# Explorations in Summer Camp in CT: Prioritizing non-coding mutations as potential cancer drivers
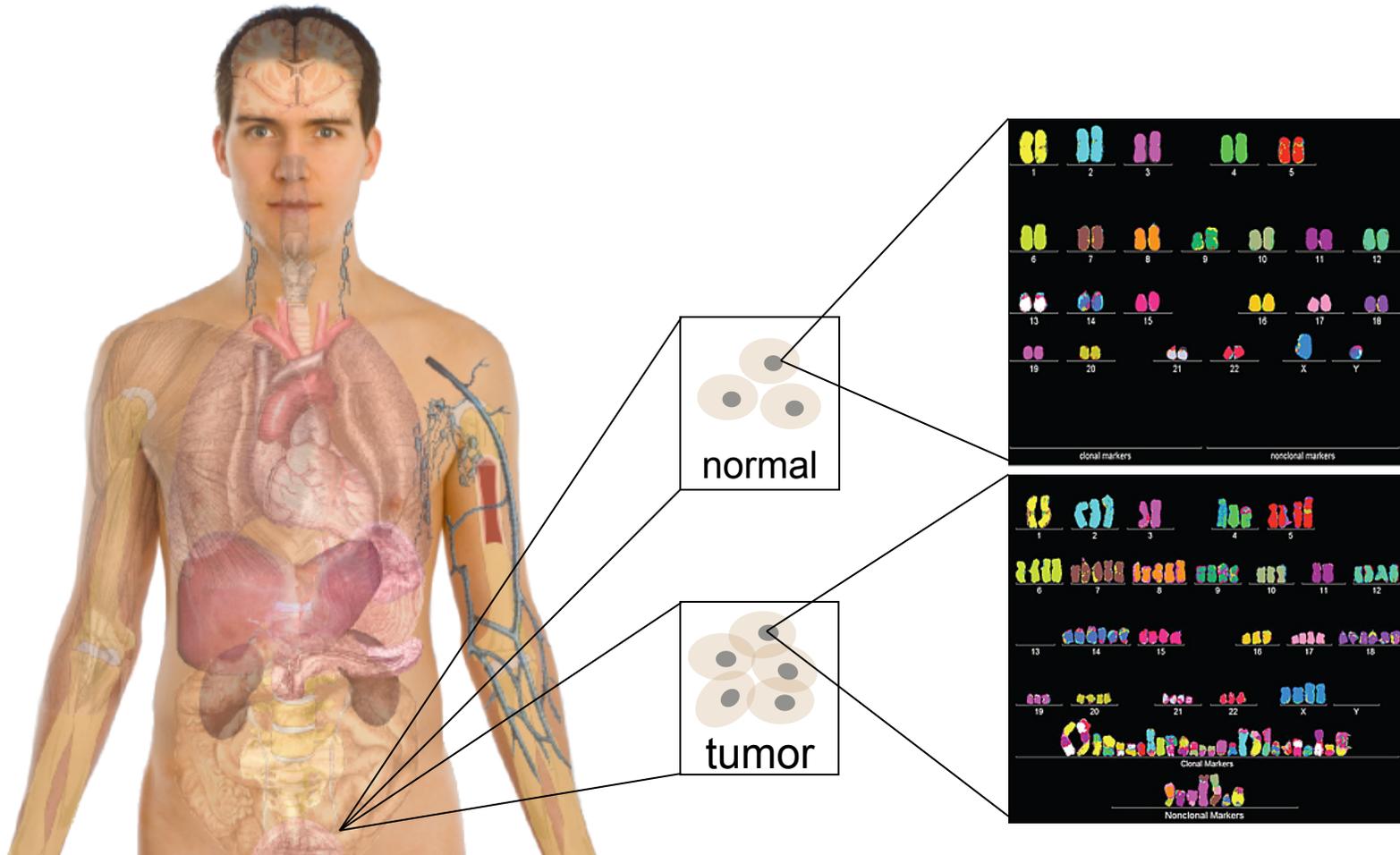
Mark Gerstein

# Personal Genomics
# as a Gateway into Biology

Personal genomes soon will become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.
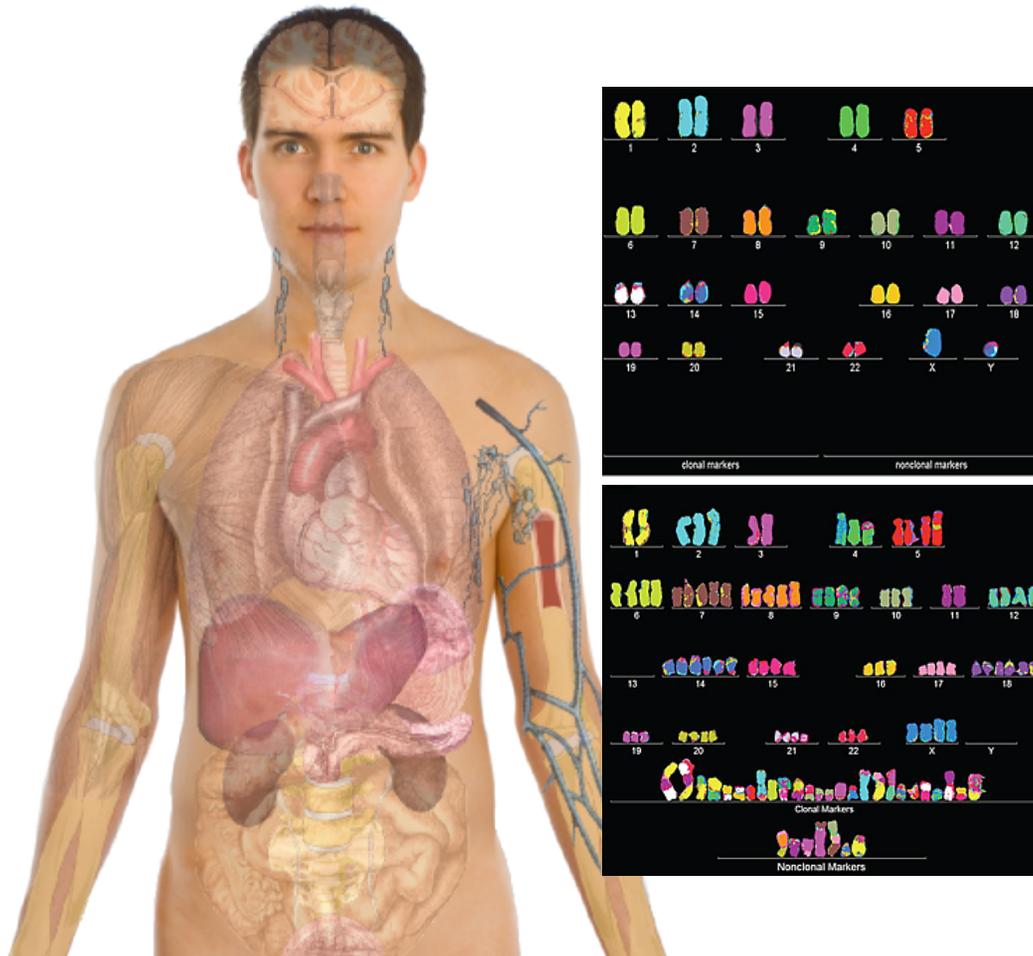
# Personal Genomics
# as a Gateway into Biology

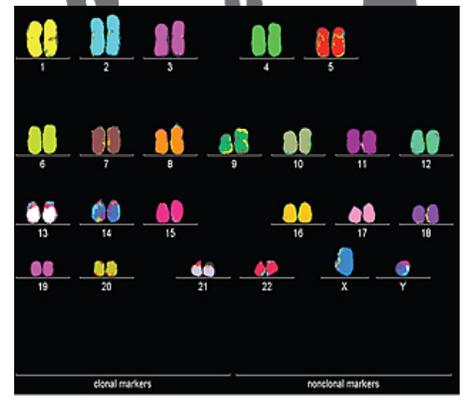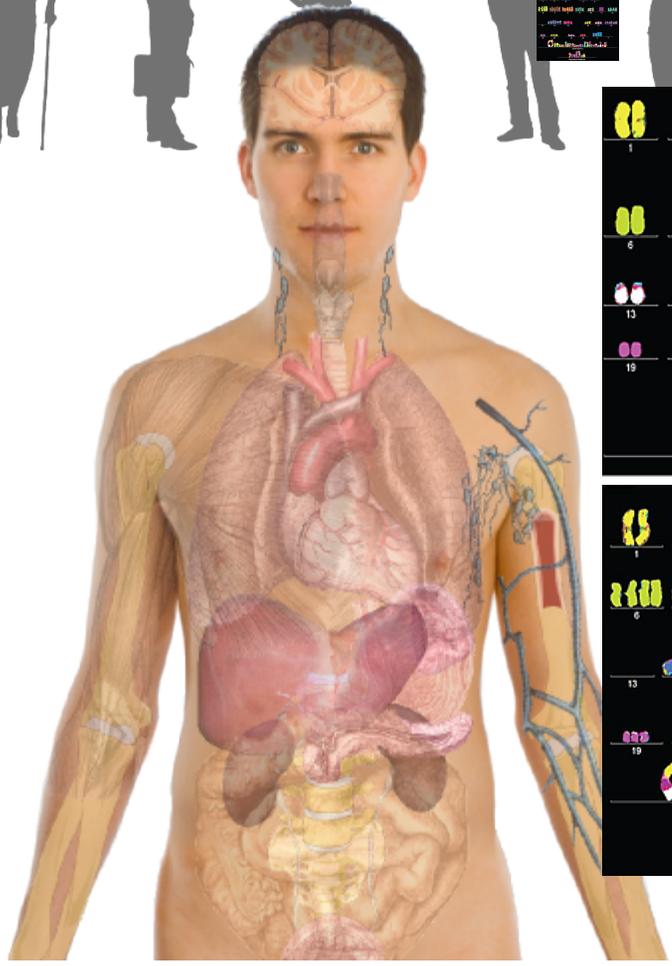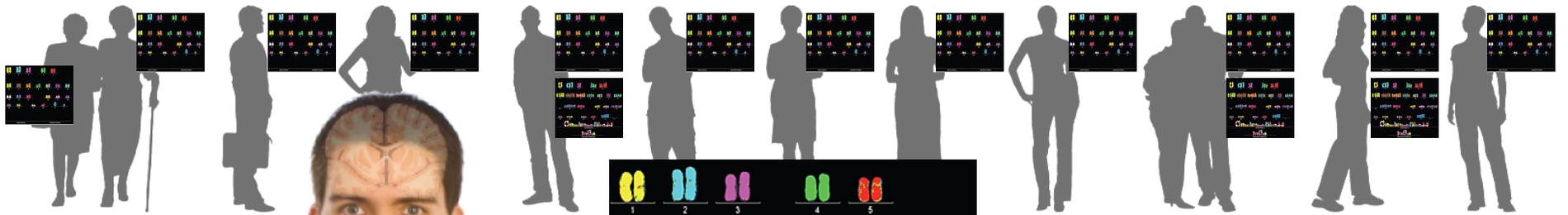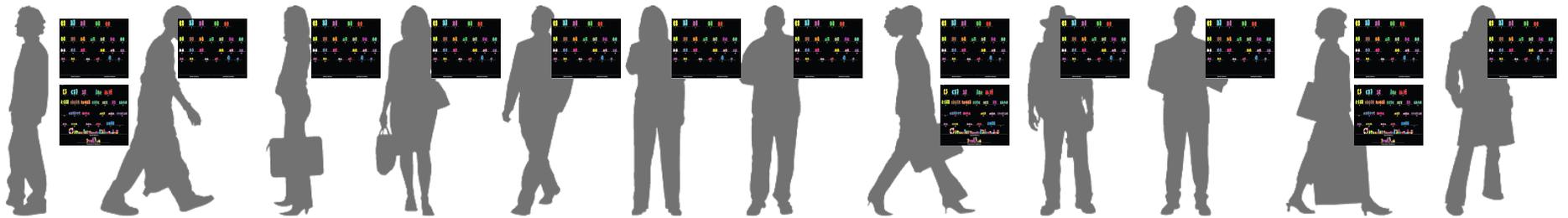Personal genomes soon will become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.

# Where is Waldo?

## (Finding the key mutations in ~3M Germline variants & ~5K Somatic Variants in a Tumor Sample)

# Non-coding Annotations: Overview

Most of cancer genomics has focused on mutations in non-coding regions – ie the exome
There are several collections of information "tracks" related to non-coding features, perhaps of use

Sequence features, incl. **Conservation**

**Functional Genomics**
Chip-seq (Epigenome & seq. specific TF)
and ncRNA & un-annotated transcription

- **Finding Non-coding Regions Sensitive to Mutations**
  - **1st Level Linear Annotation: Regulatory Sites**
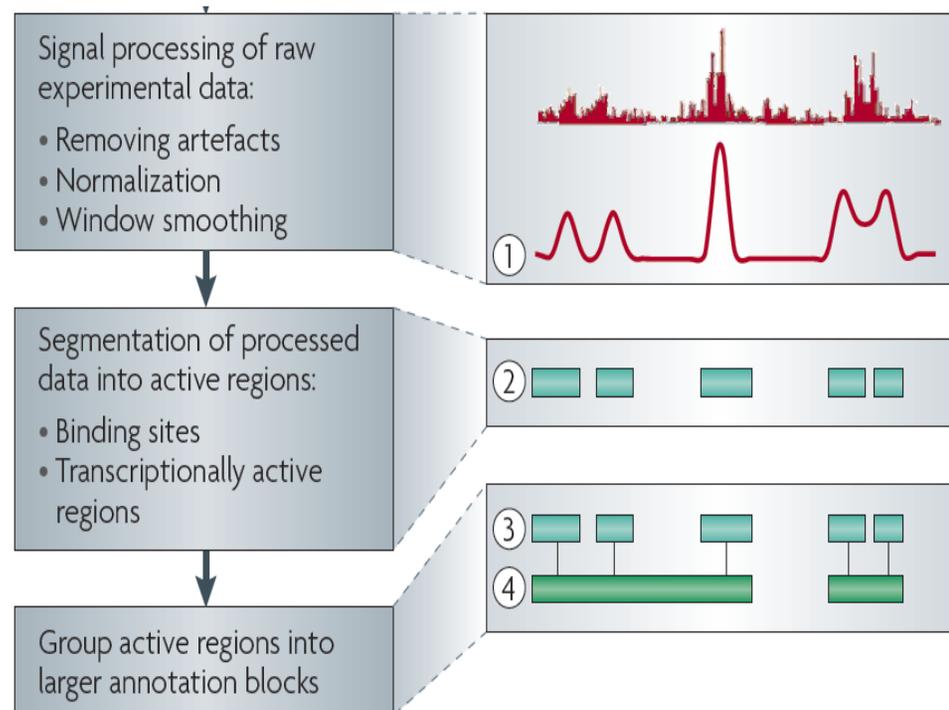    - Multi-scale "site" calling (with Music)
    - Finding small number of sites particularly sensitive to mutations
  - **2nd Level Network Annotation**
    - Building a network from the linear annotation
    - More connectivity = more constraint => highlights hubs

- **Using this to Interpret Alterations in Cancer**
  - **LARVA: to find recurrently mutated annotations**
    - Need to correct for overdispersion in bionomial
    - Use beta-bin parameterized according to replication timing
  - **FunSeq software tool for mutation prioritization**
    - Systematically weighting all the features, for non-coding prioritization

# Summarizing the Signal:
# "Traditional" ChipSeq Peak Calling

- Generate & threshold the signal profile to identify candidate target regions
  - Simulation (PeakSeq),
  - Local window based Poisson (MACS),
  - Fold change statistics (SPP)

ChIP

Threshold

Potential Targets

Normalized Control

- Score against the control

Significantly Enriched targets

# Now an update: "PeakSeq 2" => MUSIC

[Rozowsky et al. ('09) *Nat Biotech*]
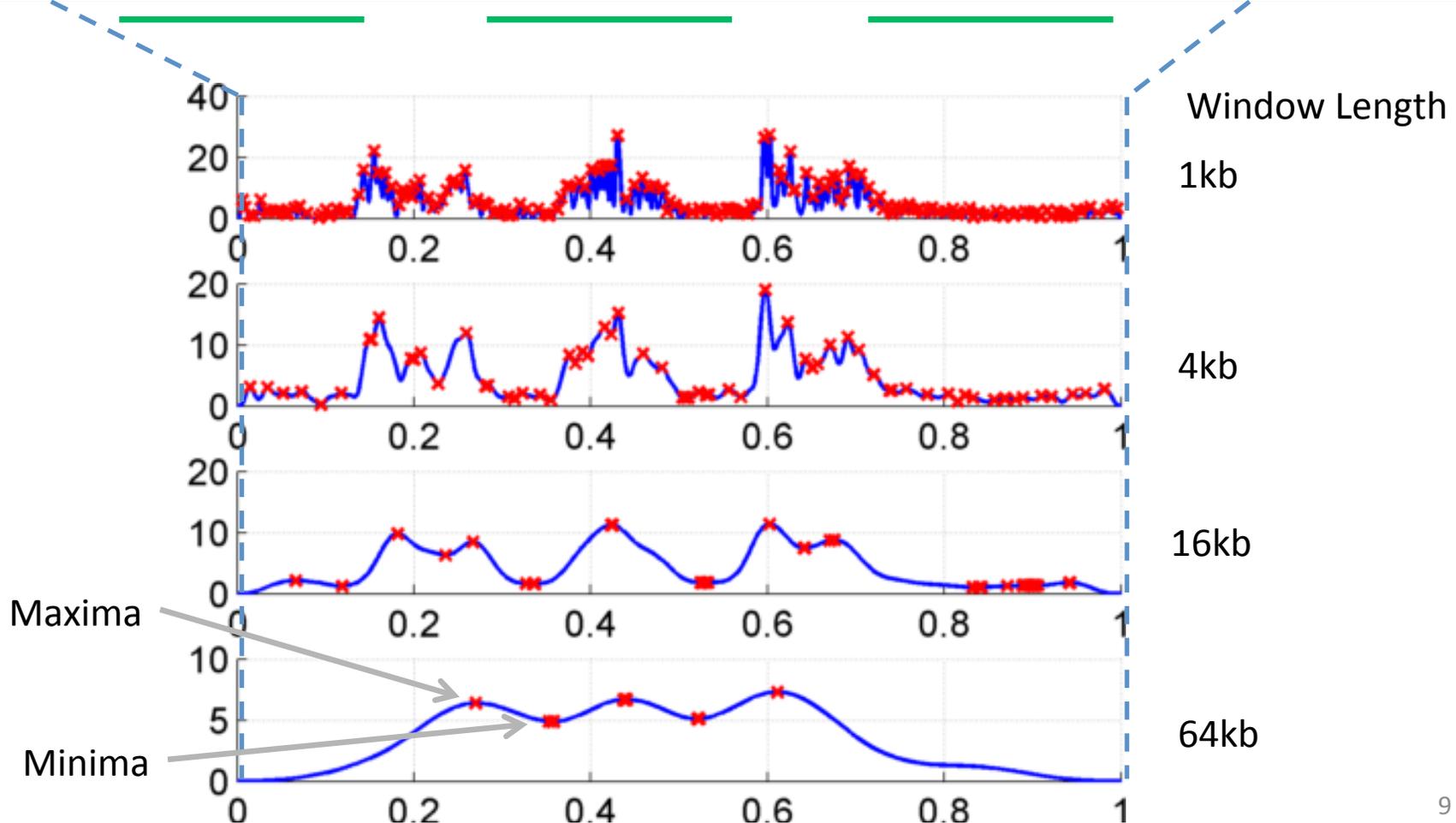
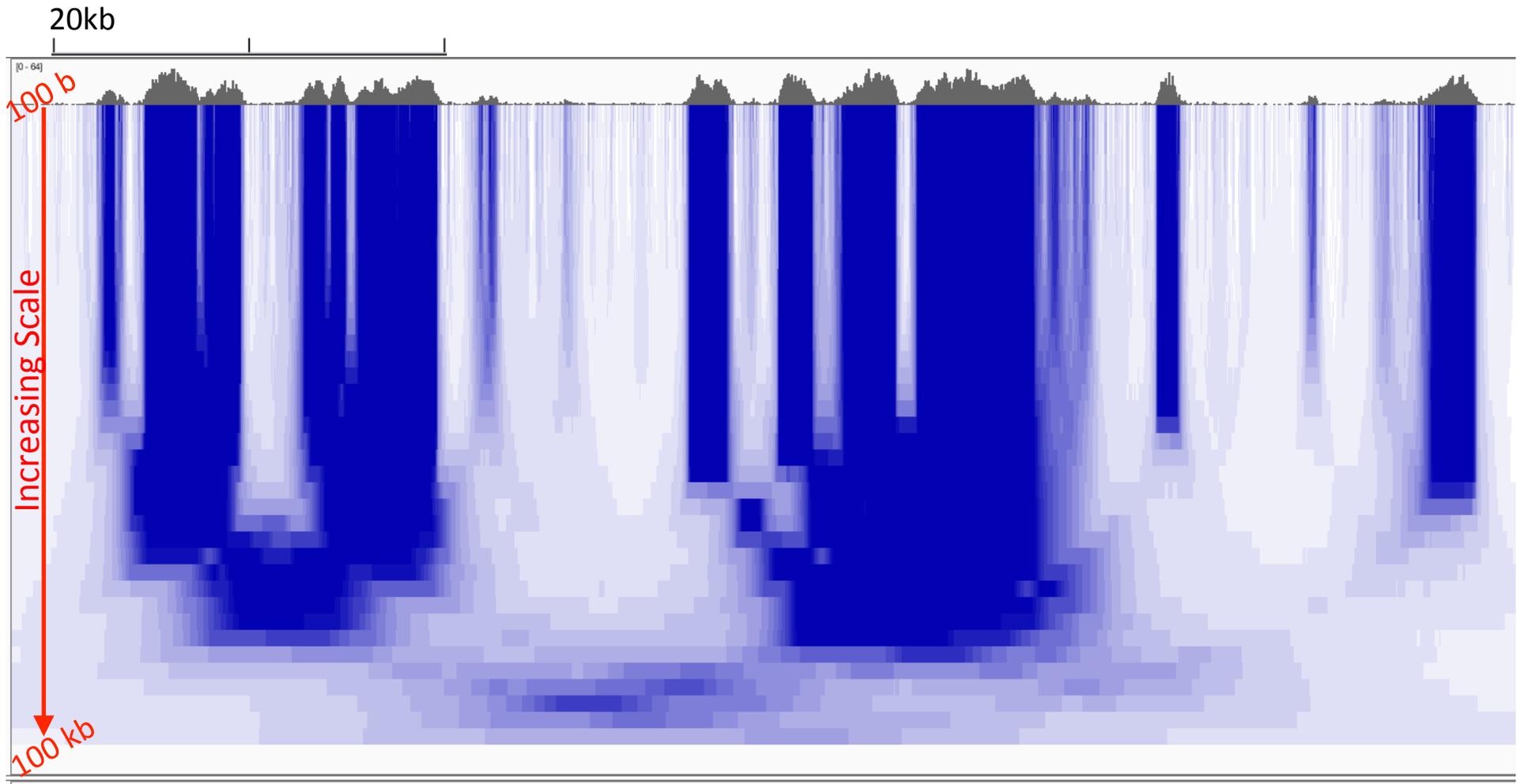# Multiscale Analysis, Minima/Maxima based Coarse Segmentation



*Harmanci et al, Genome Biology 2014, MUSIC.gersteinlab.org*

Window Length

1kb

4kb

16kb

64kb

Maxima

Minima

# Multiscale Decomposition



20kb

100 b

Increasing Scale

100 kb

[0 - 64]

[Harmanci *et al, Genome Biol.* ('14)]

# Multiscale Decomposition



[Harmanci *et al, Genome Biol.* ('14)]

11

# Finding "Conserved" Sites in the Human Population:

## Negative selection in non-coding elements based on Production ENCODE & 1000G Phase 1



**Broad Categories**

Coding

Genomic Avg

Enhancer

(Non-coding RNA) ncRNA

(DNase I hypersensitive sites) DHS

(Transcription factor binding sites) TFBS
- TFSS (TFSS: Sequence-specific TFs)
- General
- Chromatin

Pseudogene

0.56   0.58   0.60   0.62   0.64   0.66   0.68

Fraction of rare SNPs

**Depletion of Common Variants in the Human Population**

- Broad categories of regulatory regions under negative selection
- Related to:

ENCODE, *Nature*, 2012
Ward & Kellis, *Science*, 2012
Mu et al, *NAR*, 2011

[Khurana et al., *Science* ('13)]

# Differential selective constraints among specific sub-categories



**A** Broad Categories

| Category | SNPs |
|---|---|
| Genomic Avg | 27M SNPs |
| Coding | 0.27M |
| Missense | 0.15M |
| Synonymous | 0.12M |
| UTR | 0.4M |
| Enhancer | 1.4M |
| DHS | 4.8M |
| TFSS | 3.7M |
| General | 0.8M |
| Chromatin | 1.2M |
| Pseudogene | 57K |
| ncRNA | 38K |

Fraction of rare SNPs

**B** Specific Categories

TF Families (motifs)

Coding, HMG, Forkhead, bZIP, STAT, MADs-box, NR, Homeodomain, p53, IPT/TIG, ZNF, ETS, HLH, AP2, wHTH, CBF-NFY

Sub-categorization possible because of better statistics from 1000G phase 1 v pilot

[Khurana et al., *Science* ('13)]

# Defining Sensitive non-coding Regions

~0.4% genomic coverage (~ top 25)

~0.02% genomic coverage (top 5)

**A** Broad Categories

**B** Specific Categories

TF Families (motifs)

Fraction of rare SNPs

Start **677** high-resolution non-coding categories; Rank & find those under strongest selection

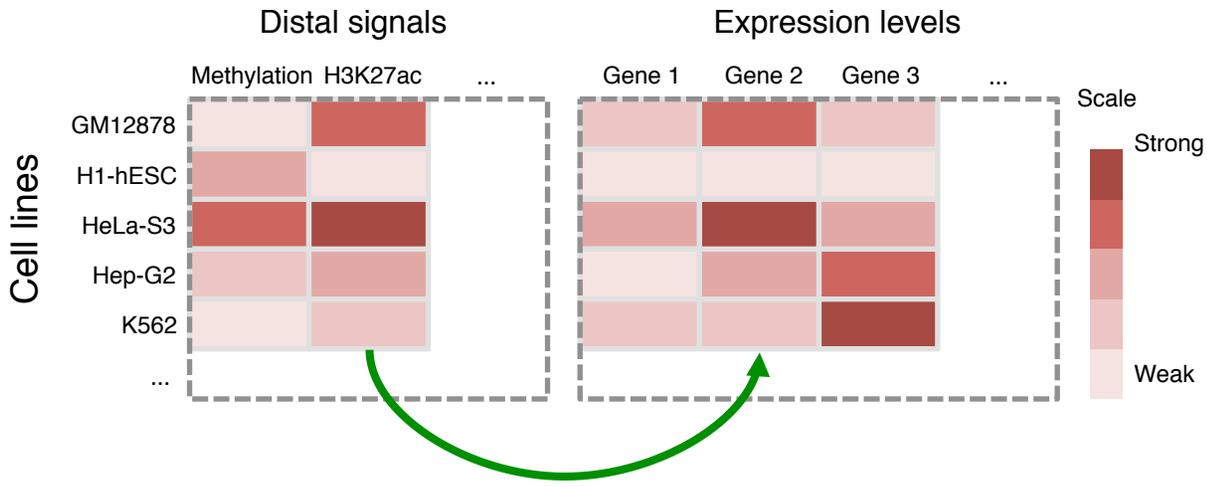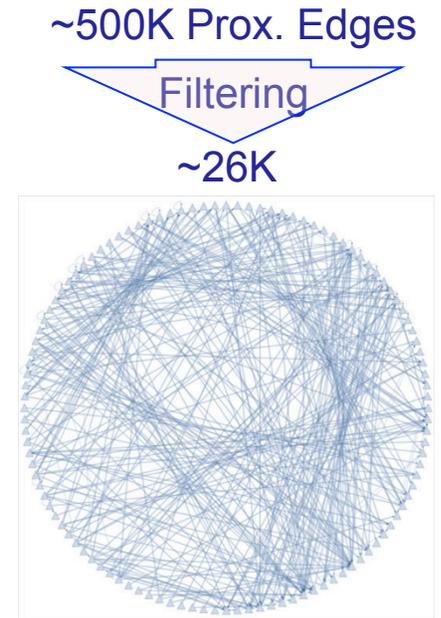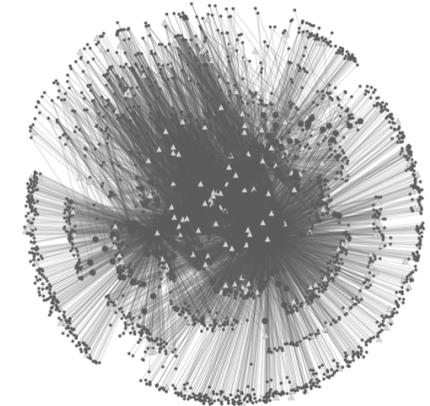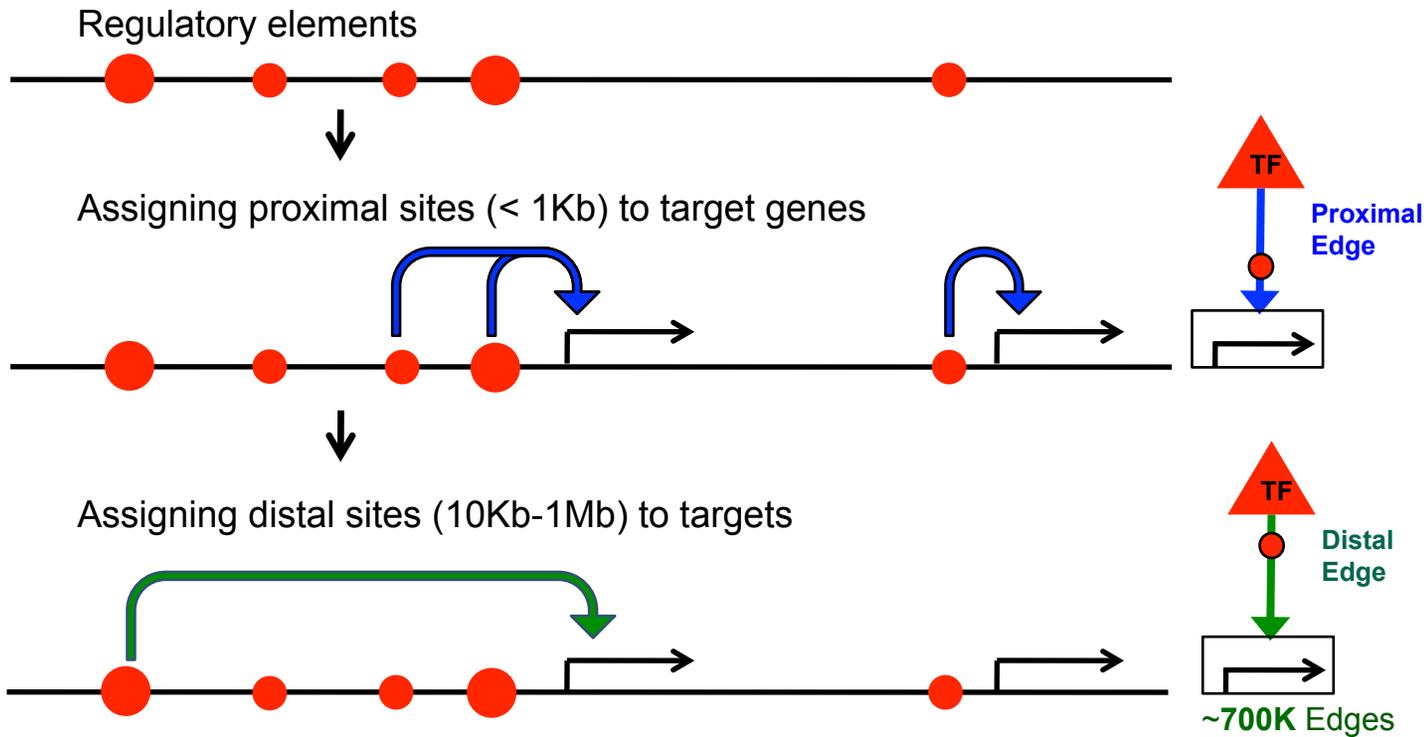Sub-categorization possible because of better statistics from 1000G phase 1 v pilot

[Khurana et al., *Science* ('13)]

[Khurana et al., *Science* ('13)]

- **Finding Non-coding Regions Sensitive to Mutations**
  - **1st Level Linear Annotation: Regulatory Sites**
    - Multi-scale "site" calling (with Music)
    - Finding small number of sites particularly sensitive to mutations
  - **2nd Level Network Annotation**
    - Building a network from the linear annotation
    - More connectivity = more constraint => highlights hubs
- **Using this to Interpret Alterations in Cancer**
  - **LARVA: to find recurrently mutated annotations**
    - Need to correct for overdispersion in bionomial
    - Use beta-bin parameterized according to replication timing
  - **FunSeq software tool for mutation prioritization**
    - Systematically weighting all the features, for non-coding prioritization

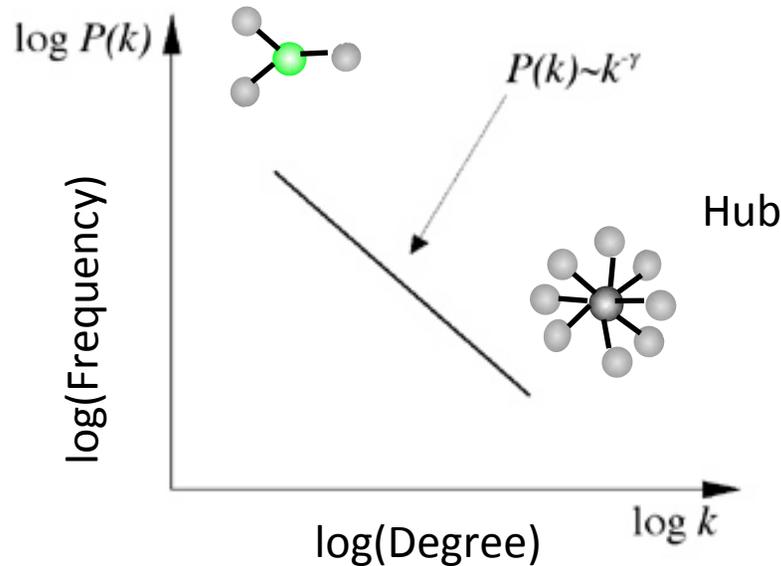# Relating Non-coding Annotation to Protein-coding Genes via Networks

[ Cheng et al., *Bioinfo.* ('11), Gerstein et al., *Nature* ('12), Yip et al., *GenomeBiology* ('12), Fu et al., *GenomeBiology*('14) ]

Regulatory elements

Assigning proximal sites (< 1Kb) to target genes

**TF**

**Proximal Edge**

Assigning distal sites (10Kb-1Mb) to targets

**TF**

**Distal Edge**

~**700K** Edges

~500K Prox. Edges

Filtering

~26K

## Distal signals

Methylation  H3K27ac  ...

## Expression levels

Gene 1  Gene 2  Gene 3  ...

Cell lines
- GM12878
- H1-hESC
- HeLa-S3
- Hep-G2
- K562
- ...

Scale
- Strong
- Weak

Connecting Distal Elements via **Activity Correlations**.

Other strategies to create linkage incl. eQTL and Hi-C. Much in recent Epigenomics Roadmap.
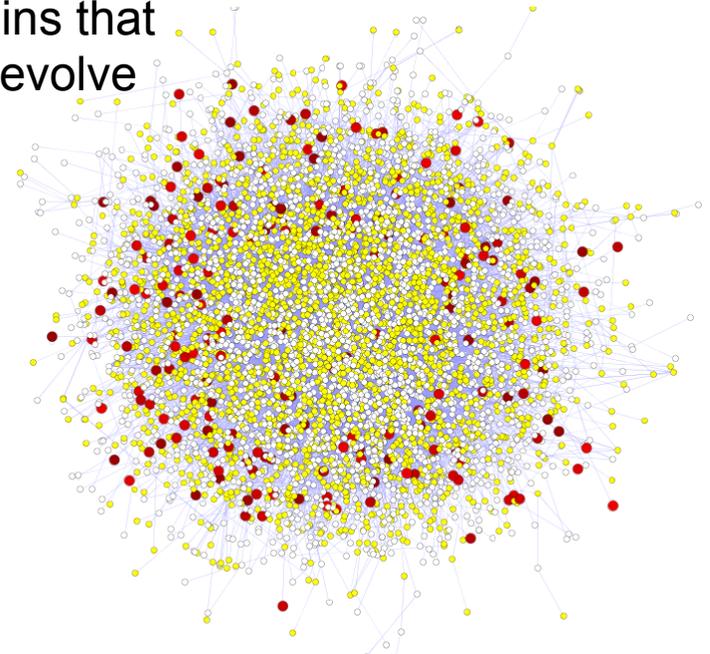
## Power-law distribution



log $P(k)$

$P(k) \sim k^{\gamma}$

log(Frequency)

Hub

log(Degree)

log $k$

## Hubs Under Constraint: A Finding from the Network Biology Community

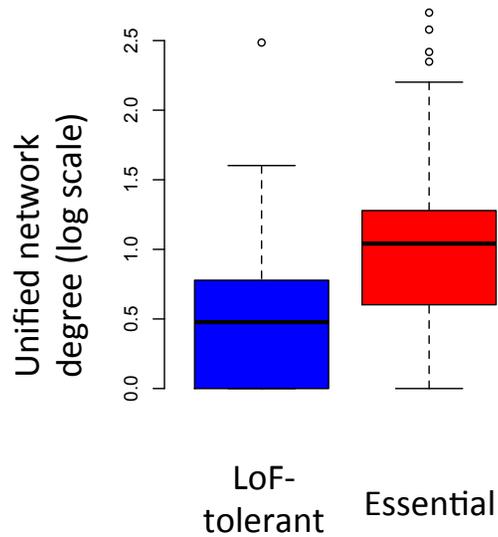- High likelihood of positive selection
- Lower likelihood of positive selection
- Not under positive selection
- No data about positive selection

[Nielsen et al. *PLoS Biol.* (2005), HPRD, Kim et al. PNAS (2007)]

- **More Connectivity, More Constraint:** Genes & proteins that have a more central position in the network tend to evolve more slowly and are more likely to be essential.

- This phenomenon is observed in **many organisms & different kinds of networks**
  - **yeast PPI** - Fraser et al ('02) Science, ('03) BMC Evo. Bio.
  - **Ecoli PPI** - Butland et al ('04) Nature
  - **Worm/fly PPI** - Hahn et al ('05) MBE
  - **miRNA net** - Cheng et al ('09) BMC Genomics
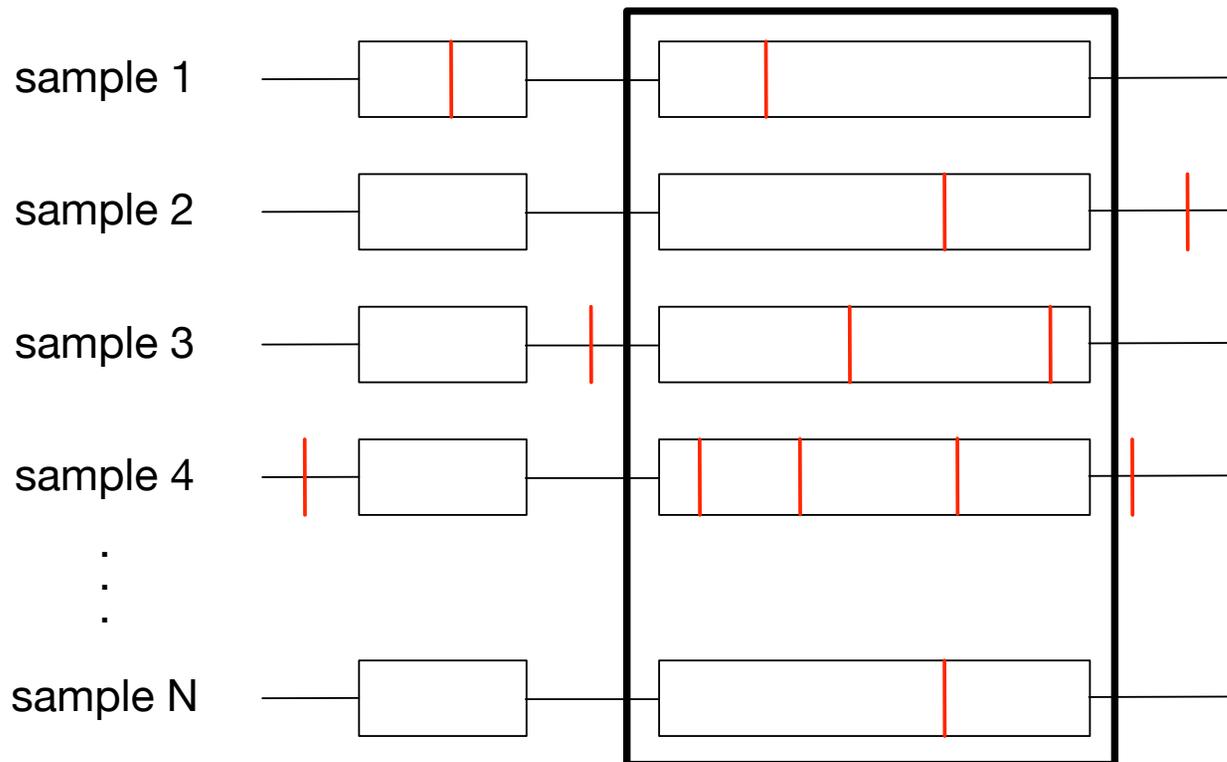
# Regulatory Hubs
# are more Essential



Unified network degree (log scale)

LoF-tolerant    Essential

Proximal Regulatory Network

Total degree (IN + OUT) (log scale)

Wilcoxon pvalue=1.29e-2

LoF-tolerant    Essential

LoF-tolerant genes    Essential genes

Size of nodes scaled by total degree

[Khurana et al., *PLOS Comp. Bio.* '13]

- **Finding Non-coding Regions Sensitive to Mutations**
  - **1st Level Linear Annotation: Regulatory Sites**
    - Multi-scale "site" calling (with Music)
    - Finding small number of sites particularly sensitive to mutations
  - **2nd Level Network Annotation**
    - Building a network from the linear annotation
    - More connectivity = more constraint => highlights hubs
- **Using this to Interpret Alterations in Cancer**
  - **LARVA: to find recurrently mutated annotations**
    - Need to correct for overdispersion in bionomial
    - Use beta-bin parameterized according to replication timing
  - **FunSeq software tool for mutation prioritization**
    - Systematically weighting all the features, for non-coding prioritization
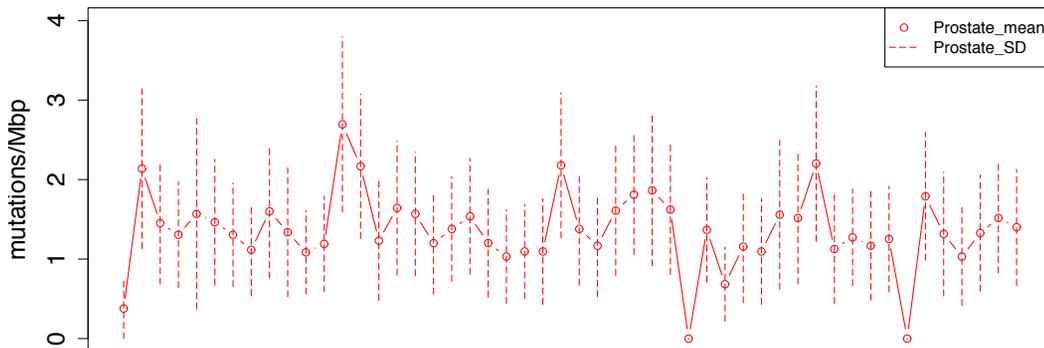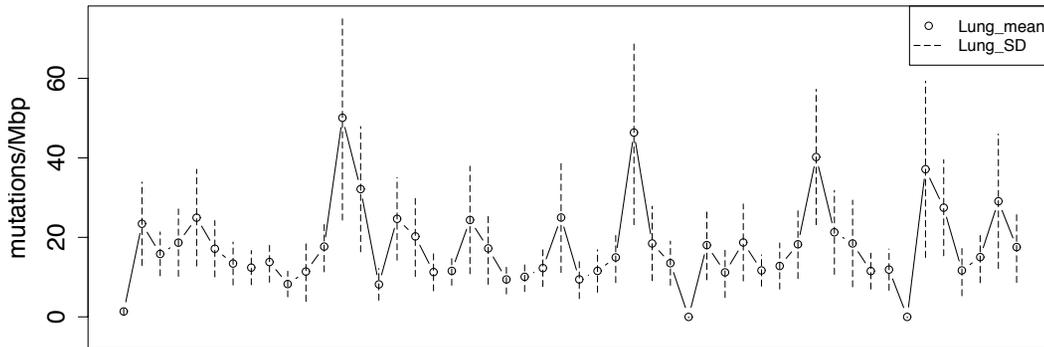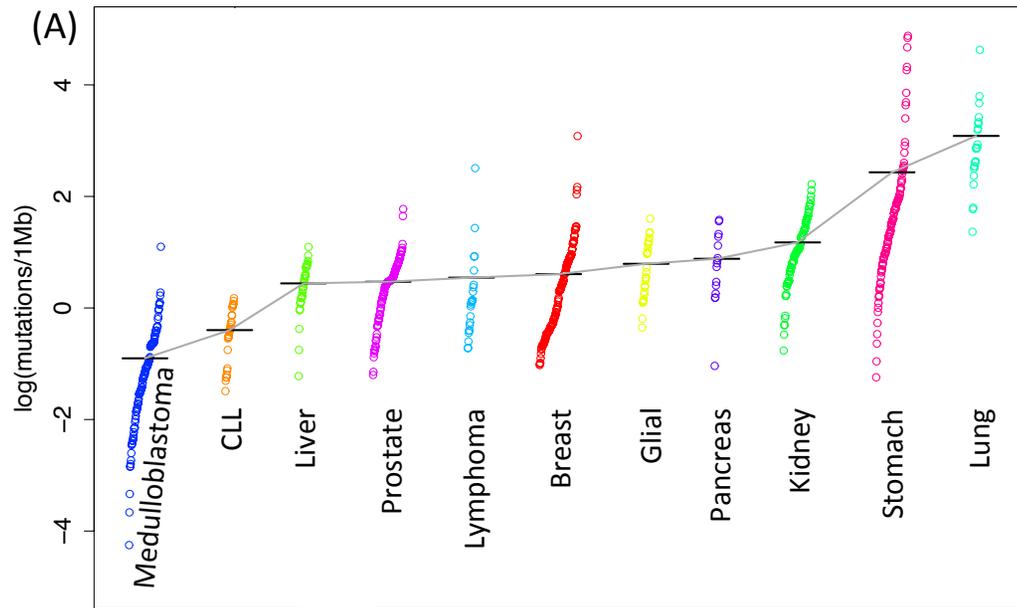
# LARVA

- Somatic single nucleotide variant (SNV) data
  - Which functional noncoding genome elements are hotspots for SNVs across multiple samples?
    - Are they mutated more than expected from neutral mutation processes?



[Lochovsky et al. NAR ('15, in press)]

# Cancer Somatic Mutational Heterogeneity

- The distribution of variants throughout the genome indicates high mutation rate heterogeneity between samples of the same cancer type, and on many other levels

- **Goal:** Develop a model for the whole genome background somatic mutation distribution in cancer to identify potential noncoding cancer driving elements

- LARVA: <u>La</u>rge-scale <u>A</u>nalysis of <u>R</u>ecurrent <u>V</u>ariants in noncoding <u>A</u>nnotations

[Lochovsky et al. NAR ('15, in press)]

# Cancer Somatic Mutation Modeling

- We tested 3 models evaluating the significance of a mutation burden of a genome element

- Suppose there are $k$ genome elements. For element $i$, define:
  - $n_i$: total number of nucleotides in $i$
  - $x_i$: the number of mutation within element $i$
  - $p$: the probability of observing a mutation in each position
  - $R$: The replication timing tenth percentile of $i$

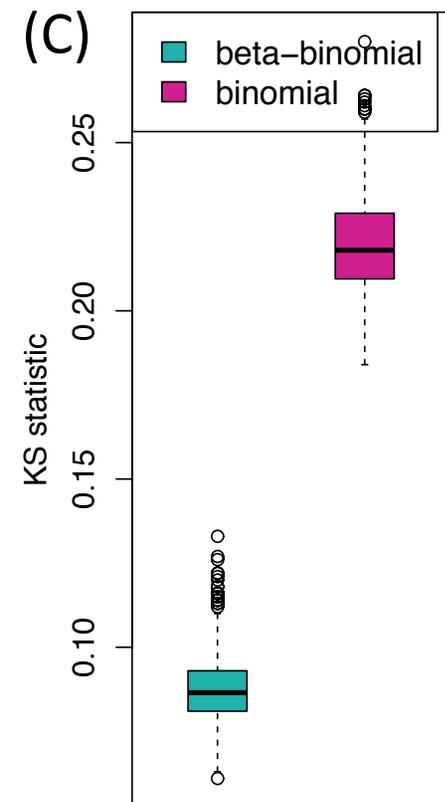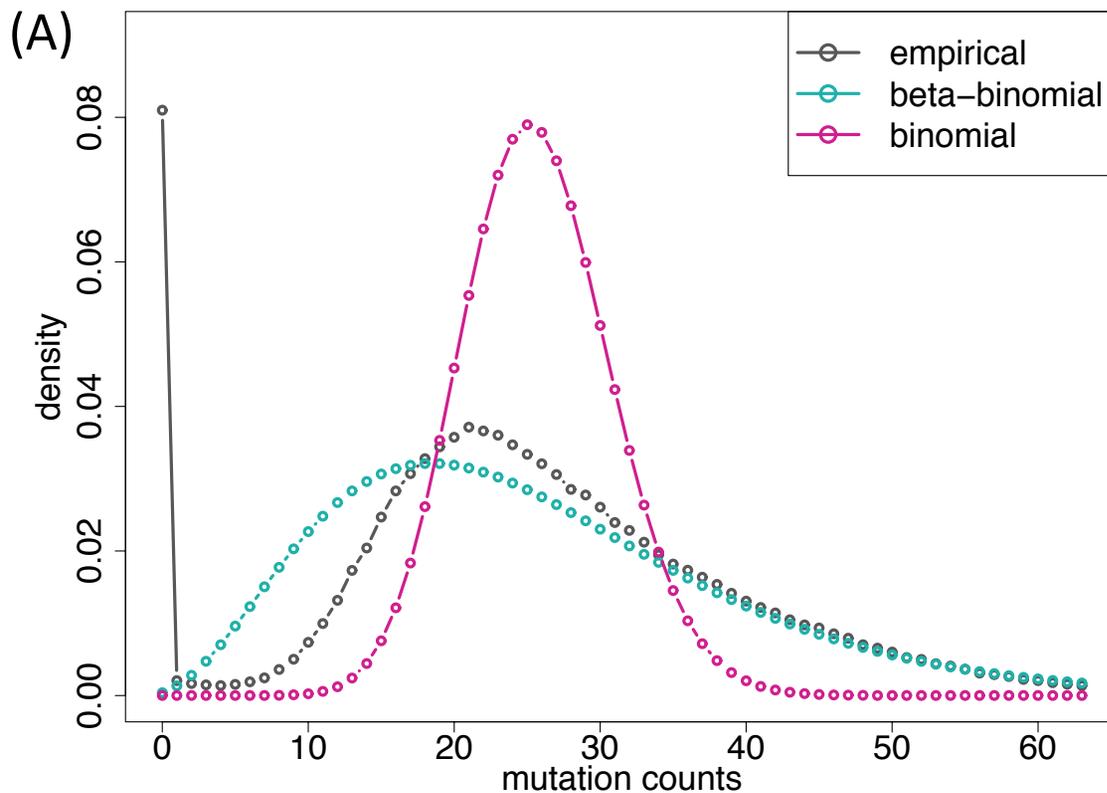| |
|---|
| **Model 1: Constant Background Mutation Rate (Model from Previous Work[1])** <br> $x_i : Binomial(n_i, p)$ |
| **Model 2: Varying Mutation Rate** <br> $x_i \mid p : Binomial(n_i, p)$ <br><br> $p : Beta(\mu, \sigma)$ |
| **Model 3: Varying Mutation Rate with Replication Timing Correction** <br> $x_i \mid p : Binomial(n_i, p)$ <br><br> $p : Beta(\mu \mid R, \sigma \mid R)$ <br><br> $\mu \mid R, \sigma \mid R : \text{constant within the same } R \text{ bin}$ |

[Lochovsky et al. NAR ('15, in press)]

1. Weinhold, N., Jacobsen, A., Schultz, N., Sander, C. & Lee, W. Genome-wide analysis of noncoding regulatory mutations in cancer. *Nature Genetics* **46,** 1160–1165 (2014).
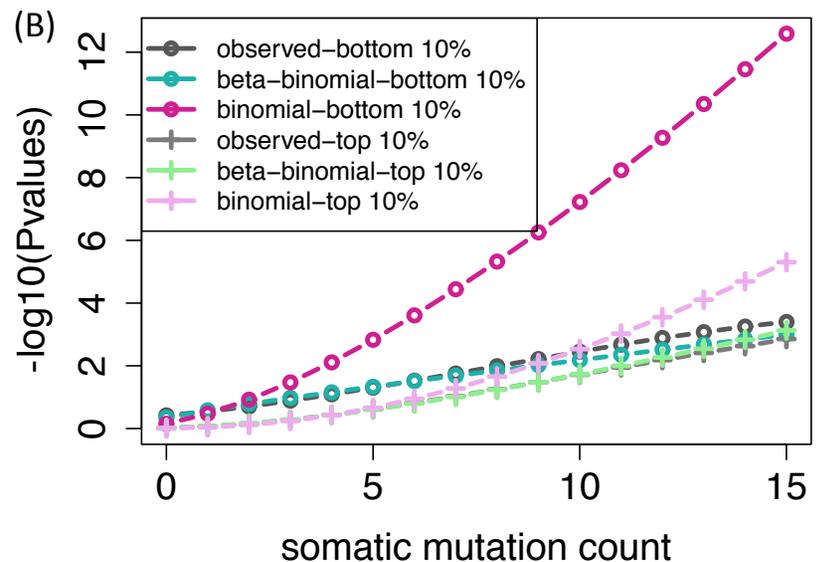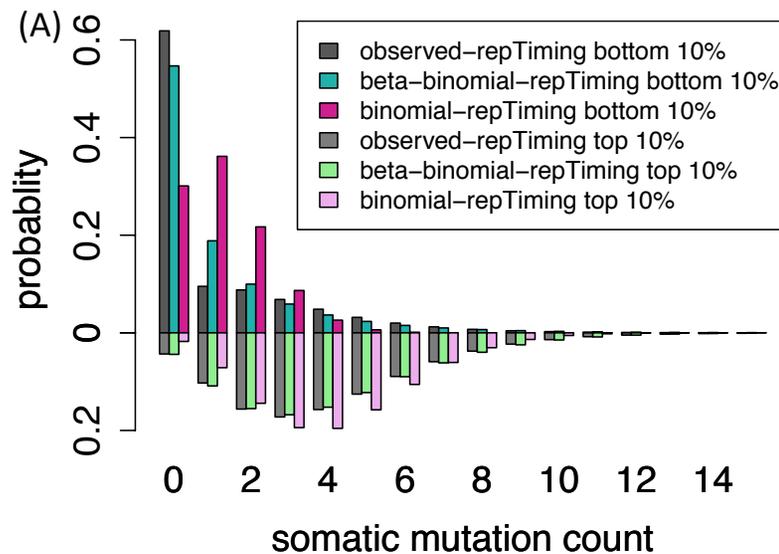
# LARVA Model Comparison

- Comparison of mutation count frequency implied by the binomial model (model 1) and the beta-binomial model (model 2) relative to the empirical distribution

- The beta-binomial distribution is significantly better, especially for accurately modeling the over-dispersion of the empirical distribution



[Lochovsky et al. NAR ('15, in press)]
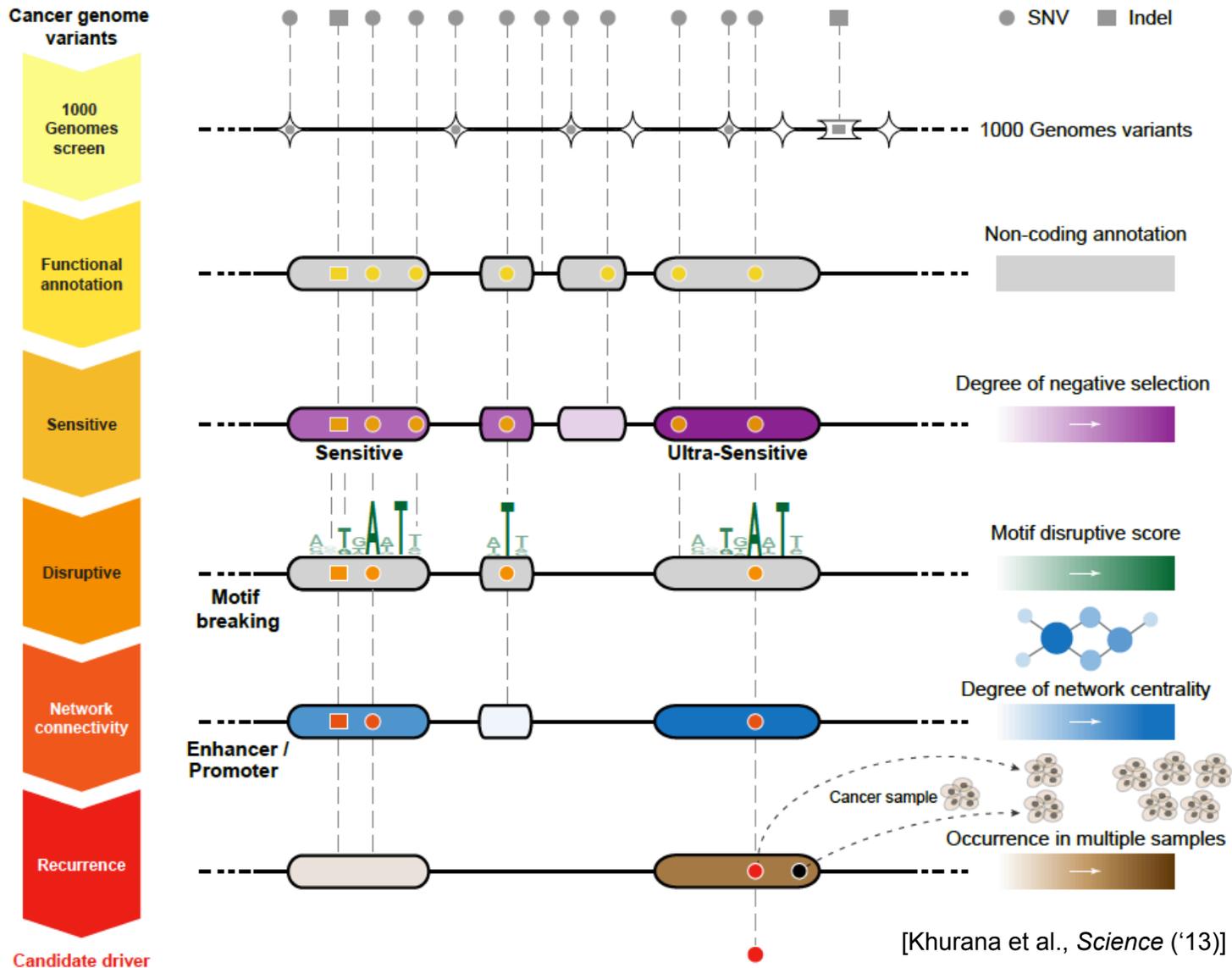
# LARVA Model Comparison

- Demonstrate that adding the DNA replication timing correction (model 3) further improves the beta-binomial model (model 2)

- Top 10% of replication timing bins requires little correction

- Bottom 10% of replication timing bins requires massive correction

- Demonstrate that the number of significant p-values is inflated under the binomial model

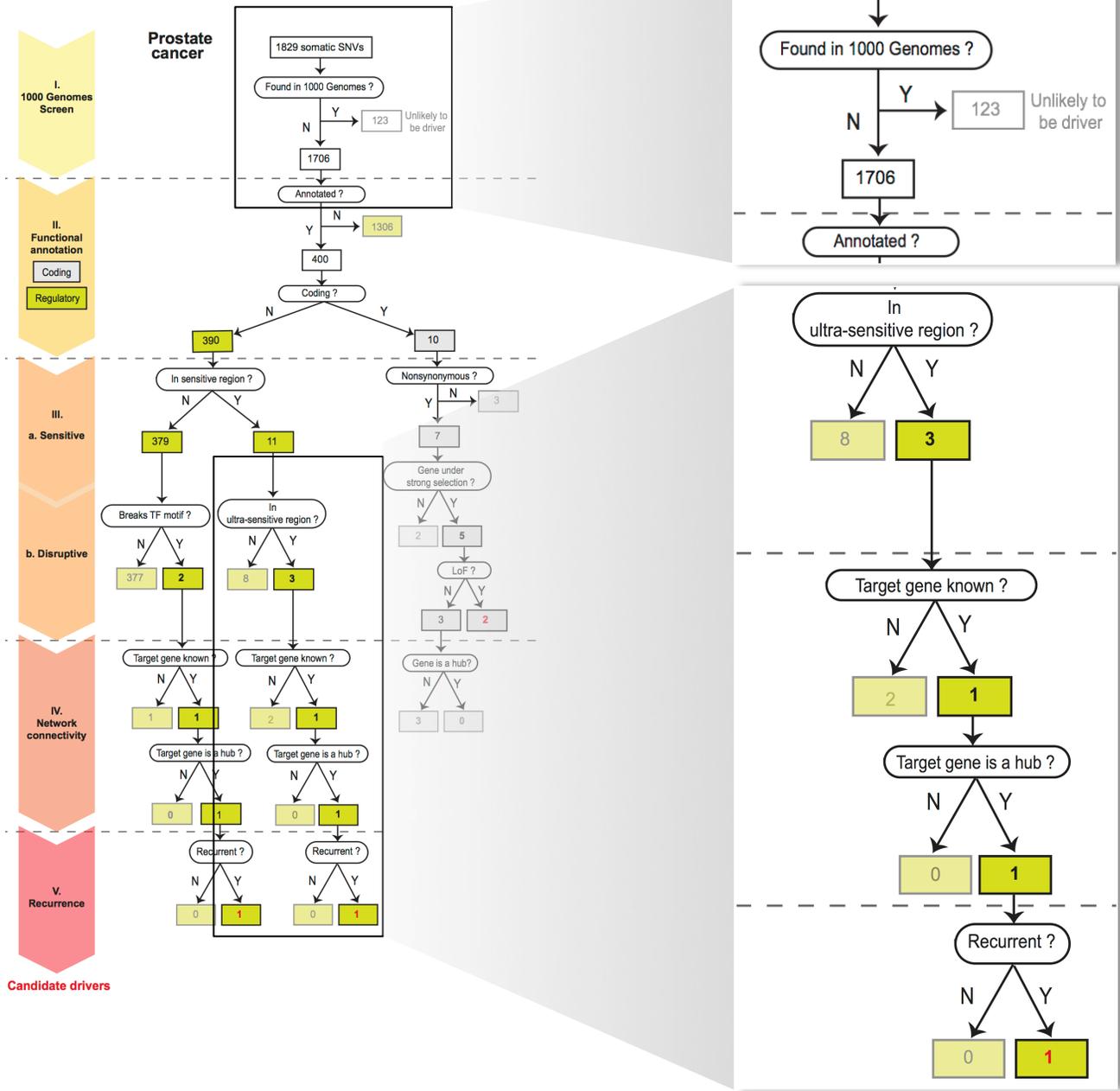- Neither the empirical or beta-binomial models exhibit this inflation



(A) Legend: observed–repTiming bottom 10%; beta–binomial–repTiming bottom 10%; binomial–repTiming bottom 10%; observed–repTiming top 10%; beta–binomial–repTiming top 10%; binomial–repTiming top 10%. Axes: probablity vs somatic mutation count.

(B) Legend: observed–bottom 10%; beta–binomial–bottom 10%; binomial–bottom 10%; observed–top 10%; beta–binomial–top 10%; binomial–top 10%. Axes: -log10(Pvalues) vs somatic mutation count.

[Lochovsky et al. NAR ('15, in press)]

- **Finding Non-coding Regions Sensitive to Mutations**
  - **1st Level Linear Annotation: Regulatory Sites**
    - Multi-scale "site" calling (with Music)
    - Finding small number of sites particularly sensitive to mutations
  - **2nd Level Network Annotation**
    - Building a network from the linear annotation
    - More connectivity = more constraint => highlights hubs
- **Using this to Interpret Alterations in Cancer**
  - **LARVA: to find recurrently mutated annotations**
    - Need to correct for overdispersion in bionomial
    - Use beta-bin parameterized according to replication timing
  - **FunSeq software tool for mutation prioritization**
    - Systematically weighting all the features, for non-coding prioritization

# Identification of non-coding candidate drivers amongst somatic variants: Scheme



[Khurana et al., *Science* ('13)]

# Flowchart for 1 Prostate Cancer Genome
## (from Berger et al. '11)

Site integrates user variants with large-scale context

**Data Context**

**Variant Prioritization**

Weighted scoring scheme

Highlighting variants

User Cancer Variants

Variant Reports

FunSeq.gersteinlab.org

[Fu et al., GenomeBiology ('14)]

28

- Feature weight

  - Weighted with mutation patterns in natural polymorphisms

    (features frequently observed weight less)

  - entropy based method

| | HOT region | ▬ |
| --- | --- | --- |
| | Sensitive region | ▬ |
| | Polymorphisms | │ |

Genome

- Feature weight

  - Weighted with mutation patterns in natural polymorphisms

    (features frequently observed weight less)

  - entropy based method

HOT region

Sensitive region

Polymorphisms

Genome

$$p = \frac{3}{20}$$

[Fu et al., GenomeBiology ('14)]

■ Feature weight

- Weighted with mutation patterns in natural polymorphisms

(features frequently observed weight less)
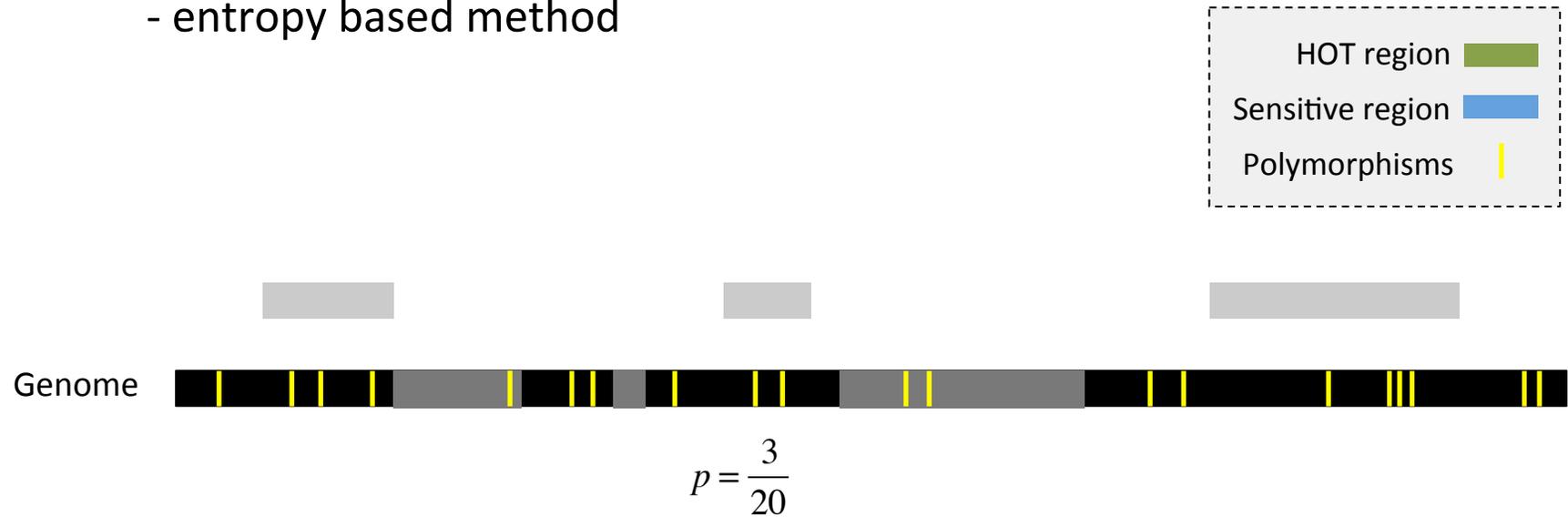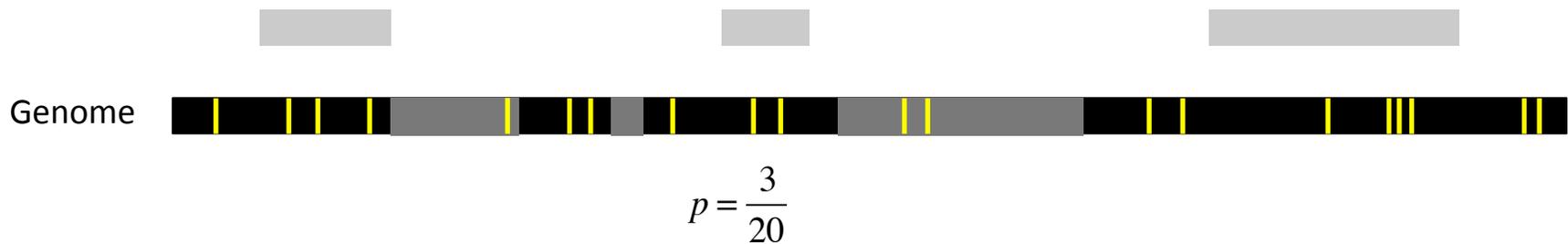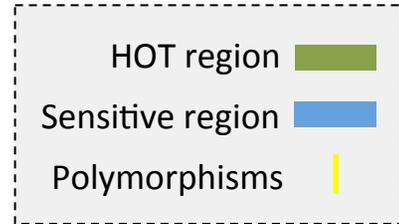
- entropy based method

| | |
|---|---|
| HOT region | ▬ |
| Sensitive region | ▬ |
| Polymorphisms | │ |

Genome

$$p = \frac{3}{20}$$

*Feature weight:* $w_d = 1 + p_d log_2 p_d + (1 - p_d) log_2 (1 - p_d)$

$p \uparrow \quad w_d \downarrow \qquad$ *p = probability of the feature overlapping natural polymorphisms*

*For a variant:* $Score = \sum w_d \;\; of\; observed\; features$

31

[Fu et al., GenomeBiology ('14)]

- **Finding Non-coding Regions Sensitive to Mutations**
  - **1st Level Linear Annotation: Regulatory Sites**
    - Multi-scale "site" calling (with Music)
    - Finding small number of sites particularly sensitive to mutations
  - **2nd Level Network Annotation**
    - Building a network from the linear annotation
    - More connectivity = more constraint => highlights hubs
- **Using this to Interpret Alterations in Cancer**
  - **LARVA: to find recurrently mutated annotations**
    - Need to correct for overdispersion in bionomial
    - Use beta-bin parameterized according to replication timing
  - **FunSeq software tool for mutation prioritization**
    - Systematically weighting all the features, for non-coding prioritization

# Cancer Prioritzation Acknowledgements

~50 people ← ~1000 "authors"

### Functional Interpretation Subgroup


A THOUSAND GENOMES

## Yale

**Ekta Khurana, Yao Fu, Jieming Chen,**
**Xinmeng Mu**, Lucas Lochovsky,
Arif Harmanci, Alexej Abyzov,
Suganthi Balasubramanian, Cristina Sisu,
Declan Clarke, Mike Wilson

## Sanger

**Vincenza Colonna**, Yali Xue,
**Chris Tyler-Smith**

## Cornell

Steven Lipkin, Jishnu Das, Robert Fragoza, Xiaomu Wei, **Haiyuan Yu**

Andrea Sboner, Dimple Chakravarty, Naoki Kitabayashi, Vaja Liluashvili, Zeynep H. Gümüş, **Mark A. Rubin**

## US, UK, Switzerland….

**Hyun Min Kang, Tuuli Lappalainen,** Kathryn Beal, Daniel Challis, Yuan Chen, Laura Clarke, Fiona Cunningham, Emmanouil T. Dermitzakis, Uday Evani, Paul Flicek, Erik Garrison, Javier Herrero, Yong Kong, Kasper Lage, Daniel G. MacArthur, Gabor Marth, Donna Muzny, Tune H. Pers, Graham R. S. Ritchie, Jeffrey A. Rosenfeld, Fuli Yu, Richard Gibbs

33

## Acknowledgements

# Info about content in this slide pack

- General PERMISSIONS

  – This Presentation is copyright Mark Gerstein, Yale University, 2014.

  – Please read permissions statement at http://www.gersteinlab.org/misc/permissions.html .

  – Feel free to use slides & images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).

  – Paper references in the talk were mostly from Papers.GersteinLab.org.

- For SeqUniverse slide, please contact Heidi Sofia, NHGRI

- PHOTOS & IMAGES. For thoughts on the source and permissions of many of the photos and clipped images in this presentation see http://streams.gerstein.info .

  – In particular, many of the images have particular EXIF tags, such as  kwpotppt , that can be easily queried from flickr, viz: http://www.flickr.com/photos/mbgmbg/tags/kwpotppt