# Trameur: A 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. 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

Treebanks, or parsed text corpora that annotate syntactic or semantic sentence structure, are widely appreciated in language studies. These corpora are getting more and more complex and consist of several layers of annotations. Such multifaceted data collections present many challenges for the development of specific NLP (Natural Language Processing) tools, statistical modelling and data visualisation techniques. In terms of statistical text analysis, most research has been carried out using either unannotated texts or tagged texts, annotated with parts of speech only (Lebart et al., 1998). As for corpora with richer annotations such as clauses, grammatical functions and dependency links, recent research projects have been focused on the development of treebank querying methods with fast search algorithms (Cunningham et al., 2002; Mírovský and Ondruška, 2002; Götze and Dipper, 2006). However, for annotated corpus analysis, it is often necessary to precisely locate different types of linguistic schemes under study, with simultaneous access to multiple corpus layers and their interactions statistics.

To face these challenges, we develop an integrated system for statistical analysis of annotated text data called Trameur (Fleury, 2013a). Trameur manages multiple layers of linguistic annotations and allows statistical exploration of complex linguistic features and embedded dependency relations. Additionally, Trameur incorporates advanced NLP processing and text mapping features. This framework has been developed to allow multiple corpus analyses within a single graphical user interface, accessible to any corpus linguist, without extensive programing 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).

## 2.1 Data structures for annotated text processing

Several software packages dedicated to multi-level text processing apply pipeline architecture to support reproducible analyses of text annotations stored in specific formats (Eberle et al., 2012; Eckart et al., 2012). We provide here a description of the specific data structures used by Trameur (text segmentation, partition, statistical tables, etc.) to implement incremental textual resources for treebanks.

In Trameur, a formal representation of text segmentation relies upon two systems of units: text containers and contents (or items) (Fleury, 2013a). Text containers reflect the structure of different corpus parts, while contents represent systems of textual units (lexemes, graphemes, etc.) and their frequency variations which can be examined on a statistical basis. Based on computerized procedures, the systems of containers and contents allow automatic counts in the form of large statistical tables. These tables can be further processed by quantitative methods to detect salient points of statistical distribution in the corpus and build the premises of variation analysis.

Containers and contents are essential elements of double-tracking text segmentation process: (1) Cutting text into containers (parts, textual zones, sections, chapters, paragraphs, sentences, etc.) allows creating a system of text spans that can be further compared on a quantitative basis. (2) Text segmentation into items (or contents) isolates distinct textual units within a text corpus. This text segmentation is closely interconnected with the other two operations: annotation of items (attributing a distinct tag, such as lemma or grammatical category, to each isolated item), and typing of isolated items (judged similar because of their intrinsic nature or on the basis of corpus annotations). Thus, each segmented text is a succession of isolated items that can always be reunited to restore the initial text sequence. Each item is further attached to a given type.

After text segmentation in Trameur, the items are numbered to create a system of coordinates, where each item is spotted by its corresponding sequence number. In Trameur, this system is called Thread (in French: *trame*). The same system of coordinates is used to define and locate text spans (or textual zones) formed by a series of consecutive items between positions $x_1$ and $x_2$, or by a certain number of textual zones of this type. The definition of a Thread allows describing various systems of textual zones corresponding to the text containers (corpus parts, sections, chapters, paragraphs, sentences, etc.). The descriptions of the containers are united within a specific data structure called Frame.

The Thread/Frame structure allows defining Selection as any object corresponding to the selected subset of Thread items used in the corpus exploration. There are several types of Selection: (1) Selection corresponding to the contents (certain Thread items, such as occurrences of the same token, lemma, or items corresponding to the same regular expression operated on one or several annotation levels); (2) Selection corresponding to certain containers (textual zones, paragraphs, sections, etc.) composed of contiguous item sequences, selected according to their positions in the text; (3) Selection resulting from highlighting operations performed by a human expert on the basis of some complex criteria, possibly difficult to describe on a formal basis; (4) A textual zone discovered following the results of a statistical calculation; (5) A sub-set of textual units selected following a statistical calculation that highlights a specific distribution of items (contents) within a corpus.

The following demonstration of these principles is based on a data sample from the Rhapsodie corpus of spoken French annotated for prosody and syntax (Gerdes et al., 2012). The data file is available online in tabular form and can be displayed in a spreadsheet.[1] The data set used in this demonstration comes from Rhapsodie v1 micro-syntax annotations. This first release is limited to 13 annotations.[2] The corpus data is structured as follows: all Rhapsodie texts are first identified by codes (column 1). Each text is further divided into separate units (column 2) and segmented into tokens (column 3). The remaining columns display linguistic annotations of the tokens, including microsyntax (rection, dependency and constituancy) (Fleury, 2013b).

Figure 1 shows the Rhapsodie data set transformed from the spreadsheet format into a Trameur base file using regular expressions. The data structure is composed of two parts: (1) a Thread, which is a list of items with position identifiers; (2) a Frame, which is a list of corpus partitions defined on the Thread. Each partition has a name and a list of named constituents identified through their first and

---

[1] The Rhapsodie corpus is available for download from: `http://www.projet-rhapsodie.fr`
[2] Rhapsodie v3 is the latest online release with micro-syntax, macro-syntax and prosody annotations. It is currently under study in a number of on-going research projects that will be revealed in forthcoming publications.

last token positions on the Thread. Thus, each annotated token from the Rhapsodie corpus becomes an item identified by its position on the Thread.

Dependency relations among items are annotated as follows: RELATION(TARGET), where: RE-LATION is a character string corresponding to the name of the relation; TARGET is a numerical value of the position identifier on the Thread (see Figure 1). All linguistic annotations from the Rhapsodie sample file are preserved in the Trameur base file. These annotations can be displayed and edited in the graphical user interface (see Figure 2): a-00001:Token, a-00002 Lemma, a-00003 P.O.S, a-00004 Mode, a-00005 Tense, a-00006 Person, a-00007 Number, a-00008 Gender, a-00009 Type_rection(Gov_rection), a-00010 Type_para(Gov_para), a-00011 Type_inher(Gov_inher), a-00012 Type_junc(Gov_junc), a-00013 Type_junc-inher(Gov_junc-inher).



Figure 1: Rhapsodie data in a Trameur base file.

## 2.2    Exploration of dependency relations



Figure 2: Trameur. Graph of a double dependency relation.

Dependency relations emerging from linguistic annotations can be explored, filtered and displayed using dependency graphs, text maps with context return, concordances, co-occurrence statistics, etc. All these features with related screenshots are presented in the on-line user guide.[3]

Figure 2 displays a sample graph of a double dependency 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 Type_rection(Gov_rection) is ROOT. For each item in context, a pop-up window displays all annotation levels with corresponding tags (or "no tag" signs <—>) and their frequencies (Freq) on a given annotation level.

Figure 3 displays the nodes of the same graph in a concordance window. A selection tool is used to highlight annotation n°3 in concordance lines (value = B_adv).
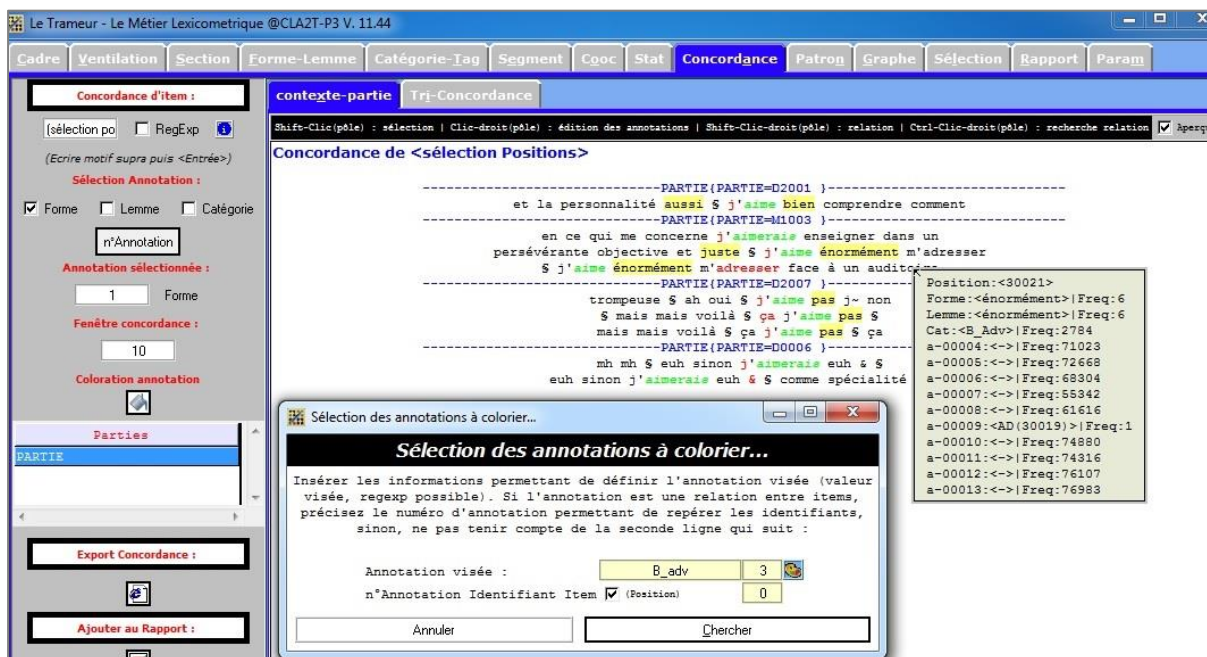
---

Figure 3: Trameur. Concordance window with selected dependency relations.

Figure 4 displays a co-occurrence graph for the lemma "vouloir" constrained by the dependency RELATION set by the regular expression SUB|AD (SUB or AD).
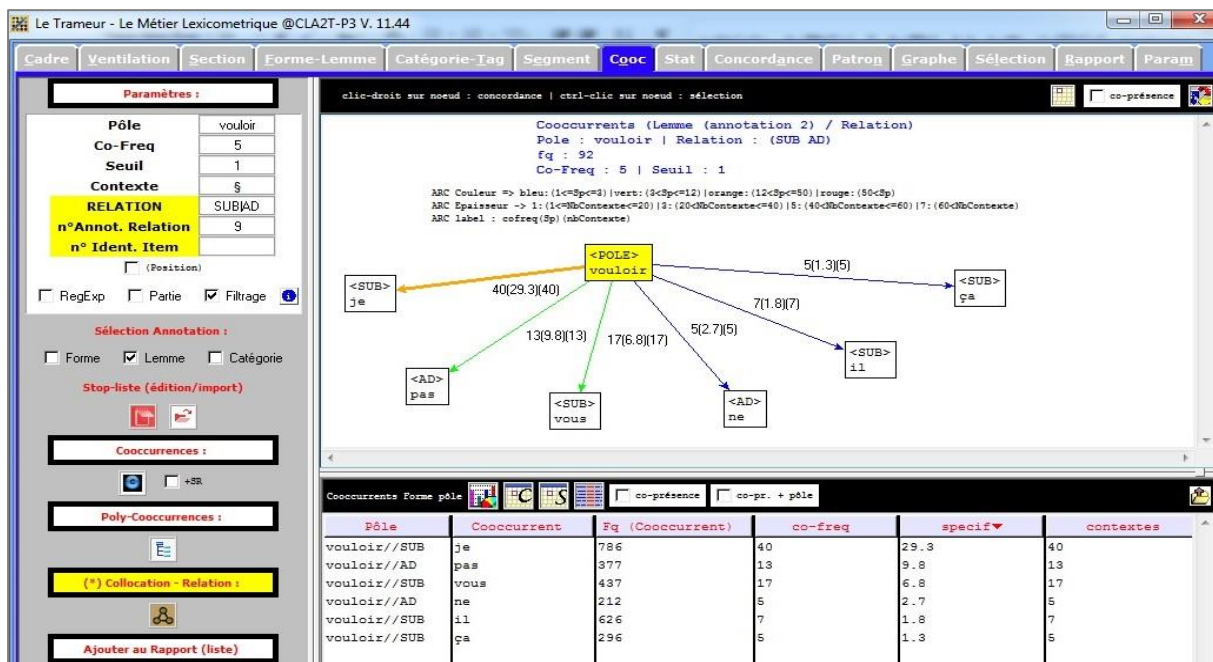


Figure 4: Trameur. Co-occurrence graph constrained by selected dependency relations.

## 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: PDT2.0[4], GATE[5], ANNIS[6], Macaon[7]. 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 of dependency relations within a single graphical user interface.

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

In Trameur, quantitative analyses are applied to textual resources, built upon a Thread and an evolutionary Frame. These analyses can progressively enrich these resources with new text divisions or annotations, without deleting previous information. Following this approach, it becomes possible to implement incremental textual resources for treebanks.

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

# References

Sonia Branca-Rosoff, Serge Fleury, Florence Lefeuvre 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`

Hamish Cunningham, Diana Maynard, Kalina Bontcheva and Valentin Tablan. 2002. *GATE: A Framework and Graphical Development Environment for Robust NLP Tools and Applications. Proceedings of the 40th Anniversary Meeting of the Association for Computational Linguistics (ACL'02)*. Philadelphia, July 2002.

Kurt Eberle, Kerstin Eckart, Ulrich Heid and Boris Haselbach, 2012. *A Tool/Database Interface for Multi-Level Analyses. Proceedings of the eighth conference on International Language Resources and Evaluation (LREC 2012)*. Istanbul (Turkey), May 2012.

Kerstin Eckart, Arndt Riester and Katrin Schweitzer, 2012. *A Discourse Information Radio News Database for Linguistic Analysis. Linked Data in Linguistics*. Springer.

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` (v1 base) `http://www.tal.univ-paris3.fr/trameur/bases/rhapsodie2trameur-v3.pdf` (v3 base)

Kim Gerdes, Sylvain Kahane, Anne Lacheret, Arthur Truong and Paola Pietrandrea. 2012. *Intonosyntactic data structures: The Rhapsodie treebank of spoken French. Proceedings of the Linguistic Annotation Workshop, COLING 2012*. Jeju, Republic of Korea, July 2012.

Michael Götze and Stefanie Dipper. 2006. *ANNIS: Complex Multilevel Annotations in a Linguistic Database. Proceedings of the 5th Workshop on NLP and XML: Multi-Dimensional Markup in Natural Language Processing (NLPXML-2006)*. Trento (Italy), April 2006.

Ludvic Lebart, André Salem and Lisette Barry. 1998. *Exploring Textual Data*. Kluwer Academic Publishers.

Jiří Mírovský and Roman Ondruška. 2002. *NetGraph System: Searching through the Prague Dependency Treebank. Prague Bulletin of Mathematical Linguistics*, 77, MFF UK, Prague, Czech Republic, Prague, 2002.

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), 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, June 2014.

---

[4] PDT (The Prague Dependency Treebank 2.0): `ufal.mff.cuni.cz/pdt2.0/`
[5] GATE (General Architecture for Text Engineering): `http://gate.ac.uk/gate/`
[6] ANNIS (ANNotation of Information Structure): `http://www.sfb632.uni-potsdam.de/annis/`
[7] MACAON Project: `http://macaon.lif.univ-mrs.fr`