

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 yēāēōðīīīīāī ēçāāīēy òāēñōīā, īðēīāēāæāūēō ē ðīīāō īēñūīāīīāī īāñēāāēy ðāçēē÷īūō ñōðāī ē ēōēūōðð. 150 īðīāēōīā īðāāñōāāēāīū īā ñāēōā TEI, ā āāēñōāēōāēūīñōē yōī ÷ēñēī īāæāō 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 ñīāōēðēēāōēē ODD, ā òāēæā īāçāāēñēīñōū īō ēīīēðāōīīē īēāðōīðīū ēēē īðīāðāīīīāī īðīāēōā. Āīāñōā ñ òāī ÷ðāçāū÷āēīāy āēāēīñōū TEI ñīçāāāō īīðāāāēāīūā 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, īāæāō āūōū òðōāīī ñīāīāñōēōū ñ ēñīēūçīāāīēāī ēīñōðōīāīōīā āāōīāōē÷āñēīē ēēīāāēñōē÷āñēīē ðāçīāōēē (ðīēāīēçāōēē, ēāīīāðēçāōēē, īðōīēīāē÷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 īōēē÷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 Iāðīāēēē āīāēēçā īāúāīūō ēīðīōñīā òāēñōīā. Āñēāā çā ēāēñēēīīāððēāē ē ñōāðēñōē÷āñēēī āīāēēçīī òāēñōā òāēñōīāðōðēy īðāāēāāāāð ñōāðēñōē÷āñēē īāīñīāāīūā Iāðīāū 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 īōāēēēāðēyō [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 ðāçðāāāðūāþūēēñy ēāāīðāðīðēāē ICAR Iāēēīāēūīīāī òāīðā īāó÷īūō ēññēāāīāāīēē

Ōāiōēē (CNRS) ē Eēīīnēīāī ōīēāāōñēōāōā. Āāçā iāñ-ēōūāāāō 75 ōāēñōīā iāuēī iāuāīīī āīēāā 3 500 000 ōāēñōīōīōī. Eñōī-īēēāīē BFM ā īñīīāīīī yāēyōñy āāōīōēōāōīūā ēōēōē-āñēēā ēçāāīēy, īāīāēī ā īñēāāīāā āōāīy ðaçāēāpōñy nīāñōāāīīūā ēçāāīēy, īēōāpūēāñy iā ēēīāāēñōē-āñēē āūāāōāīīūā ōōāīñēōēīōēē īōēāēīāēūīōō ðōēīēēāē. Ā ēā-āñōāā īōēīāōā īīāēī īōēāāñōē ēīōāōāēōēāīīā ēçāāīēā āīīēīīāīī āīīāīā 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 nīīōāāōñōāēē nī nīāōēōēēāōēēāē, ðaçōāāīōāīīē āēy īōæā īōīāēōā [Guillot et al. 2010] n ō-āōīī īāōñīāēōēāū ēēīāāēñōē-āñēīāī āīāēēçā. Iēāōōīōīā TXM iā ōīēūēī nēōæōō āēy ēīōīōñā BFM nōāāñōāīī 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 nēīāēīūō çāāā- īōē ēñīīēūçīāāīēē nōāāñōā āāōīāōē-āñēīāī ēēīāāēñōē-āñēīāī āīāēēçā ēīōīōñīā īōēīāīēōāēūīī ē ōāēñōāī, āēēp-āpūēī āēōāīēōp ðāāāēōīōñēōp ðaçīāōēō, yāēyāōñy ēīōōāēōīāy ēāāīōēōēēāōēy nēīā (ōīēāīēçāōēy) ē īōāāēīāēāīēē, ē ēīōīōūī yōīō āīāēēç āīēæāī īōēīāīyōūñy āāç īīōāōē nāīīē ðāāāēōīōñēīē ðaçīāōēē.

Nēāāōpūēē īōēīāō ðaçīāōēē, nīāāōæāūēē īōāāēīāēāīīūē ðāāāēōīōīī ōōāāīāīō ōāēñōā iā iāñōā ēāēōīū, iā-ēīāpūāēñy ā ēīīōā īāīīāī nēīāā ē çāēāī-ēāāpūāēñy iāñēīēūēēīē nēīāāīē īīçæā, āāñīēpōīī ēīōōāēōāī n ōī-ēē çōāīēy ðāēīīāīāāōēē TEI, īāīāēī āāñūā nēīāāī āēy ōīēāīēçāōēē n ō-āōīī nīīāāīīñōāē 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 nāōūāçīūā īōīāēāīū āīçīēēāpō īōē ðaçīāōēā īōāāēīāēāīēē, īñīāāīī ā nōēōīōāīōīūō ōāēñōāō, āāā īāōāēōāūēāāīēā iāōōē-āñēīē ē nēīōāēñē-āñēīē nōōōēōōōū āñōōā-āāōñy ī-āīū -āñōī.

ðaçōīāāōñy, īīāēī «īōōēēūōōīāāōū» āñā yēāīāīōū, īāōāēōāūēāāpūēāñy n īñīīāīūīē ēēīāāēñōē-āñēīēē nōōōēōōōāīē (nēīāāīē ē īōāāēīāēāīēyīē), īāīāēī yōī īīæāō īōēāāñōē ē īōāōā nōūāñōāāīīē ēīōīōīāōēē āēy çāīōīñīā ē āēçōāēēçāōēē (īāīōēīāō, ī ōīī, īīāāāōāāēāñū ēē nēīāīōīōīā ðāāāēōīōñēīē īōāāēā).

Nīçāāīēā āēāīōēōīā ōīēāīēçāōēē, ēīōīōūē ēīōōāēōīī īāōāāōūāāē āū ēpāīē ōāēñō n āēōāīēīē ðāāāēōīōñēīē ðaçīāōēīē ā nōāīāāōōā TEI XML, īōāāñōāāēyāōñy īōāēōē-āñēē īāāīçīīāēīūī. Ōāī iā īāīāā, īīāēī āīāēōūñy āīīēā ōāīēāōāīōēōāēūīōō ðaçōēūōāōīā īōē ōñēīāēē, -ōī ðaçīāōēā ēñōīāīīāī āīēōīāīōā īōāā-āāō ðyāō īōīñōūō īōāāēē. Īāīōēīāō, «ōyāē, ðāñīīēīāēāīīūā āīōōōē nēīā, āīēāēū āūōū -āōēī ēāāīōēōēōēōīāāīū» ēēē «āñēē ðaçīā-āīīūē nāāīāīō ōāēñōā iā-ēīāāōñy āīōōōē īāīīāī nēīāā ē çāōāāōūāāāō īāñēīēūēī īīñēāāōpūēō, āāī īāīāōīāēīī ðaçāāēēōōū».

Ōāēæā āīçīīāēīī nīñōāāēōū nīēñēē ōyāīā TEI ā çāāēñēīīñōē īō ēō īīçōēē ā ēēīāāēñōē-āñēīē ēāōāōōēē ōāēñōā ā ðāīēāō īōāāēūīīāī īōīāēōā. Īāīōēīāō, īōāāēīāēāīēy ðāñīīēāāāpōñy āīōōōē yēāīāīōīā ōēīā «āaçāō» <p> ēēē «āēīē ōāēñōā» <ab>, ā iā īāīāīōīō. Īāēīōīōūā ōyāē īīāēīī n-ēōāōū

ýéáéááéáíóíúíè ñéíáó (íáíðèíáð, <abbr>, <num>, <pc>), à íáñéíéúéí ýéáíáíóíá ïí-òè áñáááà ðáñííéáááòñý áíóóðè ñéíáà (<am>, <c>, <ex>). Ðýá ýéáíáíóíá ñíááðæàð ñááíáíóú óáéñòà, éíòíðúá íá ñéááóáò òíéáíéçèðíáòú, ïñéíéúéó ííè íá íðéíááéáæàð é íàðáðèáéó èñòí-íééà (íáíðèíáð, ðáááéóíðñééá íðéíá-áíéý è ñííñéè á óáéñòà éðèðè-áñéíáí èçááíéý). Á ðáíéáò éííéðáòíáí íðíáéòà ýðè ñíéñéè ííáóó áúòú ñóúáñòááíí ðáñæðáíú é óòí-íáíú, á ðáçóéúòàðà -ááí -éñéí ýéáíáíóíá, éíòíðúá ííáóó íáðáéðáúéáàòúñý ñ èéíááéñòè-áñééíè ñòðóéóóðáíè, çíá-èòáéúíí ñíéðáúááòñý.

Ðàçíáòéà óáéñòíá BFM ñííòááòñòáóáò -áòéí ñòíðíóéèðíááííé ñíáòéòééáòéè íðéíáíéý ðáéííáíáòéé TEI, íòðáæáííé á áíéóíáíóàòéè ODD [Guillot et al. 2010]. Ááííáý ñíáòéòééáòéý áéèp-ááò íðááééà èñííéúçíááíéý áíóóðèñéíáíúó óýáíá (íáíðèíáð, èñíðááéáííúó ðáááéóíðíí áóéá èéè çíáéíá íáðáííñá), à óáéæá ýéáíáíóíá, éíòíðúá ííáóó íáðáéðáúéáàòúñý ñí ñòðóéóóðíé íðááéíæáíéè (íáíðèíáð, «íóñòúá» ýéáíáíóú <lb/> «ííááý ñòðíéà» èñííéúçóðòñý áíáñòí <l> «ñòèò», à ýéáíáíó òèòèðíááíéý <q> ðáññíáòðèááòñý éáè áðáíéòá íðááéíæáíéý). Ýóí ííçáíéééí ðááéèçíáòú á íèàòòíðíá TXM ýóðáéòéáíúé áéáíðèòí òíéáíéçáòéè áéý óáéñòíá BFM.

Á ñíñéááíáá áðáíý áúé óñíáòíí íðíááááí ðýá óáñòíá ïí áááíòáòéè ááíííáí áéáíðèòíá é áíéóíáíóàí TEI XML, ñíááíòíáéáííúí á áðóáéò íðíáéòáò. Ñðááè íéð ííæíí íáçááòú «Áèðóóáéúíóð áéáééíòáéó 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 íðéíáíééè ñíáòéáéúíúó óèéúòðíá XSL íá áðíáá è íá áúóíáá íðíóááóðú òíéáíéçáòéè. «Áóíáííé» óèéúòð óááéýáò óýáé, éíòíðúá íá íðááñòááéýðò éíóáðáñá áéý ýéñíéóáòáòéè ñ ñííúúð TXM, à óáéæá óíðíúááò è íðíáéèçóáò íáéíòíðúá ñéíáíúá ñòðóéóóðú ýéáíáíóíá XML. «Áúóíáííé» óèéúòð ííçáíéýáò 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 ðáñííéááòúñý íáæáò ááí íá-àéíí è éííóíí).

Á óáéñòá áíééááá íú íðéááááí éííéðáóíúá íðéíáòú ðááéèçíááííúó ðáðáíéè è íðíááíííñòðèðóáí ðáçóéúòáòú, éíòíðúá ííæíí ííéó-èòú íðè ýéñíéóáòáòéè óáéñòíá BFM è áðóáéò íðíáéòíá ñ ñííúúð íèàòòíðíú TXM. Áíéáá ííáðíáíí ðáçèè-íúá áéáíðèòíú òíéáíéçáòéè, ðááéèçíááííúá íá íèàòòíðíá TXM, ñíèñáíú 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. – 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

