

## UNDERSTANDING DISTAL TRANSCRIPTIONAL REGULATION FROM SEQUENCE, EXPRESSION AND INTERACTOME PERSPECTIVES

ARVIND RAO\*, DAVID J. STATES†  
and ALFRED O. HERO, III‡

*Bioinformatics, University of Michigan  
Ann Arbor, MI 48109, USA  
\*[ukarvind@umich.edu](mailto:ukarvind@umich.edu)  
†[dstates@umich.edu](mailto:dstates@umich.edu)  
‡[hero@umich.edu](mailto:hero@umich.edu)*

JAMES DOUGLAS ENGEL  
*Cell and Developmental Biology  
University of Michigan  
Ann Arbor, MI 48109, USA  
[engel@umich.edu](mailto:engel@umich.edu)*

Received 20 July 2009  
Revised 17 October 2009  
Accepted 17 October 2009

Gene regulation in eukaryotes involves a complex interplay between the proximal promoter and distal genomic elements (such as enhancers) which work in concert to drive precise spatio-temporal gene expression. The experimental localization and characterization of gene regulatory elements is a very complex and resource-intensive process. The computational identification of regulatory regions that confer spatiotemporally specific tissue-restricted expression of a gene is thus an important challenge for computational biology. One of the most popular strategies for enhancer localization from DNA sequence is the use of conservation-based prefiltering and more recently, the use of canonical (transcription factor motifs) or *de novo* tissue-specific sequence motifs. However, there is an ongoing effort in the computational biology community to further improve the fidelity of enhancer predictions from sequence data by integrating other, complementary genomic modalities.

In this work, we propose a framework that complements existing methodologies for prospective enhancer identification. The methods in this work are derived from two key insights: (i) that chromatin modification signatures can discriminate proximal and distally located regulatory regions and (ii) the notion of promoter-enhancer cross-talk

\*Corresponding author.

(as assayed in 3C/5C experiments) might have implications in the search for regulatory sequences that co-operate with the promoter to yield tissue-restricted, gene-specific expression.

*Keywords:* Nephrogenesis; random forests; transcriptional regulation; transcription factor binding sites (TFBS); *GATA* genes; comparative genomics; functional genomics; tissue-specific genes; network analysis; directed information; heterogeneous data integration.

## 1. Introduction

Understanding the mechanisms underlying regulation of tissue-specific gene expression remains a challenging question. While all mature cells in the body have a complete copy of the human genome, each cell type only expresses those genes it needs to carry out its assigned task. This includes genes required for basic cellular maintenance (often called “housekeeping genes”) and those genes whose function is specific to the particular tissue type that the cell belongs to. Gene expression by way of transcription is the process of generation of messenger RNA (mRNA) from the DNA template representing the gene. It is the intermediate step before the generation of functional protein from messenger RNA. During gene expression, transcription factor (TF) proteins are recruited at the proximal promoter of the gene as well as at sequence elements (enhancers/silencers) which can lie several hundreds of kilobases from the gene’s transcriptional start site (Figs. 1 and 2).

It is hypothesized that the collective set of transcription factors that drive (regulate) expression of a target gene are cell, context-and tissue-dependent.<sup>1,2</sup> Some of these TFs are recruited at proximal regions such as the promoter of the gene, while others are recruited at these distal regulatory regions. There are several (hypothesized) mechanisms for promoter-enhancer interaction through protein interactions between TFs recruited at these elements during formation of the transcriptional complex.<sup>3</sup> A commonly accepted mechanism of distal interaction, during regulation, is looping,<sup>4-6</sup> shown in Fig. 2, wherein intervening DNA between the enhancer and promoter is “looped out” to facilitate the interaction between the TFs of the promoter and the enhancer, leading to formation of the transcriptional complex.

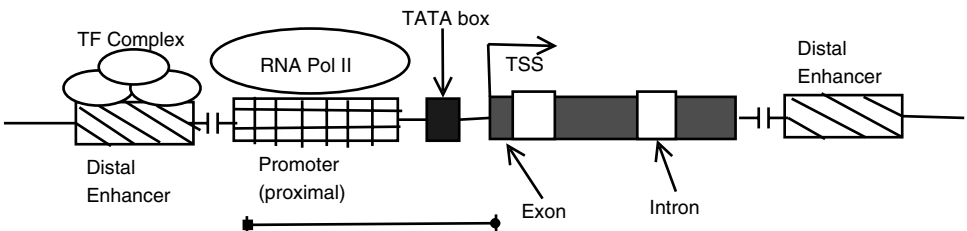


Fig. 1. Schematic of transcriptional regulation. Sequence motifs at the promoter and the distal regulatory elements together confer specificity of gene expression via TF binding.

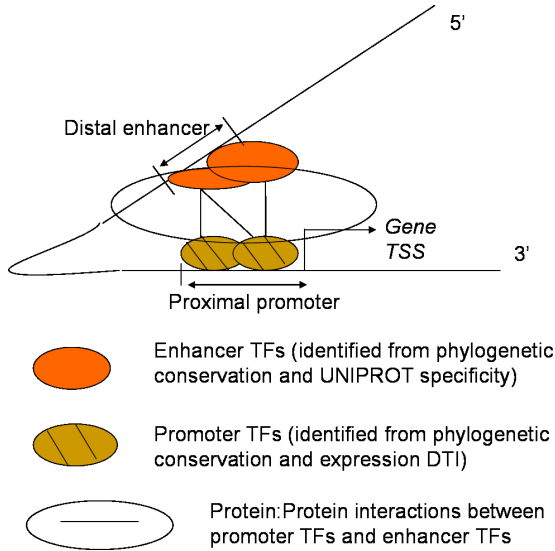


Fig. 2. Distal enhancer-promoter interaction via looping is mediated via protein interactions during TF complex formation. The set of TFs that are putatively recruited at the proximal promoter and distal enhancer can be found from sequence and expression data. Evidence of protein-interaction between proximal and distal TFs can be found from interaction databases.

An important challenge in current biology is to understand *where* functional regulatory elements (like enhancers) are located for a gene of interest. Given the complexity of the regulatory process, there are several instances wherein the enhancer for a gene is located hundreds of kilobases from the gene it regulates.<sup>5,7,8</sup> One of the typical experimental approaches to localize a gene-specific enhancer is via bacterial artificial chromosome (BAC) trap assays.<sup>9,10</sup>

Thereafter, using conservation and TFBS-based criteria, smaller genomic sequences (1 – 2kb) are isolated for subsequent transgenic analysis. However, even short genomic regions can have several conserved sequence elements (CSEs) worthy of experimental testing (e.g. ~ 120 CSEs surpass a 70% sequence conservation in a 45kb *human-mouse* aligned region, neighboring *Gata2*). Since an experimental analysis of each of these several regions is clearly unfeasible, there is a need for the use of principled methodologies that could potentially reduce this large list of enhancer candidates to a much shorter *high-confidence* list for experimental validation.

Since the main problem of interest is the prospective discovery of enhancers in a pre-specified sequence region, it would seem imperative to explore modalities that supplement conservation and TFBS criteria to reduce false positives. In this work, we explore three such modalities that emerge from functional genomic assays (from several recent independent studies as well as from the ENCODE project). These three modalities reveal some interesting new features of regulatory regions that are potentially of great use in discriminating gene-specific enhancers versus

other neutral regions. We note that there are promoter-independent enhancers too, and that their computational study has been far more principled.<sup>1,11</sup> However, their study is outside the scope of this study where we focus on gene-specificity in addition to tissue-specificity.

Understanding the characteristics of such regulatory regions entails several aspects:

- (1) Do regulatory regions like promoters and enhancers have any interesting *sequence properties* that depend on their tissue-specificity of gene expression? Such properties can be examined based on their individual sequences or their epigenetic preferences. A common approach is the identification of canonical or *de novo* tissue-specific motif-signatures<sup>2,13</sup> for such elements, and has been applied quite extensively. In this work, we focus on the sequence-specificity and epigenetic preferences of tissue-specific distal regulatory regions (enhancers) versus proximal regulatory regions (promoters).
- (2) To reduce the large number of false positives that arise from sequence comparisons alone, we appeal to a mechanistic insight from biology. For long-range transcriptional regulation to be possible, there has to be an enhancer-promoter interaction during formation of the tissue-specific, gene-specific transcriptional machinery. Literature suggests that such interaction is mediated by protein-protein interactions between promoter TFs and enhancer TFs after looping along the chromosomal length.<sup>4,6,14</sup> This insight (Fig. 2) leads to two further questions:
  - Which TFs bind the promoter and the putative enhancer(s)?
  - Does this resultant “interaction graph” between enhancer and promoter TFs have any special *structural* characteristic that can discriminate functional non-coding regulatory regions from non-functional ones?

The primary goal of answering the questions above is to build an enhancer discovery program that can localize tissue-restricted gene-specific enhancers in a given chunk of genome sequence (within a  $\sim 200$  kb genomic window, as obtained from BAC trap strategies<sup>10</sup>). These questions will help us understand the nature of distal regulatory regions and provide a way to complement existing approaches in enhancer localization<sup>11,12</sup> to achieve lower false positive rate and higher experimental efficacy.

As a case study to answer these questions, we examine the distal regulation of *Gata2* regulation in the developing kidney. *Gata2* is a gene belonging to the GATA family of transcription factors (*GATA1-6*), and binds the consensus –WGATAR– motif on DNA. It is located on mouse chromosome 6, and plays an important role in mammalian hematopoiesis, nephrogenesis and CNS development, with important phenotypic consequences. The study of long-range regulatory elements that effect *Gata2* expression has been ongoing for several years now.

Recently, Khandekar and co-workers<sup>10</sup> reported the characterization of two enhancer elements, conferring urogenital-specific (UG) expression of *Gata2*, between 80–150 kb downstream from the *Gata2* transcription start site on chromosome 6. In this experiment, four regions were selected for transgenic analysis based on sequence identity and TF motif matches. However, only two of these were functional.

Based on the insights from the various individual studies since and the ENCODE project, outlined above, we asked if it might now be possible to explain the behavior of these four regions along these new modalities (tissue-specificity, epigenetic signatures and TF-interaction graphs), thereby enabling the proposal of a *framework for promoter-specific enhancer discovery from sequence*.

## 2. Rationale and Data Sources

The overall schematic of distal transcriptional regulation via looping is shown in Fig. 2. This schematic and the discussion in Sec. 1 suggests the decomposition of the regulatory process along three main modalities: sequence, expression and interactome. Our main goal in this paper is to understand urogenital enhancer potential of these four UG sequence candidates<sup>10</sup> from these three perspectives. These attributes are discussed below:

- (1) **Sequence perspective:** To build motif signatures underlying kidney-specific enhancer activity, it would be ideal to have a database of known, previously characterized, urogenital (UG) enhancers so that we could learn the sequence preferences of such tissue-specific regulatory regions. However, due to the unavailability of such data, we instead utilize kidney-specific promoter sequences (like in Refs. 2, 13, or 15). Apart from this approach, we also examine a public dataset of histone-modified sequences of regulatory regions to find motif-signatures of genomic elements that are potentially enhancers. Though this data source is not kidney-specific, we observe that these epigenetic signatures have a strong, discriminative association with distal regulatory regions.
  - *Chromatin marks in known regulatory elements:* The ENCODE project suggests that mono-methylation of the lysine 4 residue of Histone *H3* is associated with enhancer (or distal regulatory) activity<sup>16</sup> whereas tri-methylation of *H3K4* and *H3* acetylation are associated with promoter activity. Using this set of *H3K4me1*, *H3K4me3* and *H3ac* sequences, we aim to find sequence motifs that are indicative of such epigenetic preferences during transcription. Though such epigenetic data is available for five different cell lines, we choose the HeLa cell line data because of its widespread use as a model system to understand transcriptional regulation *in vitro* in the laboratory.

For simplicity, we find the frequencies of six-nucleotide long motifs in the *H3K4me1* and *H3K4me3/H3ac* sequences. Then, we build a random forest (RF) classifier to discriminate monomethylated-*H3K4* sequences from trimethylated-*H3K4*/acetylated-*H3* sequences based on motif occurrence. We note that even though this data is not kidney cell-specific, it has favorable

specificity and sensitivity characteristics. The motifs thus obtained are putatively associated with epigenetic properties of proximal and distally located regulatory regions (such as enhancers), and are predictive of the regulatory potential of new sequences (Sec. 8).

- *Promoters of kidney-specific genes:* A catalog of kidney-specific mouse promoters is available from the GNF SymAtlas (<http://symatlas.gnf.org/>). This database contains a list of annotated genes and their expression in several tissue types, including the kidney. Since the proximal promoter of such kidney-specific genes harbors the transcriptional machinery for gene regulation, their sequences putatively have motifs that are associated with kidney-specific expression. Additionally, promoters that are spatio-temporally expressed during kidney development are also analyzed (MGI: <http://www.informatics.jax.org/>). The GNF dataset profiles mostly adult tissue-types. Since our goal is to study enhancer activity during nephrogenesis, we focus on genes expressed between day *e*10 and *e*12 in the developing kidney — such a list is obtained from the MGI database.

Without loss of generality, we use six-nucleotide motifs (hexamers) for characterizing these sequences. This is based on the observation that most transcription factor binding motifs have a 5–6 nucleotide core sequence with degeneracy at the ends of the motif. A similar strategy was introduced in Refs. 17 and 18. The main difference in our approach from such previous work is that differential hexamer analysis was done for the same class of sequences, and the statistical nature of the “test-set” is, by design, similar to the training set. That is, in Ref. 17, differential hexamers are found between known Cis-Regulatory Modules (CRMs) and non-CRMs, and used for the prediction of new CRMs from sequence. On the other hand, Ref. 18 deals with finding hexamer features of known promoters and using them to predict new promoters from sequence. In our case, however, we do not have enhancer sequence data (equivalent to CRMs) and we are using promoter sequence-data for the prediction of enhancers (CRMs) instead. Thus, the nature of the test sequence is very different. We demonstrate that our approach is partially useful in the discovery of putative enhancers from sequence. Also, the presented motif-finding approach does not depend on motif length and can be scaled, depending on biological knowledge.

We set up the motif discovery as a feature extraction problem from these tissue-specific promoter sequences and then build a random forest (RF) classifier to classify new sequences into tissue-specific and non tissue-specific categories based on these identified sequence features (motifs). Based on the motifs derived using a RF classifier algorithm we are able to accurately classify more than 95% (training-error rate) of tissue-specific genes based upon their upstream promoter region sequences alone. Since promoters are non-coding regulatory regions, the derived motifs can be putatively used to find kidney-specificity of other non-coding regions genome-wide (Sec. 9).

- (2) **Expression perspective:** There is limited expression data for the developing mouse kidney, mainly due to small tissue yield at such early time points. For this study, we use microarray expression data from a public repository of kidney microarray data (<http://genet.chmcc.org>, Ref. 19, <http://spring.imb.uq.edu.au/>).<sup>20</sup> Each of these resources contain expression data profiling kidney development from about day 10.5 dpc to the neonate stage. Such expression data can be mined for potential regulatory influence between upstream TF genes and *Gata2*.<sup>21,22</sup>
- *Inference of TF effectors at the promoter region:* The TFs putatively recruited at the proximal promoter are identified using the directed information (DTI) metric, that uses gene-expression (mRNA-level) influence in addition to phylogenetic conservation of the corresponding binding site. We have earlier shown that DTI is a good predictor of gene influence and can be used to infer transcriptional regulatory networks.<sup>22</sup>
  - *Inference of TF effectors at each non-coding region:* At the distal enhancer, it is believed that there is recruitment of tissue-specific transcription factors that co-operate with the basal transcriptional machinery (at the promoter) to direct tissue-specific gene expression.<sup>23,2</sup> Whereas phylogeny and expression-based influence metrics can yield high confidence candidates for promoter TFs, a similar analysis for enhancers is not possible, because of higher order effects.<sup>2,13</sup> To this end, the only way to search for putative enhancer TFs is to combine phylogeny with tissue-specific annotation (from UNIPROT or MGI). Hence, every transcription factor, whose motif is conserved at a non-coding (putative enhancer) region and is tissue-specific in annotation, is considered a likely candidate TF at that non-coding region.
- (3) **Interactome perspective:** The identification of phylogenetically conserved effector TFs at the promoter (identified via DTI), as well as those that are phylogenetically conserved at the putative enhancer candidate regions, lead to the exploration of protein-interactions (PPI) between these TFs, during distal enhancer-promoter interaction (Sec. 10). The STRING database (<http://string.embl.de>) integrates various experimental modalities (genomic context, high-throughput experiments such as co-immunoprecipitation, co-expression and literature) to maintain a list of organism-specific functional protein-association networks that is amenable to such exploration.

In this work, the above perspectives are examined in the context of the urogenital enhancers identified in Ref. 10. We aim to show that each of these modalities (tissue-specificity motifs, epigenetic signatures and TF-interaction graphs) has a predictive value for the identification of enhancers and the integration of these heterogeneous perspectives can lead to potential reduction in false positive rate during large-scale enhancer discovery, genome-wide. Such analyses can also be examined

in the context of new studies.<sup>24,25</sup> To date, there has been no comprehensive study for summarizing these various heterogeneous data sources to understand the characteristics of such regulatory regions.

### 3. Validation/Biological Application

As suggested in Sec. 1, we use the recently identified *Gata2* urogenital (UG) enhancers to validate our computational approach. All the data sources (and their analysis) are therefore going to be focused on the developing kidney.

The experimental characterization of these enhancers was done as follows. Based on BAC transgenic<sup>10</sup> studies, the approximate location of the urogenital enhancer(s) of *Gata2* were localized to a 70 kilobase region on chromosome 6. Using interspecies conservation plots, four elements were selected for transgenic analysis in the mouse. These were designated UG1, 2, 3 and 4. After a lengthy and resource-intensive experimental effort, two out of these four non-coding elements, *UG2* and *UG4*, were found to be true UG enhancers. Our goal is to find preferences at the sequence, expression and interactome level that can explain these experimental observations: i.e. that *UG2,4* are *Gata2*-specific urogenital enhancers and *UG1,3* are not urogenital enhancers for *Gata2*.

It is easy to see the utility of such an “*enhancer discovery*” methodology, since this can be applied also to other genes of interest. Given the complexity of 1% of the genome, made possible by the ENCODE project, the search for functional elements genome-wide is going to be an important and challenging exercise.

### 4. Organization

With a view to understanding the discriminating characteristics of transcriptional regulatory regions, the first part of this paper (Secs. 5–8) addresses identification of motif signatures representative of transcriptional control from kidney-promoter and epigenetically marked sequence sets. The second part of this work (Secs. 10.1–10.2) integrates phylogeny and expression data to find regulatory TFs at the proximal promoter and enhancer(s) of *Gata2*. Using the notion of TF interactions between enhancer and promoter, we examine if protein-interaction data (Sec. 10.3) can offer supporting evidence for the observed *in-vivo* behavior of the four *Gata2* candidate sequences. Classifiers are designed to discriminate regulatory versus non-regulatory regions based on these three modalities (kidney-specific motifs, epigenetic signatures and TF-interaction graphs). Finally, a probabilistic combination of these classifiers is done to obtain a validation (Sec. 11) of the *Gata2* UG enhancer (UGE) candidates (*UG1–4*). Sections 12 and 13 conclude the paper.

### 5. Sequence Data Extraction and Pre-Processing

Before proceeding to motif identification, a matrix of motif–chromatin-sequence correspondences is created. In this matrix, the counts of hexamer (six-nucleotide)



Table 1. The ‘motif count matrix’ for a set of histone-modified sequences. The first column is their genomic locations along the chromosome, the next two columns are hexamer quantile labels, and the last column is the corresponding sequence class label (+1/−1).

Sequence	AAAATA	AAACTG	Class
chr2:41410492-41411867	2	1	+1
chr6:41654502-41654782	4	2	+1
chr3:41406971-41408059	1	1	−1
chr2:41665970-41667002	2	3	+1
chr4:41476956-41478365	1	2	−1
chrX:41783327-41784532	1	2	+1

motif occurrence in the ‘ $H3K4me1$ ’ and ‘ $H3K4me3/H3ac$ ’ regions is obtained using sequence parsing (*R* package: ‘*seqinr*’). The motif length of six is not overly restrictive, and can be changed based on biological insight. A Welch t-test is then performed between the relative counts of each hexamer in the two epigenetic-modification categories (‘ $H3K4me1$ ’ and ‘ $H3K4me3/H3ac$ ’) and the top 1000 hexamers with  $\overrightarrow{p\text{-value}} \leq 10^{-6}$  are selected. This set of discriminating hexamers is designated ( $\overrightarrow{\mathbf{H}} = H_1, H_2, \dots, H_{1000}$ ). This procedure resulted in two hexamer-gene co-occurrence matrices, one for the ‘ $H3K4me1$ ’ (or +1) class of dimension  $N_{train,+1} \times 1000$  and the other for the ‘ $H3K4me3/H3ac$ ’ (or −1) class — dimension  $N_{train,-1} \times 1000$ . Here  $N_{train,+1}$  is the matrix of  $H3K4me1$  sequences corresponding to distal regulatory regions.  $N_{train,-1}$  is the set of ‘ $H3K4me3/H3ac$ ’ sequences that are associated with proximal promoters.

This dataset is obtained from the Sanger ENCODE database ([http://www.sanger.ac.uk/Post Genomics/encode/data-access.shtml](http://www.sanger.ac.uk/Post_Genomics/encode/data-access.shtml)) and contains 298 sequences that undergo modification ( $me1/me3/ac$ ) in histone ChIP assays. 140 of these correspond to  $H3K4me1$  (enhancers) and 158 correspond to  $H3K4me3/H3ac$  marks (promoters).

### 5.1. Kidney-specific promoter sequences

The Novartis foundation tissue-specificity atlas [<http://symatlas.gnf.org/>], has a compendium of genes and their corresponding tissues of expression. Genes have been profiled for expression in about 25 tissues, including adrenal gland, brain, dorsal root ganglion, spinal chord, testis, pancreas, liver, etc. Considering these diversity of tissue-types, one concern with the interpretation of this data is the variability in expression across tissue-types. To address this concern, we take a fairly stringent approach — if a gene is expressed in less than three tissue types, it is annotated tissue-specific (‘*ts*’), and if it is expressed in more than 22 tissue types, it is annotated non-specific (‘*nts*’). Based on this assignment, we find a list of 86 genes that are tissue-specific as well as have kidney expression (MGI: <http://www.informatics.jax.org/>). For these kidney-specific genes, we extract their promoter sequences from the ENSEMBL database (<http://www.ensembl.org/>),

using sequence 2000bp upstream and 1000bp downstream up to the first exon relative to the transcriptional start site reported in ENSEMBL (release 37).

## 6. Motif-Class Correspondence Matrices

From the above,  $N_{train,+1} \times 1000$  and  $N_{train,-1} \times 1000$  dimensional count matrices are available for the chromatin-modified sequences. Before proceeding to the feature (hexamer motif) selection step, the counts of the  $M = 1000$  hexamers in each training sample are normalized to account for variable sequence lengths. In the co-occurrence matrix, let  $gc_{i,k}$  represent the absolute count of the  $k^{th}$  hexamer,  $k \in 1, 2, \dots, M$  in the  $i$ th chromatin-sequence. Then, for each sequence  $g_i$ , the quantile labeled matrix has  $X_{i,k} = l$  if  $gc_{i, [\frac{l-1}{K}M]} \leq gc_{i,k} < gc_{i, [\frac{l}{K}M]}$ ,  $K = 4$ . Matrices of dimension  $N_{train,+1} \times 1001$ ,  $N_{train,-1} \times 1001$  for the specific and non-specific training samples are now obtained. Each matrix contains the quantile label assignments for the 1000 hexamers ( $X_i, i \in (1, 2, \dots, 1000)$ ), as stated above, and the last column would have the corresponding class label ( $Y = -1/+1$ ). Having constructed two groups of sequences for analysis — enhancer-associated (*'H3K4me1'*) and promoter-associated (*'H3K4me3/H3ac'*) — we seek to find the smallest set of hexamer motifs that are most discriminatory between these two classes. Towards this goal, we use random forest classifiers (RF)<sup>26</sup> for finding such a discriminative hexamer subset.

Based on the above strategy for epigenetically marked sequences, we follow the same procedure, from sequence extraction, parsing and quantization to obtain hexamer-promoter counts for the kidney-specific gene promoter sequences. As an illustration, we show a representative matrix (Table 2).

## 7. Random Forest Classifiers

A random forest (RF) is an ensemble of classifiers obtained by aggregating (bagging) several classification trees.<sup>26</sup> Each data point (represented as an input vector) is classified based on the majority vote gained by that vector across all the trees of

Table 2. The ‘motif count matrix’ for a set of gene-promoters. The first column is their ENSEMBL gene identifiers, the next two columns are hexamer quantile labels, and the last column is the corresponding gene’s class label (+1/−1).

Ensembl Gene ID	AAAAAA	AAATAG	Class
ENSG00000155366	1	1	+1
ENSG000001780892	4	3	+1
ENSG00000189171	1	2	−1
ENSG00000168664	4	3	−1
ENSG00000160917	2	1	−1
ENSG00000176749	1	1	−1
ENSG00000006451	3	2	+1

the forest. Each tree of the forest is grown in the following way:

- A bootstrapped sample (with replacement) of the training data is used to grow each tree. The sampling for bootstrapped data selection is done individually at each tree of the forest.
- For an  $M$ -dimensional input vector, a random subspace of  $m$  ( $\ll M$ )-dimensions is selected, and the best split on this subspace is used to split the node. This is done for all nodes of the tree.

During the training step, before sampling by replacement, one-third of the cases are kept “out of the training bag.” This OOB (out-of-bag) data is used to obtain an unbiased estimate of the classification error as trees are added to the forest. It is also used to get estimates of variable importance.

Several interesting insights into the data are available using random forest analysis. The variable importance plot yields the variables that are most discriminatory for classification under the ‘ensemble of trees’ classifier. This importance is based on two measures: ‘Gini index’ and ‘decrease in accuracy.’ The Gini index is an entropy-based criterion which measures the purity of a node in the tree, while the other metric simply looks at the relative contribution of each variable to the accuracy of the classifier. For our studies, we use the ‘randomForest’ package for R. The classifier performance on the individual data and the related diagnostics are mentioned under Sec. 8.

## 8. Random Forests on Chromatin-Modified Sequences

We train the RF classifier on the set of 298 chromosome sequences that have varying chromatin modifications associated with them (i.e. *H3K4me1/me3*, and *H3ac*), as mentioned in Sec. 2. These are derived from the HeLa cell line and are not necessarily context-specific for kidney development. However, given the widespread use of this cell line for transcriptional studies, we aim to find if the motifs associated with regulatory elements are indeed predictive of enhancer activity.

Before proceeding to motif identification, we check for possible sequence bias (such as GC-nucleotide composition) between these two classes of chromatin modified sequences. If there is a significant bias, then the motifs turn out to be just GC-rich sequences that are not very biologically informative for determination of regulatory potential. The GC composition of these two classes of sequences is represented in Fig. 3. As can be seen, the average GC composition is the same and there is no such sequence bias that would skew the discovery and subsequent interpretation of these epigenetic motifs. The performance of the histone-RF classifier is explained in the context of the classifier combination in Sec. 11.

The motifs obtained from the random forest analysis indicate the “sequence-preferences” of regulatory elements that are nucleosome-free in HeLa cells (Fig. 4). We analyze the performance of these classifiers on the 4 UG candidate regions, mentioned previously. In both cases, *UG2* – 4 are classified as enhancers, whereas

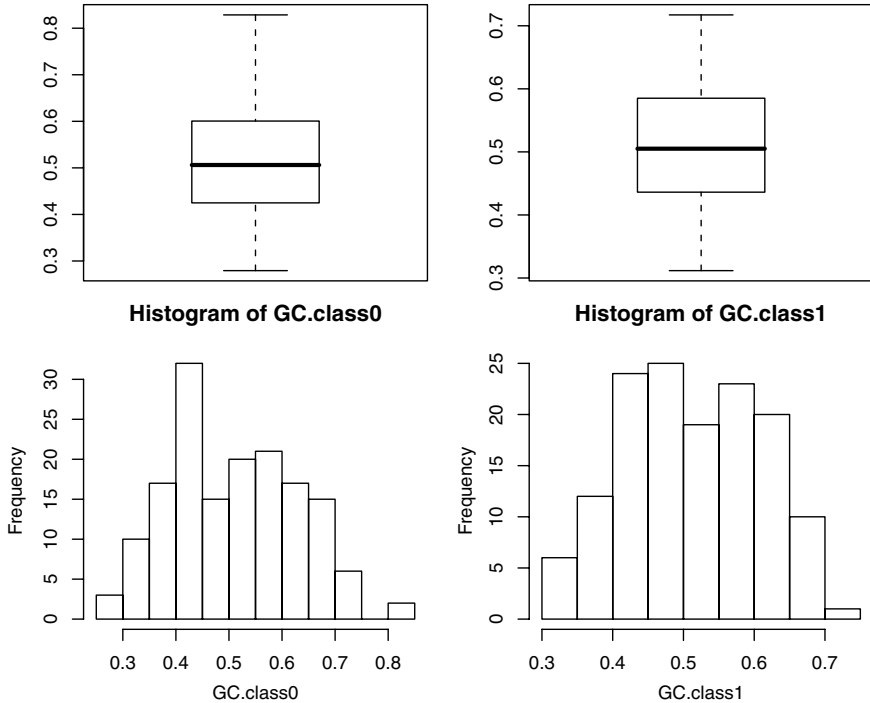


Fig. 3. GC plots for sequence bias in *H3K4me1* histone sequences versus *H3K4me3* and *H3ac* sequences. We observe that there is no significant bias in GC content.

*UG1* is correctly classified as not being regulatory. Additionally, a control set of “promoter-independent” enhancers derived from the Mouse Enhancer database<sup>1</sup> was also classified as enhancers based on these chromatin-sequence motif signatures. This high prediction accuracy in spite of non-specificity of cell context (*HeLa* cell line) is very interesting and has potentially high predictive value.

## 9. Random Forests on Kidney-Specific Promoters

In this section, we aim to find discriminating sequence motifs between a set of kidney-specific promoters and housekeeping promoters with a goal to find sequence motifs underlying kidney-specific regulation. The kidney enriched dataset has 86 genes that are assigned to a tissue specific class and have higher than mean expression in the kidney. For the purpose of training and testing, we consider the set of housekeeping genes identified from the ‘*nts*’ class and reported in literature.<sup>27,28</sup> There are almost 1500 genes in the housekeeping gene (‘*nts*’) set. Since this would lead to unbalanced predictions during classifier training, we use a stratified sampling approach<sup>29</sup> to select a sample size that reduces this effect (the sampling itself is done with a prior on the relative sizes of the two classes). Here, the set of (−1) promoter-sequences are taken to be of the same size as the (+1) class. Using this

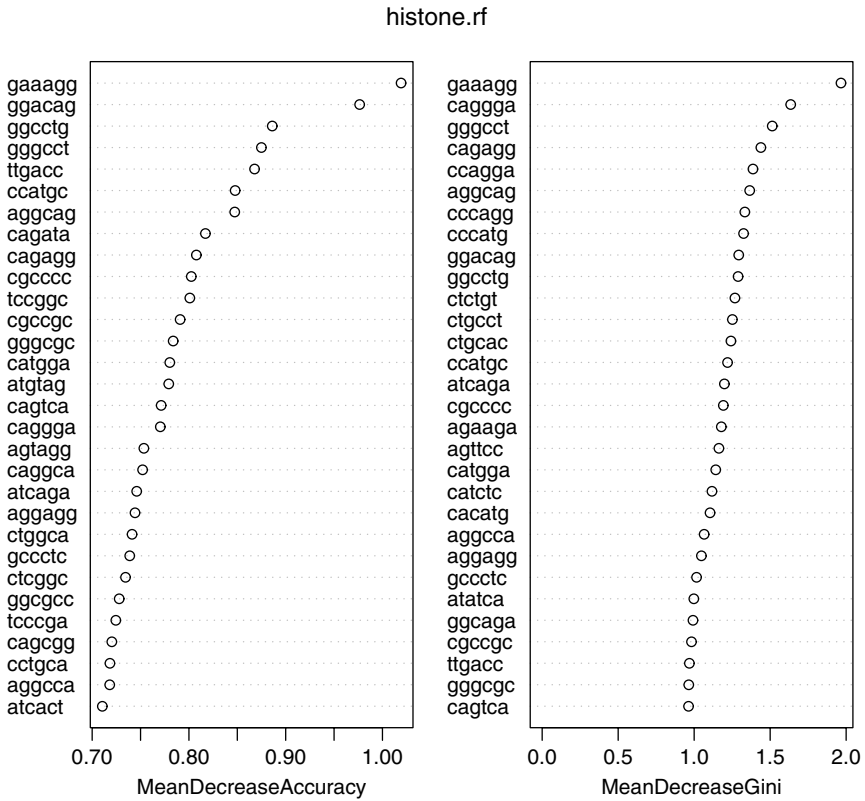


Fig. 4. Top hexamers which can discriminate between *H3K4me1* histone sequences versus *H3K4me3* and *H3ac* sequences.

approach, we obtain a training-error classification accuracy of > 95% on the kidney enriched tissue-specificity data set. Before proceeding to motif identification, it is necessary to check for possible sequence bias (GC composition) between the two classes of promoters (kidney-specific versus housekeeping). The GC composition of these two classes of sequences is represented in Fig. 5. We note that though only a subset of ‘nts’ gene-promoters were used during the RF analysis, we show the GC-composition for the entire class of ‘nts’ sequences for completeness. As can be seen, the average GC composition is the same. The ROC space representation and variable importance plot for the overall classification is indicated below (Fig. 10, represented by (·) and Fig. 6, respectively). The confusion matrices are all explained in the context of the classifier combination in Sec. 11.

To address a related question, we examine if the top ranked hexamers in the kidney dataset correspond sequence-wise to known transcription factor binding sites. Using the publicly available Opossum tool (<http://www.cisreg.ca/cgi-bin/oPOSSUM/opossum/>) or MAPPER (<http://bio.chip.org/mapper/>), we found

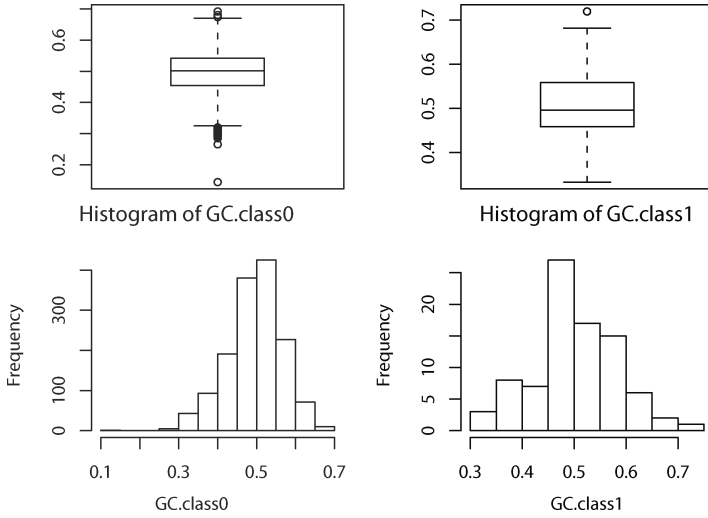


Fig. 5. GC plots for sequence bias in kidney-specific versus housekeeping promoters.

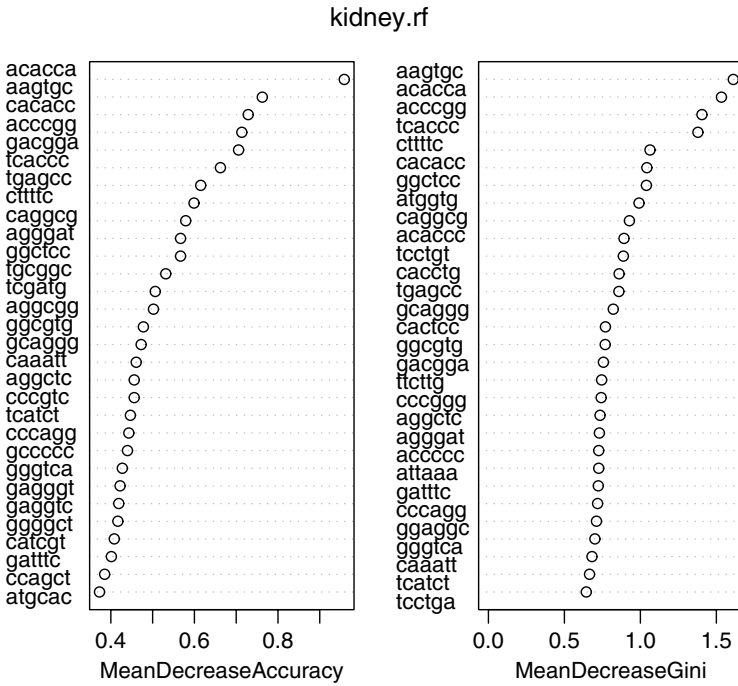


Fig. 6. Top hexamers which can discriminate between kidney-specific and housekeeping genes.

several interesting transcription factors to map to these motifs, such as *Nkx*, *ARNT*, *c-ETS*, *FREAC4*, *NFAT*, *CREBP*, *E2F*, *HNF4A*, *Pax2*, *MSX1*, *SP1*, several of which are kidney-specific. Though this is highly consistent with the tissue-specificity of the dataset, the functional relevance of these sites remains to be experimentally validated.

## 10. PPI Between Promoter and Enhancer TFs

In order to understand the nature of interactions between the enhancer and promoter TFs (Fig. 2), we decouple the overall regulation problem into three parts:

- (1) Identification of putative TF effectors at the promoter (Sec. 10.1),
- (2) Identification of enhancer TFs (Sec. 10.2), and
- (3) Examination of the interaction-graph formed between enhancer-TFs and promoter TFs (Sec. 10.3).

The key question that is explored in the following sections is: having identified the set of tissue-specific TFs that might putatively bind the promoter and the candidate regulatory regions, does the *structure of the bipartite TF-interaction graph* (across the promoter TFs and the enhancer TFs) reveal any interesting features that distinguish the functional *UG2, 4* regions from the non-functional *UG1, 3* regions?

### 10.1. TF effector identification at promoter and enhancer

*Promoter TF identification:* TFs that regulate basal transcription at the promoter can be identified from phylogenetic conservation or co-expression studies. In this approach, the promoter sequence (here, the *Gata2* promoter) is aligned across multiple species and the TFBS motifs that are conserved in the multiple alignment are considered to be putative effectors of gene regulation. Such sequence-based approaches have been examined in literature.<sup>2,13</sup>

Since the list of putative TFs (identified above) that potentially bind at the promoter is still large, there have been efforts to incorporate gene-expression data to reduce the set of potential TF effectors. In this scenario, if the gene corresponding to the conserved TF has a high expression-level influence on *Gata2* expression, then that TF has stronger evidence for being a potential regulator.<sup>21</sup>

Recently, we introduced the directed information (DTI) as a metric to infer expression-level influence between any putative transcription factor (TF) gene and a target gene (such as *Gata2*).<sup>22</sup> This seeks to integrate sequence and expression data into the determination of relationships between transcription factors and their target-genes. All additional details (performance on synthetic data, other biological data and comparison with other metrics) are available in Ref. 22. Information-based measures have enabled the investigation of non-linear gene relationships in the presence of measurement noise.<sup>21</sup> An important point to note is that unlike mutual information, the DTI is a *directed* metric that enables the determination

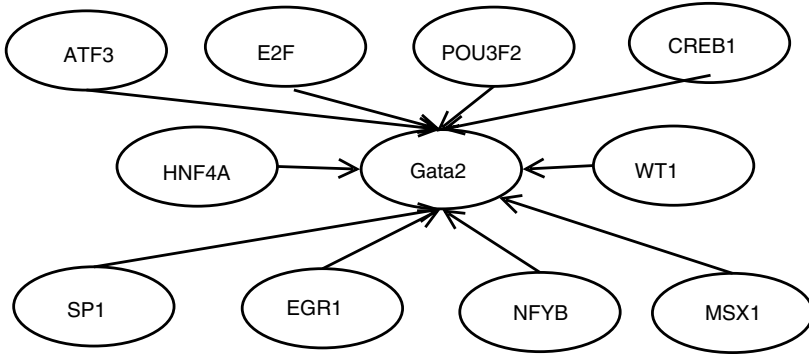


Fig. 7. Putative upstream TFs using DTI for the *Gata2* gene.

of the strength, significance and direction of gene influence. For *Gata2*, this list of effectors is listed in Fig. 7.

### 10.2. Enhancer TF identification

In Sec. 10.1, we examined the identification of promoter TFs using phylogenetic sequence conservation of TFBS motifs in conjunction with expression level influence using DTI. The next key step towards determining the structure of promoter-enhancer TF interactions is the identification of enhancer-TFs. As has been alluded to earlier, there is no method to precisely infer which transcription factors bind a certain regulatory element during long-range gene regulation. Thus, we appeal to a traditional approach of finding tissue-specific transcription factors that are phylogenetically conserved at any potential regulatory region<sup>11,2</sup> (one caveat, however, is that conservation is not a very reliable predictor of TF binding<sup>30,31</sup>). This is consistent with earlier observations that enhancers recruit tissue-specific transcription factors during the formation of the overall transcriptional machinery during gene expression, whereas promoters recruit components of the basal transcriptional machinery.<sup>23,2,13,4</sup>

To ascertain the tissue-specificity of each TF that putatively binds a regulatory element (identified via phylogenetic conservation), we examine that TF's annotation in the UNIPROT or MGI database.

### 10.3. Enhancer-promoter distal interaction via protein-protein interactions — a graph-based analysis

Using the notion of protein-protein interaction (PPI) mediating long-distance interactions between promoters and enhancers during looping,<sup>3,14,32</sup> we explore the interactome to look for within-group and between-group interactions in the promoter-TF and the enhancer-TF groups.



The interaction-graphs (e.g. Fig. 8) are obtained in the following manner:

- One part of the graph (hollow circles) corresponds to the TF effector group at the promoter. These  $V_p$  TFs are identified based on phylogenetic conservation, tissue-specificity and directed information (Sec. 10.1).
- The other part of the graph (filled circles) corresponds to the  $V_e$  tissue-specific TFs group at the enhancer, identified based on phylogeny and tissue-specificity annotation (Sec. 10.2).
- The interaction-graph is defined by the vertices  $V = (V_p \cup V_e)$ , and the edges  $E = e_{i,j}$ ,  $i, j \in (1, 2, \dots, |V_p \cup V_e|)$ . Each bidirectional edge  $E = (e_{i,j})$  is derived from an annotated interaction between TFs  $i$  and  $j$ , based on an interaction database. These edges describe both within-group TF interactions as well as between-group interactions. These interactions are obtained from the STRING (<http://string.embl.de/>) and MiMI (<http://mimi.ncibi.org/MiMI/home.jsp>) databases, both of which contain data derived from multiple sources, such as yeast-2-hybrid screens, literature etc.

It would be of great value to use a catalog of gene-specific and tissue-specific regulatory regions (with all possible transcription factors) from which to find such interaction characteristics. However, such a repository does not yet exist. In this section, we use a few examples (*Gata3* OVE, *Gata3* KE, *Fgf* OVE, *Mecp2* F21/F6, *Shh* FE) of known tissue-specific and gene-specific regulatory elements from literature, as a positive training set. For the negative training set, we consider the set of regions that were reportedly investigated in these transgenic experiments but did not yield gene-specific regulatory activity.

We have presented a preliminary analysis of enhancer-promoter TF interaction-graphs for some genomic elements with known regulatory or non-regulatory activity<sup>6,33,8,34</sup> in Table 3. The table represents the listing of some of the structural attributes of these interaction-graphs, following analysis methods from literature.<sup>35</sup>

Table 3. The first column is the various regulatory and non-regulatory elements from literature, the next column corresponds to its class label (+1/ - 1). The subsequent columns correspond to the attributes of the overall TF-interaction graph (both within-group and between-group interactions).

Sequence	Class	Clustering coefficient	Characteristic path length	Heterogeneity	Centralization	Density
<i>Mecp2</i> F21 <sup>6</sup>	+1	0.208	2.824	0.668	0.184	0.133
<i>Mecp2</i> F6 <sup>6</sup>	-1	0	1.75	0.342	0.067	0.145
<i>Gata3</i> OVE <sup>8</sup>	+1	0.036	2.254	0.779	0.359	0.154
<i>Gata3</i> KE <sup>8</sup>	+1	0.409	2.0	0.813	0.684	0.216
<i>Gata3</i> NE1 <sup>8</sup>	-1	0.383	2.131	1.139	0.757	0.15
<i>Gata3</i> NE2 <sup>8</sup>	-1	0.458	2.013	0.872	0.699	0.203
<i>Fgf10</i> OVE <sup>34</sup>	+1	0.313	2.433	0.72	0.323	0.133
<i>Shh</i> FE <sup>33</sup>	+1	0.394	2.312	0.797	0.49	0.175

A deeper analysis of other graph topology metrics and their relation to functional enhancer activity is a topic of future interest.

- *Clustering coefficient*: The clustering coefficient of a node is always a number between 0 and 1. The network clustering coefficient is the average of the clustering coefficients for all nodes in the network.
- *Characteristic path length*: The characteristic path length denotes the average shortest-path distance of the graph. This gives the expected distance of any two connected nodes in the graph and is a global indicator of network-connectivity.
- *Heterogeneity*: Network heterogeneity denotes the coefficient of variation of the degree distribution.
- *Centralization*: This refers to the overall connectivity (cohesion) of the graph. It indicates how strongly the graph is organized around its most central point(s).
- *Density*: It shows how densely the network is populated with edges (i.e. how “close-knit” an empirical graph is). A network which contains no edges and solely isolated nodes has a density of 0, whereas the density of a clique (completely connected graph) is 1.

The above-mentioned network properties (as well as clustering coefficients, number of connected components etc.) are examined for the overall interaction-graphs for the reported enhancers from literature. A logistic regression reveals that low values of heterogeneity, characteristic path length and centralization are strong predictors of potential enhancer activity. All of these attributes point to the decentralized, homogenous and somewhat tighter connectivity of the interaction-graphs for true enhancers. We note that the OOB error rate of the RF here is about 20%. The quality of this classifier can be expected to improve as we obtain more data (gene-specific regulatory regions) from which to extract features.

We now examine the interaction-graphs for the test set, i.e. the four *Gata2* UGE candidates. For illustration, we only show the largest connected component of the inter-group edges for each interaction graph (Fig. 8). For comparison, we have also shown the interaction graphs between the UG candidates and the promoters proximal to the *Gata2* promoter (*Rpn1*, *Rab7*, *Eefsec*, *Dnajb8*). We observe that the interaction densities for the proximal promoters are very low in comparison to the density for *Gata2*, and that there are no strongly connected components in the interaction graphs for the *Gata2* proximal promoters, suggesting a high specificity of interaction between the true UG candidates (UG2/UG4) and the *Gata2* promoter, Fig. 9, and a corresponding low specificity of interaction for the proximal promoters. This demonstrates the utility of this approach to the resolution of promoter-specificity of enhancer action.

This figure indicates a very interesting property of the real enhancers vis-a-vis the other conserved elements. We see that the TF effectors for *Gata2* such as *SP1*, *POU3F2* (identified in the TF effector network above, Fig. 7), are involved in cross-element interactions at the protein level, between the promoter and true enhancers (UG2/4). However, the network linkage in the elements that showed no enhancer

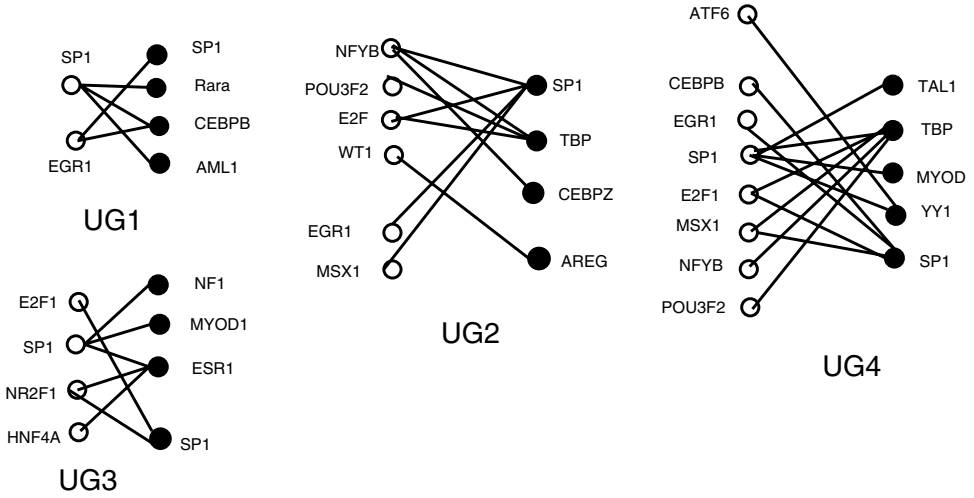


Fig. 8. Protein-protein interaction between putative *Gata2* TFs (hollow circles) and putative UG element TFs (filled circles). Note: This only shows the connections between two groups for one of the connected components. For our analysis, we consider both *intra*- and *inter*-group connections. From <http://string.embl.de/>

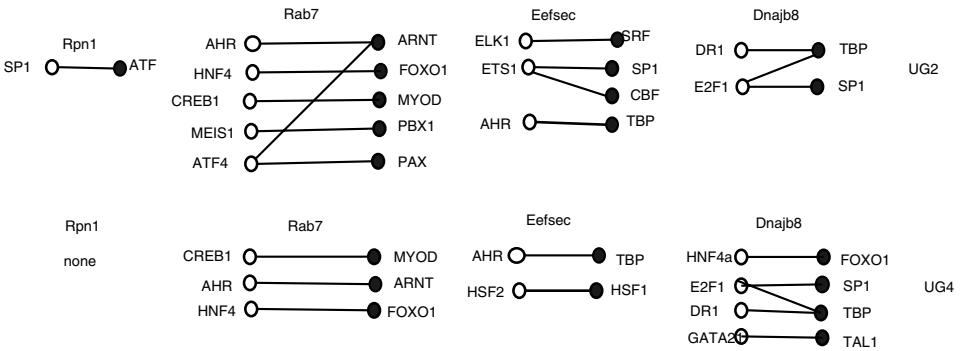


Fig. 9. Protein-protein interaction between *Gata2* proximal promoter TFs (hollow circles) and putative UG element TFs (filled circles). Here, we observe that the interaction densities are very low, in comparison to Fig. 8. From <http://string.embl.de/>

activity is very sparse, suggesting low crosstalk between promoter and enhancer. Also, the TFs at the enhancer nodes (dark circles) have a more uniform degree distribution in the functional elements *UG2/4* as compared to the non-functional ones. Both these observations suggest lower heterogeneity and centralization of such functional interaction-graphs. Thus, the extent of TF crosstalk is a potential discriminator of possible enhancer function. This shows that superimposing such PPI information along with sequence and expression data helps reduce the number of false positives while integrating various aspects of distal regulation.

## 11. Heterogeneous Data Integration and Validation on GATA2 UGE Candidate Sequences

As mentioned previously, the primary goal of the framework developed above is to understand the behavior of known regulatory elements along different genomic modalities. To validate their predictive potential, we demonstrate their application to predicting the behavior of the experimentally-verified *Gata2* UG enhancer candidates (which is our test set). Here we combine the results of the individual classifiers (kidney-promoter RF, histone RF and interactome-RF) to obtain an integrated prediction that a candidate sequence is an enhancer. For combining heterogeneous classifiers, we use a “probabilistic belief fusion” approach.

The framework involves combining the ‘beliefs’ of the individual classifiers to obtain a combined belief of prediction. To compute the belief of each classifier we start by examining the confusion matrices for each of the classifiers (promoter RF, histone-RF and graph-RF), following Ref. 38. Since each of the classifiers are random forests, we can obtain their OOB error estimates through these confusion matrices. For the graph-RF, this confusion matrix is as below,

$$\mathbf{CM}_{\text{graph-RF}} = \begin{pmatrix} \text{Class} & -1 & 1 & \text{class.error} \\ -1 & 4 & 1 & 0.20 \\ 1 & 1 & 4 & 0.20 \end{pmatrix},$$

thereby yielding an OOB error estimate of  $\sim 20\%$ .

Similarly, we have,

$$\mathbf{CM}_{\text{promoter-RF}} = \begin{pmatrix} \text{Class} & -1 & 1 & \text{class.error} \\ -1 & 67 & 19 & 0.22 \\ 1 & 10 & 76 & 0.12 \end{pmatrix},$$

thus yielding an OOB error estimate of  $\sim 17\%$ .

$$\mathbf{CM}_{\text{histone-RF}} = \begin{pmatrix} \text{Class} & -1 & 1 & \text{class.error} \\ -1 & 134 & 24 & 0.15 \\ 1 & 21 & 119 & 0.15 \end{pmatrix},$$

yielding an OOB error estimate of  $\sim 15\%$ .

As can be seen, these classifiers have fairly good sensitivity and specificity characteristics. However, we note that each of the modalities might be imbalanced in class membership in the original study, and so might not be as generalizable. This is expected to improve as more training data for these classifiers becomes available (especially for the graph-RF case). Moreover, these are three complementary data sources and can be effectively combined to improve detection reliability. Since they are trained on very different modalities, they can be assumed to be independent. It can also be seen that this method of belief combining is applicable to as many modalities ( $K$ ) as necessary to the biological problem of interest, and hence is truly scalable.

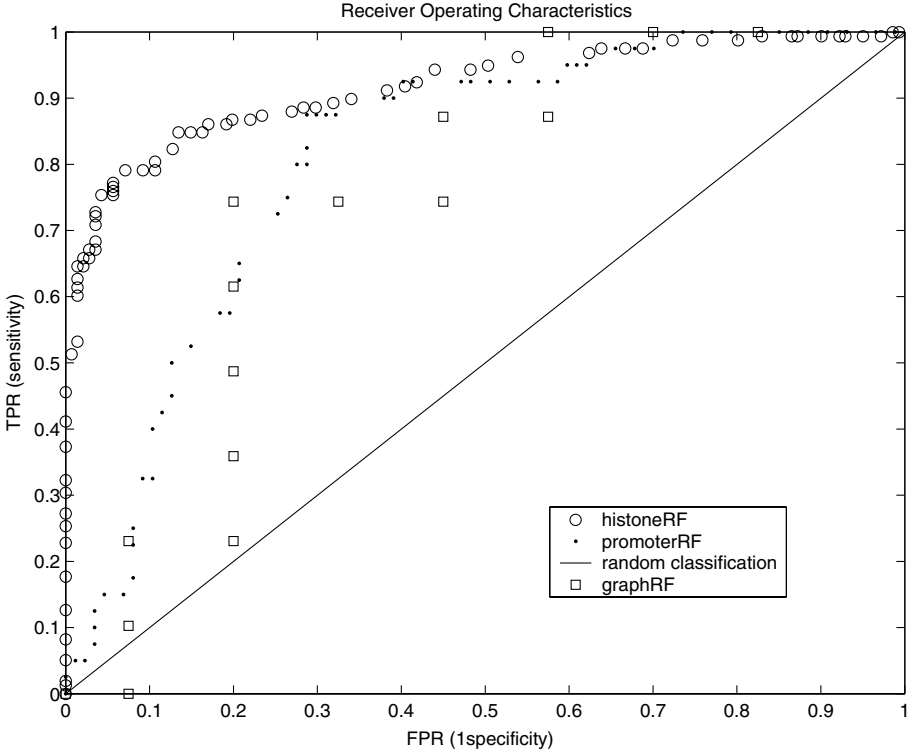


Fig. 10. Representation of the three RF classifiers in ROC space (RF-histone in (○), promoter-RF in (•), and graph-RF in (□)). The diagonal line is the classification by random chance.

Let each classifier be characterized by its decision function  $e_k(x) = j_k$  that maps a data point ( $x$ ) to the class ‘ $j$ ’, for  $k = 1, 2, \dots, K$  and  $j_k \in (-1, 1)$ . Here,  $K = 3$ , and  $J = 2$  classes.

The belief of the  $k$ th classifier is defined as

$$bel_k(x \in C_i | e_k(x) = j_k) = P(x \in C_i | e_k(x) = j_k)$$

The overall belief,  $bel(i)$ , is computed using Bayes rule

$$bel(i) = P(x \in C_i) \cdot \frac{\prod_{k=1}^K P(x \in C_i | e_k(x) = j_k)}{\prod_{k=1}^K P(x \in C_i)}$$

$$bel(C_i) = \frac{\prod_{k=1}^K P(x \in C_i | e_k(x) = j_k)}{\sum_{i=1}^J \prod_{k=1}^K P(x \in C_i | e_k(x) = j_k)}$$

Table 4. Combined belief generation during heterogeneous classifier integration. The last column represents the combined belief (probability that the UG candidate sequence is an enhancer) as a result of integrating the promoter-RF, histone-RF and graph-RF predictions.

Sequence	True Class	Promoter RF prediction $e_1(x)$	Histone RF prediction $e_2(x)$	Interaction-graph RF prediction $e_3(x)$	P(Class = +1) (Overall Belief)
<i>Gata2</i> UG1	-1	-1	-1	-1	0.0067
<i>Gata2</i> UG2	+1	+1	+1	+1	0.989
<i>Gata2</i> UG3	-1	-1	+1	-1	0.432
<i>Gata2</i> UG4	+1	+1	+1	+1	0.9875

Note:  $J = 2$  and  $K = 3$ . Depending on the belief value  $bel(i)$ , the decision rule ( $E(x)$ ) for classifying data point  $x$  is

$$E(x) = j, \quad \text{if } bel(j) = \max_i bel(i), \quad \text{or,}$$

$$E(x) = j, \quad \text{if } bel(j) = \max_i bel(i), \quad \text{and } bel(j) \geq \alpha,$$

where  $0 < \alpha \leq 1$ , with  $\alpha$  being a threshold.

We now show the output classes of each of the three classifiers as well as the combined belief on the *Gata2* UG enhancer candidates in Table 4. More specifically, for the first row in Table 4, the overall belief equation above becomes

$$bel(ug1 = +1) = \frac{\prod_{k=1}^K P(ug1 = +1|e_k(x) = j_k)}{\prod_{k=1}^K [P(ug1 = +1|e_k(x) = j_k)] + \prod_{k=1}^K [P(ug1 = -1|e_k(x) = j_k)]}$$

$$= \frac{\prod_{k=1}^K (1 - prec_{n,k})}{[\prod_{k=1}^K (1 - prec_{n,k}) + \prod_{k=1}^K prec_{n,k}]}$$

Here,  $prec_{n,k} = \frac{TN_k}{TN_k + FN_k}$ . Similarly,  $prec_{p,k} = \frac{TP_k}{TP_k + FP_k}$ . These are the negative and positive precision values respectively, for the  $k$ th classifier. These rates are obtained from the corresponding confusion matrices shown above. This approach is followed for each of the *UG1* – 4 elements (Table 4).

If we set a threshold of  $\alpha = 0.50$  or  $0.90$ , we would get *UG2* and *UG4* as the true enhancers (100% accuracy). However, for a choice of  $\alpha = 0.40$ , *UG3* is predicted to be an enhancer in spite of it being declared a member of the (-1) class by the graph-RF. This choice of threshold thus determines the performance of the combined classifier (just like in any other hypothesis-testing scenario). We note that at the present time, there is no known repository of promoter-specific regulatory elements to carry out such graph-analysis on each element.

Under the  $\alpha = 0.40$  case, however, the results are not to be interpreted as a 25% error rate since the nature of the test set (*Gata2* UG enhancers) are very different from the training data of each modality (histone sequences are for a different cell-context; and interaction-graphs are obtained over different genes). The fact that we are getting such good prediction in spite of the training sets being so different

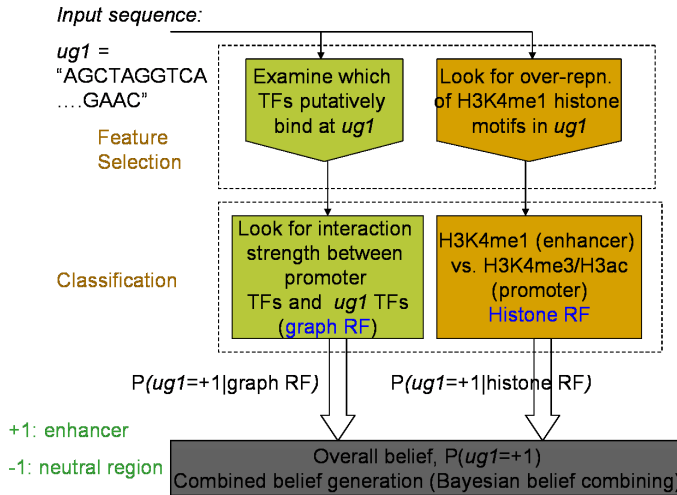


Fig. 11. Data integration across multiple modalities and combined belief generation. Given a query CSE sequence (*ug1*), and the tissue of interest, the classifier can use the chromatin-modification and interaction-graph modalities to output a combined belief for *ug1* to be an enhancer.

is a strong point in favor of examining and integrating these data sources. The test-error rates are given by the OOB error estimates of the individual classifiers.

## 12. Summary of Approach

In this work, we have shown that,

- Tissue-specificity motifs are useful for regulatory element identification. In spite of an unbiased search, they discover motifs that are regulatory (such as TFBS motifs) and potentially predictive.
- Chromatin modification motif signatures are predictive of regulatory element location. These point to the cell-specific epigenetic preferences of distally located regulatory regions.
- Promoter and enhancer TFs that are putatively recruited during gene (*Gata2*) regulation can be identified using a combination of phylogenetic conservation, expression data, and tissue-specificity annotation.
- Effector TFs at the gene proximal promoter have high network linkage with enhancer TFs in the case of functional enhancers. The TF interaction-graphs of truly functional elements are seen to have a lower centralization, characteristic path length and heterogeneity, suggesting higher crosstalk during formation of the transcription factor complex.

These diverse perspectives (based on sequence, expression and interactome data) shed some light on the sequence and mechanistic preferences of true regulatory

regions interspersed genome-wide. It is to be noted that this model is data-driven and needs further validation to correspond directly with the biology of transcription.

### 13. Conclusions

The novelty of the proposed work spans several areas. Firstly, data sources that are relevant to understanding the mechanism of gene regulation (with *Gata2* as an example) have been identified. We have developed methods that reconcile the behavior of known regulatory elements along each of these modalities. The utilization of histone-modified sequences and their exploration for sequence motifs are indicative of epigenetic preferences and nucleosome-occupancy patterns. This has not been explored before for the prediction of distal regulatory regions. The use of DTI as a metric to infer putative TF to target-gene influence is a recent one that serves to integrate phylogenetic TFBS conservation with expression data. Finally, the utilization of graph-based analysis techniques to understand the “structure” of the TF interaction-graph between enhancer and promoter helps us understand true enhancer behavior from a mechanistic viewpoint. The probabilistic combination of multiple classifiers (each deriving from a unique data resource) aims to reconcile the behavior of existing enhancers along multiple modalities. We hope to demonstrate that a principled integration of non-overlapping genomic modalities can be used to interpret the context and specificity of gene regulation.

### 14. Future Work

Some key elements directly emerge for guiding future research. As already alluded to in the motif-signature procedure, specific expression data corresponding to stages and tissues of interest would greatly improve the specificity of regulatory element prediction. Furthermore, as histone modification maps for related cell lines are generated, the false positive rate of prediction would decrease, thereby improving accuracy. Several other learning paradigms can be introduced into this setting since we are learning from structured data. Also, methods in joint classifier and feature optimization might likely improve the accuracy of predictions. Additionally, methods that analyze the grammar of these cis-regulatory regions (LREs) and look for motif position, spacing and orientation will be of great utility.

At the expression level, methods for supervised network inference would have a great impact on the discovery of TF effectors. Rapid advances have been made in this area and their relevance to the biological context of the problem has become very principled. At the interactome level, the work presented here can be extended to the investigation of graph-clusters for weighted interaction-graphs. The weighted edges are obtained from the confidence of the individual data sources, as well as the number of species over which that particular edge is conserved.<sup>35,36</sup> Such analysis enables the discovery of subgraphs of various degrees of inter-connectedness, thereby discovering functional “graph-motifs.”



An important point to note here is that there is currently no resource for promoter-specific enhancer data genome-wide. However, as various high throughput experiments become more prevalent, we can look forward to using these methods for precision-recall analysis on such public repositories.

## Acknowledgments

We thank Ms. Swapna Jayaraman for useful discussions about network analysis. We are also very grateful to the two anonymous reviewers for their help in revising the conference manuscript.

## References

1. Pennacchio LA, Ahituv N, Moses A, Prabhakar S, Nobrega M, Shoukry M, Minovitsky A, Dubchak I, Holt A, Lewis K, Plazer-Frick I, Akiyama J, DeVal S, Afzal V, Black B, Couronne O, Eisen M, Visel A, Rubin EM, *In vivo* enhancer analysis of human conserved non-coding sequences, *Nature* **444**(7118):499–502, 2006.
2. MacIsaac KD, Fraenkel E, Practical strategies for discovering regulatory DNA sequence motifs, *PLoS Comput Biol* **2**(4):e36, 2006.
3. Petrascheck M, Escher D, Mahmoudi T, Verrijzer CP, Schaffner W, Barberis A, DNA looping induced by a transcriptional enhancer *in vivo*, *Nucleic Acids Res* **33**(12):3743–3750, 2005.
4. Simonis M, Kooren J, de Laat W, An evaluation of 3C-based methods to capture DNA interactions, *Nature Methods* **4**(11):895, 2007.
5. Fraser P, Transcriptional control thrown for a loop, *Curr Opin Genet Dev* **16**(5):490–495, 2006.
6. Liu J, Francke U, Identification of cis-regulatory elements for MECP2 expression, *Hum Mol Genet* **15**(11):1769–1782, 2006.
7. Lakshmanan G, Liew KH, Lim KC, Gu Y, Grosveld F, Engel JD, Karis A, Localization of distant urogenital system-, central nervous system-, and endocardium-specific transcriptional regulatory elements in the GATA-3 locus, *Mol Cell Biol* **19**:1558–1568, 1999.
8. Hasegawa SL, Moriguchi T, Rao A, Kuroha T, Engel JD, Lim KC, Dosage-dependent rescue of definitive nephrogenesis by a distant Gata3 enhancer, *Dev Biol* **301**(2):568–577, 2007.
9. Lee EC, Yu D, Martinez de Velasco J, Tessarollo L, Swing DA, Court DL, Jenkins NA, Copeland NG, A highly efficient Escherichia coli-based chromosome engineering system adapted for recombinogenic targeting and subcloning of BAC DNA, *Genomics* **73**:56–65, 2001.
10. Khandekar M, Suzuki N, Lewton J, Yamamoto M, Engel JD, Multiple, distant Gata2 enhancers specify temporally and tissue-specific patterning in the developing urogenital system, *Mol Cell Biol* **24**(23):10263–10276, 2004.
11. Pennacchio LA, Loots GG, Nobrega MA, Ovcharenko I, Predicting tissue-specific enhancers in the human genome, *Genome Res* **17**(2):201–211, 2007.
12. Hallikas O, Palin K, Sinjushina N, Rautiainen R, Partanen J, Ukkonen E, Taipale J, Genome-wide prediction of mammalian enhancers based on analysis of transcription-factor binding affinity, *Cell* **124**(1):47–59, 2006.
13. Kreiman G, Identification of sparsely distributed clusters of cis-regulatory elements in sets of co-expressed genes, *Nucleic Acids Res* **32**(9):2889–2900, 2004.

14. Blackwood E, Kadonaga J, Going the distance: A current view of enhancer action, *Science* **281**:60–63, 1998.
15. Mayer H, Bilban M, Kurtev V, Gruber F, Wagner O, Binder BR, de Martin R, Deciphering regulatory patterns of inflammatory gene expression from interleukin-1-stimulated human endothelial cells, *Arterioscler Thromb Vasc Biol* **24**(7):1192–1198, 2004.
16. Heintzman ND, Stuart RK, Hon G, Fu Y, Ching CW, Hawkins RD, Barrera LO, Van Calcar S, Qu C, Ching KA, Wang W, Weng Z, Green RD, Crawford GE, Ren B, Distinct and predictive chromatin signatures of transcriptional promoters and enhancers in the human genome, *Nat Genet* **39**(3):311–318, 2007.
17. Chan BY, Kibler D, Using hexamers to predict cis-regulatory motifs in Drosophila, *BMC Bioinformatics* **6**:262, 2005.
18. Hutchinson GB, The prediction of vertebrate promoter regions using differential hexamer frequency analysis, *Comput Appl Biosci* **12**(5):391–398, 1996.
19. Stuart RO, Bush KT, Nigam SK, Changes in gene expression patterns in the ureteric bud and metanephric mesenchyme in models of kidney development, *Kidney International* **64**(6):1997–2008, 2003.
20. Challen G, Gardiner B, Caruana G, Kostoulias X, Martinez G, Crowe M, Taylor DF, Bertram J, Little M, Grimmond SM, Temporal and spatial transcriptional programs in murine kidney development, *Physiol Genomics* **23**(2):159–171, 2005.
21. Margolin AA, Nemenman I, Basso K, Wiggins C, Stolovitzky G, Dalla Favera R, Califano A, ARACNE: An algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context, *BMC Bioinformatics* **7**(Suppl 1):S7, 2006.
22. Rao A, Hero AO, States DJ, Engel JD, Using directed information to build biologically relevant influence networks, *Proc Computational Systems Bioinformatics (CSB)*, 2007.
23. Kleinjan DA, van Heyningen V, Long-range control of gene expression: Emerging mechanisms and disruption in disease, *Am J Hum Genet* **76**(1):8–32, 2005.
24. Carvajal JJ, Keith A, Rigby PW, Global transcriptional regulation of the locus encoding the skeletal muscle determination genes Mrf4 and Myf5, *Genes Dev* **22**(2):265–76, 2008.
25. Landry JR, Bonadies N, Kinston S, Knezevic K, Wilson NK, Oram SH, Janes M, Piltz S, Hammett M, Carter J, Hamilton T, Donaldson IJ, Lacaud G, Frampton J, Follows G, Kouskoff V, Gttgens B, Expression of the leukemia oncogene Lmo2 is controlled by an array of tissue-specific elements dispersed over 100 kb and bound by Tal1/Lmo2, Ets, and Gata factors, *Blood* **113**(23):5783–5792, 2009.
26. Breiman L, Random forests, *Machine Learning* **45**(1):5.32, 2001.
27. Eisenberg E, Levanon EY, Human housekeeping genes are compact, *Trends Genet* **19**(7):362–365, 2003.
28. Farr D, Bellora N, Mularoni L, Messeguer X, Alb MM, Housekeeping genes tend to show reduced upstream sequence conservation, *Genome Biol* **8**(7):R140, 2007.
29. Liaw A, Wiener M, Classification and Regression by randomForest, *R News* **2**:18–22, 2002.
30. McGaughey DM, Vinton RM, Huynh J, Al-Saif A, Beer MA, McCallion AS, Metrics of sequence constraint overlook regulatory sequences in an exhaustive analysis at Pox2b, *Genome Res* **18**(2):252–260, 2008.
31. Odom DT, Dowell RD, Jacobsen ES, Gordon W, Danford TW, MacIsaac KD, Rolfe PA, Conboy CM, Gifford DK, Fraenkel E, Tissue-specific transcriptional regulation has diverged significantly between human and mouse, *Nat Genet* **39**(6):730–732, 2007.

32. Gilbert SF, *Developmental Biology*, Sinauer Associates Inc., Publishers Sunderland, Massachusetts, 1997.
33. Lettice LA, Heaney SJ, Purdie LA, Li L, de Beer P, Oostra BA, Goode D, Elgar G, Hill RE, de Graaff E, A long-range Shh enhancer regulates expression in the developing limb and fin and is associated with preaxial polydactyly, *Hum Mol Genet* **12**(14):1725–1735, 2003.
34. Ohuchi H, Yasue A, Ono K, Sasaoka S, Tomonari S, Takagi A, Itakura M, Moriyama K, Noji S, Nohno T, Identification of cis-element regulating expression of the mouse Fgf10 gene during inner ear development, *Dev Dyn* **233**(1):177–187, 2005.
35. Bader GD, Hogue CW, An automated method for finding molecular complexes in large protein interaction networks, *BMC Bioinformatics* **4**:2, 2003.
36. Sharan R, Suthram S, Kelley RM, Kuhn T, McCuine S, Uetz P, Sittler T, Karp RM, Ideker T, Conserved patterns of protein interaction in multiple species, *Proc Natl Acad Sci USA* **102**(6):1974–1979, 2005.
37. Assenov Y, Ramrez F, Schelhorn SE, Lengauer T, Albrecht M, Computing topological parameters of biological networks, *Bioinformatics* **24**(2):282–284, 2008.
38. Xu L, Krzyzak A, Suen CY, Methods of combining multiple classifiers and their applications to handwriting recognition, *IEEE Transactions on Systems, Man and Cybernetics* **22**(3):418–435, 1992.

**Arvind Rao** received his Bachelor of Engineering degree (with distinction) in Electronics and Communications from Bangalore University, India in 2001. In 2003, he received the Master of Science in Engineering degree from the Electrical and Computer Engineering department at the University of Texas at Austin, with a specialization in Communications, Networks and Systems. He earned an A.M. in statistics from the University of Michigan in 2007 and was a Rackham Predoctoral Fellow. For his doctoral work at the University of Michigan, he worked on understanding long-range transcriptional regulation in higher eukaryotes. Currently, he is a Lane Fellow in Computational Biology in the School of Computer Science, Carnegie Mellon University. His research interests lie at the intersection of signal processing, machine learning, experimental and computational systems biology.

**David J. States** received his B.A. (magna cum laude) from Harvard College (1975) and his M.D. and Ph.D. from Harvard University (1983) in Biophysics. He was a resident in Internal Medicine at the University of California, San Diego and a Clinical Associate at the National Heart, Lung and Blood Institute. In 1988, he joined the National Center for Biotechnology Information, and in 1992, moved to Washington University in St. Louis to be the Director of the Institute for Biomedical Computing and in 2001, he became a Professor of Human Genetics and founding Director of the Bioinformatics Program at the University of Michigan in Ann Arbor. In 2008, Dr. States moved to the School of Health Science Information at the University of Texas in Houston to establish a Center for Systems Biology and Bioinformatics. His recent research interests have been in alternative splicing in cancer, the analysis of genome regulation and data integration in molecular biology. He is a Fellow of the American College of Medical Informatics (ACMI). He was a founding member of the Board of Directors of the International Society for Computational Biology and

served as Treasurer from 1997 to 2000. He chaired organizing committees for the Intelligent Systems in Molecular Biology Conference in 1996 and 2005.

**Alfred O. Hero III** received his B.Sc. (summa cum laude) from Boston University (1980) and his Ph.D. from Princeton University (1984), both in Electrical Engineering. Since 1984, he has been with the University of Michigan, Ann Arbor, where he is the R. Jamison and Betty Professor of Engineering. His primary appointment is in the Department of Electrical Engineering and Computer Science and he also has appointments, by courtesy, in the Department of Biomedical Engineering and the Department of Statistics. In 2008, he was awarded the Digiteo Chaire d'Excellence, sponsored by Digiteo Research Park in Paris, located at the Ecole Supérieure d'Electricité, Gif-sur-Yvette, France. He is a Fellow of the Institute of Electrical and Electronics Engineers (IEEE) and several of his research articles have received best paper awards. He received the IEEE Signal Processing Society Meritorious Service Award (1998), and the IEEE Third Millennium Medal (2000). Alfred Hero was the President of the IEEE Signal Processing Society (2006–2008) and sits on the Board of Directors of the IEEE (2009–2011).

Alfred Hero's recent research interests have been in detection, classification, pattern analysis, and adaptive sampling for spatio-temporal data. Of particular interest are applications to network security, multi-modal sensing and tracking, biomedical imaging, and genomic signal processing.

**James Douglas Engel** received his Ph.D. in biophysical chemistry at the University of Oregon in 1975 with P.H. von Hippel. He was a Helen Hay Whitney Fellow with N. Davidson and T. Maniatis at Caltech from 1975 to 1978, and then joined the faculty at Northwestern University, where he became the Owen L. Coon Professor and Assoc. Director for Basic Sciences of the Robert H. Lurie Comprehensive Cancer Center. He moved to the University of Michigan School of Medicine in 2002, and was endowed the first G. Carl Huber Chair in Developmental Biology and Chair of the Department of Cell and Developmental Biology. Dr. Engel's lab is devoted to deciphering, both experimentally and theoretically, the transcriptional regulatory networks that lead to correctly modulated gene expression during embryonic development. Dr. Engel is a Fellow of the American Association of Arts and Sciences and an Editor of *Molecular and Cellular Biology*.