## Slide 1

*Personal Transcription Factor Binding Site Mutations Point to Personal Medical Histories*

One GENOME to Rule Them All

**Gill Bejerano**

DEVBIO · STANFORD COMPUTER SCIENCE · STANFORD SCHOOL OF MEDICINE Pediatrics

http://bejerano.stanford.edu — 1

## Slide 2

**Genome = Genes + Gene Regulation**

CIS REGULATION

| Type | # in genome | % of genome |
|---|---|---|
| genes | 20,000 | 2% |
| ncRNA | 20,000 | 1% |
| cis elements | 1,000,000 | >10% |

- Encode causality
- Disease susceptibility
- Driver sequences
- Alter cell state
- Key for evolution

Atomic event – transcription factor binding

http://bejerano.stanford.edu — 2

## Slide 3

**Disease Associated tag SNPs**

- Over 15,000 distinct tag SNPs in the GWAS Catalog
- 80-90% far away from (linkage with) gene exons
- Are most gene cis regulatory?
- Are they near genes with common functionality?

GWAS Catalog Growth

2008     2011     2016

http://bejerano.stanford.edu — 3

## Slide 4

**Cis-reg enrichments: GREAT.stanford.edu**

½ million job submissions, 700+ references, established defaults

Gene transcription start site  
$\pi$ Function ('abnormal cardiac output')  
Gene regulatory domain  
Cis-reg rich region set

$p_\pi$ = 0.33 of genome annotated with $\pi$  
$n$ = 6 genomic regions  
$k_\pi$ = 5 genomic regions hit annotation

GREAT = Genomic Regions Enrichment of Annotations Tool

$P = \mathrm{Pr}_{binom}(k_\pi \geq 5 \mid n=6, p_\pi=0.33)$

[McLean et al, Nature Biotech, 2010]

http://bejerano.stanford.edu/ — 4

## Slide 5

**Cis-reg enrichment: GREAT.stanford.edu**

½ million job submissions, 700+ references, established defaults

Gene transcription start site  
$\pi$ Function ('abnormal cardiac output')  
Gene regulatory domain  
Cis-reg rich region set

Advantages of GREAT:
1. Accounts for both proximal and distal binding sites
2. Variable length gene regulatory domains
3. Multiple hits next to same gene add significance
4. Extensive body of knowledge (16,000 functions)

$p_\pi$ = 0.33 of genome annotated with $\pi$  
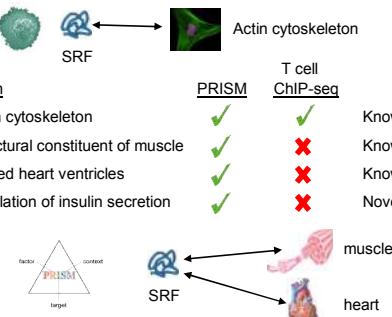$n$ = 6 genomic regions  
$k_\pi$ = 5 genomic regions hit annotation

GREAT = Genomic Regions Enrichment of Annotations Tool

$P = \mathrm{Pr}_{binom}(k_\pi \geq 5 \mid n=6, p_\pi=0.33)$

[McLean et al, Nature Biotech, 2010]

http://bejerano.stanford.edu/ — 5

## Slide 6

**Unlinked GWAS SNPs → GREAT**

http://bejerano.stanford.edu — 6

## Unlinked GWAS SNPs → GREAT

## Individual Genomes → GREAT ?



Associate to nearest gene(s)

Do individuals carry pathway specific cis-reg mutation load?
Which binding sites?
• Biochemically active ≠ functional
• Under purifying selection
From which context?
   The genome, which codes for *all* contexts.

## TF Motif Library (PBM+ChIP+SELEX)



motifs for 657 TFs

## Predict conserved binding sites

. = same as human

• We in fact allow:
   • Imperfect matches
   • Binding site / alignment wobble
• Take measures against alignment fragmentations.
• Predict efficiently.
• Improve state of the art using "Excess conservation" scoring

## Compare *everything* to shuffled motifs & weed!



SRF motif
shuffle #1
shuffle #2
shuffle #3
shuffle #4
…
shuffle #10

[Wenger et al, Genome Research, 2013]

## PRISM vs. ChIP-seq → GREAT



SRF — Actin cytoskeleton

| Term | PRISM | ChIP-seq |
|---|---|---|
| actin cytoskeleton | ✓ | ✓ (T cell, Known) |

2

## PRISM vs. ChIP-seq → GREAT



Actin cytoskeleton
SRF

| Term | PRISM | T cell ChIP-seq | |
|---|---|---|---|
| actin cytoskeleton | ✓ | ✓ | Known |
| structural constituent of muscle | ✓ | ✗ | Known |
| dilated heart ventricles | ✓ | ✗ | Known |
| regulation of insulin secretion | ✓ | ✗ | Novel |

SRF → muscle
→ heart

http://bejerano.stanford.edu
13

## PRISM vs. ChIP-seq → GREAT



Actin cytoskeleton
SRF

| Term | PRISM | T cell ChIP-seq | |
|---|---|---|---|
| actin cytoskeleton | ✓ | ✓ | Known |
| structural constituent of muscle | ✓ | ✗ | Known |
| dilated heart ventricles | ✓ | ✗ | Known |
| regulation of insulin secretion | ✓ | ✗ | Novel |

Every known function is supported
by dozens or hundreds of *novel* binding sites.

Note: Sensitivity vs Specificity

http://bejerano.stanford.edu
14

## GWAS SNPs: Predict upstream regulator



Huang et al. A prostate cancer susceptibility allele at 6q22 increases RFX6 expression by modulating HOXB13 chromatin binding. Nat. Genet. 2014 Feb;46(2):126–135.
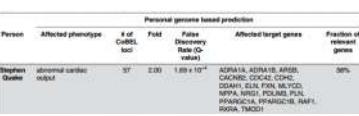
http://bejerano.stanford.edu
15

## Personal Deleterious Binding Sites



COBELs = Conserved Binding Site Eroding Loci

Individuals w public medical records

http://bejerano.stanford.edu
16

## COBELs → GREAT



http://bejerano.stanford.edu
17

## COBELs → GREAT



http://bejerano.stanford.edu
18

## COBELs → GREAT

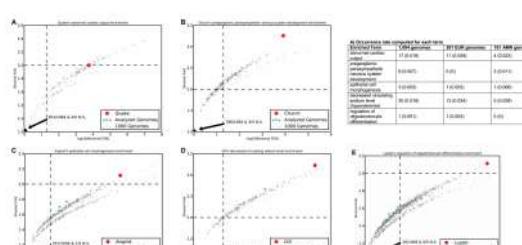## COBELs → GREAT

## COBELs → GREAT

## COBELs → GREAT

## Randomize COBELs

- Replace every CoBEL with a random binding site prediction for the same transcription factor of same affinity and similar cross-species conservation.

- Using 10,000 random control sets, the likelihood of obtaining the functions reported in Table 1 as top prediction due to bias in the distribution of binding sites in the genome is low (Quake P = 3 x 10−4, Church P = 5.7 x 10−3, Angrist P = 4.8 x 10−3, Gill P = 1 x 10−4, Lupski P = 1.9 x 10−3, and combined P = 1.6 x 10−15).

- Significance remains high when we relax the requirement to recover each exact same term with matching any one of a broader group of 12–60 related functions as a top prediction (Quake P = 1.1 x 10−3, Church P = 1.3 x 10−2, Angrist P = 7.7 x 10−3, Gill P = 7.4 x 10−3, Lupski P = 6.5 x 10−3, and combined P = 5.2 x 10−12).

## KGP as Controls

4

## Randomize Medical Histories

- Define an **association matrix linking enrichment and medical history**, with the phenotypes observed in the five individuals as rows, and top enriched terms in all as columns. A cell in the matrix would be marked "true" only where the enriched term (of any individual) is thought to be related to the etiology of the phenotype (of any individual).
- One instance of this matrix was filled by a **medical doctor** based on their medical knowledge and training and another instance was independently filled using a **literature survey**. The objective was to compute the chance of associating a set of five individuals with random medical histories with the observed enrichments using one of the two association matrices as the "gold" association.
- We generated 1,000 sets of five individuals with random medical histories composed of similar disease profiles and assessed the likelihood of being able to associate them with enrichments. Successfully linking five random individuals with enrichments was highly significant using the association matrix generated by the **medical doctor (P = 3.0 x 10−3)** and by the matrix generated by **literature survey (P = 3.0 x 10−2)** suggesting our links between enrichment and medical histories are not just a function of the listed histories.
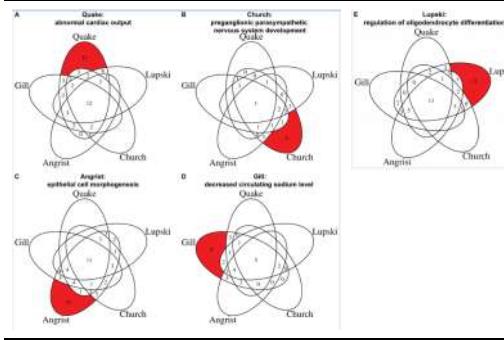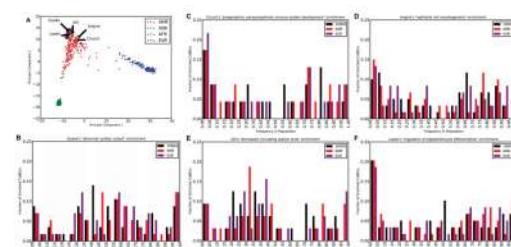
http://bejerano.stanford.edu 25

## COBELs ≠ GWAS SNPs or HGMD



http://bejerano.stanford.edu 26

## Most Predictive COBELs are Private



http://bejerano.stanford.edu 27

## Contributions from Common & Rare



Subset to 1% freq in KGP → lose all enrichments
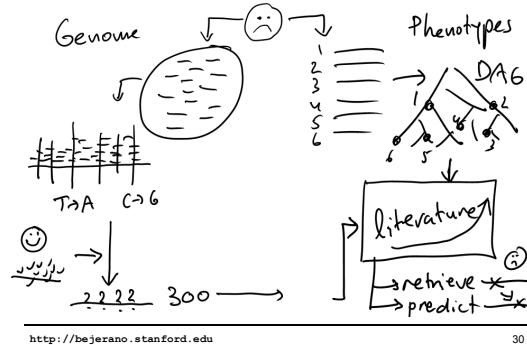
http://bejerano.stanford.edu 28

## Summary

- We define likely deleterious events as personal variants that erode the affinity of human conserved binding sites.
- When the set of all such events is probed for lying next to gene sets of particular function or phenotype, we repeatedly get a solid match between top genomic prediction and self reported medical summary.
- Top genomic predictions are eroded at both gene and gene set level.
- The variants we highlight appear to be part of the mutational load pre-disposing individual lineages to different diseases.
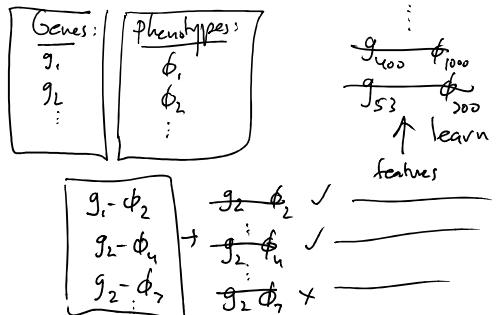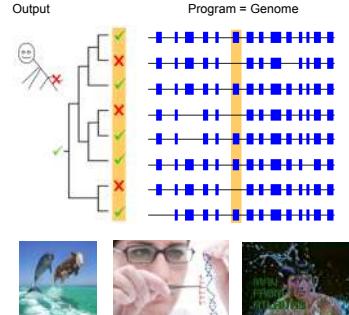
COBELs

http://bejerano.stanford.edu 29

## Other Lab Interests: 1) Solve Patient



http://bejerano.stanford.edu 30

5

## 2) Automate Patient Solving

## 3) Discover Mammalian Adaptations

Output          Program = Genome

## Kudos

COBELs: (PLoS Comp Bio, 2016)
**Harendra Guturu**, Sandeep Chinchali, Shoa Clarke

PRISM: (Genome Research, 2013)
Aaron Wenger, Shoa Clarke, Harendra Guturu, Jenny Chen, Bruce Schaar, Cory McLean

GREAT: (Nature Biotechnology, 2010)
Cory McLean, Dave Bristor, Michael Hiller, Shoa Clarke, Bruce Schaar, Craig Lowe, Aaron Wenger

Bejerano Lab past & present          The Organizers