

# Ìðĩãããĩũ èẽĩããẽñòè-ãñẽĩé ðãçĩãòèè è ãĩãèèçã ýããèòðĩĩĩũõ èðèè-ãñẽõ èçããĩèè òãẽñòĩã

Áãõĩð Áãããñãé Ìèòãããĩãè- Èããðãĩòũãã  
05.08.2012 ã.  
Ìñãããããã ããĩããããã 27.08.2012 ã.

Summary. In this paper we consider some problems of automatic linguistic annotation and analysis of textual heritage documents encoded according to the TEI XML guidelines. TEI XML is a popular standard for encoding electronic editions of textual heritage documents as it allows highly customizable semantically-oriented markup independent of a particular platform or software. TEI is aimed at facilitating data exchange and interoperability. However, rich editorial markup including various readings and interpretations at various levels of linguistic hierarchy may be a serious challenge if one wants to apply NLP (natural language processing) tools to such an edition. Based on the example of the Base de Français Médiéval Old French corpus and on the electronic edition of the Queste del saint Graal, we will discuss the solutions to these problems that are implemented in the TXM platform import modules.

## Ñòãĩããðò ýããèòðĩĩĩé

ðãçĩãòèè òãẽñòĩã ã òĩðĩãòã XML ã ñĩĩòããòñòãèè ñ ðãẽĩããããòèýĩè Ìãããóããðĩãĩé Èĩèèããòèãũ ñ Èããèðĩããĩèp Òããñòĩã TEI (Text Encoding Initiative, <http://www.tei-c.org>) Ìðĩããðãè ã ññãããĩèã ãĩãũ ãĩñòãòĩ-ĩĩ èèðĩèĩã ðãñĩðĩñòðãĩãĩã ã ããããñòè ýããèòðĩĩĩãĩ èçãããĩèý òãẽñòĩã, Ìðĩãããããããũèõ è ðĩããò Ìèñũĩãĩĩã ããñãããèý ðãçèè-ĩũò ñòðãĩ è èóèũòòð. 150 Ìðĩãèòĩã Ìðããñòããèãĩũ ã ñãèòã TEI, ã ãããñòãèòããũĩñòè ýòĩ -èñẽĩ Ìããò ãũòũ çĩã-èòããũĩĩ ãóèũòè. Ìðããĩòũããñòããĩè ðãçĩãòèè ã ñòãĩããðòã TEI ýããýðòñý ãã ñĩðã ãã òãããèũĩĩ ðãçðããĩòãĩĩòð òãĩðèð ñòðòèèòòðũ òãẽñòã è ãĩèòããĩòã, èãããèñòũ Ìãðñĩããèèçãòèè è ãããĩòãòèè è èĩĩèðãòĩĩò Ìãòãðèããèò çã ñ-ãò Ìãòèũĩĩé Ìðããĩèçãòèè è ñĩãòèããũĩĩã Ìãðãĩèçã ñĩãòèèèããòèè ODD, ã òãèãã Ìãçããèñèĩñòũ Ìò èĩĩèðãòĩĩé Ìããòðĩũ èèè Ìðĩãðããĩĩã Ìðĩãèòã. Ìããñòã ñ òãã ðãçãũ-ããĩãý ãããèñòũ TEI ñĩçãããò ñĩðããããããĩũã òðããĩñòè ããý ðãçðããĩòèè Ìðĩãðããĩũũò ñðããñòã Ìãðããĩòèè, ãããèèçã è Ìããèèããòèè òãẽñòĩã, ðãçĩã-ããĩũò ã ýòĩ ñòãĩããðòã, Ìñĩããĩĩ ãñèè ðã-ũ èããò Ì ñðããñòããò «èèðĩèĩã Ìðĩèèèý», Ìðããĩãçãã-ããĩũò ããý èñĩũèçãããĩèý ããã ðãĩè Ìãããèũĩĩ ãçýòãĩ Ìðĩãèòã. Ì ðããñòĩñòè, «ãèòããèòð» òèèĩèĩãè-ãñèòð ðãçĩãòèò, ò-èòũããðũòð ðãçĩã-ðãĩèý è ããðèããòũ èĩòãðĩðãòãòèè òðããããĩòã òãèñòã ã ðãçĩũò òðĩãĩýò èãðãðòèè ýçũèĩãũò ñòðòèèòòð, Ìããò ãũòũ òðòãĩ ñĩãããñòèòũ ñ èñĩũèçãããĩèã èĩñòðòããĩòã ããòũããòè-ãñẽĩé èẽĩããèñòè-ãñẽĩé ðãçĩãòèè (òĩèããĩèçãòèè, èããããòèçãòèè, Ìðòĩèĩãè-ãñẽĩé èãòããããðèèè è ò.ĩ.).

Á Ìãðãã ããèèããã Ìũ Ìðããñòãããè Ìãòããèèò ñĩãããããèè (Ìðĩãèèçãòèè) òèèĩèĩãè-ãñẽĩé ðãçĩãòèè òãẽñòĩã Áãçũ ñðãããããããããããã òðããòóçñèĩã ÿçũèã BFM (Base de Français Médiéval, <http://bfm.ens-lyon.fr>) ã Ìðĩããñãã èõ çããðóçèè Ìã Ìããòòĩðĩò TXM ñ òãèũð èõ ããèũããèøããã èẽĩããèñòè-ãñẽĩã ãããèèçã. Ìũ òãèãã ðãññĩòðè ðãçòèèòãòũ ñũòãã ñ ãããããòèè ãããĩé Ìãòãèèè è òãèñòã ãðòãèò Ìðĩãèòã, ðãçĩã-ããĩũ Ìã ññĩãã ðããèãããããòèè TEI, Ì Ìðèãããããũè Ìèèè-ũã Ìò BFM ðãðããĩèý ã ðýãã èèð-ããũò ããý èẽĩããèñòè-ãñẽĩã ãããèèçã ãñĩãèòã ðãçĩãòèè.

Òããñòããòðèý (textométrie) ãĩçãèèã èãé Ìãò-ĩĩã Ìããããããããã ãã Òðããòèè ã 1980-ã ããũ. Á ãã ðããèò ãũèè ðãçðããããããã ÿòããèèããããã Ìãòãèèè ãããèèçã Ìããããããã ãããããããã ðããñòããããã. Áñããã çã èããèèèãããðèãã è ñòãèñòè-ãñẽĩé ãããèèçã ðããñòã òããñòããããã Ìðãããããããã ñòãèñòè-ãñẽè Ìããããããããããã Ìãòããã è èĩñòðòãããã ãããèèçã ããý ðãçèè-ũũò ãããããèòãðĩũò Ìãóè.

TXM &ndash; ýòĩ Ìããòèũãã Ìããòòĩðã ñ Ìèèðũòũ ãñòããããã ãããã, èãðããã ñãã-ãòããò òóãèèè ðãçèè-ũũò ðãããã ðãçðããããããããã Ìðĩãðãã òããñòãããòðè-ãñẽĩã ãããèèçã. Ìã Ìðããñòãããããã Ìããã Ìããèãããã òããñòãããòðè-ãñẽĩã èĩñòðòãããããðèý, èñĩũèçãòðããã ñããðããããããã èĩðĩòãããã òãòãããããèè (Unicode, XML, TEI, NLP). Ìãðãããããã ããããããããã Ìããòòĩðã TXM Ìðããñòãããããã ã Ìããèèããòèýò [Heiden 2010; Heiden et al. 2010; Pincemin et al. 2010], ã òãèãã Ìã ñãèòã <http://textometrie.ens-lyon.fr/?lang=en>.

Áãçã ñðãããããããããããã òðããòóçñèĩã ÿçũèã (BFM) &ndash; ýòĩ èĩðãã òããñòãã ñòãðã-è ñðãããããããããããã ÿçũèã (IX &ndash; XV ãã.), ã Ìãñòãããããã ãðããã ðãçðãããããããããããã ãããããããããããã ICAR Ìãèèããããããããã òããããã Ìãò-ũũò èññããããããããã

Ōāiōēē (CNRS) ē Ēēīīnēīāī ōīēāāōñēōāōā. Āāçā iāñ-ēōūāāāō 75 ōāēñōīā iāuēī iāuāīīī āīēāā 3 500 000 ōāēñōīōīōī. Ēñōī-īēēāīē BFM ā īñīīāīīī yāēyōñy āāōīōēōāōīūā ēōēōē-āñēēā ēçāāīēy, īāīāēī ā īñēāāīāā āōāīy ðaçāēāpōñy ñīāñōāāīīūā ēçāāīēy, īēōāpūēāñy iā ēēīāāēñōē-āñēē āūāāōāīīūā ōōāīñēōēīōēē īōēāēīāēūīūō ðōēīēñāē. Ā ēā-āñōāā īōēīāōā īīāēī īōēāāñōē ēīōāōāēōēāīīā ēçāāīēā āīīēīīāīī āīīāīā XIII ā. «Iīēñēē Nāyōīāī Āōāāēy» («La Queste del saint Graal») īīā ðāāāēōēēē Ē. Iāōēāēēī-īēçy [Queste 2011].

N īāy 2012 āīāā āīñōōī ē BFM īñōūāñōāēyāōñy īñōāāñōāīī īōōāēā <http://txm.bfm-corpus.org/bfm>, īñōōīāīīīāī iā īēāōōīōīā TXM.

Āñā ōāēñōū BFM ðaçīā-āīū ā ōīōīāōā XML iā īñīīāā ōāēīīāīāōēē TEI, ā ñīīōāāōñōāēē ñī ñīāōēōēēāōēēāē, ðaçōāāīōāīīē āēy íōæā īōīāēōā [Guillot et al. 2010] ñ ō-āōīī īāōñīāēōēēē ēēīāāēñōē-āñēīāī āīāēēçā ēīōīōñīā īōēīāīēōāēūīī ē ōāēñōāī, āēēp-āpūēī āēōāīēōp ðāāāēōīōñēōp ðaçīāōēō, yāēyāōñy ēīōōāēōīāy ēāāīōēōēēāōēy ñēīā (ōīēāīēçāōēy) ē īōāāēīāēāīēē, ē ēīōīōūī yōīō āīāēēç āīēæāī īōēīāīyōūñy āāç īīōāōē ñāīīē ðāāāēōīōñēīē ðaçīāōēē, āūyāēāīēy īōyīīē ðā-ē).

Īāīīē ēç iāēāīēāā ñēīāēīūō çāāā- īōē ēñīīēūçīāāīēē ñōāāñōā āāōīāōē-āñēīāī ēēīāāēñōē-āñēīāī āīāēēçā ēīōīōñīā īōēīāīēōāēūīī ē ōāēñōāī, āēēp-āpūēī āēōāīēōp ðāāāēōīōñēōp ðaçīāōēō, yāēyāōñy ēīōōāēōīāy ēāāīōēōēēāōēy ñēīā (ōīēāīēçāōēy) ē īōāāēīāēāīēē, ē ēīōīōūī yōīō āīāēēç āīēæāī īōēīāīyōūñy āāç īīōāōē ñāīīē ðāāāēōīōñēīē ðaçīāōēē.

Nēāāōpūēē īōēīāō ðaçīāōēē, ñīāāōæāūēē īōāāēīāēāīīūē ðāāāēōīōīī ōōāāīāīō ōāēñōā iā iāñōā ēāēōīū, iā-ēīāpūāēñy ā ēīīōā īāīīāī ñēīāā ē çāēāī-ēāāpūāēñy iāñēīēūēēēē ñēīāāīē īīçāā, āāñīēpōīī ēīōōāēōāī ñ ōī-ēē çōāīēy ðāēīīāīāāōēē TEI, īāīāēī āāñūā ñēīāēāī āēy ōīēāīēçāōēē ñ ō-āōīī īñīāāīīñōāē yçūēā XML (çāīōāō «īāōāēōāūēāāīēy» yēāīāīōīā, ðēñē īīyāēāīēy īōīāāēīā iāæāō ōyāāīē ē ōāēñōīāūīē ōçēāīē īōē īāōāāīōēā):

en<supplied>tra a  
cheval en la</supplied> sale une mout bele damoisele

Āūā āīēāā ñāōūāçīūā īōīāēāīū āīçīēēāpō īōē ðaçīāōēā īōāāēīāēāīēē, īñīāāīī ā ñōēōīōāīōīūō ōāēñōāō, āāā īāōāēōāūēāāīēā iāōōē-āñēīē ē ñēīōāēñē-āñēīē ñōōōēōōōū āñōōā-āāōñy ī-āīū -āñōī.

ðaçōīāāōñy, īīāēī «īōōēēūōōīāāōū» āñā yēāīāīōū, īāōāēōāūēāāpūēāñy ñ īñīīāīūīē ēēīāāēñōē-āñēēēē ñōōōēōōōāīē (ñēīāāīē ē īōāāēīāēāīēyīē), īāīāēī yōī īīāāō īōēāāñōē ē īōāōā ñōūāñōāāīīē ēīōīōīāōēē āēy çāīōīñīā ē āēçōāēēçāōēē (īāīōēīāō, ī ōīī, īīāāāōāāēāñū ēē ñēīāīōīōīā ðāāāēōīōñēīē īōāāēā).

Ñīçāāīēā āēāīōēōīā ōīēāīēçāōēē, ēīōīōūē ēīōōāēōīī īāōāāōūāāē āū ēpāīē ōāēñō ñ āēōāīēīē ðāāāēōīōñēīē ðaçīāōēīē ā ñōāīāāōōā TEI XML, īōāāñōāāēyāōñy īōāēōē-āñēē īāāīçīīāēīūī. Ōāī iā īāīāā, īīāēī āīāēōūñy āīīēā ōāīēāōāīōēōāēūīūō ðaçōēūōāōīā īōē ōñēīāēē, -ōī ðaçīāōēā ēñōīāīīāī āīēōīāīōā īōāā-āāō ðyāō īōīñōūō īōāāēē. Īāīōēīāō, «ōyāē, ðāñīīēīāēāīīūā āīōōōē ñēīā, āīēāēīū āūōū -āōēī ēāāīōēōēōēōīāāīū» ēēē «āñēē ðaçīā-āīīūē ñāāīāīō ōāēñōā iā-ēīāāōñy āīōōōē īāīīāī ñēīāā ē çāōāāōūāāāō īāñēīēūēī īīñēāāōpūēō, āāī īāīāōīāēīī ðaçāāēēōū».

Ōāēæā āīçīīāēī ñīñōāāēōū ñīēñēē ōyāīā TEI ā çāāēñēīīñōē īō ēō īīçōēē ā ēēīāāēñōē-āñēīē ēāōāōēē ōāēñōā ā ðāīēāō īōāāēūīīāī īōīāēōā. Īāīōēīāō, īōāāēīāēāīēy ðāñīīēāāāpōñy āīōōōē yēāīāīōīā ōēīā «āaçāō» <p> ēēē «āēīē ōāēñōā» <ab>, ā iā īāīāīōīō. Īāēīōīōūā ōyāē īīāēīī ñ-ēōāōū

ýéáéááéáíóíúíè ñéíáó (íáíðèíáð, <abbr>, <num>, <pc>), à íáñéíéúéí ýéáíáíóíá ïí-òè áñáááá ðáñĩíéááááòñý áíóóðè ñéíáá (<am>, <c>, <ex>). Ðýá ýéáíáíóíá ñíááðæáð ñááíáíóú óáéñòá, éíòíðúá íá ñéááóáò òíéáíéçèðíááòú, ïñéíéúéó ííè íá íðéíááéáæáð é íàðáðèáéó èñóí-íééá (íáíðèíáð, ðáááéóíðñééá íðéíá-áíéý è ñííñéè á óáéñòá éðèðè-áñéíáí èçááíéý). Á ðáíéáó éííéðáòíáí íðíáéòá ýè ñíéñéè ííáóó áúòú ñóúáñòááíí ðáñæéðáíú é óòí-íáíú, á ðáçóéúòáðá -ááí -éñéí ýéáíáíóíá, éíòíðúá ííáóó íáðáéðáúéááòúñý ñ èéíááéñòè-áñééíè ñòðóéóóðáíè, çíá-èòáéúíí ñíéðáúááòñý.

Ðàçíáðéá óáéñòíá BFM ñííòááòñòáóáò -áòéí ñòíðíóéèðíááííé ñíáòéòééáòéè íðéíáíáíéý ðáéñíáíááòéé TEI, íòðáæáííé á áíéóíáíóáòéè ODD [Guillot et al. 2010]. Ááííáý ñíáòéòééáòéý áéèp-ááò íðááééá èñííéúçíááíéý áíóóðèñéíáíúó óýáíá (íáíðèíáð, èñíðááéáííúó ðáááéóíðíí áóéá èéè çíáéíá íáðáííñá), à óáéæá ýéáíáíóíá, éíòíðúá ííáóó íáðáéðáúéááòúñý ñí ñòðóéóóðíé íðááéíæáíéè (íáíðèíáð, «íòñòúá» ýéáíáíóú <lb/> «ííááý ñòðíéá» èñííéúçóáòñý áíáñòí <l> «ñòèò», à ýéáíáíó òèðèðíááíéý <q> ðáññíáòðèááòñý éáé áðáíéòá íðááéíæáíéý). Ýóí ïíçáíéééí ðááéèçíááòú á íèàòóíðíá TXM ýóðáéòéáíúé áéáíðèòí òíéáíéçáòéè áéý óáéñòíá BFM.

Á ïñéááíáá áðáíý áúé óñíáòí íðíááááí ðýá óáñòíá ïí áááíòáòéè ááííáí áéáíðèòíá è áíéóíáíóáí TEI XML, ïíááíòíáéáííúí á áðóáéò íðíáéòáò. Ñðááè íéò ííæíí íáçááòú «Áèðóóáéúíóá áéáéèíòáéó áóíáíéñòíá» (<http://www.bvh.univ-tours.fr>) è èçááíéá -áðííáééíá ðííáíá «Áóááð è Íáèpðá» Ápñòááá Õéíááðá (<http://dossiers-flaubert.ish-lyon.cnrs.fr>). Áááíòáòéý ñíñòíéò á íðéíáíáíéè ñíáòéáéúíúó óèéúòðíá XSL íá áðíáá è íá áúóíáá íðíóááóðú òíéáíéçáòéè. «Áóíáííé» óèéúòð óááéýáò óýáé, éíòíðúá íá íðááñòááéýáò éíóáðáñá áéý ýéñíéóáòáòéè ñ ïííúúá TXM, à óáéæá óíðíúááò è íðíáéèçóáò íáéíòíðúá ñéíáíúá ñòðóéóóðú ýéáíáíóíá XML. «Áúóíáííé» óèéúòð ïíçáíéýáò èñíðááéòú ðýá íðéáíé, éíòíðúó íðáéòè-áñéè íááíçííæíí èçááæáòú á íðíóáññá íáðáé-ííé òíéáíéçáòéè (íáíðèíáð, ááéáíéá ñéíáá íðè íáðáííñá á éííóá ñòðáíéòú, éííáá óáéúé ðýá óýáíá ííæáò ðáñĩíéááòúñý íáæáò ááí íá-àéíí è éííóíí).

Á óáéñòá áíééááá íú íðéááááí éííéðáóíúá íðéíáðú ðááéèçíááííúó ðáðáíéè è íðíááíííñòðèðóáí ðáçóéúòáòú, éíòíðúá ííæíí ííéó-èòú íðè ýéñíéóáòáòéè óáéñòíá BFM è áðóáéò íðíáéòíá ñ ïííúúá íèàòóíðíú TXM. Áíéáá ïíáðíáíí ðáçèè-íúá áéáíðèòíú òíéáíéçáòéè, ðááéèçíááííúá íá íèàòóíðíá TXM, ïíéñáíú á [Heiden 2010].

Ñíèñíé èèòáðáòóðú

Guillot, C., Heiden, S., Lavrentiev, A., Bertrand, L. Manuel d'encodage XML-TEI des textes de la Base de Français Médiéval, Lyon: Équipe BFM, 2010. &ndash; Ááðáñ á Éíóáðíáò [http://bfm.ens-lyon.fr/IMG/pdf/Manuel\\_Encodage\\_TEI.pdf](http://bfm.ens-lyon.fr/IMG/pdf/Manuel_Encodage_TEI.pdf).

Heiden, S. The TXM Platform: Building Open-Source Textual Analysis Software Compatible with the TEI Encoding Scheme. 24th // Pacific Asia Conference on Language, Information and Computation / Ed. Kiyoshi Ishikawa Ryo Otoguro. Institute for Digital Enhancement of Cognitive Development, Waseda University, 2010. P. 389-398.

Heiden, S., Magué, J.-P., Pincemin, B. TXM : Une plateforme logicielle open-source pour la

textométrie &ndash; conception et développement // Statistical Analysis of Textual Data - Proceedings of 10th International Conference JADT 2010 / Bolasco, S. et al. (Eds.). &ndash; Rome: Edizioni Universitarie di Lettere Economia Diritto, 2010. &ndash; P. 1021-1032. &ndash; Àäðñ à Èíòðíáò  
[http://www.ledonline.it/ledonline/JADT-2010/allegati/JADT-2010-1021-1032\\_025-Heiden.pdf](http://www.ledonline.it/ledonline/JADT-2010/allegati/JADT-2010-1021-1032_025-Heiden.pdf).

Pincemin, B., Heiden, S., Lay, M.-H., Leblanc J.-M. and Viprey, J.-M. Fonctionnalités textométriques: Proposition de typologie selon un point de vue utilisateur. // Statistical Analysis of Textual Data - Proceedings of 10th International Conference JADT 2010 / Bolasco, S. et al. (Eds.). &ndash; Rome: Edizioni Universitarie di Lettere Economia Diritto, 2010. &ndash; P. 341-353. &ndash; Àäðñ à Èíòðíáò  
[http://www.ledonline.it/ledonline/JADT-2010/allegati/JADT-2010-0341-0354\\_023-Pincemin.pdf](http://www.ledonline.it/ledonline/JADT-2010/allegati/JADT-2010-0341-0354_023-Pincemin.pdf).

Queste del saint Graal. Édition numérique interactive / Ed. Marchello-Nizia, Ch. &ndash; Lyon: Équipe de la BFM, 2011. &ndash; Àäðñ à Èíòðíáò <http://txm.bfm-corpus.org/txm>.