



---

*Traitement Automatique des Langues Naturelles*  
(TALN) <sup>1</sup>

Actes de la 29e Conférence sur le Traitement Automatique des Langues Naturelles.  
Volume 1 : conférence principale

---

Yannick Estève, Tania Jiménez, Titouan Parcollet, Marceley Zanon Boito (Éds.)

Avignon, France, 27 juin au 1<sup>er</sup> juillet 2022

---

1. <https://taln2022.univ-avignon.fr>



Avec le soutien de



## Préface

### Bienvenue à TALN-RECITAL 2022

La 29ème édition de la conférence TALN et 24ème édition du RECITAL se déroulent cette année à Avignon, organisée conjointement par le Laboratoire d'Informatique et Systèmes (LIS) et le Laboratoire Informatique d'Avignon (LIA) sous l'égide de l'Association pour le Traitement Automatique des Langues (ATALA).

La thématique spéciale de la conférence est : « Vers un TAL inclusif. »

Nous accueillons les participants sur le campus Hannah Arendt d'Avignon Université du lundi 27 juin au vendredi 1er juillet 2022.

Après deux années minées par les contraintes sanitaires, nous sommes heureux de pouvoir organiser cet événement en mode présentiel, avec plus de 150 participants.

Le premier jour, lundi, est consacré à deux ateliers thématiques : « DEFT Fouille de texte » dédié cette année à la correction automatique de copies d'étudiants ; et la première édition de l'atelier « TAL et Humanités Numériques » qui applique les thèmes du TAL aux problématiques des corpus SHS.

Ce même jour se tiendra le tutoriel « Quelques étapes souvent omises dans la préparation de corpus ». Nous remercions les organisateurs de ces ateliers et du tutoriel pour leurs propositions, leur animation et l'organisation de ces événements.

Les quatre jours qui suivent associent dans des séances communes la conférence TALN et la conférence jeunes chercheurs RECITAL.

Le dernier jour, vendredi 1er juillet, nous avons le plaisir d'accueillir une table ronde autour de laquelle industriels et académiques échangeront autour des thèmes de la conférence. À la suite de cette table ronde, les démonstrations et posters permettent de continuer ses échanges.

Nous accueillons deux conférencières invitées : Seza Doğruöz, Professeure de la Faculté de philosophie et lettres à l'Université de Gand et Teresa Lynn, Research Fellow à The ADAPT Centre, Dublin City University. Nous sommes très fiers qu'elles aient accepté notre invitation.

TALN2022 et RECITAL2022 auront permis aux chercheurs de présenter leurs travaux sous forme de communications orales, de posters et de démonstrations. 62 papiers au total ont été soumis à TALN, 47 papiers (76%) ont été acceptés : 35 en présentation orale et 12 en poster. En ce qui concerne la conférence RECITAL, 15 articles ont été soumis, et 9 articles ont été acceptés.

Nous tenons à remercier chaleureusement les membres des comités de programme de TALN et de RECITAL, qui ont fait un travail formidable, ainsi que le comité permanent de la conférence (CPERM) et son président pour l'aide dans l'organisation de l'événement et vos conseils éclairés.

Un grand merci à nos sponsors : Avignon Université, Aix-Marseille Université, NAVER LABS, MOBIDYS et ORKIS, sans lesquels cette édition de TALN - RECITAL n'aurait pas pu avoir lieu.

Yannick Estève, Tania Jiménez, Titouan Parcollet, Marcelly Zanon Boito

## Comités

TALN 2022 est organisé conjointement par le Laboratoire d’Informatique et Systèmes (LIS) et le Laboratoire Informatique d’Avignon (LIA) sous l’égide de l’Association pour le Traitement Automatique des Langues (ATALA).

### Comité local d’organisation

- Yannick Estève, LIA
- Tania Jiménez, LIA
- Teva Merlin, LIA
- Antoine Caubrière, LIA
- Arthur Amalvy, LIA
- Bassam Jabaian, LIA
- Carlos Ramisch, LIS
- Corinne Fredouille, LIA
- Gaëlle Laperrière, LIA
- Jarod Duret, LIA
- Juan-Manuel Torres-Moreno, LIA
- Marcelly Zanon Boito, LIA
- Natalia Tomashenko, LIA
- Paul-Gauthier Noé, LIA
- Pierre Jourlin, LIA
- Salima Mdhaffar, LIA
- Stéphane Huet, LIA
- Sylvie Ros, LIS
- Titouan Parcollet, LIA
- Jean-François Bonastre, LIA

### Comité de programme

- Benoit Favre, Aix-Marseille Université
- Caio Corro, Université Paris-Saclay
- Caroline Brun, Naver Labs Europe
- Delphine Bernhard, LiLPa, Université de Strasbourg
- Géraldine Damnati, Orange Labs
- Guillaume Wisniewski, LLF, Université de Paris
- Magalie Ochs, LSIS
- Marie Candito, Université Paris 7 / INRIA
- Maud Ehrmann, EPFL, DHLAB
- Natalia Grabar, STL CNRS Université Lille 3
- Philippe Langlais, Université de Montréal
- Philippe Muller, IRIT, Université Toulouse
- Thomas François, Université Catholique de Louvain
- Yannick Estève, LIA
- Yannick Parmentier, LORIA — Université de Lorraine

### Comité de Lecture TALN

- Andon Tchechmedjiev, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales
- Aurélie Névél, Université Paris-Saclay, CNRS, LISN

- Béatrice Daille, Laboratoire d’Informatique Nantes Atlantique (LINA)
- Cédric Lopez, Emvista
- Céline Alec, Université de Caen-Normandie
- Charles Teissède, Synapse Développement
- Christine Jacquin, LS2N
- Cyril Grouin, LIMSI, CNRS, Université Paris-Saclay
- Damien Nouvel, INaLCO
- Denis Maurel, University Francois Rabelais Tours
- Didier Schwab, Université Grenoble Alpes
- Dominique Estival, Western Sydney University
- Emmanuel Morin, Université de Nantes, LS2N
- Eric Laporte, Université Gustave Eiffel
- François Yvon, LIMSI/CNRS et Université Paris-Sud
- Gaël Dias, Normandie University
- Gaël Guibon, Télécom Paris - SNCF
- Gwénolé Lecorvé, Orange
- Jean-Philippe Prost, Aix-Marseille Université
- Jian-Yun Nie, Université de Montréal
- Juan-Manuel Torres-Moreno, Laboratoire Informatique d’Avignon — Avignon Université
- Karèn Fort, Sorbonne Université
- Laurent Besacier, Laboratoire d’Informatique de Grenoble
- Ludovic Tanguy, CLLE-ERSS
- Luka Nerima, Université de Genève
- Mathieu Dehouck, INRIA
- Matthieu Constant, Université de Lorraine, ATILF, CNRS
- Matthieu Labeau, Telecom Paris
- Maxime Amblard, Université de Lorraine
- Maximin Coavoux, CNRS, Université Grenoble Alpes
- Mikaela Keller, Université Lille 3 - Inria
- Nabil Hathout, CNRS
- Nathalie Camelin, LIUM — Université du Maine
- Nicolas Dugué, LIUM — Université du Maine
- Nicolas Hernandez, Université de Nantes — LINA CNRS UMR 6241
- Olivier Ferret, CEA List
- Olivier Hamon, Syllabs
- Pascal Amsili, Sorbonne Nouvelle
- Pascale Sébillot, IRISA
- Patrice Bellot, Aix-Marseille Université - CNRS (LSIS)
- Peggy Cellier, IRISA/INSA Rennes
- Philippe Blache, CNRS & Université de Provence
- Remi Cardon, CENTAL
- Richard Moot, CNRS (LIRMM) & Université de Montpellier
- Salima Mdhaffar, LIUM
- Solen Quiniou, LINA — Université de Nantes
- Stéphane Huet, Laboratoire Informatique d’Avignon — Avignon Université
- Sylvain Kahane, Modyco, Université Paris Ouest Nanterre & CNRS
- Sylvain Pogodalla, LORIA/INRIA Lorraine
- Thierry Charnois, LIPN CNRS University of PARIS 13
- Thierry Hamon, Université Paris-Saclay, CNRS, LIMSI & Université Sorbonne Paris Nord
- Véronique Moriceau, IRIT-CNRS
- Vincent Claveau, IRISA - CNRS

- Xavier Tannier, Sorbonne Université, INSERM, LIMICS
- Yves Bestgen, F.R.S-FNRS et UCL
- Yves Lepage, Waseda University

## Présentations invitées

- **A. Seza Doğruöz**

Professeure — Faculté de philosophie et lettres, Université de Gand

**Titre / Title :** "Multilingualism 101" for Computational Linguists : Challenges & Opportunities for Multilingual & Inclusive Language Technologies

**Résumé / Abstract :** A multilingual is someone who speaks more than one language for communication in his/her daily life. Multilingualism is commonly practiced at the individual or societal levels in most parts of the world (e.g., Europe, Asia, Africa). However, current language technologies are usually built around monolingual assumptions which do not always reflect the linguistic variation in real world communication. The first goal of this talk is to familiarize the computational linguistics community with the key concepts around multilingualism. Secondly, I will discuss the challenges and opportunities for building linguistically diverse, inclusive and multilingual language technologies through giving examples from high and low resource languages and their speakers/users.

**Short bio :** A. Seza Doğruöz is a faculty at Ghent University. Her research focuses on analyzing multilingual language use and linguistic variation between humans and developing language technologies accordingly. Her publications cover the themes related to Computational Sociolinguistics, Multilingualism & Linguistic Variation, Curating Language Resources (especially for low resource languages) and Open Domain Dialog Systems. Examples of her recent publications are :

- Doğruöz, A.S., Sitaram, S. (2022). Language Technologies for Low Resource Languages : Sociolinguistic and Multilingual Insights. LREC'22/SIGUL.
- Doğruöz, A.S., Sitaram, S., Bullock, B.E., Toribio, A.J. (2021). "A Survey of Code-switching : Linguistic and Social Perspectives for Language Technologies", Proceedings of the Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021), Bangkok, Thailand. (Nominated for the Best Paper Award)
- Doğruöz, A.S., Skantze, G. (2021). "How open are the conversations with open-domain chatbots? A proposal for speech-event based evaluation". Proceedings of the 22nd Annual Meeting of Special Interest Group on Discourse and Dialogue (SIGDIAL2021), Singapore.

- **Teresa Lynn**

Research Fellow — The ADAPT Centre, Dublin City University

**Titre / Title :** The challenges of working with low-resourced languages in NLP : An Irish story.

**Résumé / Abstract :** Language technology underpins many of the applications and platforms that enable our digitally enhanced lives (virtual assistants, search engines, translation tools, spell-checkers, language learning tools, etc.). Yet these advances do not benefit all languages equally. Due to a lack of sufficient language technologies, users often need to revert to using another (better supported) language. Such a language shift plays a major role in the risk of digital extinction, i.e. an eventual decline in language use due to lack of technological support.

This talk focuses on Irish — an official EU language, and considered a low-resourced language in terms of digital support. "Low-resourced" not only means that there is a severe lack of speech and language applications available for Irish speakers to use, but it also means that the fundamental tools and language

resources required to build these technologies are also lacking. Irish is also in a precarious position while it competes alongside the most technologically supported language in the world — English. I will talk about our work at Dublin City University where we are taking some steps towards addressing this risk of digital extinction through the development of NLP tools and resources for Irish.

**Short bio :** Dr. Teresa Lynn is a Research Fellow at the ADAPT Centre in Dublin City University. She was awarded her PhD in 2015 from both DCU and Macquarie University Sydney, and was a recipient of the 2014-2015 Enterprise Ireland Fulbright Award, carrying out research at Saint Louis University, Missouri, USA.

Her main research interests lie in developing tools and resources for Irish language technology. She is the principal investigator on the GaelTech project, funded by the Irish Government Department of the Gaeltacht, which covers various research topics in Irish NLP. She is also a core member of the European Language Equality project and Ireland's National Anchor Point for the ELRC (European Language Resource Coordination), overseeing national data collection for Irish machine translation. Other projects include ELRI (European Language Resource Infrastructure), the National Digital Plan for Irish and the Universal Dependencies Project. Her research covers treebank development, syntactic parsing, social media NLP and multiword expressions. Teresa also worked in industry for several years, namely in the areas of localisation, NLP and machine translation.

## Table des matières

<b>I</b>	<b>Travaux originaux</b>	<b>1</b>
	<b>Abstraction ou hallucination ? État des lieux et évaluation du risque pour les modèles de génération de résumés automatiques de type séquence-à-séquence</b>	<b>2</b>
	<i>Eunice Akani, Benoit Favre, Frederic Bechet</i>	
	<b>Choisir le bon co-équipier pour la génération coopérative de texte</b>	<b>12</b>
	<i>Antoine Chaffin, Vincent Claveau, Ewa Kijak, Sylvain Lamprier, Benjamin Piwowarski, Thomas Scialom, Jacopo Staiano</i>	
	<b>Décodage guidé par un discriminateur avec le Monte Carlo Tree Search pour la génération de texte contrainte</b>	<b>27</b>
	<i>Antoine Chaffin, Vincent Claveau, Ewa Kijak</i>	
	<b>Détection d'anomalies textuelles à base de l'ingénierie d'invite</b>	<b>42</b>
	<i>Yizhou Xu, Kata Gábor, Leila Khouas, Frédérique Segond</i>	
	<b>Détection des influenceurs dans des médias sociaux par une approche hybride</b>	<b>54</b>
	<i>Kevin Deturck, Damien Nouvel, Namrata Patel, Frederique Segond</i>	
	<b>Étiquetage ou génération de séquences pour la compréhension automatique du langage en contexte d'interaction ?</b>	<b>64</b>
	<i>Rim Abrougui, Géraldine Damnati, Johannes Heinecke, Frédéric Béchet</i>	
	<b>Etude des stéréotypes genrés dans le théâtre français du XVIIe au XIXe siècle à travers des plongements lexicaux</b>	<b>74</b>
	<i>Alexandra Benamar, Cyril Grouin, Meryl Bothua, Anne Vilnat</i>	
	<b>FENEC : un corpus équilibré pour l'évaluation des entités nommées en français</b>	<b>82</b>
	<i>Alice Millour, Yoann Dupont, Alexane Jouglar, Karën Fort</i>	
	<b>Filtrage et régularisation pour améliorer la plausibilité des poids d'attention dans la tâche d'inférence en langue naturelle</b>	<b>95</b>
	<i>Duc Hau Nguyen, Guillaume Gravier, Pascale Sébillot</i>	
	<b>Génération de question à partir d'analyse sémantique pour l'adaptation non supervisée de modèles de compréhension de documents</b>	<b>104</b>
	<i>Elie Antoine, Jeremy Auguste, Frederic Bechet, Géraldine Damnati</i>	
	<b>Identification of complex words and passages in medical documents in French</b>	<b>116</b>
	<i>Kim Cheng Sheang, Anaïs Koptient, Natalia Grabar, Horacio Saggion</i>	
	<b>Impact du français inclusif sur les outils du TAL</b>	<b>126</b>
	<i>Cyril Grouin</i>	
	<b>Investigating associative, switchable and negatable Winograd items on renewed French data sets</b>	<b>136</b>
	<i>Xiaoou Wang, Olga Seminck, Pascal Amsili</i>	
	<b>Langues par défaut ? Analyse contrastive et diachronique des langues non citées dans les articles de TALN et d'ACL</b>	<b>144</b>

*Fanny Ducel, Karën Fort, Gaël Lejeune, Yves Lepage*

<b>Le projet FREEM : ressources, outils et enjeux pour l'étude du français d'Ancien Régime</b>	<b>154</b>
<i>Simon Gabay, Pedro Ortiz Suarez, Rachel Bawden, Alexandre Bartz, Philippe Gambette, Benoît Sagot</i>	
<b>Mesures linguistiques automatiques pour l'évaluation des systèmes de Reconnaissance Automatique de la Parole</b>	<b>166</b>
<i>Thibault Bañeras Roux, Mickaël Rouvier, Jane Wottawa, Richard Dufour</i>	
<b>Modèle-s bayés-ien-s pour la segment-ation à deux niveau-x faible-ment super-vis-é-e</b>	<b>174</b>
<i>Shu Okabe, François Yvon</i>	
<b>Ré-ordonnancement via programmation dynamique pour l'adaptation cross-lingue d'un analyseur en dépendances</b>	<b>183</b>
<i>Nicolas Devatine, Caio Corro, François Yvon</i>	
<b>Remplacement de mentions pour l'adaptation d'un corpus de reconnaissance d'entités nommées à un domaine cible</b>	<b>198</b>
<i>Arthur Amalvy, Vincent Labatut, Richard Dufour</i>	
<b>RésuméSVD : Un outil efficace et performant pour le résumé de texte non supervisé</b>	<b>206</b>
<i>Gabriel Shenouda, Christophe Rodrigues, Aurélien Bossard</i>	
<b>Stratégies d'adaptation pour la reconnaissance d'entités médicales en français</b>	<b>215</b>
<i>Tiphaine Le Clercq de Lannoy, Romaric Besançon, Olivier Ferret, Julien Tourille, Frédérique Brin-Henry, Bianca Vieru</i>	
<b>Un algorithme d'analyse sémantique fondée sur les graphes via le problème de l'arborescence généralisée couvrante</b>	<b>226</b>
<i>Alban Petit, Caio Corro</i>	
<b>Une chaîne de traitement pour prédire et appréhender la complexité des textes pour enfants d'un point de vue linguistique</b>	<b>236</b>
<i>Delphine Battistelli, Aline Etienne, Rashedur Rahman, Charles Teissèdre, Gwénolé Lecorvé</i>	
<b>Une étude statistique des plongements dans les modèles transformers pour le français</b>	<b>247</b>
<i>Loïc Fosse, Duc-Hau Nguyen, Pascale Sébillot, Guillaume Gravier</i>	
<b>Vers la compréhension automatique de la parole bout-en-bout à moindre effort</b>	<b>257</b>
<i>Marco Naguib, François Portet, Marco Dinarelli</i>	
<b>II Traduction de soumissions en cours à des conférences internationales</b>	<b>269</b>
<b>Adaptation au domaine de modèles de langue à l'aide de réseaux à base de graphes</b>	<b>270</b>
<i>Merieme Bouhandi, Emmanuel Morin, Thierry Hamon</i>	
<b>Annotation d'expressions polylexicales verbales en arabe : validation d'une procédure d'annotation multilingue</b>	<b>280</b>
<i>Najet Hadj Mohamed, Cherifa Ben Khelil, Agata Savary, Iskander Keskes, Jean Yves Antoine, Lamia Hadrich Belguith</i>	

<b>CLISTER : Un corpus pour la similarité sémantique textuelle dans des cas cliniques en français</b>	<b>287</b>
<i>Nicolas Hiebel, Karën Fort, Aurélie Névéol, Olivier Ferret</i>	
<b>COMFO : Corpus Multilingue pour la Fouille d’Opinions</b>	<b>297</b>
<i>Lamine Faty, Khadim Drame, Edouard Ngor Sarr, Marie Ndiaye, Yoro Dia, Ousmane Sall</i>	
<b>Classification automatique de questions spontanées vs. préparées dans des transcriptions de l’oral</b>	<b>305</b>
<i>Iris Eshkol-Taravella, Angèle Barbedette, Xingyu Liu, Valentin-Gabriel Soumah</i>	
<b>Décontextualiser des plongements contextuels pour construire des thésaurus distributionnels</b>	<b>315</b>
<i>Olivier Ferret</i>	
<b>Évaluation comparative de systèmes neuronal et statistique pour la résolution de coréférence en langage parlé</b>	<b>325</b>
<i>Maëlle Brassier, Théo Azzouza, Jean-Yves Antoine, Loïc Grobol, Anaïs Lefevre-Halftermeyer</i>	
<b>Extraction d’informations de messages aéronautiques (NOTAMs) avec des modèles de langue appris de façon auto-supervisée</b>	<b>335</b>
<i>Alexandre Arnold, Fares Ernez, Catherine Kobus, Marion-Cécile Martin</i>	
<b>Fine-tuning de modèles de langues pour la veille épidémiologique multilingue avec peu de ressources</b>	<b>345</b>
<i>Stephen Mutuvi, Emanuela Boros, Antoine Doucet, Adam Jatout, Gaël Lejeune, Moses Odeo</i>	
<b>French CrowS-Pairs : Extension à une langue autre que l’anglais d’un corpus de mesure des biais sociétaux dans les modèles de langue masqués</b>	<b>355</b>
<i>Aurélie Névéol, Yoann Dupont, Julien Bezançon, Karën Fort</i>	
<b>Identification des Expressions Polylexicales dans les Tweets</b>	<b>365</b>
<i>Nicolas Zampieri, Carlos Ramisch, Irina Illina, Dominique Fohr</i>	
<b>L’importance des entités pour la tâche de détection d’événements en tant que système de question-réponse</b>	<b>374</b>
<i>Emanuela Boros, Jose Moreno, Antoine Doucet</i>	
<b>Les représentations distribuées sont-elles vraiment distribuées ? Observations sur la localisation de l’information syntaxique dans les tâches d’accord du verbe en français</b>	<b>384</b>
<i>Bingzhi Li, Guillaume Wisniewski, Benoît Crabbé</i>	
<b>Mieux utiliser BERT pour la détection d’évènements à partir de peu d’exemples</b>	<b>392</b>
<i>Aboubacar Tuo, Romaric Besançon, Olivier Ferret, Julien Tourille</i>	
<b>Preuve de concept d’un bot vocal dialoguant en wolof</b>	<b>403</b>
<i>Elodie Gauthier, Papa Séga Wade, Thierry Moudenc, Patrice Collen, Emilie De Neef, Oumar Ba, Ndeye Khoyane Cama, Cheikh Ahmadou Bamba Kebe, Ndeye Aissatou Gningue, Thomas Mendo’O Aristide</i>	
<b>Tâches Auxiliaires Multilingues pour le Transfert de Modèles de Détection de Discours Haineux</b>	<b>413</b>
<i>Arij Riabi, Syrielle Montariol, Djamé Seddah</i>	

<b>Tâches auxiliaires pour l’analyse biaffine en graphes de dépendances</b>	<b>424</b>
<i>Marie Candito</i>	
<b>Un jeu de données pour répondre à des questions visuelles à propos d’entités nommées en utilisant des bases de connaissances</b>	<b>434</b>
<i>Paul Lerner, Olivier Ferret, Camille Guinaudeau, Hervé Le Borgne, Romaric Besançon, Jose Moreno, Jesús Lovón-Melgarejo</i>	
<b>III Résumés d’articles internationaux</b>	<b>445</b>
<b>Encouraging Neural Machine Translation to Satisfy Terminology Constraints.</b>	<b>446</b>
<i>Melissa Ailem, Jingshu Liu, Raheel Qader</i>	
<b>L’Attention est-elle de l’Explication ? Une Introduction au Débat</b>	<b>447</b>
<i>Adrien Bibal, Remi Cardon, David Alfter, Rodrigo Wilkens, Xiaoou Wang, Thomas François, Patrick Watrin</i>	
<b>Quand être absent de mBERT n’est que le commencement : Gérer de nouvelles langues à l’aide de modèles de langues multilingues</b>	<b>450</b>
<i>Benjamin Muller, Antonios Anastasopoulos, Benoît Sagot, Djamé Seddah</i>	
<b>IV Prises de position</b>	<b>452</b>
<b>Evaluation of Automatic Text Simplification : Where are we now, where should we go from here</b>	<b>453</b>
<i>Natalia Grabar, Horacio Saggion</i>	

Première partie  
Travaux originaux