

Processamento de Linguagens

MiEI (3ano)

Trabalho Prático nº 2 (YACC)

Ano lectivo 19/20

1 Objectivos e Organização

Este trabalho pratico tem como principais **objectivos**:

- aumentar a experiência de uso do ambiente Linux, da linguagem imperativa C (para codificação das estruturas de dados e respectivos algoritmos de manipulação), e de algumas ferramentas de apoio à programação;
- rever e aumentar a capacidade de escrever *gramáticas independentes de contexto (GIC)*, que satisfaçam a condição LR(), para criar Linguagens de Domínio Específico (DSL);
- desenvolver processadores de linguagens segundo o método da *tradução dirigida pela sintaxe*, suportado numa *gramática tradutora (GT)*;
- utilizar *geradores de compiladores* como o par *flex/yacc*

Para o efeito, esta folha contém 6 enunciados, dos quais deverá resolver um escolhido em função do número do grupo (NGr), usando a fórmula $exe = (NGr \% 6) + 1$.

Notem que podem escolher outro enunciado por troca com outro grupo que esteja de acordo.

Neste TP, que se pretende que seja resolvido com agilidade, os resultados pedidos são simples e curtos. Aprecia-se a imaginação/criatividade dos grupos ao incluir outros processamentos!

Deve entregar a sua solução **até Domingo dia 7 de Junho**. O ficheiro zipado a entregar (com relatório, ficheiros Lex, Yacc e C) deve ter o nome 'p119TP2NA11-NA12-NA13' em que NALx é o numero de cada aluno do grupo — *perto da data de entrega serão dadas indicações precisas sobre a forma de submissão e calendário de avaliações a decorrer ao longo de Junho*.

O **relatório** a elaborar, deve ser claro e, além do respectivo enunciado, da descrição do problema, das decisões que lideraram o desenho da solução e sua implementação (incluir a especificação FLEX e Yacc em Apêndice), deverá conter exemplos de utilização (textos fontes diversos e respectivo resultado produzido). Como é de tradição, o relatório será escrito em L^AT_EX (ver template no Bb).

2 Enunciados

Para sistematizar o trabalho que se lhe pede em cada uma das propostas seguintes, considere que deve, em qualquer um dos casos, realizar a seguinte lista de tarefas:

1. Especificar a gramática concreta da linguagem de entrada.
2. Desenvolver um reconhecedor léxico e sintáctico para essa linguagem com recurso ao par de ferramentas geradoras Flex/Yacc.
3. Construir o gerador de código que produza a resposta solicitada. Esse gerador de código é construído associando acções semânticas de tradução às produções da gramática, recorrendo uma vez mais ao gerador Yacc.

Conteúdo

1	Objectivos e Organização	1
2	Enunciados	1
2.1	Reverse Engineering dum Dicionário financeiro	2
2.2	DSL para Cadernos de Anotações em HD	3
2.3	Conversor de Pug ⁺⁺ para HTML	4
2.4	Conversor toml2json	5
2.5	Conversor de GEDCOM	6
2.6	Visualizador de árvores genealógicas	7

2.1 Reverse Engineering dum Dicionário financeiro

Em `natura.di.uminho.pt/~jj/pl-20/TP2` está o texto de um dicionário EN-PT. Considere o seguinte extrato de um dicionário de negócio e finanças EN-PT:

```
1 dispute:
2   labour -          conflito (m) trabalhista
3 dissolution        dissolução (f)
4 distribution:
5   - costs           custos de distribuição
6   channels of -     canais (mpl) de distribuição
7   physical - management controle (m) da distribuição física
8 distribuidor        distribuidor (m), atacadista (m)
9 diversification:    diversificação (f)
10  - strategy         estratégia (f) de diversificação
11 product -          diversificação (f) de produtos
```

Este ficheiro de partida foi resultado da conversão automática PDF-TXT- Pretendemos agora fazer *reverse engineering* de modo a tirar o máximo partido da informação.

No contexto deste TP pede-se que:

- Escreva uma gramática que cubra tanto quanto possível as entradas do dicionário.
- Construa um processador (flex, yacc) que converta o ficheiro num formato mais tratável—ver abaixo um exemplo de saída pretendida.
- Detete e trate os casos de erro sintático, nas situações em que o formato da entrada seja incoerente com a gramática criada acima; nessas circunstâncias, avise a linha onde ocorreu o erro, e continue o processamento.

```
1 EN labour dispute
2 +base dispute
3 PT conflito (m) trabalhista
4
5 EN dissolution
6 PT dissolução (f)
7
8 EN distribution costs
9 +base distribution:
10 PT custos de distribuição
11
12 EN channels of distribution
13 +base distribution:
14 PT canais (mpl) de distribuição
```

```

12 EN physical distribution management
13 +base distribution:
14 PT controle (m) da distribuição física

15 EN distributor
16 PT distribuidor (m)
17 PT atacadista (m)

18 EN diversification
19 PT diversificação (f)

20 EN diversification strategy
21 +base diversification
22 PT estratégia (f) de diversificação

23 EN product diversification
24 +base diversification
25 PT diversificação (f) de produtos

```

2.2 DSL para Cadernos de Anotações em HD

Para guardar informação no domínio das *Humanidade Digitais*, pretendemos criar um formato e um processador de *cadernos de anotações*.

Cada *caderno* é composto por diversos pares (`doc`, `triplos`). Cada documento inclui um tópico principal (conceito), um título e um texto. Os triplos seguem a notação turtle (versão simplificada), o que corresponde a (`sujeito`, `relação`, `objecto`) em que `sujeito` e `predicado` são de novo conceitos.

O texto a seguir é um fragmento de um *caderno de anotações*, conforme pretendido.

```

1 === :OnorioL
2 @tit: Onório de Lima

3 # Infância no Brasil

4 Onório de Lima viveu ...

5 # Percurso cultural no Porto

6 Entre 1873 e 1930 Onório de Lima viveu no Porto, Cedofeita e foi...

7 ...parágrafos separados por linhas em branco

8 @triplos:
9 :OnórioL
10   a      :Person ,
11         :Mecenas_cultural ;
12   :Name  "Onório de Lima" ;
13   :birth_place :Brasil ;
14   :viveu_em :Porto,
15           :Gerês ;
16   :img "OnórioDL1.jpg" .

17 :OnórioL
18   :pai_de :OndinaL .

19 :CarolinaDA
20   :filho_de :Ontinal;

```

```

21   :name "Carolina Dias Antunes".
22   :casou_com :JorgeDomingos

23 === :JorgeDomingos
24 @tit: António Jorge Domingos
25 ...
26 @triplos:
27 ...

```

No contexto deste TP pede-se que:

- Escreva uma gramática para definir a DSL de *caderno de anotação*, acima descrita.
- Construa um processador (flex, yacc) que reconheça e valide cadernos de anotação, escritos na DSL definida acima, gerando um site HTML de acordo com as seguintes regras:
 - crie uma página por cada elemento (conceito) – tópico, sujeito, ou objeto (:OnórioL, :Person, :Mece-nas_cultural, :Brasil, etc) que agrupe a informação recolhida, e que hiperlink com tudo o que for possível.
 - inclua na página HTML de um sujeito :x uma imagem "la.jpg" sempre que encontrar no *caderno* um triplo da forma (:x :img "la.jpg").

Extra: permita a existência de uma secção *meta* onde se indiquem triplos que definam relações inversas (Ex (:filho_de :inverseOf :pai_de)) de modo a ser possível a inferência de triplos inversos.

2.3 Conversor de Pug⁻⁻ para HTML

Pug (inicialmente designado por Jade) é um *template engine* que oferece uma notação simples para escrever código HTML reutilizável (permite ter partes do documento que serão preenchidas com valores extraídos de bases de dados ou outras fontes) e que gera HTML puro para mostrar essa informação em páginas Web.

Poe exemplo, o documento a seguir, escrito em notação Pug

```

1  html(lang="en")
2    head
3      title= pageTitle
4      script(type='text/javascript').
5        if (foo) bar(1 + 5)
6    body
7      h1 Pug - node template engine
8      #container.col
9        if youAreUsingPug
10         p You are amazing
11        else
12         p Get on it!
13        p.
14          Pug is a terse and simple templating language with a
15          strong focus on performance and powerful features

```

será transformado no documento HTML abaixo:

```

1  <html lang="en">
2    <head><title></title>
3      <script type="text/javascript">
4        if (foo) bar(1 + 5)</script>
5    </head>
6    <body>
7      <h1>Pug - node template engine</h1>
8      <div class="col" id="container">

```

```

9         <p>Get on it!</p>
10         <p>Pug is a terse and simple templating language with a
11 strong focus on performance and powerful features</p>
12     </div>
13 </body>
14 </html>

```

Neste trabalho pretende-se que:

- Escreva uma gramática que cubra um subconjunto da linguagem Pug à sua escolha (e claramente identificado no relatório de desenvolvimento do TP). Para fazer a escolha veja a documentação em:
<https://pugjs.org/language/tags.html>
<https://pugjs.org/language/plain-text.html>
<https://pugjs.org/language/attributes.html>
- Construa um processador (flex, yacc) que reconheça e valide os documentos Pug⁺⁺, escritos na sintaxe definida acima, gerando o HTML correspondente.

2.4 Conversor toml2json

TOML¹ é uma linguagem simples, fácil de ler e escrever, para descrever estruturas complexas (usada frequentemente em ficheiros de configuração), desenhada para ser mapeada diretamente e sem ambiguidades para dicionários generalizados. TOML, cujo nome significa *Tom's Obvious, Minimal Language* devido ao seu criador *Tom Preston-Werner*, é usada em muitos projetos de software sendo equiparável a JSON ou YAML.

Pode ver mais detalhes em: <https://github.com/toml-lang/toml>

Para ilustrar esta proposta de trabalho, considere o seguinte exemplo:

```

1  title = "TOML Example"
2
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 = [ 8001, 8001, 8002 ]
11 connection_max = 5000
12 enabled = true
13
14 [servers]
15 [servers.alpha]
16 ip = "10.0.0.1"
17 dc = "eqdc10"
18
19 [servers.beta]
20 ip = "10.0.0.2"
21 dc = "eqdc10"
22
23 # Line breaks are OK when inside arrays
24 hosts = [
25     "alpha",
26     "omega"
27 ]

```

¹Consulte <https://en.wikipedia.org/wiki/TOML>

No contexto deste TP pede-se que:

- Escreva uma gramática que cubra um subconjunto da linguagem TOML à sua escolha (e claramente identificado no relatório de desenvolvimento do TP).
- Construa um processador (flex, yacc) que reconheça e valide estruturas/dicionários escritos na DSL TOML² definida acima, gerando o JSON correspondente.

2.5 Conversor de GEDCOM

A Genealogia tem desde sempre atraído muitas pessoas. Investigadores, historiadores, membros da nobreza, pessoas comuns curiosas ou que, por um motivo ou outro, precisam de construir a sua árvore genealógica.

Ao longo dos últimos anos, o formato GEDCOM ganhou popularidade e uma alargada comunidade de utilizadores, o que faz com que seja o formato de eleição para a transmissão e intercâmbio de informação nesta área.

Na internet é possível encontrar muitos recursos e informação sobre este formato.

Neste projeto, pretende-se que o grupo de trabalho desenvolva um conversor de GEDCOM 5.5 para XML/HTML de modo a poder mostrar de imediato numa página Web a árvore genealógica especificada na notação em causa.

Sugerem-se as seguintes etapas para o desenvolvimento deste projeto:

1. Estudar o formato, sua estrutura e alguns casos de estudo (prepare os seus próprios exemplos). Para isso pode consultar os seguintes links:

<http://homepages.rootsweb.com/~pmcbride/gedcom/55gcch1.htm> – Contem uma descrição do formato juntamente com um rascunho da gramática;

<http://www4.di.uminho.pt/~jcr/AULAS/didac/ontologias/BibleFamilyTree.ged.txt> – Dataset com as famílias da Bíblia;

<http://www4.di.uminho.pt/~jcr/AULAS/didac/ontologias/GreekGods.ged.txt> – Dataset com as famílias da mitologia grega;

<http://www4.di.uminho.pt/~jcr/AULAS/didac/ontologias/RomanGods.ged.txt> – Dataset com as famílias da mitologia romana;

<http://www4.di.uminho.pt/~jcr/AULAS/didac/ontologias/Royal92.ged.txt> – Dataset com as famílias da realeza europeia.

2. Escrever uma gramática para o GEDCOM 5.5
3. Definir o formato de saída em XML/HTML². Como fonte de inspiração, apresenta-se a seguir um possível formato de saída³:

```
1 <genoa>
2   <essoa>
3     <id>I1</id>
4     <nome>Victoria Hanover</nome>
5     <título>Queen of England</título>
6     <sexo>F</sexo>
7     <dataNasc>1819-05-24</dataNasc>
8     <localNasc>Kensington,Palace,London,England</localNasc>
9     <dataÓbito>1901-01-22</dataÓbito>
10    <pai>I133</pai>
11    <mae>I138</mae>
12  </essoa>
13  ...
14 </genoa>
```

4. Construir um processador (flex, yacc) que reconheça e valide árvores genealógicas escritas na DSL GEDCOM definida acima, gerando o XML/HTML pretendido de acordo com o passo anterior.

²Peça orientação ao professor JCRamalho.

³Veja pormenores em <http://www4.di.uminho.pt/~jcr/AULAS/didac/ontologias/royalFamilies.xml>

2.6 Visualizador de árvores genealógicas

A Genealogia tem desde sempre atraído muitas pessoas. Investigadores, historiadores, membros da nobreza, pessoas comuns curiosas ou que, por um motivo ou outro, precisam de construir a sua árvore genealógica.

Nas aulas de Processamento e Representação de Conhecimento definiu-se uma linguagem para descrever ontologias (notação ilustrada abaixo) e desenhou-se uma ontologia muito simples para a representação de famílias. Essa ontologia é formada pelos *conceitos* *Pessoa* e *Sexo*, e pelas *relações* *temPai* e *temMãe*.

Neste projeto, pretende-se que o grupo de trabalho desenvolva um conversor de ontologias escritas nessa linguagem para grafos escritos em DOT, que depois podem ser visualizados em GraphViz ou convertidos em SVG para serem usados em páginas Web.

A seguir apresenta-se um breve excerto de uma das ontologias que serão usadas como material de teste e análise:

```
1 #####
2 #   Individuals
3 #####
4
5 ### http://www.di.uminho.pt/prc/ontologies/2020/prc-genoa#Abilio_Silva_Ramalho
6 :Abilio_Silva_Ramalho rdf:type owl:NamedIndividual ,
7                        :Pessoa ;
8                        :temMae :Custodia_Azevedo_1867 ;
9                        :temPai :Manuel_Silva_Ramalho_1866 .
10
11 ### http://www.di.uminho.pt/prc/ontologies/2020/prc-genoa#Albina_Esteves_Araujo_1910
12 :Albina_Esteves_Araujo_1910 rdf:type owl:NamedIndividual ,
13                               :Pessoa ;
14                               :temMae :Maria_Araujo_1884 ;
15                               :temPai :Henrique_Luiz_Araujo_1867 .
16 ...
```

Seguem-se algumas explicações sobre a notação usada acima:

- Linhas iniciadas por "#" são comentários e devem ser ignoradas;
- Uma ontologia nesta forma é uma lista de triplos: (sujeito, predicado, objeto);
- O sujeito é sempre um identificador (URI) dum indivíduo;
- O predicado é sempre um identificador (URI) duma relação ou propriedade;
- O objeto pode ser um identificador (URI) dum indivíduo ou conceito, ou ainda um valor escalar: string, número, ...
- Um triplo normal termina por ".";
- Um triplo terminado por "," significa que o próximo triplo está abreviado e vai partilhar o mesmo sujeito e predicado;
- Um triplo terminado por ";" significa que o próximo triplo está abreviado e vai partilhar o mesmo sujeito.

Sugerem-se as seguintes etapas para o desenvolvimento deste projeto:

1. Estudar o formato, sua estrutura e alguns casos de estudo (prepare os seus próprios exemplos). Para isso pode consultar o seguinte link:

<http://www4.di.uminho.pt/~jcr/AULAS/didac/ontologias/mini-familia.ttl> – Contem o exemplo completo do excerto acima;

Outros – O docente irá fornecer mais exemplos.

2. Escrever uma gramática para a notação ilustrada para este tipo de ontologias.
3. Estudar e definir o formato de saída: <https://www.graphviz.org/> ou https://graphviz.gitlab.io/_pages/pdf/dotguide.pdf
4. Construir um processador (flex, yacc) que reconheça e valide ontologias genealógicas escritas na DSL definida acima, gerando o DOT pretendido de acordo com o passo anterior.
5. Seria interessante incluir relações de parentesco inferidas a partir da informação dada: `temAvô`, `temIrmão`, `temTio...`