

# Extraction automatique de règles pour la détermination de types de relations sémantiques dans les constructions génitives en français

H. Guenoune<sup>1,2</sup>, M. Lafourcade<sup>1,2</sup>

<sup>1</sup> Université de Montpellier, France

<sup>2</sup> Laboratoire d'Informatique, de Robotique et de Microélectronique de Montpellier, LIRMM

{hani.guenoune, mathieu.lafourcade}@lirmm.fr

## Résumé

Cette étude concerne les relations sémantiques portées par les entités sous forme génitive « A de B ». Après identification des types sémantiques pertinents, nous construisons, à l'aide d'une IA générative, un corpus annoté. Nous proposons un algorithme de découverte des règles permettant de sélectionner la relation entre A et B. Ces règles correspondent à la sélection dans une base de connaissances du voisinage adéquat d'un terme donné. Soit « désert d'Algérie », portant la relation de lieu, le terme désert identifié comme lieu géographique et Algérie comme pays. Ces contraintes aboutissent par calcul à une règle permettant de sélectionner la relation de lieu.

## Mots-clés

Génitif - Entités polylexicales - Relations sémantiques.

## Abstract

We are interested in the semantic relations conveyed by polylexical entities in the postnominal prepositional noun phrases form "A de B" (A of B). After identifying a relevant set of semantic relations types, using generative AI, we build a collection of phrases, for each relation type identified. We propose an algorithm for creating rules allowing the selection of the relation between A and B in noun phrases of each type. Rules consist in selecting from a knowledge base the appropriate neighborhood of a given term. For the phrase "désert d'Algérie" carrying the location relation, "désert" is identified as a geographical location, and "Algérie" as a country. Constraints are used to learn rules for selecting the location relation for this type of example.

## Keywords

Genitive, Postnominal construction, Semantic relations.

## 1 Introduction

Au-delà de la nécessité d'identification des entités polylexicales pour l'analyse automatisée du langage, il est important pour diverses applications, de cerner la nature des rapports qui lient les différents composants des termes polylexicaux. Nous nous intéressons au cas génitif du complément de nom « de N » (construction « post-nominale

» [1]). En d'autres termes, les mots composés construits à travers l'emploi de la préposition "de" introduisant un complément syntaxique à une tête nominale (« A de B », où A et B sont des noms). Nous cherchons dans ce travail à identifier de manière automatique la relation sémantique entre les termes A et B dans les formes « A de B » (et les variantes, « A d'B », « A du B », etc.). Une telle démarche peut servir à une interprétation plus riche de contenus textuels, et permet de mettre au point des systèmes aboutissant à des représentations sémantiques connexes. Parmi les applications dont bénéficieraient une telle étude, nous pouvons citer la tâche de question-réponse (*Question Answering* ([4, 6]) qui nécessite une représentation sémantique suffisante du texte et des rapports entre les entités qui y sont mentionnées. Ou encore, la tâche de résolution des anaphores déclenchées par un déterminant possessif consistant en une transformation de formes post-nominales en syntagmes anaphorisés (« le vélo de Julie → son vélo ») et dont la résolution se base sur des contraintes sur la nature des relations entre l'anaphore et son antécédent [5].

Dans le cadre spécifique à notre projet, ces efforts sont aussi menés dans une perspective de consolidation d'une base de connaissances de sens commun, passant notamment par l'identification des types de relations sémantiques dont l'intégration s'avère la plus intéressante pour nos différents mécanismes d'inférences ainsi qu'aux applications de Traitement Automatique du Langage Naturel (TALN) qui exploitent la base de connaissances. Un autre moyen d'améliorer la qualité de la base de connaissances est de développer un système de classification servant d'*outil de contrôle*. L'analyse de la justesse des résultats d'un tel outil apporterait des éclaircissements sur la qualité globale des connaissances utilisées. Pour que cela soit possible, l'accent doit être mis sur l'*explicitabilité* des méthodes à utiliser. Il est en effet essentiel de disposer des explications des résultats du système afin d'être en mesure d'identifier les lacunes potentielles et mettre en évidence les moyens appropriés de consolider la base de connaissances.

Dans « A de B », la tête nominale (A) joue un rôle important dans le sens qu'entretient le syntagme avec son complément (B). Les types de noms liés syntaxiquement par la préposition dans une construction génitive conditionnent

la relation sémantique qui les lie. La distinction citée dans [1, 2, 3], différencie l’usage des noms de *classes d’entités* (*sortal nouns*) des *noms relationnels* (*relational nouns*). La différence dans leurs définitions est analogue à celle des prédicats unaires et binaires de la logique de premier ordre. Certains noms ne prennent leurs sens qu’en étant rapportés à exactement deux arguments. L’exemple cité est celui des noms de famille [*père, mère, sœur, frère...*]. La relation entretenue entre le nom et son complément dans la phrase « la mère de Lucie » se rapporte directement à la sémantique de la tête nominale (le nom *relationnel* “*mère*”). Un syntagme à tête nominale relationnelle permet donc une *interprétation lexicale* du type de relation sémantique et s’oppose aux cas de lecture pragmatique qui résultent de l’emploi de certains noms de classe (*sortal nouns*) et qui nécessitent de recourir à des informations extra-lexicales afin d’identifier la nature du lien sémantique entre les termes du syntagme. Parler du « *nuage de Marie* » implique d’introduire des éléments d’ordre pragmatique afin d’être en mesure d’interpréter le type de relation (*le nuage qu’elle regardait, qu’elle dessinait, dont elle a rêvé, etc.*).

Les sens portés par ce type de syntagme sont donc variés, quand bien même les efforts d’interprétation automatiques réduisent bien souvent le type de relations sémantiques à la relation d’appartenance/possession (un élément d’explication serait l’importance de la relation de possession et son rôle prépondérant dans les typologies sémantiques standards considérées dans les travaux de *TALN*).

Au-delà du cadre typologique des têtes nominales, les noms utilisés (qu’ils soient relationnels ou de classe) introduisent une multitude de types de relations possibles entre les deux termes. Ce travail cherche à étudier la nature des liens sémantiques dans cette configuration, et vise à proposer une *typologie sémantique* de tels syntagmes polylexicaux. En outre, nous proposons un algorithme d’apprentissage, servant de support à un système de classification des types dans les syntagmes polylexicaux en complément de nom (avec « *de* »).

Les sens figurés étant difficilement accessibles au calcul, ce travail n’aborde pas la détermination du sens global de la forme « *A de B* » quand celle-ci a également acquis un sens idiomatique (*par exemple* : « *homme de paille / écran de fumée* »).

Cet article portera donc sur plusieurs points :

**Proposition d’une typologie des relations sémantiques** dans les syntagmes en complément de nom, puis production d’un corpus associatif entre des exemples de constructions génitives et les types de relations qui y correspondent. Les données sont collectées en exploitant *avec prudence* une IA générative, puis validées manuellement. Une partie de ce corpus servira de données d’entraînement et l’autre d’ensemble de test à un système d’apprentissage de règles de classification ;

**Présentation de *GRASP-it*** (*Genitive Relations And Semantic Pattern Identification Tool*), un algorithme d’apprentissage calculant des règles de décision pour le/les types de relations probables ;

**Évaluation de la qualité des règles produites** en implémentant un deuxième algorithme de classification des types sémantiques dans les formes « *A de B* ».

Nous commençons par une présentation des ressources utilisées dans le cadre de ce projet, en l’occurrence la base de connaissances pour laquelle nous souhaitons apporter des améliorations à travers cette classification des types sémantiques, puis les données ayant servi au développement de *GRASP-it*. Nous donnons, dans cette section, quelques exemples de chaque type de relation résultant de l’emploi des constructions en « *de N* ». Nous décrivons ensuite le mécanisme d’apprentissage mis en place afin de synthétiser, - automatiquement et sous forme de couples de contraintes -, les rapports sémantiques entre les termes de chaque type présenté dans la section qui précède. Pour finir, nous proposons une évaluation de la qualité des règles produites, réalisée à travers une procédure d’identification des types de relations appliquée à une portion du corpus.

## 2 Ressources

Ce projet a nécessité l’emploi de ressources externes pour mener à bien l’apprentissage réalisé par l’algorithme *GRASP-it*, mais également pour la mise au point d’heuristiques de classification et l’évaluation du système réalisé.

### 2.1 Base de connaissances utilisée

L’étude et les algorithmes développés sont soutenus par une base de connaissances issue de la dernière version de la collection de données du projet *JeuxDeMots* (*JDM*) (datée du 11 février 2024).

*JeuxDeMots* (*JDM*) [9] est un réseau lexico-sémantique représenté par un graphe orienté. Les nœuds du graphe sont des termes ou des informations de nature symbolique, tandis que les arcs indiquent des relations typées, pondérées et potentiellement annotées entre les nœuds. Le graphe aborde la polysémie lexicale en spécifiant des sens hiérarchiques “*raffinements*”, où un sens spécifique est affilié au sens général du terme. *JDM* repose sur des outils, principes et concepts pratiques (*e.g.* la notion de raffinement, la diversité des types sémantiques, des relations inverses telles que *r\_isa* et *r\_hypo*, ainsi qu’une série d’outils collaboratifs).

Le réseau *JDM* est conçu pour être utilisé comme support de connaissances pour les solutions d’IA (analyse sémantique de texte, raisonnement, prise de décision, résumé automatique, etc.). Un système de pondération et d’évaluation symbolique (annotation de méta-informations, *e.g.* *rare*, *pertinent*, *non-pertinent*, etc.) a été mis en œuvre pour faciliter le parcours du graphe et son exploitation.

Au 1er février 2024, *JDM* contenait environ 540 millions de relations entre plus de 9 millions de termes et 24 millions de nœuds [8].

Un des objectifs centraux de ce travail est d’enrichir notre base de connaissances avec des informations sémantiques, en particulier celles concernant les relations dans les syntagmes nominaux génitifs. Cela contribuera principalement aux tâches d’analyse textuelle et à l’extraction de connaissances, en particulier. En effet, lorsqu’on rencontre dans

un texte une forme génitive « A de B », il est souhaitable d’avoir des outils d’identification de la (ou des) relation(s) entre A et B.

## 2.2 Corpus de constructions génitives

Nous proposons un corpus de petite taille servant à l’apprentissage des règles de détermination de type sémantique et à leur évaluation. Ce corpus, mis à la disposition de la communauté, peut également être vu comme un point de départ à la création de collections de constructions génitives de plus grande envergure. En effet, en dépit de l’importance des petits corpus [10], en deçà de plusieurs milliers d’exemples, il se révèle difficile de destiner ces données à une exploitation par des procédures gourmandes en ressources telles que des algorithmes d’apprentissage neuronal. Toutefois, nous souhaitons inscrire cet effort dans un projet à plus long terme dans lequel des procédures d’augmentation de données peuvent être mises en place, tels que des mécanismes automatiques d’enrichissement sémantique ou encore une complétion par annotation manuelle.

Dans ce qui suit, nous détaillons le protocole d’acquisition et de validation des données, puis listons les types sémantiques identifiés pour les relations entre les éléments de la paire de syntagmes dans une construction génitive. Afin d’éviter qu’un quelconque biais soit introduit, nous avons choisi de collecter des données à partir d’une source indépendante de la base de connaissances *JDM*.

### 2.2.1 Typologie sémantique

Dans le Tableau 1, nous listons les types sémantiques que nous avons choisis de considérer. Nous apportons pour chacun d’eux une explication et quelques exemples, ainsi que le type de relation correspondant dans *JDM* (la relation sémantique avec l’orientation appropriée, les relations dont le nom a la forme ‘*r\_x-I*’ étant la relation converse de ‘*r\_x*’).

Type de relation	Relation <i>JDM</i>
Conséquence (Co) : <i>Le terme A est une conséquence de (est causé par) B.</i>	<i>r_has_causatif</i>
<i>dégâts de la tempête - retards de la circulation</i>	
Caractérisation (Ca) : <i>Le terme A est une propriété ou le nominalisation d’un adjectif pouvant qualifier B.</i>	<i>r_has_property-I</i>
<i>sournoiserie du politicien - sagesse du vieillard</i>	
Matière/composition (M) : <i>Le terme A est composé de ou est de la matière B.</i>	<i>r_objet&gt;matière</i>
<i>cuillère de bois - trône de fer</i>	
Origine (O) : <i>Le terme A est originaire de B.</i>	<i>r_lieu&gt;origine</i>
<i>vin de France - café du Brésil</i>	

Topic (T) : <i>Le terme A a pour thème (ou sujet) le terme B.</i>	<i>r_topic</i>
<i>restaurant de sushis - film d’horreur</i>	
Dépeiction (D) : <i>Le terme A est une représentation du terme B.</i>	<i>r_depict</i>
<i>peinture d’un paysage - photo d’une famille</i>	
Holonymie (H) : <i>Le terme A fait partie de B.</i>	<i>r_holo</i>
<i>coque du bateau - écaille du poisson</i>	
Lieu (L) : <i>Le terme peut avoir pour lieu le terme B.</i>	<i>r_lieu</i>
<i>tour de Pise - sahara d’Algérie</i>	
Agent (A) : <i>Le terme A est la nominalisation d’une action dont l’acteur est le terme B.</i>	<i>r_processus_agent</i>
<i>travail de l’ouvrier - cours du professeur</i>	
Patient (P) : <i>Le Terme A est la nominalisation d’une action que le terme B subit.</i>	<i>r_processus_patient</i>
<i>travail du bois - ouverture de la porte</i>	
Instrument (I) : <i>Le terme A est instrument de l’action B ou d’une action que le terme B subit.</i>	<i>r_processus_instr-I</i>
<i>clé d’ouverture - clé de la porte</i>	
Possession (Po) : <i>A est possédé par B</i>	<i>r_own-I</i>
<i>fusil du soldat - vélo du cycliste</i>	
Quantification (Q) : <i>A sert de mesure à B.</i>	<i>r_quantificateur</i>
<i>brin d’herbe - minute d’attente</i>	
Lien social (LS) : <i>A tient un rôle de ‘A’ vis-à-vis de B.</i>	<i>r_social_tie</i>
<i>avocat d’une femme battu - chef du groupe</i>	
Auteur/créateur (AC) : <i>A est produit par B.</i>	<i>r_product_of</i>
<i>portrait de Van Gogh - gâteau du pâtissier</i>	

TABLE 1 – Liste des types sémantiques considérés dans les entités « A de B » et leur correspondances avec les types de relations sémantiques dans le réseau *JDM*.

Il est à noter que cette liste est celle que nous avons arrêtée en tant que typologie de base de notre étude. Les choix concernant la granularité et le nombre de types ont été faits de sorte à aligner cette typologie avec les exigences des ressources et des outils que nous utilisons ainsi que les be-

soins des applications dans lesquelles cette typologie est exploitée. La typologie adoptée en tant qu'état de l'art pour l'étude des modificateurs de noms est donnée dans [7]. Il ne s'agit, dans notre cas, aucunement d'une liste exhaustive de tous les types de relations possibles entre les termes A et B d'une forme « *A de B* ». Quelques cas peuvent s'ajouter à cette liste. Il est également possible de *spécifier/généraliser* certains types de manière à ce qu'ils correspondent plus ou moins précisément à des cadres théoriques différents du nôtre. Notamment dans le cas d'une exploitation avec une autre base de connaissance (que JDM), définissant un ensemble de types de relations différent. Parmi ces cas, nous pouvons mentionner les exemples portant des relations sémantiques temporelles *absolues* (portées par des noms de classe), « *repas de midi - brise du matin - bus de nuit* », ou encore des lien *relatifs* spatiaux et temporels (portés par des noms relationnels) « *milieu/droite/gauche de la pièce - bas de page* ». Un autre cas est celui des nominations : « *Théorème de Pythagore - Rôle de Wallace - Kappa de Fleiss* » qui pourrait faire l'objet d'une catégorie à part entière. Nous choisissons pour les cas mentionnés de les inclure dans des types de sémantique similaire, par exemple, les deux premiers cas sont inclus dans le type *topic*, tandis que le cas des nominations, lui, est classé parmi les instances de *auteur/créateur* (même s'il n'est pas systématiquement question de création).

### 2.2.2 Collecte de données et validation

Pour chaque type de relation sémantique ci-dessus, nous avons fait appel à une IA générative<sup>1)</sup> pour créer un ensemble d'exemples. Nous avons choisi de construire un corpus de forme génitives à partir d'une source indépendante de la base de connaissances JDM. La première raison est de s'assurer de ne pas introduire de biais entre les données du corpus (servant à l'apprentissage et au test de notre algorithme) avec la base de connaissances utilisée pour l'extraction des attributs représentant la sémantique des termes (extraction des *signatures*, voir 3). Une seconde raison est la nécessité d'annoter un syntagme « *A de B* » issu de JDM avec la (ou les) relations sémantiques attendues entre A et B. Or cette information n'est pas systématiquement disponible dans JDM alors qu'elle peut être demandée à l'IA générative. Un travail manuel important deviendrait donc nécessaire et limitant. D'aucun pourrait penser que la base de connaissances pourrait contenir une ou plusieurs relations sémantiques entre les termes A et B, mais rien ne garantirait qu'il s'agisse des mêmes relations que celles attendues (introduites par le syntagme).

Concrètement, nous avons fait le choix de limiter chaque type considéré à 80 exemples, dont 50 sont consacrés à l'apprentissage et 30 servent à l'évaluation de l'algorithme présenté dans la section 3. La stratégie de construction des requêtes émises à l'agent conversationnel a été différente selon les types de relations. En effet, il s'est avéré plus ou moins difficile, selon les cas, d'obtenir des exemples satisfaisants. Pour les types où les exemples générés étaient

peu exploitables, nous avons choisi d'orienter le modèle par l'exemple. Nous avons procédé en donnant une dizaine d'exemples rédigés par nos soins, puis en expliquant les points communs au niveau des relations sémantiques sous-jacentes. Ces explications se sont étalées sur plusieurs *tours de parole* avec l'agent conversationnel.

Quand bien même cette démarche itérative a permis d'obtenir les exemples du type recherché, elle présente néanmoins l'inconvénient « *d'influencer* » excessivement les réponses produites par le *chatbot*, ce qui mène à un ensemble d'entités polylexicales de faible diversité (forte adéquation aux exemples proposés à l'agent). Il a par conséquent été nécessaire, dans un but de diversification, à partir d'une première génération d'exemples, d'insister sur le fait de réitérer la génération en cherchant à la diversifier. Cette stratégie a été répétée jusqu'à ce que nous ayons considéré l'ensemble comme convaincant. Toutefois, les ensembles contenaient environ 10% d'exemples mal classifiés ou dupliqués et sont également restés imparfaits concernant leur diversité. Nous avons donc procédé à une validation manuelle de l'ensemble des exemples produits par le *chatbot*. Précisément, la validation a consisté en un remplacement des cas de duplication et des entités trop similaires, ainsi qu'un reclassement des exemples mal classés.

### 2.2.3 Mise en forme et post-traitement des données

La collection produite inclut des exemples de morphologie variable. En effet, en ce qui concerne la présence ou non de déterminant du complément du nom (terme B), on pourrait supposer qu'une étape de normalisation morphologique serait intéressante. Il s'avère néanmoins que ce critère constitue un marqueur morpho-syntaxique pouvant se révéler crucial pour une classification, raison pour laquelle nous faisons le choix de ne pas opérer de transformation au niveau morphologique ou lexical. Toutefois, notons qu'une polysémie des termes du corpus peut éventuellement être observée. L'exploitation du corpus nécessitera donc de préparer les données, en l'occurrence une phase de désambiguïsation sémantique devra être menée afin de sélectionner les sens appropriés des termes.

## 3 Présentation de GRASP-it

L'algorithme *GRASP-it* (*Genitive Relations And Semantic Pattern Identification Tool*) vise à produire un ensemble de paires de contraintes pour chaque type de relation en se basant sur les données d'entrée. Ces contraintes sont fondées sur les types sémantiques de la tête nominale et du complément. Elles peuvent être considérées comme une synthèse des attributs sémantiques alignée avec le contenu d'une base de connaissances. Le but de cet ensemble de contraintes est de guider un processus de classification des relations sémantiques dans les syntagmes nominaux génitifs.

Un autre objectif est de produire des contraintes "interprétables" qui peuvent être facilement lues et expliquées. Dans le cas général, la première étape de *GRASP-it* implique de stocker, pour chaque exemple d'un certain type, des informations sémantiques qui pourraient permettre de classer

1. LLM GPT. Version du modèle : gpt-4-0613, datée du 13-06-2023.

l'exemple dans le type pertinent :

- *Hyperonymes des termes A et B (H)* : L'objectif est de capturer, aussi précisément que possible, les "types" sémantiques des deux termes. Un hyperonyme est un terme (entité lexicale) dans JDM accessible via la relation *r\_isa*.
- *Cible des types de relation (TRT)* : Une sélection de types de relation conduisant au terme. Par exemple, un terme fréquemment ciblé par la relation de localisation est considéré, par cette approche, comme une localisation. Cela permet de renforcer la pertinence de cette classe sémantique pour un terme spécifique. La sélection des relations conduisant aux termes peut être vue comme un moyen de compléter la liste des hyperonymes pour un terme donné.
- *Type sémantique standard (SST)* : À travers la relation *\_INFO-SEM*, le type standard associe un terme lexicalisé à un type ontologique (conceptuel) standard.

Le résultat de cette étape est un ensemble de paires pondérées, appelées ici *signatures* des termes A et B. Le nombre de paires à cette étape correspond au nombre d'exemples pour chaque type, qui, dans le cas de notre corpus (portion d'entraînement), s'élève à 50 unités de la forme « *A de B* ». Une signature est définie comme un ensemble non ordonné de symboles. Chaque symbole prend une valeur d'une entrée spécifique de JDM. Par exemple, la signature *s* associée au terme « *véhicule* » serait la suivante.

```
s(véhicule) = {
véhicule, transport urbain, partie de
l'espace, Transport urbain, mode de
transport, instrument, lieu, transport,
moyen, machine, moyen de transport,

r_isa, r_hypo, r_has_part, r_holo,
r_agent, r_patient, r_lieu, r_instr,
r_carac-1, r_lieu-1, r_action_lieu,
r_mater>object, r_processus>agent,
r_own, r_is_instance_of,

_INFO-SEM-SUBST, _INFO-SEM-THING-
ARTEFACT, _INFO-SEM-PLACE, _INFO-SEM-
THING-CONCRETE, _INFO-SEM-PLACE-HUMAN
}
```

Pour des raisons de clarté, nous avons divisé les symboles en trois blocs : Hyperonymes, TRT et SST. Il convient de noter que la signature d'un terme contient le terme lui-même, dans le but de capturer des instances qui sont des hyponymes du terme signé. En plus de la nécessité de son explicabilité, cette représentation des termes est conçue pour être contrôlable du point de vue de son contenu et de sa taille. Cela permet à la méthode *GRASP-it* d'être adaptable aux exigences variables de l'application pour laquelle elle est utilisée.

La deuxième étape vise à agréger les règles de chaque type pour traiter l'ensemble complet par généralisation. Comme

le montre (1), nous définissons une règle *R* comme une paire de contraintes *sL* et *sR* (qui sont des signatures, respectivement gauche et droite correspondant aux termes A et B) et un type de relation sémantique *rt*.

$$R : \langle s_L, s_R, rt \rangle \quad (1)$$

L'agrégation est une opération de fusion de deux règles et est définie dans (2).

$$Fusion(R1, R2) = \langle s_{1L} \cup s_{2L}, s_{1R} \cup s_{2R}, rt \rangle^2 \quad (2)$$

Une fusion de deux règles signifie que les contraintes qu'elles associent respectivement sont suffisamment similaires pour être représentées par une seule paire de contraintes. Formellement, comme une signature *s* peut être vue comme un vecteur, nous avons adopté la similarité cosinus (produit scalaire divisé par le produit des normes), notée *sim*. Deux signatures sont considérées comme suffisamment similaires lorsque leur valeur *sim* est supérieure à un seuil de 0.5 (établi empiriquement). La signature fusionnée est la somme vectorielle des deux signatures (ce qui correspond à l'union ensembliste).

Une paire produite par une ou plusieurs fusions successives est considérée comme plus générale et fiable qu'une paire qui n'a pas subi de fusion. La fiabilité est donc une mesure de la couverture des exemples du type et est calculée en attribuant un poids à la paire de contraintes correspondant au nombre de fusions effectuées pour arriver à la forme finale de la paire. À la sortie de cette étape, un ensemble de paires de contraintes plus ou moins agrégées, avec une cardinalité d'au plus deux fois le nombre d'exemples du type considéré, est attribué à chaque type de relation considéré. L'idée derrière la fusion des règles est que le résultat de fusions successives est une règle qui représente de manière appropriée un grand ensemble d'exemples d'un certain type. Une règle fusionnée peut donc être considérée comme un *modèle généralisé* pour un type de relation donné. Un type de relation peut être associé à plusieurs modèles. Un bon modèle associera de manière appropriée le type de relation entre deux termes A et B dans un syntagme génitif.

## 4 Evaluation de *GRASP-it*

Dans cette section nous présentons les conditions dans lesquelles notre évaluation a été réalisée et finissons par l'analyse des scores de performance de notre système.

### 4.1 Préparation des données

Une phase minimale de préparation des données a été entreprise avant d'appliquer l'algorithme de classification (détaillé dans la section 4.2). Afin de réussir à prendre en compte l'intégralité du corpus lors de l'entraînement ainsi que de la classification, nous avons identifié deux tâches principales qu'il est nécessaire de réaliser au préalable.

2. Notons que seules des règles ayant le même *rt* sont à fusionner.

**Identification des mots composés :** Les instances de syntagmes nominaux contenant plusieurs prépositions "de", comme dans "lunettes de soleil de marque - détecteur de fumée de protection", posent le problème du choix de la bonne préposition de séparation. Cela a une influence directe sur l'identification des termes A et B, et par conséquent sur le type de relation approprié à identifier. Nous avons procédé à l'identification de ces instances en vérifiant leur existence dans la base de connaissances. Dans "lunettes de soleil de marque", les candidats pour les termes A et B seraient "(A : lunettes, B : soleil de marque)" et "(A : lunettes de soleil, B : marque)". Dans le premier cas, l'inexistence du terme B dans JDM nous permet d'assigner les valeurs du dernier candidat à A et B. Dans le cas où les deux candidats résultent en des termes A et B connus, la préposition de séparation *prévue*<sup>3</sup> est désignée manuellement.

**Entités nommées génériques :** Les instances contenant des entités nommées telles que le prénom ou le nom de famille d'une personne ne sont bien représentées dans notre base de connaissances que lorsque ces entités sont renommées ou relève de connaissances communes/culture populaire (par exemple, "Coca-Cola" - "Lucie"). De ce fait, il est important de prendre en compte les instances inexistantes en effectuant des transformations qui remplacent le nom par un autre (du même type) que nous savons bien représenté dans la base de connaissances.

## 4.2 Algorithme d'application

Afin d'évaluer les couples de contraintes sémantiques produit par l'algorithme d'apprentissage, nous mettons en place un processus de validation qui cherche à vérifier la satisfaction de ces contraintes, l'objectif étant d'identifier les types sémantiques tirés de la portion du corpus n'ayant pas été impliquée dans le calcul des contraintes. Il s'agit donc de 450 exemples répartis équitablement sur les 15 types de relations possibles.

### 4.2.1 Critères de décision

En théorie, l'approche de validation des contraintes est basée sur une recherche de similarité des types sémantiques des termes du syntagme en entrée et les termes à partir desquels le système GRASP-it a été entraîné. L'idée est donc d'induire une *identité de type*<sup>4</sup> si les termes nouveaux sont suffisamment proches de la sémantique synthétisant les entrées du processus d'apprentissage. En pratique, la recherche est menée en calculant une similarité entre les termes A et B de l'entrée et la signature respective dans toutes les règles  $\langle sL, sR, rt \rangle_i$  apprises par GRASP-it. Les deux similarités obtenues (pour le terme A avec  $sL$  et B avec  $sR$ ) sont agrégées par une moyenne arithmétique. Par conséquent, la similarité entre une forme « A de B » et une règle (paire de contraintes)  $\langle sL, sR, rt \rangle$  est donnée dans la formule 3.

$$\frac{1}{2}(sim(s(A), s_L) + sim(s(B), s_R)) \quad (3)$$

3. Celle permettant d'identifier des termes A et B ayant le type de relation annoté.

4. Supposer, par induction, que deux termes sont du même type sémantique

Il convient de noter qu'une signature pour un terme et une contrainte d'une règle partageant une structure identique. Une réponse positive est renvoyée pour le type le mieux classé en termes de valeurs de similarité moyenne avec chaque paire. Notons qu'afin que le classement puisse être fait, la procédure de vérification est menée une fois que les couples de contraintes sont calculés pour tous les types considérés.

Il reste possible d'avoir uniquement recourt aux règles qui sont soit le résultat d'une fusion, soit qui n'ont pas été fusionnées. Autrement dit, nous pourrions exclure les règles qui ont participé à la création d'une nouvelle règle. Nous prévoyons que cette réduction de l'ensemble de règles entraînerait un temps d'exécution plus court sans fortement dégrader les résultats. Cet aspect fait l'objet d'une évaluation (voir *Expérience 3*).

### 4.2.2 Traits extra-sémantiques

La détection de certains types dépend plus ou moins de marqueurs extra-sémantiques, tels que l'utilisation ou non de déterminant, l'utilisation d'entités nommées, ou encore la définitude des compléments de nom. Un exemple illustrant cette dépendance est le syntagme "photo de famille" par opposition à "photo d'une famille" qui, en raison de la présence ou non d'un déterminant, conditionne l'interprétation du lien sémantique (resp., *topic* et *dépicition*).

De telles heuristiques ne font pas partie de la composante principale de notre solution, tant nous souhaitons mettre en évidence l'implication des règles sémantiques *spécifiquement*. Néanmoins, étant donné que la représentation des termes (signature) consiste en un ensemble de symboles, la solution proposée semble également adaptée pour prendre en compte des informations extra-sémantiques. Cela revient à inclure les traits pertinents dans les signatures. Par conséquent, et afin de confirmer l'intuition de l'utilité de marqueurs extra-sémantiques, nous avons procédé à l'intégration du trait de définitude des compléments (B) à nos représentations. Cela constitue une étude distincte dans la section suivante<sup>5</sup> (voir *Expérience 2*).

## 4.3 Protocole d'évaluation

Dans cette évaluation, nous menons trois expériences distinctes, visant à évaluer les aspects suivants.

**Expérience 1 :** Nous cherchons à évaluer individuellement le gain en performance permis par chaque trait sémantique inclus dans les signatures. Dans cette expérience, nous prenons pour référence de base (*baseline*) l'inclusion des *hyperonymes* (H) uniquement. L'idée est de procéder de manière contrastive et d'analyser les cas que nous pourrions classer avec succès en ajoutant séparément les traits TRT et SST.

**Expérience 2 :** L'inclusion de marqueurs morphologiques pourrait se révéler bénéfique pour le processus de classification. Sans concevoir d'heuristiques élaborées *ad hoc* pour ces traits, nous expérimentons les effets induits

5. Une classification aboutie devra néanmoins faire appel à des règles *ad hoc*, éventuellement plus élaborées, traitant le cas des marqueurs relevant de la morphologie, par exemple.

par la construction des signatures avec des traits extra-sémantiques, à savoir le trait de *définitude* et la présence ou non d'un déterminant pour le terme B. Nous utilisons une approche symbolique simple expliquée dans 4.4.2.

**Expérience 3 :** Ici, nous nous intéressons à estimer la valeur globale d'efforts computationnels supplémentaires (efforts se traduisant par des coûts de calcul par exemple). Le système développé étant également destiné à être utilisé comme composant de tâches spécifiques de *TALN* (les plus critiques d'entre elles étant l'extraction de relations à partir de textes ou la résolution d'anaphores), il est important d'étudier la faisabilité d'intégrer *GRASP-it* dans de tels systèmes. Dans ces systèmes appliqués, les exigences pour les sous-tâches concernent souvent le temps d'exécution. Par conséquent, cette expérience est conçue pour étudier le gain/perte de performance par rapport au temps de calcul dans deux paramétrages différents.

Afin de maintenir une équivalence entre le nombre d'exemples pour chaque classe, nous ne considérons pas les cas de classification multiple dans cette évaluation. Compte tenu de la méthode de création du corpus, un seul type est associé à chaque exemple. Les cas pouvant être classés dans plus d'un type devront par conséquent être renseignés à la main. Indépendamment de la justesse des prédictions supplémentaires (fortuites/incidentelles), ces exemples n'étaient pas prévus comme appartenant à cette classe supplémentaire et ne sont donc pas comptés dans le nombre d'instances. De plus, le nombre de ces cas de classifications multiples n'est pas prévisible, ni uniformément réparti entre les types que nous considérons.

Selon les besoins d'exigence de l'évaluation, une alternative indulgente ne nécessitant pas d'annotations supplémentaires peut être envisagée. Il serait alors possible de considérer comme correctes non seulement les instances classifiées dans le type annoté dans le corpus de test, mais également celles pour lesquelles le type attendu se serait retrouvé en deuxième position (ou en position  $n$ ) dans le classement des types par valeurs de similarités, si (et seulement si) sa valeur est proche de celle calculée pour le type le mieux classé. En d'autres termes, il s'agirait de considérer la sortie de l'algorithme comme une classification dans  $n$  classes différentes lorsque celles-ci auraient obtenu les  $n$  meilleures valeurs de similarité et que l'écart maximal entre les  $n$  meilleures valeurs aura été jugé négligeable (lorsque la  $n$ -ième meilleure valeur de similarité est à moins de 5% de la première, avec  $n = 2$ , par exemple, serait une prédiction d'appartenance dans 2 classes).

#### 4.4 Scores de performance

Dans ce qui suit, étant donné un type de relation, nous considérons la précision (P) d'une classe comme la proportion d'exemples pour lesquels la classe est correctement prédite par rapport à toutes les instances prédites comme appartenant à cette classe. Le rappel (R) représente le rapport d'exemples pour lesquels la classe est correctement prédite à toutes les instances réelles de cette classe.

##### 4.4.1 Expérience 1 - Traits sémantiques

Une approche basée uniquement sur la sémantique de A et B donne les résultats illustrés dans le Tableau 2. Afin d'évaluer séparément le gain de performance permis par chaque trait inclus dans la signature des termes, nous rapportons les résultats de 4 configurations contrastives des paramètres de *GRASP-it*.

Configuration	P (%)	R (%)	F1
<i>H</i>	67,3	65,9	0,653
<i>H+SST</i>	70,4	69,7	0,691
<i>H+TRT</i>	77,6	77	0,767
<b><i>H+TRT+SST</i></b>	<b>78</b>	<b>77,3</b>	<b>0,772</b>

TABLE 2 – Moyennes de précision P (%), rappel R (%), et score F1 permis par les différentes configurations de *GRASP-it*.

Les éléments combinés pour construire chacune des configurations représentent les informations stockées lors du calcul des signatures (donc, lors de l'apprentissage des règles) : Hyperonymes (*H*), Cible pour les types de relation (*TRT*), et Types sémantiques standard (*SST*).

Tout d'abord, en tant que *baseline*, les hyperonymes, étant des traits lexico-sémantiques, conduisent à un score *F1* moyen de 0,653, ce qui est satisfaisant compte tenu de la rareté potentielle des hyperonymes de certains termes (par exemple, nous attirons l'attention sur les termes A des types *Caractérisation (Ca)*, *Agent (A)* et *Patient (P)* étant tous des entités abstraites pour lesquelles il est délicat d'identifier un hyperonyme lexical). Deuxièmement, nous observons que les traits *SST* et *TRT* conduisent à des améliorations conséquentes, le gain non linéaire ayant naturellement tendance à devenir moins important à mesure que les scores s'améliorent. Les traits *H* et *TRT* d'une part et les *SST* d'autre part, sont complémentaires car ils répondent chacun à des besoins spécifiques de description. Les *SST* apportent la couverture des typologies standards (par exemple, en fournissant le type *\_INFOSEM-THING-ABSTRACT* qui aide dans le scénario discuté ci-dessus), tandis que les *TRTs* et les *H* (étant des entrées terminologiques de nature lexicale) apportent une granularité plus fine rendue possible par l'abondance de la terminologie.

Avec un score *F1* de 0,772, la configuration la plus favorable est celle combinant tous les traits sémantiques. Le Tableau 3 rapporte les résultats de la configuration *H+TRT+SST* pour chaque type de relation sémantique considéré.

Nous observons des résultats relativement élevés, bien que des disparités existent en fonction du type à identifier. En particulier, le faible rappel pour la relation *Caractérisation (Ca)* peut être attribué à sa représentation limitée dans la base de données (environ 10000 relations comparées à la relation *Holonymie (H)*, qui en compte plus de 10 millions), entraînant une faible proportion d'exemples d'apprentissage correctement annotés. Il en va de même pour

Type	P	R	F1
Origine (O)	100	86	0,92
Lien social (LS)	83	100	0,91
Holonymie (H)	78	86	0,82
Quantification (Q)	82	80	0,81
Agent (A)	71	93	0,81
Dépiation (D)	88	73	0,80
Matière (M)	78	83	0,80
Instrument (I)	77	80	0,78
Lieu (L)	84	70	0,76
Topic (T)	68	86	0,76
Patient (P)	74	76	0,75
Auteur (AC)	76	73	0,74
Conséquence (Co)	76	63	0,69
Possession (Po)	65	63	0,64
Caractérisation (Ca)	71	50	0,59
<b>Moyenne</b>	<b>78</b>	<b>77,3</b>	<b>0,772</b>

TABLE 3 – Pourcentages (%) de Precision (P), Rappel (R), et scores F1 pour la meilleure configuration sémantique du système (*GRASP-it*<sub>(H+TRT+SST)</sub>), pour chaque type de relation considéré.

les exemples de test (jusqu'à la moitié des cas). De plus, dans le cas de (*Ca*), les cas génériques sont également peu représentés et souvent associés à un sens du terme qui n'est pas une propriété (voir Section 4.5). Il est à noter qu'il y a un rappel maximal pour le type de *lien social (LS)* porté par des têtes nominales *relationnelles*, ce qui permet une interprétation lexicale de ce type et est bien synthétisé par les contraintes créées. Le type d'*Origine (O)* est précis en raison de son petit ensemble de règles générales (un grand nombre de règles ont pu être fusionnées), mais échoue pour les exemples qui ne sont pas bien représentés dans le corpus. Ce cas est intéressant car les instances du type *Origine* sont presque inexistantes dans *JDM* (29 relations), les hyperonymes (H) ayant facilité la synthèse de règles efficaces.

À bien des égards, le protocole suivi cherche à évaluer le modèle sans complaisance; en effet, les scores devraient être interprétés comme une limite basse à améliorer à travers divers traitements de processus de classification utilisant les règles de *GRASP-it*. Il convient de noter l'absence d'heuristiques morpho-syntaxiques (extra-sémantiques). De plus, les instances qui ont été considérées comme erronées (selon le corpus) mais qui sont en réalité également valides pour le type prédit sont comptabilisées comme des erreurs. Une évaluation à classes multiples ferait évoluer les scores de chaque type à la hausse<sup>6</sup>, elle nécessiterait néanmoins une annotation manuelle.

#### 4.4.2 Expérience 2 - Définitude

Nous cherchons, en plus des traits sémantiques discutés précédemment, à *retenir* le patron morpho-syntaxique du syntagme à décrire afin de prendre en considération le caractère défini ou indéfini du complément de nom (B). Ceci

6. À titre d'exemple, le score F1 du type le moins bien classé (*Ca*) passerait à 0,76

se fait en incluant dans la signature du terme B deux symboles distincts correspondant à la présence (*resp.* absence) d'un déterminant défini ou indéfini (*Det, NoDet, Def, NoDef*). Un cas particulier concerne les entités nommées où, en dépit d'un *NoDet*, l'attribut *Def* est "forcé". Quelques exemples :

- *chat du rabbin* => *Det + Def*
- *écran de cinéma* => *NoDet + NoDef*
- *tableau de Chagall* => *NoDet + Def*

Configuration	P	R	F1
<i>H+TRT+SST</i>	78	77,3	0,772
<b><i>H+TRT+SST+DEF</i></b>	<b>80,3</b>	<b>79,8</b>	<b>0,795</b>
Origine (O)	100	90	0,95
Lien social (LS)	85	100	0,92
Holonymie (H)	81	100	0,90
Instrument (I)	88	80	0,84
Quantification (Q)	86	83	0,84
Matière (M)	83	83	0,83
Dépiation (D)	85	80	0,82
Lieu (L)	85	76	0,80
Agent (A)	67	96	0,79
Auteur (AC)	79	76	0,77
Topic (T)	70	86	0,77
Patient (P)	72	70	0,71
Conséquence (Co)	76	63	0,69
Possession (Po)	76	63	0,69
Caractérisation (Ca)	71	50	0,59

TABLE 4 – Comparaison des scores de performance entre la configuration exclusivement sémantique et celle incluant le trait de définitude (*GRASP-it*<sub>(H+TRT+SST+DEF)</sub>), pour chaque type de relation sémantique considéré.

Tel que présenté dans le Tableau 4, prendre en compte le trait de définitude améliore le score F1 global (par rapport à la configuration exclusivement sémantique). Néanmoins, nous observons une certaine variabilité dans l'amélioration, et dans deux cas particuliers (agent et patient), une diminution des performances. La raison en est que le trait de définitude n'est pas typique d'un type unique, mais plutôt d'un sous-ensemble de types. Son inclusion aide lors de la décision entre deux types pour lesquels le trait de définitude est décisif. Au contraire, cela apporte une certaine confusion entre les types dans le même sous-ensemble de règles (pour lesquels la définitude n'est pas décisive).

#### 4.4.3 Expérience 3 - Élagage des règles

Dans le tableau 5, nous mettons en évidence les différences en termes du nombre de règles appliquées (#r) et des temps d'exécution (T) pour les paramètres *Trim* et *No Trim*, qui correspondent respectivement à un ensemble complet de règles et à un ensemble réduit. L'ensemble de règles réduit contient uniquement les règles qui n'ont pas été utilisées pour une opération de fusion. Notons que (T) est la durée pour 450 instances de test.

Le temps d'exécution dépend linéairement du nombre de règles et de la taille des signatures. Plus il y a de fusions,

Configuration	P	R	F1	#r	T (s)
<i>Trim</i>	79,6	77,6	0,77	49	<b>25.42</b>
<i>No Trim</i>	80,36	80	0,798	1384	92.78

TABLE 5 – Effets de la réduction de règles sur les scores de performance et les temps d’exécution.

plus les signatures sont longues avec toutefois une asymptote sur la longueur. La configuration d’élagage (*Trim*) conduit à une amélioration considérable du temps de calcul (temps divisé un peu moins de 4 fois) avec une légère diminution de la qualité. Cela signifie que, conformément à l’intuition de départ, les systèmes qui nécessitent des réponses dans des délais particulièrement courts pourraient bénéficier de la réduction de l’ensemble de règles sans subir une dégradation significative de la performance globale.

#### 4.5 Analyse des cas d’échec

Pour rappel, ce travail fait également office d’outil de contrôle du contenu de la base de connaissances, dans le sens où les manques éventuels sont mis en évidence par les cas d’échec de classification de l’algorithme et permettent de consolider les types de relations appropriés.

Parmi les cas d’échec, nous rapportons qu’autour de 75% des occurrences sont directement dus à la polysémie du terme A et/ou du terme B (confondus). Le terme A dans « *richesse du royaume* » est considéré à tort dans le sens d’un *bien* > *objet* et non de la propriété d’*abondance*, ce qui a mené le système à sa classification (fausse) dans le type *lieu* (un objet pouvant se trouver dans un lieu). La relation sémantique correcte dans ce cas est plutôt *caractéristique*.

En écartant les cas de classes multiples (comme « *ombre d’un arbre* » ou encore « *travail du réalisateur* » classés respectivement dans *caractéristique* et *auteur*, et prédits par le système comme *dépiction* et *agent*), le reste des erreurs peuvent être expliquées par différentes raisons, parmi elles : *Les défauts de connaissances* (représentant des manques dans la base) comme dans « *restaurant de cuisine végétalienne* » (où le terme B n’était pas inclus dans JDM). Ensuite, quelques *défauts de compétence* (dûs au manque de traitements) expliquent un petit nombre de cas. En l’occurrence, la gestion de la dispersion des attributs sémantiques dans JDM à travers les différentes variantes morphologiques d’un terme : il arrive, pour des raisons intentionnelles/valables ou suite à un manque de connaissance, que certaines relations n’aient pas été propagées à toutes les dérivations d’un terme donné. En ce qui concerne notre système, ce défaut peut être observé particulièrement au niveau du trait *type sémantique standard* (SST). Dans « *liste de films* » devant être classé dans le type *Quantification*, les types standards du lemme *film* étaient manquants dans sa variante au pluriel. Bien que ce cas d’échec ait été utile en signalant l’état de la base en ce qui concerne ce type de relation, une phase de normalisation (en l’occurrence, une lemmatisation) aurait permis de contourner ce cas de dispersion.

L’analyse des cas d’échecs peut mener à l’identification d’un défaut de connaissance au niveau d’un exemple test ou d’un exemple d’apprentissage. Dans le second cas, l’impact sera plus important dans la mesure où il influencera négativement tous les exemples qui relèvent de ce type de relation.

## 5 Conclusion et perspectives

Nous avons présenté les résultats d’une recherche portant sur l’analyse automatique des formes « *A de B* ». Les contributions de ce travail ont consisté à définir une typologie non-exhaustive des relations sémantiques entre les termes A et B, puis à produire un petit corpus d’exemples annotés par ces types. Enfin, nous avons proposé un algorithme d’apprentissage de règles d’ordre sémantique sous la forme de couples de contraintes agrégés. L’objectif est d’une part, de réussir une classification automatique des types pouvant être utilisée lors de l’analyse sémantique de textes, et d’autre part, d’employer l’algorithme dans des démarches de consolidation de la base de connaissances utilisée dans le système présenté. On notera que le système proposé ne dépend pas d’une base de connaissances spécifique.

Les enseignements apportés par ce projet se définissent sur les plans théorique et appliqué. En premier lieu, nous identifions le défi que pose la possibilité d’inclusion des exemples dans plusieurs types différents. Ensuite, même si une identification du type repose fortement sur la sémantique, cette dernière ne constitue pas l’unique critère de décision de classification. Ceci signifie qu’un algorithme performant devra inclure une série de traitements, notamment de préparation des données, mais aussi d’analyse (identification de mots composés, gestion de la polysémie, normalisation des entrées, etc). En outre, comme pour toute construction de corpus, le problème de l’équilibre de la distribution et de la couvertures des traits/modèles se pose. Ce problème est d’autant plus marqué pour une collection relativement petite de données.

La question de la couverture est difficile : intuitivement on pourrait penser qu’il faut augmenter le nombre d’exemples. Cependant la couverture des types de relations dans les formes génitives suit une loi de puissance, c’est-à-dire qu’il y a un grand nombre de cas spécifiques dans la longue traîne. Ces cas sont difficilement calculables car ils peuvent correspondre à des formes figurées et sont du point de vue des signatures souvent des exemples non-prototypiques (des quasi hapax). De plus, les formes prototypiques correspondant au début de la distribution en loi de puissance sont souvent nombreux (par exemple, un nombre impressionnant de « *<animal> de <lieu>* »). Multiplier les exemples d’apprentissage ayant quasiment les mêmes signatures (représentés par la même règle *modèle*) n’a que peu d’effet, sur le long terme, sur la qualité de l’apprentissage. De surcroît, étant donné la contrainte de coût de calcul posée par le besoin d’intégrer la solution dans des systèmes appliqués, il serait contre-productif de rajouter des exemples relevant d’un modèle déjà connu. On peut cependant envisager une

stratégie qui exploiterait un critère de fusion (règle fusionnable ou pas) via un apprentissage incrémental.

L'approche symbolique adoptée présente l'avantage d'être facilement explicable, ses résultats dépendent aussi bien des types considérés, de la richesse de la connaissance du monde utilisée, que de la qualité des données (justesse et pertinence mais également l'équilibre *signaux faibles/forts* qui sont autant d'enjeux de représentation des cas discutés ci-dessus).

Le fait que les règles ne soient pas mutuellement exclusives (et dans certains cas quasi-identiques) représente une difficulté qui ne pourrait être levée, - même pour un locuteur humain -, à moins de disposer d'un contexte suffisant, par exemple : *présentation de l'élève - portrait de Van Gogh - travail du ciment*, pouvant aussi bien être classés dans les types : [Agent (A) et/ou Patient (P)] - [Agent (A) et/ou Dé-piction (D)] - [Agent (A) et/ou Patient (P)], resp.).

Étant donné la complexité et la difficulté inhérentes à de nombreux cas, la question de la détection automatique des relations sémantiques dans ce type de syntagmes reste ouverte et sujette à une exploration continue. En perspective de ce travail, nous pouvons évoquer les tâches suivantes :

- Enrichissement du corpus par l'application de notre algorithme : Une collection plus grande permettrait d'employer des méthodes diverses notamment celles nécessitant de très grands nombres d'instances telles que des approches d'apprentissage neuronal. L'idée ici est de générer un ensemble plus large d'exemples à partir des termes de la base de connaissances qui sont liés par les relations appropriées. Le résultat devra être validé manuellement afin de garantir, d'une part, la correction du syntagme nominal généré, et d'autre part, sa pertinence. Produire "*atome du garçon*" pour la classe de *holonymie* ou de *composition*, bien que correct en théorie, ne semble pas très pertinent. Les annotations de méta-informations relevant de la *pertinence* des relations dans *JDM* pourraient aider dans cette démarche ;
- Extension de la typologie des relations sémantiques considérées dans le processus d'apprentissage à davantage de types, et leur hiérarchisation afin de permettre des exploitations plus ou moins précises qui répondraient aux exigences des diverses ressources et outils pouvant être utilisés ;

## Distribution

Le corpus ainsi que l'algorithme sont accessibles sur le lien en bas de page<sup>7</sup>. À des fins de démonstration et d'expérimentation, un exemple d'implémentation de *GRASP-it* est disponible. Plusieurs paramètres sont proposés, notamment l'inclusion ou non de chacun des traits discutés dans cet article (*H*, *TRT*, *SST* et *DEF*), la réduction des règles et le

réglage du seuil de fusion sont également possibles. L'entraînement et le test peuvent être effectués aussi bien avec le corpus proposé, qu'avec une autre collection de données. De plus, il n'est pas nécessaire que les types sémantiques soient restreints à ceux définis dans cette étude (plus ou moins de types peuvent être définis et utilisés sur le démonstrateur web). Cependant, étant donné que les signatures construites pour cette implémentation sont basées sur notre base de connaissances et sa structure (la structure *JDM*), nous devons noter que si une autre collection de données est utilisée, il est nécessaire de s'assurer que les exemples sont conformes aux exigences discutées dans la section *Préparation des données* 4.1.

## Références

- [1] C. Barker, C. Maienborn, K. Von-Heusinger, P. Portner. Possessives and relational nouns. *Semantics-noun phrases and verb phrases*, pp. 177-203, 2019.
- [2] S. Löbner. Concept types and determination. *Journal of semantics*, pp. 279-333, 2011.
- [3] J. De Bruin and R. Scha. The interpretation of relational nouns. *26th Annual Meeting of the Association for Computational Linguistics*, pp. 25-32, 1988.
- [4] A. Ben Abacha. Recherche de réponses précises à des questions médicales : le système de questions-réponses MEANS. *PhD thesis*, Paris 11, 2012.
- [5] H. Guenoune. Résolution des anaphores dans la communication électronique médiée : heuristiques et apport d'informations de sens commun. *PhD thesis*, Université de Montpellier, 2022.
- [6] P. Kapanipathi, I. Abdelaziz, S. Ravishankar, S. Rouskos, A. Gray, R. Astudillo, M. Chang, C. Cornelio, S. Dana, A. Fokoue. Leveraging abstract meaning representation for knowledge base question answering. *Association for Computational Linguistics*, 2021.
- [7] V. Nastase and S. Szpakowicz. Exploring noun-modifier semantic relations. *Fifth international workshop on computational semantics (IWCS-5)*, pp. 285-301, 2003.
- [8] M. Lafourcade and N. Le Brun. Apport du jeu pour la construction de connaissances : le projet JeuxDeMots. *Technologie et innovation*, Vol. 8, pp. 4, 2023.
- [9] M. Lafourcade. Making people play for Lexical Acquisition with the JeuxDeMots prototype. *SNLP'07 : 7th international symposium on natural language processing*, p 7, 2007.
- [10] F. Landragin. Étude de la référence et de la coréférence : rôle des petits corpus et observations à partir du corpus MC4. Bases, corpus et langage-UMR 6039. *Corpus*, Vol. 2, pp.18 , 2018.

7. <https://www.jeuxdemots.org/rezo-GEN1.php>