

CardioVINEdb: a data warehouse approach for integration of life science data in cardiovascular diseases

Benjamin Kormeier, Klaus Hippe, Thoralf Töpel and Ralf Hofestädt

Bioinformatics Department
Bielefeld University
Universitätsstraße 25
D-33501 Bielefeld
Germany
bkormeie@techfak.uni-bielefeld.de

Abstract: One of the major challenges in bioinformatics is to integrate and manage data from different sources as well as experimental microarray data and present them in a user-friendly format.

Therefore, we present CardioVINEdb, a data warehouse approach developed to integrate and explore life science data. The data warehouse architecture provides a platform-independent web interface that can be used with any common web browser. A monitor component controls and updates the data from the different sources to guarantee up-to-dateness. In addition, the system provides a visualization component for interactive graphical exploration of the integrated data based on networks of biological objects.

1 Introduction

Large amounts of high dimensional biological data are generated from different high-throughput experiments and from literature. The rapidly growing number of databases and data types poses the challenge of integrating the heterogeneous data, especially in biology. Currently there are about 1170 important molecular biology databases [1].

Thus, the challenge is to capture, model, integrate and analyze the data in a consistent way to provide a new and deeper insight into complex biological systems.

The huge quantity of information generated in life sciences is dispersed in many databases and repositories. Diverse integration approaches for molecular biological data sources have been developed. These systems are based on different data integration techniques. Several systems like Atlas [2], BioWarehouse [3], Columba [4], Systomonas [5] and Reactome [6] have been developed to integrate and present heterogeneous biological data. But most of these systems are not platform-independent and are implemented in different programming languages. Those systems take more time to install and make it difficult to determine whether data is up-to-date due to lack of logging information. Therefore, the key task is not only to integrate and manage the data scattered in different sources by creating a data warehouse, but also to create a method for visually representing the integrated data in a simple way to understand the underlying biological complexity with ease.

2 Design and Implementation

Based on the CardioWorkBench EU project we implemented a platform-independent data warehouse system that integrates multiple heterogeneous data sources into a local database enriched with protein microarrays from human smooth muscle cells that are related to cardiovascular diseases. Based on our VINEdb [8] information system we extended CardioVINEdb with more data sources, better data warehouse infrastructure including monitoring and microarray data. In addition, we upgraded the visualization components and web pages for better navigation and exploration. To ensure maximum up-to-dateness of the integrated data, we developed a data warehouse infrastructure including a monitor component. Furthermore, the common web-based user interface provides a visualization component that allows interactive exploration of the integrated data.

The CardioVINEdb system architecture consists of a 4-layer architecture that is illustrated in Figure 1. The source layer contains the multiple data sources BRENDA, EMBL, GO, IntAct, KEGG, MINT, OMIM, PubChem, SCOP, Transfac, Transpath and UniProt. In addition to the public available databases, we integrate experimental microarray data of human smooth muscle cells that are associated to cardiovascular diseases. Most of these databases provide parseable flat files that can be processed by our data warehouse infrastructure. A monitor component that is part of the integration layer controls the different data sources. It recognizes changes in the original sources and starts download if files changed. In a defined cycle the parser will be activated to start the ETL (Extraction-Transform-Load) process. ETL means that data is extracted from the source data, transformed into the data warehouse schema and loaded into the data warehouse. Data marts for specific analysis applications can easily be constructed by the database layer, i.e. the data warehouse.

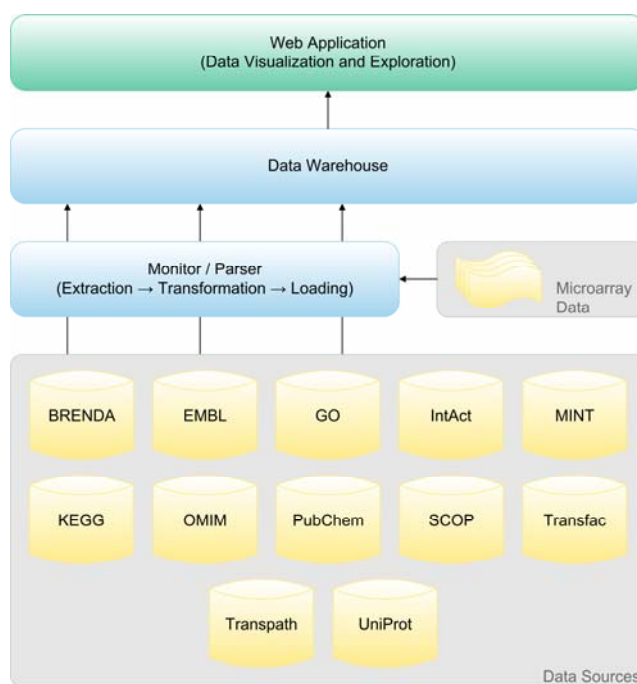


Figure 1: Schematic representation of the CardioVINEdb 4-layer system architecture from the original heterogeneous data sources to the web application layer.

The web-based graphical user interface of CardioVINEdb is implemented with JavaServer Pages (JSP) and runs on an Apache Tomcat web server. Each data entry has detailed information and a further link to the original data source. A general search engine allows the user to find information of interest spanning multiple domains, such as proteins, enzymes, genes, compounds etc. Additionally, each domain has its own specific search engine to find required information for research.

For better understanding the relationships between the biological object, the network-based visualization enables intuitive and comfortable exploration of the integrated data. The images or the graphical representations of the relationships between the entities of the data warehouse created using JUNG at runtime are dynamic and interactive allowing further exploration. JUNG (<http://jung.sourceforge.net/>) is a Java-based library that provides classes to describe graphs, nodes and edges with additional layout preferences. Thus, a graph is generated by JUNG according to the entities selected by the user. The system produces a PNG image file with the graphical visualization of biological objects in different domains and their linkage. Finally, this image is embedded in the HTML pages and displayed by the web browser. The dynamic component using a Java Applet works in the same manner, but in this case the graph is directly generated and displayed in the applet. For more interactive navigation and exploration the applet has a zoom function, different graph layouts and a picking function to move and select nodes within the graph. Therefore the applet is embedded in the HTML pages and can be displayed by the web browser if Java Runtime Environment is installed on the computer.

The database management for integrated data is realized in MySQL using Java Database Connectivity (JDBC). The core of the data warehouse infrastructure, with the name BioDWH [7], is completely implemented in Java, which ensures platform independence of the operating system. Therefore, it could be used separately from the web interface. BioDWH is a bioinformatics data warehouse software kit that integrates biological information from multiple public life science data sources into a local RDBMS. It provides up-to-date integrated knowledge, platform and database independence. This data warehouse infrastructure is available for interested scientific users as a SourceForge project (<http://sourceforge.net/projects/biodwh/>).

3 Summary

A major challenge in life sciences is the integration of heterogeneous data. Apart from the difficulty to facilitate the study of such data within the biological context, a fundamental problem is to represent and make the available knowledge accessible. CardioVINEdb provides integrated data from different popular life science databases and microarray data related to cardiovascular diseases from an EU project in a homogeneous web-based system. The system enables intuitive search of integrated life science data, simple navigation to related information as well as visualization of biological domains and their relationships. Equipped with a monitor component that updates the integrated data in defined update circles, it enables a way for presenting complex biological data in a very user-comprehensible manner. CardioVINEdb is available at <http://agbi.techfak.uni-bielefeld.de/CardioVINEdb/>.

4 Acknowledgements

This work is supported by Sixth Framework Programme Priority 1 Life Sciences, Genomics and Biotechnology for Health: „CardioWorkBench - Drug Design for Cardiovascular Diseases: Integration of in Silico and in Vitro Analyses“ (Proposal/Contract no.: PL 018671).

References

- [1] M. Y. Galperin and G. R. Cochrane. Nucleic Acids Research annual Database Issue and the NAR online Molecular Biology Database Collection in 2009. *Nucleic Acids Research*, 37(Database issue):D1-D4, 2009.
- [2] S. P. Shah, Y. Huang, T. Xu, M. M. S. Yuen, J. Ling, and B. F. F. Ouellette. Atlas - a data warehouse for integrative bioinformatics. *BMC Bioinformatics*, 6:34, 2005.
- [3] T. J. Lee, Y. Pouliot, V. Wagner, P. Gupta, D.W. J. Stringer-Calvert, J. D. Tenenbaum, and P. D. Karp. BioWarehouse: a bioinformatics database warehouse toolkit. *BMC Bioinformatics*, 7:170, 2006.
- [4] S. Trißl, K. Rother, H. Müller, T. Steinke, I. Koch, R. Preissner, C. Frömmel, and U. Leser. Columba: an integrated database of proteins, structures, and annotations. *BMC Bioinformatics*, 6:81, 2005.
- [5] C. C. Choi, R. Münch, S. Leupold, J. Klein, I. Siegel, B. Thielen, B. Benkert, M. Kucklick, M. Schobert, J. Barthelmes, C. Ebeling, I. Haddad, M. Scheer, A. Grote, K. Hiller, B. Bunk, K. Schreiber, I. Retter, D. Schomburg, and D. Jahn. SYSTOMONAS - an integrated database for systems biology analysis of *Pseudomonas*. *Nucleic Acids Research*, 35(Database issue):D533-D537, 2007.
- [6] G. Joshi-Tope, M. Gillespie, I. Vastrik, P. D'Eustachio, E. Schmidt, B. de Bono, B. Jassal, G.R. Gopinath, G.R. Wu, L. Matthews, S. Lewis, E. Birney and L. Stein. Reactome: a knowledgebase of biological pathways. *Nucleic Acids Research*, 33(Database issue):D428-D432, 2005.
- [7] T. Töpel, B. Kormeier, A. Klassen and R. Hofestädt. BioDWH: A Data Warehouse Kit for Life Science Data Integration. *Journal of Integrative Bioinformatics*, 5(2):93, 2008.
- [8] S. Hariharaputran, T. Töpel, B. Brockschmidt and R. Hofestädt. VINEdb: a data warehouse for integration and interactive exploration of life science data. *Journal of Integrative Bioinformatics*, 4(3):63, 2007. Online Journal: http://journal.imbio.de/index.php?paper_id=63