

## **The new direction in bioinformatics: Integrative data mining for genomics and proteomics**

Paul Bertone<sup>1</sup> and Mark Gerstein<sup>2\*</sup>

<sup>1</sup>Department of Molecular, Cellular, and Developmental Biology, <sup>2</sup>Department of Molecular Biophysics and Biochemistry and Department of Computer Science, Yale University, New Haven, CT 06520 USA.

### **Abbreviations**

ChIP Chromatin immunoprecipitation  
EST Expressed sequence tag  
ORF Open reading frame  
PCA Principle component analysis  
SAGE Serial analysis of gene expression  
SOM Self-organizing map  
SVM Support vector machine

**Keywords:** Bioinformatics, microarray expression analysis, machine learning, database integration

\* Corresponding author  
Telephone: (203) 432-6105  
Fax: (360) 838-7861  
E-mail: [Mark.Gerstein@yale.edu](mailto:Mark.Gerstein@yale.edu)

## **Abstract**

Biological research is becoming increasingly database driven, motivated, in part, by the advent of large-scale functional genomics and proteomics experiments such as those comprehensively measuring gene expression. These provide a wealth of information on each of the thousands of proteins encoded by a genome. Consequently, a challenge in bioinformatics is integrating databases to connect this disparate information, and performing large-scale studies to collectively analyze many different data sets. This approach represents a paradigm shift away from traditional single-gene biology, and often involves statistical analyses focusing on the occurrence of particular features (e.g. folds, functions, interactions, pseudogenes, or localization) in a large population of proteins. Moreover, the explicit application of machine learning techniques can be used to discover trends and patterns in the underlying data. We give several examples of these in a genomic context: Clustering methods to organize microarray expression data, support vector machines to predict protein function, Bayesian networks to predict subcellular localization, and decision trees to optimize target selection for high-throughput proteomics.

## **Biological research is database-oriented**

Databases have defined the information structure of molecular biology for over a decade, archiving thousands of protein and nucleotide sequences and three-dimensional structures. As large-scale genomics and proteomics move to the forefront of biological research, the role of databases has become more significant than ever. The current landscape of biological databases includes large public archives, such as GenBank, DDBJ, and EMBL for nucleic acid sequences [1], PIR and SWISS-PROT for protein sequences [2], and the Protein Data Bank for three-dimensional protein structure coordinate sets [3]. Another source of sequence data is dbEST [4], a division of GenBank storing expressed sequence tags (ESTs) from cell lines, which provide information about gene expression in various tissues. Databases such as these have been steadily accumulating gene sequences and protein structures for more than a decade, which are submitted on a per-instance basis from disparate laboratories in the biological sciences community.

In addition to these general biomolecular data repositories, specialized systems have been developed which extend its interpretation by providing a context for individual sequences and structures. The SCOP, CATH, and FSSP [5] databases classify proteins based on structural similarity, Pfam and ProtoMap [6] identify families of proteins based on sequence homology, while PartsList and GeneCensus [7] give dynamic reports on the occurrence of protein families in various genomes. Databases have also been developed

to provide comprehensive access to sequence, expression, and functional data for all the known genes of specific model organisms [8].

#### *Integration of databases and large-scale surveys*

In addition to sequence and structure databases, many diverse experimental data sets have been compiled that focus on various aspects of protein function. However, these need to be related with other data in order to be placed in a useful context, as individual bits of information are relatively meaningless outside of this wider scope. For example, an experiment measuring high mRNA copy number for the human gene J00068 carries little intrinsic value until one associates this fact with the gene's functional classification, learning that it codes for the protein actin, an abundant component of skeletal muscle.

The scale of genomics research precludes the traditional single-molecule approach in biology, where all the experimental knowledge about a given protein could be obtained through reading the literature. Thus, integrative database analysis and data mining are essential aspects of modern biology, enabling the synthesis of large-scale portraits of genome function.

In practice, data mining often takes the form of statistical surveys, in the sense of demographic censuses, of different populations of genes and proteins. The goal in this analysis is to identify certain outstanding features possessed by a given population. For instance, large-scale surveys have been used to characterize the features of pseudogenes (inactive protein-coding regions) in the worm genome [9], correlate the levels of protein

expression with subcellular localization [10], examine the relationship between protein structure with function [11-12], and measure the composition of protein folds in complete genomes [13-14]. Bioinformatics researchers conducting integrative database surveys face the challenge of merging genomics data from many different sources into a common framework. The development of standards allowing federated databases to interoperate is essential.

### **Functional genomics and microarray technology**

With the introduction of complete genome sequences, much comprehensive functional analysis has been done on the proteins encoded by an organism's genome. The sequence determination phase of genomics offers only a fraction of the analytical possibilities. Many functional experiments can be performed using that data, and countless opportunities exist to relate those experimental results with other properties via integrative database analysis.

Among the experimental techniques available for genome-wide analysis are gene disruption [15], two-hybrid studies [16], large-scale proteomics [17], silicone elastomer protein chips [18], serial analysis of gene expression (SAGE) [19], and various DNA microarray technologies. Of these, microarrays have become particularly popular due to the highly parallel nature of the experiments, and the inherent homogeneity of the data captured. Using array technologies, immobilized oligonucleotide fragments having known sequences are exposed to fluorescence-labeled DNA probes, and the signals

corresponding to hybridized fragments are scanned and quantified (Figure 1A). A number of these systems have been devised, allowing the simultaneous interrogation of thousands of genes in a single experiment. Due to the recent sequencing of complete genomes, laboratories are using these array technologies to generate expression data on a scale that most researchers would have considered nearly impossible just a few years ago.

Microarray technologies include cDNA arrays [20] and high-density oligonucleotide systems, such as GeneChips [21] and intergenic DNA arrays [22]. Different approaches are designed to capture different types of information. GeneChips and cDNA arrays enable genome-wide expression monitoring by measuring mRNA copy number under various cellular conditions, compiling an expression profile for many genes at incremental time points. cDNA microarrays measure gene expression as the difference in signal strength between reverse-transcribed mRNA levels from cells under various conditions. This difference is given relative to a reference value, which yields a ratio for every target gene under investigation. Conversely, GeneChips measure the absolute expression levels of mRNA transcripts, in number of copies per cell.

Intergenic DNA arrays are an emerging technology, involving the construction of high-density oligonucleotide chips that are probed with DNA fragments captured via chromatin immunoprecipitation (sometimes referred to as ChIP chips). DNA-binding proteins, complexed with their cognate recognition sequences, are sequestered from cell extracts with antibody-bound beads. The DNA is then isolated and used to probe the array, hybridizing to elements containing complementary nucleotide sequences. In

experiments where transcription factor-bound DNA is used to interrogate a microarray containing intergenic sequences, sites of transcription factor binding can be identified across an entire genome [22]. Evidence of putative gene expression accumulates when both cDNA and ChIP-derived probes are hybridized to identical arrays, and comparisons are made between differential mRNA levels and upstream sites of transcription factor binding. Combined with existing expression analysis methods, this location data can be used to assemble a more comprehensive picture of large-scale genetic function.

## **Machine learning approaches to genomic data analysis**

### *Unsupervised learning and clustering*

A general problem in data analysis is how to structure information into meaningful taxonomies or categories. This issue is of great importance when trying to infer relationships in diverse biological data sets. Statistical methods for finding trends and patterns in experimental results have played a large role in their interpretation. Principal component analysis (PCA) can be an effective method of identifying the most discriminating features in a data set. This technique usually involves finding two or three linear combinations of the original features that best summarize the types of variation in the data. If much of the variation is captured by these two or three most significant principal components, class membership of many data points can be observed.

Several approaches to this problem employ unsupervised learning to find these categories; that is, no *a priori* information is required, and generally, no feedback is given

to the model to adjust its performance. Unsupervised learning enables pattern discovery by organizing data into clusters, using recursive partitioning methods. The  $k$ -means algorithm is a popular instance-based method of cluster analysis. The algorithm partitions data into a predetermined number of categories as instances are examined, according to a distance measure (e.g. Euclidean). Category centroids are fixed at random positions when the model is initialized, which can affect the clustering outcome.

The self-organizing feature map (SOM) [23] consists of a neural network whose nodes move in relation to category membership. As with  $k$ -means, a distance measure is computed to determine the closest category centroid. Unlike  $k$ -means, this category is represented by a node with an associated weight vector. The weight vector of the matching node, along with those of neighboring nodes, are updated to more closely match the input vector. As data points are clustered and category centroids are updated, the positions of neighboring nodes move in relation to them. The number of network nodes which constitute this neighborhood typically decrease over time.

While both algorithms require the number of clusters to be empirically chosen, the SOM algorithm overcomes some limitations of  $k$ -means by imposing global relationships between clusters, thereby improving interpretability. Like PCA, the SOM is capable of reducing high-dimensional data into a 1- or 2- dimensional representation. The algorithm produces a topology-preserving map, conserving the relationships among data points. Thus, although either method may be used to effectively partition the input space into clusters of similar data points, the SOM can also indicate relationships between clusters.



### *Supervised learning and classification*

Analysis of large data sets which contain diverse information often involves the explicit application of supervised learning. This generally involves dividing the data set features into two categories: *predictors*, or features in a data set that are relevant for learning; and the *response variable*, or property to be classified. Given the heterogeneous information presented by integrated bioinformatics databases, a number of potential correlations between predictors and response variables may be discovered, depending on the subset of features used and the classifications sought after.

Machine learning can be applied to a wide variety of biological information, both to partition the data into categories and to classify previously unseen examples. Supervised learning is conducted in two phases, training and testing the classifier model. Using this strategy, the data set is divided into two mutually exclusive sets. The first set is used to train the model, where correct classifications/responses of the input examples are known *a priori*. This information is used to improve the performance of the model and reduce the classification error rate, a process which incrementally adjusts an  $n$ -dimensional hyperplane that serves to partition the data set into categories. Afterward, unseen instances in the test set are classified according to the partitioning established during the training phase.

### *Classification performance vs. ease of interpretation*

For biological research applications, interpretability of results is a key factor in selecting a particular machine learning method. By deciphering the mapping between points in feature space and learned categories, one can apply classification results in a practical context. The ability to decipher this mapping depends largely upon the internal representation used by a particular algorithm. For example, decision trees are quite useful in this respect, because rules can be extracted from the tree that discriminate amongst classes. Backpropagation neural networks, while useful classifiers, employ a numerical weight matrix to establish nonlinear mappings that can be difficult to interpret.

### **Biological data mining applications**

#### *Organizing microarray data*

Clustering algorithms are being applied to microarray data sets with increasing regularity; they are often incorporated into microarray image analysis software, and are therefore frequently used to visualize local and global relationships among hybridization signals captured by the array. Currently, hierarchical clustering is the most popular technique employed for microarray data analysis [24]. Hierarchical methods are based on building a distance matrix summarizing all the pairwise similarities between expression profiles, and then generating cluster trees (also called dendrograms) from this matrix. Genes which appear to be co-expressed at various time points are positioned close to one another in the tree, whose branch lengths represent the degree of similarity between expression profiles.

Agglomerative methods such as  $k$ -means and self-organizing maps have also been used to cluster microarray data. The SOM algorithm has been applied to identify cell differentiation patterns and distinguish between cancer cell types via array expression profiling [25]. Using both hierarchical and neighborhood analysis techniques, genes exhibiting similar mRNA levels tend to cluster together, revealing patterns in gene expression data (Figure 1B).

While useful visualization tools, the analytical utility of either type of clustering approach is limited in terms of interpreting microarray results. These techniques serve to group data points based on changes in mRNA levels under various cellular conditions.

Although cluster membership of related data points can indicate co-expression or co-repression of the genes they represent, the existence of common regulatory pathways based on partitioning expression profiles remains speculative in the absence of external information. Additionally, there is no reason to believe that co-expressed genes are evolutionarily related as hierarchical ordering might imply.

#### *Partitioning phenotype data with multiple clustering techniques*

Where appropriate, clustering methods can be combined to establish a multi-tiered partitioning of a data set. An example of this approach is the tandem use of  $k$ -means clustering with hierarchical clustering as applied to genome-wide phenotypic analysis [15]. Yeast cells containing transposon insertions were sorted with the  $k$ -means algorithm based on growth conditions, to produce clusters of genes with similar phenotypes. The same data was then ordered hierarchically, producing trees of experimental assays,

clustering those which produce transformants with shared phenotypes (Figure 1C). Using this method, assays can be identified that are capable of selecting functionally related proteins.

#### *Classifying gene function with support vector machines*

An example of a supervised learning method applied to functional genomics data is the use of support vector machines (SVM) to classify yeast ORFs into several functional categories by correlating mRNA expression with prior knowledge of gene function [26]. The SVM maps an  $n$ -dimensional input space onto a higher-dimensional feature space, simultaneously transforming a nonlinear class boundary into a simple hyperplane. The new feature space is embodied as a set of nonlinear combinations of the original features. To apply the SVM for gene classification, a set of examples was assembled containing genes of known function, along with their corresponding microarray expression profiles. The SVM was then used to predict the functions of uncharacterized yeast ORFs based on the expression-to-function mapping established during training. Supervised learning techniques appear to be ideal for this type of functional classification of microarray targets, where sets of positive and negative examples can be compiled from genomic sequence annotations.

#### *Predicting subcellular localization with Bayes' rules*

Bioinformatics data mining often involves the aggregation of proteomic data from multiple sources, then using machine learning techniques to predict various macromolecular properties from these features. One such application of this strategy was

the assignment of thousands of proteins in the yeast genome to subcellular compartments [27]. In this study, gene expression levels from a number of SAGE, GeneChip and cDNA microarray experiments were accrued and cross-referenced with protein localization information from public databases. Each protein was first assigned a default probability of being localized to one of several compartments. These were then updated using a Bayesian system that incorporates a number of external features to arrive at probabilities for the localization of each protein (Figure 2A). Training the model on a set of examples for which this information is known made possible the prediction of the localization of unknown proteins. This analysis also enabled the elucidation of the relationships between expression and compartmental localization, with cytosolic proteins being highly expressed, and nuclear and membrane-related proteins exhibiting lower expression levels (Figure 2B).

#### *Using decision trees to optimize high-throughput proteomics*

The future of biological databases will likely include systems that are specifically designed to capture data in a format that facilitates retrospective analysis. This requires careful standardization of experimental parameters, thereby rendering the data amenable to subsequent computation. As discussed above, computing techniques are readily applied to microarray results, partially due to the inherent homogeneity of these data sets. Similarly, a goal in the fusion of databases with data mining applications is to standardize highly divergent data sets as much as possible, making them accessible to machine learning algorithms.

Recently, a combined approach was developed to facilitate distributed collaboration among many laboratories using a shared database system, and subsequently predict macromolecular properties from structural proteomics data via decision tree analysis [28]. Experimental parameters for protein cloning, purification, biophysical characterization, NMR, and X-ray crystallization were recorded in a database. These values were standardized across many laboratories and institutions, providing a common data format for retrospective analysis. Using the database to generate training sets for supervised learning, decision trees [29] were used to classify proteins as either soluble or insoluble, based on features of their amino acid sequences. Useful rules relating these features with protein solubility were then determined by tracing the paths through the decision trees (Figure 3). Protein solubility strongly influences whether a given protein is a feasible target for structure determination, so the ability to predict this property can be a valuable asset in the optimization of high-throughput projects.

## **Conclusion**

As with biological sciences in general, the emphasis of computational biology has changed in recent years from the characterization of individual molecules to the analysis of genome-wide expression profiles and proteomics data sets. In contrast to traditional single-gene experimentation, comprehensive functional data sets are being produced, the scale of which has increased with the tandem sequencing of entire genomes and the advent of microarray technologies. An essential part of this work is the integration of a wide variety of experimental data, to assemble a larger picture of biological function and

frame pertinent information within a meaningful context. These composite data sets are conducive to extensive computational analysis, and present new opportunities for data mining.

Machine learning is a useful tool for interpreting genomic information. Both supervised and unsupervised approaches can often be used to analyze the same kinds of data, depending on the desired result and the range of features available. While supervised learning can be used to classify examples according to a given response variable, unsupervised learning can aid in revealing previously unknown relationships without *a priori* information. However, meaningful clustering relationships are often difficult to discern. Although supervised models require an explicit training and testing regime, they can be used to predict the class of unknown examples based on previous learning by approximating a target function or discovering classification rules. This capability can allow investigators to interpret data classifications and potentially elucidate functional properties.

Bioinformatics calls for a mosaic of computing techniques to facilitate the effective organization and interpretation of data generated by functional genomics projects. Large-scale experiments, such as those performed with microarrays, yield large homogenous data sets that are well suited for computational analysis. Consequently, new types of databases have been created to handle this information [30]. These information systems are designed to provide management and web-based retrieval of microarray data. Currently under debate is whether data sets should be normalized as they are entered into

a particular system, and to what degree the details of experimental conditions should be recorded. While the microarray intensity values are well suited to relational table storage, experimental parameters vary widely across individual projects and currently defy uniform standardization.

The data storage and analysis requirements of genomics research will likely promote the development of hybrid approaches that merge database design with computational methods. To facilitate this, it will be necessary to establish consistent formats for database interoperation, which will add greatly to the utility of experimental results in terms of potential data mining applications.

## References

1. Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Rapp BA, and Wheeler DL (2000). Genbank. *Nucl Acids Res* **28**:15-18; Tateno Y, Miyazaki S, Ota M, Sugawara H, and Gojobori T (2000). DNA bank of Japan (DDBJ) in collaboration with mass sequencing teams. *Nucl Acids Res* **28**:24-26; Baker W, van der Broek A, Camon E, Hingamp P, Sterk P, Stoesser G, and Tuli MA (2000). The EMBL nucleotide sequence database. *Nucl Acids Res* **28**:19-23.
2. Barker WC, Garavelli JS, Huang H, McGarvey PB, Orcutt B, Srinivasarao GY, Xiao C, Yeh LS, Ledley RS, Janda JF, Pfeiffer F, Mewes HW, Tsugita A, and Wu C (2000). The Protein Information Resource (PIR). *Nucl Acids Res* **28**:41-44; Bairoch A, and



- Apweiler R (2000). The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucl Acids Res* **28**:45-48.
3. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, and Bourne PE (2000). The Protein Data Bank. *Nucl Acids Res* **28**:235-242.
4. Boguski MS, Lowe TM, and Tolstoshev CM (1993). dbEST - database for "expressed sequence tags". *Nat Genet* **4**:332-333.
5. Murzin AG, Brenner SE, Hubbard T, and Chothia C (1995). SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J Mol Biol* **247**:536-540; Orengo CA, Michie AD, Jones S, Jones DT, Swindells MB, and Thornton JM (1997), CATH- A Hierarchic Classification of Protein Domain Structures. *Structure* **8**:1093-1108; Holm L and Sander C (1998). Touring protein fold space with Dali/FSSP. *Nucl Acids Res* **26**:316-319.
6. Bateman A, Birney E, Durbin R, Eddy SR, Howe KL, and Sonnhammer EL (2000). The Pfam protein families database. *Nucl Acids Res* **28**:263-266; Yona G, Linial N, and Linial M (2000). ProtoMap: Automatic classification of protein sequences and hierarchy of protein families. *Nucl Acids Res* **28**:49-55.
7. Qian J, Stenger B, Wilson CA, Lin J, Jansen R, Krebs W, Alexandrov V, Echols N, Teichmann S, Park J, Yu H, and Gerstein M (2000). PartsList: a web-based system for dynamically ranking protein folds based on disparate attributes (in press).
8. Cherry JM, Adler C, Ball C, Chervitz SA, Dwight SS, Hester ET, Jia Y, Juvik G, Roe T, Schroeder M, Weng S, and Botstein D (1998). SGD: Saccharomyces Genome Database. *Nucl Acids Res* **26**:73-80; Mewes HW, Frishman D, Gruber C, Geier B, Haase D, Kaps A, Lemcke K, Mannhaupt G, Pfeiffer F, Schuller C, Stocker S, and Weil B

- (2000). MIPS: A database for genomes and protein sequences. *Nucl Acids Res* **28**:37-40;
- Stein L, Sternberg P, Durbin R, Thierry-Mieg J, and Spieth J (2001). WormBase: Network access to the genome and biology of *Caenorhabditis elegans*. *Nucl Acids Res* **29**:82-86;
- Gelbart WM, Crosby M, Matthews B, Rindone WP, Chillemi J, Russo Twombly S, Emmert D, Ashburner M, Drysdale RA, Whitfield E, Millburn GH, de Grey A, Kaufman T, Matthews K, Gilbert D, Strelets V, and Tolstoshev C (1997). FlyBase: a *Drosophila* database. The FlyBase consortium. *Nucl Acids Res* **25**:63-66.
9. Harrison P, Echols N, and Gerstein M (2001). Digging for dead genes: An analysis of the characteristics of the pseudogene population in the *C. elegans* genome. *Nucl Acids Res* (in press).
10. Drawid A, Jansen R, and Gerstein M (2000). Genome-wide analysis relating expression level with protein subcellular localization. *Trends Genet* **16**:426-430 .
11. Hegyi H and Gerstein M (1999). The relationship between protein structure and function: a comprehensive survey with application to the yeast genome. *J Mol Biol* **288**:147-164.
12. Thornton JM, Orengo CA, Todd AE, and Pearl FM (1999). Protein folds, functions, and evolution. *J Mol Biol* **293**:333-342.
13. Gerstein M (1998). Patterns of protein-fold usage in eight microbial genomes: A comprehensive structural census. *Proteins* **33**:518-534.
14. Wolf YI, Brenner SE, Bash PA, and Koonin EV (1999). Distribution of protein folds in the three superkingdoms of life. *Genome Res* **9**:17-26.
15. Ross-Macdonald P, Coelho PS, Roemer T, Agarwal S, Kumar A, Jansen R, Cheung KH, Sheehan A, Symoniatis D, Umansky L, Heidtman M, Nelson FK, Iwasaki H, Hager

- K, Gerstein M, Miller P, Roeder GS, and Snyder M (1999). Large-scale analysis of the yeast genome by transposon tagging and gene disruption. *Nature* **402**:413-418.
16. Uetz P, Giot L, Cagney G, Mansfield TA, Judson RS, Knight JR, Lockshon D, Narayan V, Srinivasan M, Pochart P, Qureshi-Emili A, Li Y, Godwin B, Conover D, Kalbfleisch T, Vijayadamodar G, Yang M, Johnston M, Fields S, and Rothberg JM (2000). A comprehensive analysis of protein-protein interactions in *Saccharomyces cerevisiae*. *Nature* **403**:623-7.
17. Christendat D, Yee A, Dharamsi A, Kluger Y, Savchenko A, Cort JR, Booth V, Mackereth CD, Saridakis V, Ekiel I, Kozlov G, Maxwell KL, Wu N, McIntosh LP, Gehring K, Kennedy MA, Davidson AR, Pai EF, Gerstein M, Edwards AM, and Arrowsmith CH (2000). Structural proteomics of an archaeon. *Nat Struct Biol* **7**:903-909.
18. Zhu H, Klemic JF, Chang S, Bertone P, Casamayor A, Klemic KG, Smith D, Gerstein M, Reed MA, and Snyder M (2000). Analysis of yeast protein kinases using protein chips. *Nat Genet* **26**:283-289.
19. Velculescu VE, Zhang L, Zhou W, Vogelstein J, Basrai MA, Bassett DE Jr., Heiter P, Vogelstein B, and Kinzler KW (1997), Characterization of the yeast transcriptome. *Cell* **88**: 243-251.
20. Shalon D, Smith SJ, and Brown PO, (1996), A DNA microarray system for analyzing complex DNA samples using two-color fluorescent probe hybridization. *Genome Res* **6**: 639-645.
21. Lockhart DJ, Dong H, Byrne MC, Follettie MT, Gallo MV, Chee MS, Mittmann M, Wang C, Kobayashi M, Horton H, et al. (1996). Expression monitoring by hybridization to high-density oligonucleotide arrays. *Nat Biotechnol* **14**:1675-1680.

22. Iyer VR, Horak CE, Scafek CS, Botstein D, Snyder M, and Brown PO (2001). Genomic binding sites of the yeast cell-cycle transcription factors SBF and MBF. *Nature* **409**:533-538.
23. Kohonen T (1984). *Self-organization and Associative Memory*. Berlin: Springer-Verlag.
24. Eisen MB, Spellman PT, Brown PO, and Botstein D (1998). Cluster analysis and display of genome-wide expression patterns. *Proc Natl Acad Sci USA* **95**:14863-14868.
25. Tamayo P, Slonim D, Mesirov J, Zhu Q, Kitareewan S, Dmitrovsky E, Lander E, and Golub T (1999). Interpreting patterns of gene expression with self-organizing maps: Methods and application to hematopoietic differentiation. *Proc Natl Acad Sci USA* **96**: 2907-2912; Golub TR, Slonim DK, Tamayo P, Huard C, Gaasenbeek M, Mesirov JP, Coller H, Loh ML, Downing JR, Caligiuri MA, Bloomfield CD, and Lander ES (1999). Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring. *Science* **286**:531-537; Ross DT, Scherf U, Eisen MB, Perou CM, Rees C, Spellman P, Iyer V, Jeffrey SS, Van de Rijn M, Waltham M, Pergamenschikov A, Lee JCF, Lashkari D, Shalon D, Myers TG, Weinstein JN, Botstein D, and Brown PO (2000). Systematic variation in gene expression patterns in human cancer cell lines. *Nat Genet* **24**:227-235; Perou CM, Jeffrey SS, Van de Rijn M, Rees CA, Eisen MB, Ross DT, Pergamenschikov A, Williams CF, Zhu SX, Lee JCF, Lashkari D, Shalon D, Brown PO, and Botstein D (1999). Distinctive gene expression patterns in human mammary epithelial cells and breast cancers. *Proc Natl Acad Sci USA* **95**:9212-9217; Alizadeh AA, Eisen MB, Davis RE, Ma C, Lossos IS, Rosenwald A, Boldrick JC, Sabet H, Tran T, Yu X, Powell JJ, Yang L, Marti GE, Moore T, Hudson J Jr., Lu L, Lewis DB, Tibshirani R,

Sherlock G, Chan WC, Greiner TC, Weisenburger DD, Armitage JO, Warnke R, Levy R, Wilson W, Grever MR, Byrd JC, Botstein D, Brown PO, and Staudt LM (2000). Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. *Nature* **403**:503-511.

26. Brown MPS, Grundy WN, Lin D, Cristianini N, Sugnet CW, Furey TS, Ares M Jr., and Haussler D (2000). Knowledge-based analysis of microarray gene expression data by using support vector machines. *Proc Natl Acad Sci USA* **97**:262-267.

27. Drawid A and Gerstein M (2000). A Bayesian system integrating expression data with sequence patterns for localizing proteins: comprehensive application to the yeast genome. *J Mol Biol* **301**:1059-1075.

28. Bertone P, Kluger Y, Lan N, Zheng D, Edwards AM, Arrowsmith CH, Montelione GT, and Gerstein M (2001). SPINE: An integrated tracking database and data mining approach for prioritizing feasible targets in high-throughput structural proteomics (in press).

29. Quinlan JR (1993). *C4.5: Programs for Machine Learning*. Morgan Kaufman.

30. National Center for Biotechnology Information ([ncbi.nlm.nih.gov/geo](http://ncbi.nlm.nih.gov/geo)); Ermolaeva O, Rastogi M, Pruitt KM, Schuler GD, Bittner ML, Chen Y, Simon R, Meltzer P, Trent JM, and Boguski MS (1998). Data management and analysis for gene expression arrays. *Nat Genet* **20**:19-23; Aach J, Rindone W, and Church GM (2000). Systematic management and analysis of yeast gene expression data. *Genome Res* **10**:431-445.

Figure 1. A) Scanned image of a cDNA microarray. To construct the array, DNA fragments of known sequence are immobilized to each array element. Messenger RNA transcripts from cell extracts are reverse-transcribed and labeled with the fluorophore cyanine-5, serving as hybridization probes. Total genomic DNA is labeled with cyanine-3 and hybridized to all the array elements as a negative control. The two signals are observed separately with a laser using different excitation wavelengths (633 and 543nm, respectively), and the emissions are scanned and quantified via densitometry. Changes in gene expression are measured as the normalized ratios between mRNA level and background signal for each array element. The image is false-colored after scanning; here, genes that are up-regulated are shown in red, while genes that are down-regulated appear green. In this way, differential expression of thousands of genes can be quantified simultaneously. B) Hierarchical clustering techniques applied to microarray expression data, adapted from [24]. In this software display, microarray intensity values are depicted as colored blocks, where rows represent individual genes, and columns correspond to incremental time points when successive experiments were performed. Thus, each row constitutes an expression profile for a given gene. Red blocks indicate high mRNA levels, green blocks indicate low levels. The data was organized via hierarchical clustering to group genes with related expression profiles. C) Double-clustering of disruption phenotype data using  $k$ -means and hierarchical methods [15]. Both examples of hierarchical clustering use the Pearson coefficient to measure similarity between gene expression profiles. Given two normalized expression ratio profiles, the Pearson correlation coefficient is given by the dot product

$$R_{ij} = \frac{1}{N-1} \mathbf{X}_i \bullet \mathbf{X}_j$$

where  $N$  is the number of elements represented by the profiles  $X_i$  and  $X_j$ . The measurement expression ratio profile  $x$  can be used to compute a  $Z$ -score from the normalized profile  $X$  with the following equation

$$X(k) = \frac{x(k) - x_{avg}}{\sigma_x}$$

where  $X(k)$  and  $x(k)$  are the  $k$ th components of the profiles  $X$  and  $x$  respectively,  $x_{avg}$  is the average and  $\sigma_x$  is the standard deviation of values in  $x$ . The correlation coefficient matrix  $R$  can then be computed for a group of genes, where each matrix element  $R_{ij}$  corresponds to the Pearson correlation coefficient between genes  $i$  and  $j$ . The average correlation coefficient  $R_{avg}$  is found by averaging the elements of  $R$ , including the diagonal. This statistic indicates the overall similarity of the expression profiles in a group of genes.

Figure 2. A) Bayesian system for predicting subcellular localization [27]. Charts depict the state vector for a protein  $m$ , representing the distribution of probabilities of protein  $m$  localizing to various subcellular compartments. Bar graphs illustrate the vectors for each feature considered. The prior state for a given protein is sequentially updated with Bayes' rule using feature vectors that represent different attributes of the protein (e.g., NLS = nuclear localization signal sequence, GLY = glycosylation site). The predicted compartment for each protein is indicated by the probability distribution of the final state vector. B) Total predicted compartmental populations for the yeast genome. Estimates are calculated according to the distribution of a population vector, obtained by summing the probability state vectors of all the proteins considered.



Figure 3. Integrated database and machine learning system for optimizing high-throughput proteomics, adapted from [28]. Experimental data for protein cloning, purification, biophysical characterization, and structure determination are uploaded from many distributed workstations to a central database server via Internet. These are stored as standardized values and used as features for data mining analyses. A number of response variables can be predicted in this manner. Decision trees were built with the C4.5 algorithm [29] to sort proteins based on solubility, using training sets compiled from the database. The model was trained using features derived from the protein sequences, such as the composition of various amino acid groups (e.g., aliphatic =  $C(\text{AIGLV})$ , aromatic =  $C(\text{FWY})$ ), secondary structure features, presence of hydrophobic regions, sequence complexity, etc. Ellipses represent the decision tree nodes, filled to indicate the number of soluble versus insoluble proteins at each vertex. These values appear to the right and left of each node, respectively. Rules that discriminate between classes were then extracted from decision trees by examining the paths from root to leaf nodes that lead to correct classifications. In this example, it was found that proteins without a hydrophobic stretch (Hphobe), with fewer than 25% of their residues in short low-complexity sequences (CPLX(s)), and having less than 6% total leucine composition (C(L)) were likely to be soluble.