

DISCOVERING REGULATORY MODULES FROM HETEROGENEOUS INFORMATION SOURCES

TIJL DE BIE, PIETER MONSIEURS, KRISTOF ENGELEN, BART DE MOOR
K.U.Leuven, ESAT-SCD, Kasteelpark Arenberg 10, 3001 Leuven, Belgium

NELLO CRISTIANINI
U.C.Davis Dept. of Statistics, 360 Kerr Hall, One Shields Ave, CA 95616, US

KATHLEEN MARCHAL
*Centre of Microbial and Plant Genetics and ESAT-SCD, K.U. Leuven
Kasteelpark Arenberg 20, 3001 Leuven-Heverlee, Belgium*

We present a method for inference of transcriptional modules from heterogeneous data sources. It allows identifying the responsible set of regulators in combination with their corresponding DNA recognition sites (motifs) and target genes. Our approach distinguishes itself from previous work in literature because it fully exploits the knowledge of three independently acquired data sources: ChIPchip data; motif information as obtained by phylogenetic shadowing; and gene expression profiles obtained using microarray experiments. Moreover, these three data sources are dealt with in a new and fully integrated manner. By avoiding approaches that take the different data sources into account sequentially or iteratively, the transparency of the method and the interpretability of the results are ensured. Using our method on biological data demonstrated the biological relevance of the inference.

1. Introduction

Nowadays, data representative of different cellular processes are being generated at large scale. Based on these *omics* data sources, the action of the regulatory network that underlies the organism's behavior can be observed.

Whereas until recently bioinformatics research was driven by the development of methods that deal with each of these data sources separately, the focus is now shifting drastically towards integrative approaches dealing with several data sources simultaneously. Indeed, technological and biological noise in the individual data sources is often so prohibitive and unavoidable that standard methods are bound to fail. Then only a combined use of *heterogeneous* and *independently acquired* information sources can help to solve the problem. Furthermore, these different points of view on the biological system allow gaining a holistic insight into the network studied. Therefore, the integration of heterogeneous data is an important, though non-trivial, challenge of current bioinformatics research.

In this study we focus on 3 types of *omics* data that give independent information on the composition of transcriptional modules, the basic building blocks of transcriptional networks in the cell: *ChIPchip* data (*chromatin immunoprecipitation on arrays*) provides information on the direct physical interaction between a regulator and the upstream regions of its target genes; *motif* information as obtained by phylogenetic shadowing describes the DNA recognition sites of these regulators; and *gene expression profiles* obtained using microarray experiments describe the expression behavior in the conditions tested. By integrating these three data sources, we aim at identifying the concerted action of regulators that elicit a characteristic expression profile in the conditions tested, the target genes of these regulators, and the DNA binding sites recognized by these regulators, thus fully specifying the relevant regulatory modules.

Previous successful approaches to integrative analyses in bioinformatics can be found in the class of *kernel methods* [6,5] and methods based on *graphical models* [4,2,3]. Still, to our knowledge, no successful attempts to solve the problem of module inference exploiting all 3 independently acquired ChIPchip, motif and expression data have been made so far. Furthermore, most existing approaches that exploit the availability of heterogeneous data sources proceed in a *sequential* or an *iterative* way (see e.g. [14] for simultaneous detection of motifs and clustering of expression data, e.g. [7] for an iterative approach using ChIPchip and expression data, and e.g. [13] for simultaneous motif detection and analysis of ChIPchip data). In this paper, we present an approach that is different in spirit from previous methods, taking the different data sources into account in a highly *concurrent* way. The performance of the algorithm was demonstrated using the Spellman dataset [10] as a benchmark.

2. Materials and Algorithms

2.1. Data sources

As microarray benchmark set the Spellman dataset was used [10], which contains 77 experiments describing the dynamic changes of 6178 genes during the yeast cell cycle. The profiles were normalized (subtracting the mean of each profile and dividing by the standard deviation across the time points) and stored in a gene expression data matrix further denoted by A with a row for each gene expression profile and a column for each condition.

Genome-wide location data performed by Lee et al. [8] were downloaded from http://web.wi.mit.edu/young/regulator_network. These contain data on the binding of 106 regulators to their respective target genes in rich medium. The

ChIPchip data matrix (further denoted by \mathbf{R}) used in our study consists of *one minus* the p-values obtained from combined ratio's between immunoprecipitated and control DNA (see [8]). Thus, a large value (close to one) indicates that the regulator is probably present.

The motif data used in this study were obtained from a comparative genome analysis between distinct yeast species (phylogenetic shadowing) performed by Kellis et al. [9]. The authors describe the detection of 72 putative regulatory motifs in yeast. These motifs, available online as regular expressions, were transformed into the corresponding probabilistic representation (weight matrix): for each motif, the 20 *Saccharomyces cerevisiae* genes in which the motif was most reliably detected according to the scoring scheme of Kellis et al. [9] were selected. The intergenic sequences of these genes were subjected to motif detection based on Gibbs sampling [MotifSampler,12]. If the statistically overrepresented motif in this set of putatively coexpressed genes corresponded to the motif that was detected by the comparative motif search of [9] the motif model was retained. As such 53 of the 71 motifs could be converted into a weight matrix. This weight matrix was subsequently used to screen all intergenic sequences for the presence of the respective regulatory motifs using MotifLocator [11]. Absolute scores were normalized [11]. As the score distribution of the motif hits depends on the motif length and the degree of conservation of the motif, the distribution of the normalized scores differs between motifs. Therefore, normalized scores were converted into percentile values. This allows for an unbiased choice of the thresholds on the motif quality parameter in the algorithm. The matrix containing these percentile values is the motif data matrix \mathbf{M} that will be used in this work.

2.2. Module construction algorithm

The aim of the method is to find regulatory modules, based on the gene expression, ChIPchip, and Motif data matrices as specified above.

A module is fully specified by the set of genes it regulates (denoted by an index set \mathbf{g} , pointing to the relevant set of rows of \mathbf{R} , \mathbf{M} and \mathbf{A}), in addition to the set of regulators (corresponding to the columns with indices in a set called \mathbf{r} in the ChIPchip matrix \mathbf{R}) and motifs (corresponding to the columns with indices \mathbf{m} in the Motif matrix \mathbf{M}) that are responsible for the regulation of these genes. The goal of our method is to come up with regulatory modules specified in this way, by fully exploiting the heterogeneous data sources available.

We note that the principles behind the method developed here are based on ideas similar to those that laid the foundations for the Apriori algorithm, originally developed in the database community [15].

Seed construction. This is the main step of the algorithm, and allows the construction of a good guess (or *seed*) of the modules. The idealized goal of this step is to find a set of genes \mathbf{g} , that have the same expression profile, and such that there exist sufficiently large sets of regulators \mathbf{r} and of motifs \mathbf{m} that are entirely present in all these genes. Since in practice it is not known exactly in which intergenic regions a certain motif occurs or where a regulator binds, we have to resort to the score matrices \mathbf{R} and \mathbf{M} . Furthermore, the expression profiles \mathbf{A} of genes in a module will only be approximately equal, and possibly only in a set of conditions, so we relax this constraint to requiring a strong correlation instead of equality between them.

Formally, then the task to solve is:

Find all maximal gene sets \mathbf{g} for which there exist an \mathbf{r} of size $|\mathbf{r}|=r_{\min}$ and a set \mathbf{m} of size $|\mathbf{m}|=m_{\min}$, such that the following 3 constraints are satisfied:

1. $\mathbf{R}(i,j) > t_r$ for all $i \in \mathbf{g}$ and $j \in \mathbf{r}$
2. $\mathbf{M}(i,j) > t_m$ for all $i \in \mathbf{g}$ and $j \in \mathbf{m}$
3. $\text{corr}(\mathbf{A}(i,:), \mathbf{A}(j,:)) > t_a$ for all $i, j \in \mathbf{g}$

where r_{\min} , m_{\min} and thresholds t_r , t_m and t_a are parameters of the method.

Here, a maximal set \mathbf{g} is defined as a set that cannot be extended with another gene without violating one or more of these constraints. In the following, we will use the term *valid set* for a gene set \mathbf{g} that satisfies these constraints.

Clearly it is computationally impossible to tackle this problem with a naive approach: the number of gene sets is exponentially large in the number of genes in the dataset, which is prohibitive even for the smallest genomes. However, it is trivial to verify that:

Observation 1: When a gene set does not satisfy the constraints, none of its supersets satisfy the constraints.

This means that we can build up the maximal sets incrementally, starting with valid sets of size one, and gradually expanding them. Concretely, the (already less naive) algorithm would then look like*:

* Notationally, we will use \mathbf{L}^i to denote the list containing all valid gene sets with i genes. For an individual valid gene set we will use a bold face \mathbf{g}_k^i , with a superscript i to specify that it is an element of \mathbf{L}^i and thus contains i genes, and with a subscript k to distinguish it from the other gene sets in \mathbf{L}^i . The x -th gene in this gene set is denoted as $g_k(x)$, for brevity without superscript.

- For all single genes, check if they satisfy constraints 1 and 2 (constraint 3 is trivially satisfied for singleton gene sets). Make a list \mathbf{L}^1 of all singleton gene sets that contain such a valid gene.
- Set $i = 2$.
- While $\text{size}(\mathbf{L}^{i-1}) \neq 0$
 - For $k=1:\text{size}(\mathbf{L}^{i-1})$, expand set $\mathbf{g}_k^{i-1} = \{g_k(1), g_k(2), \dots, g_k(i-1)\} \in \mathbf{L}^{i-1}$ once for each gene g that is not yet contained in \mathbf{g}_k^{i-1} . Put the thus expanded sets $\{g_k(1), g_k(2), \dots, g_k(i-1), g\}$ that satisfy the 3 constraints (to be verified in \mathbf{R} , \mathbf{M} and \mathbf{A}), in a list \mathbf{L}^i .
- Set $i = i+1$.

Notice that following this strategy, a gene set can be constructed in different ways, by adding the genes to it in a different ordering (i.e. in different iterations i). This can be avoided by adding a gene to a gene set \mathbf{g}_k^{i-1} only whenever its row number g is larger than that of all other genes already in \mathbf{g}_k^{i-1} . Thus for every $\mathbf{g}_k^i = \{g_k(1), g_k(2), \dots, g_k(i)\} \in \mathbf{L}^i$ we always have that $g_k(x) < g_k(y)$ for $x < y$.

Additionally, in this way we can easily keep the list \mathbf{L}^i of gene sets \mathbf{g}_k^i sorted as well, where the sorting is carried out first according to the first added gene and last according to the last added gene. More formally: \mathbf{g}_k^i precedes \mathbf{g}_l^i in \mathbf{L}^i if and only if $g_k(\text{argmin}_{x:(g_k(x) \neq g_l(x))}) < g_l(\text{argmin}_{x:(g_k(x) \neq g_l(x))})$ (this ordering of the list \mathbf{L}^i is indeed a total ordering relation.)

Still the number of expanded gene sets can be huge in every iteration: each of the gene sets \mathbf{g}_k^{i-1} in \mathbf{L}^{i-1} must be expanded by all genes $g > g_k(i-1)$, after which the validity has to be checked by looking at the matrices \mathbf{R} , \mathbf{M} and \mathbf{A} . This can still be too expensive. However, we can exploit the converse of Observation 1:

Observation 2: Whenever a gene set satisfies the constraints, all of its subsets satisfy the constraints.

Using this so-called *hereditary property* of the constraint set, in some cases we can conclude *a priori*—i.e. without checking in \mathbf{R} , \mathbf{M} and \mathbf{A} —if an extended gene set of size i can possibly be valid or not: we simply have to check if all of its size $i-1$ subsets belong to \mathbf{L}^{i-1} . Only if this is the case, we still have to access the data in \mathbf{R} , \mathbf{M} and \mathbf{A} ; if it is not the case, we know without further investigation that the extended subset is invalid.

Specifically, assume that we expand the gene set $\mathbf{g}_k^{i-1} = \{g_k(1), g_k(2), \dots, g_k(i-2), g_k(i-1)\} \in \mathbf{L}^{i-1}$ with g , leading to $\{g_k(1), g_k(2), \dots, g_k(i-2), g_k(i-1), g\}$. Then, since for a valid size i set each of its size $i-1$ subsets must be contained in \mathbf{L}^{i-1} ,

also $\{g_k(1), g_k(2), \dots, g_k(i-2), g\}$ must be contained in \mathbf{L}^{i-1} . In other words: there has to be a $\mathbf{g}_l^{i-1} = \{g_l(1), g_l(2), \dots, g_l(i-2), g_l(i-1)\} \in \mathbf{L}^{i-1}$ for which $g_k(x) = g_l(x)$ for $x \leq i-2$, and $g = g_l(i-1)$. This can efficiently be ensured *constructively*, by exploiting the fact that the list \mathbf{L}^{i-1} , and all \mathbf{g}_k^{i-1} themselves are sorted. Indeed, thanks to this, all gene sets \mathbf{g}_k^{i-1} that have the first $i-2$ genes in common occur consecutively in \mathbf{L}^{i-1} . Therefore, to expand \mathbf{g}_k^{i-1} with an additional gene, we only have to screen the list \mathbf{L}^{i-1} starting at \mathbf{g}_{k+1}^{i-1} and move forward in \mathbf{L}^{i-1} for as long as the first $i-2$ genes are equal to $g_k(1), g_k(2), \dots$ and $g_k(i-2)$. For every gene set \mathbf{g}_l^{i-1} screened in this way, read the last gene $g_l(i-1)$ and append it to \mathbf{g}_k^{i-1} , thus resulting in a candidate gene set of size i , potentially to be appended to \mathbf{L}^i . To find out whether this candidate gene set is valid indeed, one still has to check the constraints explicitly. However, thus constructively exploiting the hereditary property, the number of queries to \mathbf{R} , \mathbf{M} and \mathbf{A} is drastically reduced. Note that this strategy also ensures that \mathbf{L}^i is sorted automatically.

Module validation. In some cases the first step described above is not sufficient for adequate module inference. There are three reasons for this:

First, the seed construction method can be rather conservative in recruiting genes, since each of the genes in the module has to satisfy all 3 of the constraints. Therefore, in a second step, we calculate the mean of the expression profiles of the seed modules found in the first step, further called the *seed profile*. Then we can additionally recruit all genes with a high correlation with the seed profile to be incorporated in the module. In order to determine an optimal threshold value for this correlation, we compute the enrichment of each of the motifs and regulators in the genes that have an expression profile that achieves this threshold correlation with the seed profile. The logarithm of the p-value of the enrichment is then plotted as a function of this threshold (Figure (1)), and the threshold can be chosen such that this value is minimal.

Second, sometimes it is undesirable to a priori decide how many motifs and regulators we want in the module, or it may be difficult to choose the thresholds t_r , t_m and t_a (even though experiments show little dependence on these). Then one can first use the seed construction algorithm requiring only 1 regulator and motif, and with stringent thresholds, after which again the enrichment of all motifs and regulators can be plotted as a function of the correlation threshold with the mean profile of the seed module. For each such seed profile, the corresponding enrichment plot will visually hint at the number of motifs and regulators (namely the number of significantly enriched motifs and regulators).

Third, similarly, the enrichment plot allows excluding false positive motifs or regulators: when they are selected in step 1, but appear not to be enriched in the validation step, they are considered as a false positives and discarded.

To calculate the enrichment, we first calculate the mean score of the module for the particular motif or regulator. Note that the mean score of a module by random gene selection is approximately Gaussianly distributed (central limit theorem), with mean equal to the mean over all genes, and variance equal to the overall variance divided by the size of the module. Thus, we can calculate the enrichment as the logarithm of the p-value based on a Gaussian approximation.

Note that the p-values have been computed based on profiles that have been obtained from the data, such that they do not have a rigorous probabilistic interpretation here. Hence, we can only use them as explained above.

2.3. Calculating overrepresentation of functional classes

Functional categories for each gene were obtained from MIPS [18]. Functional enrichment of the modules was calculated using the hypergeometric distribution [16], which assigns to each functional class a p-value.

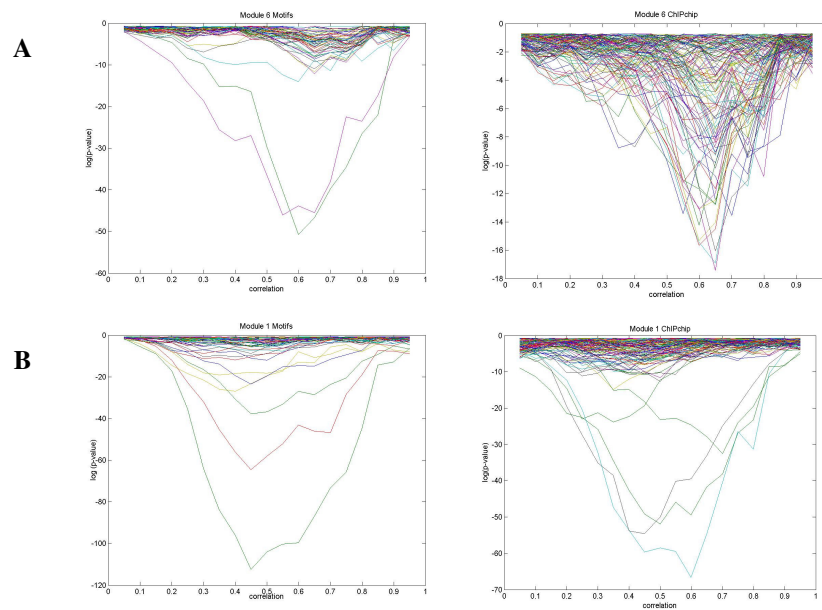


Figure 1: Two examples (A and B) of the module validation step for two seed profiles: on the left, the logarithms of the p-values are plotted for all motifs as a function of the correlation threshold, on the right similarly for the regulators. Panel A shows the results for a false positive prediction (module 6): the regulators (right figure) of the identified seed module turn out not to be significantly overrepresented in genes correlated with the seed profile. In panel B the results are displayed for the positive example described in the text.

3. Results

Cell cycle related modules. To test the reliability of our method, we used the well-studied Spellman dataset as benchmark. The analysis we performed using our two-step algorithm is illustrated by elaborating on the detection of the cell cycle related module 1. Using the seed detection step we searched for modules of genes having at least 1 common motif (1M) in their intergenic sequences and 1 common regulator (1R) showing a small p-value in the ChIPchip data, and of which the expression profiles were mutually correlated with a minimal correlation of 0.7. This seed identification step then predicts several potential modules, and for each of them a seed profile can be calculated. For each of these modules we performed the module validation step.

Fig. 1A (right figure) shows how this validation step allows one to visually detect that the regulator associated with this module is probably a false positive.

In Fig 1B, using the parameter settings of 1M/1R, we identified a potential seed module containing regulator 98 (Swi4) and motif M_11 (known as a Swi4 motif). Calculating the statistical overrepresentation of all motifs and regulators in genes correlated with the seed profile of this putative module showed that in this subset of genes indeed M_11 and Swi4 were overrepresented. The identified module seed thus is likely to be biologically relevant. These results also show that besides Swi4 and M_11, 3 additional motifs and regulators were overrepresented in subsets of genes correlated with the module seed profile, indicating the probable underestimation of the real module size. To verify whether these other regulators/motifs co-occur in the same subsets of genes and therefore comprise a larger module, we repeated the seed identification step using additional parameter settings (see Table 1 in the online supplement). From this result it appeared that we could recover a complete module consisting of the 3 overrepresented regulators (Mbp1, Swi4, Swi6) and 2 motifs (M_16, M_10) and that this module is present in genes displaying an expression profile that shows a correlation of at least 0.7 with the average seed profile. Checking the identities of the regulators and the motifs (regulators Mbp1, Swi4, Stb1 combined with the regulatory motifs Mbp1 (M_18, M_12) and Swi4 (M_11 and M67)) showed that we identified a previously extensively described regulatory module of the yeast cell cycle.

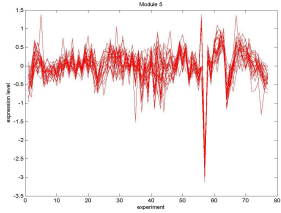
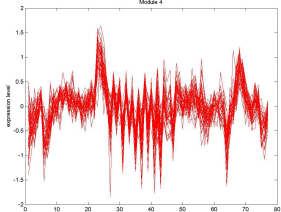
Besides this first module, 3 additional related cell cycle (Table 1) modules could be retrieved. Additional information on each of the separate modules can be found in the online supplement. Genes in the different modules showed peak expressions shifted in time relative to each other, as shown in Figure 1 of the online supplement. All of the predicted modules are conform the previously described knowledge on the cell cycle [8,17,7].

Table 1: Cell cycle related modules. Column 'R' contains the regulators, column 'M' the motifs, the column 'Functional Class: p-value' contains p-values for several functional classes, and the 'Seed Profile' column contains a plot with the expression profiles of the genes regulated by the module.

	R	M	Functional Class: p-value	Seed Profile
Module 1	Mbp1 Swi6 Swi4 Stb1	M_18 (Mbp1) M_12 (Mbp1) M_11 (Swi4) M_67 (Swi4)	10 CELL CYCLE AND DNA PROCESSING: 0 10.03 cell cycle: 2.7e-5 10.01 DNA processing: 1.3e-4 42.04 cytoskeleton: 4.2e-3	
Module 2	Swi4 Mbp1 Swi6 FKH2	M_18 (Mbp1) M_12 (Mbp1) M_11 (Swi4) M_8 (Mcm)	40 CELL FATE: 5.2e-4 40.01 cell growth / morphogenesis: 2.6e-3 43 CELL TYPE DIFFERENTIATION: 5.2e-3 43.01 fungal/microorganismic cell type differentiation: 5.2e-3 34.11 cellular sensing and response: 5.3e-3 01.05.01 C-compound and carbohydrate utilization: 6.8e-3 10.03.04.03 chromosome condensation: 9.4e-3	
Module 3	NDD1 FKH2 Mcm1	M_8 (Mcm) M_30 (Mcm)	43 CELL TYPE DIFFERENTIATION: 3.6e-3 43.01 fungal/microorganismic cell type differentiation: 3.6e-3 10.03.03 cytokinesis (cell division) /septum formation: 4.8e-3	
Module 4	Swi5 (Ace2)	M_8 (Mcm)	32.01 stress response: 3.2e-3 10.03 cell cycle: 8.7e-3	

Non cell cycle related modules. Besides the modules primarily involved in cell cycle, other modules could be identified in the Spellman dataset (see Table 2). Module 5, consisting of Fhl1, Rap1 and Yap5, involved in the regulation of ribosomal proteins was previously also identified [8]. Note that it was identified from a noise profile (i.e. a profile that does not change significantly and consistently with the cell cycles over the different time points) in this cell cycle dataset, indicating that even biological noise contains important information on regulatory networks. By our analysis we could pinpoint motif M_54 [9], as the regulatory motif correlated with this regulatory module. A second non cell cycle related module consisted of the genes regulated by the motifs M_7 and M_3 (identified as ESR1 and ESR2 [9]). For this module, related to transcription and ribosomal RNA processing only the motifs seemed informative (see module 6 in Table 2, and Figure 1A).

Table 2: Non cell cycle related modules.

	R	M	Functional Class: p-value	Seed Profile
M o d u l e 5	FKL1 Yap5 Rap1	M_54	12 PROTEIN SYNTHESIS: 0 12.01 ribosome biogenesis: 0	
M o d u l e 6	/	M_3 (ESR1) M_7 (ESR2)	11 TRANSCRIPTION: 0.000002 11.04 RNA processing: 0 11.04.01 rRNA processing: 0	

4. Discussion

We described a methodology combining ChIPchip, motif and expression data to infer complete descriptions of transcriptional modules. Our methodology consists of 2 steps. The seed construction step predicts the putative modules consisting of regulators, their corresponding motifs and the elicited expression profile. The validation step filters false positive predictions and gives further insight into the module size.

The problem is attacked in a very direct way: the integration of the data sources is achieved in a one-shot-algorithm, and requires no iteration over the different data sources. While the running time was very reasonable for all experiments carried out for this paper, it heavily depends on the parameters. The more stringent they are set, the smaller the lists L^i will be and the faster the algorithm will run. Further speed-ups are possible, but not needed for the experiments reported in this paper. Therefore we will not go into these here.

The Spellman dataset was used as a benchmark to test the performance of our method. Since this dataset and the yeast cell cycle have extensively been studied before [7,8], it is ideally suited for testing the reliability and biological relevance of the predictions. We were able to reconstruct 4 important modules known to be involved in cell cycle and also 2 non cell cycle related modules without using any prior biological knowledge or prior data reduction. These results indicate that predictions passing the module validation step are likely to be biologically relevant (no false positives present).

5. Conclusion

The 3 data types mutually agreeing with each other on the prediction of a module not only results in the most reliable predictions (as was the case for the cell cycle related modules), but also allows correlating a set of regulators with their corresponding regulatory motifs and elicited profiles in a very natural and direct way. On the other hand, because of the restricted number of experimental data yet available (chip data not known for all regulators and tested in a limited set of conditions, expression data for specific conditions not available), and the questionable quality of the motif models, the presence of a signal in 1 data type can compensate for the lack of it in another data type, allowing still to retrieve the module.

While to our knowledge this is the first time these 3 independently acquired data sources are exploited in such a concurrent way for module identification, the approach is further extendible towards any number of information sources, and in principle towards the use of other data types. The only condition for an efficient method to exist is that the constraints the gene sets have to satisfy must be hereditary. This extension will be the subject of future work.

Acknowledgments

KM and TDB are post-doctoral researcher and research assistant of the Fund for Scientific Research – Flanders (FWO–Vlaanderen) respectively. KE is a research assistant with the IWT. This work is partially supported by: 1. IWT projects: GBOU-SQUAD-20160; 2. Research Council KULeuven: GOA Mefisto-666, GOA-Ambiorics, IDO genetic networks; 3. FWO projects: G.0115.01 (microarrays/oncology), G.0413.03 (inference in bioi), G.0388.03 (microarrays for clinical use), G.0229.03 (ontologies in bioi), and G.0241.04 (Functional Genomics); 4. IUAP V-22 (2002-2006).

References

Supplementary information:

www.esat.kuleuven.ac.be/~kmarchal/Supplementary_Info_PSB2005/SuppWebsiteYeastPSB.html

- [1] **Nariai N, Kim S, et al.** 2004. *Using Protein-Protein Interactions for Refining Gene Networks Estimated from Microarray Data by Bayesian Networks*. Pacific Symposium on Biocomputing, 9: 336-347.
- [2] **Segal E, Wang H, et al.** 2003. *Discovering Molecular Pathways from Protein Interaction and Gene Expression Data*. Bioinformatics, 19(Suppl 1): 264-272.
- [3] **Segal E, Shapira M, et al.** 2003. *Module Networks: Identifying Regulatory Modules and their Condition Specific Regulators from Gene Expression Data*. Nature Genetics, 34(2): 166-176.
- [4] **Hartemink A, Gifford D, et al.** 2002. *Combining Location and Expression Data for Principled Discovery of Genetic Regulatory Networks*. Pacific Symposium on Biocomputing, 7: 437-449.
- [5] **Lauckriet G, De Bie T, et al.** 2004. *A statistical framework for genomic data fusion*. Bioinformatics, accepted.
- [6] **Vert JP and Kanehisa M.** 2003. *Extracting active pathways from gene expression data*. Bioinformatics, 19: 238ii-234ii.
- [7] **Bar-Joseph Z, Gerber GK, et al.** 2003. *Computational discovery of gene modules and regulatory networks*. Nature Biotechnology, 21(11): 1337-1342.
- [8] **Lee TI, Rinaldi NJ, et al.** 2002. *Transcriptional regulatory networks in Saccharomyces cerevisiae*. Science, 298(5594): 799-804.
- [9] **Kellis M, Patterson N, et al.** 2003. *Sequencing and comparison of yeast species to identify genes and regulatory elements*. Nature, 423(6937): 241-254.
- [10] **Spellman PT, Sherlock G, et al.** 1998. *Comprehensive Identification of Cell Cycle-regulated Genes of the Yeast Saccharomyces cerevisiae by Microarray Hybridization*. Molecular Biology of the Cell, 9: 3273-3297.
- [11] **Marchal K, De Keersmaecker S, et al.** 2004. *In silico identification and experimental validation of PmrAB targets in Salmonella typhimurium by regulatory motif detection*. Genome Biology, 5(2): R9.
- [12] **Thijs G, Marchal K, et al.** 2002. *A Gibbs sampling method to detect overrepresented motifs in the upstream regions of coexpressed genes*. Journal of Computational Biology, 9(2): 447-464.
- [13] **Liu XS, Brutlag DL et al.** 2002. *An algorithm for finding protein-DNA binding sites with applications to chromatin-immunoprecipitation microarray experiments*. Nature Biotechnology, 20(8): 835-839.
- [14] **Lapidot M and Pilpel Y.** 2003. *Comprehensive quantitative analyses of the effects of promoter sequence elements on mRNA transcription*. Nucleic Acids Research, 31(13): 3824-3828.
- [15] **Agrawal R, Imielinski T et al.** 1993. *Mining Association Rules between Sets of Items in Large Databases*. Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data, pp. 207-216.
- [16] **Tavazoie S, Hughes JD, et al.** 1999. *Systematic determination of genetic network architecture*. Nature Genetics, 22(3):281-285.
- [17] **Costanzo MC, Hogan JD, et al.** 2000. *The yeast proteome database (YPD) and Caenorhabditis elegans proteome database (WormPD): comprehensive resources for the organization and comparison of model organism protein information*. Nucleic Acids Research, 28(1): 73-76.
- [18] **Mewes HW, Amid C, et al.** *MIPS: analysis and annotation of proteins from whole genomes*. Nucleic Acids Research, 32, Database issue D41-D44.
- [19] **Friedman N.** 2004. *Inferring Cellular Networks Using Probabilistic Graphical Models*. Science, 303(5659): 799-805.