



HAL
open science

**Actes de CORIA-TALN 2023. Actes de la 30e
Conférence sur le Traitement Automatique des Langues
Naturelles (TALN)**

Christophe Servan, Anne Vilnat

► **To cite this version:**

Christophe Servan, Anne Vilnat. Actes de CORIA-TALN 2023. Actes de la 30e Conférence sur le Traitement Automatique des Langues Naturelles (TALN) : volume 6 : projets. CORIA - TALN 2023, 2023. hal-04463005

HAL Id: hal-04463005

<https://hal.science/hal-04463005>

Submitted on 16 Feb 2024

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.



Distributed under a Creative Commons Attribution 4.0 International License



*18e Conférence en Recherche d'Information et Applications,
16e Rencontres Jeunes Chercheurs en RI,
30e Conférence sur le Traitement Automatique des Langues Naturelles,
25e Rencontre des Étudiants Chercheurs en Informatique pour le
Traitement Automatique des Langues
(CORIA-TALN) ¹*

Actes de CORIA-TALN 2023.

Actes de la 30e Conférence sur le Traitement Automatique des Langues Naturelles
(TALN),
volume 6 : projets

Christophe Servan, Anne Vilnat (Éds.)

Paris, France, 5 au 9 juin 2023

1. <https://coria-taln-2023.sciencesconf.org/>

Avec le soutien de



Préface

Organisée conjointement par les laboratoires franciliens sous l'égide de l'Association francophone de Recherche d'Information et Applications (ARIA) et l'Association pour le Traitement Automatique des Langues (ATALA), la conférence CORIA-TALN-RJCRI-RECITAL 2023 regroupe :

- la 18ème Conférence en Recherche d'Information et Applications (CORIA)
 - la 30ème Conférence sur le Traitement Automatique des Langues Naturelles (TALN) ;
- ainsi que les deux conférences associées, destinées aux jeunes chercheuses et chercheurs :
- Les 16ème Rencontres Jeunes Chercheurs en RI (RJCRI)
 - la 25ème Rencontre des Etudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RECITAL)

La conférence TALN (Traitement Automatique des Langues Naturelles) est un rendez-vous annuel qui offre, depuis 1994, le plus important forum d'échange international francophone aux acteurs universitaires et industriels des technologies de la langue. Cet événement, qui accueille habituellement près de 250 participants, couvre toutes les avancées récentes en matière de communication écrite et parlée et de traitement informatique de la langue notamment la recherche et l'extraction d'information, la fouille de textes, le dialogue homme-machine, la fouille d'opinions, la traduction automatique, les systèmes de questions-réponses, le résumé automatique...

Cette année, ont été soumis 51 articles longs et 12 articles courts pour la conférence principale, dont respectivement 29 ont été acceptés pour une présentation orale (dont 2 prises de position) et 9 pour une présentation sous forme de posters. 19 présentations courtes, sous forme de posters, d'articles déjà publiés lors de conférences internationales complètent le programme de la conférence, ainsi que des démonstrations et des présentations de projets en cours. L'alternance de sessions communes entre TALN, CORIA et RJC et de sessions plus spécifiques devraient permettre de susciter des échanges fructueux.

En complément de la conférence principale, se tiennent les ateliers "Défi Fouille de Texte" (DEFT), "Atelier sur l'analyse et la recherche de textes scientifiques" (ARTS), "Humain ou pas humain ? : les nouveaux défis pour les humains" (hOUPSh) et le tutoriel "Apprentissage Profond pour le TAL français pour les débutants" (TutoriAL). Ces ateliers et tutoriel illustrent à la fois des tendances nouvelles présentes dans la communauté et des activités récurrentes.

Un grand merci à toutes celles et tous ceux qui ont soumis leurs travaux, ainsi qu'aux membres du comité de programme et aux relectrices et relecteurs pour le travail qu'ils ont accompli. Ce sont eux qui font vivre la conférence. Merci au comité d'organisation réparti sur la région parisienne, et aux sponsors qui nous ont permis d'organiser cet événement.

Christophe Servan et Anne Vilnat, co-présidents de TALN

Comités

Comité de programme

Présidents

- Christophe SERVAN
- Anne VILNAT

Membres

- Rachel BAWDEN
- Caroline BRUN
- Marie CANDITO
- Rémi CARDON
- Pascal DENIS
- Yannick ESTEVE
- Benoît FAVRE
- Amel FRAISSE
- Thomas GERALD
- Natalia GRABAR
- Lydia-Mai HO-DAC
- José MORENO
- Vassilina NIKOULINA
- Yannick PARMENTIER
- Sylvain POGODALLA
- Solène QUINIOU
- Didier SCHWAB
- Iris TARAVELLA-ESHKOL

Comité d'organisation

- Marie CANDITO
- Thomas GERALD
- José MORENO
- Benjamin PIWOWARSKI
- Christophe SERVAN
- Laure SOULIER
- Anne VILNAT

Table des matières

Les jeux de données en compréhension du langage naturel et parlé : paradigmes d’annotation et représentations sémantiques	1
<i>Rim Abrougui</i>	
Projet Gender Equality Monitor (GEM)	21
<i>Gilles Adda, François Buet, Sahar Ghannay, Cyril Grouin, Camille Guinaudeau, Lufei Liu, Aurélie Névéol, Albert Rilliard, Uro Rémi</i>	
CEN-CENELEC JTC 21 : La standardisation en TALN au service du règlement européen sur l’IA	22
<i>Lauriane Aufrant</i>	
TEMITALC : Text Mining et TAL pour Analyser le Langage des Cachalots	23
<i>Jose Coch, Olivier Adam</i>	
muDialBot, vers l’interaction humain-robot multimodale pro-active	26
<i>Fabrice Lefèvre, Timothée Dhaussy, Bassam Jabaian, Ahmed Njifenjou, Virgile Sucal</i>	
Projet ANR MALIN : MANuels scoLaires INclusifs	30
<i>Olivier Pons, Isabelle Barbet, Jérôme Dupire, Valérie Grembi, Camille Guinaudeau, Céline Hudelot, Caroline Huron, Elise Lincker, Vincent Mousseau, Léa Pacini</i>	
PROPICTO : Développer des systèmes de traduction de la parole vers des séquences de pictogrammes pour améliorer l’accessibilité de la communication	32
<i>Didier Schwab</i>	
Recherche d’information conversationnelle	36
<i>Laure Soulier, Pierre Erbacher, Thomas Gerald, Hanane Djeddal, Jian-Yun Nie, Philippe Preux</i>	
Autogramm : développement simultané de treebanks et de grammaires à partir de corpus	37
<i>Kahane Sylvain, Santiago Herrera, Bruno Guillaume, Kim Gerdes</i>	

Les jeux de données en compréhension du langage naturel et parlé : paradigmes d’annotation et représentations sémantiques

Rim Abrougui^{1, 2}

(1) Orange Innovation, Lannion, France

(2) Aix-Marseille Université, LIS UMR 7020, Marseille, France

rim.abrougui@orange.com

RÉSUMÉ

La compréhension du langage naturel et parlé (NLU/SLU) couvre le problème d’extraire et d’annoter la structure sémantique, à partir des énoncés des utilisateurs dans le contexte des interactions humain/machine, telles que les systèmes de dialogue. Elle se compose souvent de deux tâches principales : la détection des intentions et la classification des concepts. Dans cet article, différents corpora SLU sont étudiés au niveau formel et sémantique : leurs différents formats d’annotations (à plat et structuré) et leurs ontologies ont été comparés et discutés. Avec leur pouvoir expressif gardant la hiérarchie sémantique entre les intentions et les concepts, les représentations sémantiques structurées sous forme de graphe ont été mises en exergue. En se positionnant vis à vis de la littérature et pour les futures études, une projection sémantique et une modification au niveau de l’ontologie du corpus MultiWOZ ont été proposées.

ABSTRACT

The Challenges of Spoken Language Understanding Datasets : A Study on Annotations and Semantic Representations

Natural and Spoken Language Understanding (NLU/SLU) covers the problem of extracting and annotating the meaning structure from user utterances in the context of human/machine interaction, such as dialogue systems, consisting oftenly of two main tasks : intent detection and slot filling. In this paper, different SLU corpora were studied at a formal and semantic level : their different annotation formats (flat and structured) and ontologies were compared and discussed. With their expressive power maintaining the semantic hierarchy between intents and slots, graph semantic representations were highlighted. In line with the literature and for future studies, a semantic projection and a modification of the ontology of the MultiWOZ corpus were proposed.

MOTS-CLÉS : Compréhension du langage, ontologies, représentation sémantique à plat (BIO), représentation sémantique structurée (graphe).

KEYWORDS: Language Understanding, ontologies, flat semantic representation (BIO), structured semantic representation (graph).

1 Introduction

La compréhension du langage naturel et parlé est un sujet d’étude important dans le cadre des interactions homme-machine. Le domaine comprend plusieurs niveaux d’étude, mais actuellement la tâche de compréhension est principalement axée sur la compréhension de la sémantique globale

des requêtes des utilisateurs et sur l'identification des concepts génériques des mots-clés, à savoir, la détection des intentions et l'identification des concepts (Tur & De Mori, 2011).

Bien qu'il existe plusieurs approches de représentations sémantiques, la majorité des méthodes se base sur la représentation à base de frames sémantiques en utilisant des modèles supervisés pour la classification et l'étiquetage de séquence. Ces modèles qui sont basés sur des réseaux de neurones utilisant des modèles de langage pré-entraînés, ont obtenu des performances élevées sur plusieurs jeux de données SLU qui sont représentés avec un schéma à plat (Béchet & Raymond, 2019) (cf. figure 3).

Cependant, les interactions dans une conversation en conditions réelles sont beaucoup plus complexes. Afin de relever ces défis, il est nécessaire d'avoir d'un côté des corpus d'apprentissage dotés de représentations sémantiques complexes et contextuelles, et de l'autre côté des schémas d'annotation capables de prendre en compte les représentations sémantiques hiérarchiques. De plus, pour construire des modèles de SLU plus robustes et pour les comparer de manière plus précise, il est primordial d'unifier les ensembles de données existants. Cette unification permettra de diversifier les domaines et de fournir plus de connaissances aux systèmes de compréhension du langage qui pourront par conséquent apprendre plusieurs structures sémantiques.

Il existe plusieurs études complètes sur la compréhension du langage naturel et parlé, telles que (Weld *et al.*, 2022) et (Qin *et al.*, 2021), mais cet article explorera plus la question des jeux de données en étudiant leurs ontologies et formats d'annotation et de représentation sémantique. Dans le cadre de la problématique d'unification de toutes les ressources publiques d'apprentissage et d'évaluation des systèmes NLU/SLU, nous avons comparé les différentes ontologies et schémas d'annotation. Nous mettons ainsi en exergue le potentiel des représentations sémantiques structurées qui peuvent être utilisées plus facilement avec le développement des modèles de génération du langage. En s'intéressant particulièrement aux représentations en graphe qui préservent la hiérarchie et le lien sémantique entre les différents labels, nous proposerons une projection sémantique et une modification au niveau de l'ontologie du corpus MultiWOZ2.3.

La présentation des jeux de données et la comparaison de leurs ontologies sont présentées dans la section 2, alors que l'étude des schémas d'annotation et des représentations sémantiques structurées ainsi que nos perspectives sont exposées dans la section 3.

2 Exploration des jeux de données en compréhension du langage

Les jeux de données jouent un rôle crucial dans l'avancement de la recherche dans le domaine de la compréhension du langage. Ils permettent en effet d'entraîner, d'évaluer et de comparer les systèmes SLU. Les corpora disponibles sont variés et peuvent couvrir plusieurs domaines. La façon dont les annotations sont représentées et le choix des labels reflètent une certaine variation au niveau des ontologies ce qui rend difficile l'unification de ces jeux de données. Les schémas de projection et des représentations sémantiques choisis affectant la qualité d'annotation peuvent être aussi différents. Certains corpora sont des conversations complètes multi-domaines et/ou multi-intentions, tandis que d'autres ne contiennent que des simples requêtes. Malgré cette diversité, un corpus large avec des conversations complètes entièrement annotées en logique SLU selon une ontologie générique applicable à toutes les données est difficile à trouver, ce qui rend les travaux sur l'exploitation de l'histoire conversationnelle et du contexte plus difficile.

Nous pouvons trouver également plusieurs types de collectes de ces corpus. L’une de ces méthodes est appelée la méthode «Wizard-Of-OZ» où les participants interagissent en temps réel avec un système qu’ils croient être autonome, mais il est contrôlé en réalité par un opérateur humain invisible. Cette méthode est plus fréquente puisqu’elle permet de collecter des données de manière contrôlée, avec une grande variété de scénarios de dialogues. Les transcriptions de la parole à l’aide de la reconnaissance automatique de la parole (ASR), sont un autre moyen de collecter des données. La qualité des transcriptions va dépendre des systèmes ASR mais cette méthode permet de collecter d’une façon plus rapide les énoncés dans des conditions réelles et qui vont être vérifiés et annotés manuellement ou d’une manière semi-automatique. Nous trouvons enfin des méthodes plus automatisées qui consistent à synthétiser les conversations à l’aide de modèles de langage ou de modèles de génération de texte. Cette méthode est la plus rapide mais elle est artificielle car elle ne peut pas refléter toutes les variations linguistiques dans la parole humaine et peut avoir des problèmes d’hallucination. Cette section présentera certains jeux de données publiés pour la tâche SLU, ainsi que leurs ontologies et leurs méthodes d’annotation.

2.1 Les ensembles de données pour les tâches de compréhension du langage naturel et parlé

Jeux de données	Langues	Modes	#énoncés/dialogues	#labels
ATIS (Hemphill <i>et al.</i> , 1990)	en	requêtes	4978	17 intentions 84 slots
Frames (Asri <i>et al.</i> , 2017)	en	dialogues	1369	20 actes 16 slots
Massive (FitzGerald <i>et al.</i> , 2022)	multilingues 51 langues	requêtes	19521 par langue	60 intentions 55 slots
MEDIA (Devillers <i>et al.</i> , 2004)	fr	dialogues	1250	83 slots 19 spécifieurs
mTOD (Schuster <i>et al.</i> , 2019)	multilingues 3 langues	requêtes	43000	12 intentions 11 slots
mTOP (Li <i>et al.</i> , 2020)	multilingues 6 langues	requêtes	10000	117 intentions 78 slots
MultiDoGo (Peskov <i>et al.</i> , 2019)	en	dialogues	15000 annotés	85 intentions 73 slots
MultiWOZ (Budzianowski <i>et al.</i> , 2018)	en	dialogues	10438	32 actes 27 slots
M2M (Shah <i>et al.</i> , 2018)	en	dialogue	3000	15 actes 12 slots
SNIPS (Coccke <i>et al.</i> , 2018)	en	requêtes	14484	7 intentions 39 slots
TOP (Gupta <i>et al.</i> , 2018)	en	requêtes	44783	25 intentions 36 slots
The restaurant-8K dataset (Coope <i>et al.</i> , 2020)	en	requêtes	8198	5 slots 5
VocaDOM (Portet <i>et al.</i> , 2019)	fr	requêtes	4610	7 intentions 12 slots

TABLE 1 – Tableau récapitulatif des ensembles de données en NLU/SLU

Il existe de nombreux corpora SLU disponibles publiquement, chacun ayant ses propres caractéristiques et domaines d’application. Les tableaux 1 et 2 synthétisent les caractéristiques de ces données et nous détaillons ci-dessous chaque corpus.

Jeux de données	Multi-Domains	Multi-Intents	Inter-Domains	Annot. contextuelles	Annot. à plat	Annot. structurée ou semi-structurée
ATIS (Hemphill <i>et al.</i> , 1990)					X	
Frames (Asri <i>et al.</i> , 2017)				X		X
Massive (FitzGerald <i>et al.</i> , 2022)	X				X	
MEDIA (Devilleers <i>et al.</i> , 2004)					X	
mTOD (Schuster <i>et al.</i> , 2019)	X				X	
mTOP (Li <i>et al.</i> , 2020)	X				X	X
MultiDoGo (Peskov <i>et al.</i> , 2019)	X	X				X
MultiWOZ (Budzianowski <i>et al.</i> , 2018)	X	X	X	X		X
M2M (Shah <i>et al.</i> , 2018)						X
SNIPS (Couccke <i>et al.</i> , 2018)	X				X	
TOP (Gupta <i>et al.</i> , 2018)	X					X
The restaurant-8K dataset (Coope <i>et al.</i> , 2020)				X	X	
VocaDOM (Portet <i>et al.</i> , 2019)					X	

TABLE 2 – Tableau récapitulatif des caractéristiques des jeux de données en NLU/SLU

2.1.1 Jeux de données avec des requêtes simples

1. Ressources mono-lingues :

- (a) **ATIS** (Hemphill *et al.*, 1990) : Le corpus ATIS (Air Travel Information System) est l'un des corpus SLU les plus utilisés. Il contient des informations sur des compagnies aériennes et des commandes pour réserver des vols. La première version a été collectée suivant l'approche «Wizard-Of-OZ», mais les auteurs ont utilisé des transcriptions ASR pour les autres versions (Dahl *et al.*, 1994). Les premières annotations de ce corpus ont été effectuées à l'aide d'une requête SQL, puis ont été transférées au niveau global sous forme d'intentions, ainsi qu'au niveau mot sous forme des concepts (Béchet & Raymond, 2018). Ce corpus est en anglais et contient 4978 énoncés annotés en frame avec 17 intentions et 84 concepts.
- (b) **Frames** (Asri *et al.*, 2017) : Il s'agit d'un corpus avec des interactions humain-humain en anglais et qui contient des informations sur les réservations d'hôtel. Il a été publié pour encourager les recherches sur les systèmes conversationnels textuels. La notion de «mémoire» et l'exploitation de l'histoire conversationnelle ont été les premières questions abordées, où les auteurs ont rajouté à la tâche NLU, la tâche de "suivi des frames sémantiques". Des références et des identifiants des frames sémantiques ont été donc rajoutés aux annotations en acte de dialogue et en slot-valeur. Ce corpus peut être utilisé dans les tâches de compréhension et dans les tâches de suivi d'état de dialogue. 1369 dialogues ont été collectés suivant l'approche «Wizard-Of-OZ». Les énoncés au niveau utilisateurs et systèmes ont été annotés avec 20 types d'acte de dialogue et 16 types de slot.
- (c) **MEDIA** (Devilleers *et al.*, 2004) : Le corpus MEDIA contient des informations touristiques en français. Il a été collecté suivant l'approche «Wizard-Of-OZ». 1250 dialogues

ont été transcrits et annotés manuellement suivant une ontologie très riche au niveau des énoncés de l'utilisateur. Nous retrouvons 83 concepts de base regroupés dans un dictionnaire sémantique et 19 "spécifieurs" pouvant leur être associés. En lien avec le corpus MEDIA, PORT-MEDIA (Lefevre *et al.*, 2012) a été publié en français et en italien. Les méthodes de collecte étaient similaires en rajoutant des scénarios de dialogue variés. Les annotations étaient semi-automatiques à partir des systèmes de compréhension entraînés sur MEDIA suivies d'une vérification et correction manuelle. La version italienne a été générée par des traductions automatiques de la version française. Les questions de la complexité sémantique et les différents niveaux hiérarchiques ont été creusées et traitées sur une base linguistique solide.

- (d) **SNIPS** (Coucke *et al.*, 2018) : SNIPS est un corpus en anglais qui contient plusieurs domaines différents. Les données ont été collectées par des transcriptions ASR suivi d'une annotation et vérification manuelles. Il contient 7 intentions et 39 concepts. La même version de ce corpus a été publiée en d'autres langues comme le français et l'allemand. SNIPS a été l'origine d'un autre corpus, **Almawave SLU** (Bellomaria *et al.*, 2019), le premier ensemble de données en italien pour les expériences SLU. Il a été généré d'une manière semi-automatique en traduisant les énoncés et les labels et en remplaçant les entités ouverts (comme les noms de restaurants et des livres) par des références italiennes. La vérification et la correction manuelles a été effectuée pour les énoncés.
- (e) **TOP** (Gupta *et al.*, 2018) : Ce corpus a été publié pour étudier les problématiques sémantiques plus complexes. Les auteurs ont introduit une représentation hiérarchique appelée "*Task Oriented Parsing*" (TOP) pour les systèmes de dialogue basés sur des intentions et des concepts. Les énoncés ont été collectés par *crowd-sourcing* et annotés par deux annotateurs, un troisième annotateur peut intervenir en cas de désaccord. 44783 annotations avec 25 intentions et 36 slots ont été obtenues. Une version étendue du corpus avec 6 domaines supplémentaires a été publié dans **TOPv2** (Chen *et al.*, 2020).

The restaurant-8K dataset (Coope *et al.*, 2020) : Pour renforcer le travail d'extraction de concepts dans le cadre de dialogues, ce jeu de donnée qui comprend des conversations d'un système de réservation de restaurant, a été introduit. 8198 énoncés d'utilisateurs réels interagissant avec un système de dialogue déployé dans le domaine de la réservation de restaurants ont été annotés d'une manière contextuelle indiquant quels concepts ont été demandés par le système. Les réponses de système ne sont pas incluses dans l'ensemble des données, et il n'y a que 5 concepts.

VocaDOM (Portet *et al.*, 2019) : Afin de soutenir les tâches dans le cadre des systèmes "*Smart Home*" comme l'identification du locataire, la reconnaissance de la parole et les tâches SLU, ce corpus a été publié en rassemblant des interactions dans des conditions réelles de 11 participants dans une maison intelligente, la méthode «*Wizard-Of-OZ*» était la base de ce protocole. Les enregistrements ont été transcrits et annotés manuellement par des intentions et des concepts. Au total, le corpus contient 4610 énoncés en français étiquetés par 7 intentions et 12 concepts. Dans (Desot *et al.*, 2018), un jeu de donnée synthétiques du même domaine a été généré automatiquement à partir du corpus VocaDOM.

2. Ressources multi-lingues :

- (a) **Massive** (FitzGerald *et al.*, 2022) : Massive est un corpus multilingue qui contient des requêtes appartenant à 18 domaines. Sa publication a été motivée par le manque des jeux de données en plusieurs langues pour évaluer les modèles multilingues. Le corpus

SLURP (Bastianelli *et al.*, 2020), publié pour développer un assistant robotique personnel à domicile et pour des expériences SLU *End-to-end*, est la version d'origine du Massive. La version du corpus SLURP disponible publiquement est textuelle en anglais, elle a été collectée par les travailleurs de «Mechanical Turk» (AMT) et annotée manuellement au niveau "scénario" (domaine), "action" (Intention) et "entités" (concepts). Des traducteurs professionnels ont traduit les énoncés du corpus en 51 langues et ont également vérifié les frontières des concepts sur les tokens, aboutissant à la création du corpus Massive. Ce dernier est considéré comme une grande source des intentions (60) et de concepts (55) vu la diversité des domaines.

- (b) **mTOD** (Schuster *et al.*, 2019) : Il s'agit d'un ensemble de données multilingue qui permet d'étudier les méthodes d'apprentissage par transfert inter-linguistique. Ce corpus offre l'opportunité d'étudier les modèles sémantiques inter-langues et constitue le premier ensemble de données parallèles pour une tâche d'étiquetage de mots qui a été annoté selon les mêmes guides d'annotation dans plusieurs langues. Les auteurs ont collecté 43000 énoncés en anglais dans les domaines *ALARM*, *REMINDER*, et *WEATHER*, et ils ont demandé à des anglophones natifs de proposer des labels d'intentions utilisées par deux annotateurs pour étiqueter les énoncés et les valeurs par des concepts. Cette annotation a été vérifiée ensuite par un troisième annotateur. Des locuteurs natifs en espagnol et en thaï ont traduits les énoncés qui ont été aussi annotés par deux annotateurs. Il contient au total 12 types d'intentions et 11 concepts
- (c) **mTOP** (Li *et al.*, 2020) : En creusant la même problématique de la sémantique compositionnelle mise en évidence dans le corpus TOP (Gupta *et al.*, 2018) et en suivant la même logique de représentation hiérarchique, le corpus mTOP, a été publié. Cet ensemble de données est le premier qui contient des représentations sémantiques compositionnelles qui permettent l'annotation des requêtes imbriquées. Il a été publié avec les deux versions d'annotation : une plate et une autre compositionnelle. Les auteurs ont commencé par collecter une version en anglais des données, suivant la même approche dans (Gupta *et al.*, 2018), qui est traduite ensuite par des traducteurs professionnels. Le corpus mTOP est plus grand que TOP où nous avons 100.000 exemples avec 6 langues différentes, 11 domaines, 117 intentions et 78 concepts. En outre, une version parlée (**STOP**) a été publiée dans (Tomassello *et al.*, 2023) à partir de **TOPv2** pour encourager les recherches sur les approches end-to-end tout en focalisant sur les problématiques des requêtes compositionnelles.

2.1.2 Jeux de données conversationnels

1. Ressources mono-lingues :

- (a) **MultiDoGo** (Peskov *et al.*, 2019) : Cet ensemble de données a été collecté par «crowd-sourcing» dans le cadre du progrès des assistants virtuels et du manque des données pour leur développement. Ce corpus est composé par 81000 conversations, dont 15000 ont été annotées avec 6 domaines différents, 85 intentions et 73 slots. Dans le corpus disponible publiquement, les énoncés systèmes ne sont pas annotés. L'article présentant ces données a mis en valeur la possibilité d'avoir des multi-intentions en montrant que les annotations ont été réalisées par des experts selon deux niveaux : au niveau des tours des dialogues et au niveau des phrases, afin de garder l'ordre entre les énoncés coordonnés et leurs intentions. Toutefois, dans la version publiée, nous ne trouvons pas souvent cette illustration et nous pouvons même perdre le lien entre les différentes

intentions et leurs concepts.

- (b) **MultiWOZ** (Budzianowski *et al.*, 2018) : Il s'agit d'un corpus de dialogue multi-domaines en anglais, à grande échelle, souvent utilisé pour plusieurs tâches, notamment le suivi de l'état du dialogue, la politique de dialogue et les tâches de génération de dialogue. Il a été collecté à partir de la méthode «Wizard-Of-OZ» via un «crowd-sourcing». La première version de ce corpus a été publiée dans le but de faciliter la construction de systèmes de dialogue supervisés. Chaque énoncé dans les dialogues est annoté avec une séquence d'acte de dialogue. Cependant, la première version comporte des erreurs d'annotations, surtout au niveau de l'utilisateur, puisqu'elles ont été effectuées automatiquement à partir des annotations système (Eric *et al.*, 2019). Plusieurs versions ont été produites pour corriger ces erreurs et simplifier le format des annotations. Les versions MultiWOZ2.2 (Zang *et al.*, 2020) et MultiWOZ2.4 (Ye *et al.*, 2021) sont mieux adaptées aux tâches de suivi de l'état du dialogue, tandis que la version MultiWOZ2.3 (Han *et al.*, 2020) a des annotations utilisateur plus précises. Le corpus contient 7 domaines, 32 actes de dialogue et 27 slots.
- (c) **M2M** (Shah *et al.*, 2018) : M2M est une fusion de 2 données contenant des dialogues en anglais pour la réservation des restaurants et des tickets de cinémas. Les méthodes de collecte et de l'annotation ont été réalisées d'une manière automatique où un développeur de dialogue fournit un scénario et les chatbots génèrent des tours de conversations en les annotant par des actes de dialogue et par des slots. Ce processus était répété jusqu'à la fin des tours de dialogue soit par un acte "bye" soit en atteignant un seuil maximum de tours. Au total, nous avons 3000 dialogues annotés avec 15 actes de dialogues et 12 slots.

Dans la partie suivante, nous allons nous focaliser sur l'analyse de certains ensembles de données au niveau de leurs ontologies et schémas sémantiques.

2.2 Ontologies et Annotations

L'annotation sémantique repose généralement sur des ontologies basées sur des connaissances linguistiques permettant de définir des liens hiérarchiques entre les entités (Ma *et al.*, 2009). Dans les tâches de compréhension du langage, les ontologies permettent de décrire le lien sémantique entre les domaines, les intentions ou les actes de dialogue et leurs concepts (Loos, 2006). Ces ontologies se varient selon les schémas et les règles d'annotation des corpus, mais parfois, nous pouvons trouver la même logique surtout lorsque le schéma d'annotation est limité à un cadre sémantique simple. En d'autres termes, dans des jeux de données comme le cas d'ATIS (Hemphill *et al.*, 1990), de SNIPS (Coucke *et al.*, 2018), de Massive (FitzGerald *et al.*, 2022) et de mTOD (Schuster *et al.*, 2019), chaque énoncé est labélisé par une seule intention, la notion de multi-intentions ou le croisement entre les domaines (inter-domaines) sont donc absents pour ces corpora.

En plus, les données citées peuvent avoir des domaines en commun et donc des similarités au niveau des concepts. Par exemple, le corpus ATIS n'a qu'un seul domaine (réservation de vol) qui peut se croiser avec le domaine "airline" dans le corpus MultiDoGo (Peskov *et al.*, 2019). SNIPS a aussi plusieurs domaines en commun avec TOP (Gupta *et al.*, 2018) («Restaurant, Weather, Music»...). Cependant, les différentes façons d'exprimer les intentions et le choix des concepts entraînent une grande diversité au niveau des ontologies, ce qui rend leur unification assez problématique. Les intentions dans ATIS sont des noms simples, par contre dans mTOD et SNIPS elles sont composées

d'un acte de dialogue avec le domaine («Show_alarms», «BookRestaurant»). Les actes de dialogue dans MultiWOZ sont aussi composés par le domaine associé à l'acte («Hotel-Inform»), mais le choix des actes est associé à la sémantique derrière les concepts plus qu'au sens global des énoncés. Autrement dit, les concepts comme «Phone» et «Car» pour le domaine «Taxi» sont toujours associés à l'acte «Request».

Les concepts peuvent être aussi variés au niveau de leur composition, ils peuvent se représenter comme une seule entité («country») ou une entité composée («party_size_description»). Nous avons remarqué aussi que dans toutes ces données les concepts sont plus compositionnels et reliés à leurs domaines. Dans le corpus Massive, on a par exemple les slots «sport_type», «drink_type» et «alarm_type», à l'antipode de MultiWOZ qui n'a que le concept «Type» partagé par les domaines «Attraction» et «Hotel». Dans le corpus ATIS, les concepts présentent un niveau plus haut au niveau de sa composition où les prépositions peuvent faire partie des slots (comme «from» dans «fromloc.city_name»), ou nous pouvons même trouver des concepts imbriqués (le slot «depart_time.period_of_day» est composé par le slot «time» et le slot «period_of_day»).

Le corpus MEDIA (Devillers *et al.*, 2004) est par ailleurs assez particulier au niveau de son schéma d'annotation. L'ontologie a été basée sur un niveau sémantique haut qui essaie de relever le lien hiérarchique des labels sémantiques. Les frames sémantiques ne se composent pas que par des paires de slot-valeur, mais aussi par un spécifieur qui définit les relations entre les entités. Une annotation des modes a été aussi effectuée et attachée aux concepts. Ainsi, la représentation hiérarchique a été recomposée par la combinaison des "spécifieurs" et des concepts.

Il convient également de noter que ces ensembles de données varient considérablement en termes de complexité linguistique. Dans l'article (Bechet *et al.*, 2022), divers phénomènes linguistiques qui peuvent impacter les performances des systèmes SLU sont observés. Certains corpora ne reflètent pas les caractéristiques des interactions dans des conditions réelles, tandis que d'autres sont plus difficiles pour les modèles. L'approche proposée pour évaluer la qualité et comparer les corpora, comme décrite dans (Bechet *et al.*, 2022) et (Béchet & Raymond, 2019), peut contribuer à la question d'unification des données.

En ce qui concerne la méthode de projection des frames sémantiques, le format d'annotation "BIO" à plat (cf. tableau 3) est souvent le paradigme le plus utilisé, notamment pour les données avec de simples requêtes, afin de faire un étiquetage de séquence facilement avec les modèles pré-entraînés à base de Transformers de type Bert (Devlin *et al.*, 2018). Néanmoins, ce paradigme est limité si nous voulons passer à la compréhension des énoncés plus complexes en exploitant le contexte de la conversation. Nous avons remarqué les limites de ce paradigme avec le corpus MultiWOZ où un seul énoncé peut avoir plusieurs domaines ou plusieurs intentions (cf. 4). Les annotations d'origine de ce corpus se représentent sous un format plus structuré en format json où nous pouvons trouver le lien entre les différents concepts et leurs intentions. Une suggestion d'une annotation à plat a été proposée dans (Lee *et al.*, 2019), mais les difficultés des annotations en contexte, notamment pour les concepts sans informations d'empan, n'ont pas été entièrement résolues. Ces entités sont justement définies comme des "slots de catégories", où leurs valeurs se trouvent dans les énoncés précédents, ou bien implicitement dans le sens global de l'énoncé. Les slots "Parking" et "Post" dans la table 4 sont un exemple des concepts de catégories.

Les différents paradigmes d'annotation et les différentes hiérarchies sémantiques des ontologies mettent en valeur la difficulté de représenter les données d'une manière structurée où le contexte peut-être bien exploité. Dans la section suivante, cette question sera creusée où nous allons nous

Enoncés	[CLS]	I'm	traveling	to	dallas	from	philadelphia
Annotation	Flight	O	O	O	B-toloc.city_name	O	B-fromloc.city_name

TABLE 3 – Annotation à plat en BIO du corpus ATIS : "B" pour «Begining», "I" pour «Inside» et "O" pour «Outside»

Enoncés	Annotation en acte de dialogue
◇ I won't have a car, so parking isn't important	"Hotel-Inform": [{"Parking", "no"}]
◇ Can I have the postcode for the attraction, I also need a Taxi	"Attraction-Inform": [{"Post", "?"}], "Taxi-Inform": [{"none", "none"}]

TABLE 4 – Exemples d'énoncés annotés dans le corpus MultiWOZ2.3

Enoncés	
take grandma Jane off the call	
Annotation	[IN:update_call [SL:contact_removed [IN:get_contact [SL:type_relation grandma] [SL:contact Jane]]]]

TABLE 5 – Exemple d'annotation dans le corpus mTOP

intéresser aux différentes motivations qui sous-tendent l'utilisation de représentations structurées des étiquettes sémantiques des données pouvant être une solution de certaines limites des annotations à plat.

3 Projections des annotations et représentations sémantiques structurées dans le NLU/SLU

Nous avons présenté dans la section précédente les différences et les similitudes entre les jeux de données utilisés pour les tâches de compréhension du langage au niveau de leurs ontologies et leurs schémas d'annotation. Dans cette section, nous allons retracer d'une manière générale l'historique des différentes représentations sémantiques utilisées en compréhension du langage, ainsi que les enjeux liés aux représentations actuelles dans le contexte des systèmes de dialogue. Nous étudierons ensuite les projections structurées en illustrant les problématiques liées aux différents schémas à plat, ainsi que le potentiel des formats structurés, en particulier ceux basés sur la méthode des graphes.

3.1 Représentations sémantiques

Les interprétations sémantiques peuvent être considérées comme un processus de traduction, réalisé par un parseur sémantique, entre les mots d'une phrase et les représentations sémantiques du langage, comme montré dans (Dinarelli, 2010). Les représentations sémantiques permettent la modélisation des énoncés et de leurs interprétations sémantiques pour les machines. Elles peuvent être sous forme de logique formelle (Zettlemoyer & Collins, 2012), des frames sémantiques (Dinarelli *et al.*, 2009) ou encore des graphes sémantiques (Banarescu *et al.*, 2013).

Avec le progrès des systèmes de dialogue et des modèles modernes basés sur des approches de

statistiques et de probabilités, les interprétations sémantiques en SLU sont principalement basées le plus souvent sur l'identification des intentions ou des actes de dialogue et des slots. En outre, ces annotations notamment au niveau des slots sont très similaires aux annotations par frame sémantique, comme dans FrameNet (Baker *et al.*, 1998), dans la mesure où la représentation SLU est basée sur des attributs, qui sont des unités sémantiques, instanciées par séquences de mots. Contrairement aux frames, les attributs en SLU n'ont pas besoin d'explicitier les relations sémantiques entre les éléments de la phrase (Dinarelli, 2010). Au niveau concepts, les étiquettes consistent à identifier les éléments clés de l'énoncé, comme les entités nommées, les dates ou les destinations. Quant aux intentions, les annotations sont utilisées pour aider les systèmes SLU à comprendre le but global de l'utilisateur. Nous pouvons trouver ainsi des annotations en acte de dialogue pour mieux représenter les intentions. Par ailleurs, des recherches ont été menées pour représenter les attributs sémantiques du dialogue sous forme d'une ontologie abstraite, générique et structurée en exploitant les représentations AMR pour l'analyse sémantique du dialogue. (Bonial *et al.*, 2020).

Il est important ainsi de réfléchir à la question de la projection des informations sémantiques pour les traduire en entrées exploitables par les modèles. L'approche la plus courante pour la tâche de compréhension du langage repose sur les approches de l'étiquetage de séquence. La projection à plat en format BIO facilite cette tâche notamment pour les approches jointes pour prédire les intentions et les concepts simultanément. Comme il est montré dans l'exemple 3, ces approches associent généralement l'intention au token de classification globale [CLS] et détectent les concepts sur chaque token concerné avec une étiquette B ou I lorsque cela est possible.

Bien que la projection basée sur le format BIO soit utile pour l'utilisation des modèles à base de Transformers, elle ne permet pas de tirer parti des structures sémantiques hiérarchiques et imbriquées. Ainsi, des recherches sur des annotations plus structurées ont été étudiées, notamment avec l'utilisation des approches de séquence à séquence (*seq2seq*). Les chercheurs dans (Li *et al.*, 2020) ont proposé une représentation composée découplée (cf. 5) pour représenter des intentions imbriquées dans les slots. De même les expériences dans (Hu *et al.*, 2022) présentent le NLU comme une tâche de génération des graphes composés par des nœuds pour les labels. Dans la section suivante, nous allons montrer les différentes motivations et illustrations des schémas de représentation structurée.

3.2 Les représentations structurées des frames sémantiques

Selon (Devillers *et al.*, 2004), la représentation sémantique des annotations d'un corpus a été définie comme un moyen pour représenter les frames sémantiques d'une manière générique et complète selon la tâche mais qui permet aussi l'annotation des corpora larges d'une manière simple. Par conséquent, les schémas de représentation à plat ont été le centre des annotations dans la majorité des corpora publics. Cependant, les représentations hiérarchiques sont plus expressives et permettent le lien entre les sous-structures (Tur & De Mori, 2011).

Dans la section 2.2 nous avons présenté quelques tentatives de projection des annotations d'origine du corpus multi-domaines MultiWOZ au paradigme à plat dans (Lee *et al.*, 2019). Comme nous remarquons dans l'exemple 6, le schéma de représentation a repris l'idée de compositionnalité des labels notamment pour les concepts de catégorie, où l'ensemble de l'acte de dialogue, le concept et sa valeur ont été projetés au niveau global [CLS]. Il est important de souligner cependant que cette projection demeure limitée si l'on veut prédire des représentations plus complexes.

En outre dans (Gupta *et al.*, 2018) le paradigme à plat a été remis en question pour les requêtes

Enoncés	[CLS]	I	need	parking
Annotation	Hotel-Inform+Parking*yes	O	O	O

TABLE 6 – Annotation à plat en BIO du corpus MultiWOZ

plus complexes. En effet le corpus TOP est modélisé d’une manière compositionnelle qui autorise les intentions imbriquées, il s’agit en effet d’une représentation hiérarchique similaire aux arbres syntaxiques. Vu la complexité de la représentation et pour faciliter l’utilisation du même schéma en plusieurs langues, (Aghajanyan *et al.*, 2020) ont proposé une extension de la représentation compositionnelle en représentation découplée qui a été utilisée dans (Li *et al.*, 2020). La figure 1 illustre la projection : le premier niveau de l’arbre correspond à l’intention, qui peut inclure un ou plusieurs concepts. Ces derniers peuvent à leur tour comporter des intentions ou une séquence de mot comme une valeur. En somme, les auteurs ont démontré que cette projection est un compromis entre le paradigme traditionnel à plat et la représentation en logique formelle. De même, les expériences faites dans (Cheng *et al.*, 2020) s’inscrivent dans le même cadre de la sémantique compositionnelle en affirmant que la compositionnalité peut simplifier la compréhension pour faciliter la tâche de suivi d’état de dialogue.

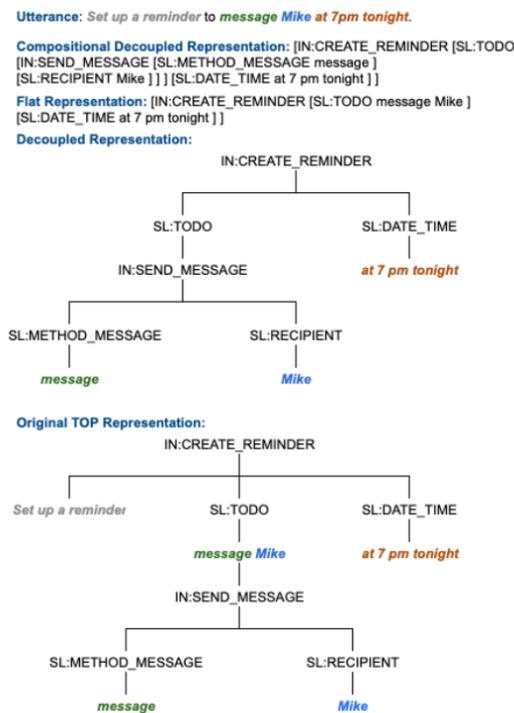


FIGURE 1 – Représentation compositionnelle découplée vs représentation plate dans mTOP (Li *et al.*, 2020)

De surcroît, les limitations du paradigme à plat et les motivations d’une représentation plus structurée ont été discutées dans (Hu *et al.*, 2022). L’article a proposé «DMR» une représentation en graphe qui comprend des nœuds d’Intention, de slots, des opérateurs indiquant les coréférences et la conjonction, et des mots-clés pour quelques éléments spéciaux en sémantique comme la négation. Une définition d’une nouvelle ontologie du domaine «Fast Food» du corpus MultiDoGo et une ré-annotation structurée qui lie les tours de dialogue permettant les annotations en contexte ont été l’objet de cet article. Des notions comme la "quantification" et "les adjectifs modificateurs" ont été ainsi soulignées.

Un exemple de leur représentation est dans la figure 2.

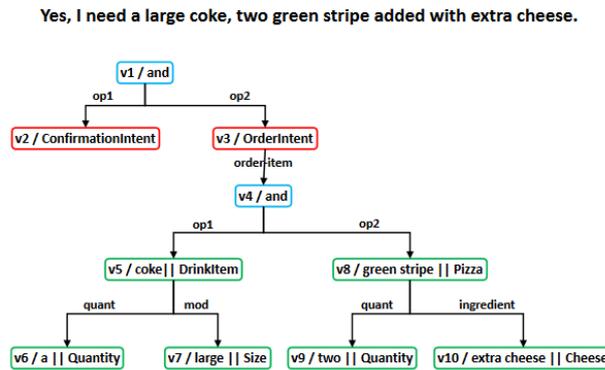


FIGURE 2 – Représentation en graphe dans DMR (Hu *et al.*, 2022)

Une proposition similaire à (Hu *et al.*, 2022) a été suggérée dans (Abrougui *et al.*, 2022). Cette proposition illustre des projections effectuées sur les annotations d’origine du corpus MultiWOZ2.3 sans qu’il soit nécessaire de modifier l’ontologie. Tel qu’indiqué dans la figure 3 Les actes de dialogue et les slots-valeurs sont transposés comme des nœuds, tandis que les slots sont représentés sous forme d’arcs reflétant la hiérarchie entre les labels. Les cas complexes, tels que les multi-intentions ou les intentions imbriquées peuvent être projetés dans ce format grâce à l’encodage Penman (Kasper, 1989), également utilisé dans les analyses AMR (Banarescu *et al.*, 2013).

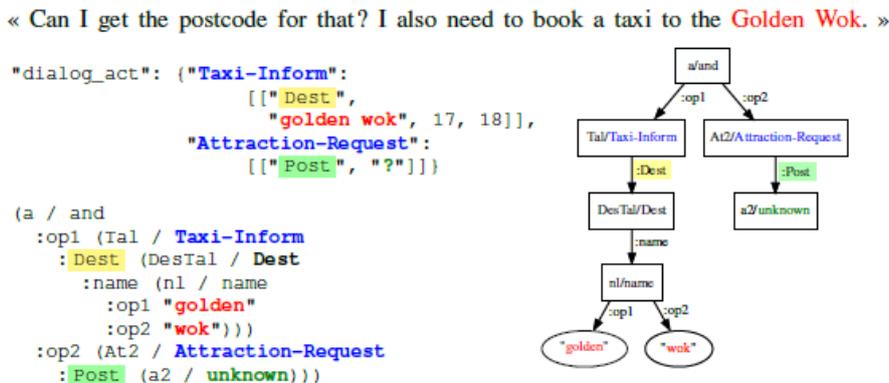


FIGURE 3 – Représentation en graphe et encodage Penman dans (Abrougui *et al.*, 2022)

Avec le développement des modèles de réseaux de neurones et l’essor des approches *seq2seq*, les représentations structurées pour les annotations sont devenues facilement exploitables par les systèmes NLU. Dans le même contexte, nous présentons nos perspectives et projets de recherche dans la partie suivante.

3.3 NLU avec une annotation structurée sur le corpus MultiWOZ

Les schémas d’annotation structurés et hiérarchiques offrent la possibilité d’étudier d’une façon plus approfondie une sémantique complexe et compositionnelle. Nous avons choisi donc de travailler sur le corpus MultiWOZ présentant des défis intéressants en raison de sa complexité (multi-domaines,

croisement entre les différentes intentions et concepts). Nous avons choisi de travailler en particulier sur la version 2.3, car elle offre des annotations utilisateur plus précises que les autres versions. Cependant, nous avons constaté que les annotations de cette version comportent des erreurs selon les règles de la logique NLU. Pour remédier à cela, notre objectif est de corriger et d'enrichir l'ontologie du corpus et de proposer un modèle de représentation qui permet une annotation contextuelle plus précise liant les tours de dialogue entre eux.

3.3.1 Corrections des annotations

Nous avons constaté dans la section 2 que les premières annotations au niveau utilisateurs de MultiWOZ ont été générées automatiquement. Bien que les auteurs dans MultiWOZ2.3 aient apporté des corrections, des précisions manquent encore. Nous avons donc commencé à examiner les conversations en identifiant une liste d'expressions, tels que «I want to travel from» ou «I want to arrive by», qui permettent de sélectionner un ensemble d'énoncés à vérifier et à corriger. Notre objectif dans cette étape est de garantir la cohérence entre les actes de dialogue, les concepts et les valeurs sans modifier l'ontologie de base.

En d'autres termes, l'exemple dans le tableau 7 doit être corrigé comme "Hotel-Request": [{"Internet", "free"}] puisque l'utilisateur demande une information spécifique sur l'accès gratuit à Internet. Toutefois, la combinaison de la valeur "free" avec l'acte "Request" n'existe pas dans l'ontologie du corpus. Les actes dans MultiWOZ sont choisis en fonction du type des concepts, (slots-valeur en cas des concepts de catégorie), plutôt que du sens global de l'énoncé. Ainsi dans cette ontologie, la valeur de catégorie "free" est associée au concept "Price", qui est lui-même associé à l'acte "Inform". Cependant, le slot "Price" n'est utilisé que pour définir les prix des hôtels et des restaurants. Dans le cas du slot "Internet", on a 4 valeurs principales : "yes" "no" "dontcare" associées à l'acte "Inform" et la valeur "?" associée à l'acte "Request". Etant donné que l'énoncé dans l'exemple 7 est une simple demande d'information, nous nous limitons donc à corriger l'annotation en "Hotel-Request": [{"Internet", "?"}]. Les travaux de correction sont actuellement en cours et leur évaluation fera l'objet d'une publication future.

Stratégie	Recherche de l'expression "do they offer free wifi ?"
Type de correction	semi-automatique
Annotation d'origine	"Hotel-Request": [{"Internet", "free"}]
Correction	"Hotel-Request": [{"Internet", "?"}]

TABLE 7 – Exemple de correction dans MultiWOZ2.3

3.3.2 Projection structurée

Dans (Abrougui *et al.*, 2022) nous trouvons une projection structurée des intentions, concepts et valeurs du MultiWOZ sans effectuer aucune modification (cf. figure 3). L'avantage de cette représentation est sa capacité de projeter les jeux de données standards comme ATIS, et les requêtes les plus complexes comme dans TOP et dans MultiWOZ. Dans cette étape nous avons deux étapes. Tout d'abord, nous visons à rajouter de nouveaux labels à l'ontologie du corpus comme l'acte "Confirmation", ou comme les adverbes représentés par un concept de catégorie, comme il est

montré dans l'exemple 8. Nous réfléchissons aussi à la question si l'acte doit être associé ou dissocié du domaine.

Enoncé	Can you also give me some information about Finches Bed and Breakfast? We 're thinking of staying there .
Annotation d'origine	"Hotel-Inform": [{"Name", "finches bed and breakfast"}]
Proposition d'annotation	"Hotel": [{"Inform-Name", "finches bed and breakfast"}, {"Request-Info", "?"}, {"Modifier", "maybe"}]

TABLE 8 – Proposition d'une nouvelle annotation du corpus MultiWOZ2.3

Notre objectif dans un second temps est de reprendre cette structure en rajoutant des annotations de corréférence en liant les antécédents et les anaphores avec les variables utilisées dans la notation Penman. La figure 4 ci-dessous qui présente un exemple fictif illustre ce projet. En examinant l'énoncé, nous remarquons la présence de deux actes de dialogues distincts qui partagent le slot "Area" clairement indiqué par l'utilisation de l'adverbe "there". La représentation sémantique de cette structure avec un format à plat serait particulièrement difficile. En outre, il convient de souligner la problématique de l'implicite soulevée par l'expression "I have a car", qui fait référence au concept de catégorie "Parking" et à sa valeur normalisée "yes". Afin de faciliter la représentation de ces entités et de leurs liens, l'utilisation du format penman basé sur les variables (comme "a1" dans la figure) s'avère être un choix judicieux.

Notre objectif est en effet de faire une tâche NLU mais en couvrant toutes les informations possibles, comme l'annotation de la corréférence.

"I want a restaurant in Bastille, and I also need a hotel there, I have a car."

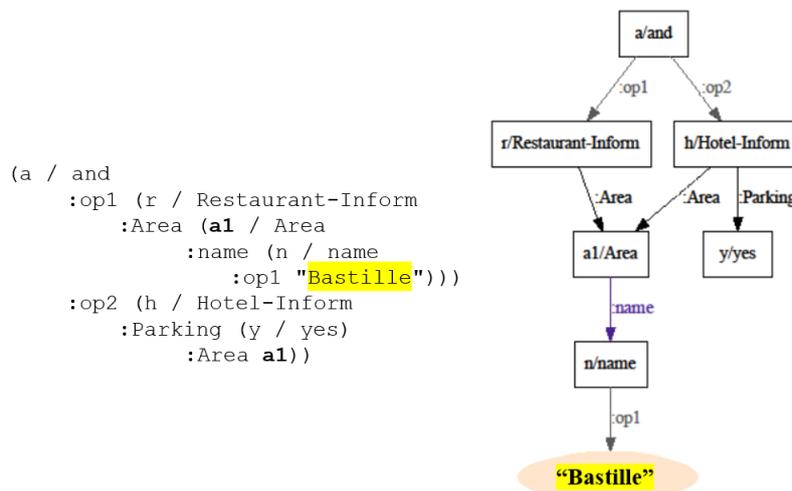


FIGURE 4 – schéma structuré pour la corréférence

3.3.3 Sémantique des labels

Il est vrai que les schémas de représentation structurés sont considérés comme des outils efficaces pour représenter les données avec des informations sémantiques complexes, mais la sémantique

sous-jacente des étiquettes est tout aussi importante. Les auteurs dans (Athiwaratkun *et al.*, 2020) ont souligné cette question en expliquant que les modèles de langage génératifs offrent un moyen naturel d’incorporer le sens des étiquettes dans les tâches de compréhension. Ils ont représenté la sortie en une séquence augmentée qui contient la séquence d’entrée avec leurs labels. Ils ont nettoyé les labels de tous les symboles et les ont représentés sous leur forme de langage naturel (par exemple, "AddToPlaylist" est devenu "Add To Playlist") afin de mieux exploiter les capacités des modèles de langage génératifs à comprendre le langage naturel.

Par ailleurs, l’unification des ensembles de données SLU reste un défi dans ce domaine, car les ontologies et les noms des concepts varient considérablement. Néanmoins, si des noms de concepts sont communs et que les modèles génératifs sont capables de comprendre le langage naturel, il est possible d’unifier ces jeux de données plus facilement sans effectuer de nombreuses modifications au niveau de leurs ontologies. Nous avons testé cette hypothèse en examinant les données SNIPS et MultiWOZ.

En effet, SNIPS a un domaine en commun avec MultiWOZ, qui est "Restaurant". Il existe également le concept "Name" associé à cinq catégories différentes (tels que "restaurant_name", "movie_name" et "location_name"). Nous avons fusionné les deux en deux étapes : **(1)** la première consiste à fusionner les deux données sans changer les ontologies. **(2)** La deuxième consiste à changer les concepts de SNIPS en une seule entité ("movie_name" devient "name") et les intentions en une forme "Domaine-Acte" ("SearchScreeningEvent" devient "Event-Search"). En ce qui concerne MultiWOZ, nous avons étendu les concepts en leur forme de langage naturel ("Dest" devient "Destination"), nous les avons mis tous en minuscules et nous avons appris un modèle mT5 (Xue *et al.*, 2020) qui prend en entrée les énoncés et génère les représentations structurées comme indiqué dans la figure 3, codées en format Penman.

Les performances globales dans le tableau 9 montrent que la fusion des deux corpora n’a pas engendré de résultats significatifs. Bien qu’une légère baisse ait été observée pour SNIPS au niveau de l’accuracy globale, MultiWOZ n’a pas vraiment changé. Nous avons ensuite examiné les performances au niveau des concepts composés comportant le mot "name" dans leur nom et simplifiés au format MultiWOZ.

Le tableau 10 présente les performances des F1-mesures au niveau Intention(slot,valeur) correspondant aux slots modifiés et leurs domaines respectifs.

Nous constatons une amélioration au niveau des résultats pour le slot "movie_name" et en particulier pour "restaurant_name" suite à la modification de l’ontologie de SNIPS (MS-S^O), avec une augmentation de 16%. Ce label partage en effet le domaine et le nom du concept avec MultiWOZ, et il semble que le modèle génératif a bien capturé la sémantique derrière.

En revanche, les performances pour "location_name" ont diminué. Les résultats pour "object_name" ont également diminué avec l’intention "SearchScreeningEvent" mais ils s’améliorent avec 1% avec l’intention "RateBook". Tout cela montre que l’association entre les différents labels affectent leurs significations, et il est possible pour les modèles génératifs de les prédire si on peut les représenter d’une manière cohérente. L’augmentation significative de certains concepts met en évidence l’importance de l’unification des ontologies. En effet, lorsque les énoncés et leurs labels partagent une sémantique commune, et sont bien définis et unifiés sous le même label, cela peut contribuer à renforcer les performances des systèmes N/SLU.

Nous envisageons d’approfondir cette approche et étudier précisément les ontologies en exploitant à la fois les représentations structurées des frames sémantiques et le potentiel des modèles génératifs

pour l'unification des jeux de données en compréhension du langage naturel et parlé.

SNIPS			
	S-S	MS-S	MS-S ^O
F1 intention	98,1	97,8	98,6
F1 (concept,valeur)	95,0	94,8	94,9
F1 Intent(concept,valeur)	94,7	94,6	94,6
Accuracy global	88,7	87,8	88,0
MultiWOZ 2.3			
	M-M	MS-M	MS-M ^O
F1 intention	96,2	96,2	96,3
F1 (concept,valeur)	94,7	94,9	94,9
F1 Intent(concept,valeur)	94,1	94,2	94,3
Accuracy global	87,6	87,7	87,5

TABLE 9 – Résultats des expériences sur la sémantique des labels : apprentissage et test sur SNIPS (S-S), apprentissage et test sur MultiWOZ (M-M), apprentissage sur MultiWOZ et SNIPS sans modifications de l'ontologie et test sur SNIPS (MS-S), apprentissage sur MultiWOZ et SNIPS sans modifications de l'ontologie et test sur MultiWOZ (MS-M), apprentissage sur MultiWOZ et SNIPS avec modifications de l'ontologie et test sur SNIPS (MS-S^O), apprentissage sur MultiWOZ et SNIPS avec modifications de l'ontologie et test sur MultiWOZ (MS-M^O)

	S-S	MS-S	MS-S ^O
SearchScreeningEvent+movie_name (event-search+name)	86,7	77,1	89,6
SearchScreeningEvent+location_name (event-search+name)	97,9	97,9	91,7
SearchCreativeWork+object_name (work-search+name)	85,9	85,6	82,9
AddToPlaylist+entity_name (music-add+name)	74,6	80,0	71,9
RateBook+object_name (book-rate+name)	95,0	97,5	96,3
BookRestaurant+restaurant_name (restaurant-book+name)	73,3	83,9	89,7

TABLE 10 – F1 mesure niveau Inent(concept,valeur) pour les concepts composés par le slot "name" : format d'origine (format modifié)

4 Conclusion

La compréhension du langage naturel et parlé est une tâche fondamentale dans les systèmes de dialogue. Les deux tâches connues visent à comprendre les commandes de l'utilisateur et ses interactions avec un agent robotique. Dans cet article, nous avons présenté les différents jeux de données utilisés dans ce domaine et nous avons comparé leurs ontologies et leurs schémas d'annotation. Les représentations sémantiques structurées et la méthode de graphe ont été mises en valeur pour leur potentiel à refléter la hiérarchie sémantique entre les frames sémantiques et à exploiter le contexte. La sémantique des labels a également fait l'objet d'une discussion dans nos projets de recherche basés fondamentalement sur des expériences sur le corpus conversationnel MultiWOZ. Dans le but d'unifier les données et d'exploiter mieux le contexte pour construire des systèmes robustes, nous envisageons d'approfondir nos études des ontologies et des représentations structurées tout en exploitant l'histoire conversationnelle.

Références

- ABROUGUI R., DAMNATI G., HEINECKE J. & BÉCHET F. (2022). Étiquetage ou génération de séquences pour la compréhension automatique du langage en contexte d'interaction ? In *Traitement Automatique des Langues Naturelles (TALN 2022)*, p. 64–73 : ATALA.
- AGHAJANYAN A., MAILLARD J., SHRIVASTAVA A., DIEDRICK K., HAEGER M., LI H., MEHDAD Y., STOYANOV V., KUMAR A., LEWIS M. *et al.* (2020). Conversational semantic parsing. *arXiv preprint arXiv :2009.13655*.
- ASRI L. E., SCHULZ H., SHARMA S., ZUMER J., HARRIS J., FINE E., MEHROTRA R. & SULEMAN K. (2017). Frames : a corpus for adding memory to goal-oriented dialogue systems. *arXiv preprint arXiv :1704.00057*.
- ATHIWARATKUN B., SANTOS C. N. D., KRONE J. & XIANG B. (2020). Augmented natural language for generative sequence labeling. *arXiv preprint arXiv :2009.13272*.
- BAKER C. F., FILLMORE C. J. & LOWE J. B. (1998). The berkeley framenet project. In *COLING 1998 Volume 1 : The 17th International Conference on Computational Linguistics*.
- BANARESCU L., BONIAL C., CAI S., GEORGESCU M., GRIFFITT K., HERMJAKOB U., KNIGHT K., KOEHN P., PALMER M. & SCHNEIDER N. (2013). Abstract meaning representation for sembanking. In *Proceedings of the 7th linguistic annotation workshop and interoperability with discourse*, p. 178–186.
- BASTIANELLI E., VANZO A., SWIETOJANSKI P. & RIESER V. (2020). Slurp : A spoken language understanding resource package. *arXiv preprint arXiv :2011.13205*.
- BÉCHET F. & RAYMOND C. (2018). Is atis too shallow to go deeper for benchmarking spoken language understanding models ? In *InterSpeech 2018*, p. 1–5.
- BÉCHET F. & RAYMOND C. (2019). Benchmarking benchmarks : introducing new automatic indicators for benchmarking spoken language understanding corpora. In *Interspeech*.
- BECHET F., RAYMOND C., HAMANE A., ABROUGUI R., MARZINOTTO G. & DAMNATI G. (2022). Can we predict how challenging spoken language understanding corpora are across sources, languages, and domains ? In *Conversational AI for Natural Human-Centric Interaction : 12th International Workshop on Spoken Dialogue System Technology, IWSDS 2021, Singapore*, p. 33–45 : Springer.
- BELLOMARIA V., CASTELLUCCI G., FAVALLI A. & ROMAGNOLI R. (2019). Almwave-slu : A new dataset for slu in italian. *arXiv preprint arXiv :1907.07526*.
- BONIAL C., DONATELLI L., ABRAMS M., LUKIN S., TRATZ S., MARGE M., ARTSTEIN R., TRAUM D. & VOSS C. (2020). Dialogue-amr : abstract meaning representation for dialogue. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, p. 684–695.
- BUDZIANOWSKI P., WEN T.-H., TSENG B.-H., CASANUEVA I., ULTES S., RAMADAN O. & GAŠIĆ M. (2018). Multiwoz—a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. *arXiv preprint arXiv :1810.00278*.
- CHEN X., GHOSHAL A., MEHDAD Y., ZETTMLOYER L. & GUPTA S. (2020). Low-resource domain adaptation for compositional task-oriented semantic parsing. *arXiv preprint arXiv :2010.03546*.
- CHENG J., AGRAWAL D., ALONSO H. M., BHARGAVA S., DRIESEN J., FLEGO F., GHOSH S., KAPLAN D., KARTSAKLIS D., LI L. *et al.* (2020). Conversational semantic parsing for dialog state tracking. *arXiv preprint arXiv :2010.12770*.

- COOPE S., FARGHLY T., GERZ D., VULIĆ I. & HENDERSON M. (2020). Span-convert : Few-shot span extraction for dialog with pretrained conversational representations. *arXiv preprint arXiv :2005.08866*.
- COUCKE A., SAADE A., BALL A., BLUCHE T., CAULIER A., LEROY D., DOUMOIRO C., GISSELBRECHT T., CALTAGIRONE F., LAVRIL T. *et al.* (2018). Snips voice platform : an embedded spoken language understanding system for private-by-design voice interfaces. *arXiv preprint arXiv :1805.10190*.
- DAHL D. A., BATES M., BROWN M. K., FISHER W. M., HUNICKE-SMITH K., PALLETT D. S., PAO C., RUDNICKY A. & SHRIBERG E. (1994). Expanding the scope of the atis task : The atis-3 corpus. In *Human Language Technology : Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994*.
- DESOT T., RAIMONDO S., MISHAKOVA A., PORTET F. & VACHER M. (2018). Towards a french smart-home voice command corpus : Design and nlu experiments. In *Text, Speech, and Dialogue : 21st International Conference, TSD 2018, Brno, Czech Republic, September 11-14, 2018, Proceedings 21*, p. 509–517 : Springer.
- DEVILLERS L., MAYNARD H., ROSSET S., PAROUBEK P., MCTAIT K., MOSTEFA D., CHOUKRI K., CHARNAY L., BOUSQUET C., VIGOUROUX N. *et al.* (2004). The french media/evalda project : the evaluation of the understanding capability of spoken language dialogue systems. In *LREC*.
- DEVLIN J., CHANG M.-W., LEE K. & TOUTANOVA K. (2018). Bert : Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv :1810.04805*.
- DINARELLI M. (2010). *Spoken language understanding : from spoken utterances to semantic structures*. Thèse de doctorat, University of Trento.
- DINARELLI M., QUARTERONI S., TONELLI S., MOSCHITTI A. & RICCARDI G. (2009). Annotating spoken dialogs : from speech segments to dialog acts and frame semantics. In *Proceedings of SRS� 2009, the 2nd Workshop on Semantic Representation of Spoken Language*, p. 34–41.
- ERIC M., GOEL R., PAUL S., SETHI A., AGARWAL S., GAO S. & HAKKANI-TUR D. (2019). Multiwoz 2.1 : Multi-domain dialogue state corrections and state tracking baselines. *arXiv :1907.01669*.
- FITZGERALD J., HENCH C., PERIS C., MACKIE S., ROTTMANN K., SANCHEZ A., NASH A., URBACH L., KAKARALA V., SINGH R. *et al.* (2022). Massive : A 1m-example multilingual natural language understanding dataset with 51 typologically-diverse languages. *arXiv preprint arXiv :2204.08582*.
- GUPTA S., SHAH R., MOHIT M., KUMAR A. & LEWIS M. (2018). Semantic parsing for task oriented dialog using hierarchical representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, p. 2787–2792, Brussels, Belgium : Association for Computational Linguistics. DOI : [10.18653/v1/D18-1300](https://doi.org/10.18653/v1/D18-1300).
- HAN T., LIU X., TAKANOBU R., LIAN Y., HUANG C., WAN D., PENG W. & HUANG M. (2020). Multiwoz 2.3 : A multi-domain task-oriented dialogue dataset enhanced with annotation corrections and co-reference annotation. *arXiv :2010.05594*.
- HEMPHILL C. T., GODFREY J. J. & DODDINGTON G. R. (1990). The atis spoken language systems pilot corpus. In *Speech and Natural Language : Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990*.
- HU X., DAI J., YAN H., ZHANG Y., GUO Q., QIU X. & ZHANG Z. (2022). Dialogue meaning representation for task-oriented dialogue systems. *arXiv preprint arXiv :2204.10989*.

- KASPER R. T. (1989). A flexible interface for linking applications to penman’s sentence generator. In *Speech and Natural Language : Proceedings of a Workshop Held at Philadelphia, Pennsylvania, February 21-23, 1989*.
- LEE S., ZHU Q., TAKANOBU R., LI X., ZHANG Y., ZHANG Z., LI J., PENG B., LI X., HUANG M. *et al.* (2019). Convlab : Multi-domain end-to-end dialog system platform. *arXiv preprint arXiv :1904.08637*.
- LEFEVRE F., MOSTEFA D., BESACIER L., QUIGNARD M., CAMELIN N., FAVRE B., JABAIAN B., BARAHONA L. M. R. *et al.* (2012). Leveraging study of robustness and portability of spoken language understanding systems across languages and domains : the portmedia corpora. In *The International Conference on Language Resources and Evaluation*.
- LI H., ARORA A., CHEN S., GUPTA A., GUPTA S. & MEHDDAD Y. (2020). Mtop : A comprehensive multilingual task-oriented semantic parsing benchmark. *arXiv preprint arXiv :2008.09335*.
- LOOS B. (2006). Scaling natural language understanding via user-driven ontology learning. In *Proceedings of the Third Workshop on Scalable Natural Language Understanding*, p. 33–40.
- MA Y., AUDIBERT L. & NAZARENKO A. (2009). Ontologies étendues pour l’annotation sémantique. In *20es Journées Francophones d’Ingénierie des Connaissances*, p. 205–216.
- PESKOV D., CLARKE N., KRONE J., FODOR B., ZHANG Y., YOUSSEF A. & DIAB M. (2019). Multi-domain goal-oriented dialogues (multidogo) : Strategies toward curating and annotating large scale dialogue data. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, p. 4526–4536.
- PORTET F., CAFFIAU S., RINGEVAL F., VACHER M., BONNEFOND N., ROSSATO S., LECOUTEUX B. & DESOT T. (2019). Context-aware voice-based interaction in smart home-vocadom@ a4h corpus collection and empirical assessment of its usefulness. In *2019 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCOM/CyberSciTech)*, p. 811–818 : IEEE.
- QIN L., XIE T., CHE W. & LIU T. (2021). A survey on spoken language understanding : Recent advances and new frontiers. *arXiv preprint arXiv :2103.03095*.
- SCHUSTER S., GUPTA S., SHAH R. & LEWIS M. (2019). Cross-lingual transfer learning for multilingual task oriented dialog. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics : Human Language Technologies, Volume 1 (Long and Short Papers)*, p. 3795–3805, Minneapolis, Minnesota : Association for Computational Linguistics. DOI : [10.18653/v1/N19-1380](https://doi.org/10.18653/v1/N19-1380).
- SHAH P., HAKKANI-TÜR D., TÜR G., RASTOGI A., BAPNA A., NAYAK N. & HECK L. (2018). Building a conversational agent overnight with dialogue self-play. *arXiv preprint arXiv :1801.04871*.
- TOMASELLO P., SHRIVASTAVA A., LAZAR D., HSU P.-C., LE D., SAGAR A., ELKAHKY A., COPET J., HSU W.-N., ADI Y. *et al.* (2023). Stop : A dataset for spoken task oriented semantic parsing. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, p. 991–998 : IEEE.
- TUR G. & DE MORI R. (2011). *Spoken language understanding : Systems for extracting semantic information from speech*. John Wiley & Sons.
- WELD H., HUANG X., LONG S., POON J. & HAN S. C. (2022). A survey of joint intent detection and slot filling models in natural language understanding. *ACM Computing Surveys*, **55**(8), 1–38.

XUE L., CONSTANT N., ROBERTS A., KALE M., AL-RFOU R., SIDDHANT A., BARUA A. & RAFFEL C. (2020). mt5 : A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv :2010.11934*.

YE F., MANOTUMRUKSA J. & YILMAZ E. (2021). Multiwoz 2.4 : A multi-domain task-oriented dialogue dataset with essential annotation corrections to improve state tracking evaluation.

ZANG X., RASTOGI A., SUNKARA S., GUPTA R., ZHANG J. & CHEN J. (2020). Multiwoz 2.2 : A dialogue dataset with additional annotation corrections and state tracking baselines. In *WNLPC AI, ACL'20*.

ZETTLEMOYER L. S. & COLLINS M. (2012). Learning to map sentences to logical form : Structured classification with probabilistic categorial grammars. *arXiv preprint arXiv :1207.1420*.

Projet Gender Equality Monitor (GEM)

Gilles Adda¹ François Buet¹ Sahar Ghannay¹ Cyril Grouin¹
Camille Guinaudeau¹ Lufei Liu¹ Aurélie Névéol¹ Albert Rilliard¹ Rémi Uro¹
(1) Université Paris-Saclay, CNRS, LISN, 507 rue du Belvédère, 91400 Orsay
prenom.nom@lisn.fr

RÉSUMÉ

Le projet ANR Gender Equality Monitor (GEM) est coordonné par l'Institut National de l'Audiovisuel (INA) et vise à étudier la place des femmes dans les médias (radio et télévision). Dans cette soumission, nous présentons le travail réalisé au LISN : (i) étude diachronique des caractéristiques acoustiques de la voix en fonction du genre et de l'âge, (ii) comparaison acoustique de la voix des femmes et hommes politiques montrant une incohérence entre performance vocale et commentaires sur la voix, (iii) réalisation d'un système automatique d'estimation de la féminité perçue à partir des caractéristiques vocales, (iv) comparaison de systèmes de segmentation thématique de transcriptions automatiques de données audiovisuelles, (v) mesure des biais sociétaux dans les modèles de langue dans un contexte multilingue et multi-culturel, et (vi) premiers essais d'identification de la publicité en fonction du genre du locuteur.

ABSTRACT

Gender Equality Monitor (GEM) Project

The ANR Gender Equality Monitor (GEM) project is coordinated by the Institut National de l'Audiovisuel (INA) and aims to study the place of women in the media (radio and television). In this submission, we present the work done at LISN : (i) diachronic study of acoustic voice characteristics as a function of gender and age, (ii) acoustic comparison of the voices of female and male politicians showing an inconsistency between vocal performance and comments on the voice, (iii) realisation of an automatic system for estimating perceived femininity from vocal characteristics, (iv) comparison of systems for thematic segmentation of automatic transcripts of audiovisual data, (v) measurement of societal biases in language models in a multilingual and multicultural context, and (vi) first attempts at speaker gender identification of advertising.

MOTS-CLÉS : Traitement de la parole, Modèles acoustiques, Modèles neuronaux.

KEYWORDS: Speech Processing, Acoustics Models, Neural Networks.

Remerciements

Ce travail est financé par l'ANR (financement ANR-19-CE38-0012).

CEN-CENELEC JTC 21 : La standardisation en TALN au service du règlement européen sur l'IA

Lauriane Aufrant
Inria, Rocquencourt, France
prenom.nom@inria.fr

RÉSUMÉ

Cette contribution présente les travaux du comité européen de standardisation de l'IA en matière de TALN. Le comité CEN-CENELEC JTC 21 a été mandaté par la Commission européenne pour développer les standards techniques permettant la mise en application du futur règlement européen sur l'IA : performance, robustesse, transparence, etc. Dans ce contexte, le TALN a été identifié comme un volet spécifique de l'IA, méritant ses propres outils, critères et bonnes pratiques. Ce constat a mené au développement d'une feuille de route ambitieuse incluant plusieurs projets de standardisation en TALN.

À ce jour, un premier travail d'inventaire et de définition des tâches de TALN a déjà été initié, et la rédaction d'un standard sur les métriques d'évaluation débute. Ces travaux ont aussi été l'occasion d'une réflexion plus large sur les besoins en standardisation du TALN, incluant une taxonomie des méthodes et des travaux sur les formats d'annotation et l'interopérabilité.

ABSTRACT

CEN-CENELEC JTC 21: NLP standardization in support of the AI Act

This contribution presents the NLP-related work of the European committee for AI standardization. The CEN-CENELEC JTC 21 committee has been mandated by the European Commission to develop the technical standards that will enable the application of the upcoming AI Act : accuracy, robustness, transparency, etc. In that context, NLP has been identified as a specific field of AI, that warrants its own tools, criteria and best practices. This observation has led to developing an ambitious roadmap with several NLP standardization projects.

To date, work for an inventory and definition of NLP tasks is already ongoing, and the writing of a standard on evaluation metrics is starting. Those efforts have also been an opportunity for a broader reflection on the standardization needs of NLP, including a taxonomy of methods, as well as considerations for annotation formats and interoperability.

MOTS-CLÉS : règlement européen sur l'IA ; standards ; métriques d'évaluation ; interopérabilité.

KEYWORDS: AI Act; standards; evaluation metrics; interoperability.

TEMITALC : Text Mining et TAL pour Analyser le Langage des Cachalots

José Coch¹ Olivier Adam²

(1) Service NLP, Dassault Systèmes, 10, place de la Madeleine, 75008 Paris, France

(2) Institut d'Alembert Sorbonne Université, 4 place Jussieu 75005 Paris France

jose.cochdiyacovo@3ds.com, olivier.adam@sorbonne-universite.fr

RESUME

Les cachalots sont des cétacés qui communiquent par des clics organisés en séquences appelées "codas". Elles sont numérisables relativement facilement : il existe en effet des corpus de transcriptions de conversations.

Une collaboration interdisciplinaire entre le Service NLP de Dassault Systèmes et l'équipe Bioacoustique de Sorbonne Université, a initié un projet d'application des techniques du TAL au cas des cachalots. Ses premiers résultats sont exposés dans les Actes de TextMine'23.

Le projet utilise le logiciel Proxem, qui permet de construire des modèles de langue à partir des corpus à analyser.

TEMITALC couvre les points suivants :

- Analyse des propriétés formelles du langage,
- Identification des corrélations entre des éléments non linguistiques et des éléments du langage.

Il bénéficie d'un financement de Dassault Systèmes et de Sorbonne Université. Sa fin est prévue pour décembre 2024.

Nos résultats vont contribuer à décrire le sophistiqué langage d'une espèce non-humaine.

ABSTRACT

TEMITALC: Text Mining and NLP to Analyze the Language of Sperm Whales

Sperm whales are cetaceans that communicate through clicks organized in sequences called "codas". They can be digitized relatively easily: there are indeed corpora of transcriptions of conversations.

An interdisciplinary collaboration between the NLP Service of Dassault Systèmes and the Bioacoustics team of Sorbonne University, initiated a project to apply NLP techniques to the case of sperm whales. Its first results are set out in the Proceedings of TextMine'23.

The project uses the Proxem software, which makes it possible to build language models from the corpora to be analyzed.

TEMITALC covers the following points:

- Analysis of the formal properties of this language,
- Identification of correlations between non-linguistic elements and language elements.

It receives funding from Dassault Systèmes and Sorbonne University. Its end is scheduled for December 2024.

Our results will help to describe the sophisticated language of a non-human species.

MOTS-CLES : zoolinguistique ; text-mining ; interdisciplinaire ; ordre des mots ; corrélations sémantiques.

KEYWORDS: zoolinguistics; text-mining; interdisciplinary; word order; semantic correlations

1. TEMITALC : Text Mining et TAL pour Analyser le Langage des Cachalots

Les cachalots (*Physeter macrocephalus*) sont les plus grands des cétacés à dents. Comme tous les cétacés, ils communiquent notamment par des émissions vocales. Les cachalots produisent des clics au cours de leurs activités vitales et leurs interactions sociales. Certains de ces sons sont organisés en séquences temporelles, appelées « codas ». Depuis plus d'une dizaine d'années, des échanges audio ou « conversations » entre cachalots sont enregistrés dans de nombreux endroits dans le monde, par exemple dans l'Océan Pacifique, dans les Caraïbes et dans l'Océan Indien. La particularité des échanges vocaux entre cachalots fait que ces codas sont numérisables relativement facilement. Ainsi, il existe des corpus de transcriptions de conversations en particulier venant des origines géographiques citées.

Durant 2022, une collaboration entre le Service NLP de Dassault Systèmes et l'équipe Bioacoustique de Sorbonne Université, basée sur les enregistrements sonores collectés et mis à disposition par Longitude 181 et Label Bleu Production, nous a permis d'initier un projet d'application des techniques de Text Mining et Traitement Automatique du Langage à l'étude du langage des cachalots. Nous avons exposé les premiers résultats du projet dans un article publié dans les Actes de l'atelier TextMine'23 de la conférence EGC'2023 concernant un corpus de cachalots résidents au large de l'île Maurice et identifiés individuellement.

Nous utilisons dans ce projet le logiciel Proxem Studio, qui a la particularité de pouvoir être appliqué sans modèle de langue préalable car il peut construire des modèles de langue à partir des corpus à analyser.

L'objectif du projet couvre les points suivants :

- Optimiser et automatiser la transcription en codas des échanges audio entre cachalots,
- Analyser les propriétés formelles du langage des cachalots : mettre en évidence que l'ordre entre codas a une importance, et découvrir s'il est possible de décrire une proto-syntaxe de ce langage,

- Mettre au point un référentiel d'éléments non linguistiques (comportements sociaux, données démographiques, relations familiales) et identifier des codas ou des séquences de codas montrant une corrélation avec ces éléments non linguistiques, et in fine, avancer des hypothèses sur la fonction de certaines codas ou séquences de codas,
- Etudier les corrélations entre les participants à chaque conversation et les codas émis afin de déterminer si des codas ou séquences de codas peuvent être associées à des individus.

Le projet bénéficie d'un financement de Dassault Systèmes et de Sorbonne Université. La fin du projet est prévue pour décembre 2024.

Nos résultats vont contribuer ainsi à décrire le sophistiqué langage d'une espèce non-humaine.

Références

ADAM, O., A. YERNAUX, M. SAUVÊTRE, J. NGOSSO, G. NUEL, M. HAFFNER-TRINH, R. TROUSSIER, Z.-L. GUILLERM, L. PICON, L. BARLUET, J. MACKY, L. BARLUET DE BEAUCHESNE, V. KUHN, F. DELFOUR, V. SARANO, H. VITRY, A. PREUD'HOMME, R. HEUZEY, J.-L. JUNG, ET F. SARANO (2020). Study of behaviours and emitted codas during sperm whale social interactions. *e-Forum Acusticum 2020*, Dec 2020, Lyon, France. pp.3225-3227.

ANDREAS, J., G. BEGUS, M. BRONSTEIN, R. DIAMANT, D. DELANEY, S. GERO, S. GOLDWASSER, D. GRUBER, S. HAAS, P. MALKIN, R. PAYNE, G. PETRI, D. RUS, P. SHARMA, P. TØNNESEN, A. TORRALBA, D. VOGT, ET R. WOOD (2021). Cetacean translation initiative: a roadmap to deciphering the communication of sperm whales. arXiv (2104.08614)

BERMANT, P., M. BRONSTEIN, R. WOOD, S. GERO, ET D. GRUBER (2019). Deep machine learning techniques for the detection and classification of sperm whale bioacoustics. *Scientific Reports* 9, 1–10.

BOSSHARD, A., M. LEROUX, N. LESTER, B. BICKEL, S. STOLL, ET S. TOWNSEND (2022). From collocations to call-ocations: using linguistic methods to quantify animal call combinations. *BEHAVIORAL ECOLOGY AND SOCIOBIOLOGY* ; Heidelberg.

COCH, J., BERKENBAUM, L., ADAM, O. (2023). Une contribution du Text-mining à la connaissance du langage des cachalots. In *Actes de TextMine, atelier de la conférence EGC (Extraction et Gestion des Connaissances)*, pp. 15-32. Lyon, France, janvier 2023.
<https://textmine.sciencesconf.org/data/pages/TextMine23.pdf>

DOH, YANN & ECALLE, BEVERLEY & DELFOUR, FABIENNE & PANKOWSKI, CYPRIEN & COZANET, GILDAS & BECOUARN, GUILLAUME & OVIZE, MARION & DENIS, BERTRAND & ADAM, O.. (2023). Performance Assessment of the Innovative Autonomous Tool CETOSCOPE© Used in the Detection and Localization of Moving Underwater Sound Sources. *Journal of Marine Science and Engineering*. 11. 960. 10.3390/jmse11050960.

μDialBot: Multi-party perceptually-active situated DIALog for human-roBOT interaction

Fabrice Lefèvre, Timothée Dhaussy, Ahmed Ndouop Njifenjou,
Virgile Sucal, Bassam Jabaian
LIA, Avignon Université, France
{fabrice.lefevre}@univ-avignon.fr

RÉSUMÉ

Dans *muDialBot* (ANR-20-CE33-0008-01), notre ambition est d'incorporer pro-activement des traits de comportements humains dans la communication parlée. Nous projetons d'atteindre une nouvelle étape de l'exploitation de l'information riche fournie par les flux de données audio et visuelles venant des humains. En particulier en extraire des événements verbaux et non-verbaux devra permettre d'accroître les capacités de décision des robots afin de gérer les tours de parole plus naturellement et aussi de pouvoir basculer d'interactions de groupe à des dialogues en face-à-face selon la situation. Récemment on a vu croître l'intérêt pour les robots compagnons capable d'assister les individus dans leur vie quotidienne et de communiquer efficacement avec eux. Ces robots sont perçus comme des entités sociales et leur pertinence pour la santé et le bien-être psychologique a été mise en avant dans des études. Les patients, leurs familles et les professionnels de santé pourront mieux apprécier le potentiel de ces robots, dans la mesure où certaines limites seront rapidement affranchies, telles leur capacité de mouvement, vision et écoute afin de communiquer naturellement avec les humains, au-delà de ce que permettent déjà les écrans tactiles et les commandes vocales. Les résultats scientifiques et technologiques du projet seront implémentés sur un robot social commercial et seront testés et validés avec plusieurs cas d'usage dans le contexte d'une unité d'hôpital de jour.

ABSTRACT

In *muDialBot* (ANR-20-CE33-0008-01) our ambition is to actively incorporate human-behavior cues in spoken human-robot communication. We intend to reach a new level in the exploitation of the rich information available with audio and visual data flowing from humans when interacting with robots. In particular, extracting highly informative verbal and non-verbal perceptual features will enhance the robot's decision-making ability such that it can take speech turns more naturally and switch between multi-party/group interactions and face-to-face dialogues where required. Recently there has been an increasing interest in companion robots that are able to assist people in their everyday life and to communicate with them. These robots are perceived as social entities and their utility for healthcare and psychological well being for the elderly has been acknowledged by several recent studies. Patients, their families and medical professionals appreciate the potential of robots, provided that several technological barriers would be overcome in the near future, most notably the ability to move, see and hear in order to naturally communicate with people, well beyond touch screens and voice commands. The scientific and technological results of the project will be implemented onto a commercially available social robot and they will be tested and validated with several use cases in a day-care hospital unit.

MOTS-CLÉS : projet ANR, interaction humain-robot, pro-activité, multi-partie.

KEYWORDS: ANR project, human-robot interaction, pro-activity, multi-party.

1 Motivations

Companion robots able to assist people in their everyday life and to communicate with them :

social assistive robots

How to improve acceptance ?

Several technological barriers to overcome :

- better ability to move, see and hear in order to naturally communicate with people in various configurations (face-to-face, multi-party, close, distant etc)
- pro-active use of perceptions
- audio-visual strategies to handle interactions

2 Ambitions

New level in exploitation of rich information available with audio and visual data flowing from humans when interacting with robots :

- fusing highly informative **verbal and non-verbal perceptive features** to enhance the robot's decision-making ability such that it can
- switch between **multi-party/group interactions and face-to-face dialogues** where required, and
- take speech turns **more naturally**



Tested and validated with several **use cases in a day-care hospital** unit. Large-scale data collection, complement in-situ tests, to fuel further researches.

Project mainly **funded by ANR and labeled by the Pole SCS**. More information on : <https://www.pole-scs.org/projets/mudialbot/>

3 Partners

 AVIGNON UNIVERSITÉ	LIA, COORDINATOR • Human-machine vocal interactions, decision-making learning
 LABORATOIRE HUBERT CURIE	Lab Hubert Curien • Image analysis
 Inria	INRIA Perception • Audio-visual scene analysis
 ERM	ERM • Robotic engineering, software integration, data management
 ASSISTANCE PUBLIQUE HÔPITAUX DE PARIS	AP-HP/Hopital Broca • Healthcare application, psychological survey

4 Organization

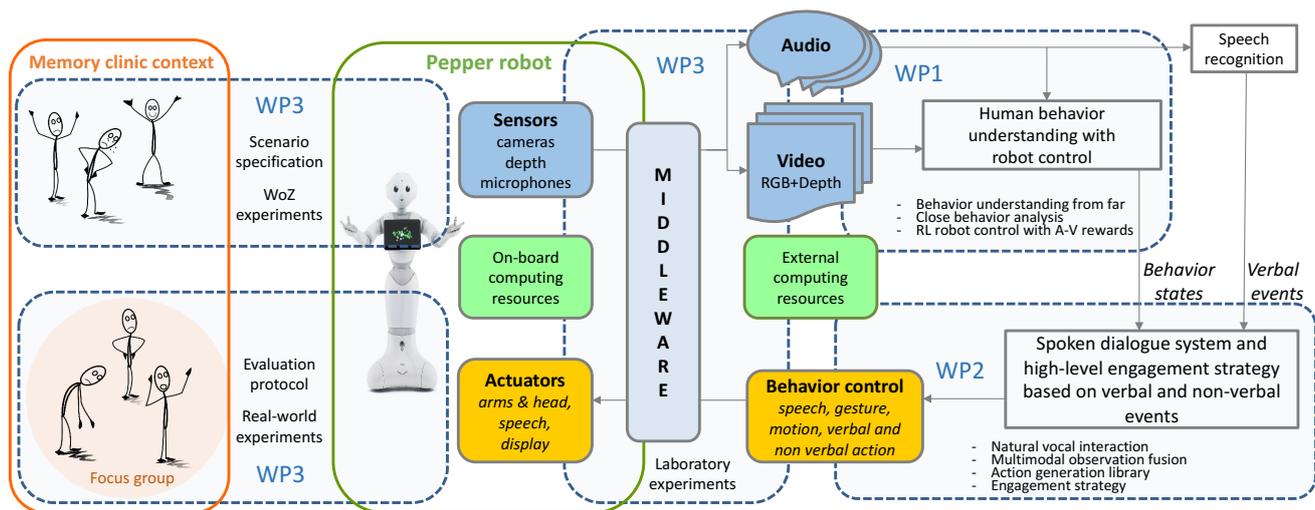
Methodology followed in μ DialBot at the same time mainstream, in line with the classical methodology to handle **new machine learning-based approaches**, and pioneering, with **novel online learning techniques** to reduce the requirement for initial data collection to its minimum. Relying on deep learning techniques which have shown their efficiency in other domains, a new formalism dedicated to **pro-active perceptually-based control of a conversational robot** will be proposed.



General methodology of the project consists of two operational building blocks for :

- the estimation of non-verbal states and
- the learning of an event-guided strategy.

Integrated on the robotic platform through a software abstraction layer. A first round of WoZ experiments to test the proposed models, then the complete system gradually introduced in the true clinical context, using a well-defined protocol.



3 main work-packages implemented :

— **WP1 Human behavior understanding with robot control**

Objectives : to develop methods and algorithms to extract HBU cues from audio and visual data.

— **WP2 Spoken dialogue system and high-level engagement strategy based on verbal and non-verbal events**

Objectives : to develop the natural vocal interaction ability of the robot.

— **WP3 Specifications, integration on robotic platform, iterative and final evaluation of the human-robot interactions in a memory clinic**

Objectives : to define the experimental protocol, specify laboratory and real-world experimental protocols, and conduct “Wizard of Oz” (WoZ) experiments.

Projet MALIN : MANuels scoLaires INclusifs

Olivier Pons^{1,*} Isabelle Barbet^{1,*} Jérôme Dupire^{1,*} Valerie Grembi^{7,*}
Camille Guinaudeau^{2,3,**} Céline Hudelot^{4,†} Caroline Huron^{5,6,‡}
Vincent Mousseau^{4,†} Élise Lincker^{1,*} Léa Pacini^{1,5,*}

(1) Cedric, CNAM, Paris, France

(2) Japanese French Laboratory for Informatics, CNRS, NII, Tokyo, Japon

(3) Université Paris-Saclay, Orsay, France

(4) MICS, CentraleSupélec, Université Paris-Saclay, Orsay, France

(5) SEED, Inserm, Université Paris Cité, Paris, France

(6) Learning Planet Institute, Paris, France

(7) Le cartable fantastique Paris, France

*prenom.nom@lecnam.net, **nom@nii.ac.jp

†prenom.nom@centralesupelec.fr, ‡prenom.nom@cri-paris.org,

*valerie.grembi@cartablefantastique.fr

RÉSUMÉ

L'école joue un rôle essentiel dans la vie des enfants. La restriction de la participation à l'école en raison d'un handicap réduit la qualité de vie. Une difficulté est l'inaccessibilité des manuels scolaires systématiquement utilisés en France pour accompagner les apprentissages. Notre projet vise à les rendre accessibles aux élèves en situation de handicap en innovant pour automatiser leur adaptation. Il s'appuie sur le croisement d'expertises médicale, pédagogique et de psychologie cognitive d'une part, d'expertises en interactions/interfaces homme-machine, accessibilité numérique, traitement de la langue et en conception de systèmes intelligents, d'autre part. Il s'agira de concevoir une plate-forme qui, en partant d'un manuel au format PDF (ou EPUB), mettra en oeuvre, via des modèles structurels et sémantiques du manuel, les adaptations et interfaces qui sont aujourd'hui principalement faites manuellement par les organismes de transposition. Ce travail est financé par l'ANR (financement ANR-21-CE38-0014).

ABSTRACT

Inclusive textbooks

School plays an essential role in children's lives. Restricting participation in school due to disability reduces quality of life. One difficulty is the inaccessibility of school textbooks, which are systematically used in France to support learning. Our project aims to make them accessible to students with disabilities, by innovating to automate their adaptation. It is based on a combination of medical, pedagogical and cognitive psychology expertise on the one hand, and expertise in human-computer interaction/interfaces, digital accessibility, language processing and intelligent systems design on the other. The aim is to design a platform which, starting from a manual in PDF (or EPUB) format, will implement, via structural and semantic models of the manual, the adaptations and interfaces which are today mainly done manually by transposition organizations. This work is financed by the ANR (funding ANR-21-CE38-0014).

MOTS-CLÉS : adaptation de manuels scolaires, accessibilité.

KEYWORDS: textbook adaptation, accessibility.

PROPICTO : Développer des systèmes de traduction de la parole vers des séquences de pictogrammes pour améliorer l’accessibilité de la communication

Lucía Ormaechea^{1,2}, Pierrette Bouillon², Maximin Coavoux¹, Emmanuelle Esperança-Rodier¹, Johanna Gerlach², Jérôme Goulian¹, Benjamin Lecouteux¹, Cécile Macaire¹, Jonathan Mutal², Magali Norré^{2,3}, Adrien Pupier¹, Didier Schwab¹ et Hervé Spechbach⁴

(1) Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP*, LIG, 38000 Grenoble, France

(2) FTI, Université de Genève, Genève, Suisse

(3) CENTAL, ILC, Université Catholique de Louvain, Louvain-la-Neuve, Belgique

(4) Hôpitaux Universitaire de Genève, Genève, Suisse

RÉSUMÉ

PROPICTO est un projet financé par l’Agence nationale de la recherche française et le Fonds national suisse de la recherche scientifique, qui vise à créer des systèmes de traduction de la parole vers des pictogrammes avec le français comme langue d’entrée. En développant de telles technologies, nous avons l’intention d’améliorer l’accès à la communication pour les patients non francophones et les personnes souffrant de troubles cognitifs.

ABSTRACT

PROPICTO : Developing Speech-to-Pictograph Translation Systems to Enhance Communication Accessibility

PROPICTO is a project funded by the French National Research Agency and the Swiss National Science Foundation, that aims at creating Speech-to-Pictograph translation systems, with a special focus on French as an input language. By developing such technologies, we intend to enhance communication access for non-French speaking patients and people with cognitive impairments.

MOTS-CLÉS : Pictogrammes, Corpus, Communication Alternative et Augmentée.

KEYWORDS: Pictographs, Corpora, Augmentative and Alternative Communication.

1 Introduction

Les dispositifs de communication alternative et augmentée (CAA) jouent un rôle de plus en plus important auprès des personnes en situation de handicap et de leurs proches. Cependant, l’utilisation de ces technologies (*i.e.*, tableaux de communication ou médias électroniques) peut être fastidieuse (Vaschalde *et al.*, 2018). Pour surmonter ce problème, nous partons de l’hypothèse que les systèmes de traduction de la parole vers des pictogrammes (S2P) peuvent être utiles aux utilisateurs de CAA. En outre, nous pensons qu’ils peuvent améliorer l’accessibilité des services de santé pour les patients qui ne parlent pas la langue locale. Le développement de tels outils nécessite des recherches approfondies

dans plusieurs domaines du traitement de la langue et de la parole (TALP). Nous présentons ici un projet de recherche visant à créer des systèmes qui traduisent automatiquement le français parlé en pictogrammes.

Lancé début 2021, PROPICTO¹ (l’acronyme signifie *PRO*jection de la parole vers des *PIC*Togrammes) est un projet franco-suisse de quatre ans, financé à la fois par l’Agence nationale de la recherche française² et par le Fonds national suisse de la recherche scientifique.³ Il est mené en collaboration entre le Département de technologie de la traduction de l’Université de Genève et le Groupe d’étude pour la traduction automatique et le traitement automatisé des langues et de la parole, une équipe du Laboratoire d’informatique de Grenoble.

Dans le cadre de ce projet, nous examinerons plusieurs domaines liés au TALP, à savoir la reconnaissance vocale, l’analyse syntaxique, la simplification des phrases et la génération de pictogrammes. En intégrant cette série de tâches dans différents flux de travail (en fonction du scénario cible), nous proposons de nouveaux systèmes de traduction automatique multimodaux qui convertissent le langage parlé en unités pictographiques. En utilisant cette approche, nous avons l’intention de répondre aux besoins sociétaux et de communication dans les milieux suivants : (1) *handicap*, où un individu cherche à communiquer avec une personne souffrant d’un trouble cognitif, et (2) *médicaux*, où une barrière linguistique existe entre le patient et le praticien.

2 Vue d’ensemble

PROPICTO vise à améliorer la facilité d’utilisation des dispositifs de CAA en tirant parti de solutions basées sur le TALP pour une plus grande accessibilité. Nous concevons de nouvelles méthodes et de nouveaux corpus afin de permettre la transcription directe d’énoncés parlés en séquences de pictogrammes, avec un *objectif général* comme ARASAAC,⁴ et un *objectif spécifique* (*i.e.*, SantéBD <https://santebd.org/> pour les concepts liés à la santé). Le projet devra relever deux défis majeurs :

- La *carence de corpus parallèles de pictogrammes vocaux*, qui constitue un obstacle important à la mise en œuvre de l’apprentissage automatique de pointe (en particulier de bout en bout) ;
- La *nécessité d’une évaluation humaine et automatique approfondie* afin d’évaluer la compréhensibilité des séquences de sortie auprès de divers groupes cibles.

Pour mieux les aborder, nous adopterons une approche en cascade pour notre flux de travail de traitement S2P, qui sera adapté en fonction de l’objectif visé. Ainsi, une première approche privilégiera une stratégie *basée sur les concepts* pour traiter la génération de pictogrammes et sera intégrée dans une architecture S2P à vocation médicale,⁵ composée d’un système de reconnaissance automatique de la parole (ASR) et d’un module neuronal de conversion texte-UMLS⁶ qui définira les pictogrammes à produire et la syntaxe. Une autre stratégie de génération de pictogrammes reposera sur une approche basée sur les mots et sera précédée des étapes suivantes (comme le montre la figure 1) : ASR, analyse des dépendances et simplification des phrases.

1. <https://www.propicto.unige.ch/>

2. <https://anr.fr/en>

3. <https://www.snf.ch/en>

4. <https://arasaac.org/>

5. Pour plus de détails sur cette architecture, se référer à [Mutal et al. 2022](#)

6. Cet acronyme fait référence aux concepts du système de langage médical unifié (UMLS).

L'utilisation d'une approche en cascade est motivée par l'avantage attendu d'une phase par rapport à la suivante. En outre, elle permet d'assurer une meilleure explicabilité du modèle. Notre deuxième proposition d'architecture multimodale partira d'un module ASR, s'appuyant sur des modèles Wav2Vec2.0 à la pointe de la technologie. La tâche d'analyse en dépendance sera traitée avec un analyseur syntaxique de bout en bout dont l'entrée est le signal brut d'un énoncé donné.

L'utilisation du signal brut au lieu des transcriptions nous permet d'utiliser l'information prosodique pour mieux prédire les frontières syntaxiques (Pupier *et al.*, 2022).

L'extraction d'une représentation syntaxique peut à son tour fournir des informations clés pour une simplification plus efficace au niveau de la phrase. La réduction de la complexité linguistique de la transcription d'entrée devrait faciliter l'étape suivante, au cours de laquelle la traduction en pictogrammes sera également régie par des règles de grammaire expertes⁷.

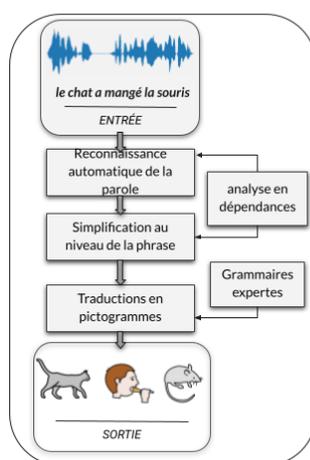


FIGURE 1 – Vue d'ensemble de l'architecture en cascade Parole vers pictogrammes utilisant une approche basée sur *les mots*.

3 Contributions

PROPICTO mettra à la disposition de la communauté scientifique des méthodes et des ressources permettant une traduction du français parlé vers des pictogrammes. Les licences seront aussi permissives que possible et conformes à celles des ensembles pictographiques utilisés. En outre, plusieurs prototypes destinés à différents publics cibles seront mis en production à la fin du projet : (1) dans les situations d'urgence aux Hôpitaux Universitaires de Genève, une démonstration est disponible sur : <https://propicto.demos.unige.ch/pictoDrClient/translate/> et (2) dans des institutions pour enfants et adultes souffrant de handicaps multiples. Ils seront testés en conditions réelles et évalués à l'aide de méthodes humaines et automatiques.

Références

MUTAL J., BOUILLON P., NORRÉ M., GERLACH J. & ORMAECHEA GRIJALBA L. (2022).

7. expressions multi-mot, temps des verbes, noms propres

A Neural Machine Translation Approach to Translate Text to Pictographs in a Medical Speech Translation System – The BabelDr Use Case. In *Proc. Association for Machine Translation in the Americas*, p. 252–263.

PUPIER A., COAVOUX M., LECOUTEUX B. & GOULIAN J. (2022). End-to-End Dependency Parsing of Spoken French. In *Proc. Interspeech 2022*, p. 1816–1820. DOI : [10.21437/Interspeech.2022-381](https://doi.org/10.21437/Interspeech.2022-381).

VASCHALDE C., TRIAL P., ESPERANÇA-RODIER E., SCHWAB D. & LECOUTEUX B. (2018). Automatic Pictogram Generation from Speech to Help the Implementation of a Mediated Communication. In *Proc. Swiss Centre for Barrier-Free Communication 2018*, p. 97–101.

Recherche d'information conversationnelle

Laure Soulier^{1,2}, Pierre Erbacher¹, Thomas Gerald², Hanane Djeddal¹, Jian-Yun Nie³, Preux Philippe⁴

(1) Sorbonne Université, CNRS, ISIR, F-75005 Paris, France.

(2) Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique, 91400, Orsay, France.

(3) University of Montreal, Montreal, Canada

(4) Univ. Lille, CNRS, Inria, Centrale Lille, UMR 9189 – CRIStAL, Lille, France
first.last @ {sorbonne-universite.fr, lisn.fr, inria.fr},
nie@iro.umontreal.ca

RÉSUMÉ

Le projet ANR JCJC SESAMS s'intéresse depuis 2018 au paradigme désormais actuels des systèmes de recherche d'information conversationnels. L'objectif est de formaliser des modèles de recherche d'information capables de fluidifier les interactions avec les utilisateurs pendant une session de recherche. Nous abordons différents enjeux : la prise en compte d'une conversation en langage naturel en contexte d'une recherche d'information, la génération d'interactions permettant de clarifier les besoins en information, la génération de réponse en langage naturel, ainsi que l'apprentissage continu pour s'adapter aux nouveaux besoins des utilisateurs. Nous présenterons dans ce poster ces différents enjeux et les contributions associées. Nous pourrions également discuter les perspectives de recherche dans ce domaine suite au développement récents des gros modèles de langue.

ABSTRACT

Search-oriented conversational systems

Since 2018, the ANR JCJC SESAMS project has focused on the now trendy paradigm of conversational information retrieval systems. The objective is to formalize information retrieval models able of smoothing interactions with users during a search session. We address different issues : taking into account a natural language conversation in the context of a search session, generating interactions to clarify information needs, generating natural language responses, and continuously learning to adapt to new user needs. In this poster, we will present these different issues and the associated contributions. We will also discuss the research perspectives in this area following the recent development of large language models.

MOTS-CLÉS : Recherche d'information, système conversationnel, modèles de langue, clarification des requêtes, apprentissage continu.

KEYWORDS: Information retrieval, conversational system, language models, query clarification, continual learning.

Remerciements. Ce travail est financé par l'ANR JCJC SESAMS (ANR-18- CE23-0001).

Autogramm : développement simultané de treebanks et de grammaires à partir de corpus

Sylvain Kahane¹ Santiago Herrera¹ Bruno Guillaume² Kim Gerdes³

(1) Modyco, Université Paris Nanterre, CNRS

(2) Sémagramme, Loria, Inria Nancy - Grand Est, Université de Lorraine

(3) LISN, Université Paris-Saclay, CNRS

sylvain@kahane.fr, s.herrera@parisnanterre.fr,
bruno.guillaume@inria.fr, kim@gerdes

RÉSUMÉ

Ce projet de recherche vise à créer de nouveaux treebanks en dépendance pour des langues sous-dotées, en unifiant autant que possible leur développement avec celui de grammaires descriptives quantitatives. Nous présenterons notre chaîne de traitement et de développement de treebanks et nous discuterons du type de grammaire que nous voulons extraire. Enfin, nous examinerons l'utilisation de ces ressources en typologie quantitative.

ABSTRACT

Autogramm : Simultaneous development of treebanks and corpus-driven grammars

This research project aims to create new dependency treebanks for low-resource languages, unifying as far as possible their development with that of quantitative descriptive grammars. We will present our processing pipeline and discuss the type of grammar we want to extract. Finally, we will examine the use of these resources in quantitative typology.

MOTS-CLÉS : Autogramm, treebanks, extraction de grammaires, typologie quantitative.

KEYWORDS: Autogramm, treebanks, grammar extraction, quantitative typology.

1 Introduction

Les études linguistiques comparatives nécessitent des corpus de haute qualité qui représentent au mieux la variation et la diversité des langues du monde, avec des annotations suffisamment riches pour en extraire des descriptions grammaticales, et suffisamment comparables pour permettre des études contrastives et typologiques.

Le projet de recherche *Autogramm*¹ ici présenté vise à répondre à ces besoins, du moins en partie, en créant de nouveaux treebanks syntaxiques pour plus de 15 langues sous-dotées et en unifiant autant que possible leur développement avec celui de grammaires descriptives quantitatives. La production de corpus annotés et de grammaires descriptives est certainement complémentaire et leur développement simultané pourrait permettre de réduire le temps de travail et d'améliorer la qualité des

1. Autogramm est un projet financé par l'Agence Nationale de la Recherche (ANR-21-CE38-0017). Pour accéder à la liste complète des participants, y compris les linguistes et les langues sur lesquelles ils travaillent, ainsi que les outils développés, voir <https://autogramm.github.io/>.

deux ressources. De plus, les grammaires basées sur des corpus encodent facilement des informations quantitatives, permettant, par exemple, la hiérarchisation des observations grammaticales et leur comparaison.

Plus précisément, nous présenterons notre chaîne de développement de treebanks en dépendances, en utilisant les schémas d’annotation *Universal Dependency* (UD) (Nivre *et al.*, 2020; de Marneffe *et al.*, 2021) et *Surface Syntactic* UD (SUD) (Gerdes *et al.*, 2018, 2019a). Ensuite, nous discuterons du type de grammaire que nous voulons extraire à partir de treebanks et nous présenterons les prochaines étapes de notre travail qui tend vers les études comparatives et vers une typologie quantitative.

2 Chaîne de traitement et nouvelles ressources

Le projet rassemble une équipe diversifiée, comprenant des linguistes de terrain spécialisés en langues peu décrites et des experts en annotation de corpus. Une chaîne de traitement a été mise en place pour que chacun puisse contribuer au développement de ces ressources pour chacune des langues étudiées.

En général, le processus commence par la transformation des gloses interlinéaires (IGT), souvent utilisées par les linguistes de terrain, en un pre-treebank, sans perdre les informations qu’elles contiennent (la segmentation, les traits morpho-syntaxiques, les gloses, etc.). Cela implique un travail avec le linguiste qui consiste en sélectionner et normaliser les informations pertinentes. L’annotation syntaxique peut alors se faire au niveau des mots ou des morphes (e.g. Kahane *et al.*, 2021). Nous utilisons l’outil d’annotation en ligne `ArboratorGrew`, qui propose un système de bootstrapping syntaxique (Guibon *et al.*, 2020; Peng *et al.*, 2022) : l’analyseur syntaxique peut être entraîné avec le travail déjà effectué pour annoter automatiquement le reste du corpus, autant de fois que nécessaire. En parallèle, nous construisons des grammaires pour chacune des langues (voir section 3).

Plusieurs treebanks sont actuellement en cours de développement pour les langues suivantes : Amdo Tibetan (sino-tibétain), dialectes arabes (marocain, égyptien, tunisien ; sémitique), bambara (mandingue), breton (celte), gbara (oubanguien), haïtien (créole), hausa (chadique), salar (turc), sungwadia (austroasiatique), tuwari (papoue), vietnamien (austroasiatique), yali (papoue), ye’kwana (caribe), etc. Un treebank pour le Beja (Kahane *et al.*, 2021) et un pour le Zaar (Caron, 2015) ont déjà été développées en utilisant une approche similaire et publiées dans la base de données d’UD.

3 Extraction de grammaires à partir de corpus

Il existe un grand nombre des travaux utilisant différents formalismes et différentes stratégies pour extraire de la manière la plus automatique possible les grammaires et les propriétés typologiques des corpus annotés. La plupart des méthodes produisent des grammaires formelles à partir de ressources linguistiques, telles que les IGT, en utilisant des connaissances grammaticales externes et élaborées à la main (voir Bender *et al.*, 2002; Zamaraeva *et al.*, 2022; Howell & Bender, 2022). Ces grammaires ne contiennent généralement pas d’informations quantitatives, bien que le fait de disposer de telles données permet d’obtenir une description fine de la langue étudiée et de classer les descriptions extraites en fonction de leur importance au sein d’un corpus. D’autres systèmes d’extraction parviennent à encoder des informations quantitatives (e.g. Blache *et al.*, 2016), mais le nombre de règles extraites restent élevé et la forme des règles est limitée. En outre, certaines propriétés

ne sont encodées qu’au niveau de leurs constructions. Par exemple, ces grammaires indiqueront si chaque construction a une tête en position finale, mais pas si la langue étudiée est une langue à tête finale.

Notre objectif est d’extraire des descriptions grammaticales facilement interprétables par n’importe quel utilisateur et, puisque la tâche est réalisée à partir de données annotées, nous cherchons à associer à chacune d’entre elles des informations quantitatives. Cependant, contrairement aux grammaires analysées, nous visons à classer les observations grammaticales en fonction de leur importance dans un corpus et à obtenir des grammaires de différentes tailles en fonction de la manière dont les règles extraites sont classées et regroupées. Les données quantitatives peuvent également être utilisées pour mettre en évidence les relations qui existent entre les différentes propriétés afin d’expliquer certains phénomènes et d’en découvrir d’autres qui seraient autrement passés inaperçus (voir [Bresnan et al. \(2007\)](#) dans une étude classique sur l’alternance dative ; plus récemment, [Chaudhary et al. \(2020\)](#) et [Chaudhary \(2022\)](#) ont extrait des règles d’accord, d’ordre, ainsi que des cas à l’aide de treebanks). Dans ce but, nous nous concentrons sur la fréquence des phénomènes observés et sur d’autres mesures continues ([Levshina, 2019](#) et [Gerdes et al., 2019b](#)).

Nous travaillons actuellement sur différents systèmes d’extraction de grammaires. Dans le cadre du projet, nous avons déjà développé une première méthode qui permet d’extraire avec succès des motifs grammaticaux à partir de treebanks et de les classer en fonction de leur importance statistique dans le corpus ([Herrera et al., 2022](#)). Plus précisément, nous calculons la probabilité d’obtenir la distribution observée de certains motifs apparentés à partir d’une hypothèse d’indépendance. Plus cette probabilité est élevée, plus le motif est significatif.

Enfin, en interagissant avec les descriptions et les grammaires extraites, le linguiste travaillant sur une langue peu décrites, sera en mesure de vérifier si les caractéristiques (spécifiques de la langue) choisies pour annoter le corpus représentent bien la grammaire de la langue en question.

4 Typologie quantitative

Les descriptions grammaticales, du moins en ce qui concerne les propriétés universelles, sont exprimées à l’aide du même jeu d’étiquettes et du même formalisme en dépendance. Cela signifie que les grammaires que nous extrayons sont des grammaires comparables qui permettent des analyses comparatives d’une même observation entre différentes langues et différents corpus. De cette façon, nous pouvons déterminer précisément ce qui est particulier à une langue par rapport à d’autres, sans limiter l’étude typologique à une liste préétablie d’observations, qui peuvent pour certaines (familles de) langues se révéler impertinentes.

Lorsque l’on travaille avec des informations quantitatives, on peut également comparer des observations en termes de valeurs continues plutôt que de valeurs discrètes. Contrairement aux bases de données importantes et fondamentales (WALS ([Dryer & Haspelmath, 2013](#)), APiCS ([Michaelis et al., 2013](#)), ValPal ([Hartmann et al., 2013](#)), entre autres), il sera possible d’expliciter dans quelle mesure une caractéristique typologique est présente dans un corpus spécifique et dans quelle mesure elle diffère de celle d’autres langues. Ce faisant, nous travaillons dans le cadre de la typologie quantitative (cf. [Cysouw, 2005](#)), en suivant une nouvelle perspective d’étude encore à explorer ([Futrell et al., 2015](#); [Gerdes et al., 2021](#)).

Nous explorons des méthodes d’échantillonnage et de comparaison pour trouver des similitudes et

des différences entre les corpus en utilisant les observations extraites, tout en cherchant de nouvelles métriques autres que la fréquence. Nous construisons une base de données typologique contenant les observations quantitatives collectées. Il est à noter que ces méthodes pourraient également être utilisées pour détecter d'autres variations dans la langue, telles que les variations sociolinguistiques et diachroniques.

5 Obstacles et perspectives

Les différents objectifs du projet se heurtent évidemment à plusieurs obstacles, dont les suivants : l'explosion combinatoire des variables lorsque l'on tente d'extraire des modèles grammaticaux des banques d'arbres, des résultats incongrus dus à différentes interprétations du schéma d'annotation et aux particularités du corpus, et des échantillons linguistiques déséquilibrés qui empêchent les études typologiques cohérentes. Une partie du projet consiste à étudier ces problèmes afin de contribuer à la linguistique théorique et à la documentation linguistique.

Références

- BENDER E. M., FLICKINGER D. & OEPEN S. (2002). The grammar matrix : An open-source starter-kit for the rapid development of cross-linguistically consistent broad-coverage precision grammars. In *COLING-02 : Grammar Engineering and Evaluation*.
- BLACHE P., RAUZY S. & MONTCHEUIL G. (2016). MarsaGram : an excursion in the forests of parsing trees. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, p. 2336–2342, Portorož, Slovenia : European Language Resources Association (ELRA).
- BRESNAN J., CUENI A., NIKITINA T. & BAAYEN R. (2007). *Predicting the Dative Alternation*, In *Cognitive foundations of interpretation*, p. 69–94. KNAW : Amsterdam.
- CARON B. (2015). Zaar grammatical sketch. In ASDT, Éd., *Corpus-based Studies of Lesser-described Languages. The CorpAfroAs corpus of spoken AfroAsiatic languages*. Amsterdam-Philadelphia : John Benjamins.
- CHAUDHARY A. (2022). *Automatic Extraction and Application of Language Descriptions for Under-Resourced Languages*. Thèse de doctorat, Carnegie Mellon University. DOI : [10.1184/R1/21708035.v1](https://doi.org/10.1184/R1/21708035.v1).
- CHAUDHARY A., ANASTASOPOULOS A., PRATAPA A., MORTENSEN D. R., SHEIKH Z., TSVETKOV Y. & NEUBIG G. (2020). Automatic extraction of rules governing morphological agreement. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, p. 5212–5236, Online : Association for Computational Linguistics. DOI : [10.18653/v1/2020.emnlp-main.422](https://doi.org/10.18653/v1/2020.emnlp-main.422).
- CYSOUW M. (2005). Quantitative methods in typology (quantitative methoden in der typologie). In R. KÖHLER, G. ALTMANN & R. G. PIOTROWSKI, Éd., *Quantitative Linguistik / Quantitative Linguistics - Ein internationales Handbuch / An International Handbook*, p. 554–557. DeGruyter.
- DE MARNEFFE M.-C., MANNING C. D., NIVRE J. & ZEMAN D. (2021). Universal Dependencies. *Computational Linguistics*, **47**(2), 255–308. DOI : [10.1162/coli_a_00402](https://doi.org/10.1162/coli_a_00402).

- DRYER M. S. & HASPELMATH M., Éds. (2013). *WALS Online*. Leipzig : Max Planck Institute for Evolutionary Anthropology.
- FUTRELL R., MAHOWALD K. & GIBSON E. (2015). Quantifying word order freedom in dependency corpora. In *Proceedings of the Third International Conference on Dependency Linguistics (Depling 2015)*, p. 91–100, Uppsala, Sweden : Uppsala University, Uppsala, Sweden.
- GERDES K., GUILLAUME B., KAHANE S. & PERRIER G. (2018). SUD or surface-syntactic Universal Dependencies : An annotation scheme near-isomorphic to UD. In *Proceedings of the Second Workshop on Universal Dependencies (UDW 2018)*, p. 66–74, Brussels, Belgium : Association for Computational Linguistics. DOI : [10.18653/v1/W18-6008](https://doi.org/10.18653/v1/W18-6008).
- GERDES K., GUILLAUME B., KAHANE S. & PERRIER G. (2019a). Improving surface-syntactic Universal Dependencies (SUD) : MWEs and deep syntactic features. In *Proceedings of the 18th International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2019)*, p. 126–132, Paris, France : Association for Computational Linguistics. DOI : [10.18653/v1/W19-7814](https://doi.org/10.18653/v1/W19-7814).
- GERDES K., KAHANE S. & CHEN X. (2019b). Rediscovering greenberg’s word order universals in UD. In *Proceedings of the Third Workshop on Universal Dependencies (UDW, SyntaxFest 2019)*, p. 124–131, Paris, France : Association for Computational Linguistics. DOI : [10.18653/v1/W19-8015](https://doi.org/10.18653/v1/W19-8015).
- GERDES K., KAHANE S. & CHEN X. (2021). Typometrics : From implicational to quantitative universals in word order typology. *Glossa : a journal of general linguistics*, **6**(1). DOI : [10.5334/gjgl.764](https://doi.org/10.5334/gjgl.764).
- GUIBON G., COURTIN M., GERDES K. & GUILLAUME B. (2020). When collaborative treebank curation meets graph grammars. In *Proceedings of The 12th Language Resources and Evaluation Conference*, p. 5293–5302, Marseille, France : European Language Resources Association.
- HARTMANN I., HASPELMATH M. & TAYLOR B., Éds. (2013). *The Valency Patterns Leipzig online database*. Leipzig : Max Planck Institute for Evolutionary Anthropology.
- HERRERA S., KAHANE S. & GUILLAUME B. (2022). Extraction de règles de grammaire à partir de treebanks : développement d’un outil et premiers résultats. In L. BECERRA, B. FAVRE, C. GARDENT & Y. PARMENTIER, Éds., *Journées Jointes des Groupements de Recherche Linguistique Informatique, Formelle et de Terrain (LIFT) et Traitement Automatique des Langues (TAL)*, p. 93–98, Marseille, France : CNRS. HAL : [hal-03846825](https://hal.archives-ouvertes.fr/hal-03846825).
- HOWELL K. & BENDER E. (2022). Building analyses from syntactic inference in local languages : An hpsg grammar inference system. *Northern European Journal of Language Technology*, **8**. DOI : [10.3384/nejlt.2000-1533.2022.4017](https://doi.org/10.3384/nejlt.2000-1533.2022.4017).
- KAHANE S., VANHOVE M., ZIANE R. & GUILLAUME B. (2021). A morph-based and a word-based treebank for Beja. In *Proceedings of the 20th International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2021)*, p. 48–60, Sofia, Bulgaria : Association for Computational Linguistics.
- LEVSHINA N. (2019). Token-based typology and word order entropy : A study based on universal dependencies. *Linguistic Typology*, **23**(3), 533–572. DOI : [doi:10.1515/lingty-2019-0025](https://doi.org/10.1515/lingty-2019-0025).
- MICHAELIS S. M., MAURER P., HASPELMATH M. & HUBER M., Éds. (2013). *APiCS Online*. Leipzig : Max Planck Institute for Evolutionary Anthropology.
- NIVRE J., DE MARNEFFE M.-C., GINTER F., HAJIČ J., MANNING C. D., PYYSALO S., SCHUSTER S., TYERS F. & ZEMAN D. (2020). Universal Dependencies v2 : An evergrowing multilingual treebank collection. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, p. 4034–4043, Marseille, France : European Language Resources Association.

PENG Z., GERDES K. & GUILLER K. (2022). Pull your treebank up by its own bootstraps. In L. BECERRA, B. FAVRE, C. GARDENT & Y. PARMENTIER, Édts., *Journées Jointes des Groupements de Recherche Linguistique Informatique, Formelle et de Terrain (LIFT) et Traitement Automatique des Langues (TAL)*, p. 139–153, Marseille, France : CNRS. HAL : [hal-03846834](https://hal.archives-ouvertes.fr/hal-03846834).

ZAMARAIEVA O., CURTIS C., EMERSON G., FOKKENS A., GOODMAN M., HOWELL K., TRIMBLE T. & BENDER E. M. (2022). 20 years of the grammar matrix : cross-linguistic hypothesis testing of increasingly complex interactions. *Journal of Language Modelling*, **10**(1), 49–137. DOI : [10.15398/jlm.v10i1.292](https://doi.org/10.15398/jlm.v10i1.292).

