**Biosphere: the interoperation of web services in microarray cluster analysis**

Kei-Hoi Cheung[1,2,*], Remko de Knikker[1], Youjun Guo[1], Guoneng Zhong[1], Janet Hager[3,4], Kevin Y. Yip[5], Albert K.H. Kwan[5], Peter Li[6], and David W. Cheung[5]

[1]Center for Medical Informatics, [2]Department of Genetics, [3]Keck Biotechnology Resource Laboratory, [4]Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, USA, [5]Department of Computer and Information Systems, Pokfulam Road, University of Hong Kong, Hong Kong, [6]School of Computing Science, University of Newcastle, NE1 7RU, UK.

## ABSTRACT

**Summary**: The growing use of DNA microarrays in biomedical research has led to the proliferation of analysis tools. These software programs address different aspects of analysis (e.g., normalization and clustering within and across individual arrays) as well as extended analysis methods (e.g., clustering, annotation and mining of multiple data sets). Microarray data analysis therefore typically requires the interoperability of multiple software programs involving different analysis types and methods. Such interoperation is often hampered by the heterogeneity involved in the software tools (which may function by implementing different interfaces and using different programming languages). To address this problem, we employ the SOAP-based web service approach that provides a uniform programmatic interface to these heterogeneous software components. To demonstrate this approach in the microarray context, we have created a web server application, **Biosphere**, which interoperates a number of web services that are geographically widely distributed. These web services include a **clustering web service**, that is a suite of different clustering algorithms for analyzing microarray data, **XEMBL** developed at EBI for retrieving EMBL sequence data, and three gene annotation web services: **getGO**, **getHAPI**, and **getUMLS**. While **getGO** allows retrieval of Gene Ontology (GO) annotation, the other two web services retrieve annotation from the biomedical literature that is indexed based on the Medical Subject Headings (MeSH) terms. With these web services, the application allows the users to do the following: i) cluster gene expression data using seven different algorithms, ii) visualize the clustering results statistically grouped in color, and iii) retrieve sequence, annotation and citation data for the genes of interest.

**Availability**: The web server application and its web services described in WSDL can be accessed at: http://rook.cecid.hku.hk:8280/BiosphereServer or https://rook.cecid.hku.hk:8443/BiosphereServer (this works for browsers such as Internet Explorer 6 or higher that supports secure data encryption)

**Contact**: kei.cheung@yale.edu

## INTRODUCTION

Interpretation of microarray gene expression data utilizes many different types of computational and statistical analyses methods (e.g., normalization, clustering, etc). Despite the fact that a large repository of software tools have been developed and made

---

[*] To whom correspondence should be addressed.

available to the microarray community (Dudoit et al. 2003), researchers are confronted with the following problems.

1. Identifying and locating available tools may be difficult since they are widely scattered and developed by geographically disperse groups.
2. Some types of microarray analysis may require large amounts of memory or computational power especially when large quantities of data are involved. Researchers may not have adequate local computational resources to perform such analyses.
3. Tools available by download are often written in different languages and may not be cross-platform-compatible.
4. Some tools provide a web server interface that allows users to submit their data over the Internet for processing. While this type of interface is friendly to biological users, it is difficult for bioinformatics developers to access programmatically.
5. It is often difficult to orchestrate a specific set of bioinformatics tools from diverse sources for a particular application due to their incompatibility and heterogeneity.

To address these problems, we employ web services, which are implemented based on the XML-based Simple Object Access Protocol (SOAP), as a uniform model to provide a common programmatic interface to heterogeneous software applications in the bioinformatics context. There have been a growing number of SOAP web services available for bioinformatics use. This includes XEMBL (Wang et al. 2002), Soaplab (http://industry.ebi.ac.uk/soaplab), DDBJ Soap Server (http://xml.ddbj.nig.ac.jp/soapp.html), and so on. As pointed out by Stein (Stein 2002), web services have the potential of unifying the "bioinformatics nation". We demonstrate how web services can be integrated to create a single-entry portal that is broadly accessible to researchers for data analysis involving multiple programs and databases.

**WEB SERVICES APPLICATION**
A wide variety of clustering algorithms have been used to cluster genome-wide expression patterns (Eisen et al. 1998). In addition to traditional methods, new approaches to clustering gene expression data have been developed (e.g., HARP (Yip et al. 2004)). It has become important to evaluate and compare these methods. In (Datta et al. 2003), a number of algorithms (including hierarchical, K-means, Diana, fanny, and model-based methods) were compared in terms of their performance in gene clustering. To access and use these clustering programs may not be straightforward since they are scattered and available in different formats (some of them are platform-specific). For example, Eisen's cluster program (http://rana.lbl.gov/EisenSoftware.htm) can only be run on Windows. EPCLUST (http://ep.ebi.ac.uk/EP/EPCLUST/) is a web server application allowing clustering, visualization and analysis of gene expression data, but it lacks a programmatic interface, making automatic interoperation with other programs difficult. To demonstrate how these issues are addressed, we have built a web server application, **Biosphere**, which interoperates a number of web services. **Biosphere** allows different algorithms to be used to cluster gene expression datasets uploaded by individual users via the Web. The algorithms include both hierarchical and non-hierarchical methods, which

are programmatically accessible through a "cluster" web service. In addition, we use the following web services to retrieve annotation and sequences of the clustered genes.

1. **getHAPI:** This is the web service counterpart of HAPI (Masys et al. 2001), which extracts MeSH terms represented by the Unified Medical Language System (UMLS) concept IDs from the biomedical literature.
2. **getUMLS:** This web service builds a hierarchy of UMLS concepts for a set of concept IDs.
3. **getGO:** This web service retrieves genome annotation (human only) from the Gene Ontology database (http://www.geneontology.org/).
4. **XEMBL:** This web service was developed at EBI (http://www.ebi.ac.uk), which allows retrieval of sequence data from EMBL.

Extensibility is an important property of web-service based analyzing systems. By defining standard service interfaces through Web Service Description Language (WSDL), new algorithms and modules (e.g. data preprocessing and charting) can easily be implemented by any parties and added to an existing system. Additional data mining algorithm implementations can be plugged into the system easily. The clustering and annotation web services are implemented and deployed using Axis that is the Java implementation of SOAP on Apache. The underlying programs are written in Java. To increase performance, the clustering programs run on a two-node cluster (with each node being an Intel Pentium 4 processor with 1 GB memory and 40 GB hard disk space). In addition, the system notifies the users by e-mail when clustering is done. This feature is especially useful for clustering large datasets. For security, Biosphere provides the option of using the Secure Sockets Layer (SSL) to encrypt data transmitted over the Internet.

Fig. 1 shows the colored cluster output display generated based on the Diana algorithm asking for genes (not samples) to be clustered into six groups based on expression ratio (together with the GO and HAPI annotation). This represents a subset of data generated from a human microarray experiment using 2 µgmRNAfrom HeLa and HL60 cell lines that was amino-allyl indirect cy-dye labeled using oligo dT primer and reverse transcription. The array data input for Biosphere are three technical hybridization replicates (ITIK_HUCOMP3V2C-319U, ITIK_HUCOMP3V2C-319L, ITIK_HUCOMP3V2C-320U). The array used (HUCOMP3V2C-300-400) is a composite array of both cDNA and oligonucleotide targets. The color difference of a gene across different samples is totally dependent on the expression values of the gene in the input data file. A perfect cluster is a group of genes that have zero distances from each other, i.e., in each sample, all the genes have the same expression value.

## REFERENCES

Datta, S. and S. Datta (2003). Comparisons and validation of statistical clustering techniques for microarray gene expression data. Bioinformatics **19**(4):459-466.

Dudoit, S., R. Gentleman and J. Quackenbush (2003). Open source software for the analysis of microarray data. Biotechniques **Suppl**:45-51.

Eisen, M., P. Spellman, P. Brown and D. Botstein (1998). Cluster analysis and display of genome-wide expression patterns. PNAS **95**(25):14863-8.

Masys, D., J. Welsh, J. L. Fink, M. Gribskov, I. Klacansky and J. Corbeil (2001). Use of keyword hierarchies to interpret gene expression patterns. Bioinformatics **17**(4):319-26.

Stein, L. (2002). Creating a bioinformatics nation. Nature **417**:119-20.

Wang, L., J. Riethoven and A. Robinson (2002). XEMBL: distributing EMBL data in XML format. Bioinformatics **18**(8):1147-8.

Yip, K.Y., D.W. Cheung, M.K. Ng and K.H. Cheung (2004). Identifying projected clusters from gene expression profiles. IEEE International Symposium on Bio-Informatics and Biomedical Engineering (BIBE), pp. 259-66.

**Fig. 1.** Colored display of Diana cluster output together with GO annotation (top right), HAPI results (middle right), and XEMBL sequences (bottom right) for the selected genes indicated by the yellow radio buttons next to them.