versão PLY do Python. Tendo em conta ao que nos foi proposto acreditamos ter feito um bom trabalho, embora tenhamos consciência que a nossa aplicação ficou longe de ser perfeita, uma vez que não implementamos os estados do ply e fazemos poucas verificações a nível do código em Python.