

Thursday, October 11, 2007

Welcome / Registration / Poster Set Up – 5:00-5:20pm

5:30 pm	Keynote: <u>Michael Eisen</u>	22
6:00 pm	<u>Reijman/Shamir</u> : Evolution of TF combinations	56
6:15 pm	<u>Cordero/Stormo</u> : Motifs using Phylo & 3D struct.....	18
6:30 pm	<u>Berger/Bulyk</u> : Motifs using PBMs	8

Poster Session I – 7:00-9:00pm

Authors of odd-numbered posters present 7pm - 8pm

Authors of even-numbered posters present 8pm - 9pm

Friday, October 12, 2007

Breakfast – 8am

9:00 am	Keynote: <u>George Church</u>	14
9:30 pm	<u>Ernst/Bar-Joseph</u> : Dynamic regulatory maps.....	4
9:45 am	<u>Fan/Liu</u> : Bayesian Cell Cycle.....	23
10:00 am	<u>Bais/Vingron</u> : Nucleosome Depletion & TFBS	3

Coffee / Snacks / Fruit Break – 10:15-10:45am

10:45 am	<u>Palin/Ukkonen</u> : KA stats for CRMs.....	52
11:00 am	<u>Ward/Bussemaker</u> : Alignment free TFBS prediction.....	70
11:15 am	<u>Ray/Xing</u> : Multi-resolution phylo shadowing.....	75
11:30 am	Keynote: <u>Ewan Birney</u>	10

Lunch Break / Networking Opportunities – 12-1pm

1:00 pm	<u>Stark/Kellis</u> : Regulatory motifs and targets in 12 Flies	37
1:15 pm	<u>Jaeger/Bulyk</u> : CRM targets.....	34
1:30 pm	<u>Yeger-Lote/Fraenkel</u> : Perturbations Expression Binding ...	76
1:45 pm	Keynote: <u>Erin O'Shea</u>	51

Poster Session II – 2:15-4:00pm

2:30-3:15pm - Authors of even-numbered posters present

3:15-4:00pm - Authors of odd-numbered posters present

4:00 pm	Keynote: <u>Nir Friedman</u>	26
4:30 pm	<u>Sriswasdi/Ge</u> : Interactome network motifs	27
4:45 pm	<u>Zhang/Zhang</u> : Splicing motif targets.....	80
5:00 pm	<u>Xiao/Burge</u> : Co-evolution splicing networks	74

Sunset Break – 5:15-6:15pm

Friday, October 12, 2007 – evening session

6:15 pm	<u>Kumar/Young</u> : Polycomb switch differentiation.....	38
6:30 pm	<u>Thurman/Stam</u> : DNase chromatin domains map	64
6:45 pm	<u>Xi/Weng</u> : DNase tissue-specific structures	71
7:00 pm	Keynote: <u>Michael Levine</u>	41

Conference Reception – 7:30-9:30pm

Jazz band and Hors-d'Oeuvres

Saturday, October 13, 2007

Breakfast – 8am

9:00 am	Keynote: <u>David Bartel</u>	5
9:30 am	<u>Ebert/Lauffenberger</u> : miRNA sponge inhibitors.....	21
9:45 am	<u>Yousef/Benos</u> : miRNA upstream regulators.....	6
10:00 am	<u>Kheradpour/Kellis</u> : miRNA characterization	37

Coffee / Snacks / Fruit Break – 10:15-10:45am

10:45 am	<u>Zheng/Burge/Sharp</u> : Novel short RNAs	82
11:00 am	<u>Kapranov/Gingeras</u> : Promoter assoc. RNAs.....	36
11:15 am	<u>O'Neil</u> : Ultra-conserved centromere virus RNA.....	50
11:30 am	<u>Bernick/Lowe</u> : RNA-mediated silencing Archaea.....	9

Lunch Break / Networking Opportunities – 11:45-12:45pm

12:45 pm	Keynote: <u>Trey Ideker</u>	33
1:15 pm	<u>Mogno/Gardner</u> : Temporal logic Ecoli sugar reg.....	46
1:30 pm	<u>Lim/Califano</u> : Druggable apoptosis regulators	43
1:45 pm	Keynote: <u>James Collins</u>	17

Closing remarks + adjourn – 2:15 pm

An integrative strategy for characterizing transcriptional regulatory networks involved in *Drosophila* embryonic mesoderm development

Anton Aboukhali^{1,5}, Brian W. Busser⁶, Savina A. Jaeger¹, Aditi Singhanía⁶, Michael F. Berger^{1,4}, Stephen S. Gisselbrecht¹, Alan M. Michelson⁶, Martha L. Bulyk¹⁻⁴

¹*Division of Genetics, Department of Medicine;* ²*Department of Pathology;* ³*Harvard/MIT Division of Health Sciences and Technology (HST); Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115;* ⁴*Harvard University Graduate Biophysics Program, Cambridge, MA 02138;* ⁵*MIT Aeronautics & Astronautics Department, Cambridge, MA 02139;* ⁶*National Heart, Lung and Blood Institute, National Institutes of Health, Bethesda, MD 20892.*

Elucidating transcriptional regulatory networks underlying progressive cell fate determination in embryonic development is of fundamental biological importance. Myogenesis in the *Drosophila* embryo is a powerful model for such investigation. Our approach is integrative and iterative, combines wet lab and computational methods, and is generalizable to other systems. First, the genes expressed in embryonic mesodermal cells are identified via flow cytometry cell purification, genetic perturbation, genome-wide expression profiling, and resolution of spatiotemporal patterns by *in situ* hybridization (ISH). New myoblast transcription factors (TFs) are identified via expression profiles, and a high-throughput *in vitro* protein binding microarray (PBM) technology is used to determine their binding specificities. Previously identified and experimentally-derived findings are integrated within a computational framework to predict: (a) novel cis-regulatory modules (CRMs) associated with coexpressed genes (via PhylCRM algorithm, which quantifies motif clustering and evolutionary conservation) (b) mapping motif combinations (MCs) to their target genes (via Lever algorithm, which evaluates the enrichment of MCs within predicted CRMs in the noncoding sequences surrounding genes in a collection of gene sets) (*Warner et al, in review*). Predicted CRMs, MCs, and individual motifs are validated *in vivo*. Findings are refined by iterative application of our strategy to improve its specificity/sensitivity. Finally, a search is imposed for shared novel motifs within candidate and validated CRMs to further classify them and reveal their cellular specificity transcriptional codes.

Our prior studies developed genomic methods for finding potentially relevant gene sets (*Estrada et al, PLoS Genetics, 2006*). We identified the combination of Ets+Twi+Tin TF sites as a core code for expression of some genes in muscle founder cells (FCs) (*Philippakis et al, PLoS Comp Bio, 2006*), while recognizing that more information is needed to direct expression to specific and overlapping FC subsets. We are currently tackling the problem of specificity pertaining to the distinct effects of multiple homeodomain (HD)-containing TFs with similar binding specificity in partially distinct groups of FCs. For instance, we derived motifs for the HD TFs Slou, Ap, Msh & Lbl using PBMs. We then defined and ISH-validated many of their downstream targets. To validate the functional significance of the Slou site, we knocked it out in the *ndg* and *mib2* enhancers and noted specific alterations in FC gene expression. Finally, we computationally derived a potential regulatory code of Ets+Twi+Tin+Slou for a set of Slou-responsive FC genes.

The oncogene IKBKE is overexpressed in only two breast cancer subtypes: the Her2+ subtype with lymphocytic infiltration and the basal-like subtype

Gabriela Alexe*, Nilay Sethi*, Lyndsay Harris, Shridar Ganesan and Gyan Bhanot

Affiliations: GA: The Broad Institute of MIT and Harvard, Cambridge, MA 02142; USA; GB, GB: The Simons Center for Systems Biology, IAS, Princeton, NJ 08540 USA; NS: Robert Wood Johnson University Hospital/UMDNJ, New Brunswick, NJ 08903 USA & Dept. of Mol. Bio., Princeton University, Princeton, NJ 08540 USA; LH: Yale Cancer Center, Yale University, New Haven, CT 06519 USA; GB: BioMaPS Institute & Dept. of BME, Rutgers University, Piscataway, NJ 08854 USA; GB, SG: Cancer Institute of New Jersey, New Brunswick, NJ 08903 USA.

In a recent study [1], Boehm et al showed that IKBKE, a kinase in the NF- κ B pathway, is an oncogene over-expressed in a subset of breast cancers. They also found an allelic gain of chr 1q. Using microarray data from Wang et al [2] and Richardson et al [3], we have recently shown that [4] node negative breast cancers treated only with surgery and radiation separate into eight subtypes. One of the HER2+ subtypes (HER2₊) was characterized by a lymphocytic infiltrate [4, 5] correlated with improved natural history (89 % vs 52% 10 year disease free survival) compared to the other (HER2_{NI}). Using an independent HER2+ dataset [5], we verified [4] the correlation between HER2₊ and a lymphocytic infiltrate. We also found that Basal-like tumors separate into two subtypes (BA₁ and BA₂), with BA₁ characterized by upregulation of the interferon pathway, suggesting that this subtype also elicits a differential immune response compared to the BA₂ subtype. Our Basal-like subtypes correlate well with the X chromosomal alterations identified in such tumors in a recent study [3].

In this paper we study the hypothesis that upregulation of IKBKE (and the NF- κ B pathway [1]) is correlated with immune system activation.

Using three independent datasets [2,3,6] we found that IKBKE is up-regulated only in the Basal-like (BA₁ more than BA₂) and HER2₊ subtypes. Interestingly, IKBKE was shown [1] to activate the interferon pathway in some breast cancers, which correlates with the immune signature of the BA₁ subtype. Mapping the gene expression data to chromosomes, we found that HER2₊ exhibits amplicons in the chr1q32 region (FLJ13310, CR2, CR1, IL10, IL19, IL24, FAIM3, PIGR, YOD1, PFKFB2, C4BPA) in the vicinity of the IKBKE gene, as well as in the chr1q23 region. The latter region includes interferon gamma, immunoglobulin and several antigen genes (FCRL2, KIRREL, ELL2, CD1D, CD1A, CD1C, CD1B, CD1E, SPTA1, MNDA, PYHIN1, IFI16, AIM2, IGSF4B, FY, FCER1A).

Our results enrich the analysis in [1] by showing that in IKBKE activates the NF- κ B pathway only Basal-like and HER2₊ subtypes. They also suggest the possibility of a causal connection between immune activation in BA₁ and HER2₊ [4] and IKBKE/ NF- κ B activation. We are in the process of validating end extending our results in-vivo using breast cancer tissue microarray using cell lines MCF-7, MDA-MB-453/ MCF-10A etc as case/controls for IKBKE activation. Results from this analysis will be presented if available.

1. Boehm JS, Zhao JJ, et al. Cell. 2007 Jun 15;129(6):1065-79.
2. Wang Y, Klijn JG, et al. Lancet. 2005 Feb 19-25;365(9460):671-9.
3. Richardson AL, Wang ZC, et al. Cancer Cell. 2006 Feb;9(2):121-32.
4. Alexe G, Dalgin GS, et al. under review in Cancer Research.
5. Harris LN, You F, et al. Clin Cancer Res. 2007 Feb 15;13(4):1198-207.
6. Ivshina AV, George J, et al. Cancer Res. 2006 Nov 1; 66(21):10292-301.

Combining the prediction of nucleosome-binding sequences and transcription factor binding sites.

Abha Singh Bais, Ho-Ryun Chung, Dennis Kostka, Hugues Richard and Martin Vingron

Department of Computational Molecular Biology, Max Planck Institute for Molecular Genetics, Ihne str. 63-73, Berlin 14195, Germany.

Recently there has been a surge in experimental and computational efforts to characterize nucleosomal positions along the chromatin.

Increasing evidence points to a significant contribution of DNA sequence characteristics to nucleosomal positioning in vivo. Based on their sequence composition, certain DNA sequences are more likely to be bound by nucleosomes. In a recent work, Segal *et al.* employ characteristics of experimentally found nucleosome-bound sequences to predict nucleosome positions genome-wide in yeast. A probabilistic model of nucleosome-DNA interaction, akin to a position-weight matrix for transcription factor binding sites (TFBSs), is used to detect putative nucleosome positions.

The binding of nucleosome is implicated in hindering the binding of transcription factors by limiting sequence accessibility. It is hypothesized that functional binding sites of a transcription factor usually lie in nucleosome free regions. Hence, employing the knowledge of nucleosome positioning may aid in identifying putative binding sites that are more likely to be functional, while at the same time restricting the search space. Based on this hypothesis, we study in how far the combined knowledge of nucleosome as well as transcription factor binding preferences help in identifying putative locations of both. We propose a hidden markov model framework that employs the proposed probabilistic model of nucleosome-DNA interaction of Segal *et al.* along with the standard position-weight matrix for TFBSs to predict both binding site and nucleosome positions.

References: Segal *et al.* A genomic code for nucleosome positioning. *Nature*, 442, 772-778.

Reconstructing dynamic regulatory maps

Jason Ernst¹, Oded Vainas², Christopher T. Harbison³, Itamar Simon², Zoltan N. Oltvai⁴, Naftali Kaminski⁵ and Ziv Bar-Joseph¹

1. Machine Learning Department, School of Computer Science, Carnegie Mellon University, 5000 Forbes Ave. Pittsburgh, PA,

2. Dept. Molecular Biology, Hebrew University Medical School, Jerusalem, Israel

3. Whitehead Institute, Nine Cambridge Center, Cambridge, Massachusetts USA.

4. Department of Pathology, University of Pittsburgh, Pittsburgh, PA, USA

5 Simmons Center for Interstitial Lung Disease, University of Pittsburgh Medical School, Pittsburgh, PA, USA..

Even simple organisms have the ability to respond to internal and external stimuli. This response is carried out by a dynamic network of protein–DNA interactions that allows the specific regulation of genes needed for the response. We have developed a novel computational method termed the Dynamic Regulatory Events Miner (DREM) to reconstruct such networks. DREM uses an Input-Output Hidden Markov Model to model these regulatory networks while taking into account their dynamic nature. Our method works by identifying bifurcation points, places in the time series where the expression of a subset of genes diverges from the rest of the genes. These points are annotated with the transcription factors regulating these transitions resulting in a unified temporal map. Applying DREM to study yeast response to stress we derive dynamic models that are able to recover many of the known aspects of these responses. Predictions made by our method have been experimentally validated by PCR and ChIP-chip experiments leading to new roles for Ino4 and Gcn4 in controlling yeast response to stress. The temporal cascade of factors reveals common pathways and highlights differences between master and secondary factors in the utilization of network motifs and in condition specific regulation. We have also used DREM to reconstruct the Aerobic-Anaerobic switch network in *e. coli* and validated many of the predictions using new microarray data. Recently we have begun applying DREM to reconstruct networks in mammalian cells. Using motif data from comparative genomics studies and new time series expression data we reconstructed networks for the progression of fibrosis (an often deadly lung disease) and for immune response. The models recovered known and novel aspects regarding the regulation of these processes, some of which are being tested in follow up experiments.

MicroRNAs

David Bartel

Department of Biology, MIT/Whitehead Institute/HHMI, Cambridge, MA

MicroRNAs are small endogenous RNAs that can guide the posttranscriptional repression of protein-coding genes. We have been using molecular and computational approaches to find microRNAs in plants and animals, identify the messages that they repress, and then investigate their functions during development, oncogenesis, and other processes. This talk will feature recent results, such as the analyses of high-throughput sequencing of small RNAs, computational and experimental results revealing the widespread impact of miRNAs on mRNA expression and evolution, methods for identifying mRNAs that are most effectively repressed by miRNAs, and the biological roles of particular miRNA-target interactions.

Small but essential: let-7d miRNA in a transcriptional regulatory network that determines the fate of epithelial cells

Hanadie Yousef, David L. Corcoran, Daniel Handley, Kusum Pandit, Isidore Rigoutsos, Oliver Eickelberg, Naftali Kaminski, Panayiotis V. Benos

Department of Computational Biology, University of Pittsburgh, 3501 Fifth Avenue, Pittsburgh, Pennsylvania, USA and Department of Medicine, University of Pittsburgh, Pittsburgh, Pennsylvania, USA.

Transforming growth factor β (TGF β) is implicated in cancer, fibrosis, and embryonic development. SMAD3 mediates TGF β dependent dramatic effects on cellular phenotype by directly binding to DNA. microRNAs (miRNAs) are small (22-61 bp long), non-coding RNAs that downregulate their target genes *via* base complementarity to their mRNA molecules. So far our knowledge about the transcriptional regulation of miRNA genes is limited. Considering the significant role of TGF β in carcinogenesis, development and response to injury we hypothesized that TGF β also regulates miRNA expression. To address this hypothesis we analyzed upstream sequences of intergenic miRNA genes in vertebrates. Impressively, the conservation of their upstream sequences resembles the conservation of the promoters of the protein coding genes, supporting the notion that similar mechanisms and factors may regulate the expression of both types. The similarity of the conservation patterns is even more striking in the first 500 bp. Discriminative motif analysis of the 1 kb upstream region of the miRNA genes known to be expressed in the lung identified three clusters of genes with differences in the predicted regulatory sites. SMAD3 was predicted to regulate let-7d and miR-10a intergenic miRNA genes. Using ChIP and EMSA assays on A549 cells we confirmed the direct association of SMAD3 to the promoters of these two miRNAs. We determined that let-7d expression invariably decreases in TGF β treated cells and that inhibition of let-7d increases HMGA2, a TGF β inducible molecule critical to epithelial mesenchymal. In summary, we have shown that TGF β regulates miRNA expression through a regulatory cascade. The impressive down-stream effects of this regulation suggest that TGF β regulation of miRNAs has critical and profound effects on cellular and –potentially– organ phenotype.

Familial Binding Profiles Made Easy

Shaun Mahony, Philip E. Auron, [Panayiotis V. Benos](#)

Department of Computational Biology, University of Pittsburgh, 3501 Fifth Avenue, Pittsburgh, Pennsylvania, USA and Department of Biological Sciences, Duquesne University, Pittsburgh, Pennsylvania, USA.

Background.

Transcription factor (TF) proteins regulate the expression of their target genes by recognizing highly similar DNA binding sites in the promoters of these genes. The evolution of their DNA binding preferences has been the subject of many recent studies. Generalized motifs or *familial binding profiles* (FBPs) have been generated using semi-manual methods and have been used as priors to improve motif prediction algorithms or for the classification of newly discovered motifs.

Results.

STAMP (*Similarity, Tree-building, and Alignment of Motifs and Profiles*) is a comprehensive framework to facilitate the classification and evolutionary analysis of DNA motifs. We developed a new metric to determine the optimal number of clusters in a phylogenetic tree of DNA motifs and we used it to build an improved set of FBPs. STAMP offers the equivalent of BLAST and CLUSTALW for DNA motifs and its uses include the prediction of a TF that binds to a given motif and meta-analysis of motif finder outputs.

URL: <http://www.benoslab.pitt.edu/stamp>

High-Resolution Characterization of Mammalian Transcription Factor Binding Specificities Using Protein Binding Microarrays

Michael F. Berger^{1,4}, Gwenael Badis-Breard⁶, Andrew R. Gehrke¹, Shaheynoor Talukder⁵, Anthony A. Philippakis^{1,3,4}, Savina A. Jaeger¹, Esther Chan⁶, Genita Metzler¹, Anastasia Vedenko¹, Olga Botvinnik¹, Wen Zhang¹, Hanna Kuznetsov¹, Chi-Fong Wang¹, David Coburn⁶, Quaid D. Morris⁶⁻⁸, Timothy R. Hughes^{6,9}, and Martha L. Bulyk¹⁻⁵

¹Division of Genetics, Department of Medicine, ²Department of Pathology, and ³Harvard/MIT Division of Health Sciences & Technology (HST); Brigham & Women's Hospital and Harvard Medical School, Boston, MA, USA. ⁴Harvard University Biophysics Program, Cambridge, MA, USA. ⁵Broad Institute, Cambridge, MA, USA. ⁶Banting and Best Institute of Medical Research, and ⁷Departments of Computer Science, ⁸Electrical and Computer Engineering, and ⁹Medical Genetics and Microbiology, University of Toronto, Toronto, ON, Canada.

Sequence-specific protein-DNA interactions play a critical role in establishing complex temporal and spatial patterns of gene expression. A full understanding of these regulatory programs requires detailed knowledge of the DNA binding specificities of the transcription factors (TFs) involved. However, the binding sites for the vast majority of TFs are currently poorly characterized or completely unknown. To meet this need, we are using compact, universal protein binding microarrays (PBMs) to create a comprehensive dataset of the high-resolution *in vitro* binding specificities of mouse TFs. The uniform coverage of *k*-mer sequence variants on our microarrays enables us to simultaneously assay the binding preferences of any TF for all possible binding sites [Berger, Philippakis, *et al.*, *Nature Biotechnology*, 24:1429-1435 (2006)]. We have cloned the DNA-binding domains from over 1,000 distinct mouse TFs and have optimized strategies for protein expression both in *E. coli* and by *in vitro* transcription and translation. To date, we have successfully characterized the binding specificities of more than 350 of these TFs in high resolution. Using these data, we are able to compare the binding profiles among divergent TFs, both within a single structural class and across different families. Interestingly, we have identified several TFs that share the same highest-affinity binding sites but differ in their weaker affinity targets. We have developed a computational method to represent the preferences for individual *k*-mers by a position weight matrix (PWM), but we note that the measured binding specificity is often better captured by combinations of distinct PWMs. The majority of TFs studied show evidence for a significant "secondary" PWM due to interdependence between nucleotide positions or, in some cases, potentially due to alternate binding conformations. We have also begun to examine the relationship between amino acid sequence and DNA binding specificity, with a particular emphasis on the homeodomain structural class. Finally, we expect the TF binding data presented here to facilitate the more accurate prediction of *cis*-regulatory modules and inference of *cis*-regulatory codes [Warner, Philippakis, Jaeger, *et al.*, unpublished] as we try to interpret the intricate tissue-specific patterns of gene expression in mammals.

RNA-mediated silencing in Archaea

David L Bernick, and Todd M Lowe

Biomolecular Engineering, University of California, Santa Cruz, California, USA

Epigenetic control systems provide a mechanism to alter gene expression without modification of the underlying genome. Some of these systems likely came about in response to genomic compromise by invasive elements such as transposons, plasmids and viruses. RNA interference (RNAi) is a class of epigenetic control that recognizes the invading element and specifically silences the effect of the invasion. This system is widespread in eukaryotes but has so far not been identified within prokaryotes. One of the components of RNAi, the Argonaute protein, has been identified and verified among a few members of the euryarchaea (*Pyrococcus furiosus*, for example), in a bacterium, *Aquifex aeolicus*, and putative sequence homologs in five other bacterial phyla have been described. This suggests either an unrecognized prokaryotic function for this protein or that this versatile epigenetic control mechanism may be active in prokaryotes as well.

Within the eukaryotes, arrays of genomically encoded sequences sharing sequence homology to endogenous transposons have been identified at specific loci (flamenco). An RNA interference mechanism, employing the Piwi subfamily of the Argonaute proteins, utilizes this array as its source of complimentary sequence. An analogous prokaryotic array, the so called CRISPR array (clustered regularly interspaced short palindromic repeats), is found in roughly half of the sequenced eubacteria and in all sequenced archaea. Acquired sequences maintained within the CRISPR array have been shown to confer immunity against invading phage. Some archaeal species have over thirty instances of the array with some arrays containing over fifty unique sequence elements.

Using 454 sequencing of small RNA extracted from total RNA from exponentially growing *Pyrococcus furiosus* (Pfu) cells, we show that a class of small RNA (20-28 bases) exist that are antisense to the 5' end of Pfu transposases. The region involved spans the 5' UTR and initial coding sequence of the transposase transcript as has been seen with the bacterial *IS10* element. Additionally, sequences associated with CRISPR arrays were found (size range 18-28 bases), some of which are antisense to endogenous genes.

The mechanism that might make use of these small RNA remains unknown. Pfu has been shown to encode an Argonaute protein. We now have experimental evidence that Pfu also expresses small antisense RNA complimentary to endogenous genes. Future work will examine the possibility that these findings are related and that they may support an RNA-mediated silencing mechanism in this non-eukaryotic life form.

ENCODE: Understanding our genome

Ewan Birney

Sanger Center

The ENCODE pilot project has provided an unprecedented view into how the human genome functions. With over 33 groups contributing a variety of functional experiments on a targetted 1% of the human genome, the ENCODE project has revealed a far more complex, intercalated transcription, association with transcription factors and histone modifications and surprising correlations with evolutionary conservation.

I will present this work and discuss some on the consequences of this new understanding into human genome function.

A Method for Inferring Extracellular Environment from Gene Expression Profiles Using Metabolic Flux Balance Models

Aaron Brandes, Desmond S. Lun, Jeremy Zucker, and James E. Galagan

Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA

The programs of metabolic regulation used by the bacterium *E. coli* vary with carbon source, and result in differences in metabolic flux. For example the gene *aceA* is highly upregulated for growth on acetate resulting in significant flux through the glyoxylate shunt, in contrast to growth on glucose which results in virtually no flux through the shunt. The goal of this research is to couple metabolic models with gene expression sets resulting from the organism's program of regulation in order to identify the corresponding nutrient source. Genes key for making the differentiation can be identified. They are candidates for hierarchical regulation, and could be used to help identify regulatory networks. The algorithm developed here also has the potential to allow the use of expression data as probe for challenging environmental conditions such as *M. tuberculosis* in a phagosome.

Using a model of *E. coli* central carbon metabolism and gene expression data obtained by Liu et al. [1] for *E. coli* grown separately on six different carbon sources—glucose, glycerol, succinate, L-alanine, acetate, and L-proline we used a novel method to create a set of six corresponding flux cones. Each flux cone represents the limitations imposed on metabolic fluxes, assuming that relative enzyme activities between nutrient conditions are related to the expression levels. We then measured the distance between these flux cones and optimal solutions to a flux balance problem to successfully identify five of the six nutrient sources from the corresponding gene expression values.

1. Liu M, Durfee T, Cabrera JE, Zhao K, Jin DJ, et al. (2005) Global transcriptional programs reveal a carbon source foraging strategy by *Escherichia coli*. *J Biol Chem* 280: 15921-15927.

Characterization of genome-wide DNA binding by the replication initiator and transcription factor DnaA in response to perturbation of replication

Adam M. Breier and Alan D. Grossman

Department of Biology, Building 68-530, MIT, Cambridge, MA, USA 02139

Growing cells use multiple mechanisms to coordinate DNA replication, cell division, and development. In the *Bacillus subtilis*, an overlapping set of genes is affected by perturbations in replication elongation or initiation independently of the well-known RecA-mediated SOS response to DNA damage. Many of these genes are thought to be regulated by DnaA, the highly conserved replication initiation protein and transcription factor. We have characterized the genome-wide binding profile of DnaA using chromatin immunoprecipitation and genomic microarrays. Significant DnaA binding was observed at many sites around the chromosome, including upstream of at least 13 operons which responded transcriptionally to replication perturbation and contained sequence matches to the DnaA binding motif. At four loci, including the origin of replication, DnaA binding increased strongly in response to perturbation of replication; each of these loci contain multiple matches to the DnaA binding consensus sequence. The activity of DnaA as a transcriptional regulator appears to be controlled in part by modulation of its affinity for DNA. Current work is focused on the mechanisms by which DnaA senses replication status.

Selective optimization of codon usage detected by information theory based codon bias measurements

Zehua Chen, and Jeffrey Chuang

Biology Department, Boston College, Chestnut Hill, MA, USA

Most amino acids are encoded by multiple synonymous codons. In the standard genetic code, the synonymous codon groups contain one to six codons that differ mostly at the third codon position. It is widely known that different codons of the same synonymous group, and different synonymous groups as well, are not used with equal frequencies, and that the biased codon usage differs within and across genomes. Several factors, including GC content, mRNA expression and stability, tRNA abundance and translation efficiency, have been proposed to account for the codon usage biases.

Information theory has been widely used in various biological systems, especially in transcription factor binding site modeling and prediction. In this study, novel information measurements were developed to quantify amino acid usage bias (R_m) and synonymous codon usage bias (R_w) for individual genes or genomes. By applying the information analysis to 487 bacterial genomes with different GC contents, we find that R_m and R_w correlate closely with the background information at the first and third codon positions, respectively, while correlations with the second position are different. The base composition at second position is less variable, which defers control of amino acid usage and synonymous codon bias to the first and third positions.

When correlating R_m and R_w with gene expression data within five species, the individual amino acid biases ($R_m(i)$) and synonymous biases ($R_w(i)$) give highly variant correlations, but some conserved correlation patterns are shared in even distantly related species. The highly variant correlations of individual biases may be a result of selective optimization of synonymous codon groups coding for lower cost amino acids.

We are currently processing microarray datasets for more species. We aim to compute gene expression index from cross-laboratory and crossplatform microarray expression datasets for more than 50 species. Then we can investigate correlation between codon bias and gene expression in a larger scale.

New genomic tools for reading, writing & regulation.

George Church

Harvard University

2nd Generation Sequencing platforms typically employ polymerase-colony cyclic Sequencing by ligation (SbL) or by polymerase (SbP). We have explored these via an "open-source" approach to software, hardware, wetware, and ELShare (<http://arep.med.harvard.edu/Polonator/>). Applications of polonies include finding rare mutations in evolved microbial and cancer genomes, RNA quantitation, genome rearrangements, chromosome conformation capture (5C) and exon resequencing.

The Personal Genome Project (PGP) aims to correlate medical history, identifiable traits, and allele-specific regulatory genomics on PGP cell lines (now at Coriell with no commercial or privacy restrictions). Automated homologous allele replacement has been developed to allow exploration of promoter combinatorials and new genetic codes for multi-virus resistance.

A nonlinear dynamical model to infer transcriptional regulatory networks from time dependent expression data

Adriana Climescu-Haulica¹ and Michelle Quirk²

¹ Laboratoire Jean Kuntzmann, Grenoble 38041 France

adriana.climescu@imag.fr

² Los Alamos National Laboratory, Los Alamos, New Mexico, U.S.A.

A foreground level of understanding the regulatory mechanism is to decipher the learning machinery enabling the transcriptional program to adapt in time as the cell progresses through development or undergoes environmental changes. An actual way to trace dynamic features on the relationship between genes and their regulators is to analyze time dependent microarray gene expression data obtained in pertinent conditions. The similarity of the pattern for two gene expression time series do not implies necessarily a causal regulatory relationship between them. Recent studies [3, 2] show that this type of relationship could be inferred from a mathematical model which fit the transcription pattern of a specific target gene using a set of putative regulators. In this perspective we propose a novel model derived as a nonlinear stochastic differential equation which has two novel particularities: a) it considers a kinetic interaction model for the decay rate accounting for the mRNA degradation; this leads to a nonlinear term on the stochastic differential equation modeling the target gene mRNA; b) for each target gene it proposes a choice for the prototype of the regulatory function between the beta sigmoid function, designed to keep track of the local temporal patterns of the target gene regulators, and the sigmoid function, shaped around statistical parameters; this feature accommodates partially the variability of the regulatory pattern from one gene to another. Applied to the expression measurements of the mRNA levels of *Saccharomyces cerevisiae* [1] this model improves the fitting for 43% more genes with respect to the results from [2] and 32 % for results from [3], respectively. We provide also the analysis of time dependent gene expression measurements on *Drosophila melanogaster* embryogenesis [4]. For this data set we obtained a good fit for almost half of a set of 3418 genes analyzed.

1. Spellman, P.T. et al. (1998) Comprehensive identification of cell cycle-regulated genes of the yeast *S.cerevisiae* by micro-array hybridization, *Mol. Biol. Cell*, 9, 3273-3297

2. Chen, K-C. et al.(2005), A stochastic differential equation model for quantifying transcriptional regulatory network in *S.cerevisiae*, *Bioinformatics*, 21, 2883-2890

3. Climescu-Haulica, A. et al. (2007) A stochastic differential equation model for transcriptional regulatory networks, *BMC Bioinformatics*. 8(Suppl 5): S4.

4. www.fruitfly.org

Riboswitches, RNA conformational switches and prokaryotic gene regulation

Eva Freyhult^a, Vincent Moulton^b, Peter Clote^c

a Linnaeus Centre for Bioinformatics, Uppsala University, 75124 Uppsala, Sweden, eva.freyhult@lcb.uu.se, *b* School of Computing Sciences, University of East Anglia, Norwich, NR4 7TJ, UK, vincent.moulton@cmp.uea.ac.uk, *c* Department of Biology, Boston College, Chestnut Hill, MA 02467, USA, clote@bc.edu. This work is funded in part by NSF DBI-0543506.

Metabolite-sensing 5'-UTR (untranslated regions) of certain mRNAs, called riboswitches, have been discovered to undergo a conformational change upon ligand-binding, which thereby can up- or down-regulate the corresponding protein product. For instance, upon the binding of nucleotide guanine, the G-box riboswitch in the 5' UTR of the XPT gene of *Bacillus subtilis* undergoes a conformational change to create a terminator loop, thereby prematurely terminating transcription of the XPT gene. Since XPT is involved in guanine metabolism, this is an example of negative autoregulation by a riboswitch. Although riboswitches have been postulated to be an ancient genetic regulatory system, first developed in bacteria, the remarkable discovery of Cheah et al. in Nature 2007 suggests that eukaryotes may have co-opted riboswitches to control alternative splicing of genes.

Here we describe a new algorithm RNAbor (Freyhult, Moulton, Clote Bioinformatics 2007) which gives information on possible conformational switches by computing the Boltzmann probability of structural neighbors of a given RNA secondary structure. A secondary structure T of a given RNA sequence s is called a d -neighbor of S if T and S differ by exactly δ base pairs. RNAbor computes the number (N_{δ}), the Boltzmann partition function (Z^{δ}) and the minimum free energy (MFE^{δ}) and corresponding structure over the collection of all d -neighbors of S . This computation is done simultaneously for all $\delta \leq m$, in run time $O(m^2 n^3)$ and memory $O(mn^2)$, where n is the sequence length. We apply RNAbor for the detection of possible RNA conformational switches, and compare RNAbor with an existent switch detection method. We also provide examples of how RNAbor can at times improve the accuracy of secondary structure prediction

Biology by Design: Integrating Synthetic Biology and Systems Biology

J.J. Collins

*Center for BioDynamics and Dept. of Biomedical Engineering
Boston University*

Many fundamental cellular processes are governed by genetic programs which employ protein-DNA interactions in regulating function. Owing to recent technological advances, we describe how techniques from nonlinear dynamics and molecular biology are being utilized to model, design and construct synthetic gene regulatory networks. Importantly, engineered gene networks represent a first step towards logical cellular control, whereby biological processes can be manipulated or monitored at the DNA level. From the construction of a simple set of genetic building-block circuits (e.g., toggle switches, oscillators, etc.), one can imagine the design and construction of integrated biological circuits capable of performing increasingly elaborate functions. Moreover, synthetic gene circuits can be interfaced with the natural regulatory circuitry in cells to create, in effect, programmable cells. We discuss the implications of synthetic gene networks for biotechnology, biomedicine and biocomputing. In addition, we present integrated computational-experimental approaches that enable construction of first-order quantitative models of gene-protein regulatory networks using only steady-state expression measurements and no prior information on the network structure or function. We discuss how the reverse-engineered network models, coupled to experiments, can be used: (1) to gain insight into the regulatory role of individual genes and proteins in the network, (2) to identify the pathways and gene products targeted by pharmaceutical compounds, and (3) to identify the genetic mediators of different diseases.

Crystal Structure and phylogenetic data to infer the binding motifs for new protein

Francesca Cordero¹ and Gary D. Stormo²

¹ *Department of computer science , University of Torino, corso svizzera 185, 10149 Torino, Italy*

² *Department of Genetics , Washington University Medical School, St. Louis, MO 63110 USA*

DNA binding transcriptional regulators interpret the genome's regulatory code by binding to specific sequences to induce or repress gene expression. Accurate genome wide characterization of transcription factors binding sites is thus a necessary prerequisite to deciphering complex gene expression patterns. Positional weight matrices (PWMs) are simple motif representations that have been widely used in motif-identification algorithms. Usually, the search results using PWMs are a large number of hits many of which are false positives. This problem can arise because the data used to build a PWM are insufficient, as well as that the PWM model, which assumes independence between position in the pattern, is over simplified.

We have developed a method that infers the sequence specificity of a transcriptional factor (TF) from x-ray co-crystal structure of a protein from the same TF family bound to one particular DNA sequence and binding data from additional homologous TFs.

The principal idea is to find which positions in a protein sequence link to which positions in transcriptional factor binding site (TFBS). In a crystal structure these direct interactions are known, then we can obtain for every amino acids if it interacts or not with a nucleotide and if so, which nucleotide and its position in the TFBS.

Using BLAST we extract the homologue DNA binding domain (DBD). For every homologue TF we found the PWM that describes the TFBS. To align the different profiles we use a dynamic programming algorithm to identify high similarity regions using the ALLR statistic.

The crystal structure helps to find which amino acids interact with which position in DBD

If the amino acid changes we know which position in the DBD changes and we can infer which nucleotide interaction should also change. In this way we can predict the binding motifs for new proteins for which we don't have binding data..

Drug target prediction via Lasso regression analysis of mRNA expression compendia

Elissa J. Cosgrove^{1*}, Yingchun Zhou^{2*}, Eric D. Kolaczyk², Timothy S. Gardner¹

¹*Department of Biomedical Engineering and* ²*Department of Mathematics and Statistics, Boston University, Boston, MA 02215*

**These authors contributed equally.*

A major challenge in the development of new therapeutic drugs is the identification of the molecular targets of candidate compounds. In this study, we applied Lasso regression to a sparse simultaneous equation model (SSEM) of mRNA concentrations in order to construct a gene interaction network model from a compendium of microarray experiments. This inferred network was then used to filter data from a given experimental condition and predict its directly perturbed gene targets. We compared the performance of our algorithm to a compendium z-test model and a previously published network-based filtering approach, mode-of-action by network identification (MNI). We tested all three approaches on simulated data, two yeast microarray compendia, and one *E. coli* compendium, each composed of over 500 arrays and including deletion, mutant, overexpression, and drug perturbations. Our algorithm demonstrated an improvement in sensitivity in predicting the known target(s) of perturbations: a 30% improvement on simulated data and a 5% improvement on real data compared with the next best method. Additionally, the Lasso method is appealing as there are no unspecified algorithm parameters, unlike previously published approaches. The results of this study highlight the value of the network-based filtering approach in perturbation target identification, and provide a rigorous statistical foundation for further development of drug target prediction algorithms. Moreover, incorporation of proteomic and metabolomic data into this approach would likely enhance target prediction.

Zebrafish Gene Map and Microarray Annotations

Anthony D. DiBiase, Anhua Song, Fabrizio Ferre, Gerhard Weber, David Langenau, Leonard I. Zon, Yi Zhou

Hematology and Oncology and Stem Cell Program, Department of Medicine, Children's Hospital Boston and Harvard Medical School, 300 Longwood Ave., Boston, MA 02115

The zebrafish is a vertebrate genetic and development model system widely used to study vertebrate development and model human biology and disease. This research focuses on high-quality radiation hybrid (RH) mapping of genes and markers (18,898 to date) on the zebrafish T51 RH panel and on genome-wide, high-quality annotation of zebrafish genes (8,916 genes mapped to human orthologs to date) that are represented by Affymetrix Zebrafish Genome Array probesets. This high density RH map is being used to assist the assembly and finishing of the zebrafish genome sequencing. The established orthology between zebrafish and human genes are used to elucidate gene function and coordination and to efficiently analyze zebrafish gene expression profiles. These genome mapping and annotation efforts allow researchers to use more accurately assembled zebrafish genome sequences in research projects, quickly associate zebrafish genes with known biological functions from annotations of human known genes, and identify syntenic genomic regions across many species (e.g. teleosts). The release of this gene annotation to the research community has reduced redundant annotation effort and facilitates the use of zebrafish expression arrays for analyzing global gene expression profiles. We are continuing our annotation efforts with the aid of the international zebrafish research community as the quality and quantity of model organism sequences in repositories and their annotations improve.

Probing microRNA regulatory functions with microRNA 'sponge' inhibitors

Margaret S. Ebert, Madhu S. Kumar, Nancy Guillen, Phillip A. Sharp, Tyler Jacks, Douglas A. Lauffenburger.

Center for Cancer Research, Massachusetts Institute of Technology, E17-529B, Cambridge, MA 02139, USA.

MicroRNAs have emerged in recent years as a major class of gene regulatory molecules. They are computationally predicted to regulate thousands of mammalian genes and have been implicated in a variety of cellular processes, including those pertinent to cancer as well as normal development and tissue maintenance. Nonetheless, relatively few target predictions have been validated, and few microRNA loss-of-function phenotypes have been determined. As a means to study microRNA function in cell culture and in animal models, we developed 'microRNA sponges,' competitive inhibitors of mature microRNAs that mimic natural target mRNAs, contain multiple binding sites for a microRNA of interest, and are strongly and continuously expressed from transgenes within the cell. We are currently using sponges to probe the functions of several microRNA seed families in cancer cell lines and in mouse models of cancer. First, we have observed coordinate regulation of the G1 cyclins by the miR-16 family. How might these microRNAs promote mRNA degradation and repress translation to fine-tune the transcriptional and post-translational control of cyclin protein oscillation? Second, we are investigating a negative feedback loop linking the miR-30-5p family to the microRNA pathway component GW182. Is global microRNA processing and activity sensed and maintained through the repression of this Argonaute-interacting gene? Third, in collaboration with the Jacks lab, we are using retrovirally delivered microRNA sponges to test the effect of the let-7 family on Ras signaling and tumor growth in a xenograft model of lung cancer. Does let-7 function as a tumor suppressor by repressing the expression of oncogenic KRas and HMGA2? Fourth, in collaboration with the Lauffenburger lab, we are modulating the activity of several abundant microRNAs in hepatocellular carcinoma to test their influence on cell survival/apoptosis decisions and kinase signaling networks. Can we discover pathway interactions suggesting that specific inhibition of a microRNA could enhance the efficacy of cancer therapies?

Function and Evolution of Regions Bound by *Drosophila* Transcription Factors

Michael B. Eisen

*LBNL
Berkeley*

Identifying the genomic regions bound by sequence specific regulatory factors is central both to deciphering the complex DNA cis-regulatory code that controls transcription in metazoans and to determining the range of genes that shape animal morphogenesis. We have used whole-genome tiling arrays to map sequences bound in *Drosophila melanogaster* embryos by the six maternal and gap transcription factors that initiate anterior-posterior patterning. We find that these sequence-specific DNA binding proteins bind with quantitatively different specificities to an overlapping set of several thousand genomic regions in blastoderm embryos. The more highly-bound regions include all of the over forty well-characterized enhancers known to respond to these factors as well as several hundred putative new cis-regulatory modules clustered near developmental regulators and other genes with patterned expression at this stage of embryogenesis. In addition to these highly-bound regions, there are several thousand regions that are reproducibly bound at lower levels. However, these poorly-bound regions are, collectively, far more distant from genes transcribed in the blastoderm than highly-bound regions and are preferentially found in protein-coding sequences. We have extensively analyzed the evolution of recognition sequences in these bound regions and find little evidence for their preferential conservation. Together these observations suggest that many of these poorly-bound regions are not involved in early-embryonic transcriptional regulation and may be non-functional. I will propose that the pervasive view amongst both experimental and computational biologists that most protein-DNA interactions observed *in vivo* are functional is wrong.

Bayesian Meta-analysis for Identification of Cell cycle-regulated Genes in Fission Yeast

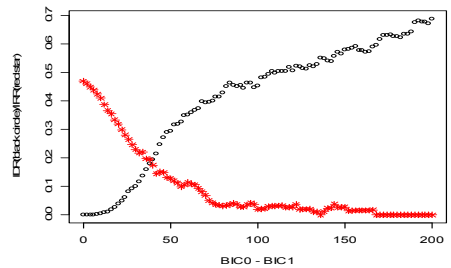
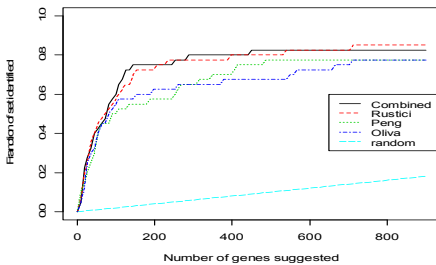
Xiaodan Fan¹, Saumyadipta Pyne², and Jun S. Liu¹

¹Statistics Dept., Harvard University, 1 Oxford Street, Cambridge, MA, USA

²Broad Institute of MIT and Harvard, 7 Cambridge Center, Cambridge, MA, USA

The effort of identifying cell cycle-related gene from gene expression data has lasted for a decade. Recently three independent research groups (Rustici *et al.*, 2004; Peng *et al.*, 2005; Oliva *et al.*, 2005) have done both elutriation and *cdc25* experiments to measure all genes' expression during the fission yeast cell cycle. Their results showed discrepancies with regard to the identity and numbers of periodically expressed genes. In this paper we introduced a hierarchical Bayesian model to integrate multiple time series experiments. A Monte Carlo Markov Chain (MCMC) approach was adapted to draw samples from the posterior distribution of the all parameters. We calculated the Bayesian Information Criterion (BIC; Schwarz, 1978) values using the Maximum A Posteriori estimates. The difference of BIC values from two competing models were then used to classify genes as either cell cycle-related gene or non-cell-cycle gene. The above approach was applied on both the combined dataset and the datasets from individual groups. The integration-driven discovery rate (Choi *et al.*, 2003) was calculated to show the gain in information from combining studies versus individual studies. The integration-driven revision rate (Stevens and Doerge, 2005) was calculated to show the amount of genes that are missed or "dropped" in meta-analysis versus separate study analysis.

Our model considered both the de-synchronization effect and the bloc-release effect. A damping sinusoidal function with quadratic trend and i.i.d. Gaussian noise was used to model each time series curve. We designed a stepwise conditional MCMC approach to explore the high-dimensional parameter space. The results showed the modified MCMC chains converged fast and our model fitted the data well. The classification performance on the benchmark set of cell cycle-regulated genes (Marguerat *et al.*, 2006) and the corresponding IDR/IRR curves are showed below. They showed the power of our Bayesian approach to pool multiple datasets is appealing.



Regulatory Hubs: Classification and Conservation

Andrew Fox, Donna K. Slonim

*Department of Computer Science, Tufts University,
161 College Ave, Medford 02155, MA, USA*

The power of comparative genomic research has grown steadily with the availability of genomic sequence and annotation for increasing numbers of organisms. A significant amount of this power is derived from and hence dependent on the fundamental assumption that relationships between key proteins are conserved between humans and model organisms. Previous studies of transcriptional regulatory networks have shown that despite the preferential conservation of functionally important proteins, regulatory network modules are poorly conserved, while individual interactions are conserved only to a limited degree.

In this study we investigate the different question of whether the high interaction degree of individual proteins is conserved with orthology. Our analysis reveals that an interaction hub in one species has an increased chance of having an ortholog of high degree in one or more other species, showing that interaction hubs are preferentially conserved. However, the effect is subtle and appears to involve only a subset of hubs.

We then investigate the nature of the genes in this subset. Incorporating gene expression and protein interaction data, we introduce the notion of classifying hubs by the variance in the number of their interactors simultaneously expressed whenever an incumbent hub is expressed. In conjunction with the traditional measure of Pearson correlation between hub and interactor expression, this variance-correlation space reveals a cluster of known regulatory hubs, including many transcription factors. This clustering allows us to identify a novel set of candidate regulatory hubs. We demonstrate that the set of hub proteins chosen in this way is significantly enriched for several important regulatory functions, including chromatin remodeling, and that these proteins show significantly higher conservation of their interaction degree.

This work provides evidence that identifying high-degree proteins relevant to particular biological processes in model organisms is more likely to shed light on human response. Our results have important implications for predicting novel protein function, identifying protein interactions, and finding key targets for therapeutic intervention.

Dynamics of replication-independent histone turnover in budding yeast.

Nir Friedman

Hebrew University, Israel

Chromatin plays roles in processes governed by different time scales. To assay the dynamic behavior of chromatin in living cells, we used genomic tiling arrays to measure histone H3 turnover in G1-arrested *Saccharomyces cerevisiae* at single-nucleosome resolution over 4% of the genome, and at lower (approximately 265 base pair) resolution over the entire genome. We find that nucleosomes at promoters are replaced more rapidly than at coding regions and that replacement rates over coding regions correlate with polymerase density. In addition, rapid histone turnover is found at known chromatin boundary elements. These results suggest that rapid histone turnover serves to functionally separate chromatin domains and prevent spread of histone states.

High-resolution mapping and characterization of open chromatin across the human genome

Terrence S. Furey¹, Alan P. Boyle¹, Sean Davis², Hennady P. Shulha³, Paul Meltzer², Zhiping Weng³, Gregory E. Crawford¹

¹*Institute for Genome Sciences and Policy, Duke University, Durham, NC 27708.*

²*Center for Cancer Research, National Cancer Institute, National Institutes of Health, Bethesda, MD 20892.*

³*Biomedical Engineering Department, Boston University, Boston, MA 02215.*

Identification of open sites in chromatin allows for identification of active genetic elements across the genome. These sites can be measured through their sensitivity to being cleaved by DNaseI. For the last 28 years, these DNaseI hypersensitive sites have been the gold standard method to identify the location of many types of regulatory elements, including promoters, enhancers, silencers, insulators, and locus control regions. However, because of the difficulty in accessing these sites on a genomic scale, there has been no whole genome mapping of DNaseI hypersensitive sites. Here, we present the first-ever DNaseI HS site map from across the entire human genome using both next-generation high throughput sequencing (Illumina and 454) and whole genome tiled microarrays (NimbleGen). Both methods by themselves have high sensitivity and specificity, but combining data from both platforms generates an even higher quality dataset. In total we find 94,797 DNaseI HS sites covering 60Mb (2.1%) of the genome. We find that all DNaseI HS sites are not equal: DNaseI HS sites at promoters of genes appear to be much more open than non-promoter HS sites on average. In addition, we find transcription factor motifs that are enriched in DNaseI HS sites provide important clues to factors that control gene expression in the cell type that we studied. The accurate and comprehensive DNaseI HS map presented here offers an unprecedented view of open chromatin structure at high resolution. The further generation of genome-wide DNaseI HS maps from a diverse set of normal and diseased human cell types, as well as from those from other species, will continue to reveal how chromatin structure differences contribute to cell-type specific gene expression and cell fate decisions. These maps will also provide a scaffold on which to combine and analyze data from ChIP-chip/ChIP-Seq experiments, gene expression data, comparative sequences from other species, and data from the HapMap consortium to further resolve the complex nature of gene regulation.

Pattern-based analysis of interactome networks in *C. elegans*

Sira Sriswasdi¹, Kesheng Liu¹, Thomas Martinez¹, Carmel Mercado¹, Pengyu Hong², and Hui Ge¹

Whitehead Institute for Biomedical Research, 9 Cambridge Center, Cambridge MA 02142;
²Department of Computer Sciences, Volen Center for Complex Systems 135 MS 018,
Brandeis University, Waltham MA 02454.

In the last few years, various interaction datasets have been generated for human and model organisms using high-throughput experimental approaches such as yeast twohybrid analysis, biochemical pull-down assays, or synthetic phenotypic screens. These large datasets can be represented as interactome networks, in which genes and proteins are depicted as nodes and interactions as edges. Given the complexity of these networks and the error-prone nature of high-throughput assays, it is not yet known how these networks can help us gain new knowledge about biological pathways at a high confidence. To address this need, we develop a computational approach that predicts undiscovered interactions and novel components of biological pathways based on interactome networks. Our prediction method is based on the concept of network motifs, which are recurrent topological patterns enriched in networks of interest as compared to randomly generated networks. Because of the high recurrence, network motifs have been implicated as basic building blocks of biological networks. From a *C. elegans* interactome network that combines protein-protein interactions and genetic interactions, we identify “complete network motifs”, which are cliques enriched in the interactome network. We also identify “incomplete network motifs”, which are patterns needing one or few edges to become complete network motifs. We hypothesize that incomplete motifs are part of complete motifs, and that the missing edges represent unidentified interactions because of the incomprehensive coverage of currently available datasets. We predict novel genetic interactions based on incomplete network motifs, and these predictions achieve a sensitivity of 60% and a specificity of 99.6% estimated by crossvalidation. The success of the predictions provides supporting evidence that motifs in interactome networks are biologically relevant. We then investigate whether such pattern-based network analysis can provide biological insights into specific pathways. We focus on *dbl-1* induced TGF- β pathway, which regulates body size in *C. elegans*. We identify all the complete and incomplete motifs containing any known components of the TGF- β pathway. We predict genes that participate in the network motifs as putative components of the pathway. In order to verify our predictions and to further characterize the interactions between the predicted genes and the known components, we plan to examine by RNAi analysis whether the predicted genes regulate body size through the TGF- β pathway. We will knock down the expression of predicted genes from mutants of known TGF- β pathway components as well as from wild-type animals, and determine whether the RNAi experiments result in any changes in body size. To this end, we have developed an image-processing program which automatically measures *C. elegans* body size from digital images. Results from our RNAi analysis should reveal epistatic relationships between the predicted genes and known pathway components and reveal whether their interactions are synergistic or antagonistic. Taken together, our approach greatly facilitates the reconstruction and knowledge expansion of biological pathways based on interactome networks.

A fast and simple motif finder utilizing rank lists of sequences

Stoyan Georgiev, Karthik Jayasurya, Sayan Mukherjee, Uwe Ohler

Institute for Genome Sciences & Policy, Duke University, Durham NC 27710, USA

We propose a fast enumerative strategy for identifying cis-regulatory elements involved in regulatory processes in the cell. In addition to a set of non-coding regulatory regions, our approach makes use of a genome-wide condition-specific ranking of genes, such as p-values resulting from chromatin-immunoprecipitation (ChIP-chip) experiments, or expression changes from microarray gene knockdown data sets. Our strategy defines gene sets by the presence of motifs in the regulatory regions, and exhaustively identifies motifs/gene sets with strong enrichments towards the top or bottom of the ranked list. Rather than explicitly modeling the large space of intergenic sequences, we execute a greedy search on the space of potential IUPAC motifs, utilizing the highly time-efficient suffix-tree data structure. We utilize the complete information in the ranked list and avoid pre-set thresholds by defining an appropriate sample statistic to guide the search.

We applied our method to the specific example scenario of integrating ChIP-chip data for motif finding in eukaryotic promoters. The resulting algorithm was validated extensively by simulation, and application on available curated yeast and human ChIP-chip datasets delivered a performance similar to or better than existing popular motif finders. Our approach provides a new look at an old problem and promises extension possibilities to a wide range of applications.

Genome-wide identification of active regulatory elements in the human genome by FAIRE

Paul Giresi¹, Vishy Iyer², and Jason Lieb¹

¹*Department of Biology and Carolina Center for Genome Sciences, University of North Carolina at Chapel Hill, 407 Fordham Hall, NC, USA*

²*Institute for Cellular and Molecular Biology, and Center for Systems and Synthetic Biology, University of Texas at Austin, 1 University Station, TX, USA*

FAIRE (Formaldehyde-Assisted Isolation of Regulatory Elements) is a simple procedure for the genome-wide isolation and identification of nucleosome-depleted regions in eukaryotic cells. Identification of "open" chromatin regions has been one of the most accurate and robust methods to identify functional promoters, enhancers, silencers, insulators, and locus control regions in mammalian cells. Here I will present an analysis of FAIRE data derived from human fibroblasts, HeLa S3 cells, and cell lines derived from two subtypes of breast cancer. In addition, I will present the initial analysis of whole-genome data derived from a 21-million probe tiling array.

Position-dependent motif characterization using non-negative matrix factorization

Lucie N. Hutchins¹, Erik L. McCarthy², Sean Murphy¹, Priyam Singh³, and Joel H. Graber¹⁻³

¹*Center for Genome Dynamics, The Jackson Laboratory, 600 Main St, Bar Harbor, ME 04609 USA*

²*Functional Genomics Program, University of Maine, Orono, MA 044 USA*

³*Bioinformatics Program, Boston University, Boston, MA 02215 USA*

Cis-acting regulatory elements are frequently constrained by both sequence content and positioning relative to a functional site, such as a splice or polyadenylation site. We describe a novel approach to regulatory motif analysis based on non-negative matrix factorization (NMF). Whereas existing pattern recognition algorithms commonly focus primarily on sequence content, our method simultaneously characterizes both positioning and sequence content of putative motifs.

We align training sequences on their common functional site and count occurrence and position of all sequence words of a given length (k -mers). The resulting positional word count (PWC) matrix is indexed in two dimensions by k -mer and relative position. Further analysis is based on NMF, a dimension reduction algorithm that models the PWC matrix as a linear superposition of a small number of underlying patterns. We demonstrate that NMF decomposes the PWC matrix into two new matrixes with intuitive interpretations, reflecting motif sequence content and positioning, respectively.

Tests on artificially generated sequences show that NMF can faithfully reproduce both positioning and content of the test motifs. Our analysis has proven to be particularly successful in distinguishing multiple elements with significant overlap in sequence content and/or positioning. NMF also identifies motifs that are missed by existing approaches. We present examples of analysis of a variety of data sets, including mRNA 3'-processing (cleavage and polyadenylation) sites, splice junctions, and transcription start sites.

Characterization of TBP binding to TATA boxes and its relationship to structural properties of TATA boxes

Hana Faiger, Marina Ivanchenko, Ilana Cohen, and Tali E. Haran

Department of Biology, Technion – Israel Institute of Technology, Haifa, Israel

TBP recognizes its target sites, TATA boxes, by recognizing their sequence-dependent structure and flexibility ("indirect readout"). Studying this mode of TATA-box recognition is important for elucidating the binding mechanism in this system, as well as for developing methods to locate new binding sites in genomic DNA.

We have probed systematically the role of the sequences flanking the MLP and E4 TATA boxes, as prototype TATA boxes having A_4 versus $(T-A)_4$ tracts, respectively, on TBP/TATA-box interactions using *in vitro* evolution methods. We searched for sequences with optimal binding affinities, as well as for those with optimal binding stabilities. We show that for TATA boxes containing alternating $(T-A)_n$ runs binding affinity and stability to TBP are significantly dependent on the nature of the sequences flanking the core TATA box, to an extent comparable to that observed when the sequences within the core TATA box itself are changed. We suggest that this is a novel form of indirect readout, and propose that structural modulation of certain TATA-boxes by their flanking sequences increases the number of different sequences that constitute a valid TATA box. This enhances the fine-tuning of gene regulation attainable at the initial stage of transcription.

We recently studied the binding properties of TBP to all consensus-like TATA boxes, as well as TBP-induced TATA-box bending. Classifying TATA boxes by their structural properties clarifies the different recognition pathways and binding mechanisms used by TBP upon binding to different TATA boxes, and reveal differences in the indirect readout of these TATA boxes by TBP. Furthermore, we show that various non-additive effects exist in TATA boxes dependent on their structural properties. Statistical analysis indicated that TATA boxes that have a context-independent cooperative structure are best described by a nearest-neighbor non-additive model, whereas TATA boxes that have a flexible context-dependent conformation cannot be described by either an additive model or by a nearest-neighbor non-additive model. We will discuss the structural and evolutionary sources of the difficulties in predicting new binding sites by probabilistic weight-matrix methods for proteins in which indirect readout is dominant.

Association Analysis of Gene Functions and Upstream Sequence Conservation in the Human Genome

Weichun Huang¹, Leping Li², Gabor Marth¹, and Bruce S Weir³

¹Department of Biology, Boston College, Chestnut Hill, MA, USA, ²Biostatistics Branch, the National Institute of Environmental Health Sciences, NIH, NC, USA, ³Department of Biostatistics, the University of Washington, Seattle, WA, USA

Cross-species sequence comparisons have shown that non-coding sequences, particularly, Conserved Non-coding Sequences (CNS), play functional roles as important as protein-coding sequences for an organism. It remains unclear, however, how sequence conservation of gene upstream sequences is related to the specific functions of genes in the human genome. In this study, with the aim of elucidating some features of human gene regulatory networks, we systematically investigated the sequence conservation of gene upstream regions and its relationships with gene functions using about 6, 000 and 1, 000 pairs of human-mouse and human-rat orthologous genes, respectively. We found that the genes with highly conserved upstream regions are significantly associated with important regulation functions such as transcription regulation, organ development control and morphogenesis. In particular, the developmental process related transcription regulators, *e.g.*, *HOX* and *POU* genes, shows extreme upstream region conservation. In contrast, the genes involved in physiological process, or encoding the components of protein complex or catalytic enzymes are significantly associated with more divergent upstream regions. Furthermore, we found that the degree of upstream region conservation is highly correlated with the density of transcription factor binding sites in the upstream regions. Our results suggest that the key developmental process-related transcription regulators are under the sophisticated control, both temporal and spatial, of their upstream regulatory regions, which are subjected to stringent evolutionary constraints, hence are highly conserved in both human and rodent species.

Protein Network Comparative Genomics

Trey Ideker

University of California San Diego

With the appearance of large networks of protein-protein and protein-DNA interactions as a new type of biological measurement, methods are needed for constructing cellular pathway models using interaction data as the central framework. The key idea is that, by comparing the molecular interaction network with other biological data sets, it will be possible to organize the network into modules representing the repertoire of distinct functional processes in the cell. Three distinct types of network comparisons will be discussed, including those to identify:

- (1) Protein interaction networks that are conserved across species
- (2) Networks in control of gene expression changes
- (3) Networks correlating with systematic phenotypes and synthetic lethals

Using these computational modeling and query tools, we are constructing network models to explain the physiological response of yeast to DNA damaging agents.

Relevant articles and links:

Yeang, C.H., Mak, H.C., McCuine, S., Workman, C., Jaakkola, T., and Ideker, T. Validation and refinement of gene regulatory pathways on a network of physical interactions. *Genome Biology* **6(7)**: R62 (2005).

Kelley, R. and Ideker, T. Systematic interpretation of genetic interactions using protein networks. *Nature Biotechnology* **23(5)**:561-566 (2005).

Sharan, R., Suthram, S., Kelley, R. M., Kuhn, T., McCuine, S., Uetz, P., Sittler, T., Karp, R. M., and Ideker, T. Conserved patterns of protein interaction in multiple species. *Proc Natl Acad Sci U S A.* **8:102(6)**: 1974-79 (2005).

Suthram, S., Sittler, T., and Ideker, T. The *Plasmodium* network diverges from those of other species. *Nature* **437**: (November 3, 2005).

<http://www.pathblast.org>

<http://www.cytoscape.org>

Acknowledgements: We gratefully acknowledge funding through NIH/NIGMS grant GM070743-01; NSF grant CCF-0425926; Unilever, PLC, and the Packard Foundation.

Systematic identification of mammalian regulatory motifs' target genes and functions

Jason B. Warner^{1,6,7}, Anthony A. Philippakis^{1,3,4,6}, Savina A. Jaeger^{1,6}, Fangxue Sherry He^{6,8}, Jolinta Lin^{5,6}, and Martha L. Bulyk^{2,3,4,6}

¹ *These authors contributed equally to this work;*

²*Department of Pathology;* ³*Harvard/MIT Division of Health Sciences and Technology (HST); Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115.*

⁴*Harvard University Graduate Biophysics Program, Harvard Medical School, Boston, MA 02115.*

⁵*Department of Biology, MIT, Cambridge, MA 02139.*

⁶*Division of Genetics, Department of Medicine;* ⁷*Current address: Applied Biosystems, Beverly, MA 01915.*

⁸*Current address: The Institute for Genomic Research, Rockville, MD 20850.*

We have developed an algorithm (“Lever”) that systematically maps metazoan DNA regulatory motifs or motif combinations to the sets of genes that they likely regulate. Lever accomplishes this by assessing whether the motifs are enriched within predicted *cis*-regulatory modules (CRMs) in the noncoding sequences surrounding genes in a collection of gene sets. When these gene sets correspond to Gene Ontology (GO) categories, the results of Lever analysis allow the unbiased assignment of functional annotations to the regulatory motifs and also to the candidate CRMs that comprise the genomic motif occurrences. We demonstrate these methods using human myogenic differentiation as a model system, for which we statistically assessed greater than 25,000 pairings of gene sets and motifs / motif combinations. These results allowed us to assign functional annotations to candidate regulatory motifs predicted previously, and to identify gene sets that are likely to be co-regulated via shared regulatory motifs. Lever represents a major genome-wide step in moving beyond the identification of putative regulatory motifs in mammalian genomes, towards understanding their biological roles. This approach is general and can be applied readily to any cell type, gene expression pattern, or organism of interest.

Cycles in the transcription regulatory network of *S. Cerevisiae*

Jieun Jeong and Piotr Berman

CSE Department, The Pennsylvania State University

Unlike the transcription regulatory network of *E. Coli*, the network of *S. Cerevisiae* (baker's yeast) contains many cycles. In the data set of Luscombe, Babu et al. We have found 51 interactions in cycles (after removing self-loops), of which 49 were in a single node set after combining node sets of overlapping cycles. We use the term LSCC (Large Strongly Connected Component) to refer this set. We observed the phenomenon of LSCC containing almost all cycles also after random rewiring. The size of LSCC and its in-component is significantly smaller than the average after random rewiring (p-value ca. 0.001). We investigated the properties of LSCC, in particular how it changes in different physiological conditions and how we can identify meaningful substructures. Three observations can be mentioned as most interesting.

LSCC is very strongly related with cell-cycle, i.e. the cycle of changes occurring in a cell as it undergoes division. This partially explains the contrast with Prokaryotic network of *E. Coli* that has no cycles or very few (dependent on the data set). The nature of the overlap between LSCC and cell cycle interactions is also interesting: 27 out of 30 of these interactions are active in all five phases of cell cycle, while only 64 of 124 interactions of cell cycle that are directed at other transcription factors are common to all five phases (and none of 135 interactions to terminal targets). Thus this cycle forms a mechanism that remains in place throughout all five phases while its function is modulated by other interactions. The second observation is that besides cell cycle, stress response (heat shock response) also has its "own" overlapping cycles, three cycles within three transcription factors.

Moreover, overlapping cycles of length 3 or 4 form a part of cell cycle interactions that are also active during sporulation. In general, members of short cycles (2 to 4 interactions) are very strongly correlated whether they are active in some conditions or not. LSCC has a very orderly layout (with few line intersections) with strongly connected subsets active together in the same conditions. The third observation is that all chains of interactions with more than two interactions proceed through LSCC which forms kind of a switchboard for the entire network. Thus we propose a variation of hierarchy of transcription factors: in-component of LSCC, LSCC, and the out-component, as well as the "egalitarian" set of ca. 1/4 of transcription factors that belong to very short chains only. Ca. 10% of the genes in the in-component and in LSCC are cancer related, while only ca. 2.5% of genes in the rest of the network are.

New Classes of RNAs Potentially Involved in Regulation of Gene Expression

Philipp Kapranov¹, Juergen Schmitz², Katalin Fejes-Toth³, Radharani Duttgupta¹, Aarron T. Willingham¹, Timofey Rozhdestvensky², Jill Cheng¹, Richard Reinhardt⁴, Sujit Dike¹, Greg Hanon³, Juergen Brosius², Thomas R. Gingeras¹

¹ Affymetrix, Inc. Affymetrix Laboratory, Santa Clara, CA, 95051

² University of Muenster, Muenster, Germany

³ Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724

⁴ Max-Planck-Institute for Molecular Genetics, Berlin, Germany

Unbiased whole-genome profiling of different types of human RNA has revealed presence of pervasive un-annotated transcription in the human genome and several novel classes of RNAs. Two of these classes are represented by short (<200nt) and long (> 200 nt) RNA species originated at the promoters and CpG islands named PASRs and PALRs (promoter-associated short/long RNAs). Several lines of evidence suggest that PASRs may be correlated with expression of the downstream genes. (1) Presence of PASRs is enriched over what would be expected by random chance at and around transcriptional start sites of genes on both strands. (2) The Density of PASRs correlates with steady-state levels of the mRNA from the associated genes. (3) PASRs are syntenically conserved between human and mouse.

Sequence analysis of PASRs combined with Northern blot analysis reveals several characteristics of these small RNAs, suggesting that they represent discrete stable, short RNAs in human and mouse cells of variable lengths. Comparative mapping of PASRs and PALRs suggests a possibility that PASRs could be made from the longer promoter-associated RNA species. Taken together these data suggest that short, stable RNAs are produced at the promoters of many mammalian genes and these RNAs could participate in regulation of gene expression (see abstract by R. Duttgupta and A. Willingham).

Systematic Discovery and Characterization of Fly microRNAs using 12 *Drosophila* Genomes

Pouya Kheradpour*, Alexander Stark*, Leopold Parts, Emily Hodges, Julius Brennecke, Gregory Hannon and Manolis Kellis

Computer Science and Artificial Intelligence Laboratory, MIT; Broad Institute of MIT and Harvard; Cold Spring Harbor Laboratory

* = contributed equally

MicroRNAs (miRNAs) are short RNA genes that direct the inhibition of target messenger-RNA expression via complementary binding sites in the targets' 3' untranslated region (UTR). Currently, miRNAs are estimated to comprise 1%–5% of animal genes, making them one of the most abundant classes of regulators. In addition, an average miRNA may regulate hundreds of genes, so that a large fraction of all genes are miRNA targets. It is thus desirable to obtain a comprehensive characterization on all miRNAs in an animal genome, especially as knowledge of the sequence alone can allow the identification of the physiologically relevant target genes.

We have used 12 *Drosophila* genomes to define structural and evolutionary signatures of miRNA hairpins, which we use for their de novo discovery. We predict more than 41 novel miRNAs, which encompass many unique families, and 28 of which we validate experimentally. We also define precise signals for the start position of mature miRNAs, which we use to correct the annotation of previously known miRNAs, often leading to drastic changes in their target spectrum. We show that miRNA discovery power scales with the number and divergence of species compared, suggesting that such approaches can be successful in human as dozens of mammalian genomes become available.

Interestingly, for some miRNAs sense and anti-sense hairpins score highly, and indeed we find mature miRNAs from both strands *in vivo*. Similarly, we find that when prediction of miRNA start is imprecise, multiple starts are indeed found to be processed *in vivo*, especially for miRNAs with few targets. Lastly, for several miRNAs, both arms of the hairpin (mature and star arms) score highly computationally, and are found to be more abundantly expressed *in vivo*. These results suggest that a single miRNA locus can produce several functional miRNA products, each with distinct functional targets. For miR-10 in particular, both arms show abundant processing, and both show highly conserved target sites in Hox genes, suggesting a possible cooperation of the two arms, and their role as a master Hox regulator.

A Polycomb Switch During Muscle Differentiation

Roshan M. Kumar¹, Stuart S. Levine¹, and Richard A. Young^{1,2,†}

¹ Whitehead Institute for Biomedical Research, 9 Cambridge Center, Cambridge, Massachusetts 02142, USA

² Department of Biology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

The polycomb group (PcG) proteins help maintain repression of key developmental regulatory genes in embryonic stem cells. How these epigenetic regulators contribute to new gene expression programs as cells differentiate is poorly understood. We report here that during myogenesis, the PcG proteins Suz12 and Bmi1 are lost from the developmental genes they occupied in muscle precursor cells (myoblasts) and appear at proliferation and metabolism genes in differentiated muscle cells (myotubes). Developmental genes vacated by PcGs retained the histone H3K27me3 modification catalyzed by PRC2 and remained repressed in myotubes, indicating that continued PcG occupancy is not necessary for maintenance of the repressed state of these genes in post-mitotic cells. Consistent with the gene occupancy data, knockdown of Suz12 in myoblasts led to the expression of markers of alternate cell fates, while differentiation of Suz12 knockdown cells led to continued expression of proliferation genes. These results reveal that PcG proteins switch from a role in repressing alternate cell fates in muscle precursors to a novel non-epigenetic function in suppressing proliferation in terminally differentiated cells.

A predictive model of the oxygen and heme regulatory network in yeast

Anshul Kundaje^a, Xiantong Xin^b, Changgui Lan^b, Steve Lianoglou^a, Mei Zhou^b, Christina Leslie^{*c}, Li Zhang^{*b},

^aDepartment of Computer Science, ^bDepartment of Environmental Health Sciences Columbia University, New York, NY, ^cComputational Biology Program, Sloan-Kettering Institute Memorial Sloan-Kettering Cancer Center, New York, NY

**To whom correspondence should be addressed: lz2115@columbia.edu or*

cleslie@cbio.mskcc.org

Mechanisms of oxygen sensing and regulation underlie many physiological and pathological processes. Due to the complexity of these mechanisms, standard biochemical and genetic experiments and conventional analysis of microarray experiments have led only to a limited understanding of the global oxygen regulatory network.

Here we present the first genomic and computational analysis to decipher the global network of oxygen and heme regulation in yeast. We identified all oxygen-regulated, heme-regulated, and Co₂⁺-inducible genes, by using microarray gene expression profiling. We then used a machine learning algorithm called MEDUSA to discover the regulatory programs mediating oxygen and heme regulation. MEDUSA integrates mRNA expression, promoter sequence, and ChIP-chip occupancy data to learn a model that accurately predicts the differential expression of target genes in held-out data. MEDUSA uses boosting, a technique from machine learning, to help avoid overfitting as the algorithm searches through the high dimensional space of potential regulators and sequence motifs.

We then used a novel margin-based score to extract significant condition-specific regulators and assemble a global map of the oxygen sensing and regulatory network. This network is constituted of approximately 60 transcriptional regulators and signal transducers, including several previously known oxygen and heme regulators, such as Hap1, Mga2, Hap4, and Upc2, as well as new predicted regulators. Many DNA motifs identified by MEDUSA are consistent with previous experimentally identified motifs.

Finally, we performed targeted biochemical validation of our model through experimental analysis of the promoter activity of OLE1, a gene that is strongly induced by hypoxia. MEDUSA identifies 5 potential regulators of the OLE1 promoter, 4 of which are predicted to upregulate OLE1 and the remaining one to downregulate OLE1 under hypoxia; only one of these regulators has been previously identified as an oxygen regulator. We measured OLE1 promoter-lacZ reporter activity under air and anoxia in wild type and mutant cells with each one of the regulator genes deleted. For all 4 predicted positive regulators, the reporter activity in anoxic mutant cells was reduced compared to that in wild type cells, consistent with the model's predictions. Also consistent with the model, deletion of the predicted negative regulator did not affect the reporter activity in anoxic cells. These experimental results support the predictive power of the MEDUSA algorithm and demonstrate its ability to generate testable and specific biological hypotheses.

Dissecting transcriptional regulation in human B lymphocytes using an integrative, context-specific network

Celine Lefebvre¹, Wei Keat Lim¹, Kai Wang¹, Pavel Sumazin¹, Katia Basso², Riccardo Dalla-Favera² and Andrea Califano¹

¹Center for Computational Biology and Bioinformatics, ²Institute for Cancer Genetics, Columbia University, 1130 St Nicholas Avenue, New York, NY 10032, USA

The availability of a genome-wide map of regulatory interactions in a human cell is a critical step towards the elucidation of normal and disease-related phenotypes. In this work we used the B Cell Interactome (BCI), an integrated network of protein-DNA (PD) and protein-protein (PP) interactions produced by our lab [1] to dissect the regulatory mechanisms controlling the Germinal Center (GC) reaction of antigen-activated B lymphocytes. The original BCI was extended using a new class of biochemically-validated, functional PP interactions predicted by the MINDY algorithm [2]. These identify proteins, such as kinases and co-factors, that post-translationally modulate the ability of a transcription factor (TF) to regulate its targets. Physical PD and PP interactions were predicted using a Naïve Bayes classification from a variety of clues, including among many others (a) gene co-expression data, (b) inferences from reverse engineering algorithms, such as ARACNE [3] and MINDY, and (c) evidence from interactions in other organisms. The context specificity of the network is supported by a large collection of 254 gene expression profiles representing 27 distinct normal and neoplastic human B cell phenotypes [3] that fully mediate the integration process.

The BCI constitutes the first integrated PD and PP interaction network for a specific human cell and provides a novel and unique tool to study regulatory processes in B cells. First, we analyzed the network for the enrichment of three-gene-motifs, i.e., fully-connected three-gene cliques typed according to the gene and interaction type. These underlie different regulatory motifs, such as the combinatorial transcriptional regulation of a target by two transcription factors or proteins participating in the same complex or pathway. Then, we combined such enriched motifs into larger network modules and studied their differential regulation within specific phenotypes. Specifically, we analyzed these larger modules to reveal activated and repressed pathways in the GC reaction, leading to the identification of 37 TFs that may play an important role in the GC formation and in its dysregulation during lymphomagenesis.

1. Lefebvre, C., et al., *A context-specific network of protein-DNA and protein-protein interactions reveals new regulatory motifs in human B cells*. Lecture Notes in Bioinformatics, 2007. **4532**: p. 42-56.
2. Wang, K., et al., *Genome-wide discovery of modulators of transcriptional interactions in human B lymphocytes*. Lecture Notes in Computer Science, 2006. **3909**: p. 348-362.
3. Basso, K., et al., *Reverse engineering of regulatory networks in human B cells*. Nat Genet, 2005. **37**(4): p. 382-90.

Whole-genome analysis of dorsal-ventral patterning thresholds in the *Drosophila* embryo

Julia Zeitlinger, Rob Zinzen, Dmitri Papatsenko, Joung-Woo Hong, Manolis Kellis, Alexander Stark, Rick Young, and Mike Levine

MIT, Whitehead Institute; Broad Institute; UC Berkeley, Dept. MCB, Center for Integrative Genomics

Dorsal is a sequence-specific transcription factor related to NF- κ B. The protein is distributed in a broad nuclear gradient in the precellular *Drosophila* embryo. This gradient controls dorsal-ventral patterning by regulating at least 50 target genes in a concentration-dependent manner. Dorsal works with two additional regulatory proteins that are encoded by genes directly regulated by the gradient, Twist and Snail. To determine how the Dorsal gradient generates diverse thresholds of gene activity, we have used ChIP-chip assays with Dorsal, Twist, and Snail antibodies. This method efficiently identified 20 known enhancers and predicted another ~50 novel enhancers associated with known or suspected dorsal-ventral patterning genes.

The analysis of ~30 different Dorsal target enhancers suggests that those mediating gene expression in response to high levels of the Dorsal gradient contain a series of disordered low-affinity Dorsal and/or Twist activator binding sites. In contrast, enhancers mediating expression in response to low levels of the gradient (the “type 2” response) contain an ordered arrangement of optimal Dorsal and Twist binding sites. This organization is highly conserved in evolution and is likely to foster cooperative occupancy of linked operator sites. We suggest that type 2 enhancers have some of the properties of a transistor: they amplify weak and unstable signals to produce a constant output of gene activity. The accurate representation of genetic circuit diagrams depends on the designation of this special subset of enhancers.

The conventional view of gene activation is that it depends on the recruitment of RNA polymerase II (Pol II) to the core promoter. However, there are a few documented examples of gene activation via Pol II elongation, whereby Pol II is bound, but paused at a fixed site just downstream of the transcription start site. It is not known to what extent this mechanism is used to establish differential patterns of gene expression in development. To investigate this issue, we performed ChIP-chip assays using antibodies directed against Pol II. At least 1,000 genes contain a paused or stalled form of Pol II prior to their activation in the *Drosophila* embryo. These include many developmental control genes, such as Hox genes and components of the FGF, Wnt, Hedgehog, TGF β , and Notch signaling pathways. Altogether, these results suggest that the regulation of Pol II elongation, not recruitment or initiation, might be a critical mechanism for gene activation in development.

An integrated approach to clustering expression trajectories in differentiating stem cells

Shenggang Li¹ Miguel Andrade² David Sanko³

¹*Dept. of Mathematics and Statistics, University of Ottawa. sli365@science.uottawa.ca*

²*Ottawa Health Research Institute. mandrade@ohri.ca*

³*Dept. of Mathematics and Statistics, University of Ottawa. sankoff@uottawa.ca*

During the differentiation of embryonic stem cells (ESC), the expression of some key genes is gradually or relatively suddenly shut down, while other genes show an initiation or increase in expression. Based on the trajectories over the first few hours of several hundred genes extracted from large arrays on the basis of significant level of expression at one or more sampling times, we address the problem of clustering time course gene expression data on differentiating ESC.

We treat a gene profile along a time course as a function composed of the following components (1) the impact of related genes or proteins in a dynamical system described by differential equations, (2) autoregressive factors representing auto correlation structure and feedback or loop effects, and (3) random components allowing for heterogeneity of individual genes in each cluster. The proposed model-based clustering method allows for more flexible characteristics by assigning both dynamical and fixed functional components into an integrated time series model with random mixtures. For the dynamical model, we propose a two-parameter hyperbolic function, which is particularly suited to the diverse types of profiles encountered in expression studies.

We use an EM algorithm to fit the model and clustering is accomplished by monitoring Bayesian posterior probabilities.

Our method is applied to the V65 mouse Embryonic Stem Cells line at 6 time points over the first 24 hour differentiation period. We derive three large clusters, which turn out to be significantly functionally distinct according to both a gene ontology analysis and the previous stem cell literature.

Assembling a compendium of druggable apoptosis regulators

Wei Keat Lim and Andrea Califano

Center for Computational Biology and Bioinformatics (C2B2) & Center for the Multiscale Analysis of Genomic and Cellular Networks (MAGNet), Columbia University, 1130 Saint Nicholas Avenue, New York, NY 10032, USA.

Negative regulation of the apoptotic machinery, due to mutations or biochemical inhibition, can facilitate aberrant cell accumulation and lead to oncogenesis. Thus, the ability to activate the apoptotic process by pharmacological intervention in specific cellular phenotypes may provide new therapeutic options in the treatment of cancer. While many studies have been successful in identifying various components of the complex apoptosis pathway, a genome-wide repertoire of apoptotic targets that are both phenotype specific and susceptible to pharmacological intervention is still an elusive goal. In this study, we adopt a reverse-engineering approach to infer the upstream regulators of a predefined panel of apoptosis-related genes in a human breast cancer context. Our analysis focuses on genes encoding druggable proteins, with the hope of identifying novel cancer therapeutic targets.

We use the MINDy algorithm¹ to identify modulators of transcription factors regulating a panel of 392 genes with a pro- and anti-apoptotic biological process annotation in the Gene Ontology database. As an input to the algorithm, we use 218 Connectivity Map microarrays from MCF-7 cells treated with a large panel of chemical perturbagens². MINDy measures the difference in mutual information (ΔI) between the expression profiles of a transcription factor and its putative apoptotic target gene(s), when a putative modulator gene is either up or down-regulated. Positive or negative ΔI values respectively indicate a gain or a loss of transcriptional regulation activity, as a function of the modulator availability. Pro- and anti-apoptotic modulators can then be inferred by performing a simple t-test on the expression of the specific apoptosis panel gene across the two subsets. Candidate modulators included in this study must: (1) have a sufficient expression range, (2) contain a druggable InterPro domain, and (3) belong to one of 5 major druggable protein families, including protein kinases, G-protein-coupled receptors, ion channels, proteases, and phosphatases.

To quantitatively assess the performance of this algorithm, we used the Gene Set Enrichment Analysis (GSEA) program. Specifically, we evaluated enrichment of druggable targets known to be pro- or anti-apoptotic in the ranked modulator list produced by the MINDy analysis. The algorithm was set to implement a weighted scoring scheme and the enrichment score significance was assessed by 1000 permutation tests. GSEA exhibits much more statistically significant enrichment in both the anti-apoptotic set (FDR<0.001) and the pro-apoptotic set (FDR<0.06), compared to the enrichment using only correlation of the targets with the apoptotic panel (FDR<0.231 and FDR<0.149, respectively). We conclude that the proposed approach provides an informative ranking of druggable, pro- and anti-apoptotic for further biochemical/therapeutic validation.

1. K. Wang, et al., *LNCS* **3909**, 348 (2006).
2. J. Lamb, et al., *Science* **313**, 1929 (2006).

DIAL: An algorithm to detect RNA 3D motifs

Fabrizio Ferre*, Yann Ponty †, William A. Lorenz †, Peter Clote †

(* *Harvard Medical School, Children's Hospital, Hematology/Oncology Department, Boston, MA 02115 fferre@enders.tch.harvard.edu*

(† *Department of Biology, Boston College, Chestnut Hill, MA 02467, USA, {lorenzwi,ponty,clote}@bc.edu. This work is funded in part by NSF DBI-0543506.*

DIAL (dihedral alignment) is a web server that provides public access to a new dynamic programming algorithm for pairwise 3-dimensional structural alignment of RNA. DIAL achieves quadratic time by performing an alignment that accounts for (i) pseudo-dihedral and/or dihedral angle similarity, (ii) nucleotide sequence similarity (iii) nucleotide base-pairing similarity.

DIAL provides access to three alignment algorithms: global (Needleman-Wunsch), local (Smith-Waterman) and semiglobal (modified to yield motif search). Suboptimal alignments are optionally returned, as well as Boltzmann pair probabilities $\Pr(a_i, b_j)$ for aligned positions a_i, b_j from the optimal alignment. If a nonzero suboptimal alignment score ratio is entered, then the semiglobal alignment algorithm may be used to detect structurally similar occurrences of a user-specified 3-dimensional motif. The query motif may be contiguous in the linear chain or fragmented in a number of noncontiguous regions.

The DIAL web server provides graphical output which allows the user to view, rotate and enlarge the 3-dimensional superposition for the optimal (and suboptimal) alignment of query to target. Although graphical output is available for all three algorithms, the semiglobal motif search may be of most interest in attempts to identify RNA motifs. DIAL is available at <http://bioinformatics.bc.edu/clotelab/DIAL>.

Peromyscus maniculatus and Mus musculus Synteny Map Reveals Divergent, Patterns of Breakpoint Reuse in Rodentia

E.E. Mlynarski¹, C.J. Obergfell¹, W. Rens², P.C.M. O'Brien², C.M. Ramsdell^{3,4}, M.J. Dewey³, M.J. O'Neill¹, R.J. O'Neill¹

Corresponding author: R. J. O'Neill rachel.oneill@uconn.edu

Presenting author: Elisabeth Mlynarski, graduate student UConn

Peromyscus maniculatus and *Mus musculus* are distantly related species, each more closely related to *Rattus norvegicus* than to each other. The genomes of both *Mus musculus* and *Rattus norvegicus* have been extensively studied, but despite the emergence of *Peromyscus* as a NIH model for biomedical research much remains unknown about its genome. Through reciprocal chromosome painting using fluorescence-labeled flowsorted chromosome probes between *Mus musculus* and *Peromyscus maniculatus* we provide the first chromosome homology map between these two species. Despite their superficial resemblance, *Peromyscus* and *Mus* have undergone extensive chromosome rearrangements since they last shared a common ancestor 25mya. Chromosome homology comparisons between *Mus*, *Rattus*, *Peromyscus* and its sister taxon *Cricetulus griseus* using GRIMM Multiple Genome Rearrangement analysis reveals a high-level of chromosome breakpoint conservation between *Rattus* and *Peromyscus* and indicates that the chromosomes of *Mus* are highly derived. This analysis identifies several conserved evolutionary breakpoints that have been reused multiple times during rodent speciation and allows for the first derivation of a high-resolution ancestral karyotype for Muroidea. The GRIMM analysis also indicates that sequence based rodent phylogenies do not accurately depict the complex evolution of genomes that has occurred within the rodent lineage.

Temporal logic, not combinatorial logic, dominates the regulation of sugar catabolic genes in E. coli

Ilaria Mogno (mogno@bu.edu) and Timothy Gardner (tgardner@bu.edu)

Department of Biomedical Engineering, Boston University

We studied the temporal response of 10 E. coli sugar metabolism operons to transient carbon starvation during diauxic shift from glucose and another sugar. Using GFP reporters to measure real time expression of the sugar operons and their transcriptional regulators, we found that these operons respond to the shift with a transcriptional burst of expression followed by a quick drop to intermediate or basal levels. Many operons showed bursting irrespective of the sugar used in the experiment, while others showed the bursting for a subset of sugars. Interestingly, the sugar specificity of the burst behavior is partially determined by the operon's regulatory network motif. Moreover, the data showed that the initial expression burst in these sugar operons is primarily driven by the activation of CRP/cAMP, not by increased activity in the sugar-specific regulatory proteins. The transcription factors at these promoters appear to operate independently across time rather than as a combinatorial logic gate. To test this hypothesis, we stimulated each of the operons with different levels of their associated sugars and cAMP to determine their logic gates. By comparing the gates to the operon response during diauxic shift, we confirmed that simultaneous high-level activation of CRP and the sugar specific regulator is not attained under physiological conditions. Rather, an initial nonspecific starvation response of the operons is driven almost exclusively by CRP, while a subsequent weaker response is driven and sustained by the sugar-specific regulator. This new vision of the sugar logic gates better clarifies the temporal dynamic of sugar metabolism in bacterial cells, and suggests that the regulatory logic observed during steady state conditions may not reflect the role served by that gate under dynamic physiological conditions.

Form, function, and information processing in small stochastic biological networks

Andrew Mugler¹, Etay Ziv², Ilya Nemenman³, and Chris H. Wiggins⁴

¹Department of Physics, ²College of Physicians and Surgeons, ⁴Department of Applied Physics and Applied Mathematics, ⁴Center for Computational Biology and Bioinformatics, Columbia University, New York, NY, USA; ³Computer, Computational and Statistical Sciences Division, ³Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, NM, USA

There has been much recent interest in the extent to which the topology of a regulatory network dictates or constrains its function. Following the experimental setup of Guet et. al. [1], we stochastically model a set of small regulatory networks, treating chemical inducers as functional “inputs” and the steady-state concentration of a reporter gene as the functional “output.” We take an information-theoretic approach, allowing the system to set various biochemical parameters that optimize signal processing ability. The advantage is twofold: we obtain a measure of network fidelity that does not require a predetermination of the network’s function(s), and we enumerate the most informative functions a network can perform.

We find that all networks studied can exhibit surprisingly high fidelity, despite the noise imposed by the biologically realistic constraints of low molecule number and low latency. Specifically, all networks can transduce more than one bit of information, showing that reporter concentration can encode more than the simple “on/off” description employed by the majority of molecular biology literature. Additionally, we find that transduction capacity is higher for network topologies in which the overall feedback is negative, a result we support analytically.

We show analytically that topology constrains a network to a particular subset of steady-state functions. Nonetheless, we find computationally that each network is able to perform a variety of high-information functions from within this subset, allowing us to quantify the extent to which a change in function depends on a change in model parameters. We present a non-parametric measure of “evolvability,” the degree to which a network can perform markedly different functions with only minimal changes in parameters, and discuss correlations of this measure with topological features.

[1] Guet CC, Elowitz MB, Hsing W, Leibler S (2002) Combinatorial synthesis of genetic networks. *Science* **296**, 1466-1470.

Generalized dynamic theory of alternative splicing regulation

Rajamanickam Murugan and Gabriel Kreiman

*Department of Neurology and Ophthalmology, Children's Hospital Boston
Harvard Medical School, MA 02115, USA*

The coding sequences of most eukaryotic genes are fragmented into many exons separated by introns. The transcription of such interrupted genes produces the primary mRNA (pre-mRNA) transcript which includes both the exons as well as the introns. Recent studies show that a majority of the genes in the higher eukaryotes are spliced alternatively and therefore significantly contribute to the structural and functional complexity of the organism. However the principles associated with exon inclusion and exclusion are not clearly understood. The inclusion or the exclusion of an exon is positively correlated with the presence or the absence of the *cis-acting* exonic enhancer (EEE) or the exonic silencer elements (ESE) (static factors). Based on the presence or absence of these elements, we can classify the exons into two major categories viz. strong and the weak exons. Most of the currently available splice-pattern prediction models and algorithms are solely based on these static factors and none of the models considers the underlying stochastic dynamics involved in the snRNP-pre-mRNA interactions (dynamical factors). Here we present a model that shows that the splice-pattern predictions based on the static factors may be valid only when the size of the pre-mRNA is small and transcription is physically uncoupled from splicing. When transcription is coupled with splicing as in most eukaryotic genes, then the dynamical factors can significantly modify the predictions. In this context, here we develop a generalized phenomenological theory of alternate splicing which predicts that the probability associated with an exon to be included in the final transcript is not only decided by the presence or absence of the *cis-acting* regulatory elements and the transcription elongation rate as predicted by the kinetic model but also by (a) the relative position of that exon in the pre-mRNA transcript and (b) the jump size associated with the one dimensional diffusion dynamics of the sn-RNPs in the process of searching for the intron-exon junctions on the pre-mRNA lattice. Detailed studies show that as the jump size increases, the distribution of the splicing patterns shifts towards the canonical type. The jump size required to generate the pure canonical transcript scales with the size of the pre-mRNA lattice in a square root manner. Our theory also predicts that under coupled conditions there exists a critical range of positions on the pre-mRNA lattice in such a way that those exons present in this range possess equal and higher probabilities to be included in the final transcript than that of those exons falling outside this critical range. This prediction is indeed consistent with the distribution of occurrences of exons in the final transcript which are computed from the exon micro array results.

An All-or-None Response in Osmolyte-Dependent Gene Expression in *Saccharomyces cerevisiae*

Gregor Neuert and Alexander van Oudenaarden

Department of Physics, Massachusetts Institute of Technology, Cambridge, MA, USA

How cells sense their environment using signal transduction pathways and respond to this information by regulating gene expression is a key problem in systems biology. Here we investigate the coupling between signal transduction and gene expression in budding yeast. We monitor the activation of the osmo-sensitive kinase Hog1 and its resulting gene expression in single yeast cells. Performing single cell time-lapse experiments, we find that, although all cells show a similar activation dynamics of Hog1, downstream gene expression of *STL1* displays an all-or-none response. We observe that a fraction of cells show no gene expression at all, where other cells show gene expression over a wide range of expression levels. This fraction is dependent on the intensity of the osmo-shock.

An ultra-conserved centromere-specific retrovirus is an epigenetic determinant of centromere function in vertebrates

Rachel J. O'Neill, Ph.D.

*Associate Professor, Genetics and Genomics
University of Connecticut*

Little is known about the function of centromere satellite RNAs in vertebrate cells or the elements controlling their transcription. Nonetheless, evidence suggests that they play a critical role in centromere identity and in the recruitment of kinetochore proteins. Our investigation of centromeres of vertebrate species spanning 300 million years of evolution has uncovered a previously unidentified conserved sequence element, a retrovirus. We have determined that this novel retrovirus is composed of ultra-conserved sequence elements, that it is transcribed bi-directionally and produces small non-coding RNAs of a novel class size. These RNAs are involved in heterochromatin recruitment kinetochore demarcation and, therefore, play an apparently domesticated role in centromere function. Furthermore, examination of the distribution of this centromeric retrovirus within vertebrates shows these sequences to be associated with centric shifts and convergent breakpoint reuse. The instability of these sequences within interspecific hybrids is directly associated with chromosome rearrangement and changes to chromatin structure. Taken together, these data indicate that this virus is not only involved in centromere function, but is also participant in chromatin remodeling involved in chromosome rearrangement. The identification of such a well-conserved, functional retrovirus, and the novel class of small RNAs they produce, forces a re-evaluation of theories of centromere function and evolution in vertebrate lineages.

Structure and Function of the Transcriptional Network Activated by the MAPK Hog1

Erin O'Shea*, Andrew Capaldi, Kristen Cook, Ying Liu, Tommy Kaplan, Naomi Habib, Aviv Regev, Nir Friedman

**Howard Hughes Medical Institute, FAS Center for Systems Biology, Departments of Molecular and Cellular Biology and Chemistry and Chemical Biology, Harvard University, Cambridge MA 02138*

To gain insight into transcriptional network structure and function we are studying the response of budding yeast to hyper-osmotic stress. A single kinase, Hog1, controls the expression of over 300 genes by modulating the activity of several transcription factors. We have determined the influence that each transcription factor has on gene expression, both alone and in combination, by conducting systematic microarray and ChIP-chip experiments in different genetic backgrounds. Computational analyses of these data have allowed us to construct a detailed model of the interactions and activities underlying this transcriptional response. Our findings reveal a dense overlapping network structure where the Hog1 dependent transcription factors act in different combinations, at different genes, to precisely define the output. This architecture allows Hog1 to control a highly diverse transcriptional program by incorporating both osmotic stress-specific and general stress transcription factors. Most remarkably, however, we find that this network structure also allows Hog1 to activate alternate expression programs in other stress conditions, by integrating information from other pathways. We demonstrate that condition specific changes in the activity of Hog1, or the general stress transcription factors, leads to a range of expression changes at the gene level, due to the variety of transcription factor combinations and differences in their combinatorial logic. Therefore, an overlapping network architecture allows the cell to create a fine-tuned response over hundreds of genes using a minimal number of signaling pathways and transcription factors, providing an explanation for its overrepresentation among biological network motifs.

A generalized Karlin-Altschul statistics for modeling conserved cis-regulatory modules

Kimmo Palin¹, Jussi Taipale², and Esko Ukkonen¹

1. *Department of Computer Science, P. O. Box 68 (Gustaf Hällströmin katu 2 B) FIN-00014 University of Helsinki, Finland*

2. *Molecular Cancer Biology Program / Biomedicum Helsinki PO Box 63 (Haartmaninkatu 8) FI-00014 University of Helsinki, Finland*

In [1], a novel model for conserved cis-regulatory modules (CRMs) was introduced and its biological significance was demonstrated experimentally. The model is based on scoring conserved patterns (clusters) of transcription factor binding sites in DNA. High scoring clusters are putative cis-regulatory modules. They can be found using a Smith-Waterman type of local alignment algorithm that evaluates the score. The score is called the EEL score according to the software tool (Enhancer Element Locator) developed in [1].

Here we improve the CRM model of [1] by developing a method to estimate the statistical significance of the EEL scores under the null hypothesis of neutral conservation. Our statistical method applies the two component extreme value distribution (TCEV) to the EEL scores of putative CRMs in neutrally evolved DNA sequences. The TCEV distribution is similar to the well-known Karlin-Altschul statistics but includes an additional component accounting for the neutral conservation due to finite evolutionary distance. The parameters of the TCEV distribution are estimated by a non-linear least-squares fit to the EEL score distribution from simulations of sequence evolution using the HKY substitution model augmented with insertions and deletions. The TCEV distribution has a remarkable fit to the observed score distribution over large range of scores and sequence lengths. We apply the statistical model to the CRMs found by comparing non-coding orthologous genomic sequences from human and mouse.

It is important to notice that our TCEV based generalization of the Karlin-Altschul statistics is not restricted to the EEL scores. It is applicable in most pairwise local alignment settings for identification of the conservation that is above the neutral one. Note that the original Karlin-Altschul statistics assumes that the sequences are completely independent. Our generalization allows the assumption that they have evolved from a common origin.

[1] Outi Hallikas, Kimmo Palin, Natalia Sinjushina, Reetta Rautiainen, Juha Partanen, Esko Ukkonen, Jussi Taipale: Genome-wide Prediction of Mammalian Enhancers Based on Analysis of Transcription-Factor Binding Affinity *CELL* 124(1), 13 January 2006, Pages 47-59.

A mechanism for bistable gene expression regulation in yeast

Saurabh Paliwal and Andre Levchenko

*Department of Biomedical Engineering, The Johns Hopkins University,
3400, N. Charles Street, Baltimore, MD, USA*

Many natural and engineered cellular biochemical networks display bistability or multistability i.e. the capacity to show two or more stable steady states in response to a single external input, thus resulting in biological switches. Such switches are inherent to several biological phenomena including cell fate decisions involving growth and differentiation, the regulation of cell-cycle oscillations as well as the basis for stochastic noise induced bimodal gene expression in individual cells.

We have previously shown that the yeast pheromone response pathway displays bimodal gene expression in a specific range of pheromone concentrations, resulting from a combination of bistability in pheromone induced gene expression and stochastic noise. Bistability is typically attributed to the presence of positive feedback loops in the signaling network. In this pathway, MAPK (Fus3 and Kss1) mediated phosphorylation of the mating specific transcription factor Ste12 followed by up-regulation of the expression of mating-induced genes results in several positive feedback loops. Additionally, the presence of unphosphorylated Kss1-mediated repression of Ste12 contributes to the bistability in the network. In light of this biological example, we have investigated the role of negative regulation in modulating the threshold properties of bistable responses. We perform a detailed analysis of a simplified mathematical model, and propose that this could be a recurring module in several other biological systems.

Efficient Discovery of Significant Co-Occurrences with an Application to Dyad Motifs Finding

Cinzia Pizzi¹, Esko Ukkonen²

¹ *Projet Helix, INRIA Rhone-Alpes and Lab. de Biometrie et Biologie Evolutive, CNRS, Univ. Lyon1,* ² *Dept. of Computer Science, Univ. Helsinki, Finland*

The computation of statistical indexes such as frequency counts, expected probabilities, and over/under-representation scores is a routine task shared among several biological problems. In fact, they often constitute the core operation of pattern discovery algorithms used to extract DNA or proteins motifs from a sequence or a set of related sequences. We are interested in the problem of computing indexes related to co-occurrences of elements in a sequence. This can be applied to discovery regulatory elements that can be modeled as dyad motifs, i.e. two conserved components that co-occur at some distance d .

Exact methods and heuristics suffer from limitations arising respectively from the exponential time and space complexity, and from the lack of guarantee of finding the optimal solution. An alternative approach to the problem is the *compact* or *conservative* approach. This is based on the notion of *maximality* to partition the search space in a set of classes, and evaluate only one representative for each class. If we consider the space of all the strings in a sequence, the classes can be identified by all strings that share the same set of starting positions. The representative maximal element is the longest string in the class. The suffix tree data structure naturally detects such representatives as the strings spelled out by any path from the root to each of its nodes. When counting co-occurrences measuring head-to-head distance up to a value d , it suffices to compute co-occurrences count only between maximal representatives. Any pair that is left out is associated with a maximal pair that has been counted. Since there are $O(n^2)$ words in a sequence of length n , but only $O(n)$ maximal representative, computing the co-occurrence count require only $O(n^2)$ time and space rather than $O(n^4)$ [Apostolico et al, 2004]. Tail-to-head distance between the two components might be more natural to consider in real application. To compute tail-to-head distance we first compute co-occurrence count between maximal for all possible head-to-head distances. This requires $O(n^3)$ time and space. Then we can compute tail-to-head distance for any pair of words in constant time.

In order to assign a significance score to each dyad motif, besides efficient counting it is necessary to study the behavior of the expectation of pairs among each class. The results of our analysis showed that either considering a background probability derived from the input strings, or from an external source, the expectation has a constant or decreasing behavior when the length of the strings increases. This implies that the score associated with the maximal element is the maximum score that can be achieved. Experiments on dyad motifs discovery where the dataset consists in 8 gene families that are regulated by zinc-bicluster transcription factors, proved our method to be effective and much faster than the exhaustive approach [VanHelden et al, 2000].

LocalMove: A web server to compute best on-lattice fits for biopolymers

Yann Ponty, Peter Clote

*Department of Biology, Boston College, Chestnut Hill, MA 02467, USA,
{ponty,clote}@bc.edu.*

LocalMove₁ is a 4500 line C++ program we have developed, that implements a Monte Carlo approach to obtain approximate on lattice-models for protein and RNA Protein DataBank (PDB) files. Lattices currently supported are the square 2D lattice, cubic 3D lattice and the face-centered cubic (FCC) lattice. Simulated annealing and greedy descent are additionally featured by the LocalMove program.

Other methods have been used to determine optimal cubic 3D on-lattice fits of C_{α} -atom traces of proteins conformations available in the PDB. Rabow and Scheraga have considered a variant of Hopfield artificial network, while Reva et al. have used self-consistent field optimization method, involving a combination of statistical mechanics of chain molecules, Monte-Carlo and dynamic programming. It seems likely that computing an optimal on-lattice fit is NP-complete, although this should at present be considered a conjecture on our part. For this reason, all currently known methods use a heuristic to approximate the best on-lattice fit. Our approach yields a significant and consistent improvement over previous works from Reva and al. Namely, preliminary results show an average 15.79% improvement on the root-mean square deviation, a classic measure of the quality of on-lattice fits.

Cipra has proven that face-centered cubic (FCC) lattice provides the tightest-packing of spheres, and Covell and Jernigan have shown that the FCC lattice is the most appropriate 3D lattice for fitting protein C_{α} -atoms as a self-avoiding walk; i.e. root mean square deviation values are smaller for the FCC than for the cubic, body-centered cubic and tetrahedral lattices. To the best of our knowledge, LocalMove is the only program providing an onlattice fit for the FCC lattice. Proteins and RNA can be modeled using the coarse-grain monomer model (C_{α} for amino acid, C_1 for RNA), as well as considering all backbone atoms. LocalMove implements a Monte-Carlo algorithm involving prefix, suffix and crankshaft moves, where all possible length $k \leq 6$ self-avoiding walks starting from the origin (0, 0, 0) are computed and used for crankshaft moves.

A webserver featuring a very large panel of options, in addition to a real-time visualization of the pseudo-folding process can be found at the URL: <http://bioinformatics.bc.edu/clotelab/localmove/index.html>

[†]This work is funded in part by NSF DBI-0543506.

Evolution and selection in yeast promoters: analyzing the combined effect of diverse transcription factor binding sites

Daniela Raijman¹, Amos Tanay² and Ron Shamir¹

¹*School of computer science, Tel Aviv University, Ramat Aviv 69978, Israel.*

²*Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 76100, Israel*

The evolutionary dynamics of regulatory regions is shaped by heterogeneous factors, notably multiple transcription factors (TFs). The interactions among these factors are manifested as evolutionary interactions between loci (epistasis). Computational modeling of evolution under epistatic interactions is difficult because such interactions invalidate the standard assumptions on the independence of evolving loci. The majority of models that are used in comparative genomics are therefore searching for conservation without explicitly formalizing the evolutionary process in terms of mutation and selection. Here we present a probabilistic model for the evolution of promoter regions given a neutral background mutation process and selection on isolated and overlapping TF binding sites. The model integrates mutation and selection at the process level, taking into account the epistasis between loci that are part of the same or overlapping binding sites, and allowing the assessment of the model's goodness of fit to the observed divergence patterns using principled likelihood computation.

Using our model, and new algorithms that support computation in this model, we examine the evolutionary dynamics in *Saccharomyces* species promoters. We study comprehensive evolutionary models, which encompass (in a single model) dozens of literature based or algorithmically derived TFs. We quantify the intensity of selection on the targets of TFs and the overall selection pattern of motifs that are targets of TFs and of motifs that are similar to such targets. Our results reveal that transcriptional interactions in yeast are evolutionary plastic with surprisingly weak selection affecting the targets of most TFs. In addition, we identify a gap between global estimates of selection on yeast promoters and the selection that our model can attribute to characterized transcriptional interactions. Our evolutionary studies therefore support the hypothesis that much of the functionality (and therefore selection) in the yeast regulatory network is encoded at loci that are not part of clearly defined, high specificity binding sites. We conjecture that rather than being a deterministic and sparse network, the yeast transcriptional program is shaped as a dense and noisy network that is continuously changing during evolution.

Inferring Interactions in Transcriptional Regulatory Networks

Anil Raj^{1,2}, Chris H. Wiggins^{1,2}

¹Department of Applied Physics and Applied Mathematics

²Department of Computational Biology and Bioinformatics

Columbia University, New York, NY, USA

Correspondence: ar2384@columbia.edu

We develop an algorithm to learn a real-valued transcription factor(TF)-gene affinity-scoring model that relies on availability of binding sites of known TFs and the whole genome of the organism. This model can be used to infer TF-gene interactions and can be used as a starting point for experimental techniques currently being used like binding of purified proteins and gene expression analysis. The TF-gene affinity-scores correlate with the frequency of binding sites of a TF p within k bases upstream of a gene g . This score matrix S_{gp} and the list of known interactions can then be used to calculate an ROC score for each value of k as a measure of model accuracy.

Experimental Results

We obtain experimentally determined TF-gene interaction data for E. Coli K12 [1] consisting of 1317 genes and 96 TFs and binding site data (motifs) for each of those TFs. We construct a binary adjacency matrix and normalize it. We also obtain the entire genome of E. Coli K12 from GenBank [2] and use it to construct a list of k -base sequences upstream of each of the 1317 genes. We then generate an adjacency matrix of Frequency of motif μ upstream of gene g and normalize it. An affinity score between genes and TFs can then be constructed (see Poster). Hence, we scale the affinity scores to give a diagonal matrix of lengths of binding sites of transcription factors.

Using the known interactions, we can generate an ROC curve by varying the discriminative threshold (see Poster). We can also calculate the ROC score for each k by calculating the AUC. We note a maximum ROC score of 0.88 for $k = 14500$ and an ROC score of 0.80 for a mere $k = 3000$. An important limitation of this algorithm, however, is the strong specificity it demands between transcription factors and binding sites. Future research will aim at constructing sets of motifs, similar to the known binding sites, to which the TFs could have real-valued binding-affinities and accounting for them as well.

References

[1] Salgado H, et. al., Nucleic Acids Res. 2006 Jan 1;34(Database issue) D394-7 (2006)

[2] National Library of Medicine. <http://www.ncbi.nlm.nih.gov/Ftp/>

Simulated Evolution of Transcription Regulatory Networks and Network Motifs

Arvind Ramanathan, Michelle Goodstein, Robbie Sedgewick, and Dannie Durand

Joint CMU/Pitt Program in Computational Biology, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh, PA, USA

Computer Science Department, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh, PA, USA

Several experimental tools are available today to analyze cellular machinery in various organisms like mass spectrometry, genome-wide chromatin immuno-precipitation and yeast two hybrid assays. The data on interactions (usually represented in a network; called a Biological Network) from several thousand proteins as well as genes are available. Much of the data gathered is known to have a high degree of error; finding network motifs, then, becomes a challenge of detecting whether the observance of a motif is significant in comparison to a null model. Studying experimental noise is an interesting challenge in that once the data becomes available, the noise cannot be separated from it. We focused instead on mathematically analyzing the impact of three different error models on the expected number of counted motifs. These models were *node removal*, where a node in the network was removed with some probability p ; *edge removal*, where an edge in the network was removed with some probability p ; and *edge reversal*, where an edge (i, j) is changed to (j, i) with some probability p . We studied the affect of these error models on three common motifs: single input modules (SIMs), bi-fans and feed forward networks. Our analysis is independent of models of the graph, and relies only on assumptions about the error. Our analysis shows that the type of error present when data is collected can definitely affect the counted number of motifs. Our work suggests that this type of analysis may allow biologists to study, if not which motifs are in error, at least which numbers seem more significant.

AhoSoft: P-value Computation for Heterotypic Clusters of Regulatory Motifs

Mireille Regnier¹, Valentina Boeva^{1,2}, Zara Kirakosyan³, Vsevolod Makeev^{4,5}

¹INRIA Rocquencourt, 78153 Le Chesnay, France

²LIX, Ecole Polytechnique (1 128 Palaiseau Cedex, France

³YSU, 377 200 Yerevan, Armenia

⁴Institute of Genetics and Selection of Industrial Microorganisms, GosNIIGenetika, 117545 Moscow, Russia

⁵Engelhardt Institute of Molecular Biology, Russian Academy of Sciences, Moscow, Russia

A eukaryotic transcription regulatory region of DNA can be usually bound by multiple regulatory proteins, which recognize specific motifs in DNA texts. Identification of multiple transcription regulation binding sites (TFBSs) clustered within a short DNA segment became an important method of predicting location of eukaryotic regulatory regions. Many factors including motif overlapping and their fuzziness may make difficult to assess the statistical significance of the observed configuration of TFBSs.

Earlier, we developed an algorithm inspired by Aho-Corasick string matching algorithm, which calculates P-value for a number of high scoring hits of different motifs to occur simultaneously within a random text of a given length. More precisely, for k motifs determined by corresponding positional weight matrices (PWM), $M_1 \dots M_k$ scoring above corresponding thresholds $T_1 \dots T_k$, or by simple word enumeration we calculate the P-value that a random text of length N contains no less than $r_1 \dots r_k$ possibly overlapping motif occurrences. Markov or Bernoulli random models can be adopted for random text. The algorithm is implemented in the AhoPro program, which can be used or downloaded through our web-site www.bioinform.genetika.ru. It was successfully tested on developmental enhancers of *D. melanogaster*. However, employing the algorithm to real world genomic studies (long motifs, long regulatory regions, low PWM thresholds) brings about a number of limitations on computer memory and computational time. Here we suggest how the efficiency of the algorithm can be increased by preprocessing the AhoCorasick tree. The new design of the algorithm solves the problem in $O(nm|H|)$ time and $O(m^2|H|)$ space complexity where m is the length of the motif, $|H|$ is the number of allowed words in the motif and n is the size of the text. The tree is built only once, and its traversal is not required. After tree processing, it becomes substituted for a system of linear equations which can be efficiently solved with a symbolic computation system.

Age, sex and genetic architecture of human gene expression

Manuel A. Rivas, Mark J. Daly, and Itsik Pe'er

*Department of Mathematics, Massachusetts Institute of Technology, 372 Memorial Drive
Cambridge, MA, USA*

Broad Institute of MIT and Harvard, Cambridge, MA, USA

*Center for Human Genetics Research, Massachusetts General Hospital, Cambridge, MA,
USA*

Department of Computer Science, Cambridge, MA, USA

Genomewide association studies and expression profiling are powerful strategies to dissect genes that impact human health. Genetic genomics holds the promise to link these strategies to dissect the heritable factors that regulate all expression-mediated processes in the cell. Indeed, coupled whole genome genotype and expression data for human lymphoblastoid cell lines have recently been analyzed by classical methods to map numerous quantitative trait loci's. As the effects of age, sex, and heritability from each parent are routinely evaluated for genetically-studied disease traits, we sought to assess these effects on transcript levels.

In this study, we use a systematic approach to derive signatures of aging, sex, and genetics of gene expression. Genes located in the Human Leukocyte Antigen (MHC) class I region of the human genome are found to be enriched across tissues as a signature of aging. Furthermore, we derive an enriched set of gender specific signature located across the autosomal chromosomes of the human genome and pursue analysis of a robust filtering of genes to test for gene set enrichment for sex. We consider public data from CEPH extended pedigrees, including expression profiles as well as microsatellite and SNP genotypes and by doing so increasing power to detect genetic effects of gene expression levels. We discover 178 cis-acting genetic variants of gene expression levels along with six genome-wide significant trans effects. Finally, we evaluate the effects of genetically mediated age signatures of gene expression and associate a trans-acting SNP variant located in LRP1B, the receptor of APOE, regulating APOE expression in an age-dependent manner. These findings will allow us to decipher and identify the genetic and non-genetic effects of gene expression that could possibly influence disease etiology.

Inferring tissue specific activity of transcription factors from predicted DNA binding affinities to core promoters

Helge G. Roeder, Thomas Manke, Martin Vingron & Stefan A. Haas

Max Planck Institute for Molecular Genetics, Ihnestrasse 73, 14195 Berlin, Germany

Transcription factors (TFs) play a central role in the regulation of genes in terminally differentiated tissues. Deciphering the machinery by which the factors gain their tissue specific activity and establish gene expression patterns in different cell types is currently among the most important topics of research in both the fields of experimental as well as computational biology.

On the computational side the target genes of a TF are typically identified based on the prediction of TF binding sites in the respective promoter sequences. This is done by scanning the potential regulatory regions using position weight matrices (PWMs) derived from experimentally verified binding sites. Using such target gene predictions together with sets of tissue specific genes derived from expression data computational studies have been able to recover a number of well known associations between TFs and tissues. Pioneering this work Wyeth Wasserman and colleagues were able to first detect the association between a small set of muscle and heart specific genes and various TFs with known muscle specific function including SRF and MYOD. Since then, large scale studies recovered a more comprehensive but still rather limited number of experimentally known TF-tissue associations such as the link between immune cell specific genes and NF κ B or liver specific genes and HNF1. While most such studies select TF targets based on the prediction of discretized TF binding sites it was shown recently that using PWMs to predict continuous binding affinities instead of discrete TF hits yields higher correlation with experimental ChIP-chip data in yeast and is more successful in ranking TF target genes (Tanay 2006, Roeder et al. 2007).

Here we present a method to detect TF-tissue associations based on the prediction of such continuous TF binding affinities to core promoters of genes specifically expressed in given tissues. To this end we use our Transcription factor Affinity Prediction program TRAP (Roeder et al. 2007) to compute TF binding affinities for 26609 Ensembl mouse promoters and 589 TRANSFAC vertebrate PWMs. EST and microarray data is used independently to generate tissue specific gene sets. In order to detect functional associations between a given TF and a given tissue we test for the enrichment of genes specifically expressed in the tissue among the genes with highest predicted affinity for the TF using a hypergeometric test. We use this procedure to first identify the optimal region for detecting tissue-TF associations and find that primarily sequences up to 200 bps upstream of the TSS confer tissue specificity. We show that predicted TF affinities together with EST or GNF expression data allows for the detection of a comprehensive number of known tissue-TF associations including factors and tissues for which no associations could be found previously. Our findings are hereby rather robust against the size of the selected promoter regions, parameter settings and the source of expression data.

Thermodynamic-based approach for miRNA target prediction

Aleksey Y. Ogurtsov, Denis V. Karamyshev, Olga V. Matveeva, Svetlana A. Shabalina

National Center for Biotechnology Information, National Institutes of Health, Building 38A, 6S604, 9000 Rockville Pike MSC 3829, Bethesda, MD 20894, USA

By combining computational and experimental approaches, it was shown that miRNAs have relatively clearly defined patterns of complementarity to the 3'UTRs of their target transcripts. The so-called 'seed' region, positions 2-8 of miRNAs from 5'end, has been determined as a key specificity feature of binding to target and requires perfect complementarity. Here we showed that feature types involved in miRNA–target 3'UTR interactions include structural, thermodynamic and position-based elements. To describe more accurately the mechanism of the seed pairing, we paid special attention to the position-based features. We developed thermodynamic-based approach for miRNA target predictions and showed that the efficiency of miRNA binding site search improved, when sequence predictors were supplemented with structural features. Our hybridization approach mainly determines sequence complementarities measured through favorable miRNA-target duplex stability. Our statistical models and an efficient algorithm that scans the hundreds and thousands of hybridization duplexes composing an output file detailing each possible miRNA/3'UTR binding site was produced in this study. Since miRNAs frequently bind imperfectly to target 3'UTRs, the stability of optimal matches of miRNA/3'UTR binding sites are discussed. Our past extensive search of 3'UTR targets for experimentally verified miRNAs using our hybridization approach has proved certain ranges for thermodynamic stability that have led to more reliable miRNA predictions.

PhyloGibbs-MP: Motif-finding, module prediction and differential motif prediction by Gibbs sampling

Rahul Siddharthan

The Institute of Mathematical Sciences, Chennai 600113, India

We present PhyloGibbs-MP[1], an update to the motif-finder PhyloGibbs[2]. PhyloGibbs is a motif-finder designed to detect binding sites for regulatory proteins in orthologous sequence from closely-related species, scoring conserved sequence by an evolutionary model, and can also run on single-species sequence. Additionally, it makes a careful assessment of the posterior probability of its predictions, via further sampling and “tracking” of a candidate set of motifs, and we demonstrated that this self-assessment is accurate on synthetic data and credible on actual biological data. PhyloGibbs-MP improves the motif-finding performance, both in speed and in accuracy, of PhyloGibbs, while extending the capabilities of PhyloGibbs in significant new directions. First, it is able to localise predictions to small, unspecified regions of large input sequence, meant to indicate *cis*-regulatory modules; it can thus simultaneously perform the tasks of CRM prediction and motif finding. We demonstrate that its performance as a CRM predictor, without prior information on known regulatory proteins and their position weight matrices (PWMs), is comparable with that of specialised CRM prediction programs that do make use of known PWMs. However, PhyloGibbs-MP can use prior PWMs to improve its predictions if necessary. Second, PhyloGibbs-MP can predict motifs that differentiate between two (or more) different groups of regulatory regions, that is, motifs that occur preferentially in one group compared to the others. This feature, that appears to be new among published programs, has clear advantages in studying contrasting gene clusters obtained, for example, from high-throughput (microarray) expression studies. There are other improvements, including output in an annotation format that can be read by the Generic Genome Browser. We briefly discuss some completed work, and work in progress, that uses PhyloGibbs-MP. As part of benchmarking for this program, we have clustered known DNase I footprints from the Flyreg database to obtain high-quality position weight matrices for 73 transcription factors in fruitfly, and have used these in analysing results of a detailed study of the regulation of a key gene in mesodermal development in *D. melanogaster*[3]. PhyloGibbs-MP can also be used successfully for motif-finding outside the context of regulatory genomics, and we will discuss a recent example[4].

Availability: <http://www.imsc.res.in/~rsidd/phylogibbs-mp/>

[1] R. Siddharthan, manuscript in preparation

[2] R Siddharthan, E D Siggia, E v Nimwegen, PloS Comp Biol 1(7):e67 (2005)

[3] K G Guruharsha, Mar-Ruiz Gomez, H.A. Ranganath, Rahul Siddharthan and K. Vijay-Raghavan, manuscript in preparation

[4] K Sanyal, R Siddharthan et al, manuscript in preparation

Rapid chromatin structural phenotyping by Digital DNaseI

Dorschner MO, Sabo PJ, Kuehn MS, Sandstrom R, Weaver M, Neph S, Haydock A, Goldy J, Thurman RT, [Stamatoyannopoulos JA](#)

Dept. of Genome Sciences, University of Washington, Seattle, WA, 98195

Chromatin accessibility plays a defining role in gene and genome regulation. In its cognate cell type, regulatory DNA becomes universally hyper-accessible relative to bulk chromatin. We developed a simple and powerful approach – Digital DNaseI – for rapidly generating high-resolution, whole-genome maps of chromatin accessibility via direct enumeration of individual DNaseI cleavage events in nuclear chromatin. We coupled Digital DNaseI with a next-generation sequencing platform to assay the chromatin phenotypes of diverse human cell types, including primary lineages. A striking feature of Digital DNaseI is the ability to ascertain completely the chromatin accessibility state of all known human canonical and alternative promoters using only a small fraction of the data generated during a single sequencing run. We find that the promoter accessibility state of a given cell type is rigidly stereotyped, with clear segregation of individual promoters into accessible and inaccessible chromatin compartments. We found surprisingly marked differences in promoter chromatin state between even closely related cell types within the same lineage. This suggests that the genome is highly plastic and retains the potential for widespread reprogramming of chromatin regulatory state during terminal cell differentiation. Digital DNaseI is rapid, inexpensive, and may be applied to capture the chromatin structural phenotype of virtually any eukaryotic cell type.

Regulatory Motif Discovery and Target Prediction in 12 *Drosophila* Genomes

Alexander Stark*, Pouya Kheradpour*, Sushmita Roy, and Manolis Kellis

MIT Computer Science and Artificial Intelligence Laboratory; Broad Institute of MIT and Harvard; University of New Mexico.

* = contributed equally

The expression of animal genes is regulated both pre- and post-transcriptionally, via subtle DNA and RNA sequence signals, forming several classes of diverse regulatory motifs. A systematic understanding of gene regulation relies on global knowledge of these motifs and the regulatory targets they define, which are often elusive due to their short lengths, several weakly-specified positions, and large distances at which they can act.

We used whole-genome alignments of 12 *Drosophila* species for the systematic discovery of regulator targets in the fly genome. We developed a phylogenetic framework that accounts for the evolutionary relationships between species, allows for motif movements, and is robust against missing data due to artifacts in sequencing, assembly or alignment. We also developed a statistical framework for evaluating robust measures of motif confidence, enabling us to seamlessly translate evolutionary conservation into confidence for each motif instance, correcting for varying motif length, composition, and background conservation.

Using this framework, we predicted targets for 83 transcription factors and 57 miRNAs, resulting in 46,525 regulatory connections. When compared to extensive genome-wide experimental data, predicted targets are of very high quality, matching and surpassing ChIP-chip, and showing higher functional enrichments even outside ChIP-bound regions, suggesting that many biologically-relevant targets are missed by traditional experimental approaches. The resulting regulatory network shows heavy targeting of known regulators, evidence of extensive regulatory feedback, reveals many intriguing novel regulatory connections, and suggests significant cooperation and redundancy between pre- and post-transcriptional regulation.

In addition, we used the 12 genomes for systematic discovery of novel motifs in promoters, enhancers, introns, 5' and 3'-untranslated regions (UTRs), protein-coding exons, and across all intergenic regions. We distinguished motifs acting at the RNA level by their strand-specific conservation, and defined regulatory motifs in protein-coding sequence based on their reading frame-independent conservation. Our search reveals many novel elements enriched in genes of common function and similar expression patterns, and depleted in ubiquitously expressed genes, similarly to known motifs. Interestingly, the most highly conserved motifs in coding sequence are complementary to known miRNA 5'ends, suggesting that coding exons contain physiologically relevant and evolutionarily selected miRNA target sites.

The inferred regulatory network provides a powerful resource for systems-level understanding of tissue formation and development in the fly. Moreover, the methods presented are general and may serve as a model for similar studies in the human genome, as dozens of new mammalian genomes become available.

A Mining Tool for Regulatory Elements in Zebrafish

TePLYuk, V¹, Hirsch, N.², Sagerstrom, C², [Straubhaar, J¹](#)

University of Massachusetts Medical School, Program in Molecular Medicine¹, Department of Biochemistry², and institutional Diabetes Endocrinology Research Center (DERC)¹.

A standard method to study developmental processes is to assay the global transcriptional expression changes that occur during development. It is assumed that the resulting expression patterns can be explained by sets of genes acting together in a network through various kinds of interactions reminiscent of electrical circuits. Furthermore, the memberships of these modules can be inferred from the identity of regulatory motifs in the promoter sequences of genes with similar expression patterns. As a concrete example of such a methodology we have built a database of regulatory motifs in zebrafish promoters we call ZAN. We used an mRNA misexpression approach to induce ectopic expression of Hoxb1b target genes in the zebrafish anterior CNS. RNA from dissected anterior CNS was then used as a probe to screen microarrays (Compugen zebrafish oligo set) for genes upregulated by Hoxb1b. The genomic regions from +20kb to -20kb relative to the TSS of each upregulated gene was obtained from the Ensembl database. These were scanned with HMM models of transcription factor binding sites and the results loaded into ZAN. The database is accessible through a webbased interface and the promoter sequences of original list of upregulated genes can be searched for common binding sites of transcription factors. We are in the process of adding phylogenetic footprinting information to ZAN and integrating motif discovery algorithms.

A genome-wide map of human chromatin domains

Thurman R, Wang H, Sabo PJ, Koch C, Dunham I, Stamatoyannopoulos JA

Dept. of Genome Sciences, University of Washington, Seattle, WA; Wellcome Trust Sanger Institute, Cambridge, UK

Recent data suggest that the human genome is partitioned into higher-order functional domains characterized by transitions in chromatin modification state, accessibility, and higher-order structure. Domain boundary elements containing the CTCF regulatory factor have been proposed to play an important role in this organization. To systematically delineate chromatin domains across the human genome, we applied a wavelet segmentation approach to partition coordinated genome-wide maps of twenty histone modifications, CTCF binding, and DNaseI sensitivity in primary human lymphoid cells. We find that abrupt transitions in chromatin state are a regular feature of the human genome landscape, and that transition regions are characterized by a unique chromatin modification and accessibility signature. Gene clusters appear to evince two basic categories of chromatin domain organization: a monolithic architecture in which all genes within a cluster are unified within the same higher-order chromatin state; and a fractured pattern in which gene-rich regions are partitioned into a sequence of sub-domains with alternating chromatin states. The domain-level patterning of many histone modifications is surprisingly redundant, enabling us to define a minimal core set of modifications sufficient to accurately capture human chromatin domain structure when combined with CTCF and DNaseI sensitivity.

Exploiting Clinical Constraints for Robust Feature Subset Selection from Microarray Data on Breast Cancer

Yasser alSafadi, Angel Janevski, Nilanjana Banerjee, and Vinay Varadan

BMI Department, Philips Research NA, 345 Scarborough Road, B. Manor, NY10510

High-throughput molecular measurements are increasingly subjected to analyses that look into clinical applications based on patterns in the datasets. One application is a selection of feature subsets (say from expression data) that can be used as diagnostic patterns for a particular category of patients. For example, van't Veer et al. detected a signature of 70 genes that selects for breast cancer patients that would benefit from adjuvant therapy.

We have previously addressed the issue of discovering diagnostic patterns, using a Genetic Algorithm (GA)-based tool for feature subset selection based on a Support Vector Machine (SVM) for classification. Using our method, we are able to identify a large number of feature subsets that classify the learning data with relatively high accuracy. Unfortunately, given the wide range of validation errors for these patterns, we conclude that many of the high-performance feature subsets retrieved by our GA (and many other powerful search tools) are spurious patterns with limited robustness. In order to avoid this form of overfitting, we extended the analysis of expression microarray data to include a constraint in the form of a clinical profile for each sample – a clinical constraint. For clinical profiles, we utilized three widely-accepted clinical prognostic indices based on patient age, tumor size, hormonal and histopathological parameters: Nottingham Prognosis Index (NPI), National Institutes of Health Consensus (NCI), and the St. Gallen Consensus Conference (STG). Furthermore, we used a “virtual index” that Eden et al. (EDN) developed using the same dataset. We utilized the clinical constraints to additionally characterize the misclassifications by the SVM and identify cases where the SVM opts for a feature subset that is not in agreement with the otherwise conservative prognostic index.

The application of clinical constraints has a detectable effect on the overall output of the GA tool. Each GA execution is characterized by the distribution of the frequency of learning errors over all feature subsets when classifying 77 cases and validation errors on additional 19 cases unseen by the GA. The application of the three clinical indices (STG, NPI and EDN) reduced the range of the validation errors from a maximum of 52.6% to 31.6%, with most of the validation errors being around 18%.

Our results indicate that adding clinical constraints to a subset discovery tool improves performance by significantly reducing overfitting of the learning data thus avoiding discovery of spurious patterns. In addition, we offer evidence that our clinical constraints method does not over-constrain the search process from discovering useful feature subsets.

A network approach to literature management

Joachim von Eichborn^a, Jose Duarte^a, Henning Stehr^a, Michael Lappe^a

^aMax Planck Institute for Molecular Genetics, Ihnestr. 63 - 73, 14195 Berlin, Germany

Literature management has increasingly become a complex issue for the modern researcher. The explosion in publications in the recent years has made matters worse while new open access publishing schemes (PLoS, BioMed Central) add to the surge of information. Keeping track of what is being published in a field of interest presents a challenge for many researchers: staying up to date with the latest developments vs. advancing own research. We present a complete pipeline for literature management, including a network representation of publications stored in a relational database for convenient querying.

The literature is a (scale free) network. We represent it as a graph $G=(V,E)$ where papers are nodes (V) and citations are edges (E). This graphic representation of the network enables the quick visualization of the dependencies between the papers (citations). Features as clusters of related papers or important papers in our network (highly cited) can be easily observed.

Predicting functional transcription factor binding through alignment-free, cross-species analysis of promoter affinity

Lucas D. Ward¹ and Harmen J. Bussemaker^{1,2}

¹Department of Biological Sciences, Columbia University, 1212 Amsterdam Ave., New York, NY 10027

²Center for Computational Biology and Bioinformatics, Columbia University, 1130 St. Nicholas Ave., New York, NY 10032

The availability of large-scale functional genomics and comparative sequence data has led to the proliferation of techniques for detecting *cis*-regulatory elements and deciphering their combinatorial logic. For the most part, current methods for analyzing mRNA expression and ChIP-chip data divide the genome into target and non-target genes. Phylogenetic footprinting approaches rely on *ad hoc* local alignment of non-coding DNA or RNA sequence. Both approaches suffer from the inability to model and detect suboptimal binding sites, which have recently been shown to be abundant and functional (Tanay, 2006).

Using *MatrixREDUCE* (Foat et al., 2006), which is based on a biophysical model for protein-DNA interaction, together with a library of position-specific affinity matrices (PSAMs), we predicted the genome-wide profile of (relative) promoter occupancy by each of a large set of transcription factors across multiple yeast species. For each factor, we then used the *minimum* predicted occupancy across each set of orthologous promoters as a predictor for the amount of *functional* binding. We demonstrate that this affinity-based, alignment-independent use of evolutionary conservation gives rise to significant improvement in predicting expression from sequence compared to the use of a single species. We also find that the profile of conserved binding affinity across promoters is significantly correlated with nucleosome depletion for many factors. The same profile allows us to detect statistically significant association with specific functional categories using T-profiler (Boorsma et al, 2005) for many factors, even when no such association can be found using a single-species. We propose novel biological functions for several factors.

Using the same general approach, we compute for each pair of transcription factors the *affinity co-conservation*, defined as the degree of correlation between their genome-wide functional affinity profiles beyond what is expected for randomized promoter sequences. The interactions thus predicted are highly enriched for known physical or functional interactions between factors. When the same approach is repeated using single-species affinity profiles, no such enrichment is detected. Organizing the co-conserved factor pairs into a network reveals several functionally coherent sub-structures.

Identification and characterization of cell type-specific and ubiquitous chromatin regulatory structures in the human genome

Hualin Xi, Hennady P. Shulha, Jane M. Lin, Teresa R. Vales, Yutao Fu, David M. Bodine, Ronald D.G. McKay, Josh G. Chenoweth, Paul J. Tesar, Terrence S. Furey, Bing Ren, Gregory E. Crawford and Zhiping Weng

The identification of regulatory elements from different cell types is necessary for understanding the mechanisms controlling cell type-specific and housekeeping gene expression. Mapping DNaseI hypersensitive (HS) sites is an accurate method for identifying the location of functional regulatory elements. We have used a high throughput method, called DNase-chip, to identify 3904 DNaseI HS sites from six cell types across 1% of the human genome. A significant number (22%) of DNaseI HS sites from each cell type are ubiquitously present among all cell types studied.

Surprisingly, nearly all of these ubiquitous DNaseI HS sites correspond to either promoters or insulator elements: 86% of them are located near annotated transcription start sites (TSS) and 10% are bound by CTCF, a protein with known enhancer blocking insulator activity. We also identified a large number of DNaseI HS sites that are cell type-specific (only present in one cell type); these regions are enriched for enhancer elements and correlate with cell type-specific gene expression as well as cell type-specific histone modifications. Finally, we find that approximately 8% of the genome overlaps a DNaseI HS site in at least one the six cell lines studied, indicating that a significant percentage of the genome is potentially functional.

Regulation of Gene Expression by Promoter Associated Short RNAs.

Radharani Duttagupta^{1,2}, Aaron T. Willingham^{1,2}, Erica Dumais¹, Ian Bell¹, Jorg Drenkow¹, Jackie Dumais¹, Philipp Kapranov¹, Jill Cheng¹, Sujit Dike¹, Thomas R. Gingeras¹

¹. Affymetrix, Inc. Affymetrix Laboratory, Santa Clara, CA, 95051

². These authors contributed equally to this work.

Unbiased transcriptional profiling of the human genome suggests that a significant portion of pervasive transcription in mammals gives rise to RNA species shorter than 200nt (*Science*, Jun 8, 2007, v.316 p.1484). Strikingly, sequence conservation analysis and genomic mapping have revealed that a subclass of these short RNAs are enriched at the 5' ends of genes and have therefore been classified as promoter associated short RNAs (PASRs) (see abstract by P. Kapranov). More recently, these studies have been extended to the mouse genome which yielded similar findings. Furthermore, syntenic conservation of this class of short RNAs in mouse and human and the correlation of the density of PASRs with the steady-state mRNA levels of the associated genes suggest a putative functional role of these RNAs in the regulation of gene expression.

In order to ascertain the cellular roles of PASRs, we have systematically perturbed a collection of these short RNAs in human and mouse cells --utilizing complimentary approaches that involve (a) exogenous "overexpression" of synthetic RNA oligonucleotides representing PASR sequences or (b) strand-specific siRNA mediated knockdown of PASRs. We identified loci where modulation of PASR levels results in modest but statistically significant changes in the expression of the associated genes. Depending on the locus, this change appears to be either repressive or activating suggesting that PASRs participate in both negative and positive mechanisms to regulate gene expression. Furthermore we find that multiple PASRs at a given loci may have similar effects on downstream gene expression suggesting that they may function in a concerted manner. The cellular consequences of these changes and the underlying mechanism(s) mediating them are currently being explored.

Study of Temporal Effects of Intermittent- and Sustained Hypoxia on Gene Expression Patterns in Rat Lungs using a Systems Biology Approach

Wei Wu¹, Nilesh B. Dave¹, Patrick J. Strollo¹, Elizabeta Kovkarova-Naumovski¹, Stefan Rytter¹, Augustine M. K. Choi¹, David Gozal², Naftali Kaminski¹

¹*Division of Pulmonary, Allergy and Critical Care Medicine, School of Medicine, University of Pittsburgh, Pittsburgh, PA, USA*

²*Division of Pediatric Sleep Medicine, Department of Pediatrics, Kosair Children's Hospital Research Institute, University of Louisville, Louisville, Kentucky, USA*

Introduction: The molecular and genetic networks underlying the lung responses to intermittent- (IH) or sustained hypoxia (SH) are not fully understood. In this work we employed a systems biology approach to study temporal effects of IH and SH on gene expression patterns in rat lungs.

Methods: Rats were exposed to room air, IH or SH for 1, 3, 7, 14 and 30 days. Gene expression profiles were obtained using Codelink arrays. To analyze these temporal data, we developed a nonparametric statistical procedure, *LRSA*, to estimate genes differentially expressed (DE) at different time points, to detect cyclically regulated genes, and to reveal temporal expression patterns of DE genes. There are two steps involved in *LRSA*: i) the time course expression data were first fitted using a nonparametric local regression smoothing model; the bandwidth parameter in the model was estimated using the leave-one-out cross-validation. Then simultaneous confidence bands were calculated for estimation of DE genes. ii) A spectral clustering algorithm was employed to identify temporal patterns of DE genes. Finally, gene ontology (GO) analysis, network and functional physiological analyses were performed to identify significant biological processes and to propose regulatory networks that underlie the lung response to IH and SH. The temporal changes in the expression of 6 genes in IH and 5 genes in SH are confirmed by real-time PCR.

Results: Using the *LRSA* procedure, we found that IH and SH induce two distinct sets of genes in rat lungs which have little in overlap. These genes exhibit different temporal expression patterns. GO analysis reveals that IH-induced genes are mostly involved in transport, homeostatic, and neurological processes, steroid hormone receptor activity, while that SH-induced genes participate mostly in anatomical structure development and immune responses. IH-activated networks suggested a role for cross-talk between steroid hormone receptors and several other proteins which play key roles in hypoxia. SH-activated networks were indicative of vascular remodeling and pulmonary hypertension.

Conclusions: Our results provide a detailed view of the distinct temporal changes in gene expression induced by IH or SH in the rat lung, and demonstrate the discovery potential in applying systems biology approaches to understanding the mechanisms that underlie the lung response to hypoxia.

Coevolutionary Networks of Splicing cis-Regulatory Elements

Xinshu Xiao¹, Zefeng Wang², Minyoung Jang¹, and Christopher B. Burge¹

¹Department of Biology, Massachusetts Institute of Technology, Cambridge, MA, USA

²Department of Pharmacology, University of North Carolina, Chapel Hill, NC, USA

Accurate and efficient splicing of eukaryotic pre-mRNAs requires coordinated recognition by *trans*-acting factors of a complex array of *cis*-acting RNA elements, including 5' and 3' splice sites (5'ss, 3'ss) and exonic and intronic splicing enhancers and silencers (ESEs, ESSs, ISEs, and ISSs). To better understand this coordinated recognition, we studied the degree to which correlated changes have occurred between different classes of regulatory elements during mammalian evolution using whole-genome alignments of human, mouse and other genomes. The coevolutionary relationships among these elements were derived using a Bayesian network inference method. Cross-exon but not cross-intron interactions between the 5'ss and 3'ss were observed in human/mouse, indicating that the exon is the primary unit of selection in mammals. Studied plants, fungi and invertebrates exhibited exclusively cross-intron interactions, suggesting that intron definition drives evolution in these organisms. In mammals, 5'ss strength and the strength of several classes of ESSs coevolved in a correlated way, while specific ESEs, including motifs associated with hTra2, SRp55 and SRp20, evolve in a compensatory manner relative to the 5'ss and 3'ss. No interactions between specific ESS or ESE motifs were observed, suggesting that elements bound by different factors are not commonly interchangeable. Genes subject to recent positive selection exhibited a higher fraction of ESE disruption, and of 5'ss changes between human and chimp, suggesting that selection for change in the encoded protein tends to be disruptive to splicing regulation.

Comparative Genomic Motif Detection via Multi-Resolution Phylogenetic Shadowing

Eric P. Xing and Pradipta R. Ray

School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213

Recent developments in genomics, computational and molecular biology and population genetics have led to a convergence of interest in Regulatory Evolution—the evolutionary aspects of gene expression regulation. Identifying transcription factor binding sites (i.e., motif) governing the temporal-spatial expression patterns of developmental genes across genomes of evolutionarily related species plays an important role in modeling the evolution of transcriptional regulatory networks. Conventional phylogenetic shadowing methods for such tasks model the evolution of genomes only at nucleotide (nt) resolution; and assume that the nt content from all taxa in aligned sites are orthologous, and sequence functions are conserved across aligned sites. These assumptions are often violated in short, degenerate functional elements such as motifs [1], and can lead to a large number of false-positive predictions in practice. Hence comparative genomic motif-finding remains a difficult challenge, especially in metazoan species where the *cis*-regulatory modules (CRMs) can be long and divergent. To address issues encountered during comparative genomic motif finding, such as presence of misidentified orthology due to imperfect alignment, incomplete orthology in aligned sites due to taxon-specific functional loss, and co-evolution of multiple sites as a unit, we propose a new method: Conditional Shadowing via Multi-resolution Evolutionary Trees, or CSMET, which employs a context-dependent probabilistic graphical model that allows a single column in a multiple alignment to be modeled by multiple evolutionary trees conditioning on the functional specifications of each row (Fig. 1). The functional specifications themselves in turn result from a phylogeny which models the evolution not of nucleotides, but of the overall functionality (e.g., functional retention or loss) of the aligned sequence segments over lineages. Combining with a hidden Markov model (HMM) that autocorrelates evolutionary rates on successive sites in the genome, CSMET offers a way to take into consideration adaptive evolution of regulatory circuitry in motif detection, which is believed to have caused site evolution to be non-*iid* and non-isotropic in *cis*-regulatory sequences across taxa [1]. Our model shares a similar objective as that of several recent comparative genomic motif detection algorithms, such as EMnEm, Footer, MONKEY, PhyloGibbs, PhyME and Shadower, but the inclusion of multi-resolution phylogenies at both sequence and functional level makes CSMET able to detect poorly conserved motifs with higher accuracy (Fig. 2), and the estimated parameters offer insights into the evolutionary dynamics at multiple levels. We compared CSMET with MoMP—an implementation of the Phylogenetic HMM (which underlies EMnEm and PhyloGibbs) tailored for supervised motif detection, and with rVISTA and MAST (which uses an HMM on single species) on both simulated data, and a carefully curated *Drosophila* dataset containing 14 developmental CRMs for 11 aligned *Drosophila* species. Annotations for motif occurrences in *D. melanogaster* of 5 gap-gene TFs – Bicoid, Caudal, Hunchback, Kruppel and Knirps – were obtained from the literature. We found that CSMET outperforms the other methods on both synthetic and real data, and identifies a number of previously unknown occurrences of motifs within and near the study CRMs. A whole genome scan of the 5 gap-gene motifs, and *in vivo* functional analysis of the newly identified motifs, is currently underway.

1. Nipam H. Patel Michael Z. Ludwig, Casey Bergman and Martin Kreitman, Evidence for stabilizing selection in a eukaryotic enhancer element. *Nature*, 403:564–567,2000.

Hidden cellular responses revealed by integrating genetic and transcriptional data

Esti Yeger-Lotem^{1,2}, Laura Riva¹, Mikko Taipale², Bhupinder Bhullar², David R. Karger³, Susan Lindquist², and Ernest Fraenkel^{1,3}

¹ MIT Department of Biological Engineering, 77 Massachusetts Ave, Cambridge, Massachusetts 02139, USA. ² Whitehead Institute for Biomedical Research, Nine Cambridge Center, Cambridge, Massachusetts 02142, USA. ³ MIT Computer Science and Artificial Intelligence Laboratory, 32 Vassar St., Cambridge, Massachusetts 02139, USA

The cellular response to a perturbation is a complex process that involves the activation and inhibition of various pathways, requires a response-specific set of proteins, and often consists of a transcriptional response. To date, efforts to identify components of these responses have largely depended on mRNA profiling experiments and genetic libraries.

Intriguingly, when we compared the outputs of these two assays across numerous genetic and chemical perturbations we found that the differentially transcribed genes were consistently very different from the genes discovered in the genetic screens. This disparity suggests that independent analysis of each method's output reveals only a partial image of the cellular response, and that an integrative approach may substantially increase the effectiveness of the results obtained by mRNA profiling experiments and genetic screens. The disparity between the two types of data has not been rigorously studied and exploited in the literature.

We developed a computational approach that integrates genetic and transcriptional data with known molecular interactions to reveal a broader image of cellular response to perturbation. The algorithm identifies a network of interactions relating the genes and proteins identified in the assays with many other genes and proteins that were not detected by these assays, and it ranks the components of this network by their importance to the response.

We assessed our approach using published data for 106 genetic and chemical perturbations. Although in most cases only limited genetic data were available, the approach detected relevant perturbed pathways for more than 75% of the perturbations.

To further evaluate our method, we conducted an mRNA profiling experiment and a genetic over-expression screen in yeast cells exposed to radicicol, a known inhibitor of Hsp90. Although the data from each assay showed no obvious connection to Hsp90, our approach detected perturbed pathways known to rely on Hsp90 and provided a mechanistic explanation for the observed data. The output of our approach may be used to uncover novel drug targets and reveal general mechanisms for cellular adaptation.

Formal Analysis of Piecewise Affine Systems under Parameter Uncertainty with Application to Gene Networks

Boyan Yordanov¹, and Calin Belta²

¹Biomedical Engineering, ²Bioinformatics, Boston University, Boston MA, USA

Discrete-time continuous-space Piecewise Affine (PWA) systems are widely used as models of gene networks because they can approximate nonlinear dynamics with arbitrary accuracy and there exist methods for their identification from experimental data. In such models, the states correspond to concentrations of species (mRNA, proteins or small molecules). The parameters for the dynamics of real gene networks rarely remain fixed so, to build more general models, parameters might be allowed to vary.

We are interested in studying properties of the trajectories of such models, expressed as temporal and logical statements over regions of the state space, corresponding to different concentrations of the species and therefore different phenotypes. Specifically, we attempt to find regions of initial conditions for the system, such that initializing trajectories there guarantees the satisfaction of a property.

We illustrate our method by analyzing a PWA model of a two-gene network (see Poster). Under fixed parameters, this system can be shown to behave as a bi-stable “toggle switch” with two modes (only gene 1 or gene 2 activated).

For this example, we find regions of state space where the system should be initialized so one of the genes becomes activated and remains activated for all future times (regions of attraction for the two stable equilibria). We show that under fixed parameters, initial conditions completely determine which mode of the system trajectories will settle in, but as uncertainty is introduced there might exist initial conditions for which the behavior of the system is uncertain. Furthermore, we show that under high uncertainty it is possible to lose one of the attracting regions completely (the system is no longer bistable).

Analysis for property “gene 1 is expressed in concentrations higher than a certain threshold (see Poster), while gene 2 is expressed in concentrations lower than a threshold”. If the system is initialized in a region shown in green all trajectories are guaranteed to satisfy the property, while if initialized in a region shown in red no trajectories will satisfy the property. No guarantees can be made about trajectories originating in regions shown in white. The model is analyzed under fixed parameters (A), with 1% uncertainty (B) or with 5% uncertainty (C).

Quantitative Analysis of the Relationship Between Gene Expression Variation and Sequence Conservation

Yue Yun, Ayo Adesanya, Rob Mitra

*Genetics Department, Washington University in St. Louis
4444 Forest Park, MO 63108*

The relationship between sequence conservation and biological function is complex, as illustrate by the recent finding that the deletion of ultra-conserved elements in mouse produced no detectable phenotype. Our experiments are designed to systematically probe the biological function of conserved non-coding elements in the yeast *S. cerevisiae*. We are interested in the following questions: 1) What fraction of conserved sequences is involved in transcriptional regulation? 2) Are regulatory sequences located in conserved regions only? 3) Are there sequences important in transcriptional regulation that are not identified as transcription factor binding sites?

By evaluating expressional variations caused by single nucleotide substitutions in the GAL1 promoter , we aim to identify specific nucleotides responsible for the GAL1 transcriptional regulation. We generated an exhaustive single nucleotide mutation library of the GAL1 promoter, corresponding to position -1 to -633 before ATG. To detect any expression changes caused by a single nucleotide substitution, we applied a dual-reporter (CFP/YFP) assay followed by flow cytometry analysis, which should allow us to detect changes in expression of less than 5%. In this system, we constructed diploid yeast strains that express both CFP and YFP proteins, with CFP under the control of wild type GAL1 promoter and YFP under the control of mutated GAL1 promoter and integrated them on the same locus of homologous chromosomes. The expression of the cells is based on the ratio of CFP and YFP fluorescent level. We have randomly selected library clones bearing different mutations and tested their expression pattern in our system. We observed significant decrease of YFP/CFP ratio in clones with mutation in conserved GAL4 binding sites, consistent with the dominant role of transcription factor Gal4 in regulating Gal1 expression. It demonstrates that our technology platform is able to detect expression changes caused by a single point mutation in the GAL1 promoter. We are currently analyzing the entire GAL1 promoter mutation library to determine all functional nucleotides. This knowledge will advance our understanding of gene expression and regulation.

IEM: An algorithm for Iterative Enhancement of Motifs using Comparative Genomics Data

Erliang Zeng, Samidh Chaterjee, Kalai Mathee, and Giri Narasimhan

*School of Computing and Information Sciences, Florida International University,
Miami, Florida, USA*

Understanding gene regulation is a key step to investigating gene functions and their relationships. Many algorithms have been developed to discover transcription factor binding sites (TFBS); they are predominantly located in upstream regions of genes and contribute to transcription regulation if they are bound by a specific transcription factor. However, traditional methods focusing on finding motifs have shortcomings, which can be overcome by using comparative genomics data that is now increasingly available.

In this paper, we propose a new algorithm called IEM to refine motifs using comparative genomics data. It differs from other earlier approaches in that no attempt is made to perform de novo detection of motifs (although that would be easy to incorporate). Instead, comparative genomics data is used to “enhance” any given motif. These motifs may have been discovered by other computational methods, or may have been identified by laboratory techniques. Thus our method leverages the best-known motif discovery methods, or utilizes the (potentially incomplete) knowledge of previous studies while incorporating newly available comparative genomics data. Another challenge in motif prediction is to develop scoring functions that reflect biological significance. In this paper, we propose a metric (weighted conservation score) to measure such biological significance.

We show the effectiveness of our techniques with several data sets. Two sets of experiments were performed with comparative genomics data on five strains of *P. aeruginosa*. One set of experiments was performed with similar data on four species of yeast. The weighted conservation score proposed in this paper is an improvement over existing motif scores.

The research described here is significant for the following reasons. First, there is a clear need to reduce the number of false positives predicted by traditional tools. Second, our method can make use of partial information (on one or more binding sites), which may be available as a result of biological experiments. Third, with the availability of high throughput gene expression techniques like Microarrays and ChIP-Chip experiments, it is possible to get sets of coexpressed genes involved in the same metabolic pathway (and, therefore, potentially coregulated). Finally, our results show that the IEM algorithm has superior ability to overcome the shortcoming of previous methods and to effectively utilize any available comparative genomics data.

Comprehensive Prediction and Characterization of pre-mRNA Targets of the Tissue-Specific Alternative Splicing Factors Fox-1/2

Chaolin Zhang^{1,2,*}, Zuo Zhang^{1,*}, John Castle⁴, Shuying Sun^{1,3}, Jason Johnson⁴, Adrian R. Krainer¹, Michael Q. Zhang¹

¹ Cold Spring Harbor Laboratory, 1 Bungtown Road, Cold Spring Harbor, NY, 11724, USA; ² Department of Biomedical Engineering, ³ Department of Molecular & Cellular Biology, The State University of New York, Stony Brook, NY 11794, USA; ⁴ Rosetta Inpharmatics, LLC (Merck & Co.), Seattle, Washington 98109, USA; *Equal contribution

Alternative splicing (AS) is a key mechanism to expand proteomic diversity and accelerate gene evolution, but the molecular basis of tissue-specific AS regulation is only partially understood. Here we used the brain- and heart-specific splicing factors Fox-1 (A2BP1) and Fox-2 (RBM9), which specifically recognize the hexamer motif UGCAUG, as a model system to investigate the role of gene-expression regulation at the splicing level. We combined genome-wide comparisons of AS events in six species (from worm to human), splicing microarray data in a large panel of human tissues, and experimental validation. Our main conclusions are as follows: 1. The hexamer element is enriched in both downstream and upstream introns flanking cassette exons, but is under-represented in exons. Consistently, intronic elements are highly conserved, whereas exonic elements are generally not conserved. 2. There are ~1000 more human cassette exons with intronic elements than expected by chance; they include 567 cassette exons with intronic elements conserved in vertebrates and many have a conserved AS pattern. 3. To validate our predictions, we used shRNAs to knock down A2BP1 or RBM9 in HeLa cells, and also generated stable cell lines overexpressing A2BP1 or RBM9. We analyzed ~70 cassette exons with conserved elements by RT-PCR; more than half of the genes expressed in HeLa cells gave a switch between exon inclusion and skipping upon overexpression or knockdown of A2BP1/RBM9. 4. Previous results suggested that downstream elements are enriched in brain-specific exons. Using splicing microarrays to study the correlation between A2BP1 expression and target-exon splicing, as well as RT-PCR, we found that downstream elements strongly enhance exon inclusion, whereas exonic elements and upstream intronic elements are repressive. 5. Gene-ontology analysis of the targets suggests specific roles of the proteins in neuromuscular functions and RNA splicing. Importantly, a number of targets have been implicated in human diseases.

A Quadratic Programming Based Motif Finding Algorithm

Yue Zhao and Gary Stormo

*Department of Genetics,
Washington University Medical School,
St Louis, MO 63110, USA*

Identification of transcription factor (TF) binding sites in unaligned DNA sequences is an important step in understanding of the regulatory networks that control gene expression. A related problem is the choice of models used to describe TF specificity. Currently, the most popular models of TF specificity are information theory based log-odds model and Match algorithm used by TRANSFAC database. A quadratic programming based model of TF specificity was recently introduced by two groups. This model correctly takes saturation in the TF/DNA binding probability into account and has been shown to outperform both log-odds model and Match algorithm. Here, we show that performance of a motif finding algorithm, CONSENSUS, is improved by the use of quadratic programming based model of TF specificity. We benchmark the new algorithm on real motifs and find its performance compares favorably to current popular methods.

Characterization of short RNA clones in the Sfmbt2 locus

Grace X.Y. Zheng^{1,2}, Joel R. Neilson¹, Christopher B. Burge^{2,3}, and Phillip A. Sharp^{1,3}

¹Center for Cancer Research; ²Computational and Systems Biology; ³Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

Short RNAs are known to control gene expression at several different levels in organisms ranging from yeast to plants to mammals. To understand the role of short RNAs in regulating the development of mammalian cells, we have comprehensively profiled short RNA expression from several distinct stages of murine T-lymphocyte development. While the majority of the clones correspond to known miRNAs, novel short RNAs constitute 10% of the cloned sequences. We mapped novel short RNA sequences to their physical location in the mouse genome. Several highly related novel sequences clustered in an intron of the uncharacterized gene *Sfmbt2* on mouse chromosome 2. The dynamics of *Sfmbt2* expression and cloning frequency of members of the cluster throughout T lymphocyte development are highly correlated, suggesting that the novel RNA sequences derived from the cluster are controlled by the *Sfmbt2* promoter. The *Sfmbt2* short RNAs are similar to known miRNAs; most of them can fold into miRNA-like hairpins, and the processing of these species is dependent on Dicer. However, while highly related miRNAs are well conserved at the 5' end, the *Sfmbt2* short RNAs share a 12-nt consensus motif spanning nucleotides 6 to 17 of most of the sequences. We are actively investigating whether these species are incorporated in the canonical RNAi pathway, and whether they play a role in the maturation of T lymphocytes.

Index of Presentations

An integrative strategy for characterizing transcriptional regulatory networks involved in *Drosophila* embryonic mesoderm development..... 1
Anton Aboukhalil^{1,5}, Brian W. Busse⁶, Savina A. Jaeger¹, Aditi Singhanian⁶, Michael F. Berger^{1,4}, Stephen S. Gisselbrecht¹, Alan M. Michelson⁶, Martha L. Bulyk¹⁻⁴

The oncogene IKBKE is overexpressed in only two breast cancer subtypes: the Her2+ subtype with lymphocytic infiltration and the basal-like subtype 2
Gabriela Alexe*, Nilay Sethi*, Lyndsay Harris, Shridar Ganesan and Gyan Bhanot

Combining the prediction of nucleosome-binding sequences and transcription factor binding sites. 3
Abha Singh Bais, Ho-Ryun Chung, Dennis Kostka, Hugues Richard and Martin Vingron

Reconstructing dynamic regulatory maps 4
Jason Ernst¹, Oded Vainas², Christopher T. Harbison³, Itamar Simon², Zoltan N. Oltvai⁴, Naftali Kaminski⁵ and Ziv Bar-Joseph¹

MicroRNAs 5
David Bartel

Small but essential: let-7d miRNA in a transcriptional regulatory network that determines the fate of epithelial cells 6
Hanadie Yousef, David L. Corcoran, Daniel Handley, Kusum Pandit, Isidore Rigoutsos, Oliver Eickelberg, Naftali Kaminski, Panayiotis V. Benos

Familial Binding Profiles Made Easy 7
Shaun Mahony, Philip E. Auron, Panayiotis V. Benos

High-Resolution Characterization of Mammalian Transcription Factor Binding Specificities Using Protein Binding Microarrays 8
Michael F. Berger^{1,4}, Gwenael Badis-Breard⁶, Andrew R. Gehrke¹, Shaheynoor Talukder⁶, Anthony A. Philippakis^{1,3,4}, Savina A. Jaeger¹, Esther Chan⁶, Genita Metzler¹, Anastasia Vedenko¹, Olga Botvinnik¹, Wen Zhang¹, Hanna Kuznetsov¹, Chi-Fong Wang¹, David Coburn⁶, Quaid D. Morris^{6,8}, Timothy R. Hughes^{6,9}, and Martha L. Bulyk¹⁻⁵

RNA-mediated silencing in Archaea 9
David L Bernick, and Todd M Lowe

ENCODE: Understanding our genome 10
Ewan Birney

A Method for Inferring Extracellular Environment from Gene Expression Profiles Using Metabolic Flux Balance Models 11
Aaron Brandes, Desmond S. Lun, Jeremy Zucker, and James E. Galagan

Characterization of genome-wide DNA binding by the replication initiator and transcription factor DnaA in response to perturbation of replication 12
Adam M. Breier and Alan D. Grossman

Selective optimization of codon usage detected by information theory based codon bias measurements	13
<u>Zehua Chen</u> , and Jeffrey Chuang	
New genomic tools for reading, writing & regulation	14
<u>George Church</u>	
A nonlinear dynamical model to infer transcriptional regulatory networks from time dependent expression data	15
<u>Adriana Climescu-Haulica</u> ¹ and Michelle Quirk ²	
Riboswitches, RNA conformational switches and prokaryotic gene regulation...	16
Eva Freyhult ^a , Vincent Moulton ^b , <u>Peter Clote</u> ^c	
Biology by Design: Integrating Synthetic Biology and Systems Biology	17
<u>J.J. Collins</u>	
Crystal Structure and phylogenetic data to infer the binding motifs for new protein	18
<u>Francesca Cordero</u> and Gary D. Stormo	
Drug target prediction via Lasso regression analysis of mRNA expression compendia.....	19
<u>Elissa J. Cosgrove</u> ^{1*} , Yingchun Zhou ^{2*} , Eric D. Kolaczyk ² , Timothy S. Gardner ¹	
Zebrafish Gene Map and Microarray Annotations	20
<u>Anthony D. DiBiase</u> , Anhua Song, Fabrizio Ferre, Gerhard Weber, David Langenau, Leonard I. Zon, Yi Zhou	
Probing microRNA regulatory functions with microRNA ‘sponge’ inhibitors	21
<u>Margaret S. Ebert</u> , Madhu S. Kumar, Nancy Guillen, Phillip A. Sharp, Tyler Jacks, Douglas A. Lauffenburger.	
Function and Evolution of Regions Bound by Drosophila Transcription Factors	22
<u>Michael B. Eisen</u>	
Bayesian Meta-analysis for Identification of Cell cycle-regulated Genes in Fission Yeast.....	23
<u>Xiaodan Fan</u> ¹ , Saumyadipta Pyne ² , and Jun S. Liu ¹	
Regulatory Hubs: Classification and Conservation	24
<u>Andrew Fox</u> , Donna K. Slonim	
Dynamics of replication-independent histone turnover in budding yeast.....	25
<u>Nir Friedman</u>	
High-resolution mapping and characterization of open chromatin across the human genome	26
<u>Terrence S. Furey</u> ¹ , Alan P. Boyle ¹ , Sean Davis ² , Hennady P. Shulha ³ , Paul Meltzer ² , Zhiping Weng ³ , Gregory E. Crawford ¹	

Pattern-based analysis of interactome networks in <i>C. elegans</i>	27
Sira Sriswasdi ¹ , Kesheng Liu ¹ , Thomas Martinez ¹ , Carmel Mercado ¹ , Pengyu Hong ² , and <u>Hui Ge</u> ¹	
A fast and simple motif finder utilizing rank lists of sequences	28
<u>Stoyan Georgiev</u> , Karthik Jayasurya, Sayan Mukherjee, Uwe Ohler	
Genome-wide identification of active regulatory elements in the human genome by FAIRE	29
<u>Paul Giresi</u> ¹ , Vishy Iyer ² , and Jason Lieb ¹	
Position-dependent motif characterization using non-negative matrix factorization	30
Lucie N. Hutchins ¹ , Erik L. McCarthy ² , Sean Murphy ¹ , Priyam Singh ³ , and <u>Joel H. Graber</u> ^{1,3}	
Characterization of TBP binding to TATA boxes and its relationship to structural properties of TATA boxes	31
Hana Faiger, Marina Ivanchenko, Ilana Cohen, and <u>Tali E. Haran</u>	
Association Analysis of Gene Functions and Upstream Sequence Conservation in the Human Genome	32
<u>Weichun Huang</u> ¹ , Leping Li ² , Gabor Marth ¹ , and Bruce S Weir ³	
Protein Network Comparative Genomics	33
<u>Trey Ideker</u>	
Systematic identification of mammalian regulatory motifs' target genes and functions	34
Jason B. Warner ^{1,6,7} , Anthony A. Philippakis ^{1,3,4,6} , <u>Savina A. Jaeger</u> ^{1,6} , Fangxue Sherry He ^{5,8} , Jolinta Lin ^{5,6} , and Martha L. Bulyk ^{2,3,4,6}	
Cycles in the transcription regulatory network of <i>S. Cerevisiae</i>	35
<u>Jieun Jeong</u> and Piotr Berman	
New Classes of RNAs Potentially Involved in Regulation of Gene Expression ..	36
<u>Philipp Kapranov</u> ¹ , Juergen Schmitz ² , Katalin Fejes-Toth ³ , Radharani Duttagupta ¹ , Aaron T. Willingham ¹ , Timofey Rozhdestvensky ² , Jill Cheng ¹ , Richard Reinhardt ⁴ , Sujit Dike ¹ , Greg Hanon ³ , Juergen Brosius ² , Thomas R. Gingeras ¹	
Systematic Discovery and Characterization of Fly microRNAs using 12 <i>Drosophila</i> Genomes	37
<u>Pouya Kheradpour</u> [*] , Alexander Stark [*] , Leopold Parts, Julius Brennecke, Emily Hodges., Gregory Hannon and Manolis Kellis	
A Polycomb Switch During Muscle Differentiation	38
<u>Roshan M. Kumar</u> ¹ , Stuart S. Levine ¹ , and Richard A. Young ^{1,2,†}	
A predictive model of the oxygen and heme regulatory network in yeast	39
<u>Anshul Kundaie</u> ^a , Xiantong Xin ^b , Changgui Lan ^b , Steve Lianoglou ^a , Mei Zhou ^b , Christina Leslie ^{*c} Li Zhang ^{*b} ,	

Dissecting transcriptional regulation in human B lymphocytes using an integrative, context-specific network 40
Celine Lefebvre¹, Wei Keat Lim¹, Kai Wang¹, Pavel Sumazin¹, Katia Basso², Riccardo Dalla-Favera² and Andrea Califano¹

Whole-genome analysis of dorsal-ventral patterning thresholds in the *Drosophila* embryo 41
Julia Zeitlinger, Rob Zinzen, Dmitri Papatsenko, Joung-Woo Hong, Manolis Kellis, Alexander Stark, Rick Young, and Mike Levine

An integrated approach to clustering expression trajectories in differentiating stem cells 42
Shenggang Li¹ Miguel Andrade² David Sanko³

Assembling a compendium of druggable apoptosis regulators..... 43
Wei Keat Lim and Andrea Califano

DIAL: An algorithm to detect RNA 3D motifs 44
Fabrizio Ferre*, Yann Ponty†, William A. Lorenz†, Peter Clote†

Peromyscus maniculatus and *Mus musculus* Synteny Map Reveals Divergent, Patterns of Breakpoint Reuse in Rodentia 45
E.E. Mlynarski¹, C.J. Obergfell¹, W. Rens², P.C.M. O'Brien², C.M. Ramsdell^{3,4}, M.J. Dewey³, M.J. O'Neill¹, R.J. O'Neill¹

Temporal logic, not combinatorial logic, dominates the regulation of sugar catabolic genes in *E. coli*..... 46
Ilaria Mogno (mogno@bu.edu) and Timothy Gardner (tgardner@bu.edu)

Form, function, and information processing in small stochastic biological networks 47
Andrew Mugler¹, Etay Ziv², Ilya Nemenman³, and Chris H. Wiggins⁴

Generalized dynamic theory of alternative splicing regulation 48
Rajamanickam Murugan and Gabriel Kreiman

An All-or-None Response in Osmolyte-Dependent Gene Expression in *Saccharomyces cerevisiae* 49
Gregor Neuert and Alexander van Oudenaarden

An ultra-conserved centromere-specific retrovirus is an epigenetic determinant of centromere function in vertebrates..... 50
Rachel J. O'Neill, Ph.D.

Structure and Function of the Transcriptional Network Activated by the MAPK Hog1..... 51
Erin O'Shea*, Andrew Capaldi, Kristen Cook, Ying Liu, Tommy Kaplan, Naomi Habib, Aviv Regev, Nir Friedman

A generalized Karlin-Altschul statistics for modeling conserved cis-regulatory modules.....	52
<u>Kimmo Palin</u> ¹ , Jussi Taipale ² , and Esko Ukkonen ¹	
A mechanism for bistable gene expression regulation in yeast	53
<u>Saurabh Paliwal</u> and Andre Levchenko	
Efficient Discovery of Significant Co-Occurrences with an Application to Dyad Motifs Finding.....	54
<u>Cinzia Pizzi</u> ¹ , Esko Ukkonen ²	
LocalMove: A web server to compute best on-lattice fits for biopolymers	55
<u>Yann Ponty</u> , Peter Clote	
Evolution and selection in yeast promoters: analyzing the combined effect of diverse transcription factor binding sites	56
<u>Daniela Rajiman</u> ¹ , Amos Tanay ² and Ron Shamir ¹	
Inferring Interactions in Transcriptional Regulatory Networks.....	57
<u>Anil Raj</u> , Chris Wiggins	
Simulated Evolution of Transcription Regulatory Networks and Network Motifs	58
<u>Arvind Ramanathan</u> , Michelle Goodstein, Robbie Sedgewick, and Dannie Durand	
AhoSoft: P-value Computation for Heterotypic Clusters of Regulatory Motifs	59
<u>Mireille Regnier</u> ¹ , Valentina Boeva ^{1,2} , Zara Kirakosyan ³ , Vsevolod Makeev ^{4,5}	
Age, sex and genetic architecture of human gene expression	60
<u>Manuel A. Rivas</u> , Mark J. Daly, and Itsik Pe'er	
Inferring tissue specific activity of transcription factors from predicted DNA binding affinities to core promoters	61
<u>Helge G. Roeder</u> , Thomas Manke, Martin Vingron & Stefan A. Haas	
Thermodynamic-based approach for miRNA target prediction	62
Aleksy Y. Ogurtsov, Denis V. Karamyshev, Olga V. Matveeva, <u>Svetlana A. Shabalina</u>	
PhyloGibbs-MP: Motif-finding, module prediction and differential motif prediction by Gibbs sampling.....	63
<u>Rahul Siddharthan</u>	
Rapid chromatin structural phenotyping by Digital DNaseI.....	64
Dorschner MO, Sabo PJ, Kuehn MS, Sandstrom R, Weaver M, Neph S, Haydock A, Goldy J, Thurman RT, <u>Stamatoyannopoulos JA</u>	
Regulatory Motif Discovery and Target Prediction in 12 Drosophila Genomes ..	65
<u>Alexander Stark</u> [*] , Pouya Kheradpour [*] , Sushmita Roy, and Manolis Kellis	
A Mining Tool for Regulatory Elements in Zebrafish.....	66
Teplyuk, V ¹ , Hirsch, N. ² , Sagerstrom, C ² , <u>Straubhaar, J</u> ¹	

A genome-wide map of human chromatin domains.....	67
<u>Thurman R</u> , Wang H, Sabo PJ, Koch C, Dunham I, Stamatoyannopoulos JA	
Exploiting Clinical Constraints for Robust Feature Subset Selection from Microarray Data on Breast Cancer.....	68
Yasser alSafadi, Angel Janevski, Nilanjana Banerjee, and <u>Vinay Varadan</u>	
A network approach to literature management	69
<u>Joachim von Eichborn</u> ^a , Jose Duarte ^a , Henning Stehr ^a , Michael Lappe ^a	
Predicting functional transcription factor binding through alignment-free, cross- species analysis of promoter affinity	70
<u>Lucas D. Ward</u> ¹ and Harmen J. Bussemaker ^{1,2}	
Identification and characterization of cell type-specific and ubiquitous chromatin regulatory structures in the human genome.....	71
Hualin Xi, Hennady P. Shulha, Jane M. Lin, Teresa R. Vales, Yutao Fu, David M. Bodine, Ronald D.G. McKay, Josh G. Chenoweth, Paul J. Tesar, Terrence S. Furey, Bing Ren, Gregory E. Crawford and <u>Zhiping Weng</u>	
Regulation of Gene Expression by Promoter Associated Short RNAs.....	72
Radharani Duttgupta ^{1,2} , <u>Aarron T. Willingham</u> ^{1,2} , Erica Dumais ¹ , Ian Bell ¹ , Jorg Drenkow ¹ , Jackie Dumais ¹ , Philipp Kapranov ¹ , Jill Cheng ¹ , Sujit Dike ¹ , Thomas R. Gingeras ¹	
Study of Temporal Effects of Intermittent- and Sustained Hypoxia on Gene Expression Patterns in Rat Lungs using a Systems Biology Approach	73
<u>Wei Wu</u> , Nilesh B. Dave, Patrick J. Strollo, Elizabeta Kovkarova-Naumovski, Stefan Ryter, Augustine M. K. Choi, David Gozal, Naftali Kaminski	
Coevolutionary Networks of Splicing cis-Regulatory Elements.....	74
<u>Xinshu Xiao</u> ¹ , Zefeng Wang ² , Minyoung Jang ¹ , and Christopher B. Burge ¹	
Comparative Genomic Motif Detection via Multi-Resolution Phylogenetic Shadowing	75
<u>Eric P. Xing</u> and Pradipta R. Ray	
Hidden cellular responses revealed by integrating genetic and transcriptional data	76
<u>Esti Yeger-Lotem</u> ^{1,2} , Laura Riva ¹ , Mikko Taipale ² , Bhupinder Bhullar ² , David R. Karger ³ , Susan Lindquist ² , and Ernest Fraenkel ^{1,3}	
Formal Analysis of Piecewise Affine Systems under Parameter Uncertainty with Application to Gene Networks.....	77
<u>Boyan Yordanov</u> ¹ , and Calin Belta ²	
Quantitative Analysis of the Relationship Between Gene Expression Variation and Sequence Conservation.....	78
<u>Yue Yun</u> , Ayo Adesanya, Rob Mitra	

IEM: An algorithm for iterative enhancement of motifs using comparative genomics data..... 79
Erliang Zeng, Samidh Chaterjee, Kalai Mathee, and Giri Narasimhan

Comprehensive Prediction and Characterization of pre-mRNA Targets of the Tissue-Specific Alternative Splicing Factors Fox-1/2 80
Chaolin Zhang^{1,2,*}, Zuo Zhang^{1,*}, John Castle⁴, Shuying Sun^{1,3}, Jason Johnson⁴, Adrian R. Krainer¹, Michael Q. Zhang¹

A Quadratic Programming Based Motif Finding Algorithm..... 81
Yue Zhao and Gary Stormo

Characterization of short RNA clones in the Sfmbt2 locus 82
Grace X.Y. Zheng^{1,2}, Joel R. Neilson¹, Christopher B. Burge^{2,3}, and Phillip A. Sharp^{1,3}

Participants List

Anton Aboukhalil	Graduate Student, MIT/BWH	anton1@mit.edu
Gabriela Alexe	Computational Biologist, Broad Institute	galaxe@broad.mit.edu
Robert Altshuler	Graduate Student, MIT	rca@csail.mit.edu
Ido Amit	Postdoctoral Fellow, Broad Institute	iamit@broad.mit.edu
Daehyun Baek	Postdoctoral Fellow, Whitehead Institute, MIT	baek@wi.mit.edu
Abha Bais	PostDoc, Max Planck Institute for Molecular Genetics	bais@molgen.mpg.de
Nilanjana Banerjee	Senior Member Research Staff, Philips Research NA	nila.banerjee@philips.com
Ziv Bar-Joseph	Assistant Professor, Carnegie Mellon University	zivbj@cs.cmu.edu
M Inmaculada Barrasa	Postdoctoral Associate, UMass Medical School	inmaculada.barrasa@umassmed.edu
David Bartel	Professor, MIT/Whitehead Institute/HHMI	dbartel@wi.mit.edu
Michael Baym	Student, MIT	baym@mit.edu
George Bell	Senior Bioinformatics Scientist, Whitehead Institute, MIT	gbell@wi.mit.edu
Panayiotis Benos	Assistant Professor, University of Pittsburgh	benos@pitt.edu
Michael Berger	Graduate Student, Harvard University	berger@fas.harvard.edu
Nynke Van Berkum	Postdoctoral Associate, UMass Medical School	nynke.vanberkum@umassmed.edu
David Bernick	Graduate Student, University of California, Santa Cruz	dbernick@soe.ucsc.edu
Ewan Birney	Principal Researcher, Sanger Center/EBI	birney@ebi.ac.uk
Alan Boyle	Graduate Student, IGSP, Duke University	alan.boyle@duke.edu
Aaron Brandes	Computational Biologist, Broad Institute	abrandes@broad.mit.edu
Adam Breier	PostDoc, MIT	abreier@gmail.com
Michael Brent	Professor, Washington University in St. Louis	brent@cse.wustl.edu
Michael Brodsky	Assistant Professor, UMass Medical School	michael.brodsky@umassmed.edu
Harmen Bussemaker	Associate Professor, Columbia University	hjb2004@columbia.edu
David Byrne	Student, Boston Univeristy	dbyrne@bu.edu
Daniel Caffrey	Scientist, Pfizer	daniel.r.caffrey@pfizer.com
Sarah Calvo	Graduate Student, Broad Institute	scalvo@broad.mit.edu
Rogério Candeias	PhD Student, MIT	candeias@mit.edu
Yuxia Cao	Assistant Professor, Boston Univeristy	yxcao@bu.edu
Andrew Capaldi	PostDoc, Harvard University	capaldi@mcb.harvard.edu
Christopher Carr	Postdoctoral Associate, MIT	chrisc@mit.edu
Scott Carter	Graduate Student, MIT	scarter@mit.edu
Luis Carvalho	Student, Brown University	carvalho@dam.brown.edu
Mia Champion	Research Publications Coordinator, Broad Institute	champion@broad.mit.edu
Lingji Chen	Senior Research Engineer, Scientific Systems Company Inc	lingji.chen@sosci.com
Zehua Chen	PostDoc, Boston College	chenzp@bc.edu
Michael Chou	PhD Student, HMS	mchou@fas.harvard.edu
Jeffrey Chuang	Assistant Professor, Boston College	chuangj@bc.edu
Victor Chubukov	Graduate Student, UC San Francisco / UC Berkeley	victor.chubukov@gmail.com
Yakov Chudnovsky	PostDoc, Whitehead Institute, MIT	chudnovs@wi.mit.edu
George Church	Professor, MIT/Harvard	arep.med.harvard.edu/gmc/email.html
Adriana Climescu-Haulica	Researcher, Laboratoire Jean Kuntzmann	adriana.climescu@imag.fr
Peter Clote	Professor of Biology, Boston College	clote@bc.edu
James Collins	Professor, Boston Univeristy	jcollins@bu.edu
Francesca Cordero	PhD Student, Computer science, University of Torino	fcordero@di.unito.it
Elissa Cosgrove	Graduate Student, Boston Univeristy	ejburk@bu.edu
Stephen Davies	Assistant Professor, University of Toronto	stephen.davies@utoronto.ca
Anthony DiBiase	Bioinformaticist, Children's Hospital/HMS	adibiase@enders.tch.harvard.edu
Yang Ding	Graduate Student, Boston College	yding1985@yahoo.com.cn
Robin Dowell	Postdoctoral Fellow, CSAIL, MIT	rdd@mit.edu
Ines Drinnenberg	Research Assistant, Whitehead Institute, MIT	drinnenb@wi.mit.edu
Alexandra Dumitriu	Student, Boston Univeristy	adumitri@bu.edu
Radharani Duttagupta	Postdoctoral Fellow, Affymetrix	radharani_duttagupta@affymetrix.com
Margaret Ebert	Graduate Student, Center for Cancer Research, MIT	ebertms@mit.edu
Joachim Von Eichborn	Visiting Scholar, Boston College	joachim.eichborn@gmail.com
Michael Eisen	Professor, UC Berkeley / LBNL	mbeisen@lbl.gov
Tom Ellis	PostDoc, Boston Univeristy	ellis@bu.edu
Eleazar Eskin	Professor, UCLA	eeskin@cs.ucla.edu
Moriah Eustice	Research Tech II, Broad Institute	meustice@broad.mit.edu
Xiaodan Fan	Graduate Student, Harvard University	xfan@fas.harvard.edu
Miaoqing Fang	Student, MIT	fangm@mit.edu

Paola Favaretto	Senior Software Developer, MathWorks	paola.favaretto@mathworks.com
Fabrizio Ferre	Research Fellow, Children's Hospital/HMS	fferre@enders.tch.harvard.edu
Andrew Fox	Graduate Student, Tufts University CS	andrew.fox@tufts.edu
Garrett Frampton	Graduate Student, Whitehead, MIT	frampton@mit.edu
Hunter Fraser	Researcher, Broad Institute	hunter@broad.mit.edu
Nir Friedman	Professor, Hebrew University	nir@cs.huji.ac.il
Guilherme Fujiwara	Master student, CSAIL, MIT	issao@mit.edu
Terrence Furey	Assistant Professor, IGSP, Duke University	terry.furey@duke.edu
David Garcia	Graduate Student, Whitehead Institute, MIT	garcia@wi.mit.edu
Hui Ge	Whitehead Fellow, Whitehead Institute, MIT	hge@wi.mit.edu
Robin Ge	Bioinformatic analyst, Whitehead Institute, MIT	rge@wi.mit.edu
Stoyan Georgiev	PhD student, Duke University	sgg4@duke.edu
David Gifford	Professor, MIT	dkg@mit.edu
Paul Giresi	Graduate Student, University of North Carolina	paulg@email.unc.edu
Amir Globerson	Postdoctoral Associate, MIT	gamir@csail.mit.edu
Suresh Gopalan	Molecular Biology, Massachusetts General Hospital	gopalan@molbio.mgh.harvard.edu
Bryan Gorman	Graduate Student, HST, Harvard/MIT	gorman@mit.edu
Joseph Gormley	Director I/CB, MBI	jg@massbiomed.org
Joel Graber	Associate Staff Scientist, The Jackson Laboratory	joel.graber@jax.org
Andrew Grimson	PostDoc, Whitehead Institute, MIT	grimson@wi.mit.edu
Yuchun Guo	Graduate Student, CSBi, MIT	yguo@mit.edu
Katherine Gurdziel	Bioinformatics Analyst, BaRC Whitehead Institute	gurdziel@wi.mit.edu
Tali Haran	Researcher, Technion	bitali@tx.technion.ac.il
Brian Haynes	Doctoral Student, Washington University in St. Louis	bch2@cec.wustl.edu
Shao-shan Huang	Graduate Student, MIT	shuang@mit.edu
Weichun Huang	Research Associate, Boston College	whduke@gmail.com
Trey Ideker	Associate Professor, UCSD	trey@bioeng.ucsd.edu
Hideo Imamura	PostDoc, Boston College	imamurah@bc.edu
Sorin Istrail	Professor, Brown University	sorin@cs.brown.edu
Savina Jaeger	Postdoctoral Fellow, Brigham & Women's/HMS	sjaeger@rics.bwh.harvard.edu
Calvin Jan	Graduate Student, Whitehead Institute, MIT	cjan@wi.mit.edu
Angel Janevski	SMRS, Philips Research NA	angel.janevski@philips.com
Jieun Jeong	Graduate Student, Pennsylvania State University	jijeong@cse.psu.edu
Maura Kak	Computational Biologist, Broad Institute	maurakak@broad.mit.edu
Philipp Kapranov	Senior Scientist, Affymetrix	philipp_kapranov@affymetrix.com
Ulas Karaoz	PhD Student, Boston University	ulask@bu.edu
Uri Keich	Assistant Professor, Cornell University	keich@cs.cornell.edu
Manolis Kellis	Assistant Professor, MIT	manoli@mit.edu
Peter Kharchenko	PostDoc, Harvard Medical School	peter.kharchenko@post.harvard.edu
Pouya Kheradpour	Student, MIT	pouyak@mit.edu
Jinkuk Kim	Graduate Student, Whitehead Institute, MIT	jinkuk@mit.edu
Niels Klitgord	PhD student, Boston University	klitgord@gmail.com
Mladen Kolar	Student, Carnegie Mellon University	meh@cs.cmu.edu
Mark Kon	Professor, Boston University	mkon@bu.edu
Andrew Krueger	PhD Candidate, Boston University	akrueger@bu.edu
Roshan Kumar	Postdoctoral Fellow, Whitehead Institute, MIT	roshan@wi.mit.edu
Anshul Kundaje	Student, Columbia University	abk2001@columbia.edu
Sooraj Kuttykrishnan	Postdoctoral Employee, Washington University in St. Louis	sooraj1.k@gmail.com
Brian Lajoie	Bioinformatician, UMass Medical School	sharon.briggs@umassmed.edu
David Lapointe	Research Computing, UMass Medical School	david.lapointe@umassmed.edu
Carolyn Lawrence	PhD Candidate, Boston University	carolynl@bu.edu
Lee Lawton	Postdoctoral Fellow, Whitehead Institute, MIT	lawton@wi.mit.edu
Celine Lefebvre	Postdoctoral Researcher, Columbia University	lefebvre@c2b2.columbia.edu
Jessica Lehoczky	PostDoc, Boston College	jlehoczk@umich.edu
Ana Leite	PhD student, Broad Institute	apleite@broad.mit.edu
Mike Levine	Professor, UC Berkeley	mlevine@berkeley.edu
Fran Lewitter	Director of Bioinformatics, Whitehead Institute, MIT	lewitter@wi.mit.edu
Jialing Li	Graduate Student, MIT	jjaling3@mit.edu
Shenggang Li	PhD Student, University of Ottawa	sli365@science.uottawa.ca
Wei Keat Lim	Student, Columbia University	wl2131@columbia.edu
Mieszko Lis	PhD Candidate, MIT	mieszko@csail.mit.edu
Kesheng Liu	Postdoctoral Associate, Whitehead Institute, MIT	kliu@wi.mit.edu

Andy Lorenz	Postdoctoral Fellow, Boston College	lorenzwi@bc.edu
Kenzie MacIsaac	Graduate Student, MIT	macisaac@mit.edu
Shaun Mahony	Research Associate, MIT	mahony@mit.edu
Erik McCarthy	Student, University of Maine	erik.mccarthy@maine.edu
Scott McCuine	Technical Associate II, Whitehead Institute, MIT	mcuine@wi.mit.edu
Rebecca McDonald	Graduate Student, Massachusetts General Hospital	rmcdonal@fas.harvard.edu
Paul McGettigan	PhD student, University College Dublin	paul.mcgettigan@ucd.ie
Leonid Mirny	Associate Professor, MIT	leonid@mit.edu
Patrycja Missiuro	Graduate Student, Whitehead/CSAIL, MIT	patrycja@mit.edu
Elisabeth Mlynarski	Graduate Student, University of Connecticut	elizabeth.mlynarski@uconn.edu
Ilaria Mogno	PostDoc, Boston University	mogno@bu.edu
Andrew Mugler	Graduate Student, Physics Dept, Columbia University	ajm2121@columbia.edu
Rajamanickam Murugan	Research Fellow, Children's Hospital/HMS	rajamanickam.murugan@childrens.harvard.edu
Giri Narasimhan	Professor, Florida International University	giri@cs.fiu.edu
Natalia Naumova	Postdoctoral Associate, UMass Medical School	natalia.naumova@umassmed.edu
Gregor Neuert	Postdoctoral Associate, Physics, MIT	gneuert@mit.edu
Daniel Newburger	Computational Technician, Brigham & Women's/HMS	daniel.newburger@gmail.com
Amir Niknejad	PostDoc, ETSU	niknejad@etsu.edu
Roland Nilsson	Postdoctoral Fellow, Massachusetts General Hospital	nilsson@chgr.mgh.harvard.edu
Michael Nodine	PostDoc, Whitehead Institute, MIT	nodine@wi.mit.edu
Noa Novershtern	Student, Broad Institute	nshefi@broad.mit.edu
Daniel O'Connell	Graduate Student, Harvard Medical School	daniel_oconnell@hms.harvard.edu
Alice Oh	Student, CSAIL, MIT	aoh@mit.edu
Michael O'Kelly	Graduate Student, MIT	mokelly@mit.edu
Rachel O'Neill	Associate Professor, University of Connecticut	rachel.oneill@uconn.edu
Erin O'Shea	Professor, Harvard University/HHMI	erin_oshea@harvard.edu
Kimmo Palin	Researcher, University of Helsinki	kimmo.palin@helsinki.fi
Saurabh Paliwal	Graduate Student, The Johns Hopkins University	spaliwal@jhu.edu
Nathan Palmer	Graduate Student, MIT	palmer@mit.edu
Hemali Patel	Graduate Student, Boston University	hemali@bu.edu
Shouyong Peng	Research Fellow, Harvard Medical School	speng2@rics.bwh.harvard.edu
Cinzia Pizzi	PostDoc, INRIA Rhone-Alpes	cinzia.pizzi@dei.unipd.it
Yann Ponty	PostDoc, Boston College	ponty@bc.edu
Shobha Potluri	Post Doctoral Fellow, Pfizer	shobha.potluri@pfizer.com
Elizabeth Purdom	PostDoc, UC Berkeley, Statistics	epurdom@stat.berkeley.edu
Saamyadipta Pyne	Postdoctoral Associate, Broad Institute	spyne@broad.mit.edu
Daniela Raijman	Student, Tel Aviv University	raijmand@post.tau.ac.il
Anil Raj	Graduate Student, Applied Math, Columbia University	anil.calvin@gmail.com
Arvind Ramanathan	Student, Carnegie Mellon University	aramanat@andrew.cmu.edu
Maria Ramirez	Assistant Professor, Boston University School of Medicine	mm Ramirez@bu.edu
Pradipta Ray	Student, Carnegie Mellon University	meh@cs.cmu.edu
Soumya Raychaudhuri	Postdoctoral Fellow, Brigham & Women's/HMS	soumya@broad.mit.edu
Christopher Reeder	Student, MIT	reeder@mit.edu
Mireille Regnier	Researcher, INRIA	mireille.regnier@inria.fr
Deborah Ritter	Student, Boston College	ritterd@bc.edu
Laura Riva	Postdoctoral Associate, MIT	lriva@mit.edu
Manuel Rivas	Student, MIT	mrivas03@mit.edu
Helge Roeder	PhD student, Max Planck Institute for Molecular Genetics	roeder@molgen.mpg.de
Graham Ruby	Graduate Student, Whitehead Institute, MIT	grahamruby@yahoo.com
Elizabeth Ryder	Associate Professor, Worcester Polytechnic Institute	ryder@wpi.edu
Svetlana Shabalina	Staff Scientist, NCBI/NLM/NIH	shabalin@ncbi.nlm.nih.gov
Oded Shaham	Graduate Student, HST, Harvard/MIT	odeds@mit.edu
Maryam Shanechi	Student, MIT	shanechi@mit.edu
Kenichi Shimada	Graduate Student, Columbia University	ks2474@columbia.edu
Suyash Shringarpure	Student, Carnegie Mellon University	meh@cs.cmu.edu
Rahul Siddharthan	Reader, Institute of Mathematical Sciences	rsidd@imsc.res.in
Mary Sifferlen	Student, Harvard University	mary.sifferlen@acm.org
Trevor Siggers	PostDoc, Brigham & Women's/HMS	tsiggers@rics.bwh.harvard.edu
Emily Singer	Biotechnology editor, Technology Review	emily.singer@nasw.org
Donna Slonim	Professor, Tufts University	donna.slonim@tufts.edu
Gustav Smith	Student, Lund University, Sweden	gustav.smith@med.lu.se
Noah Spies	Graduate Student, Whitehead Institute, MIT	nspies@mit.edu

John Stamatoyannopoulos	Assistant Professor, University of Washington	jstam@u.washington.edu
Alex Stark	PostDoc, MIT	alex.stark@mit.edu
Juerg Straubhaar	Researcher, UMass Medical School	juerg.straubhaar@umassmed.edu
Aarathi Sugathan	Research Assistant, Waxman Lab, Boston University	aarathi@bu.edu
Rui Zhen Tan	Graduate Student, MIT	ruizhentan@gmail.com
Stephanie Tauber	PhD Student, Boston Univeristy	stauber@bu.edu
Robert Thurman	Senior Scientist, University of Washington	rthurman@u.washington.edu
Sonia Timberlake	Graduate Student, Biological Engineering, MIT	soniat@mit.edu
Michael Tolstorukov	Postdoctoral Fellow, Brigham & Women's/HMS	tolstorukov@gmail.com
Giovanni Tonon	Instructor, DFCI	giovanni_tonon@dfci.harvard.edu
John Tsang	Graduate Student, Harvard/MIT	tsang@mit.edu
Alex Tsankov	Student, MIT	tsankov@mit.edu
Aviad Tsherniak	Computational Biologist, Broad Institute	aviad@broad.mit.edu
Sevin Turcan	Graduate Student, Tufts University	sevin.turcan@tufts.edu
Esko Ukkonen	Professor, University of Helsinki	ukkonen@cs.helsinki.fi
Vinay Varadan	Senior Member Research Staff, Philips Research NA	vinay.varadan@philips.com
Vikram Vijayan	Student, Harvard University	vwijayan@fas.harvard.edu
Martin Vingron	Professor, MPI for Molecular Genetics	martin.vingron@molgen.mpg.de
Qingqing Wang	Graduate Student, Systems Biology, Harvard University	wang41@fas.harvard.edu
Qingqing Wang	Student, Systems Biology, Harvard University	qingqing@mit.edu
Xiaowo Wang	PhD Student, Cold Spring Harbor Laboratory	wangx@cshl.edu
Ilan Wapinski	Student, Broad Institute	ilan@eecs.harvard.edu
Lucas Ward	Graduate Student, Columbia University	ldw2105@columbia.edu
David Waxman	Professor, Boston Univeristy	djw@bu.edu
Zhiping Weng	Associate Professor, Boston Univeristy	zhiping@bu.edu
Charlie Whittaker	Research Scientist, Center for Cancer Research, MIT	charliew@mit.edu
Winfred Williams	Associate Member, Broad Institute	wwwilliams@partners.org
Aarron Willingham	Scientist, Affymetrix	aarron_willingham@affymetrix.com
Wei Wu	Assistant Professor, University of Pittsburgh	weiwu2004@gmail.com
Zhenhua Wu	PhD student, Boston Univeristy	zhenhua@bu.edu
Zeba Wunderlich	Graduate Student, Harvard University	wunderl@fas.harvard.edu
Hualin Xi	Student, Boston Univeristy	xih@bu.edu
Xinshu (Grace) Xiao	PostDoc, MIT	xiao@mit.edu
Huangming Xie	Graduate Student, Whitehead Institute, MIT	xieh@mit.edu
Eric Xing	Assistant Professor, Carnegie Mellon University	epxing@cs.cmu.edu
Itai Yanai	PostDoc, Harvard	yanai@mcb.harvard.edu
Qiong Yang	Graduate Student, MIT	qiongy@mit.edu
Esti Yeger-Lotem	PostDoc, Whitehead, MIT	esti@wi.mit.edu
Boyan Yordanov	Graduate Student, Boston Univeristy	yordanov@bu.edu
Bingbing Yuan	Bioinformatics Scientist, Whitehead Institute, MIT	byuan@wi.mit.edu
Yue Yun	Graduate Student, Washington University in St. Louis	yyun@wustl.edu
Erliang Zeng	Student, Florida International University	ezeng001@cs.fiu.edu
Chaolin Zhang	Graduate Student, Cold Spring Harbor Laboratory	zhangc@cshl.edu
Lu Zhang	Student, Boston College	zhanglv@bc.edu
Wen Zhang	Postdoctoral Fellow, Brigham & Women's/HMS	wzhang@rics.bwh.harvard.edu
Yue Zhao	Student, Washington University in St. Louis	yue.zhao@wustl.edu
Grace Zheng	Graduate Student, MIT	gzheng@mit.edu
Wen Zhou	Graduate Student, Harvard University	wenzhou1985@gmail.com
Julie Zhu	Assistant Professor, UMass Medical School	julie.zhu@umassmed.edu