

# Visualizing and Refining Connectivity Map Query Results

November 29, 2013

Proposal for a  
Thesis in the Field of  
Biotechnology

In Partial Fulfillment of the Requirements  
for a Master of Liberal Arts Degree

Harvard University  
Extension School

Theodore Natoli  
51 Lewis Avenue, Apartment 3  
Arlington, MA 02472  
857-498-1946  
`tnatoli@fas.harvard.edu`

Proposed Start Date: November 29, 2013  
Anticipated Date of Graduation: November 29, 2013  
Thesis Director: Aravind Subramanian

## Abstract

The Connectivity Map (CMap) is a database of gene expression signatures obtained from experiments in which cultured human cells are treated with pharmacologic and genomic perturbagens. A typical use case of this database is for a researcher to query with a signature of a cell state of interest and use the matching perturbagens to develop a functional hypothesis for follow-up. Current pattern matching algorithms that perform CMap queries suffer from a universal weakness – the enormous size and richness of signatures in CMap means that a query typically generates hundreds of strong connections. These connections are hard to distinguish, thereby making prioritization difficult. I hypothesize that one mode of prioritization is to highlight query results that are highly interconnected amongst themselves over singletons. The goal of this work is to provide a web-based tool for implementing an interconnectivity-based method of query result refinement and for visualizing CMap query results in a graph layout.

## 1 The Research Problem

The CMap database is a compendium of gene expression signatures resulting from the treatment of cultured human cells with small molecule compounds (CP), short hairpin RNAs (shRNA), or over-expression constructs (OE). The utility of the CMap database is that of a gene expression search engine. Users are able to pose questions about relationships between cellular states and formulate hypotheses based on similarities or differences in the states’ gene expression signatures.

Hypotheses are generated by posing search queries into the database and examining the query results. A CMap query is a focused question in which a user inputs a gene expression signature, called the query, and computes the similarity, or connectivity, between his/her query and other signatures in

the database. Positive connectivity indicates that two signatures' expression changes are similar and vice versa. Researchers can use CMap to find connections between signatures within or external to the database. Hypotheses may be in the form of "the shRNA knockdown of gene X connects to shRNA knockdown signatures of pathway Y members, so X is probably a member of Y." Or perhaps "the signature of compound Z connects to the knockdown signature of gene X, so perhaps X is the target of Z".

Lamb et al demonstrated a more directly therapeutically relevant use of the original incarnation of CMap when they discovered that the signature of sirolimus connected strongly to a signature of dexamethasone sensitivity. Dexamethasone is a treatment for acute lymphoblastic leukemia (ALL), but many patients eventually become resistant (Tissing, Meijerink, den Boer, & Pieters, 2003). The CMap connection between sirolimus and dexamethasone sensitivity suggested that sirolimus might be effective in reversing resistance in ALL patients who had become resistant to dexamethasone. A follow-up experiment confirmed that sirolimus conferred dexamethasone sensitivity to CEM-c1 cells, a previously dexamethasone-resistant cell line (Lamb et al., 2006).

The CMap database contains over 400,000 signatures spanning over 70 cell types. Because of the large size of the CMap database, interpreting and prioritizing query results has become a difficult task. For example, accepting only the top one percent of connections yields nearly 4,000 signatures. Follow-up on such a large number of primary hits is nearly impossible in

most cases. I propose that within a set of initial query results, there will frequently exist a set or sets of signatures that are more tightly interconnected with themselves than with other signatures. These interconnected sets are more likely to be indicative of robust biological signal and should therefore be prioritized over other singleton connections. The goal of this work is to build a web-based tool to implement an algorithm to identify subsets of high interconnectivity with lists of initial query results and to visualize the relationships between these subsets in a graph layout. I propose that this tool will be useful in refining initial CMap query results into smaller, more actionable lists of connections that can be further investigated in secondary assays.

## 2 Key Terms

1. gene expression profiling:
2. gene set enrichment analysis (GSEA):
3. Kolmogorov-Smirnov (KS) statistic:
4. graphical model:
5. CMap connection:

## 3 Background

### 3.1 Gene Expression Profiling

Gene expression profiling is the simultaneous measure of the RNA transcript levels of many genes within a cell or group of cells. These measurements can help to provide insight into the cellular state or states of the cells in question. For example, if many cell-cycle genes are observed to be active, it could suggest that the cells are actively dividing. Conversely, if many apoptotic genes are active, the cells might be dying. Frequently, the goal of gene expression profiling is to identify genes that are differentially regulated between one or more sets of conditions. For example, one might measure expression in cells that have and have not been treated with a drug of interest, and then compare the resulting expression profiles to identify genes that are substantially up- or down-regulated in the treated cells relative to the untreated. Current technology, such as the microarray, allow for many such gene expression experiments to be run in parallel, enabling the comparative analysis of hundreds or thousands of expression profiles corresponding to an equal number of experimental conditions. Similarly, expression profiling can be used to identify genes differentially regulated between disease and normal states. van't Veer et al used gene expression profiling to identify a set of genes, that were predictive of breast cancer metastasis (van 't Veer et al., 2002). Because of its ability to identify such signatures, gene expression profiling is a powerful and often-used tool in contemporary biology.

### 3.2 Gene Set Enrichment Analysis (GSEA)

Gene Set Enrichment Analysis (GSEA) is an analytical approach designed to extract biological insight from gene expression data (Subramanian et al., 2005). It leverages groups of genes, called gene sets, that share some biological commonality (i.e. members of a cellular signaling pathway) and computes their enrichment, or trend towards the top or bottom, of a ranked list of genes generated by comparing expression profiles across two experimental classes (i.e. tumor vs. normal). For example, one might define many sets of genes, each corresponding to a cellular pathway. One could then rank-order all genes by their differential expression when comparing profiles of tumor vs. normal samples. Lastly, one could compute the enrichment of each pathway in the rank-ordered list to attempt to identify pathways that might be active in the particular tumor in question.

Mechanically, GSEA computes a Kolmogorov-Smirnov (KS) statistic when comparing a given gene set to a given ranked list (Hollander & Wolfe, 1975). Effectively, this amounts to walking down the ranked list, increasing a running-sum statistic when we encounter a gene in the gene set and decreasing it when we encounter genes not in the gene set. The magnitude of the increment depends on the correlation of the gene with the phenotype. The enrichment score is the maximum deviation from zero encountered in the random walk (Subramanian et al., 2005). GSEA has been used extensively for identifying coherent sets of genes that are collectively modulated under certain disease and/or experimental conditions. In fact, a

GSEA software suite and an accompanying online database exist to facilitate comparisons between novel and curated gene sets. Cite MsigDB and GSEA.

### **3.3 The Connectivity Map**

The Connectivity Map (CMap) is a database containing the gene expression signatures resulting from treating cultured cells with various chemical and genomic perturbations (Lamb et al., 2006). The purpose of CMap is to serve as a lookup table of functional annotation. These annotations might be derived by comparing signatures within the CMap database or by querying the database with externally generated signatures. The database itself can be thought of as a large matrix where each row is a gene and each column is an experiment in which a particular perturbagen was profiled under a given set of conditions (i.e. cell context, dose, treatment time, etc). The values in the matrix are differential expression measures generated by comparing the expression levels of the genes across perturbed and control states. Thus, each column of the matrix can be thought of as a given perturbagen’s expression signature.

#### **3.3.1 Computing Connections in the Connectivity Map**

A primary use of the CMap database is to compare the signatures of different perturbations and assess their similarity. Perturbagens that, when used to treat cultured cells, result in similar gene expression consequences will yield similar CMap signatures. Such signatures are said to be positively connected



in the CMap sense. Conversely, perturbagens that elicit inversely-related expression consequences are said to be negatively connected. For a given query signature  $Q$  and a reference signature  $R$ , the weighted connectivity score (WTCS) is computed by computing and integrating two KS statistics, one each for the  $n$  most up- and down-regulated genes in  $Q$ . The algorithm proceeds as follows:

1. Order  $Q$
2. compute  $ES_{up}$  as the enrichment of the  $n$  most up-regulated genes in  $R$
3. compute  $ES_{dn}$  as the enrichment of the  $n$  most down-regulated genes in  $R$
4. compute WTCS as
  - (a) 0 if  $ES_{up}$  and  $ES_{dn}$  share the same sign
  - (b)  $(|ES_{up}| + |ES_{dn}|/2)$ , where the resulting WTCS is given the sign of  $ES_{up}$

WTCS will be positive for signatures that are positively related and negative for those that are inversely related.

A common CMap use case is to select a given query signature  $Q$  from the database and compute its similarity to all other signatures. The remaining signatures can be ranked according to their connection strength with  $Q$ . The connections can be used to gain insight and form hypotheses about  $Q$ . For example, if  $Q$  is a signature of a novel compound and it connects strongly to signatures of compounds of a known pharmacological class, one

might hypothesize that the novel compound is also a member of this class. Similarly, if Q connects strongly to the knockdown signature of gene G, one might hypothesize that G is the novel compound’s target. Q need to be a signature from the CMap database. For instance, it might be the signature of some disease and one might seek connections to genes whose modulation could be causing the disease or to compounds that could have therapeutic relevance.

### 3.4 Graphical Depictions of Biological Phenomena

In the worlds of computer science and mathematics, a graph is a means by which to represent a set of objects and relationships between them. It is frequently depicted as a set of nodes, where each node represents an object. Connections, where they exist, are represented by edges (fig).

Figure 1: An example of a graph. The six nodes are connected by seven edges. All nodes other than node 6 have at least two edges.

Although they originated in the field of computer science, graphs have frequently been used as tools to model biological phenomena. Graphical models of protein interaction networks, gene expression networks, and other similar phenomena are commonplace. Friedman used graphical models to infer and visualize gene regulatory networks (?, ?). Lage et al used graphical models to characterize existing and elucidate novel protein-protein interaction networks (?, ?). Because of its widespread use and adoption, the graphical model is

an appropriate, familiar, and effective means to depict connectivity between CMap signatures.

The Long Tail Problem

### **3.5 Approach for This Work**

In this work, a graph will be used to depict the existence and strength of connections between gene expression signatures in the CMap database. Signatures will be represented as nodes and their pairwise connections by the edges. The final app might look something like Figure 2.

## **4 Methods**

### **4.1 Computing Connections**

Connections between CMap signatures will be computed using weighted connectivity score (WTCS) which is a variant of the Komolgov-Schmirnov (KS) statistic. WTCS is inspired by Subramanian’s use of KS to compute enrichment of gene sets in expression profiles (Subramanian et al., 2005).

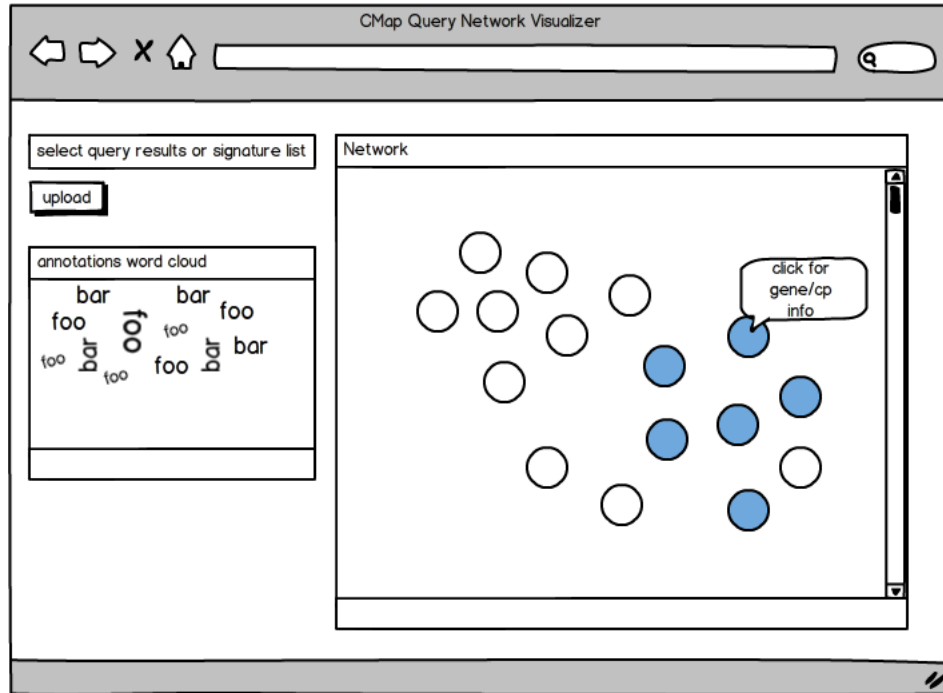


Figure 2: The App

## 4.2 Software Components

### 4.2.1 Front End: HTML & D3.js

### 4.2.2 Back End: Node.js & MongoDB

## 5 Potential Challenges

## 6 Preliminary Timeline

2013-12-01 Receive approval of proposal

**2013-12-15** Investigate and validate then notion of refining hit lists base on their interconnectivity

**2013-12-22** Implement graph visualizer in D3

**2013-12-29** Implement node selection and highlighting

**2014-01-14** Implement Node.js and MongoDB backend

**2014-02-07** Implement user-based input, uploading text files and CMap query results

**2014-03-01** Implement word-cloud generation based on selected graph nodes

**2014-03-28** Implement filtering, showing/hiding graph nodes, brushing to select nodes based on connectivity score

**2014-04-15** Implement text or spreadsheet result export

**2014-05-01** Implement linking to external sites (GeneBank, ChemBank, for clicking on graph nodes)

## References

Friedman, N. (2004, February). Inferring cellular networks using probabilistic graphical models. *Science*, *303*(5659), 799–805. Retrieved 2013-11-29, from <http://www.sciencemag.org/content/303/5659/799> (PMID: 14764868) doi: 10.1126/science.1094068

- Hollander, M., & Wolfe, D. A. (1975). Nonparametric statistical methods. *Biometrische Zeitschrift*, 17(8), 526–526. Retrieved from <http://dx.doi.org/10.1002/bimj.19750170808> doi: 10.1002/bimj.19750170808
- Lage, K., Karlberg, E. O., Strling, Z. M., Iason, P. ., Pedersen, A. G., Rigina, O., ... Brunak, S. (2007, March). A human phenome-interactome network of protein complexes implicated in genetic disorders. *Nature Biotechnology*, 25(3), 309–316. Retrieved 2013-11-29, from <http://www.nature.com.ezp-prod1.hul.harvard.edu/nbt/journal/v25/n3/full/nbt1295> doi: 10.1038/nbt1295
- Lamb, J., Crawford, E. D., Peck, D., Modell, J. W., Blat, I. C., Wrobel, M. J., ... Golub, T. R. (2006, September). The connectivity map: Using gene-expression signatures to connect small molecules, genes, and disease. *Science*, 313(5795), 1929–1935. Retrieved 2013-10-18, from <http://www.sciencemag.org/content/313/5795/1929> (PMID: 17008526) doi: 10.1126/science.1132939
- Subramanian, A., Tamayo, P., Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., ... Mesirov, J. P. (2005, October). Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America*, 102(43), 15545–15550. Retrieved 2013-10-18, from <http://www.pnas.org/content/102/43/15545> (PMID: 16199517) doi: 10.1073/pnas.0506580102

- Tissing, W. J. E., Meijerink, J. P. P., den Boer, M. L., & Pieters, R. (2003, January). Molecular determinants of glucocorticoid sensitivity and resistance in acute lymphoblastic leukemia. *Leukemia*, *17*(1), 17–25. (PMID: 12529655) doi: 10.1038/sj.leu.2402733
- van 't Veer, L. J., Dai, H., van de Vijver, M. J., He, Y. D., Hart, A. A. M., Mao, M., ... Friend, S. H. (2002, January). Gene expression profiling predicts clinical outcome of breast cancer. *Nature*, *415*(6871), 530–536. Retrieved 2013-11-10, from <http://www.nature.com/nature/journal/v415/n6871/abs/415530a.html> doi: 10.1038/415530a