



HAL
open science

Outils de fouille textuelle à partir d'annotations sémantiques

Marc Bertin, Iana Atanassova, Jean-Pierre Desclés

► **To cite this version:**

Marc Bertin, Iana Atanassova, Jean-Pierre Desclés. Outils de fouille textuelle à partir d'annotations sémantiques. 12e Conférence Internationale Francophone sur l'Extraction et la Gestion des Connaissances, EGC, Jan 2012, Bordeaux, France. hal-01954003

HAL Id: hal-01954003

<https://hal.science/hal-01954003>

Submitted on 13 Dec 2018

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.

Outils de fouille textuelle à partir d'annotations sémantiques

Marc Bertin*, Iana Atanassova*, Jean-Pierre Desclés*

*Université Paris-Sorbonne
Maison de la recherche
28 rue Serpente 75006 Paris
{ marc.bertin | iana.atanassova | jean-pierre.descles }@paris-sorbonne.fr

Résumé. Nous proposons des outils de fouille textuelle guidée par la sémantique reposant sur des annotations produites automatiquement selon des ontologies linguistiques. Cette approche offre à l'utilisateur final de nouveaux types de produits documentaires permettant une meilleure représentation du document dans le cadre d'une recherche d'informations.

1 Introduction

L'exploitation automatique d'annotations sémantiques, effectuée selon des ontologies linguistiques contenant des points de vue tels qu'hypothèse, citation, méthode, résultat, . . . permet de mettre en place de nouveaux produits documentaires, que nous appellerons *fiches de synthèse* ou *bibliographie augmentée*, offrant à l'utilisateur un aperçu structuré et synthétisé des contenus textuels correspondants à des requêtes de recherche d'informations sémantiques. Les outils de fouille textuelle que nous proposons ont pour but de fournir à l'utilisateur les briques de navigation textuelle permettant une meilleure représentation du document dans le cadre d'une recherche d'informations.

2 Système de fouille textuelle

A partir d'un corpus d'articles scientifiques, nous cherchons à extraire, par exemple, les hypothèses de travail, les méthodes utilisées, voire les conclusions, afin de répondre à un besoin spécifique : au-delà de la simple recherche par mots clés, nous voulons proposer à l'utilisateur un accès aux contenus guidé par la sémantique textuelle. Pour cela, nous avons construit de nouveaux produits documentaires, notamment des synthèses automatiques à partir d'une analyse sémantique des textes.

Les annotations sémantiques sont obtenues par le moteur Excom (Djioua et al., 2006; Al-rahabi et Desclés, 2008) qui permet d'identifier dans un texte des phrases porteuses des points de vue sémantiques. Il s'appuie sur des ressources linguistiques sous forme de marqueurs de surface et de règles d'exploration de contexte. Le moteur Excom implémente la méthode d'Exploration Contextuelle, proposée par Desclés (1997). Le système que nous avons mis en place exploite les fichiers annotés au format XML. Ces fichiers sont importés dans une base de données dédiée afin d'assurer les fonctionnalités de recherche d'informations par l'indexation du

contenu textuel et des annotations sémantiques. Le système propose un service web en utilisant la technologie Apache/PHP-Ajax/MySQL. La qualité des annotations sémantiques a été l'objet de plusieurs évaluations (Bertin, 2008, 2011).

3 Réalisations proposées

Recherche d'informations sémantique. Nous avons mis en place une interface de recherche d'informations sémantique, exploitant les points de vue annotés. L'utilisateur peut poser des requêtes selon des catégories sémantiques, associées au termes de filtrage qui l'intéressent. Par exemple, pour retrouver toutes les définitions liées à la "sémantique", il peut utiliser le point de vue définition et le terme "sémantique". Pour une approche multi-document, il est nécessaire de gérer les problèmes d'ordonnement et de similarité entre les réponses afin de fournir à l'utilisateur une information pertinente et non-redondante tout en proposant un retour au contexte. Cette problématique a été soulevée par Atanassova et al. (2008).

Fiches de synthèse. Les fiches de synthèse représentent des extraits catégorisés en utilisant les annotations sémantiques. Il s'agit d'un produit documentaire permettant de visualiser les différents types d'informations contenues dans un document ou un corpus de documents, sous forme structurée. Les catégories d'extraction sont personnalisables par l'utilisateur. Par exemple, la figure 1 présente une fiche de synthèse utilisant les points de vue *résultat*, *hypothèse*, *méthode*, *soulignement* et *opinion*.

Bibliographie augmentée. La figure 2 montre une bibliographie contenant non seulement les références mais également les catégories sémantiques de l'acte de citation ainsi que l'extrait correspondant. L'utilisateur a ainsi accès aux relations sémantiques entre les auteurs (voir les travaux de Bertin (2011)), à savoir comment ils sont cités : pour leurs résultats, pour une définition, pour leurs méthodes, ... Pour l'utilisateur, la lecture de la bibliographie se révèle éclairante sur la nature de l'article ainsi que sur l'utilisation des citations faites par l'auteur et offrant ainsi une meilleure représentation du document.

Distribution des références bibliographiques. La figure 3 montre la répartition des références bibliographiques en fonction des annotations sémantiques dans deux thèses de doctorat, l'une en philosophie¹, l'autre en informatique². L'axe des abscisses indique comme valeur le numéro du segment textuel, à savoir la phrase. L'axe des ordonnées affiche les différents points de vue de fouille sémantique. La répartition des références est relativement homogène selon les différents points de vue pour la thèse en informatique alors que celle de philosophie montre une prédominance des citations de la part de l'auteur. Cela nous conduit donc à revisiter, à travers l'expérimentation, l'approche et l'utilisation des citations et nous invite à nous interroger sur le rôle de certains indicateurs présents dans les systèmes d'évaluation.

1. Thèse de P. Gauvin, "*Les notions de mouvement et de changement en relation avec les théories contemporaines de l'aspectualité et de la cognition*".

2. Thèse de I. Atanassova, "*Exploitation informatique des annotations sémantiques d'Excom pour la recherche d'informations et la navigation*".

La démonstration soulignera l'importance à proposer des outils facilitant la navigation textuelle ainsi que la capacité du système à resituer un segment dans son contexte initial, c'est-à-dire dans le document source, montrant ainsi une rupture avec les systèmes de recherche d'information traditionnels.

Références

- Alahabi, M. et J.-P. Desclés (2008). Automatic annotation of direct reported speech in Arabic and French, according to a semantic map of enunciative modalities. In *6th International Conference of NLP, GOTAL*, Gothenburg, Sweden.
- Atanassova, I., J.-P. Desclés, A. Franchi, et F. Le Priol (2008). La plate-forme excom comme outil automatique d'annotations sémantiques des textes pour la catégorisation des informations sur le web. In *Internet : besoin de communiquer autrement*, Université St. Clément d'Ohride, Sofia, Bulgarie.
- Bertin, M. (2008). Categorizations and annotations of citation in research evaluation. In *The 21st international FLAIRS Conference*, Coconut Grove, Floride. AAAI Press.
- Bertin, M. (2011). *Bibliosématique : une technique linguistique et informatique par exploration contextuelle*. Ph. D. thesis, Université Paris-Sorbonne.
- Desclés, J.-P. (1997). *Systèmes d'exploration contextuelle*. Presses Universitaires de Caen.
- Djioua, B., J. G. Flores, A. Blais, J.-P. Desclés, G. Guibert, A. Jackiewicz, F. Le Priol, L. Nait-Baha, et B. Sauzay (2006). Excom : an automatic annotation engine for semantic information. *The 19th international FLAIRS Conference, Melbourne, Floride 1*, 285–290.

point de vue	Points de vue : point de vue (1) information (3) definition (1) prise de position (2) similitude (4) dissimilitude (4) résultat (12) méthode (2) citation (2)
74 Cette stagnation étonne Treiber & Wilcox (1984) qui l'attribuent à la complexité du matériel.	
information	
156 Le problème du nombre chez le bébé a été relancé de manière spectaculaire par les expériences de Wynn (1992) .	
198 Wynn, (2000) répond aux auteurs que ses résultats sont robustes et ont été reproduits dans 8 expériences impliquant 4 laboratoires différents du sien, mais également que les cas de non-réplication observés dans différents laboratoires sont eux-mêmes cohérents et explicables soit par le fait que les nombres en jeu étaient trop important, soit que l'on a demandé aux bébés deux « mises à jour » successives, soit enfin que l'opération ait impliqué un zéro, ce qui la conduit à conclure à la cohérence et à la reproductibilité des résultats.	
205 Dans une communication récente, Leslie (1999) fait état de plusieurs recherches avec des bébés de 11-12 mois dans lesquelles la technique est un peu différente : N objets sont présentés à l'enfant, on place un écran devant et une main en sort un ou passe simplement derrière l'écran et ressort vide.	
definition	
231 Gelman & Gallistel (1978) ont proposé le concept de numéron : représentation interne de la numérosité (ce qui est représenté), idée reprise par (Wynn, 1995).	
prise de position	
265 C'est aussi la position de Simon (1997) .	
297 Parfois même, c'est un paradigme expérimental précis qui est contesté (Haith, 1998; Cashion & Cohen (2000) .	

FIG. 1 – Fiche de synthèse produite par le système

Outils de fouille textuelle à partir d'annotations sémantiques

1 évaluer l'impact des stéréotypes dans les supports multimédia (14 annotations)

ADEN, Joëlle **cité 1 fois** : (1 **méthode**) ;
 ALSIC, Numéro 1, Vol. 9, pp. 103-128
[Résumé](#) [Bibliographie](#)

[Abric94] Abric, J.-C. (1994). Pratiques sociales et représentations. Paris : PUF.
 [Aden04a] Aden, J. (2004a). "La Presse économique anglophone dans la formation professionnelle initiale". In Cain, A. (dir.). Espace public et espace privé, enjeux et partage. Paris : L'Harmattan. pp. 255-267.

[Aden04b] Aden, J. (2004b). "Construction du sens et supports filmiques, guidage et autonomie en classe de langue". In Tardieu, C. & Pugibet, V. (dir.). Langues et cultures, les TIC, enseignement et apprentissage. Paris : Scéren-CNDP. pp. 135-145.
méthode Nous avons décrit la façon dont s'opère la construction sémantique de séquences filmiques en classe [Aden04b] .

[Cain91] Cain, A. (1991). "Stéréotypes culturels et construction de connaissances en civilisation". Les langues modernes, n° 2. pp.17-21.
information Dans ce travail, l'identité est le lieu de rencontre de l'individuel et du social [Cain91], elle se construit à partir des perceptions, de l'expérience personnelle et sociale.

FIG. 2 – Bibliographie augmentée

Summary

We propose text exploration tools based on semantic analyses using annotations which are carried out automatically according to linguistic ontologies. By this approach, the final user can access new types of documentary products, that offer a better document representation in an information retrieval context.

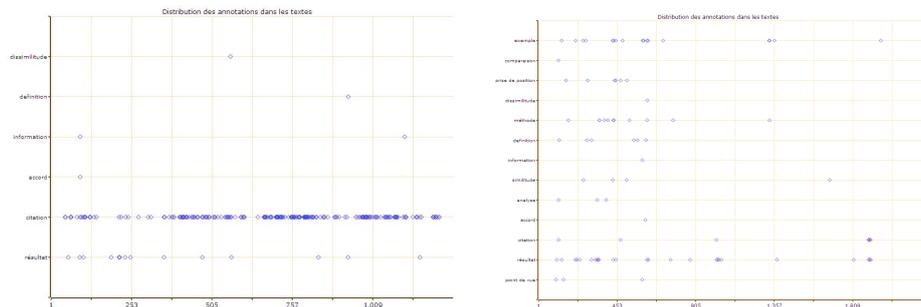


FIG. 3 – Répartition des références bibliographiques pour une thèse de philosophie (à gauche) et pour une thèse d'informatique (à droite)