



HAL
open science

Trameur: A Framework for Annotated Text Corpora Exploration

Serge Fleury, Maria Zimina

► **To cite this version:**

Serge Fleury, Maria Zimina. Trameur: A Framework for Annotated Text Corpora Exploration. Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: System Demonstrations, Dublin City University, Aug 2014, Dublin, Ireland. pp.57-61. hal-01223725

HAL Id: hal-01223725

<https://u-paris.hal.science/hal-01223725>

Submitted on 1 Mar 2023

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.

Trameur: Framework for Annotated Text Corpora Exploration

Serge Fleury

Sorbonne Nouvelle – Paris 3

SYLED-CLA2T, EA2290

75005 Paris, France

`serge.fleury@univ-paris3.fr`

Maria Zimina

Paris Diderot – Sorbonne Paris Cité

CLILLAC-ARP, EA 3967

75205 Paris, cedex 13, France

`maria.zimina@eila.univ-paris-diderot.fr`

Abstract

Corpus resources with complex linguistic annotations are becoming increasingly important in the work of language specialists with different professional backgrounds. Regardless of their computer skills, they often need to perform extensive corpus research, including Natural Language Processing (NLP), statistical modelling and data visualisation. Our software system, called Trameur, aims at making these analyses possible within a single graphical user interface. It relies upon a specific data modelling framework presented in this paper.

1 Introduction

Annotated text corpora are currently used by language specialists with different professional backgrounds. As a rule, linguists are familiar with standard text processing tools, software engineers use Natural Language Processing (NLP) tools and statisticians work with statistics software. Corpus resources are getting more and more complex and consist of several layers of annotations for multifaceted data collections. To face these challenges, we develop an integrated system for complex annotated text data analyses called Trameur (Fleury, 2013a).

Trameur manages multiple layers of linguistic annotations and allows exploring complex linguistic features and embedded dependence relations. Additionally, Trameur incorporates advanced NLP processing, statistical analysis and text mapping features. Thus, Trameur has been developed to allow multiple corpus analyses within a single graphical user interface.

This framework system was developed to make complex analyses of annotated corpora accessible to any corpus linguist, without extensive programming skills and knowledge of statistical modelling tools.

2 Corpus research with Trameur

Trameur was successfully tested for monolingual text processing within several research projects in corpus linguistics and discourse analysis (Branca-Rosoff et al., 2012; Née et al., 2012).

On-going research demonstrates its potential for processing parallel and comparable text data in distant languages (Zimina and Fleury, 2014).

2.2 Exploring dependence relations

Dependence relations emerging from linguistic annotations can be explored, filtered and displayed in several ways:

- Using graphs of dependence relations
- Using concordances and context return
- Using co-occurrence statistics
- Using visual tools (graphs) and text mapping

A detailed description of these features with related screenshots is available on-line (Fleury, 2013).

Figure 3 displays a sample graph of a double dependence relation from Rhapsodie. It is set by the regular expression OBJ|SUB (OBJ or SUB), where the relation target for the lemma is “aimer” and the annotation n°9 is ROOT:

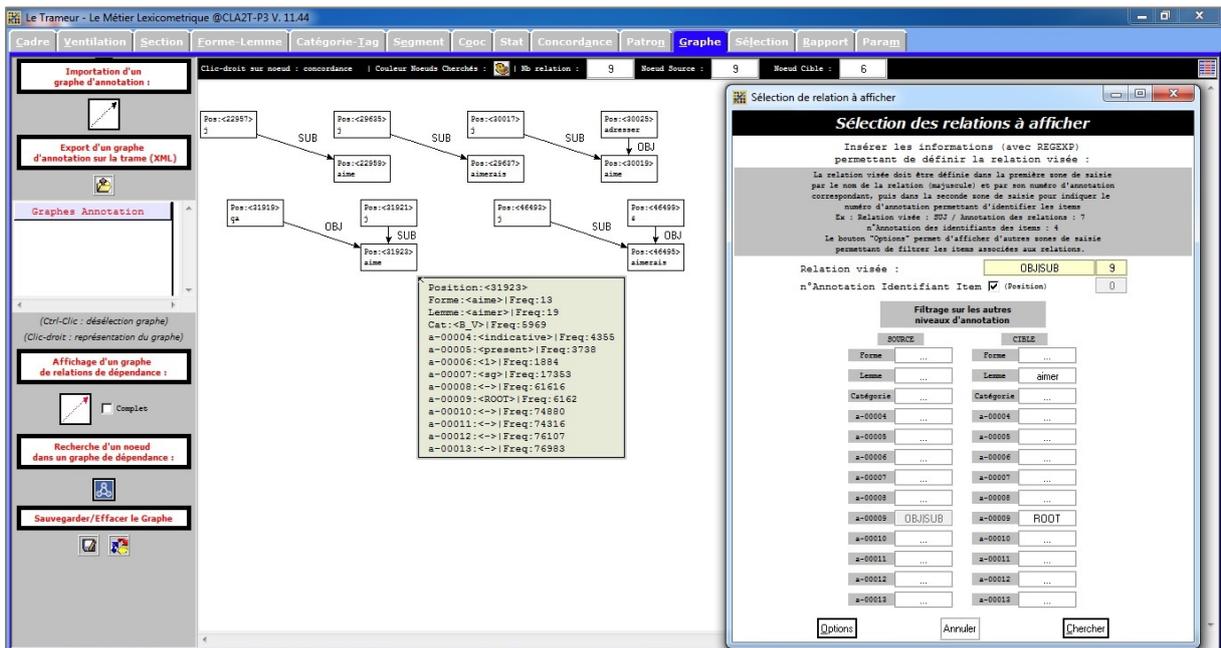


Figure 3. Graph of a double dependence relation

Figure 4 displays the nodes of the same graph in a concordance window. Selected annotations are automatically highlighted.

3 Conclusion

Trameur is a tool for exploration of complex multi-layer linguistic corpora with different annotations. It is built upon a Thread/Frame data model using XML. The software is distributed with a graphical user interface. A reference package for Windows is available for free download from the official website: <http://www.tal.univ-paris3.fr/trameur>.

Several other systems have been already developed for processing annotated corpora, for example: ANNIS² and Macaon³. However, the novelty of Trameur consists in expanding a multi-layered data model to all stages of corpus exploration, including text mapping features and statistical analysis within a single graphical user interface.

Following this integrated approach, Trameur can be used, for example, to build complex linguistic units, analyse their dependence relations and reveal their characteristic attractions using text statistics.

² ANNIS (ANNotation of Information Structure): <http://www.sfb632.uni-potsdam.de/annis/>

³ MACAON Project: <http://macaon.lif.univ-mrs.fr>

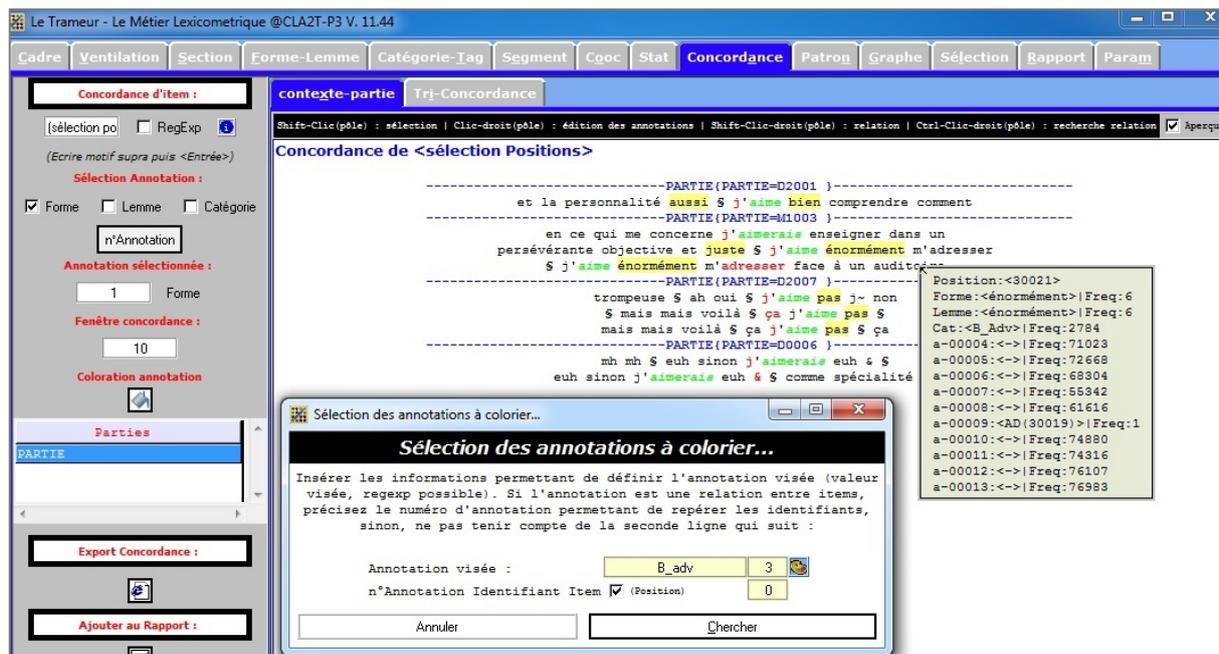


Figure 4. Access to dependence relations in a concordance window

Thus, the use of the Thread/Frame data model implemented in Trameur allows different combinations of analyses when processing multifaceted linguistic annotations. We believe that this framework offers new perspectives for corpus research.

References

- Sonia Branca-Rosoff, Serge Fleury, Florence Lefevre and Mat Pires. 2012. *Discours sur la ville. Corpus de Français Parlé Parisien des années 2000 (CFPP 2000)*. Sorbonne nouvelle – Paris 3. Online publication: <http://cfpp2000.univ-paris3.fr/CFPP2000.pdf>
- Serge Fleury. 2013a. *Le Trameur. Propositions de description et d'implémentation des objets textométriques*. Sorbonne nouvelle – Paris 3. Online publication: <http://www.tal.univ-paris3.fr/trameur/trameur-propositions-definitions-objets-textometriques.pdf>
- Serge Fleury. 2013b. *Annotations Rhapsodie pour le Trameur*. Sorbonne nouvelle – Paris 3. Online publication: <http://www.tal.univ-paris3.fr/trameur/bases/rhapsodie2trameur.pdf>
- Emilie Née, Erin MacMurray and Serge Fleury. 2012. Textometric Explorations of Writing Processes: A Discursive and Genetic Approach to the Study of Drafts. *Journées internationales d'Analyse statistique des Données Textuelles (JADT 2012)*. Liège (Belgium), 13-15 June 2012.
- Kim Gerdes, Sylvain Kahane, Anne Lacheret, Arthur Truong, and Paola Pietrandrea. 2012. *Intonsyntactic data structures: The Rhapsodie treebank of spoken French*. *Proceedings of the Linguistic Annotation Workshop, COLING 2012*. Jeju, Republic of Korea, July 2012.
- Emilie Née, Erin MacMurray and Serge Fleury. 2012. *Textometric Explorations of Writing Processes: A Discursive and Genetic Approach to the Study of Drafts*. *Journées internationales d'Analyse statistique des Données Textuelles (JADT 2012)*. Liège (Belgium), 13-15 June 2012.
- Maria Zimina and Serge Fleury. 2014. *Approche systémique de la résonance textuelle multilingue*. *Journées internationales d'Analyse statistique des Données Textuelles (JADT 2014)*. Paris, 3-6 June 2014.