

# Catégorisation d'articles scientifiques basée sur les relations sémantiques des mots-clés

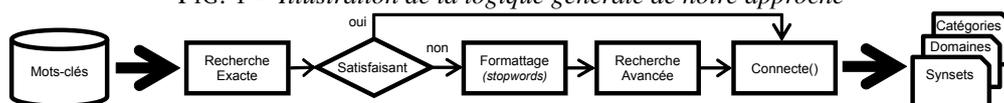
Bastien Latard<sup>\*,\*\*</sup> Jonathan Weber<sup>\*</sup>  
Germain Forestier<sup>\*</sup>, Michel Hassenforder<sup>\*</sup>

<sup>\*</sup>MIPS, Université de Haute-Alsace, Mulhouse, France

<sup>\*\*</sup>MDPI AG, Bâle, Suisse

**Introduction.** La recherche bibliographique est une étape cruciale pour tout chercheur. En effet, la connaissance des travaux existant peut faire gagner un temps précieux tant pour le choix de la méthode à adopter que pour être à jour des dernières avancées. Néanmoins, trouver des articles similaires reste une tâche compliquée et pénible autant pour les domaines étendus que réduits. Les chercheurs passent un temps considérable à chercher des travaux proches de leurs intérêts de recherche, disséminés dans 47'000 revues scientifiques appartenant à quelques 6000 éditeurs différents. Cette étape est cependant incontournable dans tout projet de recherche afin de confronter de nouvelles idées à des solutions existantes, ainsi que pour l'acquisition de connaissance à propos d'un domaine spécifique. Dans cet article, une nouvelle méthode d'extraction de connexions entre les catégories des mots-clés d'articles scientifiques est proposée. Les limites de notre approche naïve héritée de la recherche exacte ont été soulignées dans Latard et al. (2017), et cet article fournit une amélioration qui s'attaque à ce problème. Notre recherche a pour but d'intégrer les relations sémantiques dans les moteurs de recherche scientifiques afin de les rendre plus intelligents. Effectivement, en fonction du nombre de résultats renvoyés, une requête plus raffinée / étendue pourrait alors être proposée à l'utilisateur.

FIG. 1 – Illustration de la logique générale de notre approche



**Approche Proposée.** Notre approche utilise BabelNet (Navigli et Ponzetto (2012)), une base de données fusionnant lexiques sémantiques (WordNet, VerbNet) et autres bases de données collaboratives (Wikipedia et autres données Wiki). Une requête pour un terme renvoie des "entrées de dictionnaire", des synonymes, des catégories ou des domaines. Cette base de connaissance est intégrée afin d'ajouter de l'information sémantique à partir de tous les mots-clés des articles de la base de données de littérature scientifique, Scilit<sup>1</sup>. Scilit contient à ce jour les métadonnées de plus de 97 millions d'articles. La Figure 1 illustre la logique principale de notre framework. La recherche exacte est l'approche naïve de notre framework qui prend des mots-clés sans préformatage et tente de faire une recherche exacte sur BabelNet. Ses limites sont rapidement atteintes lorsqu'un article comporte des mots-clés composés (plusieurs

1. <http://www.scilit.net> – développée par MDPI (<http://www.mdpi.com>)