

### Annotating Non-coding Regions of the Human Genome

RECOMB Satellite Meeting on Regulatory Genomics Cambridge, MA 2008.10.31, 18:45-19:15

Mark B Gerstein Yale (Comp. Bio. & Bioinformatics)

#### Slides from Lectures.GersteinLab.org

(Please read permissions statement.) Paper references mostly from Papers.GersteinLab.org. See streams.gerstein.info on photos & images

> (Genome tech and Genome annotation talk, including: Seq. Sim + PeakSeq ; MSB; DART TAR classification ; TRE clustering + biplot, BoCaTFBS, sdcnvcorr [I:RECOMBSAT] )



2001: Most of the genome is not coding (only ~1.2% exon). It consists of elements such as repeats, regulatory regions, non-coding RNAs, origins of replication, pseudogenes, segmental duplications....What do these elements do? How should [IHGSC, Nature 409, 2001] they be annotated? [Venter et al. *Science* 29, 2001]



# 2007 : Pilot results from ENCODE Consortium on decoding what the bases do

- 1% of Genome (30 Mb in 44 regions)
- Tiling Arrays to assay Transcription & Binding
- Multi-organism sequencing and alignment
- Careful Annotation
- Variation Data

[IHGSC, *Nature* 409, 2001] [ENCODE Consortium, *Nature* 447, 2007]

#### How might we The Semicolon Wars annotate a human text ? Brian Hayes F YOU WANT TO BE a thorough-Every programmer going world traveler, you need to learn 6,912 ways to say "Where is the toilet, please?" That's the number of knows there is one languages known to be spoken by the peoples of planet Earth, according to true programming Ethnologue.com. If you want to be the complete polylanguage. A new one glot programmer, you also have quite a challenge ahead of you, learning all **Color** is every week the ways to say: **Function** printf("hello, world\n"); (This one is in C.) A catalog maintained a good-enough notation-for expressby Bill Kinnersley of the University of ing an algorithm or defining a data Kansas lists about 2,500 programming Lines are structure. There are programmers of my aclanguages. Another survey, compiled **Similarity** quaintance who will dispute that last by Diarmuid Piggott, puts the total even higher, at more than 8,500. And statement. I expect to hear from them. keep in mind that whereas human lan-They will argue-zealously, ardently, guages have had millennia to evolve vehemently-that we have indeed found the right programming lanand diversify, all the computer languages have sprung up in just 50 years. Even guage, and for me to claim otherwise by the more-conservative standards of is willful ignorance. The one true lanthe Kinnersley count, that means we've guage may not yet be perfect, they'll been inventing one language a week, concede, but it's built on a sound foun-[B Hayes, on average, ever since Fortran. dation and solves the main problems, Am. Sci. For ethnologists, linguistic diversity and now we should all work together is a cultural resource to be nurtured to refine and improve it. The catch, of (Jul.- Aug. and preserved, much like biodiversity. course, is that each of these friends will '06)]

All human languages are valuable; the

favor a different language. It's Lisp,

cide which end of a boiled egg to crack. This famous tempest in an egg cup was replayed 250 years later by designers of computer hardware and communications protocols. When a block of data is stored or transmitted, either the leastsignificant bit or the most-significant bit can go first. Which way is better? It hardly matters, although life would be easier if everyone made the same choice. But that's not what has happened, and so quite a lot of hardware and software is needed just to swap ends at boundaries between systems.

This modern echo of Swift's Endian wars was first pointed out by Danny Cohen of the University of Southern California in a brilliant 1980 memo, "On holy wars and a plea for peace." The memo, subsequently published in Computer, was widely read and admired; the plea for peace was ignored.

Another feud—largely forgotten, think, but never settled by truce or treaty-focused on the semicolon. In Algol and Pascal program statements have to be separated by semicolons. For example, in x := 0; y := x+1; z := 2 the semicolons tell the compiler where one statement ends and the next begins. C programs are also peppered with semi-

## **Overview of the Process of**

### **Annotation of non-coding Regions**

### • Basic Inputs

- 1. Doing large-scale similarity comparison, looking for repeated or deleted regions
- 2. Determining experimental signals for biochemical activity (e.g. transcription) across each base of genome

# Results of Analyzing Similarity Comparison

- 1. Finding large repeated or deleted blocks (e.g. CNVs) as a function of degree of similarity
  - 1. within reference human genome
  - 2. within human population
  - 3. between related organisms (e.g. mouse)
- 2. Finding smaller "exon-level" similarities (e.g. pseudogenes)



#### Determining experimental signals for biochemical activity (e.g. transcription) across each base of genome

- Development of Sequence (and Array) Technology
  - Normalizing & Scoring Signal, Correcting Artifacts, Segmenting to create Small Annotation Blocks
- Output of Production Pipelines and Surveying a Single Type of Annotation on a Large-scale
  - Clustering Small Blocks into Larger Ones, Surveying
- Integrated Analysis Connecting Different Types of Annotation
  - Building networks and beyond

## Signal Processing: Normalizing Signal and Finding Initial Annotation Blocks ("Hits")



pers. **photo**, see streams.gerstein.info

### Representative Signal from Chip-Seq



#### Genome / Genomic region



<u>Correcting</u> <u>Chip-seq Signal by</u> <u>Simulating a Non-</u> <u>uniform Genomic</u> <u>Background</u>

We developed *in silico* ChIP sequencing, a computational method to simulate the experimental outcome.

14 Lectures.GersteinLab.org

(C)

14 (c) Mark Gerstein, 2002, Yale, bioinfo.mbb.yale.edu

[Zhang et al. PLoS Comp Bio. ('08)]





#### **ChIP-sequencing simulation**

- Contrary to the common belief, the background is mildly fluctuating and contains some 'hot' spots.
- Simple uniform background model does not count for all the variation in the background and thus leads to a serious underestimation of the background noise.
- Our study demonstrates that both the genomic background of ChIP and binding sites are not uniform.
- Simulated distributions segments the actual distribution into four sections.

[Zhang et al. PLoS Comp Bio. ('08)]

### **ChIP-Seq vs Input DNA Control**



[Rozowsky et al. (submitted)]



Scored results consistent with simulation

Actual peaks at tail of power-law graph



[Rozowsky et al. (submitted)]

### Punctate Regions vs Broad Regions

#### **Representative Signal from aCGH with CNVs & Breakpoints**



LCR A

BCD

Urban et al. (2006) PNAS



- (**x**<sub>i</sub>) Observed depth of coverage counts (or array signal) as samples from PDF
- (m) Kernel-based approach to estimate local gradient of PDF
- $(\mathbf{y}_{c})$  Iteratively follow grad to determine local modes

#### Not Model-based (e.g. like HMM)

with global optimization, distr. assumption & parms. (e.g. num. of segments). Achieves discontinuity-preserving smoothing

3

### **Representative Result Showing Segmentation Based on Depth of Coverage**



MSB is not model based so can be applied equally well to pseudo-signal from coverage depth as to

CGH arrays ??

NA11995 (seq. by Sanger, MAQ mapping) chr 21 (46162500 to 46164711)

[Wang et al. Gen. Res (in press, '08)]



### Annotating a single type of signal on a large-scale: Clustering and Characterizing Binding Sites (TREs)

pers. **photo**, see streams.gerstein.info

### Clustering Binding Sites at ~50kb resolution



### Clustering Binding Sites at ~50kb resolution



26 Lectures.GersteinLab.org (c)

#### Landscape of ENCODE Transcriptional Regulatory Elements

- Analyzed 105 lists of transcriptional regulatory elements in the encode regions
- 29 transcription factors, 9 cell lines, 2 time points

◊RNA Pol2

- Histone modifications such as Ac & Me
- $\diamond \, \text{Core promoters}$
- Oromoter proximal elements
- Others such as enhancers, silencers, insulators, & response elements
- [CFTR] ENm001 <del>╶╷╏</del>╖╫<mark>╗╴╎╢┉┽╦╅┥╍┉┢┍╴╶╘┅┉┟╶┝╶╎╓╎╓┼╓╶╶╓╓╶╴╟╴╍╻┾╷╷╵┉┢┦╴</mark> [Interleukin] ENm002 [Apo] ENm003 [Chr22] ENm004 \_\_\_\_\_\_ [ChrX] ENm006 [Chr19] ENm007 [α-globin] ENm008 [β-globin] ENm009 Zhang et al. (2007) Gen. Res. [HOXA] ENm010 [IGF2/H19] ENm011 [FOXP2] ENm012 [7q21.13] ENm013 [7q31.33] ENm014 ENr121 ENr131 ENr122 ENr112 ENr132 ENr123 ENr133 ENr114 ENr221 ENr211 ENr231 ENr212 ENr222 ENr232 ENr223 ENr213 ENr233 ENr321 ENr331 ENr312 ENr332 -----ENr322 ENr323 ENr333 ENr313 111 ENr334 ENr324 1500000 bp 500000 1000000

#### Collect Total Hits for Each Factor in ~6000 Bins of 10 to 100 kb and Compare to Random Control



### **Non-random distribution of TREs**

- TREs are not evenly distributed throughout the encode regions (*P* < 2.2×10<sup>-16</sup>).
- The actual TRE distribution is power-law.
- The null distribution is 'Poissonesque.'
- Many genomic subregions with extreme numbers of TREs.



Number of TREs in a subregion

### Local TRE enrichment and depletion: Annotation of Desserts and Forests

Known Genes

KATNAL1

- Hundreds of TRE 'forests' and 'deserts' are identified in ENCODE regions.
- The entirety of *ehd1* on chromosome 11 is covered by TRE islands.
- Some of islands are located in the intergenic regions in the genome.



#### dart.gersteinlab.org/encode/tr/

Zhang et al. (2007) Gen. Res.

HMGB1

BX647267

### **Biplot to Show Overall Relationship of TFs** and Genomic Bins



15 18 22

18 13 36

24 10 32

8

9

10

а а 1.00b -0.44 .48 C 0

1





Zhang et al. (2007) Gen. Res.

- Biplot groups TFs into sequence-specific and sequence-nonspecific clusters.
  - $\diamond~$  c-Myc may behave more like a sequence-nonspecific TF.
  - H3K27me3 functions in a transcriptional regulatory process in a rather sequence-specific manner.
- Genomic Bins are associated with different TFs and in this fashion each bin is "annotated" by closest TF cluster

### TRE analysis on the microgenomic scale



### TRE analysis on the microgenomic scale



- Traditional motif learners (e.g. consensus sequences, profile methods, and HMMs) only use positive information
- ChIP-chip & Chip-seq give vast amount of negative information (regions not bound)
- Explicitly use this in constructing classifier that <u>refines</u> known positive motif seeds
- Use sequence of Alternating Decision Trees (ADTboost), which allow explicit inter-positional correlations between nucleotide positions

## Using Binding Site Regions Found by ChIP-chip to refine motifs: BoCaTFBS



[Wang et al., GenomeBiology ('06)]

#### Good performance compared to traditional motif-finders but large negative set requires training and detection cascade for efficiency and balance



45 Lectures.GeisteinLab.org

0



## Annotating a single type of signal on a large-scale: Clustering and Classifying Unannotated Transcription (TARs)

pers. photo, see streams.gerstein.info





Rozowsky et al. Genome Research (2007)







ENCODE Regions (30 Mb)

Locations of TARs

Of the approx 7,000 Novel TARs

- 955 are assigned to known genes
- 1,463 are clustered into ~200 Novel Loci

•DART Classification has been experimentally validated with some small scale experiments

- ◊ RT-PCR & Sequencing
- ◊ 18/46 (39%) confirmed by RT-PCR
- ♦ 4/5 Sequenced Products Map uniquely to correct genomic region

Rozowsky et al. Genome Research (2007)

## Example predicted structured RNAs (using RNAz)

#### Overlap of predicted structured RNAs with the union of TARs/Transfrags and the "moderate" set of sequence-constrained elements



[>700 candidate structured RNAs predicted in 1% of the reference genome] Stefan Washietl, Jakob Pedersen, Jan Korbel *et al.* (2007) *Genome Res* 17:852-864

# Analyzing Repeated Blocks in the Genome (SDs & CNVs)



55 Lectures.GersteinLab.org

(C)

pers. photo, see streams.gerstein.info

#### SEGMENTAL DUPLCATIONS AND COPY NUMBER VARIANTS ARE RELATED PHENOMENA AND HAVE BEEN CREATED BY SEVERAL DIFFERENT MECHANISMS





#### **NAHR** (Non-allelic homologous recombination)

Flanking repeat (e.g. Alu, LINE...)



#### NHEJ

(Non-homologous-endjoining)

No (flanking) repeats. In some cases <4bp microhomologies

#### PERFORM LARGE SCALE CORRELATION ANALYSIS TO DETECT REPEAT SIGNATURES OF SDs AND CNVs



[Kim et al. Gen. Res. (in press, '08), arxiv.org/abs/0709.4200v1 ]

#### SDs ARE CORRELATED WITH ALUS AND OTHER SDs



[Kim et al. Gen. Res. (in press, '08), arxiv.org/abs/0709.4200v1 ]

#### **ASSOCIATIONS ARE DIFFERENT FOR SDs AND CNVs**

#### **CNVs ARE LESS ASSOCIATED WITH** SD association with repeats **SDs THAN THE GENERAL SD TREND** 0.27 CNV 0.21 Association 0.094 0.07 with SDs Microsatellite Pseudogenes LINE Alu 0.31 <0.001 (<0.001) 0.046 0.001 0.11 **CNV** association with repeats 0.0739 0.048 0.0466 0.0006 >99% SDs\* CNVs Microsatellite **Pseudogenes** LINE Alu 0.046 0.92 <0.001 0.001

[Kim et al. Gen. Res. (in press, '08), arxiv.org/abs/0709.4200v1 ]

59



AFTER THE ALU BURST, THE **IMPORTANCE OF ALU ELEMENTS FOR GENOME** REARRANGEMENT **DECLINED RAPIDLY** 

- About 40 million years ago there was a burst in retrotransposon activity
- The majority of Alu elements stem from that time
- This, in turn, led to rapid genome rearrangement via NAHR
- The resulting SDs, could create more SDs, but with Alu activity decaying, their creation slowed



pers. **photo**, see streams.gerstein.info

## Overview of the Process of Intergenic Annotation

- Basic Inputs
  - 1. Doing large-scale similarity comparison, looking for repeated or deleted regions
  - 2. Determining experimental signals for biochemical activity (e.g. transcription) across each base of genome
- Results of Analyzing Similarity
  Comparison
  - Finding large repeated or deleted blocks (e.g. CNVs) as a function of degree of similarity
    - 1. within reference human genome
    - 2. within human population
    - 3. between related organisms (e.g. mouse)
  - 2. Finding smaller "exon-level" similarities (e.g. pseudogenes)

- Results of Processing Raw Expt. Signals
  - Signal Processing: removing artifacts, normalizing, window averaging
  - 2. Segmenting signal into larger "hits" ("Active Regions" or ARs)
  - 3. Clustering together active regions into even larger features at different length scales and classifying them
  - 4. Building networks and beyond....

### Segmenting the Raw "Signal" from Next-generation Sequencing into Usable Annotation Blocks

#### • MSB

- ◊ Mean-shift segmentation approach following grad. of PDF
- ◊ Equally applied to aCGH and depth of coverage of short reads

#### PeakSeq

- Scoring chip-seq expt relative to input control
- Simulating chip-seq expt anticipates & allows correction for nonuniformity

### First-Pass Annotation Clustering and Characterizing Novel Transcribed Regions and Groups of Binding Sites

- Deserts and Forests of Binding Activity
  - $\diamond~$  on ~50kb scale
  - Biplot gives broad separation of seq. specific and non-specific factors and associated genomic bins
- Analyzing Promotors
  - Observation BoCaTFBS: Refining binding site motifs based on the results of chIPchip experiments
- DART classification of TARs
  - $\diamond~$  1300 TARs in ~200 novel pilot ENCODE loci
    - based on expression and phylogenetic clustering

#### Analysis of Duplication in the Genome: SVs and SDs

- Large-scale analysis of existing CNVs & SDs in human genome
- SDs assoc. with Alu, pseudogenes and older SDs
- CNVs assoc. other repeats (microsat.) and not as much with SDs
- Suggestion: Alu burst 40 MYA triggered much NAHR rearrangement, then dupl. feed on itself in hotspots but now dying down and NAHR assoc. with other repeats and CNVs also from NHEJ

#### **Acknowledgements**

#### pseudogene.org, tiling.gersteinlab.org, sv.gersteinlab.org



#### **Acknowledgements**

#### pseudogene.org, tiling.gersteinlab.org, sv.gersteinlab.org



#### **Acknowledgements**

#### pseudogene.org, tiling.gersteinlab.org, sv.gersteinlab.org





### **Permissions Statement**

This Presentation is copyright Mark Gerstein, Yale University, 2007.

# Feel free to use images in it with **PROPER acknowledgement**

(via citation to relevant papers or link to gersteinlab.org).