

# TEI Lex-0: a good fit for the encoding of the Portuguese Academy Dictionary?

<sup>1</sup> Ana Salgado, <sup>1</sup> Rute Costa, <sup>2</sup> Toma Tasovac

<sup>1</sup> NOVA CLUNL, Faculdade de Ciências Sociais e Humanas, Universidade NOVA de Lisboa, Portugal

<sup>2</sup> Belgrade Center for Digital Humanities, Serbia

[anasalgado@campus.fcsh.unl.pt](mailto:anasalgado@campus.fcsh.unl.pt)

[rute.costa@fcsh.unl.pt](mailto:rute.costa@fcsh.unl.pt)

[ttasovac@humanistika.org](mailto:ttasovac@humanistika.org)



# Acknowledgements



This research has been financed by:

- Portuguese National Funding through the FCT – *Fundação para a Ciência e Tecnologia* as part of the project Centro de Linguística da Universidade NOVA de Lisboa – UID/LIN/03213/2019
- European Union’s Horizon 2020 research and innovation programme under grant agreement No 731015 (ELEXIS)



# Outline

1. Motivation and goals
2. Portuguese Academy Dictionary (DLPC)
3. TEI Guidelines for Dictionary encoding
4. TEI Lex-0
5. TEI Lex-0 encoding of the DLPC
6. Conclusions and future work

# 1. Motivation

## A new PORTUGUESE ACADEMY DICTIONARY

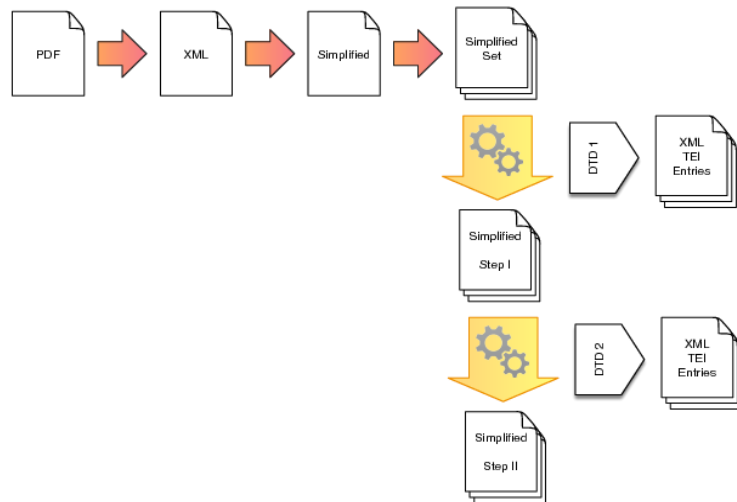
- Preservation
- Updating (of the content)
- Accessibility
- Interoperability
- Reusability
- Linguistic analysis

# 1. ... and Goals

- to refer the conversion of TEI P5 schema into TEI Lex-0
- to demonstrate the application of TEI Lex-0 to mark up 'special entries' in DLPC
- to present arguments in favour of using TEI Lex-0 dictionary encoding
- to contribute to the efforts of the TEI Lex-0 group (DARIAH-ERIC Lexical Resources group)

# 2. Portuguese Academy Dictionary

DLPC = *Dicionário da Língua Portuguesa Contemporânea (2001)*  
69 426 entries, 167 556 senses



Dicionário da Academia das Ciências de Lisboa

Páginas ▾ palavra 🔍

Condensado/Expandido

Estado: Importado

**dicionário**  
*n. m.*

1. Livro de referência em que se fornecem informações, como a categoria gramatical, as aceções, os registos, a forma correspondente noutra língua..., sobre palavras e expressões de uma língua, apresentando-as de acordo com uma ordem convencional, geralmente alfabética. *Consultou vários dicionários de língua portuguesa. Um dicionário com 50.000 entradas. Os artigos, os verbetes de um dicionário. Dicionário informatizado.*  
«O dicionário, imagem ordenada do mundo, constrói-se e desenvolve-se sobre tantíssimas palavras que viveram uma vida plena» ( Público, 6.12.1992)

**dicionário bilingue**  
o que apresenta a tradução das palavras e respectivas aceções de uma língua para outra. *Dicionário bilingue de português/francês.*

**dicionário electrónico**  
o que tem um suporte informático, com um ou vários discos compactos de grande capacidade, designado por CD-ROM.

**dicionário enciclopédico**  
aquele que, além das definições de palavras, inclui artigos desenvolvidos de carácter científico, técnico, histórico...

**dicionário inverso**  
aquele em que as palavras estão ordenadas alfabeticamente a partir do fim.

**dicionário monolíngue**  
o que apresenta a descrição do léxico de uma só língua.

**dicionário multilíngue**  
o que apresenta correspondência termo a termo entre mais de duas línguas.

2. Livro que reúne um conjunto de palavras seleccionadas de acordo com áreas temáticas, zonas geográficas em que são usadas, peculiaridades da língua... + de *medicina, pintura; +s de regionalismos, de calão, de sinónimos, de antónimos; + etimológico; + de verbos, de citações, de provérbios.*

**dicionário analógico**

1. O que parte de uma selecção de conceitos, que constituem as entradas, sob os quais agrupa o vocabulário que lhes corresponde, associando as palavras de acordo com as analogias de sentido, e que compreende um índice final onde são indexadas alfabeticamente todas as palavras constantes dos artigos, com indicação dos conceitos sob os quais figuram. *Nos dicionários analógicos, o léxico é encarado do ponto de vista da sua estruturação semântica.*

2. Aquele que, apresentando as palavras por ordem alfabética, consegue estabelecer entre elas relações de analogia, no plano do conteúdo, pela inclusão de sinónimos e antónimos e também de remissão para outros termos pertencentes aos mesmos campos semânticos. *É um dicionário ao mesmo tempo descritivo e analógico: o seu sistema de remissões leva o utilizador a descobrir palavras desconhecidas.*

3. Conjunto de palavras usadas habitualmente por um grupo social ou por uma pessoa individualmente. *Esse termo não consta do meu dicionário.*

**dicionário vivo**  
pessoa muito culta, muito erudita. ≈ **enciclopédia**  
(Do lat. medieval *dictionarium*, do lat. *dictio*, -ōnis 'palavra')

Remover Editar /db/academia/dicionário.xml

# 3. TEI Guidelines for Dictionary encoding

## TEI dictionary encoding

- TEI is a *de facto* standard in digital edition or text annotation projects, from novels, letters to music notation.
- It is frequently used in Digital Humanities as the basis for a large number of current lexicographic projects .
- TEI has a specific module for encoding dictionaries.

BUT...

- We could not find solutions in the Guidelines that covered all the microstructural elements of the dictionary.
- TEI Guidelines contain multiple encoding possibilities for the same dictionary features.

# 4. TEI Lex-0

## TEI Lex-0

Romary & Tasovac (2018).  
TEI Lex-0: A Target Format for  
TEI-Encoded Dictionaries and  
Lexical Resources

- a streamlined version of the TEI Guidelines, simplified and enhanced for regular use to improve interoperability
- given its (still) non-standard nature, it can be changed in order to accommodate relevant dictionary structures
- we have been participating in the TEI Lex-0 discussion during the Portuguese Academy Dictionary encoding: <https://github.com/DARIAH-ERIC/lexicalresources/tree/master/Schemas/TEILex0>
- the advantage of applying TEI Lex-0 lies in the fact that lexicographers and terminologists are currently making efforts to align the ongoing revision of ISO LMF with TEI (Romary, 2015)

TEI Lex-0 will not replace the Dictionaries Chapter in the TEI Guidelines. It is being discussed as a **target format** that will facilitate the integration and analysis of the existing heterogeneously encoded lexical resources.



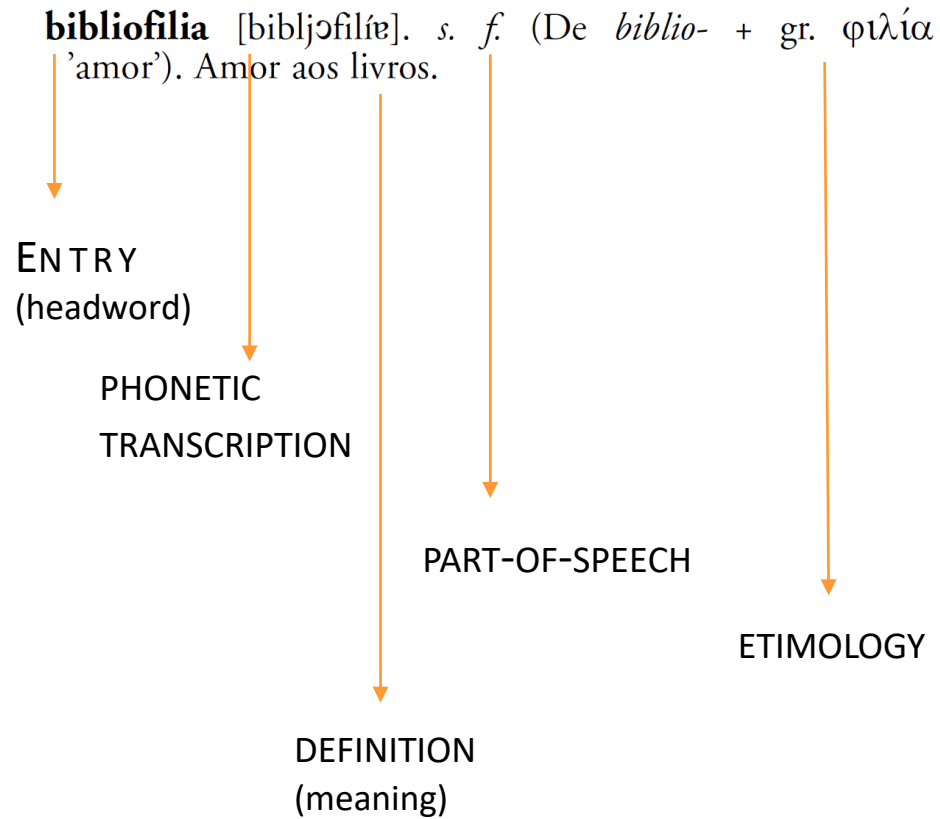
# 5. TEI Lex-0 encoding of the DLPC

## XML ESSENTIAL CHANGINGS (some examples)

Original encoding	Conversion into TEI Lex-0
<entry>	<entry xml:id=" id " xml:lang="pt">
<term>	<form type="lemma">
<gramGrp> part of speech and gender </gramGrp>	<gramGrp><gram type="pos"></gram> <Gram type="gen">f.</gram></gramGrp>
<sense>	<sense xml:id=" id ">
<quote type="example">	<cit type="example"><quote>
<syn> synonym </syn>	<xr type="synonym"><ref type="entry sense"> synonym </ref></xr>
<cit><quote> example </quote> <bibl> author ,<title> title </title>,<page> </bibl></cit>	<cit type="example"><quote> example </quote> <bibl><author> author </author> <title> title </title><citedRange> page </citedRange></bibl></cit>
<sense><def>V.<xr><ref> cross-reference </ref></xr></def></sense>	<xr><lbl>V.</lbl><ref> cross-reference </ref></xr>

# 5. TEI Lex-0 encoding of the DLPC

## Basic entry structure



*bibliofilia* [bibliophilia], DLPC (2001)

```
<entry type="monolexicalWord" xml:lang="pt" xml:id="bibliofilia">
  <form type="lemma">
    <orth>bibliofilia</orth>
    <pron>bibljɔfil'ie</pron>
  </form>
  <gramGrp>
    <gram type="pos" norm="NOUN">s.</gram>
    <gram type="gen">f.</gram>
  </gramGrp>
  <etym type="inheritance">
    <lbl>De</lbl>
    <cit type="etymon">
      <form>
        <orth>biblio-</orth>
      </form>
    <etym type="grammaticalization">
      <seg type="desc">De</seg>
      <cit type="etymon">
        <form>
          <orth extent="pref">biblio-</orth>
        </form>
      </cit> <lbl>+</lbl>
    <etym type="inheritance">
      <seg type="desc">Do</seg>
      <cit type="etymon" xml:lang="grc">
        <form><orth>φιλία</orth></form>
      </cit>
    </etym>
  </etym>
  <sense>
    <def>Amor aos livros</def><pc>.</pc>
  </sense>
</entry>
```

- entry
- @xml:id
  - @xml:lang – code (BCP 47)
- gram type element
- part-of-speech of the entry and further specifications
  - @norm attribute – values from the Universal Dependencies Part-of-Speech: <https://universaldependencies.org/u/pos/>

# 5. TEI Lex-0 encoding of the DLPC

## Mark up of ‘special entries’

- part-of-speech homonyms: *antepassado*<sup>1</sup>, *antepassado*<sup>2</sup>
- entries that have a different meaning in the plural: *antepassados*
- etymological homonyms: *cota*<sup>1</sup>, *cota*<sup>2</sup> (Bowers & Romary (2017). Deep encoding of etymological information in TEI)
- homographs: *lobo*<sup>1</sup> /ó/, *lobo*<sup>2</sup> /ô/
- spelling variants: *ouro*, *oiro*
- lexical variants: *missanga*, *miçanga*
- trademarks: *donut*<sup>®</sup>, *walkman*<sup>®</sup>

**antepassado**<sup>1</sup>, **a** [ẽtipesádu, -v]. *adj.* (De *ante-* + *passado*).

Que pertence ou viveu numa época anterior. ≈ ANTECESSOR, PREDECESSOR. ≠ DESCENDENTE, SUCESSOR.

**antepassado**<sup>2</sup> [ẽtipesádu]. *s. m.* (De *ante-* + *passado*).

**1.** Pessoa que é ascendente de outra ou outras. ≈ ASCENDENTE. ≠ DESCENDENTE. *Certos povos crêem-se descendentes de um antepassado comum. «o vaqueiro, pai do vaqueiro, o avô e outros antepassados mais antigos haviam-se acostumado a percorrer veredas, afastando o mato com as mãos.»* (G. RAMOS, *Vidas Secas*, p. 36). **2. pl.** Pessoas anteriormente ao momento actual. ≈ ANTECESSORES. ≠ VINDOUROS. *Herdámos estes costumes dos nossos antepassados. Culto dos antepassados.*

*antepassado* [ancestor, past], DLPC (2001)

```
<entry type="derivativeWord" xml:lang="pt"
xml:id="antepassado_1" n="1">
  <form type="lemma">
    <orth>antepassado</orth>
  </form>
  <form type="inflected">
    <orth>antepassado</orth>
    <pron>ẽtipes'adu</pron>
  <gramGrp>
    <gram type="gen">m.</gram>
  </gramGrp>
</form>
  <form type="inflected">
    <orth>antepassada</orth>
    <pron>ẽtipes'ade</pron>
  <gramGrp>
    <gram type="gen">f.</gram>
  </gramGrp>
</form>
  <gramGrp>
    <gram type="pos" norm="ADJ">adj.</gram>
  </gramGrp>
  <!--etc. -->
</entry>
```

```
<entry type="derivativeWord" xml:lang="pt"
xml:id="antepassado_2" n="2">
  <form type="lemma">
    <orth>antepassado</orth>
    <pron>ẽtipes'adu</pron>
  </form>
  <gramGrp>
    <gram type="pos" norm="NOUN">s.</gram>
    <gram type="gen">m.</gram>
  </gramGrp>
  <!--etc. -->
  <sense xml:id="antepassado_2.2" n="2">
    <gramGrp>
      <gram type="number">pl.</gram>
    </gramGrp>
    <!--etc. -->
  </sense>
</entry>
```

**cota**<sup>1</sup> [kóte]. *s. f.* (Do fr. ant. *cotte*). 1. *Mil.* Vestimenta ou peça da armadura usada antigamente por cavaleiros e guerreiros, sobre o busto, para os proteger dos golpes do adversário. **cota de armas** ou **cota**. 1. Vestimenta usada pelos cavaleiros sobre o arnês ou a armadura, tanto na guerra como em torneios. 2. Peça de vestuário, onde estava bordado o escudo real, usada pelos reis em actos públicos. **cota de malha(s)**, armadura defensiva, com o formato de uma camisa feita de anéis ou malhas de metal entrelaçadas, usada antigamente para proteger o corpo contra golpes de armas. **fidalgo<sup>+</sup> de cota de armas**. 2. *Mil.* Antiga peça de vestuário, usada para protecção por cavaleiros e guerreiros, que podia chegar ao joelho e apertar-se na cintura. 3. Peça de vestuário parecida com uma bata ou gibão, usada antigamente. 4. Antiga peça de vestuário feminino que consistia numa espécie de corpete. 5. *Rel.* Veste litúrgica, de cor branca, com ou sem rendas, semelhante a uma sobrepeliz mas com mangas mais curtas, envergada sobre o roquete, e usada sobretudo em Itália, por certos dignitários religiosos e acólitos.

**cota**<sup>2</sup> [kóte]. *s. f.* (De origem obscura). Lado de uma ferramenta ou de um instrumento cortante, oposto ao que está afiado, ao gume. ≈ COSTAS.

*cota* [armour, ..., back], DLPC (2001)

```
<entry type="monolexicalWord" xml:lang="pt"
xml:id="cota_1" n="1">
  <form type="lemma">
    <orth>cota</orth>
    <pron>k'ote</pron>
  </form>
  <gramGrp>
    <gram type="pos" norm="NOUN">s.</gram>
    <gram type="gen">f.</gram>
  </gramGrp>
  <etym type="inheritance">
    <seg type="desc">Do</seg>
    <cit type="etymon" xml:lang="fro">
      <form><orth>cotte</orth></form>
    </cit>
  </etym>
  <!--etc. -->
</entry>
```

```
<entry type="monolexicalWord" xml:lang="pt"
xml:id="cota_2" n="2">
  <form type="lemma">
    <orth>cota</orth>
    <pron>k'ote</pron>
  </form>
  <gramGrp>
    <gram type="pos" norm="NOUN">s.</gram>
    <gram type="gen">f.</gram>
  </gramGrp>
  <etym type="undefined">De origem obscura.
  </etym>
  <!--etc. -->
</entry>
```





**ouro, oiro** [óru], [ójru]. *s. m.* (Do lat. *aurum*). **1. Quím.** Elemento sólido (símb. Au, nº at. 79, p. at. 196,997), metal precioso de cor amarela, macio e maleável, inoxidável e não atacado pelos ácidos. *O ouro é muito usado em joalheria.* **de ouro.** 1. Diz-se de alguma coisa que tem a cor, o aspecto, o brilho desse metal precioso. *Banho de ouro.* 2. De grande valor. ≈ PRECIOSO. *O silêncio é de ouro.* **ensaiador<sup>+</sup> de ouro e prata. ouro branco.** 1. Liga deste metal que contém 20 a 50% de níquel, muito usada na ourivesaria. 2. Liga deste metal, níquel e paládio. 3. *Bras.* Algodão, considerado como riqueza agrícola. **ouro de lei.** 1. Aquele cujos quilates são determinados por disposições e que varia consoante o país. 2. Material de excelente qualidade. **ouro fino,** o que possui 24 quilates. **2.** Qualquer objecto feito desse metal. **medalha<sup>+</sup> de ouro.** **3.** Coisa muito valiosa. **a preço<sup>+</sup> de ouro. bodas<sup>+</sup> de ouro. coração<sup>+</sup> de ouro. ouro em pó.** 1. Coisa de qualidade superior. 2. Pessoa muito virtuosa, sincera, leal. **menino<sup>+</sup> de ouro. ouro negro,** petróleo, considerado como riqueza económica. **ouro coronário,** *Hist.,* coroa deste metal que davam, na Roma Antiga, aos generais vencedores. **sonho<sup>+</sup> de ouro.** **4.** Símbolo de riqueza. ≈ FORTUNA, OPULÊNCIA. **a galinha<sup>+</sup> dos ovos de ouro. mina<sup>+</sup> de ouro.** **5.** Cor amarela dourada. *Tinha cabelos claros, cor de ouro.* **6. Heráld.** Metal heráldico, amarelado, representado por uma série de pontos. **7. funç. adj. córdoba<sup>+</sup> ouro.** **8. livro<sup>+</sup> de ouro. 9. número<sup>+</sup> de ouro. 10. ouro potável,** líquido oleoso e alcoólico. *O silêncio é de ouro.* **cobrir alguém de ouro,** tratar muito bem alguém, dando-lhe tudo de bom. **nem tudo o que brilha/luz é ouro,** as coisas nem sempre são o que parecem ser. **ouro sobre azul.** 1. Coisa muito boa. 2. Reunião de duas coisas excelentes, que combinam na perfeição. **pagar a peso de ouro,** pagar muito bem. *O patrão pagava-lhe a peso de ouro.* **valer o seu peso em ouro,** ser muito valioso.

*ouro* [gold], DLPC (2001)

```
<entry type="simpleWord" xml:id="ouro" xml:lang="pt">
  <form type="lemma">
    <orth>ouro</orth>
    <pron>'oru</pron>
  </form>
  <form type="variant" xml:id="oiro" xml:lang="pt">
    <orth>oiro</orth>
    <pron>'ojru</pron>
  </form>
  <!--etc. -->
</entry>
```

**missanga** [misêŋɐ]. *s. f.* (Talvez do quimb. *missanga*, forma pl. de *mussanga*). **1.** Conjunto de pequenas contas de vidro, de cores variadas, que se usam como enfeite em colares, bordados... *Tinha um colar de missanga preta e outro de missanga branca, que costumava usar entrelaçados.* «Ao lado dormia a pequena surda-muda, dormia cheia dum suspirar e de activos sonhos em que decerto chorava a perda da sua missanga de vidro com que fazia colares e anéis.» (AGUSTINA, *Sermão*, p. 15). «o calcanhar sujo da meia saía-lhe para fora da chinela bordada a missanga» (EÇA, *Primo Basílio*, p. 33). **2.** Anel ou outro ornato feito com essas contas. **3.** *Tip.* Variedade de letra de imprensa muito miúda, correspondente aos caracteres tipográficos de corpos 4 e 5. **4.** Pessoa ou coisa insignificante, sem valor. **5.** Bagatela, bugiganga; coisa miúda.

*missanga* [glass bead], DLPC (2001)

```
<entry type="monolexicalWord" xml:id="missanga" xml:lang="pt">
  <form type="lemma">
    <orth>missanga</orth>
  </form>
  <!-- It isn't on paper; new edition information -->
  <form type="variant" xml:id="id" xml:lang="pt">
    <orth>miçanga</orth>
    <usg type="geo">
      <placeName>Brasil</placeName>
    </usg>
  </form>
  <!--etc. -->
</entry>
```



**missanga , miçanga** geo. Bras.

*n. f.*

**1.** Conjunto de pequenas contas de vidro, de cores variadas, que se usam como enfeite em colares, bordados... *Tinha um colar de missanga preta e outro de missanga branca, que costumava usar entrelaçados.*

«Ao lado dormia a pequena surda-muda, dormia cheia dum suspirar e de activos sonhos em que decerto chorava a perda da sua missanga de vidro com que fazia colares e anéis.» (AGUSTINA, *Sermão*, 15)

«o calcanhar sujo da meiasaía-lhe para fora da chinela bordada a missanga» (EÇA, *Primo Basílio*, 33)

**2.** Anel ou outro ornato feito com essas contas.

**3.** dom: **Tip.** Variedade de letra de imprensa muito miúda, correspondente aos caracteres tipográficos de corpos 4 e 5.

**4.** Pessoa ou coisa insignificante, sem valor.

**5.** Bagatela, bugiganga; coisa miúda.

(Talvez do quimb. *missanga*, forma pl. de *mussanga*)

Remover

Editar

/db/academia/missanga.xml



**donut** [dónut]. *s. m.* (Ingl.). V. *dónute*.

**dónute** [dónuti]. *s. m.* (Do ingl. *donut*). Argola de massa frita enrolada em açúcar ou com diferentes coberturas, em forma de anel ou bola.

*donut* [donut], DLPC (2001)

**walkman** [wółkmən]. *s. m.* (Ingl., da marca registada).

*Mús.* Leitor de cassetes portátil, por vezes com rádio acoplado, com auscultadores leves, para as pessoas ouvirem música quando se deslocam. «*O mesmo garoto tem um 'walkman', objecto que todos os mais novos sonham possuir*» (*DN*, 12.3.1990). Pl. walkmans.

*walkman* [walkman], DLPC (2001)

```
<entry type="simpleWord" xml:lang="en" xml:id="donut">
  <form type="lemma">
    <orth>donut</orth>
    <pron>d'ɔnut</pron>
  </form>
  <gramGrp>
    <gram type="pos" norm="NOUN">s.</gram>
    <gram type="gen">m.</gram>
  </gramGrp>
  <etym type="borrowing">
    <seg type="desc">Do</seg>
    <cit type="etymon" xml:lang="en">
      <form><orth>donut</orth></form>
      <!-- It isn't on paper; new edition information -->
      <usg type="domain">marca registada</usg>
    </cit>
  </etym>
  <!--etc. -->
</entry>
```

## 6. Conclusions and future work

- The results are useful for the discussion and definition of the TEI Lex-0 standard.
- The need to adapt and rewrite the guidelines.
- An agreement between academies and other institutions would be desirable to systematize and optimize resources that can provide a better representation of the entire European lexicographical heritage.

Thank you for your attention.

Obrigada pela vossa atenção.

anasalgado@campus.fcsh.unl.pt  
rute.costa@fcsh.unl.pt  
ttasovac@humanistika.org