



HAL
open science

Annotation micro- et macrosyntaxique manuelle et automatique de français parlé

Sylvain Kahane, Henri-José Deulofeu, Kim Gerdes, Alexis Nasr, André Valli

► **To cite this version:**

Sylvain Kahane, Henri-José Deulofeu, Kim Gerdes, Alexis Nasr, André Valli. Annotation micro- et macrosyntaxique manuelle et automatique de français parlé. Journée Floral, 2017, Orléans, France. halshs-01740461

HAL Id: halshs-01740461

<https://shs.hal.science/halshs-01740461v1>

Submitted on 19 Sep 2019

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.

Annotation micro- et macrosyntaxique manuelle et automatique de français parlé

Sylvain Kahane¹, José Deulofeu², Kim Gerdes³, Alexis Nasr², André Valli²

¹ Modyco, Université Paris Nanterre & CNRS

² LIF, Aix-Marseille Université & CNRS

³ LPP, Université Paris 3 Sorbonne Nouvelle & CNRS

1. Le contexte : le projet Orféo

Ce résumé présente les principaux choix d'annotation syntaxique fait dans le cadre du projet Orféo (ANR 2012-2017) visant à fournir un large treebank de français écrit et oral interrogeable en ligne (projet-orfeo.fr). L'annotation est réalisée selon un processus de bootstrapping usuel : un corpus d'amorçage est réalisé par une annotation manuelle, puis un analyseur syntaxique est entraîné sur ce corpus, une nouvelle portion de corpus est analysée automatiquement, puis corrigée manuellement, un nouvel analyseur est entraîné et ainsi de suite. Le corpus comprend plusieurs millions de mots et seule une partie du corpus est corrigée manuellement. Nous utilisons pour la correction manuelle l'Arborator développé par Gerdes (2013) (distribué librement et utilisable en ligne à partir de arborator.ilpga.fr). Plusieurs outils permettant d'entraîner un analyseur en dépendance sont actuellement distribués librement. Nous avons utilisé MATE (Bohnet 2010), ainsi que l'analyseur développé au LIF (Nasr et al. 2011). Le corpus d'amorçage a été développé à partir du treebank Rhapsodie, un corpus de 33 000 mots de français parlé annoté en prosodie et syntaxe distribué librement (Lacheret et al. 2014, projet-rhapsodie.fr) dont l'annotation syntaxique a été entièrement corrigée à la main, à partir d'un pré-annotation automatique réalisée avec un analyseur de l'écrit (de la Clergerie 2005), aucun analyseur pour le français parlé n'étant disponible à l'époque.

Dans ce résumé, nous ne présenterons pas davantage la chaîne de traitement, ni le découpage en unités d'analyse (énoncé, unité illocutoire, phrase, etc.), qui s'avère néanmoins une étape essentielle de l'analyse (voir par ex. Deulofeu et al. 2011, Pietrandrea et al. 2014). Nous allons nous concentrer sur les choix faits pour l'analyse syntaxique en dépendance.

Une remarque préalable essentielle : les choix d'annotation sont toujours un compromis entre diverses exigences (Gerdes & Kahane 2016). Des exigences théoriques : l'annotation doit répondre à un certain nombre de critères fixés par le cadre théorique. Des exigences pratiques liées au processus d'annotation : l'annotation doit être reproductible (accord inter-annotateur), elle doit être la plus simple possible (efficacité, rapidité), et surtout, lorsqu'elle est réalisée en grande partie automatiquement, elle doit pouvoir être propagée sur l'ensemble du corpus en minimisant les erreurs. Enfin des exigences liées à l'utilisateur final : les annotations doivent être facilement requêtées et permettre à l'utilisateur de récupérer les données qu'il souhaite étudier. Nous utilisons l'outil ANNIS (Zeldes et al. 2009) dont le langage de requête permet de décrire des configurations et d'extraire tous les énoncés dont l'arbre syntaxique contient cette configuration.

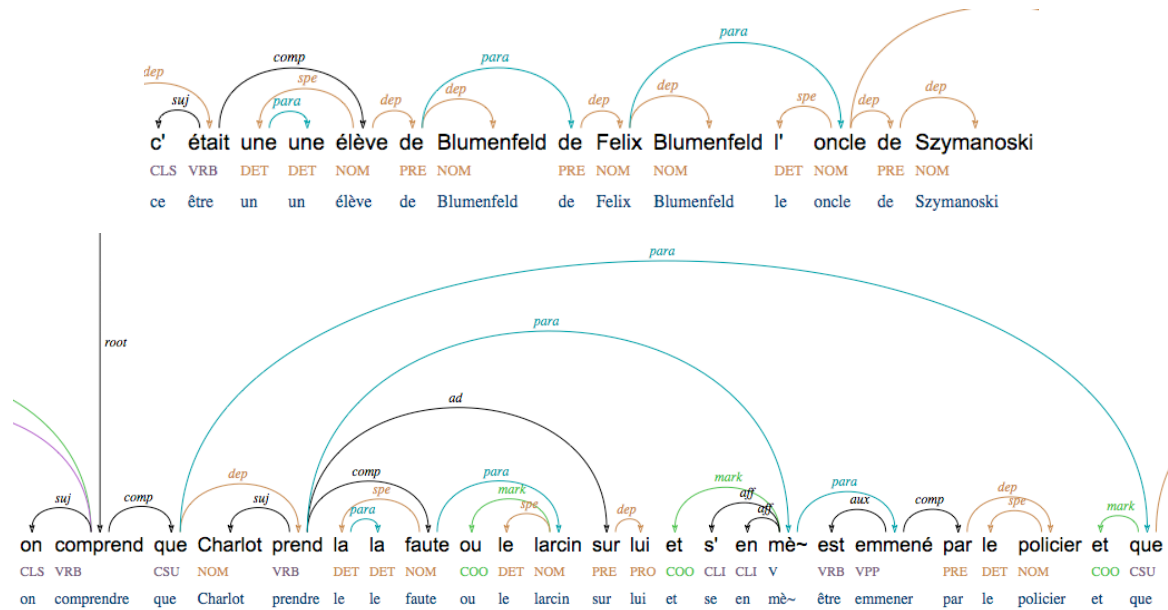


Figure 2. Liens paradigmatiques

3. Conclusion

Les premiers résultats d'entraînement d'un analyseur sur ce schéma d'annotation syntaxique sont extrêmement encourageants. Après entraînement sur seulement 60 000 mots, notre analyseur MATE affiche un f-score de 86,5% pour la reconnaissance du gouverneur d'un mot et 82,6% pour la reconnaissance du gouverneur et de la fonction, ce qui en fait d'ores et déjà un outil opérationnel pour le traitement automatique du français parlé.¹ Les f-scores pour les relations présentées ici sont : 74,0% pour *periph*, 89,9% pour *dm*, 67,3% pour *para*, 95,1% pour *mark*,

Mots-clés : syntaxe de dépendance, marqueurs de discours, listes paradigmatiques, treebank.

Références

- Blanche-Benveniste, C. (1990). Un modèle d'analyse syntaxique « en grilles » pour les productions orales. *Anuario de psicología/The UB Journal of psychology*, (47), 11-28.
- Bohnet, B. (2010, August). Very high accuracy and fast dependency parsing is not a contradiction. *ACL*, Uppsala
- Deulofeu J., Gerdes K., Kahane S., Pietrandrea P. (2010) Depends on what the French say: Spoken corpus annotation with and beyond syntactic function, *LAW IV*, ACL, Uppsala, 274-281.
- K. Gerdes, S. Kahane (2015) Non-constituent coordination and other coordinative constructions as dependency graphs, *Depling*, Uppsala.
- Gerdes K., Kahane S. (2016), Dependency Annotation Choices: Assessing Theoretical and Practical Issues of Universal Dependencies, *Proceedings of Linguistic Annotation Workshop (LAW)*, ACL, Berlin.
- Kahane, S., Pietrandrea, P. (2012). La typologie des entassements en français. *CMLF*, Lyon, 1809-1828.
- Nasr, A., Béchet, F., Rey, J. F., Favre, B., & Le Roux, J. (2011). Macaon: An NLP tool suite for processing word lattices. *ACL-HTL: Systems Demonstrations*, 86-91.

¹ Rappelons que pour l'évaluation, le parseur est entraîné sur 90% du corpus et évalué sur les 10% restant choisis aléatoirement. Le f-score est la moyenne harmonique entre la précision et le rappel.