# Bi²SoN – A Digital Library for Supporting Biomedical Research

Benjamin Köhncke
L3S Research Center
Appelstrasse 9a
30167 Hannover
Germany

koehncke@l3s.de

Sascha Tönnies
L3S Research Center
Appelstrasse 9a
30167 Hannover
Germany

toennies@l3s.de

Wolf-Tilo Balke
IFIS TU Braunschweig
Mühlenpfordtstrasse 23
38106 Braunschweig
Germany

balke@ifis.cs.tu-bs.de

## ABSTRACT
In the domain of biology a huge amount of different data sources is available. Therefore, information gathering and searching are challenging tasks. To avoid a manual assessment of all relevant data sources, their knowledge has to be integrated. The presented system focuses on all aspects needed for suitable data integration and retrieval for domain experts from the field of biology. The knowledge from different data sources is combined and further used for, e.g. synonym enrichment of the query term. The resulting prototype was presented to a group of domain experts who confirmed that the system delivers suitable results supporting the scientists by their literature search.

## Categories and Subject Descriptors
H.3.3 [**Information Systems**]: Information Storage and Retrieval – *Information Search and Retrieval.*

## General Terms
Management, Design, Human Factors

## Keywords
Architecture, Biology, Personalization, Digital Library

## 1. INTRODUCTION AND BACKGROUND
Nowadays, the directed information search is usually performed via the Internet. Independent of the field of interest a huge amount of online accessible information is available. Especially in the field of biology searching for relevant data is a challenging task. Biology is a wide domain with many different working areas. Each area has special data sources where the required knowledge is stored. In total there are more than 1000 different data sources storing biological knowledge.

Today the search for information in biology is performed by manually requesting all required data sources and manually combining the retrieved information. The combined information is subsequently used for a literature search in specialized digital libraries. The most widely-used library for the field of biomedical and life sciences information is MEDLINE (Medical Literature Analysis and Retrieval System Online) covering more than 21 million articles from around 4500 different journals. It is

compiled by the United States National Library of Medicine (NLM) and freely accessible via the PubMed interface.

The strength and weaknesses of commonly used search portals for biomedical information retrieval are compared in [1], namely PubMed, Web of Science, Scopus, and Google Scholar. The authors compared the content coverage and practical utility of these sources. They used example keyword searches to evaluate the usefulness and specific published articles to evaluate their utility in performing citation analysis. Whereas PubMed and Google Scholar offer free access to their search platform, Scopus and Web of Science are commercial products requiring access fees. The evaluation showed that PubMed is up-to-date due to frequent updates, reliable and easy to use. With GoPubMed [2] it exists a portal allowing users to explore PubMed search results with a hierarchically structured vocabulary for molecular biology, the Gene Ontology [3]. In GoPubMed an overview of the document abstracts according to the Gene Ontology is given. Therefore, the user is able to navigate through the abstracts by category. Furthermore, also more general ontology terms are shown, related to the original query that do not occur directly in the abstracts.

However, all approaches have in common that they only index the terms directly occurring in the documents. But most of the query terms have several synonyms frequently used in different documents. In [4] it was mentioned that for many proteins also different names were used even within the same article. Without taking these synonyms into account a huge amount of relevant literature is not included in the result set. Another aspect is that PubMed does not support a suitable ranking. Usually the result set is ordered by showing the newest documents first. Most of the queries deliver huge results sets making it impossible for the user to find the most important documents without a suitable ranking.

Furthermore, the user's relevance sensation is usually not only query dependent, but also related to his implicit background knowledge and the task he is currently working on. This knowledge cannot be expressed in the query. To summarize, the following additional aspects have to be considered in a retrieval system for biological literature:

- Query expansion using synonyms
- Suitable ranking functions
- Personalized access to information sources

In the course of the Bi²SoN project we will focus on these aspects to build an architecture supporting biomedical researchers with their literature search. Our aim is to support the scientist by

offering a portal with direct access to all required data sources. Beside the literature research this also includes an analysis of the required data sources for query expansion. We did a survey with domain experts from the field of infection research and found out that depending on the task the scientist is working on the different knowledgebases for query enrichment differ. We will introduce a framework integrating the knowledge from the required data sources by taking the search task and implicit knowledge of each individual user into account.

## 2. Bi²SoN ARCHITECTURE

During the project we work together with experts from the Helmholtz Centre for Infection Research (HZI). For the area of infection research it is important to understand bio-chemical reactions in cells. Therefore, a lot of mass-spectrometric experiments are accomplished leading to huge amounts of experimental data. In a first step, this data is analyzed using statistical approaches. An even more important step is the comparison of these results with other approaches known in literature. For example, it is possible that the same reactions are also known from other organisms. Hence, after the relevant entities have been determined from the experimental data, they are further used for literature search. Figure 1 gives an over view of our architecture.
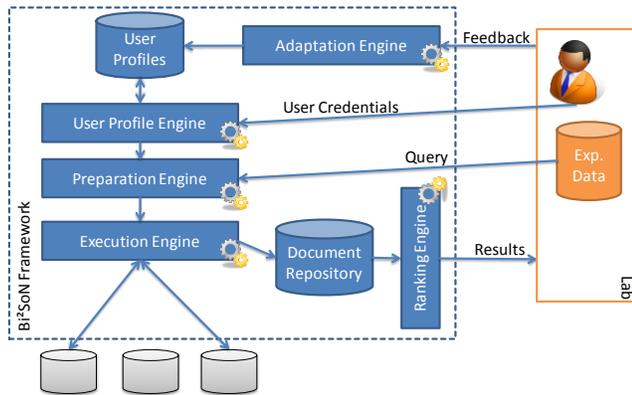


**Figure 1. Bi²SoN architecture**

The scientist analyzes the experimental data and determines the interesting entities. He identifies himself to the Bi2SoN system and submits a query term out of these entities. For each user a user profile is stored in our repository including information about preferred data sources and working tasks. The user profile engine is responsible for managing the user profiles. The profile and the query is analyzed by the preparation engine which is responsible for choosing adequate data sources. The resulting service set is transmitted to the execution engine which requests the included data sources. Each data source delivers information related to the query entity. The execution engine analyzes and merges the data from each data source and finally transmits the enriched query to the document repository. In a last step a suitable ranking function is chosen taking the user preferences defined in the profile into account. The ranked result set is delivered to the user who has the choice to give a relevance feedback to the system. The user feedback is used to adapt information in the profile, e.g. if the precision of the top-x documents is low another ranking function has to be used for this user.

We implemented the introduced architecture in a first prototype. Our document repository includes around 220000 documents from the PubMed Central corpus. Furthermore, we used external

knowledgebases to extract all genes and pathways included in each document. We combined the information from several sources, e.g. UniProt and KEGG, to allow for automatic synonym retrieval. Figure 2 shows a screenshot of our system which we presented to the domain experts from the HZI. The feedback was very positive since the system already retrieves good search results due to the integration of different data sources and synonym enrichment.
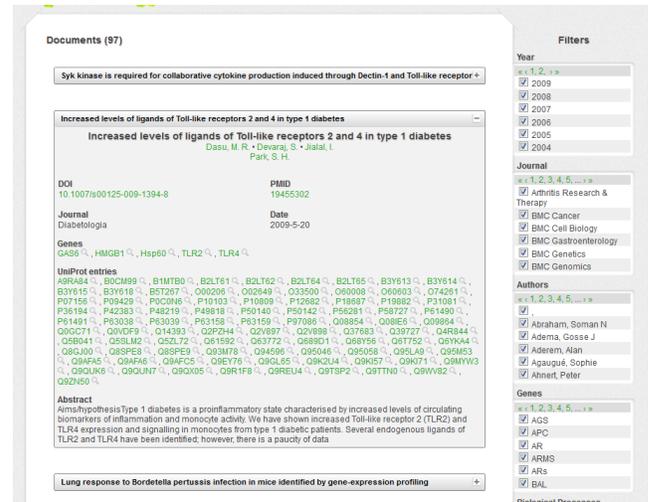


**Figure 2. Bi²SoN portal screenshot**

The user has the ability to filter the result set using bibliographic and semantic facets. The description of each document includes, among others, all genes and UniProt entries.

## 3. CONCLUSION

In the field of biology a huge amount of different data sources is available. To allow for high quality document retrieval an integration of this knowledge is essential. We have shown in a first prototype that the integration of different sources for synonym enrichment already leads to pleasant retrieval results. For future work we plan to further investigate the definition of the user profile. Beside the already included information about preferred data sources and working tasks, it would be interesting to model the implicit knowledge of a biologist. Usually he has an exact perception of what he is interested in which cannot be expressed in a query. Moreover, it is important to invent personalized ranking functions to further improve the user's search experience.

## 4. REFERENCES

[1]     M. E. Falagas et al., "Comparison of PubMed, Scopus, Web of Science, and Google Scholar: strengths and weaknesses.," *The FASEB journal*, vol. 22, no 2, 2008.

[2]     A. Doms, M. Schroeder, "GoPubMed: exploring PubMed with the Gene Ontology.," *Nucleic acids research*, vol. 33, no Web service issue, 2005.

[3]     M. a Harris et al., "The Gene Ontology (GO) database and informatics resource.," *Nucleic acids research*, vol. 32, no Database issue, 2004.

[4]     L. Hirschman et al., "Accomplishments and challenges in literature data mining for biology.," *Bioinformatics (Oxford, England)*, vol. 18, no. 12, 2002.