

Processamento de Linguagens (3º ano de Curso)

Trabalho Prático

Relatório de Desenvolvimento

Bianca Araújo do Vale
(A95835)

João Pedro Machado Ribeiro
(A95719)

Telmo José Pereira Maciel
(A96569)

28 de maio de 2023

Resumo

Neste trabalho iremos desenvolver um compilador de código *TOML* gerando o respetivo código num formato *JSON*, através do uso de um gerador de compiladores baseado em gramáticas e analisadores léxicos, através do *Lex/Yacc* do *Ply/Python*.

Conteúdo

1	Introdução	2
1.1	Conversor toml-json	2
1.2	Estrutura do Relatório	2
2	Concepção/desenho da Resolução	3
2.1	Concepção da Linguagem	3
2.2	Analizador Léxico	3
2.2.1	Tokens	3
2.3	Analizador Sintático	5
2.4	Gramática	5
2.4.1	Programa principal	6
2.4.2	Secções	7
2.4.3	Subsecções	7
2.4.4	Tabela	7
2.4.5	Nome da secção ou subsecção	8
2.4.6	Conteúdo da secção	8
2.4.7	Conteúdo	8
2.4.8	Chave	9
2.4.9	Valores	9
2.4.10	Arrays	9
2.5	Controlo de Erros	10
3	Funcionalidades implementadas	11
3.1	Conversor <i>toml</i> para <i>json</i>	11
3.1.1	Ficheiro output	11
3.2	Conversor <i>toml</i> para <i>yaml</i>	12
3.2.1	Ficheiro output	12
3.3	Conversor <i>toml</i> para <i>xml</i>	13
3.3.1	Ficheiro output	13
3.4	Funcionamento do programa	14
4	Conclusão	15

Capítulo 1

Introdução

1.1 Conversor toml-json

Área: Processamento de Linguagens

O presente relatório tem como objetivo expor o trabalho prático desenvolvido na Unidade Curricular de Processamento de Linguagens do terceiro ano da Licenciatura em Engenharia Informática.

O trabalho escolhido pelo grupo foi o **Conversor toml-json (2.6)**, sendo que iremos abordar, ao longo do relatório, o processo de resolução do respetivo trabalho.

Para a realização do respetivo projeto, colocamos em prática diversos conhecimentos adquiridos nas aulas teóricas e práticas da respetiva UC: o uso de geradores de compiladores baseados em gramáticas tradutoras, tais como o *Yacc* do *PLY* do *Python* e o gerador de analisadores léxicos *Lex*, também do *PLY* do *Python*.

1.2 Estrutura do Relatório

O presente relatório encontra-se dividido em quatro capítulos:

- O primeiro capítulo refere-se à introdução, que descreve, de forma breve, o problema e retrata o objetivo do trabalho.
- O capítulo dois refere-se à concepção/desenho da resolução, onde o grupo apresenta a estratégia que utilizou para resolver o trabalho, desde os tokens criados até à gramática desenvolvida. Ainda neste capítulo está presente o controlo de erros, onde o grupo explica o que realizou para proceder à ultrapassagem dos mesmos.
- O terceiro capítulo caracteriza-se pelas funcionalidades que o grupo implementou no trabalho prático.
- No quarto e último capítulo está contida uma reflexão crítica, por parte do grupo, acerca do trabalho realizado.

Capítulo 2

Concepção/desenho da Resolução

2.1 Concepção da Linguagem

Primeiramente, para a realização do presente trabalho prático, começamos por desenvolver um analisador léxico com a ferramenta *Lex*. Este tem como objetivo identificar os tokens (definidos pelo grupo).

Posteriormente, procedemos à realização da gramática, que irá ser apresentada com mais detalhe mais à frente.

2.2 Analisador Léxico

O analisador léxico é responsável por analisar o código-fonte caracter por caracter e transformá-lo em sequências de *tokens*, preparando o nosso código para as próximas etapas do processo de compilação.

2.2.1 Tokens

Com o objetivo de conseguirmos representar as diversas funcionalidades do nosso programa, definimos alguns tokens, que se encontram no nosso ficheiro **lex**.

Para além da vírgula (**COMMA**), dos parêntes retos (**LBRACKET**, **RBRACKET**) e da igualdade (**EQUALS**), definimos um conjunto de *tokens* que achamos necessários para a resolução do nosso problema, um compilador de linguagem *toml*.

Fora os *tokens* anteriormente referidos, temos também os seguintes:

BASIC_STRING e LITERAL_STRING: Representa todas as strings multi-line;

STRING: Representa uma string;

COMMENT: Representa um comentário;

NUMBER: Representa um número, int ou float;

DATE, TIME e DATE_TIME: Representa todos os tipos de datas e horas suportadas;

BOOLEAN: Representa qualquer palavra do tipo *bool*;

KEY: Representa as chaves dos valores, isto é, tudo que vem atrás do símbolo '=';

OBJECT e SUBOBJECT: Representa todos os objetos da nossa linguagem.

```

1 tokens = (
2     'BASIC_STRING',      'LITERAL_STRING',    'STRING',
3     'COMMENT',           'NUMBER',        'IP',
4     'DATE',              'TIME',          'BOOLEAN',
5     'KEY',               'OBJECT',        'SUBOBJECT',
6     'LBRACKET',          'RBRACKET',      'EQUALS',
7     'COMMA',             'DATE_TIME',
8 )
9
10 literals = '-.:',=[]
11
12 t_ignore = ' \t\n'
13
14 def t_BASIC_STRING(t):
15     r'\''\''\''\''((.|\\n)*?)\''\''\''
16     t.value = str(t.value[3:-3])
17     return t
18
19
20 def t_LITERAL_STRING(t):
21     r'\"\"\"\"\"\"((.|\\n)*?)\"\"\"\"\"\"
22     t.value = str(t.value[3:-3])
23     return t
24
25 def t_BOOLEAN(t):
26     r'TRUE|FALSE|true|false'
27     t.value = True if t.value == 'true' else False
28     return t
29
30 def t_IP(t):
31     r'\"\"\\d+\\.\\d+\\.\\d+\\.\\d+\"\"'
32     t.value = t.value[1:-1] # Remove aspas
33     return t
34
35 def t_DATE_TIME(t):
36     r'\\d{4}-\\d{2}-\\d{2}T\\d{2}:\\d{2}:\\d{2}(\\.\\d+)?Z?'
37     t.value = str(t.value)
38     return t
39
40 def t_DATE(t):
41     r'\\d+\\-\\d+\\-\\d+'
42     return t
43
44 def t_TIME(t):
45     r'\\d+\\:\\d+\\:\\d+'
46     return t
47
48 def t_NUMBER(t):
49     r'-?\\d+(?:\\.\\d+)?(?:[eE][+-]?\\d+)?'
50     try:
51         t.value = int(t.value)
52     except ValueError:
53         t.value = float(t.value)

```

```

54     return t
55
56 def t_COMMENT(t):
57     r'(\#\ s.*)'
58     pass
59     #return t
60
61 def t_STRING(t):
62     r'("[^"]*" | \'[^\']*\' | """"[^^"]*"""" | \'\\\'[^\']*\\\' )'
63     t.value = t.value[1:-1] # Remove aspas
64     return t
65
66 def t_KEY(t):
67     r'([a-z_\d]+)\s'
68     return t
69
70 def t_SUBOBJECT(t):
71     r'[a-zA-Z_\d]*\.[a-zA-Z_\d]+'
72     match = re.match(r'([a-zA-Z_\d]*)\.[a-zA-Z_\d]+', t.value)
73     if match:
74         t.value = match.group(2)
75     return t
76
77 def t_OBJECT(t):
78     r'[a-zA-Z_\d]+'
79     return t
80
81 t_LBRACKET = r'\['
82 t_RBRACKET = r'\]'
83 t_EQUALS = r'='
84 t_COMMA = r','

```

Após o desenvolvimento do analisador léxico, pasamos para a elaboração de um analisador sintático.

2.4 Gramática

É importante referir que se trata de uma gramática LR, ou seja, *Bottom-up*, de forma a podermos usufruir ao máximo das funcionalidades do `PLY`.

De referir, também, que todos os valores que tiverem em letra maiúscula são símbolos terminais (apresentados anteriormente no lex) e os restantes são não terminais.

```

1
2 # P1 : FT : Section_list
3 # P2 : Section_list : Section section_list
4 # P3 : | Section
5 # P4 : Section : '[' Section_name ']' Section_content
6 # P5 : | '[' Section_name ']' Subsection
7 # P6 : | Content
8 # P7 : Table : LBRACKET LBRACKET Section_name RBRACKET RBRACKET
9 # P8 : | LBRACKET LBRACKET Section_name RBRACKET RBRACKET
10 # P9 : Subsection : '[' Subsection_name ']' Section_content Subsection
11 # P10 : | '[' Subsection_name ']' Section_content
12 # P11 : Section_name : OBJECT
13 # P12 : Subsection_name : SUBOBJECT
14 # P13 : Section_content : Content Section_content
15 # P14 : | Content
16 # P15 : | COMMENT
17 # P16 : Content : Key '=' Value
18 # P17 : Key : KEY
19 # P18 : Value : BASIC_STRING
20 # P19 : | LITERAL_STRING
21 # P20 : | STRING
22 # P21 : | NUMBER
23 # P22 : | IP
24 # P23 : | DATE
25 # P24 : | TIME
26 # P25 : | DATE.TIME
27 # P26 : | BOOLEAN
28 # P27 : | Array
29 # P28 : Array : '[' Value_list ']'
30 # P29 : Value_list : Value ',' Value_list
31 # P30 : | Value

```

2.4.1 Programa principal

Primeiramente, decidimos começar o nosso programa como uma lista de secções, que por sua vez, podem-se dividir em apenas uma secção ou uma secção seguida de mais secções, de forma recursiva.

O nosso programa tem, obrigatoriamente, uma secção, pois o código toml apresentado nunca vai ser vazio.

```

1
2 def p_FT(p):
3     "FT : Section_list"
4     p[0] = p[1]
5
6 def p_Section_list1(p):
7     "Section_list : Section Section_list"
8     p[0] = {**p[1], **p[2]}
9
10 def p_Section_list2(p):

```



```
11 "Section_list : Section"
12 p[0] = p[1]
```

2.4.2 Secções

Uma secção pode ser dividida de quatro formas: nome da secção e respetivo conteúdo , nome da secção seguida de uma subsecção, apenas conteúdo (quando não temos referência a um objeto json, isto é, apenas um valor json) ou uma tabela. Aqui quando nos referimos a uma secção estamos a falar de um objeto json.

Como podemos ver pelas produções, todos os valores já estão a ser guardados num formato de dicionário para, mais tarde, facilitar a conversão.

```
1
2 def p_Section1(p):
3     "Section : LBRACKET Section_name RBRACKET Section_content"
4     p[0] = { p[2]: p[4] }
5
6 def p_Section2(p):
7     "Section : LBRACKET Section_name RBRACKET Subsection"
8     p[0] = { p[2]: p[4] }
9
10 def p_Section3(p):
11     "Section : Content"
12     p[0] = p[1]
13
14 def p_Section4(p):
15     "Section : Table"
16     p[0] = p[1]
```

2.4.3 Subsecções

Aqui quando nos referimos a subsecções estamos a falar de objetos dentro de objetos (json). Neste caso, uma subsecção é composta por um nome e o seu conteúdo. Temos também a primeira produção que vai fazer isto de forma recursiva de forma a apanhar todas as subsecções disponíveis no código.

```
1
2 def p_Subsection1(p):
3     "Subsection : LBRACKET Subsection_name RBRACKET Section_content Subsection"
4     p[0] = {p[2]: p[4]}
5     p[0].update(p[5])
6
7 def p_Subsection2(p):
8     "Subsection : LBRACKET Subsection_name RBRACKET Section_content"
9     p[0] = { p[2]: p[4] }
```

2.4.4 Tabela

Com tabela estamos a referir a tabelas aninhadas, e aqui temos duas produções, a segunda (mais simples) vai apenas retirar o nome da tabela e o seu conteúdo, já a primeira é responsável por fazer o mesmo da segunda, retirar nome e conteúdo da tabela, e ainda verificar se a tabela que vem a seguir (caso exista) tem o mesmo nome da anterior e, em caso positivo, vai apenas unir os seus conteúdos num array e guardar no respetivo dicionário.

```

1
2 def p_Table1(p):
3     "Table : LBRACKET LBRACKET Section_name RBRACKET RBRACKET Section_content Table"
4     keys_list = list(p[7].keys())
5     if p[3] == keys_list[0]:
6         p[0] = {p[3]: [p[6]]+p[7][p[3]]}
7     else:
8         p[0] = {p[3]: [p[6]]}
9         p[0].update(p[7])
10
11 def p_Table2(p):
12     "Table : LBRACKET LBRACKET Section_name RBRACKET RBRACKET Section_content"
13     p[0] = { p[3]: [p[6]] }

```

2.4.5 Nome da secção ou subsecção

Nestas produções apenas vamos retirar o nome das secções e subsecções, através do uso do respetivo *token*.

```

1
2 def p_Section_name(p):
3     "Section_name : OBJECT"
4     p[0] = p[1]
5
6 def p_Subsection_name(p):
7     "Subsection_name : SUBOBJECT"
8     p[0] = p[1]

```

2.4.6 Conteúdo da secção

Na parte do conteúdo da secção temos duas alternativas, ou este é um comentário que mais tarde vai ser ignorado ou é uma linha de conteúdo.

```

1
2 def p_Section_content1(p):
3     "Section_content : Content Section_content"
4     p[0] = {**p[1], **p[2]}
5
6 def p_Section_content2(p):
7     "Section_content : Content"
8     p[0] = p[1]
9
10 def p_Section_content3(p):
11     "Section_content : COMMENT"
12     p[0] = p[1]

```

2.4.7 Conteúdo

Uma linha de conteúdo é, obrigatoriamente, representada por uma chave e um valor.

```

1

```

```

2 def p_Content(p):
3     "Content : Key EQUALS Value"
4     p[0] = { p[1] : p[3] }

```

2.4.8 Chave

Esta produção, tal como a do nome da secção ou subsecção, apenas vai verificar o respetivo *token*.

```

1
2 def p_Key(p):
3     "Key : KEY"
4     p[0] = p[1]

```

Aqui temos a produção responsável por captar todos os valores de uma linha *toml*. A função *p_Value* vai verificar todos os *tokens* que tenham a ver com possíveis valores de uma linha de código *toml*. Temos, ainda, uma linha responsável por verificar a existência de um *Array*

2.4.9 Valores

```

1
2 def p_Value(p):
3     '''Value : BASIC_STRING
4               | LITERAL_STRING
5               | STRING
6               | NUMBER
7               | IP
8               | DATE
9               | TIME
10              | BOOLEAN
11              | Array
12              | DATE.TIME
13             '''
14     p[0] = p[1]

```

2.4.10 Arrays

Estas produções são responsáveis por verificar a existência de um *array*. A primeira vai retirar os valores que se encontram dentro de um *array* e a segunda é responsável por percorrer o *array* valor a valor.

```

1
2 def p_Array(p):
3     "Array : LBRACKET Value_list RBRACKET "
4     p[0] = p[2]
5
6 def p_Value_list(p):
7     '''Value_list : Value COMMA Value_list
8                  | Value'''
9     if len(p) > 2:
10        p[0] = [p[1]] + p[3]
11    else:
12        p[0] = [p[1]]

```

2.5 Controlo de Erros

De modo a controlarmos os erros que apareciam ao longo da realização do projeto, testámos escrever diversos exemplos de várias formas, de modo a conseguirmos que a nossa gramática ficasse mais "universal". Alguns exemplos foram extraídos da internet e também utilizamos o que nos foi dado no enunciado, efetuando o grupo várias alterações ao mesmo.

Por outro lado, de forma a controlar os erros do nosso compilador fomos criando funções que detetassem os mesmos e nos devolvessem o erro.

Ao contruir o lex, criamos a função *t_error* de forma a devolver um erro caso não consiga fazer o lexer.

```
1 def t_error(t):
2     print("Erro: " + t.value[0], file=sys.stderr)
3     t.lexer.skip(1)
```

De seguida, ainda com o mesmo pensamento do caso anterior, no yacc criamos a função *p_error* que nos devolve erro caso não consiga interpretar alguma linha do nosso código.

```
1 def p_error(p):
2     parser.success = False
3     print(f"Syntax error: {p.value}")
```

Por fim, para não estarmos a trabalhar com códigos *toml* mal formatados criamos estas linhas de código de forma a ler o ficheiro introduzido para converter e detetar logo se há erros de formatação ou não, devolvendo erro em caso afirmativo.

```
1 try:
2     parsed_data = toml.loads(data)
3     print("O arquivo TOML esta formatado corretamente.")
4
5 except toml.TomlDecodeError as e:
6     print("Erro ao decodificar o arquivo TOML:")
7     print(e)
```

Capítulo 3

Funcionalidades implementadas

Como pedido no trabalho prático, a nossa primeira funcionalidade a ser implementada foi o conversor de *toml* para *json*. No entanto, uma vez que ao longo das produções da gramática fomos guardando todos os valores num dicionário *python*, reparamos que, com isso, eram facilmente construídos outros conversores para outras linguagens, como por exemplo, *yaml* e *xml*

Em todos os exemplos que se vão apresentar de seguida, o ficheiro input é o exemplo apresentado no enunciado e *parsed_dict* é o dicionário resultante do analisador sintático.

3.1 Conversor *toml* para *json*

Esta foi a funcionalidade pedida no trabalho prático, logo foi a primeira a ser feita. Como referido anteriormente, os nossos valores estão todos armazenados num dicionário *python*, logo para converter o nosso dicionário para um ficheiro *json* apenas temos de adicionar a seguinte linha de código que vai dar *dump* ao nosso dicionário para o ficheiro *output*.

```
1 json.dump(parsed_dict, f, ensure_ascii=False)
```

3.1.1 Ficheiro output

```
1 {
2   "title ": "TOML Example",
3   "owner": {
4     "name ": "Tom Preston-Werner",
5     "date  ": "2010-04-23",
6     "time  ": "21:30:00"
7   },
8   "database": {
9     "server ": "192.168.1.1",
10    "ports  ": [
11      8001,
12      8001,
13      8002
14    ],
15    "connection_max ": 5000,
16    "enabled ": true
17  },
```

```

18     "servers": {
19         "alpha": {
20             "ip ": "10.0.0.1",
21             "dc ": "eqdc10"
22         },
23         "beta": {
24             "ip ": "10.0.0.2",
25             "dc ": "eqdc10",
26             "hosts ": [
27                 "alpha",
28                 "omega"
29             ]
30         }
31     }
32 }

```

3.2 Conversor *toml* para *yaml*

Esta foi uma funcionalidade extra que o grupo decidiu implementar. Aqui é percorrido o dicionário com a tradução do *toml*, que vai ser convertido para *yaml* com a seguinte linha de código:

```

1 yaml_data = yaml.dump(parsed_dict)

```

3.2.1 Ficheiro output

```

1
2 database:
3   'connection_max ': 5000
4   'enabled ': true
5   'ports ':
6     - 8001
7     - 8001
8     - 8002
9   'server ': 192.168.1.1
10 owner:
11   'date ': '2010-04-23'
12   'name ': Tom Preston-Werner
13   'time ': '21:30:00'
14 servers:
15   alpha:
16     'dc ': eqdc10
17     'ip ': 10.0.0.1
18   beta:
19     'dc ': eqdc10
20     'hosts ':
21       - alpha
22       - omega
23     'ip ': 10.0.0.2
24 'title ': TOML Example

```

3.3 Conversor *toml* para *xml*

Igual à funcionalidade anterior, esta foi uma funcionalidade extra que o grupo decidiu implementar. Aqui é percorrido o dicionário com a tradução do *toml*, que vai ser convertido para *xml* com as seguintes linhas de código:

```
1 def create_xml_element(parsed_dict, parent):
2     for key, value in parsed_dict.items():
3         if isinstance(value, dict):
4             element = ET.SubElement(parent, key)
5             create_xml_element(value, element)
6         else:
7             element = ET.SubElement(parent, key)
8             element.text = str(value)
9
10 root = ET.Element('root')
11
12 create_xml_element(parsed_dict, root)
13
14 # Criar a rvore XML
15 tree = ET.ElementTree(root)
```

3.3.1 Ficheiro output

```
1
2 <?xml version='1.0' encoding='utf-8'?>
3 <root>
4     <title >TOML Example</title >
5     <owner>
6         <name >Tom Preston-Werner</name >
7         <date >2010-04-23</date >
8         <time >21:30:00</time >
9     </owner>
10    <database>
11        <server >192.168.1.1</server >
12        <ports >[8001, 8001, 8002]</ports >
13        <connection_max >5000</connection_max >
14        <enabled >True</enabled >
15    </database>
16    <servers>
17        <alpha>
18            <ip >10.0.0.1</ip >
19            <dc >eqdc10</dc >
20        </alpha>
21        <beta>
22            <ip >10.0.0.2</ip >
23            <dc >eqdc10</dc >
24            <hosts >['alpha', 'omega']</hosts >
25        </beta>
26    </servers>
27 </root>
```

3.4 Funcionamento do programa

Ao correr o programa com *python3 converter.py*, o utilizador vai ser pedido para introduzir o path do ficheiro que quer converter e, de seguida, qual é a linguagem para a qual deseja converter o seu ficheiro. Neste processo o programa vai verificar se o ficheiro existe, se está num formato *toml* correto e se a linguagem para a qual o utilizador deseja converter é suportada.

Capítulo 4

Conclusão

Finalizada a realização deste projeto desenvolvido ao longo do semestre, sentimos que pudemos aprofundar os conhecimentos e aspetos adquiridos ao longo das aulas da respetiva UC, nomeadamente o uso das ferramentas *Yacc* e *Lex* do *PLY*.

Ao longo da realização do presente trabalho prático, deparamo-nos com algumas dificuldades, mas rapidamente superamos, analisando os erros e tentando perceber o porquê de estarem a acontecer.

Uma das partes mais desafiadoras deste trabalho foi a construção do analisador sintático. Durante o processo, o grupo deparou-se com vários erros e também encontrou melhores soluções para o desenvolvimento da gramática. Isto exigiu várias tentativas de escrita até chegarmos a uma forma mais correta.

Em jeito de conclusão, finalizado este trabalho, o grupo está satisfeito com o produto final, considerando ir de encontro ao que foi solicitado.