

Īōīāēāīū ēēīāāēñōē÷āñēīē ðāçīāōēē ē āīāēēçā yēāēōðīīīūō ēðēō÷āñēēō ēçāāīēē òāēñōīā

Āāōīð Āēāēñāē Iēōāēēīāē÷ Ēāāðāīōūāā
05.08.2012 ā.
Īñēāāīāā īāīīāēāīēā 27.08.2012 ā.

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.

Nōāīāāðō yēāēōðīīīē

ðāçīāōēē òāēñōīā ā òīðīāōā XML ā nīīōāāðñōāēē ñ ðāēīīāīāāōēyīē Iāæāōīāðīāīē Ēīēōēāōēāū īī Ēīāēðīāāīēp Ōāēñōīā TEI (Text Encoding Initiative, <http://www.tei-c.org>) īðēīāðāē ā īñēāāīēā āīāū āīñōāōī÷īī øēðīēīā ðāñīðīñōðāīāīēā ā īāēāñōē yēāēōðīīīāī ēçāāīēy òāēñōīā, īðēīāēāæāūēō ē ðīīāō īēñūīāīīāī īāñēāāēy ðāçēē÷īūō ñōðāī ē ēōēūōðð. 150 īðīāēōīā īðāāñōāāēāīū īā ñāēōā TEI, ā āāēñōāēōāēūīñōē yōī ÷ēñēī īīæāō āūōū çīā÷ēōāēūīī āóēūøēī. Īðāēīōūāñōāāīē ðāçīāōēē ā ñōāīāāðōā TEI yāēyþōñy āā īīðā īā òūāðāēūīī ðāçðāāīōāīōþ òāīðēþ ñōðōēōðū òāēñōā ē āīēōīāīōā, ēāāēīñōū īāðñīāēēçāōēē ē āāīōāōēē ē ēīīēðāōīīō Iāðāðēāēō çā ñ÷āð īīāōēūīīē īðāāīēçāōēē ē ñīāōēāēūīīāī īāðāīēçīā ñīāōēōēēāōēē ODD, ā òāēæā īāçāāēñēīñōū īō ēīīēðāōīīē īēāðōīðīū ēēē īðīāðāīīīāī īðīāēōā. Āīāñōā ñ òāī ÷ðāçāū÷āēīāy āēāēīñōū TEI ñīçāāāō īīðāāēāāīūā òðōāīñōē āēy ðāçðāāīōēē īðīāðāīīūō ñðāāñōā īāðāāīōēē, āīāēēçā ē īōāēēēāōēē òāēñōīā, ðāçīā÷āīīūō ā yōī ñōāīāāðōā, īñīāāīī āñēē ðā÷ū ēāāō ī ñðāāñōāāō «øēðīēīāī īðīōēēy», īðāāīāçīā÷āīīūō āēy ēñīēūçīāāīēy āīā ðāīīē īōāāēūīī āçyōīāī īðīāēōā. Ā ÷āñōīñōē, «āēōāīēōþ» òēēīēīāē÷āñēōþ ðāçīāōēō, ò÷ēōūāþþōþ ðāçīī÷ðāīēy ē āāðēāīōū ēīōāðīðāðōāōēē òðāīāīōīā òāēñōā īā ðāçīūō òðīāīyō ēāðāððēē yçūēīāūō ñōðōēōðð, īīæāō āūōū òðōāīī ñīāīāñōēōū ñ ēñīēūçīāāīēāī ēīñōðōīāīōīā āāōīāōē÷āñēīē ēēīāāēñōē÷āñēīē ðāçīāōēē (ðīēāīēçāōēē, ēāīīāðēçāōēē, īððōīēīāē÷āñēīē ēāðāāīðēçāōēē ē ò.ī.).

Ā īāøāī āīēēāāā īū īðāāñōāāēī Iāðīāēēō īīāāīðīāēē (īðīāēēçāōēē) òēēīēīāē÷āñēīē ðāçīāōēē òāēñōīā Āāçū ñðāāīāāāēīāīāī òðāīōóçñēīāī yçūēā BFM (Base de Français Médiéval, <http://bfm.ens-lyon.fr>) ā īðīōāññā ēō çāāðóçēē īā īēāðōīðīō TXM ñ òāēūþ ēō āāēūīāēøāāī ēēīāāēñōē÷āñēīāī āīāēēçā. Īū òāēæā ðāññīðōēī ðāçōēūōāðū īīūōīā īī āāāīōāōēē āāīīē Iāðīāēēē ē òāēñōāī āðōāēō īðīāēōīā, ðāçīā÷āīīūī īā īñīāā ðāēīīāīāāōēē TEI, īī īðēīāīyþūēī īòēē÷īūā īō BFM ðāðāīēy ā ðyāā ēēþ÷āāūō āēy ēēīāāēñōē÷āñēīāī āīāēēçā āñīāēōīā ðāçīāōēē.

Ōāēñōīāðōðēy (textométrie) āīçīēēēā ēāē īāó÷īīā īāīðāāēāīēā āī Ōðāīōēē ā 1980-ā āīāū. Ā āā ðāīēāð āúēē ðāçðāāīōāīū yóðāēōēāīūā īāðīāēēē āīāēēçā īāúāīūō ēīðīōñīā òāēñōīā. Āñēāā çā ēāēñēēīīāððēāē ē ñōāðēñōē÷āñēēī āīāēēçīī òāēñōā òāēñōīāðōðēy īðāāēāāāāð ñōāðēñōē÷āñēē īāīñīāāīūā īāðīāū ē ēīñōðōīāīōū āīāēēçā āēy ðāçēē÷īūō āōīāīēōāðīūō īāóē.

TXM – yōī īīāōēūīāy īēāðōīðīā ñ īōēðūōūī ēñōīāīūī ēīāīī, ēīðīðāy ñī÷āðāāð òōīēōē ðāçēē÷īūō ðāīāā ðāçðāāīōāīūō īðīāðāī òāēñōīāðōðē÷āñēīāī āīāēēçā. Īā īðāāñōāāēyāð īīāīā īīēīēāīēā òāēñōīāðōðē÷āñēīāī ēīñōðōīāīōāðēy, ēñīēūçōþūāā ñīāðāīāīūā ēīðīōñīūā òāðīīēīāēē (Unicode, XML, TEI, NLP). Īāðīāīāy ēīðīāōēy ī īēāðōīðīā TXM īðāāñōāāēāīā ā īōāēēēāōēyō [Heiden 2010; Heiden et al. 2010; Pincemin et al. 2010], ā òāēæā īā ñāēōā <http://textometrie.ens-lyon.fr/?lang=en>.

Āāçā ñðāāīāāāēīāīāī òðāīōóçñēīāī yçūēā (BFM) – yōī ēīðīōñ òāēñōīā ñōāðī- ē ñðāāīāððāīōóçñēīāī yçūēā (IX – XV āā.), ā īāñōīyūāā āðāīy ðāçðāāāðūāþþūēēñy ēāāīðāðīðēāē ICAR īāēēīāēūīīāī òāīðā īāó÷īūō ēññēāāīāāīēē

Ōāiōēē (CNRS) ē Eēīīnēīāī ōīēāāōñēōāōā. Āāçā iāñ-ēōūāāāō 75 ōāēñōīā iāuēi iāuāīīī āīēāā 3 500 000 ōāēñōīōīōī. Eñōī-īēēāīē BFM ā īñīīāīīī yāēyōñy āāōīōēōāōīūā ēōēōē-āñēēā ēçāāīēy, iāīāēī ā īñēāāīāā āōāīy ðaçāēāpōñy ñīāñōāāīīūā ēçāāīēy, īēōāpūēāñy iā ēēīāāēñōē-āñēē āūāāōāīīūā ōōāīñēōēīōēē īōēāēīāēūīūō ðōēīēñāē. Ā ēā-āñōāā īōēīāōā īīāēī īōēāāñōē ēīōāōāēōēāīīā ēçāāīēā āīīēīīāīīā ðīāīā XIII ā. «Iīēñēē Nāyōīāī Āōāāēy» («La Queste del saint Graal») īīā ðāāāēōēēē E. 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 īāōñīāēōēāū ēēīāāēñōē-āñēīāī āīāēēçā. Iēāōōīōīā TXM iā ōīēūēī ñēōæōō āēy ēīōīōñā BFM ñōāāñōāīī āīñōōīā īīēūçīāāōāēāē, īī ē īīçāīēyāō īñōūāñōāēyōū ðyā īīāōāōēē āāōīāōē-āñēīē ē īīēōāāōīāōē-āñēīē ðaçīāōēē (ā -āñōīīñōē, īōōīēīāē-āñēīē āīīōāōēē, āūyāēāīēy īōyīē 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-ēāpūāēñy iāñēīēūēēīē ñēīāāīē īīçæā, āāñīēpōīī ēīōōāēōāī ñ ōī-ēē çōāīēy ðāēīīāīāōēē TEI, iāīāēī āāñūā ñēīāēāī āēy ōīēāīēçāōēē ñ ō-āōīī īñīāīīīñōāē 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, īñīāāīīī 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-āñēīēē ñōōōēōōōāīē (ñēīāāīē ē īōāāēīāēāīēyīē), iāīāēī yōī īīæāō 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 iāīāā, 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 ðaçāāēēōōū».

Ōāēæā āīçīīāēīī ñīñōāāēōū ñīēñēē ōyāīā TEI ā çāāēñēīīñōē īō ēō īīçōēē ā ēēīāāēñōē-āñēīē ēāōāōēē ōāēñōā ā ðāīēāō īōāāēūīīāī īōīāēōā. Īāīōēīāō, īōāāēīāēāīēy ðāñīīēāāāpōñy āīōōōē yēāīāīōīā ōēīā «āççāō» <p> ēēē «āēīē ōāēñōā» <ab>, ā iā iāīāīōīō. 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 ñéááóáð òíéáíéçèðíááòú, ïñéíéúéó ííè íá íðéíááéáæáð é íàðáðèáéó èñòí-íééá (íáíðèíáð, ðáááéóíðñééá íðéíá-áíéý è ñííñéè á óáéñòá éðèðè-áñéíáí èçááíéý). Á ðáíéáó éííéðáòíáí íðíáéòá ýè ñíéñéè ííáóó áúòú ñóúáñòááíí ðáñæéðáíú è óòí-íáíú, á ðáçóéúòáðá -ááí -éñéí ýéáíáíóíá, éíòíðúá ííáóó íáðáéðáúéááòúñý ñ èéíááéñòè-áñééíè ñòðóéóóðáíè, çíá-èòáéúíí ñíéðáúááòñý.

Ðàçíáðéá óáéñòíá BFM ñííòááòñòáóáð -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 èèè çíáéíá íáðáííñá), à óáéæá ýéáíáíóíá, éíòíðúá ííáóó íáðáéðáúéááòúñý ñí ñòðóéóóðíé íðááéíæáíéè (íáíðèíáð, «íóñòúá» ýéáíáíóú <lb/> «ííááý ñòðíéá» èñííéúçóáòñý áíáñòí <l> «ñòèò», à ýéáíáíó òèòèðíááíéý <q> ðáññíáòðèááòñý éáè áðáíéòá íðááéíæáíéý). Ýóí ïíçáíéééí ðááéèçíááòú á íèàòóíðíá TXM ýóðáéòéáíúé áéáíðèòí òíéáíéçáòéè áéý óáéñòíá BFM.

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

Ñíèñíè 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. – Ááðáñ á Éíóáðíáò 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

