

Interface Web pour l'annotation morpho-syntaxique de textes

Thierry Hamon^{1,2}

(1) LIMSI, CNRS, Universit  Paris-Saclay, Bat 508, rue John von Neumann, Campus Universitaire, 91405 Orsay, France

(2) Universit  Paris 13 - Sorbonne Paris Cit , 99 avenue J.B. Cl ment, 93430 Villetaneuse, France

hamon@limsi.fr

R SUM 

Nous pr sentons une interface Web pour la visualisation et l'annotation de textes avec des  tiquettes morphosyntaxiques et des lemmes. Celle-ci est actuellement utilis e pour annoter des textes ukrainiens avec le jeu d' tiquettes Multext-East. Les utilisateurs peuvent rapidement visualiser les annotations associ es aux mots d'un texte, modifier les annotations existantes ou en ajouter de nouvelles. Les annotations peuvent  tre charg es et export es au format XML TEI, mais aussi sous forme tabul e. Des scripts de conversion de format et de chargement dans une base de donn es sont mis   disposition.

ABSTRACT

Web interface for the morpho-syntactic annotation of texts

We present a Web interface for visualizing and annotating texts with POS tags and lemma. This interface is currently used for the annotation of Ukrainian texts with the Multext-East POS tagset. The users have a fast access to the annotations associated with words from a text. They can also modify existing and add new annotations. The annotations can be loaded or exported in the TEI XML format and tabular separated format. Several scripts for loading and converting the data are also available.

MOTS-CL S : Annotation morphosyntaxique, Lemmatisation, Multext-East, Ukrainien.

KEYWORDS : Morpho-syntactic annotation, Lemmatization, Multext-East, Ukrainian.

1 Introduction

La mise au point de m thodes de TAL sur des langues peu outill es n cessite de constituer des corpus annot s manuellement. Ainsi, dans le cadre du d veloppement d'outils pour l'ukrainien, et notamment un  tiqueteur morphosyntaxique, nous souhaitons disposer de textes annot s morphosyntaxiquement avec le jeu d' tiquettes Multext-East (Erjavec, 2012). Ainsi, nous avons d velopp  une interface Web pour l'annotation des mots d'un texte avec des  tiquettes morphosyntaxiques et des lemmes, mais aussi la correction de ces informations lorsque le texte est pr -annot .

Notre objectif est d'une part, de faciliter l'utilisation du jeu d' tiquettes Multext-East, et d'autre part, de limiter les actions de l'utilisateur afin de r duire le temps d'annotation. En effet, ce jeu d' tiquette morphosyntaxique est complexe puisqu'il propose 12 cat gories grammaticales, mais aussi jusqu'  10 traits morphologiques pour les adjectifs et 11 valeurs de traits possibles pour d crire les types de pronom. Il s'agit donc de ne pr senter que les traits morphologiques et les valeurs pertinentes pour

une cat gorie donn e. L'utilisateur doit donc  galement avoir une vue synth tique des annotations et la visualisation des annotations associ es   un mot doit  tre rapide.

2 L'interface d'annotation

Afin de r pondre aux objectifs pr sent s ci-dessus, nous avons d velopp  une interface d'annotation s'appuyant sur des technologies Web (XHTML, PHP, AJAX). Une base de donn es est  galement utilis e pour le stockage des annotations et la description du jeu d' tiquettes. L'adaptation de l'interface   un autre jeu d' tiquettes ne n cessite donc que la modification de la description du jeu d' tiquettes.

Les textes au format tabul  ou XML conforme   la TEI¹ (Wittern *et al.*, 2009) pr -annot s ou non peuvent  tre charg s dans la base   l'aide de scripts Perl. Les textes annot s peuvent  tre export s dans les m mes formats.

La figure 1 pr sente une capture d' cran de l'interface d'annotation. Le document   annoter est affich  une seule fois dans l'interface et les annotations d j   associ es aux mots apparaissent dynamiquement. Une vue synth tique des annotations est propos e lorsque l'utilisateur passe avec la souris sur un mot (par exemple, **використані**   la figure 1). Les annotations peuvent  tre modifi es ou ajout es en cliquant sur le mot concern  (par exemple **приймає** visible   droite sur la figure 1).

L'interface est t l chargeable   l'adresse suivante : <https://perso.limsi.fr/hamon/Rada/index.php>.

The screenshot shows a web interface for document annotation. At the top, there is a navigation menu with links: Accueil, A-G, H-N, O-T, U-Z, 0-9, Others, and Statistiques. Below the menu, the document title is 'SUBCUINJECTIONSITES_UKR_mergeTagged'. There are buttons for 'export as tab' and 'export as XML', and a 'Statistiques' link. The main content area displays a list of text lines with annotations. A tooltip is visible over the word 'використані' in line 2, showing its morphological analysis: 'використані (4 poss. - notDisamb.)', 'використаний - Ar--pns-er', 'використаний - Ar--paslep', 'використані - Ar--pl--er', and 'використаний - Ar--pl--er'. On the right side, there is a 'Statistiques' panel for the word 'приймає (2 poss.)'. It contains radio buttons for 'приймати - Vmp1p2p' (selected), 'приймає - Vmp1p2p', and 'New tag or lemma'. Below this, the POS tag is shown as 'V m p i p 2 p -' and the lemma is 'приймати'. There are 'OK' and 'Cancel' buttons at the bottom of the panel.

FIGURE 1 : Exemple de visualisation d'un document en cours d'annotation.

Remerciements

Ce travail a  t  financ  par l'action incitative *Outiller l'Ukraine* du LIMSI-CNRS. Nous remercions  galement Natalia Grabar et Anastasiia Kuznietsova pour leur commentaires sur les premi res versions de l'interface.

¹<http://www.tei-c.org>

Références

ERJAVEC T. (2012). Multext-east : Morphosyntactic resources for central and eastern european languages. *Language Resources and Evaluation*, **46**(1), 131--142.

WITTERN C., CIULA A. & TUOHY C. (2009). The making of tei p5. *Literary and Linguistic Computing*, **24**(3), 281--296.