



Instituto Politécnico de Bragança
Engenharia Informática
Processamento de Linguagem

Trabalho nº 15
Trabalho Prático 1

Aluno 1: Gaofeng Yin a36233
Aluno 2: Hugo Gázio a36481
Professor: Paulo Matos & Célia Ramos

Maio
2020

Resumo

Uma linguagem para ser valida tem que cumprir as regras sintáticas e semânticas, porém nem todas as frases sintacticamente validas são semanticamente corretas. A gramática é um modo de descrever sintaxe de linguagens. As expressões regulares são utilizadas para descrever os símbolos válidos de uma dada linguagem.

Neste trabalho prático vamos construir um parser que faz o processamento da estrutura de dados da notação JSON que posteriormente serão necessários para próximo trabalho.

O LEX é o analisador léxico e a sua execução é chamado dentro do próprio ficheiro YACC. O YACC por sua vez valida a sequência de tokens detectados. Com YACC e LEX devemos produzir um analisador léxico e sintático que deve ser capaz de processar os dados de um ficheiro como input em formato JSON do Google Calendar.

Desenvolvimento

LEX

Iniciamos com o ficheiro.l, fazemos o include do y.tab.h em que vai estabelecer a comunicação entre YACC e LEX. Criamos as expressões regulares para identificar os valores de cada chave, vou apenas abordar algumas ER mais complexas:

A expressão regular date é capaz de conhecer datas em formato yyyy-mm-

```
3 //data em formato 1996-07-13
4 date (1|2)[0-9]+\-((0)[1-9]|1[012])\-(0[1-9]|12)[0-9]|3[01])
5 //hora em formato 23:50:46 / 06:00:00pm
6 time ([0-1]?[0-9]|2[0-3])\:[0-5][0-9]\:([0-5][0-9]|60)(pm|am)?
7 //data e hora em formato 2003-05-17 15:45:23
8 datetime {date}" "{time}
9 //expressão regular para websites
10 url (http:|https:)?(\\|/)?(www\\.)?[a-zA-Z0-9\\|\\.|\\#st]+
11 //expressão regular que identifica email
12 email [a-zA-Z@\\|\\._-0-9st]+
```

Figura 1: Expressões regulares pré definidas.

dd. O ano tem que começar com o digito 1 ou 2 para ser valido, seguido de um ou mais numero de 0-9, o mês vai de 01 até 09, mas também pode ser 10, 11 ou 12 e aplicando a mesma lógica obtemos o dia. Os respectivos yyyy mm dd é separados pelo carácter '-'. O time pode ser em 00-23h ou AM/PM.

Depois unimos o date e time para formar apenas uma expressão regular de datetime que serve para identificar a data e hora ao mesmo tempo.

A seguir definimos as regras para reconhecer as chaves e o respectivo valor:

O timeZone, hidden, date, createdon, url e entre outros chaves são palavras

```
5  "\"timeZone\": \"{string}\"      {return TIMEZONE;}
6  "\"hidden\": \"{boolean}\"      {return HIDDEN;}
7  "\"date\": \"{date}\"          {return DATE;}
8  "\"createdon\": \"{datetime}\"  {return CREATEDON;}
9  "\"url\": \"{url}\"            {return URL;}
10 [,{,}:\\[\\]"                  {return (yytext[0]);}
11 [ \\t\\n]                      ;
12 .                              {printf("Token invalido\\n");}
```

Figura 2: Identificadores de tokens.

reservadas ao qual está associado um valor dinâmico e por isso usamos a expressão regular pré definida entre chavetas ({string}) capaz validar o valor. Os símbolos da linha 10 da figura 2 são retornados como tokens. O programa não faz nada quando é encontrado uma quebra de linha, espaço e tab como está definido na linha 11 da figura 2. E por fim tudo o resto são tokens inválidos.

YACC

No YACC definimos a função main e a sequências de gramáticas validas. Os tokens que declaramos no YACC são reconhecidos no LEX e são usados para retornar o significado do input encontrado. É de denotar que os símbolos terminais estão em maiúsculas e os não terminais em minúsculas.

O start indica o ponto da partida da árvore de derivação que por derivações

```
9      //ponto da partida
10 %start program
11      //tokens para identificar terminais
12 %token KIND ETAG ID SUMMARY NAME BACKGROUNDCOLOR TYPE EMAIL LOCATION KEY
13 %token TIMEZONE TIME HIDDEN RESERVED OBJECT LINK DATE CREATEDON URL ADD_GUEST CODE
```

Figura 3: Tokens.

vai chegar até um terminal. Ao longo deste processo temos que analisar com atenção o documento de input, a forma como os dados estão agrupados e organizados.

```

17      /*programa começa com duas chavetas que abrange todo ficheiro json,
18      dentro desses chavetas estão os dados.*/
19  program      :      '{' datas '}'
20      ;
21      /*recursiva a direita.Existe pelo menos 1 ou mais objetos.*/
22  datas      :      data datas
23      |
24      ;
25      /*os dados podem ser key , array , object ou array_object conforme o dado processado*/
26  data      :      key comma
27      |      array comma
28      |      object comma
29      |      array_object comma
30      ;

```

Figura 4: Gramática.

Começamos com o program que vai derivar em datas entre as duas chavetas. considerando que um input vai ter pelo menos 1 ou mais dados, então recorremos a recursividade a direita para dizer que datas vai derivar em data datas ou nada.

Por sua vez data pode derivar em key, array , object ou array_object conforme a estrutura de dados. A comma deriva em virgula ou nada porque há

```

65      // é uma virgula ou não é nada
66  comma      :      ','
67      |
68      ;

```

Figura 5: Virgulas.

caso de dados que aparecem com virgula e outros não, e para poupar linhas de código optei por essa solução.

```

4      /*o array podem ter varios valores dentro ou array multidimensional */
5  array      :      '[]' OBJECT '[]' '[' datas ']'
6      ;
7      /*dentro de array pode ter objetos */
8  array_object:      '[]' OBJECT '[]' '[' program ']'
9      ;
10     /*um objeto*/
11  object      :      '[]' OBJECT '[]' '[' '{' datas '}'
12      ;
13     //key
14  key      :      KIND
15      |      ETAG
16      |      (...)
17      |      KEY
18      |      RESERVED
19      ;

```

Figura 6: Formato de dados.

A sintaxe do array é constituído por um OBJECT que identifica o array e dentro tem datas que podem derivar em key , um objecto ou em array outra

vez.

O array_object tem um OBJECT que o identifica e dentro é outro program que pode ter varias derivações e pode ser tudo que definimos anteriormente. O object tal como array tem a estrutura semelhante, apenas muda nas chavetas. E por fim temos o key que são as chaves e o valor associado.

A função yyerror é chamado quando algo de incorrecto é detectado e mostra

```
72 /*o yyerror para identificar a linha onde ocorreu o erro caso haja alguma */
73 int yyerror(char *msg){
74     fprintf(stderr, "ERRO(%d):%s\n",yylineno, msg);
75     variavel = 1;
76     return 0;
77 }
```

Figura 7: função yyerror.

no ecrã o tipo de erro e a linha onde foi detectado o que ajuda imenso na implementação do parser e é essencial para reconhecimento do ficheiro JSON.

Finalmente a função main, se estiver tudo correcto é mostrado na ecrã

```
78 /*caso esteja tudo bem é executado o comando de printf dentro do if*/
79 int main(){
80     yyparse();
81     if(variavel != 1){
82         printf("Processado com sucesso!\n");
83     }
84     return 0;
85 }
```

Figura 8: Função main.

”Processado com sucesso!”. Para isso criei uma variável de verificação. Se a função yyerror é executado então a variavel passa ser 1 e o printf do main não vai ser executado.

Resultado

Como podemos constatar na figura 9 os comandos yacc e flex executou sem nenhum erro e geramos o executável com sucesso. Na pasta estão 4 ficheiros de input, o input, input2 e input3 são aqueles em que os dados estão em formato JSON e foram validados com sucesso. O input4 deu erro porque o seu conteúdo não é legível, apenas foi incluído para mostrar o output quando dá erro.

```
gaofeng@gaofeng-GT70-2PE: ~/Documents/Trabalho_pratico_1_PL
gaofeng@gaofeng-GT70-2PE: ~/Documents/Trabalho_pratico_1_PL 89x17
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ yacc -d ficheiro.y
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ flex ficheiro.l
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ gcc lex.yy.c y.tab.c -o exe
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ ls
exe      ficheiro.y  input2  input4  teste  y.tab.h
ficheiro.l  input      input3  lex.yy.c  y.tab.c
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ ./exe < input
Processado com sucesso!
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ ./exe < input2
Processado com sucesso!
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ ./exe < input3
Processado com sucesso!
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$ ./exe < input4
ERRO(1):syntax error
gaofeng@gaofeng-GT70-2PE:~/Documents/Trabalho_pratico_1_PL$
```

Figura 9: Resultado final.