



## Automatically Encoding Encyclopedic-like Resources in TEI

Mohamed Khemakhem, Laurent Romary, Simon Gabay, Hervé Bohbot,  
Francesca Frontini, Giancarlo Luxardo

### ► To cite this version:

Mohamed Khemakhem, Laurent Romary, Simon Gabay, Hervé Bohbot, Francesca Frontini, et al.. Automatically Encoding Encyclopedic-like Resources in TEI. The annual TEI Conference and Members Meeting, Sep 2018, Tokyo, Japan. hal-01819505

HAL Id: hal-01819505

<https://inria.hal.science/hal-01819505>

Submitted on 20 Jun 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Automatically Encoding Encyclopedic-like Resources in TEI

**Mohamed Khemakhem<sup>1,2,3</sup>, Laurent Romary<sup>1,2,4</sup>**

{name.surname@inria.fr}

<sup>1</sup> Inria ALMAnaCH, Paris

<sup>2</sup> Centre Marc Bloch, Berlin

<sup>3</sup> Paris Diderot University, Paris

<sup>4</sup> BBAW - Berlin-Brandenburgische Akademie der Wissenschaften, Berlin

**Simon Gabay**

{surname.name@unine.ch}

Université de Neuchâtel - Institut de littérature française

**Hervé Bohbot, Francesca Frontini, Giancarlo Luxardo**

{name.surname@univ-montp3.fr}

Univ Paul Valéry Montpellier 3, CNRS, PRAXILING UMR 5267

Encyclopedic content exists in different forms of paper-based resources and remains to a very large extent not exploited given the limited alternatives to extract information from the corresponding digitized material. Where TEI has provided mechanisms for encoding such a content, recent works have yielded advanced techniques to apply these mechanisms automatically on raw resources, however clear recommendations for the encoding of such content is still lacking.

As current digitization projects concern encyclopedias, other projects are undertaken to digitize textual resources with heritage significance and with similar features in terms of layout and indexed access. In this paper we introduce new lexicographic models, customizing the TEI schema, we extend existing ones and we present their application on two categories of resources, currently under development to be made available for public use:

- Encyclopedic sections of early, out of copyright, versions of a reference French Dictionary (Petit Larousse Illustré),
- Manuscripts auction catalogues, which are old resources for referencing and semantically describing manuscripts for sale.

The definition of these models is based on the significant similarities noticed at different levels of the described information in such entry-based resources. We have employed elements from the TEI dictionaries module to enable the segmentation of morphological and semantic information encompassed by respectively `<form>` and `<sense>` elements, already extracted by dedicated existing models in GROBID-Dictionaries<sup>1</sup>'s architecture. In fact, each entry in both types of resources carries two main encyclopedic information. First, the name or the label of a concept, coupled with one or more of its extended form and coming sometimes with a brief description. Then, a second block comes to establish an exhaustive description of the concept and to

---

<sup>1</sup> <https://github.com/MedKhem/grobid-dictionaries/>

present related information such as bibliography or pricing. In addition, we have used a customized version of the `<entry>` element to allow the generic encoding of entry numbering (see second excerpt in the table below).

TEI encoding	Excerpt from Dictionary Encyclopedic Section in Petit Larousse Illustré (1948)
<pre> &lt;entry&gt;   &lt;form type="lemma"&gt;     &lt;persName&gt;ABERDEEN&lt;/persName&gt;     &lt;addName&gt;(G. H. Gordon, comte d')&lt;/addName&gt;&lt;pc&gt;,&lt;/pc&gt;     &lt;desc&gt;homme d'Etat anglais, né à Edimbourg&lt;/desc&gt;     &lt;/form&gt;     &lt;pc&gt;.&lt;/pc&gt;     &lt;sense&gt;       &lt;def&gt;Premier ministre en 1852, il conclut avec la France une alliance contre la Russie (1784- 1860)&lt;/def&gt;       &lt;/sense&gt;       &lt;pc&gt;.&lt;/pc&gt;     &lt;/entry&gt;   </pre>	<p><b>ABERDEEN</b> [<i>aberdin'</i>], v. d'Ecosse, ch.-l. de comté; port sur la mer du Nord; 170.000 h. Université.</p> <p><b>ABERDEEN</b> (G. H. Gordon, <i>comte d'</i>), homme d'Etat anglais, né à Edimbourg. Premier ministre en 1852, il conclut avec la France une alliance contre la Russie (1784-1860).</p> <p><b>ABER-VRACH</b>, fl. côtier du Finistère (Atlantique); 34 kil. Station marémotrice d'essai.</p> <p><b>ABGAR</b>, nom de huit rois d'Edesse, en Mésopotamie (132 av. J.-C.-216 apr.).</p> <p><b>ABIA</b>, roi de Juda, fils de Roboam, vainqueur de Jéroboam, roi d'Israël (957-955 av. J.-C.).</p> <p><b>ABIDJAN</b>, ch.-l. de la Côte-d'Ivoire (A.-O. F.), sur une vaste lagune navigable; 15.000 h.</p> <p><b>ABIMÉLECH</b> [<i>lèk</i>], fils de Gédéon. Il devint Juge d'Israël, après avoir fait égorguer ses frères; il établit son pouvoir sur Sichem et fut tué au siège de Thèbes, en Palestine (vers 1100 av. J.-C.).</p> <p><b>ABIRON</b>, lévite qui fut englouti dans la terre avec Coré et Dathan, tous trois révoltés contre Moïse et Aaron (<i>Bible</i>).</p>

TEI encoding	Excerpt from a Manuscripts Auction Catalogue (1889)
<pre> &lt;entry&gt;   &lt;num&gt;49&lt;/num&gt;   &lt;form type="lemma"&gt;     &lt;surName&gt;Kourakin&lt;/surName&gt;     &lt;addName&gt;(le prince Alexis B.),&lt;/addName&gt;     &lt;desc&gt; frère du précédent, homme d'Etat russe.&lt;/desc&gt;   &lt;/form&gt;   &lt;sense&gt;     &lt;pc&gt;-&lt;/pc&gt;     &lt;def&gt;       &lt;bibl&gt;Billet auto sig., en français, à M. Monférand, 1 p, in-8.&lt;/bibl&gt; &lt;num type="price"&gt;2 »&lt;/num&gt;     &lt;/def&gt;   &lt;/sense&gt; &lt;/entry&gt; </pre> <pre> &lt;entry&gt;   &lt;num&gt;54&lt;/num&gt;   &lt;form type="lemma"&gt;     &lt;surname&gt;Lassalle&lt;/surname&gt;     &lt;addName&gt;(A.-Ch.-L. de)&lt;/addName&gt;,     &lt;desc&gt;le plus brillant général de cavalerie des guerres de la République et de l'Empire, né à Metz, tué à la bataille de Wagram&lt;/desc&gt;   &lt;/form&gt;   &lt;sense&gt; .-     &lt;def&gt;       &lt;bibl&gt;L. a. s. au général Dugua; 1 p. in- f.&lt;/bibl&gt;       &lt;num type="price"&gt;10 »&lt;/num&gt;     &lt;/def&gt;     &lt;note&gt;Superbe lettre sur la campagne d'Egypte. Il profite du départ du général Desaix pour lui donner des nouvelles. Desaix lui laisse le commandement de la colonne qui doit poursuivre Mourad-Bey, et qui se compose de 400 hommes de cavalerie, 4 pièces de canon et 160 dromadaires. Le général Boyer a, dans une petite affaire, tué 10 mamelouks et 40 arabes, etc.&lt;/note&gt;     &lt;/sense&gt; &lt;/entry&gt; </pre>	<p>49 <b>Kourakin</b> (le prince Alexis B.), frère du précédent, homme d'Etat russe. — Billet aut. sig., en français, à M. Monférand, 1 p. in-8. 2 »</p> <p>50 <b>Labanoff</b> (le prince Alex.), célèbre général et écrivain russe, historien de Marie Stuart. — L. a. s., en français, 1835, 1 p. in-4. 3 »</p> <p>51 <b>Ladislas IV</b>, roi de Pologne, célèbre par ses succès contre les Russes, époux de Marie de Gonzague. — L. sig., en latin, au cardinal de Montalte; Varsovie, 1645, 1 p. in-f. 8 »</p> <p>52 <b>Lafayette</b>, illustre général.—L. a. sig. de ses initiales à M. Masclet; Washington, 13 août 1825, 1 p. 1/4 in-4. Un peu fatiguée. 15 »</p> <p>Très-curieuse lettre sur le voyage qu'il fit en Amérique, de 1824 à 1825. « C'est avec de bien tendres regrets que je quitterai cette terre américaine, le bon, grand et heureux peuple des Etats-Unis auquel je suis amalgamé depuis près d'un demi-siècle, et qui vient encore de me combler de ses bontés. J'y ai vu les miracles de l'indépendance, de la liberté, égalité et <i>self government</i>; le problème des institutions républiques a été résolu ici sur une grande échelle et jamais expérience n'a si bien réussi. » Il comptait retourner comme il était venu, sur un paquebot-poste, mais le peuple et le gouvernement en ont disposé autrement. On a donné le nom de <i>Brandywine</i> à une superbe frégate qui est chargée de le ramener en France.</p> <p>53 <b>La Roncière</b> (Emile-Clement de), fils du général, condamné pour tentative de viol.—L. a. s. aux officiers et élèves de l'école de Saumur; Paris, mai 1836, 3 p. pet. in-4. 10 »</p> <p>Très-curieuse lettre toute relative à son procès.</p> <p>54 <b>Lassalle</b> (A.-Ch.-L. de), le plus brillant général de cavalerie des guerres de la République et de l'Empire, né à Metz, tué à la bataille de Wagram.—L. a. s. au général Dugua; 1 p. in-f. 10 »</p> <p>Superbe lettre sur la campagne d'Egypte. Il profite du départ du général Desaix pour lui donner des nouvelles. Desaix lui laisse le commandement de la colonne qui doit poursuivre Mourad-Bey, et qui se compose de 400 hommes de cavalerie, 4 pièces de canon et 160 dromadaires. Le général Boyer a, dans une petite affaire, tué 10 mamelouks et 40 arabes, etc.</p>
	<p>We employed the adapted and newly defined models to automatically process and label the encyclopedic information in a cascading fashion, as introduced in Khemakhem et al. 2017 and Khemakhem et al. 2018, by relying on the same sequence labeling machine learning technique.</p>

We will present in-depth our progress in implementing morphological and semantic models and show through the results at each structuring level how such information could be uniformly encoded and automatically extracted from these resources.

**References:**

1. Mohamed Khemakhem, Luca Foppiano, Laurent Romary. Automatic Extraction of TEI Structures in Digitized Lexical Resources using Conditional Random Fields. electronic lexicography, eLex 2017, Sep 2017, Leiden, Netherlands
2. Mohamed Khemakhem, Axel Herold, Laurent Romary. Enhancing Usability for Automatically Structuring Digitised Dictionaries. GLOBALEX workshop at LREC 2018, May 2018, Miyazaki, Japan.
3. Hervé Bohbot, Francesca Frontini, Giancarlo Luxardo, Mohamed Khemakhem, Laurent Romary. Presenting the Nénufar Project: a Diachronic Digital Edition of the Petit Larousse Illustré. GLOBALEX 2018 - Globalex workshop at LREC 2018, May 2018, Miyazaki, Japan.