

Motif Discovery in Tissue-Specific Regulatory Sequences Using Directed Information

Arvind Rao, Alfred O. Hero III, David J. States, James Douglas Engel

Abstract—Motif discovery for the identification of functional regulatory elements underlying gene expression is a challenging problem. Sequence inspection often leads to discovery of novel motifs (including transcription factor sites) with previously uncharacterized function in gene expression. Coupled with the complexity underlying tissue-specific gene expression, there are several motifs that are putatively responsible for expression in a certain cell type. This has important implications in understanding fundamental biological processes such as development and disease progression.

In this work we present an approach to the identification of motifs (not necessarily transcription factor sites) and examine its application to some questions in current bioinformatics research. These motifs are seen to discriminate tissue-specific gene-promoter or regulatory regions from those that are not tissue-specific. There are two main contributions of this work: Firstly, we propose the use of directed information for such classification constrained motif discovery, and then, use the selected features with a support vector machine (SVM) classifier to find the tissue-specificity of any sequence of interest. Such analysis yields several novel interesting motifs that merit further experimental characterization. Furthermore, this approach leads to a principled framework for the prospective examination of any chosen motif to be discriminatory motif for a group of co-expressed/co-regulated genes, thereby integrating sequence and expression perspectives. We hypothesize that the discovery of these motifs would enable the large-scale investigation for the tissue-specific regulatory role of any conserved sequence element identified from genome-wide studies.

Index Terms—Directed Information, transcriptional regulation, phylogeny, Transcription factor binding sites (TFBS), tissue-specific genes.

I. 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

The material in this paper was presented in part at the IEEE Statistical Signal Processing Workshop 2007, (SSP07).

Arvind Rao and Alfred O. Hero, III are with the Departments of Electrical Engineering and Computer Science, and Bioinformatics at the University of Michigan, Ann Arbor, MI-48109, email: [ukarvind, hero] @umich.edu

David J. States is with the Departments of Bioinformatics and Human Genetics at the University of Michigan, Ann Arbor, MI-48109.

James Douglas Engel is with the Department of Cell and Developmental Biology at the University of Michigan, Ann Arbor, MI-48109.

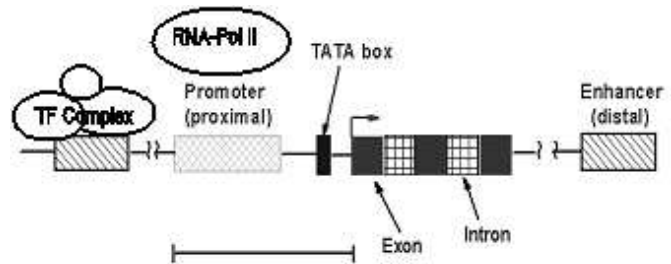


Fig. 1. Schematic of Transcriptional Regulation. Sequence motifs at the promoter and the distal regulatory elements together confer specificity of gene expression via TF binding.

representing the gene. It is the intermediate step before the generation of functional protein from messenger RNA. During gene expression (Fig. 1), 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 (TSS). The basal transcriptional machinery at the promoter coupled with the transcription factor complexes at these distal, long-range regulatory elements (LREs) are collectively involved in directing tissue-specific expression of genes.

One of the current challenges in the post-genomic era is the principled discovery of such LREs genome-wide. Recently, there has been a community-wide effort (<http://www.genome.gov/ENCODE/>) to find all regulatory elements in 1% of the human genome. The examination of the discovered elements would reveal characteristics typical of most enhancers which would aid their principled discovery and examination on a genome-wide scale. Some characteristics of experimentally identified distal regulatory elements ([21],[19]) are:

- Non-coding elements: Distal regulatory elements are non-coding and can either be intronic or intergenic regions on the genome. Hence previous models for gene finding [4] are not directly applicable. With over 98% of the annotated genome being non-coding, the precise localization of regulatory elements that underlie tissue-specific gene expression is a challenging problem.
- Distance/orientation independent: an enhancer can act from variable genomic distances (hundreds of kilobases) to regulate gene expression in conjunction with the proximal promoter, possibly via a looping mechanism [34]. These enhancers can lie upstream or downstream of the actual gene along the genomic locus.

- Promoter dependent: Since the action at a distance of these elements involves the recruitment of TFs that direct tissue-specific gene expression, the promoter that they interact with is critical.

Although there are instances where a gene harbors tissue-specific activity at the promoter itself, the role of long range elements (LREs) remains of interest, e.g: for a detailed understanding of their regulatory role in gene expression during biological processes like organ development and disease progression [18]. We seek to develop computational strategies to find novel LREs genome-wide that govern tissue specific expression for any gene of interest. A common approach for their discovery is the use of motif-based sequence signatures. Any sequence element can then be scanned for such a signature and its tissue-specificity can be ascertained [30].

Thus, our primary question in this regard is: is there a discriminating sequence property of LRE elements that determine tissue-specific gene expression - more particularly, are there any sequence motifs in known regulatory elements that can aid discovery of new elements [17]. To answer this, we examine known tissue-specific regulatory elements (promoters and enhancers) for motifs that discriminate them from a background set of neutral elements (such as housekeeping gene promoters). For this study, the datasets are derived from the following sources:

- Promoters of tissue-specific genes: Before the widespread discovery of long-range regulatory elements (LREs), it was hypothesized that promoters governed gene expression alone. There is substantial evidence for the binding of tissue-specific transcription factors at the promoters of expressed genes. This suggests that, in spite of newer information implicating the role of LREs, promoters also have interesting motifs that govern tissue-specific expression.

Another practical reason for the examination of promoters is that their locations (and genomic sequences) are more clearly delineated on genome databases (like UCSC or Ensembl). Sufficient data (<http://symatlas.gnf.org/>) on the expression of genes is also publicly available for analysis. Sequence motif discovery is set up as a feature extraction problem from these tissue-specific promoter sequences. Subsequently, a support vector machine (SVM) classifier is used to classify new promoters into specific and non-specific categories based on the identified sequence features (motifs). Using the SVM classifier algorithm, 90% of tissue-specific genes are correctly classified based upon their upstream promoter region sequences alone.

- Known long range regulatory elements (LRE) motifs: To analyze the motifs in LRE elements, we examine the results of the above approach on the Enhancer Browser dataset (<http://enhancer.lbl.gov/>) which has results of expression of ultraconserved genomic elements in transgenic mice [32]. An examination of these ultraconserved enhancers is useful for the extraction of discriminatory motifs to distinguish the regulatory elements from the non-regulatory (neutral) ones. Here the results indicate that up to 95% of the sequences can

be correctly classified using these identified motifs.

We note that some of the identified motifs might not be transcription factor binding motifs, and would need to be functionally characterized. This is an advantage of our method - instead of constraining ourselves to the degeneracy present in TF databases (like TRANSFAC/JASPAR), we look for all sequences of a fixed length.

II. CONTRIBUTIONS

Using microarray gene expression data, ([14],[37]) proposes an approach to assign genes into tissue-specific and non-specific categories using an entropy criterion. Variation in expression and its divergence from ubiquitous expression (uniform distribution across all tissue types) is used to make this assignment. Based on such assignment, several features like CpG island density, frequency of transcription factor motif occurrence, can be examined to potentially discriminate these two groups. Other work has explored the existence of key motifs (transcription factor binding sites) in the promoters of tissue-specific genes ([41],[1]). Based on the successes reported in these methods, it is expected that a principled examination and characterization of every sequence motif identified to be discriminatory might lead to improved insight into the biology of gene regulation. For example, such a strategy might lead to the discovery of newer TFBS motifs, as well as those underlying epigenetic phenomena.

For the purpose of identifying discriminative motifs from the training data (tissue-specific promoters or LREs), our approach is as follows:

- *Variable selection:* Firstly, sequence motifs that discriminate between tissue-specific and non-specific elements are discovered. In machine learning, this is a feature selection problem with features being the counts of sequence motifs in the training sequences. Without loss of generality, six-nucleotide motifs (hexamers) are used as motif features. 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 setup has been introduced in ([5], [12],[38]). The motif search space is, therefore a $4^6 = 4096$ dimensional one. The presented approach, however, does not depend on motif length and can be scaled according to biological knowledge.

For variable (motif) selection, a novel feature selection approach (based on an information theoretic quantity called **directed information** - DI) is proposed. The improved performance of this criterion over using mutual information for motif selection is also demonstrated.

- *Classifier design:* After discovering discriminating motifs using the above DI step, a SVM classifier that separates the samples between the two classes (specific and non-specific) from this motif space, is constructed.

Apart from this novel feature selection approach, several questions pertaining to bioinformatics methodology can be potentially answered using this framework - some of these are:

- Are there common motifs underlying tissue-specific expression that are identified from tissue-specific promoters and enhancers? - in this paper, an examination of motifs (from promoters and enhancers) corresponding to brain-specific expression is done to address this question.
- Do these motifs correspond to known motifs (transcription factor binding sites)? - we show that several motifs are indeed consensus sites for transcription factor binding, although their real role can only be identified in conjunction with experimental evidence.
- Is it possible to relate the motif information from the sequence and expression perspectives to understand regulatory mechanisms? - this question is addressed in Section XI.C.
- How useful are these motifs in predicting new tissue-specific regulatory elements? - this is partly explained from the results of SVM classification.

This work differs from that in ([5], [12]), in several aspects. We present the DI based feature selection procedure as part of an overall unified framework to answer several questions in bioinformatics, not limited to finding discriminating motifs between two classes of sequences. Particularly, one of the advantages is the ability to examine any particular motif as a potential discriminator between two classes. Also, this work accounts for the notion of tissue-specificity of promoters/enhancers (in line with more recent work in [20],[32],[14],[37],[16]). Also, this framework enables the principled integration of various data sources to address the above questions. These are clarified in the Results (Section: XI).

III. RATIONALE

The main approaches to finding common motifs driving tissue-specific gene regulation are summarized in ([19], [21]). The most common approach is to look for TFBS motifs that are statistically over-represented in the promoters of the co-expressed genes based on a background (binomial or Poisson) distribution of motif occurrence genomewide.

In this work, the problem of motif discovery is set up as follows. Using two annotated groups of genes, tissue-specific ('*ts*') and non-tissue specific ('*nts*'), hexamer motifs that best discriminate these two classes are found. The goal would be to make this set of motifs as small as possible - i.e. to achieve maximal class partitioning with the smallest feature subset.

Several metrics have been proposed to find features with maximal class label association. From information theory, mutual information is a popular choice [31]. This is a symmetric association metric and does not resolve the direction of dependency (i.e if features depend on the class label or vice versa). It is important to find features that induce the class label. Feature selection from data implies selection (control) of a feature subset that maximally captures the underlying character (class label) of the data. There is no control over the label (a purely observational characterization).

With this motivation, a new metric for discriminative hexamer subset selection, termed "directed information" (DI), is proposed. Based on the selected features, a classifier is used to classify sequences to tissue-specific or non-tissue-specific

categories. The performance of this DI based feature selection metric is subsequently evaluated in the context of the SVM classifier.

IV. OVERALL METHODOLOGY

The overall schematic of the proposed procedure is outlined below (Fig. 2).

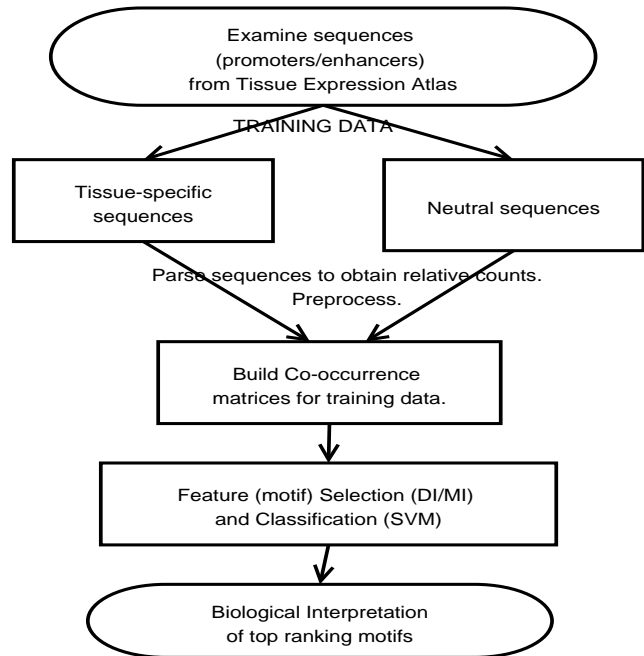


Fig. 2. An overview of the proposed approach. Each of the steps are outlined in the following sections.

Below we present our approach to find promoter-specific or enhancer-specific motifs.

V. MOTIF ACQUISITION

A. Promoter motifs:

1) *Microarray Analysis:* Raw microarray data is available from the Novartis Foundation (GNF) [<http://symatlas.gnf.org/>]. Data is normalized using RMA from the Bioconductor packages for R [cran.r-project.org/]. Following normalization, replicate samples are averaged together. Only 25 tissue types are used in our analysis including: Adrenal Gland, Amygdala, Brain, Caudate Nucleus, Cerebellum, Corpus Callosum, Cortex, Dorsal Root Ganglion, Heart, HUVEC, Kidney, Liver, Lung, Pancreas, Pituitary, Placenta, Salivary, Spinal Cord, Spleen, Testis, Thalamus, Thymus, Thyroid, Trachea, and Uterus.

In this context, the notion of tissue-specificity of a gene needs clarification. Suppose there are N genes, g_1, g_2, \dots, g_N and T tissue types (in GNF: $T = 25$), we construct a $N \times T$ tissue specificity matrix: $M = [0]_{N \times T}$. For each gene g_i , $1 \leq i \leq N$, let $g_{i,[0.5T]} = \text{median}(g_{i,k}), \forall k \in 1, 2, \dots, T$; $g_{i,k}$ being the expression level of gene ' i ' in tissue ' k '. Define, each entry $M_{i,k}$ as,

$$M_{i,k} = \begin{cases} 1 & \text{if } g_{i,k} \geq 2g_{i,[0.5T]}; \\ 0 & \text{otherwise.} \end{cases}$$

Now consider the N dimensional vector $m_i = \sum_{k=1}^T M_{i,k}$, $1 \leq i \leq N$ i.e. summing all the columns of each row. The inter-quartile range of ' m ' can be used for ' ts '/' nts ' assignment. Gene indices ' i ' that are in quartile 1 ($=3$), are labeled as ' ts ', and those in quartile 4 ($=22$), are labeled as ' nts '.

With this approach, a total of 1924 probes representing 1817 genes were classified as tissue-specific, while 2006 probes representing 2273 genes were classified as non tissue-specific. In this work, genes which are either heart-specific or brain-specific are considered. From the tissue-specific genes obtained from the above approach, 45 brain-specific gene promoters and 118 heart-specific gene promoters are obtained. As mentioned in Section II, one of the objectives is to find motifs that are responsible for brain/heart specific expression and also correlate them with binding profiles of known transcription factor binding motifs.

2) *Sequence Analysis*: Genes (' ts ' or ' nts ') associated with candidate probes are identified using the Ensembl Ensmart [<http://www.ensembl.org/>] tool. For each gene, sequence from 2000bp upstream and 1000bp down-stream upto the start of the first exon relative to their reported TSS is extracted from the Ensembl Genome Database (Release 37). The relative counts of each of the 4^6 hexamers are computed within each gene-promoter sequence of the two categories (' ts ' and ' nts ') - using the '*seqinr*' library in the R environment. A t-test is performed between the relative counts of each hexamer between the two expression categories (' ts ' and ' nts ') and the top 1000 significant hexamers ($\vec{H} = H_1, H_2, \dots, H_{1000}$) is obtained. The relative counts of these hexamers is recomputed for each gene individually. This results in two hexamer-gene co-occurrence matrices, - one for the ' ts ' class (dimension $N_{train,+1} \times 1000$) and the other for the ' nts ' class (dimension $N_{train,-1} \times 1000$). Here $N_{train,+1}$ and $N_{train,-1}$ are the number of positive training and negative training samples, respectively.

The input to the feature selection procedure is a gene promoter - motif frequency table (Table I). The genes relevant to each class are identified from tissue microarray analysis, following steps 1 and 2 above, and the frequency table is built by parsing the gene promoters for the presence of each of the $4^6 = 4096$ possible hexamers.

B. LRE motifs:

To analyze long range elements which confer tissue specific expression, the Mouse Enhancer database (<http://enhancer.lbl.gov/>) is examined. This database has a list of experimentally validated ultraconserved elements which have been tested for tissue specific expression in transgenic mice [32], and can be searched for a list of all elements which have expression in a tissue of interest. In this work, we consider expression in tissues relating to the developing brain.

Ensembl Gene ID	AAAAAA	AAAAAG	AAAAAT	AAAAACA
ENSG00000155366	0	0	1	4
ENSG000001780892	6	5	5	6
ENSG00000189171	1	2	1	0
ENSG00000168664	6	3	8	0
ENSG00000160917	4	1	4	2
ENSG00000163655	2	4	0	1
ENSG000001228844	8	6	10	7
ENSG00000176749	0	0	0	0
ENSG00000006451	5	2	2	1

TABLE I

THE 'MOTIF FREQUENCY MATRIX' FOR A SET OF GENE-PROMOTERS. THE FIRST COLUMN IS THEIR ENSEMBL GENE IDENTIFIERS AND THE OTHER 4 COLUMNS ARE THE MOTIFS. A CELL ENTRY DENOTES THE NUMBER OF TIMES A GIVEN MOTIF OCCURS IN THE UPSTREAM (-2000 TO +1000BP FROM TSS) REGION OF EACH CORRESPONDING GENE.

According to the experimental protocol, the various regions are cloned upstream of a heat shock protein promoter (*hsp68-lacz*), thereby not adhering to the idea of promoter specificity in tissue-specific expression. Though this is of concern in that there is loss of some gene-specific information, we work with this data since we are more interested in tissue expression and also due to a paucity of public promoter-dependent enhancer data .

This database also has a collection of ultraconserved elements that do not have any transgenic expression in-vivo. This is used as the neutral/background set of data which corresponds to the ' nts ' (non-tissue specific class) for feature selection and classifier design.

As in the above (promoter) case, these sequences (seventy four enhancers for brain-specific expression) are parsed for the absolute counts of the 4096 hexamers, a co-occurrence matrix ($N_{train,+1} = 74$) is built and then t-test p -values are used to find the top 1000 hexamers ($\vec{H} = H'_1, H'_2, \dots, H'_{1000}$) that are maximally different between the two classes (brain-specific and brain-non-specific).

The next three sections clarify the preprocessing, feature selection and classifier design steps to mine these co-occurrence matrices for hexamer motifs that are strongly associated with the class label. We note that though this work is illustrated using two class labels, the approach can be extended in a straightforward way to the multi-class problem.

VI. PREPROCESSING

From the above, $N_{train,+1} \times 1000$ and $N_{train,-1} \times 1000$ dimensional co-occurrence matrices are available for the tissue-specific and non-specific data, both for the promoter and enhancer sequences. Before proceeding to the feature (hexamer motif) selection step, the counts of the $M = 1000$ hexamers in each training sample need to be 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} gene. Then, for each gene 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 has the corresponding class label ($Y = -1/+1$).

VII. DIRECTED INFORMATION AND FEATURE SELECTION

The primary goal in feature selection is to find the minimal subset of features (from hexamers: \vec{H}/\vec{H}^T) that lead to maximal discrimination of the class label ($Y_i \in (-1/+1)$), using each of the $i \in (1, 2, \dots, (N_{train,+1} + N_{train,-1}))$ genes during training. We are looking for a subset of the variables ($H_{i,1}, \dots, H_{i,1000}$) which are directionally associated with the class label (Y_i). These hexamers putatively influence/induce the class label (Fig. 3). As can be seen from [33], there is considerable interest in discovering such dependencies from expression and sequence data. Following [10], we search for features (in *measurement* space) that induce the class label (in *observation* space).

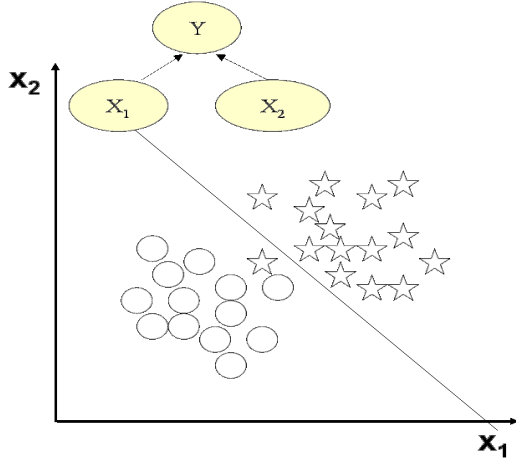


Fig. 3. Causal Feature discovery for two class discrimination, adapted from [10]. Here the variables X_1 and X_2 discriminate Y , the class label.

One way to interpret the feature selection problem is the following: Nature is trying to communicate a source symbol ($Y \in \{-1/+1\}$), corresponding to the gene class label ('nts/ts'), to us. In this setup, an encoder that extracts frequencies of a particular hexamer (H_i) maps the source symbol (Y) to $H_i(Y)$. The decoder outputs the source reconstruction \hat{Y} based on the received codeword $c_i(Y) = H_i(Y)$.

We observe that there are several possible encoding schemes $c_i(Y)$ that the encoder could potentially use ($i = 1, 2, \dots, 1000$), each corresponding to feature extraction via a different hexamer H_i . An encoder is the mapping rule $c_i : Y \rightarrow H_i$. The ideal encoding scheme is one which induces the most discriminative partitioning of the code (feature) space, for successful reconstruction of Y by the decoder. The ranking of each encoder's performance over all possible mappings yields the most discriminative mapping. This measure of performance is the amount of information flow from the mapping (hexamer) to the class label. Using mutual information as one such measure indeed identifies the best features [31], but fails to resolve the direction of dependence due to its symmetric nature $I(H_i; Y) = I(Y; H_i)$. The direction of dependence is important since it pinpoints those features that induce the class label (not vice-versa). This is necessary since these class labels

are predetermined (given to us by biology) and the only control we have is the feature space onto which we project the data points, for the purpose of classification. This loosely parallels the use of the directed edges in bayesian networks for inference of feature-class label associations [10].

Unlike mutual information (MI), directed information (DI) is a metric to quantify the directed flow of information. It was originally introduced in ([22], [23]) to examine the transfer of information from encoder to decoder under feedback/feedforward scenarios and to resolve directivity during bidirectional information transfer. Given its utility in the encoding of sources with memory (correlated sources), this work demonstrates it to be a competitive metric to MI for feature selection in learning problems. DI answers which of the encoding schemes (corresponding to each hexamer H_i) leads to maximal information transfer from the hexamer labels to the class labels (i.e. directed dependency).

The DI is a measure of the directed dependence between two vectors $X_i = [X_{1,i}, X_{2,i}, \dots, X_{n,i}]$ and $Y = [Y_1, Y_2, \dots, Y_n]$. In this case, $X_{j,i}$ = quantile label for the frequency of hexamer $i \in (1, 2, \dots, 1000)$ in the j^{th} training sequence. $Y = [Y_1, Y_2, \dots, Y_n]$ are the corresponding class labels ($-1, +1$). For a block length N , the DI is given by [23]:

$$I(X_i^N \rightarrow Y^N) = \sum_{n=1}^N I(X_i^n; Y_n | Y^{n-1}) \quad (1)$$

Using a stationarity assumption over a finite-length memory of the training samples, a correspondence with the setup in ([23], [40]) can be seen. As already known [6], the mutual information $I(X^N; Y^N) = H(X^N) - H(X^N | Y^N)$, where $H(X^N)$ and $H(X^N | Y^N)$ are the Shannon entropy of X^N and the conditional entropy of X^N given Y^N , respectively. With this definition of mutual information, the Directed Information simplifies to,

$$\begin{aligned} I(X^N \rightarrow Y^N) &= \sum_{n=1}^N [H(X^n | Y^{n-1}) - H(X^n | Y^n)] \\ &= \sum_{n=1}^N \{ [H(X^n, Y^{n-1}) - H(Y^{n-1})] - [H(X^n, Y^n) - H(Y^n)] \} \end{aligned} \quad (2)$$

Using (2), the Directed information is expressed in terms of individual and joint entropies of X^n and Y^n . This expression implies the need for higher-order entropy estimation from a moderate sample size. A Voronoi tessellation [24] based adaptive partitioning of the observation space can handle $N = 5/6$ without much complexity.

The relationship between MI and DI is given by [23],

$$\begin{aligned} \text{DI: } I(X^N \rightarrow Y^N) &= \sum_{i=1}^N I(X^i; Y_i | Y^{i-1}). \\ \text{MI: } I(X^N; Y^N) &= \sum_{i=1}^N I(X^N; Y_i | Y^{i-1}) \\ &= I(X^N \rightarrow Y^N) + I(0Y^{N-1} \rightarrow X^N). \end{aligned}$$

To clarify, $I(X^N \rightarrow Y^N)$ is the directed information from X to Y , whereas $I(0Y^{N-1} \rightarrow X^N)$ is the directed information from a (one-sample) delayed version of Y^N to X^N . From [40], it is clear that DI resolves the direction of information transfer (feedback or feedforward). If there is no feedback/feedforward, $I(X^N \rightarrow Y^N) = I(X^N; Y^N)$.

From the above chain-rule formulations for DI and MI, it is clear that the expression for DI is permutation-variant (i.e., the value of the DI is different for a different ordering of random variables). Thus, we instead find the $I_p(X^N \rightarrow Y^N)$, a DI measure for a particular ordering of the N random variables (r.v.s). The DI value for our purpose, $I(X^N \rightarrow Y^N)$ is an average over all possible sample permutations given by, $I(X^N \rightarrow Y^N) = \frac{1}{N!} \sum_{p=1}^{N!} I_p(X^N \rightarrow Y^N)$. For MI, however, $I_p(X^N; Y^N) = I(X^N; Y^N)$ because, MI is permutation-invariant (i.e., independent of r.v ordering). As can be readily observed, this problem is combinatorially complex, and hence, a monte-carlo sampling strategy (1000 trials) is used for computing $I(X^N \rightarrow Y^N)$. This is because we find that about 1000 trials yields a DI confidence interval (CI) that is only 20% more than the corresponding CI obtained from 10,000 trials of the data, a far more exhaustive number.

To select features, we maximize $I(X^N \rightarrow Y^N)$ over the possible pairs (\vec{X}, Y) . This feature selection problem for the i^{th} training instance reduces to identifying which hexamer ($k \in (1, 2, \dots, 4096)$) has the highest $I(X_k \rightarrow Y)$.

The higher dimensional entropy can be estimated using order statistics of the observed samples [24] by iterative partitioning of the observation space until nearly uniform partitions are obtained. This method lends itself to a partitioning scheme that can be used for entropy estimation even for a moderate number of samples in the observation space of the underlying probability distribution. Several such algorithms for adaptive density estimation have been proposed ([42],[27],[29]) and can find potential application in this procedure. In this methodology, a Voronoi tessellation approach for entropy estimation because of the higher performance guarantees as well as the relative ease of implementation of such a procedure.

The above method is used to estimate the true DI between a given hexamer and the class label for the entire training set. Feature selection comprises of finding all those hexamers (X_i) for which $I(X_i^N \rightarrow Y^N)$ is the highest. From the definition of DI, we know that $0 \leq I(X_i^N \rightarrow Y^N) \leq I(X_i^N; Y^N) < \infty$. To make a meaningful comparison of the strengths of association between different hexamers and the class label, we use a normalized score to rank the DI values. This normalized measure ρ_{DI} should be able to map this large range ($[0, \infty]$) to $[0, 1]$. Following [13], an expression for the normalized DI is given by:

$$\rho_{DI} = \frac{I(X_i^N \rightarrow Y^N)}{\sqrt{1 - e^{-2I(X_i^N \rightarrow Y^N)}}} = \sqrt{1 - e^{-2 \sum_{i=1}^N I(X_i^i; Y_i | Y^{i-1})}}$$

Another point of consideration is to estimate the significance of the DI value compared to a null distribution on the DI value (i.e. what is the chance of finding the DI value by chance from the N -length series X_i and Y). This is done using confidence intervals after permutation testing (Sec: VIII).

VIII. BOOTSTRAPPED CONFIDENCE INTERVALS

In the absence of knowledge of the true distribution of the DI estimate, an approximate confidence interval for the DI estimate ($\hat{I}(X^N \rightarrow Y^N)$), is found using bootstrapping [8]. Density estimation is based on kernel smoothing over the bootstrapped samples [35].

The kernel density estimate for the bootstrapped DI (with $n = 1000$ samples), $Z \triangleq \hat{I}_B(X^N \rightarrow Y^N)$ becomes,

$\hat{f}_h(Z) = \frac{1}{nh} \sum_{i=1}^n \frac{3}{4} [1 - (\frac{z_i - z}{h})^2] I(|\frac{z_i - z}{h}| \leq 1)$ with $h \approx 2.67\hat{\sigma}_z$ and $n = 1000$. $\hat{I}_B(X^N \rightarrow Y^N)$ is obtained by finding the DI for each random permutation of the X, Y series, and performing this permutation B times. As is clear from the above expression, the Epanechnikov kernel is used for density estimation from the bootstrapped samples. The choice of the kernel is based on its excellent characteristics - a compact region of support, the lowest AMISE (asymptotic mean squared error) and favorable bias-variance tradeoff [35].

We denote the cumulative distribution function (over the bootstrap samples) of $\hat{I}(X^N \rightarrow Y^N)$ by $F_{\hat{I}_B(X^N \rightarrow Y^N)}(\hat{I}_B(X^N \rightarrow Y^N))$. Let the mean of the bootstrapped null distribution be $I_B^*(X^N \rightarrow Y^N)$. We denote by $t_{1-\alpha}$, the $(1-\alpha)^{th}$ quantile of this distribution i.e. $\{t_{1-\alpha} : P([\frac{\hat{I}_B(X^N \rightarrow Y^N) - I_B^*(X^N \rightarrow Y^N)}{\hat{\sigma}}] \leq t_{1-\alpha}) = 1 - \alpha\}$. Since we need the true $\hat{I}(X^N \rightarrow Y^N)$ to be significant and close to 1, we need $\hat{I}(X^N \rightarrow Y^N) \geq [I_B^*(X^N \rightarrow Y^N) + t_{1-\alpha} \times \hat{\sigma}]$, with $\hat{\sigma}$ being the standard error of the bootstrapped distribution, $\hat{\sigma} = \sqrt{\frac{[\sum_{b=1}^B \hat{I}_b(X^N \rightarrow Y^N) - I_B^*(X^N \rightarrow Y^N)]^2}{B-1}}$; B is the number of bootstrap samples.

This hypothesis test is done for each of the 1000 motifs, in order to select the top ' d ' motifs based on DI value, which is then used for classifier training subsequently. This leads to a need for multiple-testing correction. Because the Bonferroni correction is extremely stringent in such settings, the Benjamini-Hochberg procedure [3], which has a higher false positive rate, but a lower false negative rate is used in this work.

IX. SUPPORT VECTOR MACHINES

From the top ' d ' features identified from the ranked list of features having high DI with the class label, a support vector machine classifier in these ' d ' dimensions is designed. A SVM is a hyperplane classifier which operates by finding a maximum margin linear hyperplane to separate two different classes of data in high dimensional ($D > d$) space. The training data has $N (= N_{train,+1} + N_{train,-1})$ pairs $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$, with $x_i \in \mathcal{R}^d$ and $y_i \in \{-1, +1\}$.

An SVM is a maximum margin hyperplane classifier in a non-linearly extended high dimensional space. For extending the dimensions from d to $D > d$, a radial basis kernel is used.

The objective is to minimize $\|\beta\|$ in the hyperplane $\{x : f(x) = x^T \beta + \beta_0\}$, subject to $y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i \forall i, \xi_i \geq 0, \sum \xi_i \leq \text{constant}$ [11].

X. SUMMARY OF OVERALL APPROACH

Our proposed approach is as follows. Here, the term 'sequence' can pertain to either tissue-specific promoters or LRE sequences, obtained from the GNF SymAtlas and Ensembl databases or the Enhancer Browser.

- 1) The sequence is parsed to obtain the relative counts/frequencies of occurrence of the hexamer in that sequence and to build the hexamer-sequence frequency matrix. The '*seqinr*' package in R is used for this purpose. This is done for all the sequences in the specific

(class “+1”) and non-specific (class “-1”) categories. The matrix thus has $N = N_{train,+1} + N_{train,-1}$ rows and $4^6 = 4096$ columns.

- 2) The obtained hexamer-sequence frequency matrix is pre-processed by assigning quantile labels for each hexamer within the i^{th} sequence. A hexamer-sequence matrix is thus obtained where the $(i, j)^{th}$ entry has the quantile label of the j^{th} hexamer in the i^{th} sequence. This is done for all the N training sequences consisting of examples from the -1 and +1 class labels.
- 3) Thus, two submatrices corresponding to the two class labels are built. One matrix contains the hexamer-sequence quantile labels for the positive training examples and the other matrix is for the negative training examples.
- 4) To select hexamers that are most different between the positive and negative training examples, a t-test is performed for each hexamer, between the ‘ts’ and ‘nts’ groups. Ranking the corresponding t-test p-values yields those hexamers that are most different distributionally between the positive and negative training samples. The top 1000 of these hexamers are chosen for further analysis. This step is only necessary to reduce the computational complexity of the overall procedure - computing the DI between each of the 4096 hexamers and the class label is relatively expensive.
- 5) For the top $K = 1000$ hexamers which are most significantly different between the positive and negative training examples, $I(X_k^N \rightarrow Y^N)$ and $I(X_k^N; Y^N)$ reveals the degree of association for each of the $k \in (1, 2, \dots, K)$ hexamers. The entropy terms in the directed information and mutual information expressions are found using a higher-order entropy estimator. Using the procedure of Sec: VII, the raw DI values are converted into their normalized versions. Since the goal is to maximize $I(X_k \rightarrow Y)$, we can rank the DI values in descending order.
- 6) The significance of the DI estimate is obtained based on the bootstrapping methodology. For every hexamer, a $p = 0.05$ significance with respect to its bootstrapped null distribution yields potentially discriminative hexamers between the two classes. The Benjamini-Hochberg procedure is used for multiple-testing correction. Ranking the significant hexamers by decreasing DI value yields features that can be used for classifier (SVM) training.
- 7) Train the Support Vector Machine classifier (SVM) on the top ‘d’ features from the ranked DI list(s). For comparison with the MI based technique, we use the hexamers which have the top ‘d’ (normalized) MI values. The accuracy of the trained classifier is plotted as a function of the number of features (d), after ten-fold cross-validation. As we gradually consider higher ‘d’, we move down the ranked list. In the plots below, the misclassification fraction is reported instead. A fraction of 0.1 corresponds to 10% misclassification.

Note: An important point concerns the training of the SVM classifier with the top ‘d’ features selected using DI or MI

(step 7 above). Since the feature selection step is decoupled from the classification step, it is preferred that the top ‘d’ motifs are consistently ranked high among multiple draws of the data, so as to warrant their inclusion in the classifier. However, this does not yield expected results on this data set. Briefly, a kendall rank correlation coefficient [15] was computed between the rankings of the motifs between multiple data draws (by sampling a subset of the entire dataset), for both MI and DI based feature-selection. It is observed that this coefficient is very low in both MI and DI, indicating a highly variable ranking. This is likely due to the high variability in data distribution across these multiple draws (due to limited number of data points), as well as the sensitivity of the data-dependent entropy estimation procedure to the range of the samples in the draw. To circumvent this problem of inconsistency in rank of motifs, a *median* DI/MI value is computed across these various draws and the top ‘d’ features based on the median DI/MI value across these draws are picked for SVM training [10].

XI. RESULTS

A. Tissue specific promoters

We use DI to find hexamers that discriminate brain-specific and heart-specific expression from neutral sequences. The negative training sets are sequences that are not brain or heart-specific, respectively. Results using the MI and DI methods are given below (Figs. 5 and 7). The plots indicate the SVM cross-validated misclassification accuracy (ideally 0) for the data as the number of features using the metric (DI or MI) is gradually increased. We can see that for any given classification accuracy, the number of features using DI is less than the corresponding number of features using MI. This translates into a lower misclassification rate for DI-based feature selection. We also observe that as the number of features ‘d’ is increased the performance of MI is the same as DI. This is expected since, as we gather more features using MI or DI, the differences in MI vs. DI ranking are compensated.

An important point needs to be clarified here. There is a possibility of sequence composition bias in the tissue-specific and neutral sequences used during training. This has been reported in recent work [38]. To avoid detecting GC rich sequences as hexamer features, it is necessary to confirm that there is no significant GC-composition bias between the specific and neutral sets in each of the case studies. This is demonstrated in Figs. 4, 6 and 8. In each case, it is observed that the mean GC-composition is almost same for the specific vs. neutral set. However, in such studies, it is necessary to select for sequences that do not exhibit such bias. In Figs. 6 and 8, even the distribution of GC-composition is similar among the samples. For Fig. 4, even though the distributions are slightly different, the box plots indicate similarity in mean GC-content.

Next, some of the motifs that discriminate between tissue-specific and non-specific categories for the brain promoter, heart promoter and brain enhancer cases respectively are listed in Table II. Additionally, if the genes encoding for

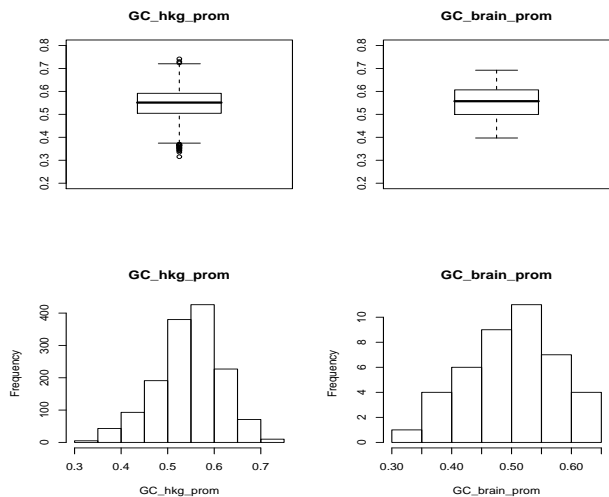


Fig. 4. GC sequence composition for brain-specific promoters and house-keeping (hkg) promoters.

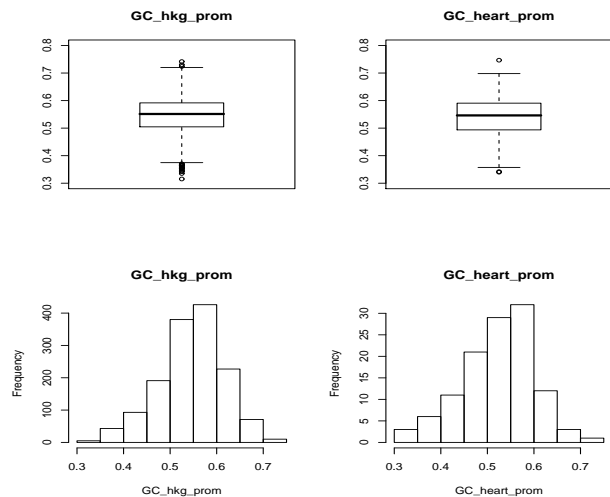


Fig. 6. GC sequence composition for heart-specific promoters and house-keeping (hkg) promoters.

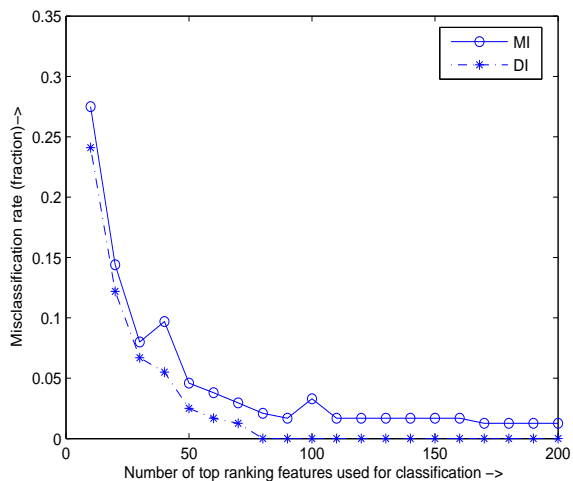


Fig. 5. Misclassification accuracy for the MI vs. DI case (brain promoter set). Accuracy of classification is ~ 0.9 i.e. 93%.

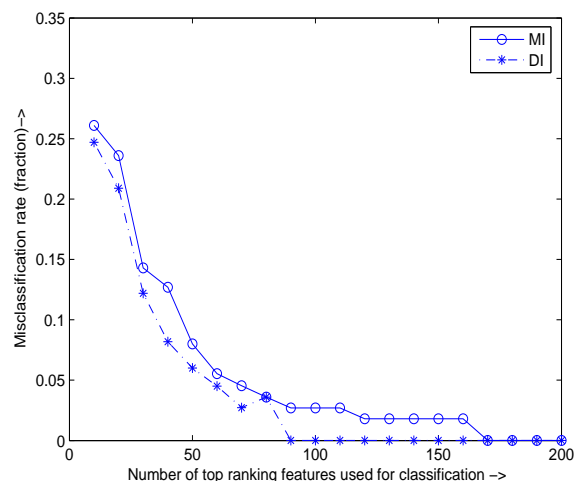


Fig. 7. Misclassification accuracy for the MI vs. DI case (heart promoter set).

these TFs are expressed in the corresponding tissue [26], a (*) sign is appended. In some cases, the hexamer motifs match the consensus sequences of known transcription factors (TF). This suggests a potential role for that particular TF in regulating expression of tissue-specific genes. This matching of hexamer motifs with TFBS consensus sites is done using the MAPPER engine (<http://bio.chip.org/mapper/>). It is to be noted that a hexamer-TFBS match does not necessarily imply the functional role of the TF in the corresponding tissue (brain or heart). However, such information would be useful to guide focused experiments to confirm their role in-vivo (using techniques such as chromatin immunoprecipitation).

As is clear from the above results, there are several other motifs which are novel or correspond to non-consensus motifs of known transcription factors. Hence, each of the identified hexamers merit experimental investigation. Also, though we identify as many as 200 hexamers in this work (please see Supplementary data), we have reported only a few due to space

constraints.

In the context of the heart-specific genes, we consider the cardiac troponin gene (*cTNT*, ENSEMBL:ENSG00000118194), which is present in the heart promoter set. An examination of the high DI motifs for the heart-specific set yields motifs with the GATA consensus site, as well as matches with the MEF2 transcription factor. It has been established earlier that GATA-4, MEF2 are indeed involved in transcriptional activation of this gene [25] and the results have been confirmed by ChIP [2].

B. Enhancer DB

Additionally, all the brain-specific regulatory elements profiled in the mouse Enhancer Browser database (<http://enhancer.lbl.gov/>), are examined for discriminating motifs. Fig. 8 shows that the two classes have similar GC-composition. Again, the plot of misclassification accuracy

Brain promoters	Heart promoters	Brain enhancers
Ahr-ARNT (*)	Pax2	HNF-4 (*)
Tcf11-MafG (*)	Tcf11-MafG (*)	Nkx2
c-ETS (*)	XBPI (*)	AML1
FREAC-4	Sox-17 (*)	c-ETS (*)
T3R-alpha1	FREAC-4	Elk1 (*)
	GATA(*)	

TABLE II

COMPARISON OF HIGH RANKING MOTIFS (BY DI) ACROSS DIFFERENT DATA SETS. THE (*) SIGN INDICATES TISSUE-SPECIFIC EXPRESSION OF THE CORRESPONDING TF GENE.

vs. number of features in the MI and DI scenarios reveal the superior performance of the DI-based hexamer selection compared to MI (Fig. 9).

In this case, the enhancer sequences are ultraconserved, thus obtained after alignment across multiple species. The examination of these sequences identified motifs that are potentially selected for regulatory function across evolutionary distances. Using alignment as a prefiltering strategy helps remove bias conferred by sequence elements that arise via random mutation but might be over-represented. This is permitted in programs like Toucan [1] and rVISTA (<http://rvista.dcode.org/>).

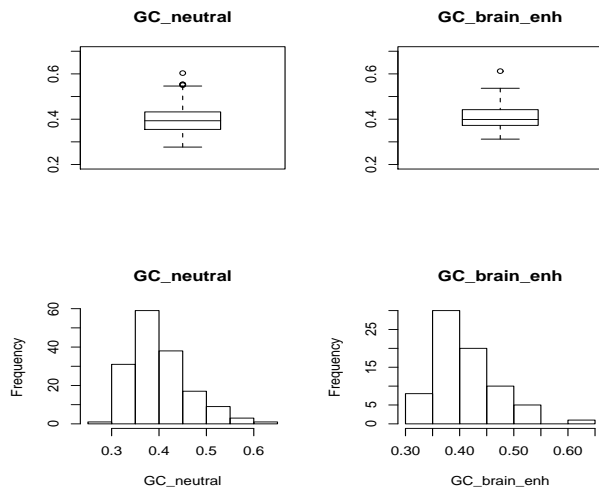


Fig. 8. GC sequence composition for brain-specific enhancers and neutral non-coding regions.

As in the previous case, some of the top ranking motifs from this dataset are also shown in Table II. The (*) signed TFs indicate that some of these discovered motifs indeed have documented high expression in the brain. The occurrence of such tissue-specific transcription factor motifs in these regulatory elements gives credence to the discovered motifs. For example, *ELK-1* is involved in neuronal differentiation [39]. Also, some motifs matching consensus sites of TEF1 and ETS1 are common to the brain-enhancer and brain-promoter set. Though this is interesting, an experiment to confirm the enrichment of such transcription factors in the population of brain-specific regulatory sequences is necessary.

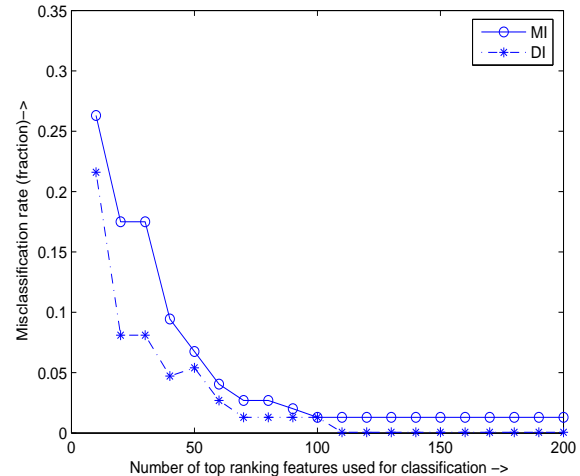


Fig. 9. Misclassification accuracy for the MI vs. DI case (brain enhancer set).

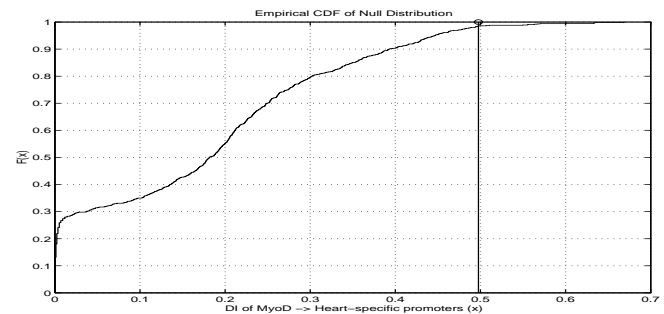


Fig. 10. Cumulative Distribution Function for bootstrapped $I(\text{MyoD motif: } CACCTG \rightarrow Y)$; Y is the class label (Heart-specific vs. Housekeeping). True $\hat{I}(CACCTG \rightarrow Y) = 0.4977$.

C. Quantifying sequence-based TF influence

A very interesting question emerges from the above presented results. What if one is interested in a motif that is not present in the above ranked hexamer list for a particular tissue-specific set? As an example, consider the case for *MyoD*, a transcription factor which is expressed in muscle and has a putative activity in heart-specific genes [28]. In fact, a variant of its consensus motif - CATTG is indeed

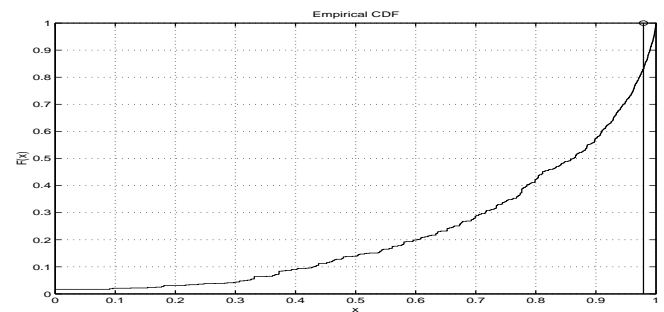


Fig. 11. Cumulative Distribution Function for bootstrapped $I(\text{Pax2 motif: } GTTCC \rightarrow Y)$; Y is the class label (UB/non-UB). True $\hat{I}(GTTCC \rightarrow Y) = 0.9792$.

in the top ranking hexamer list. The DI based framework further permits investigation of the directional association of the canonical *MyoD* motif (CACCTG) for the discrimination of heart-specific genes vs. housekeeping genes. This is shown in Fig. 10. As is observed, *MyoD* has a significant directional influence on the heart-specific vs. neutral sequence class label. This, in conjunction with the expression level characteristics of *MyoD*, indicates that the motif CACCTG is potentially relevant to make the distinction between heart-specific and neutral sequences.

Another theme picks up on something quite traditionally done in bioinformatics research - finding key TF regulators underlying tissue-specific expression. Two major questions emerge from this theme.

- 1) Which putative regulatory TFs underlie the tissue-specific expression of a group of genes?
- 2) For the TFs found using tools like TOUCAN [1], can we examine the degree of influence that the particular TF motif has in directing tissue-specific expression?

- To address the *first* question, we examine the TFs revealed by DI/MI motif selection and compare these to the TFs discovered from TOUCAN [1], underlying the expression of genes expressed on day *e14.5* in the degenerating mesonephros and nephric duct (TS22). This set has about 43 genes (including *Gata2*). These genes are available in the Supplementary data.

Using TOUCAN, the set of module TFs are combinations of the following TFs: *E47*, *HNF3B*, *HNF1*, *RREB1*, *HFH3*, *CREBP1*, *VMYB*, *GF11*. These were obtained by aligning the promoters of these 43 genes (-2000bp upstream to +200bp from the TSS), and looking for over-represented TF motifs based on the TRANSFAC/JASPAR databases.

Using the DI based motif selection, a set of 200 hexamers are found that discriminate these 43 gene promoter sequences from the background housekeeping promoter set. They map to the consensus sites of several known TFs, such as (identified from bio.chip.org/mapper/) *Nkx*, *Max1*, *c-ETS*, *FREAC4*, *Ahr-ARNT*, *CREBP2*, *E2F*, *HNF3A/B*, *NFATc*, *Pax2*, *LEF1*, *Max1*, *SP1*, *Tef1*, *Tcf11-MafG*, many of which are expressed in the developing kidney (<http://www.expasy.org/>). Moreover, we observe that the TFs that are common between the TOUCAN results and the DI based approach: *FREAC4*, *Max1*, *HNF3a/b*, *HNF1*, *SP1*, *CREBP*, *RREB1*, *HFH3* are mostly kidney-specific. Thus, we believe that this observation makes a case for finding all (possibly degenerate) TF motif searches from TRANSFAC, and filtering them based on tissue-specific expression subsequently. Such a strategy yields several more TF candidates for testing and validation of biological function.

- For the *second* question, we examine the following scenario. The *Gata3* gene is observed to be expressed in the developing ureteric bud (UB) during kidney development. To find UB specific TF regulators, conserved TF modules can be examined in the promoters of UB-specific genes. These experimentally annotated UB-specific genes are

obtained from the Mouse Genome Informatics database at <http://www.informatics.jax.org/>. Several programs are used for such analysis, like Genomatix [41] or Toucan [1]. Using Toucan, the promoters of the various UB specific genes are aligned to discover related modules. The top-ranking module in Toucan contains *AHR-ARNT*, *Hox13*, *Pax2*, *Tal1alpha-E47*, *Oct1*. Again, the power of these motifs to discriminate UB-specific and non-specific genes, based on DI, can be investigated.

For this purpose, we check if the *Pax2* binding motif (GTTCC [7]) indeed induces kidney specific expression by looking for the strength of DI between the GTTCC motif and the class label (+1) indicating UB expression (Fig. 11). This once again adds to computational evidence for the true role of *Pax2* in directing ureteric bud specific expression [7]. The main implication here is that, from sequence data, there is strong evidence for the *Pax2* motif being a useful feature for UB-specific genes. This is especially relevant given the documented role of *Pax2* ([9]) directing ureteric-bud expression of the *Gata3* gene, one of the key modulators of kidney morphogenesis. Both the *MyoD* and *Pax2* studies indicate the relevance of principled data integration using expression ([36],[26]) and sequence modalities.

D. Observations

With regard to the feature selection and classification results, in both studies (enhancers and promoters), we observe that about 100 hexamers are enough to discriminate the tissue-specific from the neutral sequences. Furthermore, some sequence features of these motifs at the promoter/enhancer emerge.

- There is higher sequence variability at the promoter since it has to act in concert with LREs of different tissue types during gene regulation.
- Since the enhancer/LRE acts with the promoter to confer expression in only one tissue type, these sequences are more specific and hence their mining identifies motifs that are probably more indicative of tissue-specific expression.

We however, reiterate that the enhancer dataset that we study uses the *hsp68-lacZ* as the promoter driven by the ultraconserved elements. Hence there is no promoter specificity in this context. Though this is a disadvantage and might not reveal all key motifs, it is the best that can be done in the absence of any other comprehensive repository.

The second aspect of the presented results highlight two important points. Firstly, the identified motifs have a strong predictive value as suggested by the cross-validation results as well as Table II. Moreover, DI provides a principled methodology to investigate any given motif for tissue-specificity as well as for identifying expression-level relationships between the TFs and their target genes, (Sec XI : C).

XII. CONCLUSIONS

In this work, a framework for the identification of hexamer motifs to discriminate between two kinds of sequences (tissue-specific promoters or regulatory elements vs non-specific

elements), is presented. For this feature selection problem, a new metric - the ‘directed information’ (DI) is proposed. In conjunction with a support vector machine classifier, this method was shown to outperform the state-of-the-art method employing undirected mutual information. We also find that only a subset of the discriminating motifs correlate with known transcription factor motifs and hence the other motifs might be potentially related to non-consensus TF binding or underlying epigenetic phenomena governing tissue-specific gene expression. The superior performance of the directed-information based variable selection suggests its utility to more general learning problems. As per the initial motivation, the discovery of these motifs can aid in the prospective discovery of other tissue-specific regulatory regions.

We have also examined the applicability of DI to prospectively resolve the functional role of any TF motif in a biological process, integrating other sources (literature, expression data, module searches).

XIII. FUTURE WORK

Several opportunities for future work exist within this proposed framework. Multiple sequence alignment of promoter/regulatory sequences across species would be a useful preprocessing step to reduce false detection of discriminatory motifs. The hexamers can also be identified based on other metrics exploiting distributional divergence between the samples of the “+1” and “-1” classes. Furthermore, there is a need for consistent high-dimensional entropy estimators within the small sample regime. A very interesting direction of potential interest is the formulation of a stepwise hexamer selection algorithm, using the directed information for maximal relevance selection and mutual information for minimizing between-hexamer redundancy [31]. This analysis is beyond the scope of this work but an implementation is available from the authors for further investigation.

XIV. ACKNOWLEDGEMENTS

The authors gratefully acknowledge the support of the NIH under award 5R01-GM028896-21 (J.D.E). We would like to thank Prof. Sandeep Pradhan and Mr. Ramji Venkataramanan for useful discussions on directed information. We are extremely grateful to Prof. Erik Learned-Miller and Dr. Damian Fermin for sharing their code for high-dimensional entropy estimation and ENSEMBL sequence extraction, respectively. We also thank the anonymous reviewers and the corresponding editor for helping us to improve the quality of the manuscript through insightful comments and suggestions.

XV. AVAILABILITY

The source code of the analysis tools (in R 2.0 and MATLAB 6.1) is available on request. Supplementary data (the set of various tissue-specific promoters, enhancers, housekeeping promoters and hexamer lists for each category) will be made available from the publisher’s website.

REFERENCES

- [1] Aerts S, Van Loo P, Thijs G, Mayer H, de Martin R, Moreau Y, De Moor B., “TOUCAN 2: the all-inclusive open source workbench for regulatory sequence analysis”, *Nucleic Acids Res.* 2005 Jul 1;33(Web Server issue):W393-6.
- [2] Azakie A, Fineman JR, He Y., “Myocardial transcription factors are modulated during pathologic cardiac hypertrophy in vivo”, *J Thorac Cardiovasc Surg.* 2006 Dec;132(6):1262-71.
- [3] Benjamini, Y. and Hochberg, Y. (1995)., “Controlling the false discovery rate: a practical and powerful approach to multiple testing”, *J. Roy. Statist. Soc. Ser. B* 57 289-300.
- [4] Burge C, Karlin S, “Prediction of complete gene structures in human genomic DNA”, *J Mol Biol* 1997, 268:78-94.
- [5] Chan BY, Kibler D., “Using hexamers to predict cis-regulatory motifs in *Drosophila*”, *BMC Bioinformatics*, 2005 Oct 27;6:262.
- [6] Cover TM, Thomas JA, *Elements of Information Theory*, Wiley- Interscience, 1991.
- [7] Dressler, G.R. and Douglas, E.C. (1992)., “Pax-2 is a DNA-binding protein expressed in embryonic kidney and Wilms tumor”, *Proc. Natl. Acad. Sci. USA* 89: 1179-1183.
- [8] Efron B, Tibshirani R.J, *An Introduction to the Bootstrap* (Monographs on Statistics and Applied Probability), Chapman & Hall/CRC, 1994.
- [9] Grote D, Souabni A, Busslinger M, Bouchard M., “Pax 2/8-regulated Gata 3 expression is necessary for morphogenesis and guidance of the nephric duct in the developing kidney”, *Development.* 2006 Jan;133(1):53-61.
- [10] I. Guyon, A. Elisseeff, “An Introduction to Variable and Feature Selection”, *Journal of Machine Learning Research* 3 1157-1182.
- [11] Hastie T, Tibshirani R, *The Elements of Statistical Learning*, Springer 2002.
- [12] Hutchinson GB., “The prediction of vertebrate promoter regions using differential hexamer frequency analysis”, *Comput Appl Biosci.* 1996 Oct;12(5):391-8.
- [13] H. Joe., “Relative entropy measures of multivariate dependence”. *J. Am. Statist. Assoc.*, 84:157164, 1989.
- [14] Kadota K, Ye J, Nakai Y, Terada T, Shimizu K., “ROKU: a novel method for identification of tissue-specific genes”, 2003. *BMC Bioinformatics.* 2006 Jun 12;7:294.
- [15] M. G. Kendall, “A New Measure of Rank Correlation”, *Biometrika*, 30 (1/2): 81-93 Jun., 1938.
- [16] 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.* 2004 Dec;24(23):10263-76.
- [17] King DC, Taylor J, Elnitski L, Chiaromonte F, Miller W, Hardison RC., “Evaluation of regulatory potential and conservation scores for detecting cis-regulatory modules in aligned mammalian genome sequences”, *Genome Res.* 2005 Aug;15(8):1051-60.
- [18] Kleinjan DA, van Heyningen V., “Long-range control of gene expression: emerging mechanisms and disruption in disease”, *Am J Hum Genet.* 2005 Jan;76(1):8-32.
- [19] Kreiman G., “Identification of sparsely distributed clusters of cis-regulatory elements in sets of co-expressed genes”, *Nucleic Acids Res.* 2004 May 20;32(9):2889-900.
- [20] Lakshmanan, G., K. H. Liew, K. C. Lim, Y. Gu, F. Grosveld, J. D. Engel, and A. Karis. 1999. “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.
- [21] MacIsaac KD, Fraenkel E., “Practical strategies for discovering regulatory DNA sequence motifs”, *PLoS Comput Biol.* 2006 Apr;2(4):e36.
- [22] H. Marko, “The Bidirectional Communication Theory - A Generalization of Information Theory”, *IEEE Transactions on Communications*, Vol. COM-21, pp. 1345-1351, 1973.
- [23] J. Massey, “Causality, feedback and directed information,” *Proc. 1990 Symp. Information Theory and Its Applications (ISITA-90)*, Waikiki, HI, Nov. 1990, pp. 303305.
- [24] Miller E, “A new class of entropy estimators for multi-dimensional densities”, *Proc. International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2003.
- [25] Murphy AM, Thompson WR, Peng LF, Jones L., “Regulation of the rat cardiac troponin I gene by the transcription factor GATA-4”, *Biochem J.* 1997 Mar 1;322 (Pt 2):393-401.
- [26] NCBI Pubmed URL: <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi>
- [27] Nemenman, F Shafee, and W Bialek. “Entropy and inference, revisited”, In TG Dietterich, S Becker, and Z Ghahramani, editors, *Advances in Neural Information Processing Systems* 14, Cambridge, MA, 2002. MIT Press.
- [28] Olson EN., “Regulation of muscle transcription by the MyoD family. The heart of the matter”, *Circ Res.* 1993 Jan;72(1):1-6.
- [29] Paninski, L. , “Estimation of entropy and mutual information”. *Neural Computation* 15: 1191-1254, 2003.
- [30] Pennacchio LA, Loots GG, Nobrega MA, Ovcharenko I., “Predicting tissue-specific enhancers in the human genome”, *Genome Res.* 2007 Jan 8;

- [31] Peng H., Long F., Ding C., "Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, No. 8, pp: 1226-1238, August 2005.
- [32] Pennacchio, L. A., 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., and Rubin, E.M. 2006., "In vivo enhancer analysis of human conserved non-coding sequences", *Nature*, 444(7118):499-502.
- [33] Proc. NIPS 2006 Workshop on Causality and Feature Selection, available at: <http://research.ihost.com/cws2006/>.
- [34] Li Q, Barkess G, Qian H., "Chromatin looping and the probability of transcription"., *Trends Genet.* 2006 Apr;22(4):197-202.
- [35] J. Ramsay, B. W. Silverman, *Functional Data Analysis* (Springer Series in Statistics), Springer 1997.
- [36] Rao A, Hero AO, States DJ, Engel JD, "Inference of biologically relevant Gene Influence Networks using the Directed Information Criterion", *Proc. of the IEEE Conference on Acoustics, Speech and Signal Processing*, 2006.
- [37] Schug J, Schuller W-P, Kappen C., Salbaum J.M., Bucan M., Stoeckert C.J. Jr., "Promoter Features Related to Tissue Specificity as Measured by Shannon Entropy"., *Genome Biology* 6(4): R33, March 2005.
- [38] Sumazin P, Chen G, Hata N, Smith A D., Zhang T, Zhang M Q., "DWE: Discriminant Word Enumerator", *Bioinformatics*, 21(1):31-38, 2005.
- [39] Vanhoutte P, Nissen JL, Brugg B, Gaspera BD, Besson MJ, Hipskind RA, Caboche J., "Opposing roles of Elk-1 and its brain-specific isoform, short Elk-1, in nerve growth factor-induced PC12 differentiation"., *J Biol Chem.* 2001 Feb 16;276(7):5189-96.
- [40] Venkataramanan, R.; Pradhan, S. S., "Source Coding With Feed-Forward: Rate-Distortion Theorems and Error Exponents for a General Source," *IEEE Transactions on Information Theory*, vol.53, no.6, pp.2154-2179, Jun. 2007.
- [41] Werner T., "Regulatory networks: Linking microarray data to systems biology"., *Mech Ageing Dev.* 2007 Jan;128(1):168-72.
- [42] Willett R, Nowak R, "Complexity-Regularized Multiresolution Density Estimation", *Proc. Intl Symp. on Information Theory, ISIT 2004*.