

# Extraction and Prediction of Biomedical Database Identifier Using Neural Networks towards Data Network Construction

**Hendrik Mehlhorn**

*Institute of Plant Genetics and Crop  
Plant Research, Germany*

**Uwe Scholz**

*Institute of Plant Genetics and Crop  
Plant Research, Germany*

**Matthias Lange**

*Institute of Plant Genetics and Crop  
Plant Research, Germany*

**Falk Schreiber**

*Institute of Plant Genetics and Crop  
Plant Research, Germany & Martin  
Luther University, Germany*

## EXECUTIVE SUMMARY

*Knowledge found in biomedical databases is a major bioinformatics resource. In general, this biological knowledge is represented worldwide in a network of thousands of databases, which overlap in content, but differ substantially with respect to content detail, interface, formats, and data structure. To support a functional annotation of lab data, such as protein sequences, metabolites, or DNA sequences, as well as a semi-automated data exploration in information retrieval environments, an integrated view to databases is essential. A prerequisite of supporting the concept of an integrated data view is to acquire insights into cross-references among database entities.*

## ***Extraction and Prediction of Biomedical Database Identifier Using Neural Networks***

*In this work, we investigate to what extent an automated construction of an integrated data network is possible. We propose a method that predicts and extracts cross-references from multiple life science databases and possible referenced data targets. We study the retrieval quality of our method and report on first, promising results.*

### **1. INTRODUCTION**

Bioinformatics is the field of science in which biology, computer science, and in particular information retrieval merge to form a single discipline. The ultimate goal of the field is to enable the discovery of new biological insights. The first step in this direction is already done. High throughput biotechnologies, like next generation sequencing, proteomics and metabolomics techniques produce a massive amount of data (Galperin & Fernandez-Suarez, 2012). But the data gathered in biology or medicine is as manifold as the biological research areas itself. If we will narrow down in this chapter the complex areas of biomedical research to molecular biology, bioinformatics attempts to model and interprets this data pathway: *genome, gene sequence, protein sequence, protein structure, protein function, cellular pathways & networks*, and *biomedical literature*. The first consequence of this revolution is the explosion of available data that biomolecular researchers have to harness and exploit (Roos, 2001) (e.g., as of March 2012, Genbank provides access to 150,000,000 DNA sequences<sup>1</sup> and in PubMed there are 2,400,000 research articles listed). The number of public available databases passed currently the number of high water mark of 1,200 (Galperin & Fernandez-Suarez, 2012).

The big players in this context are on the one hand companies like pharmaceutical or plant breeders on the other hand public or private financed research institute. Their role is either a data consumer or a data producer. In consequence there is a raising need for find, extract, merge, and synthesize information from multiple, disparate sources. Convergence of biology, computer science, and information technology will accelerate this multidisciplinary endeavor. The basic needs are formulated in Lacroix & Critchlow, 2003:

1. On demand access and retrieval of the most up-to-date biological data and the ability to perform complex queries across multiple heterogeneous databases to find the most relevant information.
2. Access to the best-of-breed analytical tools and algorithms for extraction of useful information from the massive volume and diversity of biological data.
3. A robust information integration infrastructure that connects various computational steps involving database queries, computational algorithms, and application software.

24 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: [www.igi-global.com/chapter/extraction-prediction-biomedical-database-identifier/77200](http://www.igi-global.com/chapter/extraction-prediction-biomedical-database-identifier/77200)

## Related Content

---

### Incremental Learning

Abdelhamid Bouchachia (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 1006-1012).

[www.irma-international.org/chapter/incremental-learning/10944](http://www.irma-international.org/chapter/incremental-learning/10944)

### Enhancing Web Search through Query Expansion

Daniel Crabtree (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 752-757).

[www.irma-international.org/chapter/enhancing-web-search-through-query/10904](http://www.irma-international.org/chapter/enhancing-web-search-through-query/10904)

### Count Models for Software Quality Estimation

Kehan Gao and Taghi M. Khoshgoftaar (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 346-352).

[www.irma-international.org/chapter/count-models-software-quality-estimation/10843](http://www.irma-international.org/chapter/count-models-software-quality-estimation/10843)

### Discovery of Protein Interaction Sites

Haiquan Li, Jinyan Li and Xuechun Zhao (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 683-688).

[www.irma-international.org/chapter/discovery-protein-interaction-sites/10894](http://www.irma-international.org/chapter/discovery-protein-interaction-sites/10894)

### Symbiotic Data Miner

Kuriakose Athappilly (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition* (pp. 1903-1908).

[www.irma-international.org/chapter/symbiotic-data-miner/11079](http://www.irma-international.org/chapter/symbiotic-data-miner/11079)