

# PROM-OOGLE – Data Mining and Integration of On-line Databases to Discover Gene Promoters

Dean Cheng, John Sheldon  
Department of Computing Science  
University of Alberta, Canada

{dcheng, sheldon}@cs.ualberta.ca

Marcelo Marcet-Palacios  
Department of Medicine  
University of Alberta, Canada

marcelo@ualberta.ca

Osmar R. Zaiane  
Department of Computing Science  
University of Alberta, Canada

zaiane@cs.ualberta.ca

## ABSTRACT

The vast number of on-line biological and medical databases available can be a great resource for biomedical researchers. However, the different types of data and interfaces available can be overwhelming for many biomedical researchers to learn and make effective use of. Moreover, the available resources lack needed integration. Here we focus on an important task in medical research: to provide researchers with promoter analysis for a given gene. PROM-OOGLE is a web based data mining tool that provides a means for researchers to take a gene name of interest and obtain its promoter sequence in return after automatic integration of text databases. Additionally, the program is capable of returning multiple promoters from different genes allowing researchers to study how promoters regulate genes. This tool facilitates the process of acquiring information on a promoter and may lead to interesting discoveries.

## Categories and Subject Descriptors

H.2.5 [Information Systems]: Heterogeneous Databases

## Keywords

Biomedical database integration, gene promoter analysis

## 1. INTRODUCTION

The permeation of computer science into all aspects of research is growing every day. None as much as biology where computers are aiding in the acquisition of meaningful results from data that was never thought possible. The field of bioinformatics has become a powerful area in which a tremendous potential for great research is being seen. As a result, the vast amount of on-line biological and medical databases available can be a great resource for biomedical researchers. Examples of such databases include sequence databases such as GenBank [1], literature databases such as MedLine [2], chemical databases such as PubChem [3], transcription factor databases such as TESS [4]. To fully take advantage of these databases, one has to have the proper biomedical background as well as familiarities with the different interfaces. This can be overwhelming for many medical

researchers. Moreover, the interfaces of the different on-line databases are neither standardized nor designed for integration. When the output of a query into a first database or part of it is required as input for a search in a second database, the operation needs to be done manually. This process is laborious, time consuming and prone to errors. To alleviate the above-mentioned difficulties, text mining and integration tools can be used to automate tedious searching procedures.

Developing data mining and text mining tools for biological applications is a growing field with vast potential [5]. There are many types of analysis that can be done across multiple types of data and as a result, there are many types of tools being developed. Tools range across from whole genome sequence analysis such as Vista Genome Browser [6] to information extraction from text for gene-gene/gene-disease associations such as MedMiner [7] and MedGene [8]. Our system, PROM-OOGLE, is a web content mining tool that tries to integrate sequence data and gene-promoter associations in order to help medical researchers with the important task of finding promoters of a gene. The unique functionality that PROM-OOGLE has to offer is to compare promoter sequences within a given species and to improve our understanding of promoter structure and similarities among genes that are co-regulated in similar ways.

In essence, PROM-OOGLE is an integrator that can facilitate researchers move through tedious work more efficiently and with less uncertainty with the validity of their results. It allows researchers to generate meaningful biological hypotheses with greater ease that can be tested in laboratory experiments with the aim of understanding better how gene expressions work and to eventually develop methods to counteract health problems. It is a knowledge discovery process because the tool combines parts of individual database outputs from different resources in order to generate new information that is not explicitly present in any individual database and presents it in a manner to support understanding and discovery of new knowledge regarding genes and their promoters.

## 2. MOTIVATION

### 2.1 Background Knowledge

As humans, we are essentially made up of cells. These cells contain 23 pairs of chromosomes, which are compact intertwined molecules called deoxyribonucleic acids (DNA). A DNA strand is composed of linearly linked nucleotides that are subsequently linked with one of four bases: adenine (A), thymine (T), guanine (G), and cytosine (C). On each strand of DNA, in specific positions, are genes, which contain hereditary information. Genes

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SAC'10, March 22-26, 2010, Sierre, Switzerland.

Copyright 2010 ACM 978-1-60558-638-0/10/03...\$10.00.





Fig. 8. Output from TESS database (the framed columns are targeted and extracted).

### 3.3 Describing PROM-OOGLE user interface

As hinted above, the output and input interfaces for NCBI and TESS databases are completely transparent to the user who may ignore their existence. These interfaces are managed automatically by PROM-OOGLE and only its web-based GUI is presented to the user. Figure 9 shows a snapshot of the user interface of the web based mining tool. The user interface is designed to be user-friendly for biomedical researchers. Navigating through the PROM-OOGLE website is plain and simple while it still provides various functionalities.



Fig. 9. PROM-OOGLE initial search interface (left); First gene search results (right).

Figure 9 shows, on the right, the display of returned gene matches that the user can select from after the first access to NCBI. The search for a gene name can return many hits and as such the user must ensure that PROM-OOGLE is searching the proper gene that the user has requested. The results of the gene search show the name of a possible subfamily of the gene, the short version of the name and the chromosome number. The user can then select a link to run PROM-OOGLE on the gene or can select the map element to locate the gene in the NCBI database if further information is needed pertaining to the gene. It is more important to verify that the promoter sequence returned by PROM-OOGLE is correct since wrong promoter sequences would definitely give rise to false results. During development and initial experiments, several gene names such as NOS2A, PI3K and CD8A have been verified to return biologically correct results with known correct promoter sequence.

Results from the data mining tool are returned in two formats: in the form of a gene sequence with the different bindings and directions; and in tabular form listing the promoters that can be sorted by different attributes interactively as the user sees fit. Figure 10 shows the promoter binding locations for the gene along with information pop down tables that provide more useful information regarding the gene and its binding site location (pop-up on Figure 10). The link that can be taken directly from the

name of the gene takes the researcher directly to the NCBI website which allows the user to access more information pertaining to that specific gene. This is very useful for the researcher because of the quick and simple access to pertinent information. Details of a promoter binding location, which include the sequence name, starting location, direction, length, and the sequence composition, can be displayed at will (Figure 11). It also includes the useful visual information pertaining to the exact placement of the sequence on the gene and its directions.

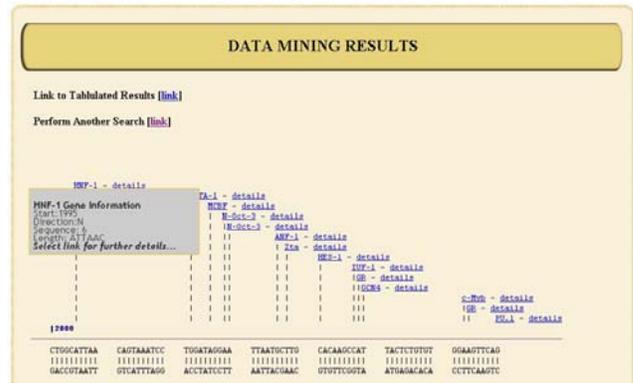


Fig. 10. Data mining results and the Basic gene information in pop down window.

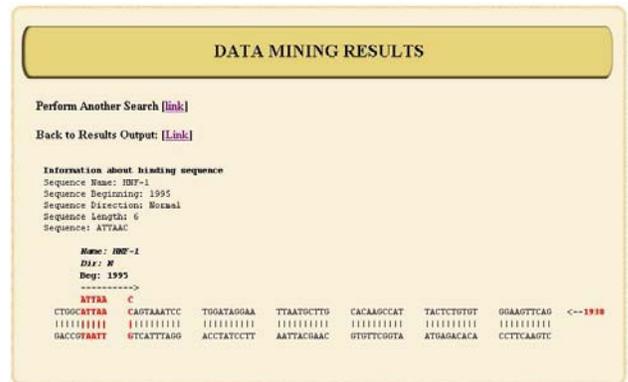


Fig. 11. Complete Gene details

ID	Factor	Beginning	Direction	Length	Sequence	Promoter
0	HNF-1	6	N	6	ATTAAC	<a href="#">[return]</a>
1	GATA-1	22	N	6	GGATAG	<a href="#">[return]</a>
2	MCFB	26	R	6	AGGAAT	<a href="#">[return]</a>
3	N-Oct-3	29	N	7	MATWAAT	<a href="#">[return]</a>
4	N-Oct-3	30	R	7	ATTWATR	<a href="#">[return]</a>
5	ANF-1	36	R	10	KTTGCACAA	<a href="#">[return]</a>
6	Zta	38	R	7	TTGCACA	<a href="#">[return]</a>
7	HES-1	41	R	6	CACAAG	<a href="#">[return]</a>
8	IUF-1	48	R	6	CAATAC	<a href="#">[return]</a>

Fig. 12. Tabular form of PROM-OOGLE results

The results in tabular form (Figure 12) can be sorted by any of the attributes. This provides the user with the ability to easily find information in the table. For example, if the user wished to find all sequences of a certain length, the user could click on "length", the list would be sorted by length, and the user would simply need to scroll down the list to the point of where the sequence length is of the desired number.

#### 4. POSSIBLE EXTENSIONS

Despite the functional correctness of our application, improvements to the tool are always possible and recommendations were collected after a first round of use by bio-medical researchers. Possible improvements include user interface functionalities that would help the user navigate through the information provided and application functionalities that would aid the user in making more informed decisions. Other improvements are more related to overall correctness of our application and require extensive user testing to ensure accuracy.

Another extension of PROM-OOGLE is to include a local database that could be used to include proteins of interest in different medical areas and alert users if the search results match some proteins in the local database. In addition, users would be able to save searches locally on their computers so that searches need not be repeated.

A final extension is the ability to run multiple searches simultaneously with results being stored for further review at a later time. This would also give the ability to perform multiple sequence alignment to find out if there are common transcription factors that are shared among a family of genes.

A concern of note is the running time of PROM-OOGLE. This is, however due to the running times of NCBI and TESS and as such we are limited to the time constraints they impose on our program.

#### 5. CONCLUSION

PROM-OOGLE is a web content mining tool and online database integrator for promoter analysis that can facilitate researchers move through tedious work more efficiently and with less uncertainty with the validity of their results by combining query results from different existing on-line databases. It allows researchers to generate meaningful biological hypotheses with greater ease that can be tested in laboratory experiments with the aim of understanding better how gene expressions work and to eventually develop methods to counteract health problems.

Even though software tools like Vista Genome Browser [6] are capable of performing phylogenetic analysis they lack the functionality of PROM-OOGLE to compare promoter sequences within a given species and improve our understanding of promoter structure and similarities among genes that are co-regulated in similar ways, therefore making PROM-OOGLE an exciting new tool for researchers.

This project is still in its first phase and there are numerous other functionalities that will be included in future developments. Preliminary feedbacks from bio-medical researchers have been very positive. This web-based tool is a promising aid to those who perform research in the area of biology and medicine.

#### 6. REFERENCES

- [1] GenBank : <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Nucleotide>.
- [2] MEDLINE: <http://medlineplus.gov/>
- [3] PubChem: <http://pubchem.ncbi.nlm.nih.gov/>
- [4] Jonathan Schug and G. Christian Overton, TESS: Transcription Element Search Software on the WWW, Technical Report CBIL-TR-1997-1001-v0.0, Computational Biology and Informatics Laboratory, School of Medicine, University of Pennsylvania, 1997, URL: <http://www.cbil.upenn.edu/tess>
- [5] D. Page and M. Craven, (2003). Biological applications of multi-relational data mining. *SIGKDD Explorations*, **5**(1): 69-79.
- [6] Frazer KA, Pachter L, Poliakov A, Rubin EM, Dubchak I. VISTA: computational tools for comparative genomics. *Nucleic Acids Res.* 2004 Jul 1;32 (Web Server issue):W273-9, URL: <http://pipeline.lbl.gov/cgi-bin/gateway2>
- [7] L. Tanabe, U. Scherf, L.H. Smith, J.K. Lee, L. Hunter and J.N. Weinstein, (1999). MedMiner: an Internet text-mining tool for biomedical information, with application to gene expression profiling. *BioTechniques Dec.* **27**:1210-1217
- [8] Y. Hu, L.M. Hines, H. Weng, D. Zuo, M. Rivera, A. Richardson, and J. LaBaer, (2003). Analysis of genomic and proteomic data using advanced literature mining. *Journal of Proteome Research*, **2**: 405-412
- [9] Wheeler, D.L., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Church, D.M., DiCuccio, M., Edgar, R., Federhen, S., Helmberg, W., Kenton, D.L., Khovayko, O., Lipman, D.J., Madden, T.L., Maglott, D.R., Ostell, J., Pontius, J.U., Pruitt, K.D., Schuler, G.D., Schriml, L.M., Sequeira, E., Sherry, S.T., Sirotkin, K., Starchenko, G., Suzek, T.O., Tatusov, R., Tatusova, T.A., Wagner, L. and Yaschenko, E. (2005) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.* **33** (Database issue):D39-45. NCBI Database: <http://www.ncbi.nlm.nih.gov>.
- [10] M. Marcet-Palacios, K. Graham, C. Cass, A.D. Befus, I. Mayers, and M.W. Radomski, (2003). Nitric Oxide and Cyclic GMP Increase the Expression of Matrix Metalloproteinase-9 in Vascular Smooth Muscle. **307**:429-436
- [11] Minematsu N, Nakamura H, Tateno H, Nakajima T, Yamaguchi K. (2001) Genetic polymorphism in matrix metalloproteinase-9 and pulmonary emphysema. *Biochem Biophys Res Commun.* Nov 23;289(1):116-9
- [12] NCBI-Entrez: <http://www.ncbi.nlm.nih.gov/Database/>
- [13] NCBI Map Viewer : [http://www.ncbi.nih.gov/map\\_search.cgi?taxid=9606](http://www.ncbi.nih.gov/map_search.cgi?taxid=9606)
- [14] TESS web form: <http://www.cbil.upenn.edu/cgi-bin/tess/tess?RQ=SEA-FR-QueryS>