

## BioMinT: the Research Assistant for Biological Text Mining

Luc Dehaspe<sup>1</sup>, Teresa K. Attwood<sup>2</sup>, Walter Daelemans<sup>3</sup>, Melanie Hilario<sup>4</sup>, Jee-Hyub Kim<sup>4</sup>, Jo Meyhi<sup>3</sup>, Alex L. Mitchell<sup>2</sup>, Johann Petrak<sup>5</sup>, Violaine Pillet<sup>6</sup>, Alexander K. Seewald<sup>5</sup>, Ioannis Selimas<sup>2</sup>, Anne-Lise Veuthey<sup>6</sup>, Marc Zehnder<sup>6</sup>

<sup>1</sup>PharmaDM, Leuven, Belgium, <sup>2</sup>School of Biological Sciences, University of Manchester, Manchester, United Kingdom, <sup>3</sup>University of Antwerp, Antwerp, Belgium, <sup>4</sup>University of Geneva, Geneva, Switzerland, <sup>5</sup>Austrian Research Institute for Artificial Intelligence (OFAI), Vienna, Austria, <sup>6</sup>Swiss Institute of Bioinformatics, Geneva, Switzerland

### Abstract

The goal of the BioMinT project is to develop a generic text mining tool that assists researchers with the extraction of relevant information from the worldwide collection of abstracts and papers, the ultimate information source to anyone trying to know more about a given DNA or protein sequence. The project is conducted by an interdisciplinary team that represents a unique combination of expertises and technologies. The biologists in the team are involved in the curation of biological databases (Swiss-Prot, PRINTS). As such they are ideally placed to provide feedback on the tool's efficiency (measured in reduction of literature screening time). They also identify and incorporate relevant biological resources (e.g., databases and ontologies). The developers in the team provide and integrate Natural Language Processing, Text Mining, and Data Mining components and adapt these components to the biomedical domain.

The core of the system is composed of an information retrieval module consisting in a meta-query engine wrapped around the PubMed server. To ensure a high recall of documents from Medline, the query is expanded with related terms. For that purpose, a new database has been developed, GPSDB<sup>(1)</sup> -Gene and Protein Synonyms DataBase- which collects gene/protein names, in a species specific way, from 14 main biological resources. A web-based search interface (freely accessible at [biomint.oefai.at](http://biomint.oefai.at)) gives access to the database: given a gene/protein name, it retrieves all synonyms for this entity and queries Medline with a set of user-selected terms.

The retrieved documents are then filtered, categorized, and ranked according to their relevance with regard to the query. A user interface provides control over each step of the query process.

Next, surviving documents are fed into the information extraction module. In this module the texts are parsed using adaptive natural language processing (NLP) techniques. The adaptation of the NLP module could be achieved by training a tokenizer and tagger on a biomedical corpus, and by adding a named-entity / concept tagger for biomedical concepts. Finally, from the parsed sentences, those are selected that deal with a user specified topic.

For the PRINTS application, the above steps have been combined in a unified system capable of taking a fingerprint, returning a set of relevant documents, extracting useful sentences and pertinent information from those sentences.

*The BioMinT project ([www.biomint.org](http://www.biomint.org)) is funded by the European Commission, contract-no. QLRI-CT-2002-02770 under the RTD programme "Quality of Life and Management of Living Resources".*

(1) V. Pillet, M. Zehnder, A.K. Seewald, A.-L. Veuthey, and J. Petrak. GPSDB: a new database for synonyms expansion of gene and protein names. Accepted in Bioinformatics.