# Personal Genomics: Managing Rapid Data Scaling through Prioritizing High-impact Variants

**Mark Gerstein, Yale**
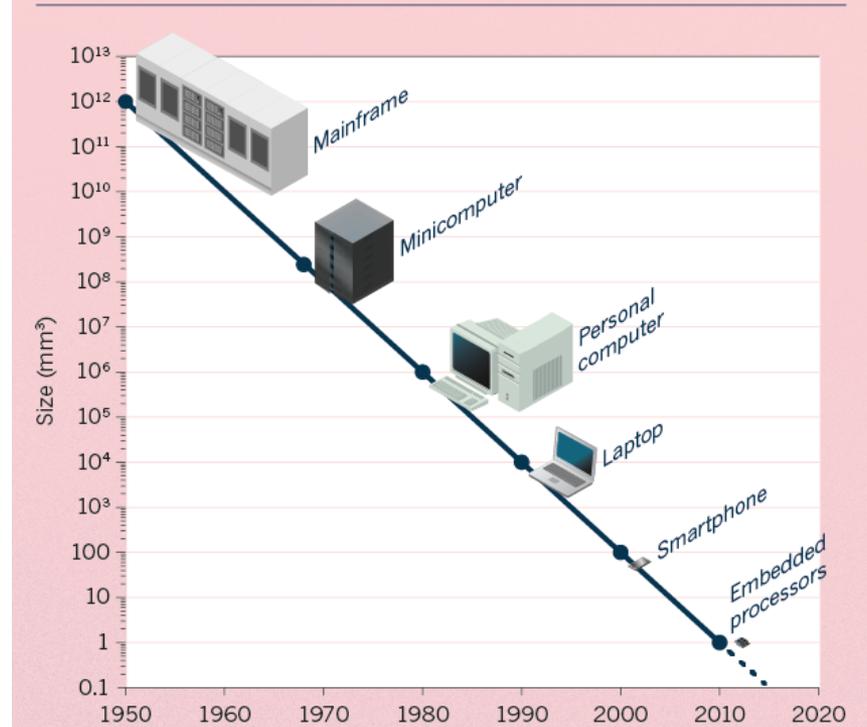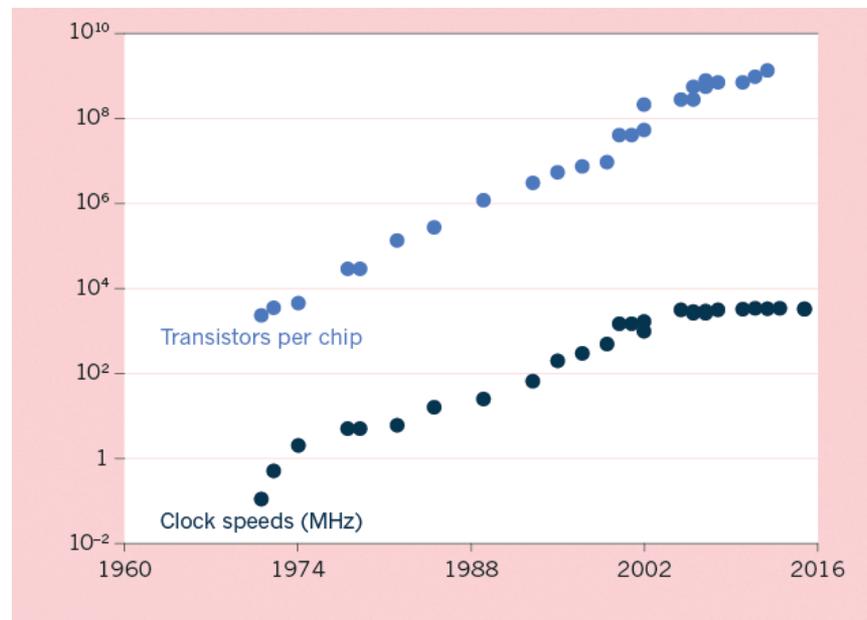
**Slides freely downloadable from Lectures.GersteinLab.org**

**& "tweetable" (via @markgerstein). See last slide for more info.**
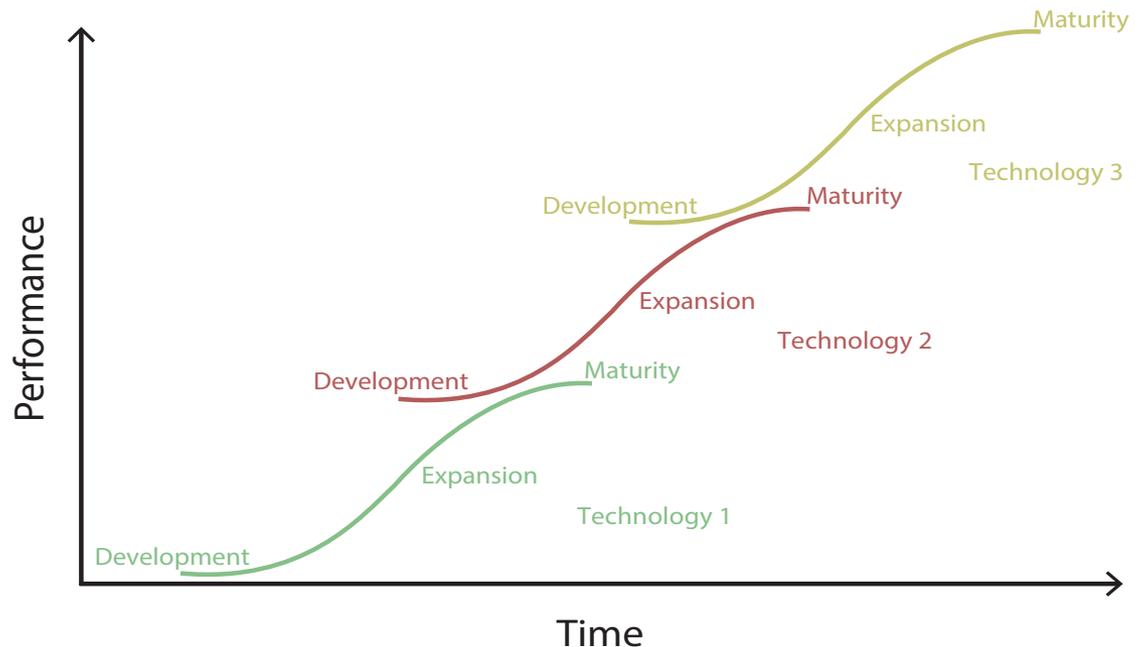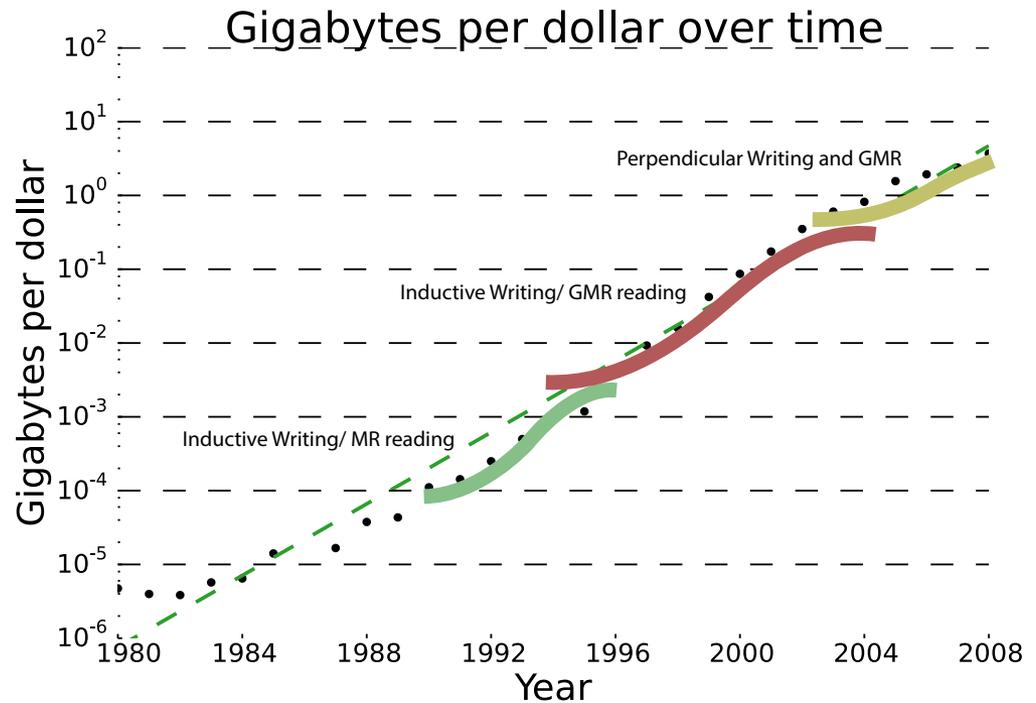
# Moore's Law: Exponential Scaling of Computer Technology



- Exponential increase in the number of transistors per chip.

- Led to improvements in speed and miniaturization.

- Drove widespread adoption and novel applications of computer technology.
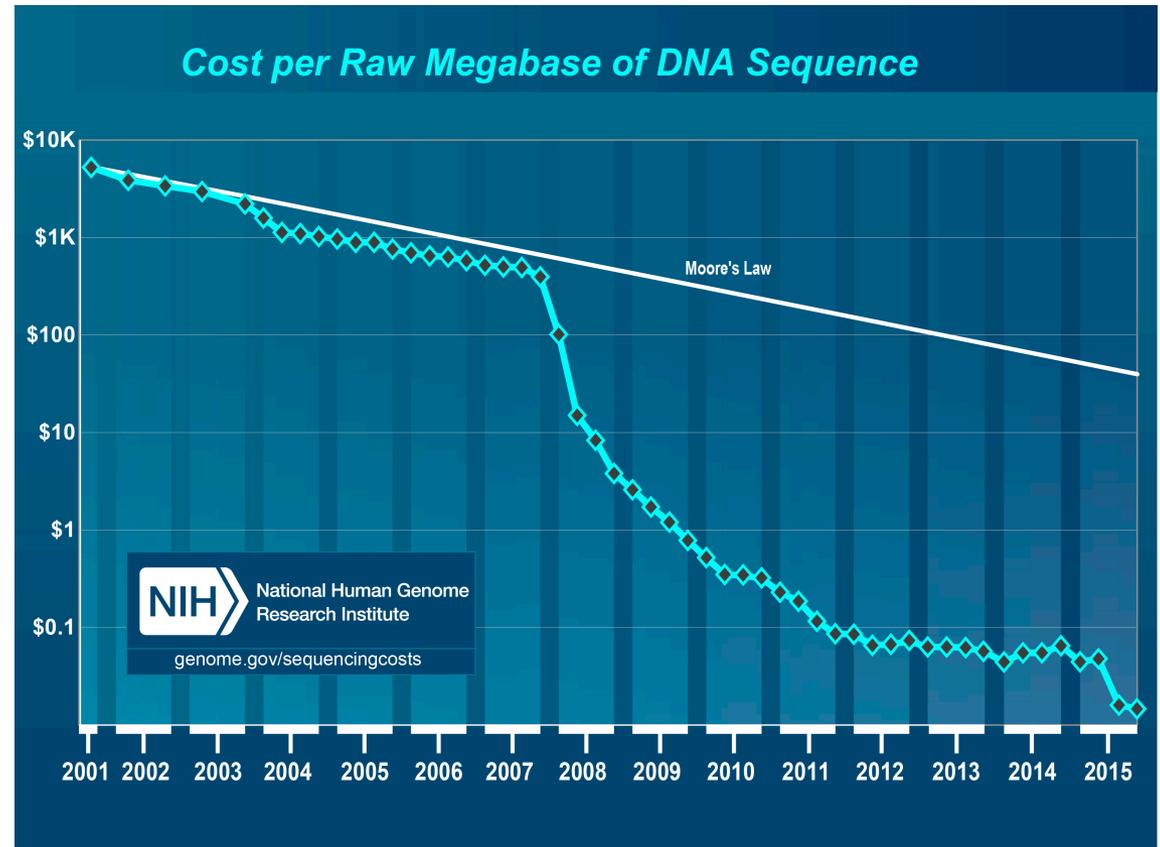
# Kryder's Law and S-curves underlying exponential growth

- Moore's & Kryder's Laws
  - As important as the increase in computer speed has been, the ability to store large amounts of information on computers is even more crucial

- Exponential increase seen in Kryder's law is a superposition of S-curves for different technologies

**Gigabytes per dollar over time**



**Time**

# Sequencing Data Explosion:
## Faster than Moore's Law for a Time (or a S-curve)
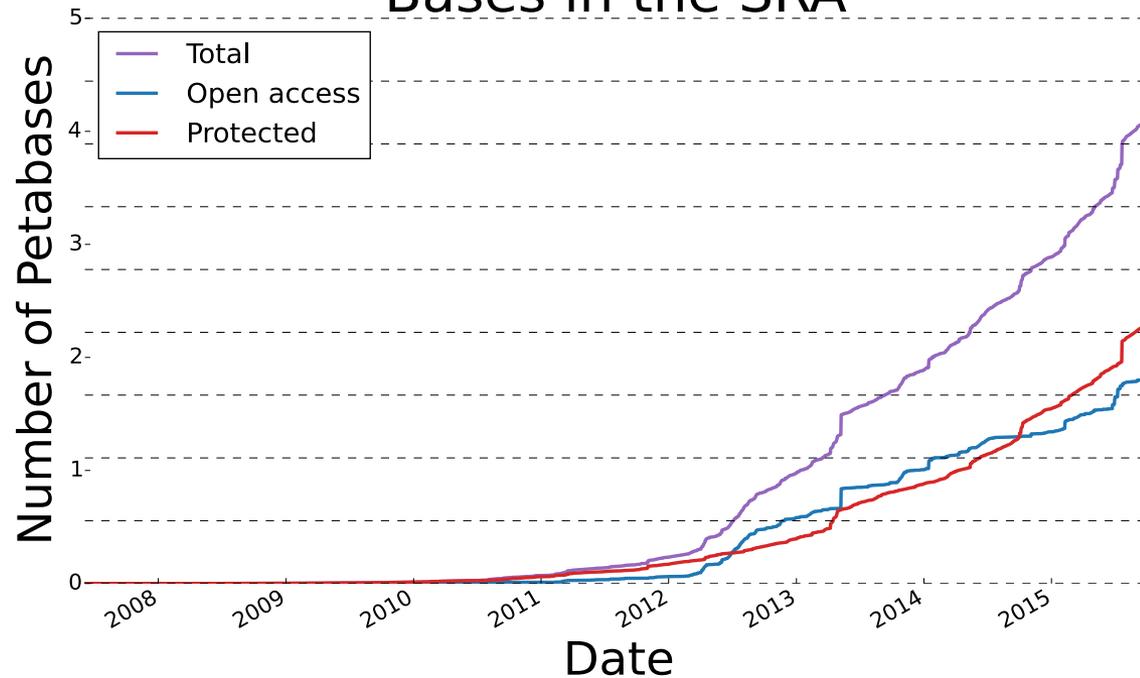
- DNA sequencing has gone through technological S-curves
  - In the early 2000's, improvements in Sanger sequencing produced a scaling pattern similar to Moore's law.

  - The advent of NGS was a shift to a new technology with dramatic decrease in cost).



Cost per Raw Megabase of DNA Sequence

Moore's Law

NIH> National Human Genome Research Institute

genome.gov/sequencingcosts

4

# Sequencing cost reductions have resulted in an explosion of data

- The type of sequence data deposited has changed as well.
  - Protected data represents an increasing fraction of all submitted sequences.

  - Data from techniques utilizing NGS machines has replaced that generated via microarray.

## Bases in the SRA



Legend:
- Total
- Open access
- Protected

X-axis: Date (2008–2015)
Y-axis: Number of Petabases (0–5)

## NIH Funding for "microarray" and "sequencing" projects



Legend:
- Microarray
- Sequencing

X-axis: Year (1990–2015)
Y-axis: Funding in USD (Billions) (0.5–3.0)

# Sequence Universe

SRA ~1 petabyte

TCGA endpoint: ~2.5 Petabytes
~1.5 PB exome
~1 PB whole genome

1000 Genomes
A Deep Catalog of Human Genetic Variation

RNASeq | Clinical
miRNA
Epigenome

**TCGA**

**2.3**

**Petabytes**
in CGHub

222 TB

68 TB

40 TB

34 TB

32 TB

29 TB

ENCODE

ARRA Autism

NHLBI ESP

GTeX

ADSP

NHGRI LSSP

National Human Genome Research Institute

Star formation
100K Genomes England

Heidi Sofia, 7-16-15

JESS3

# Increasing diversity in sequence data sources

# The changing costs of a sequencing pipeline



From '00 to ~'20,
cost of DNA sequencing expt. shifts from
the actual seq. to sample
collection & analysis

[Sboner et al. ('11), Muir et al. ('15) Genome Biology]

# The changing costs of a sequencing pipeline



From '00 to ~'20,
cost of DNA sequencing expt. shifts from
the actual seq. to sample
collection & analysis

Labor
Instrument depreciation and maintenance
Reagents and supplies
Indirect costs

[Sboner et al. ('11), Muir et al. ('15) Genome Biology]

# The changing costs of a sequencing pipeline



From '00 to ~'20,
cost of DNA sequencing expt. shifts from
the actual seq. to sample
collection & analysis

Alignment algorithms scaling to keep
pace with data generation

[Sboner et al. ('11), Muir et al. ('15) Genome Biology]

# The changing costs of a sequencing pipeline



From '00 to ~'20,
cost of DNA sequencing expt. shifts from
the actual seq. to sample
collection & analysis

Alignment algorithms scaling to keep
pace with data generation

[Sboner et al. ('11), Muir et al. ('15) Genome Biology]
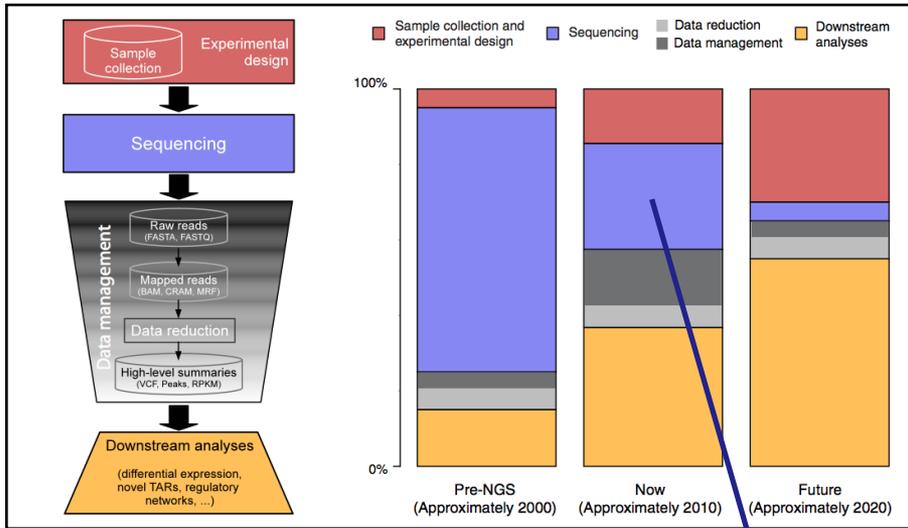
Lectures.GersteinLab.org

11

# The changing costs of a sequencing pipeline
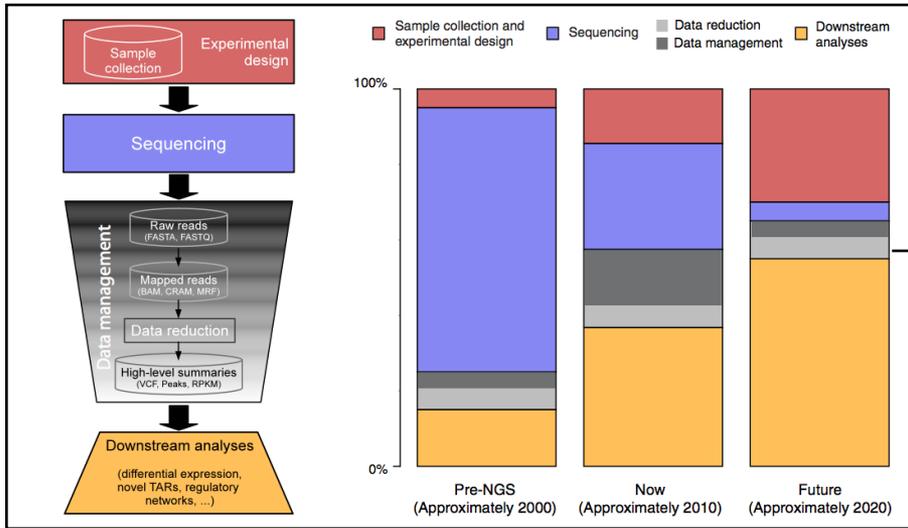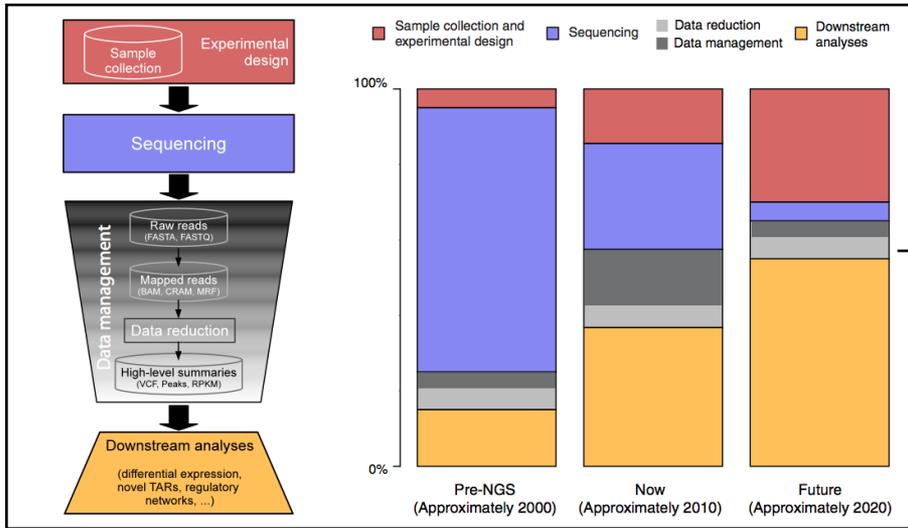


From '00 to ~'20,
cost of DNA sequencing expt. shifts from
the actual seq. to sample
collection & analysis

[Sboner et al. ('11), Muir et al. ('15) Genome Biology]

# Personal Genomics:
## Managing Rapid Data Scaling through Prioritizing High-impact Variants

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

# Human Genetic Variation

**A Cancer Genome**

**A Typical Genome**

**Population of 2,504 peoples**

## Origin of Variants

| | Coding | Non-coding |
|---|---|---|
| Germ-line | 22K | 4.1 – 5M |
| Somatic | ~50 | 5K |

## Class of Variants

| SNP | 3.5 – 4.3M |
|---|---|
| Indel | 550 – 625K |
| SV | 2.1 – 2.5K (20Mb) |
| Total | 4.1 – 5M |

| SNP | 84.7M |
|---|---|
| Indel | 3.6M |
| SV | 60K |
| Total | 88.3M |

Passenger

Driver (~0.1%)

**Prevalence of Variants**

Common

Rare* (1-4%)

Common

Rare (~75%)

* Variants with allele frequency < 0.5% are considered as rare variants in 1000 genomes project.

The 1000 Genomes Project Consortium, Nature. 2015. 526:68-74
Khurana E. et al. Nat. Rev. Genet. 2016. 17:93-108

# Finding Key Variants

## Germline

- **Common variants**
    - Can be associated with phenotype (ie disease) via a Genome-wide Association Study (GWAS), which tests whether the frequency of alleles differs between cases & controls.
    - Usually their functional effect is weaker.
    - Many are non-coding
    - Issue of LD in identifying the actual causal variant.
- **Rare variants**
    - Associations are usually underpowered due to low frequencies.
    - They often have larger functional impact
    - Can be collapsed in the same element to gain statistical power (burden tests).
    - In some cases, causal variants can be identified through tracing inheritance of Mendelian subtypes of diseases in large families.

McCarthy, M. et al. Nat. Rev. Genet. 2008. 9, 356-369, Zuk, O. et al. PNSA. 2014. Vol. 11, no. 4, MacArthur DG et al. Nature 2014. 508:469-476

# Finding Key Variants

## Somatic

- **Overall**
  - Often these can be conceptualized as very rare variants
  - A challenge to identify somatic mutations contributing to cancer is to find driver mutations & distinguish them from passengers.
- **Drivers**
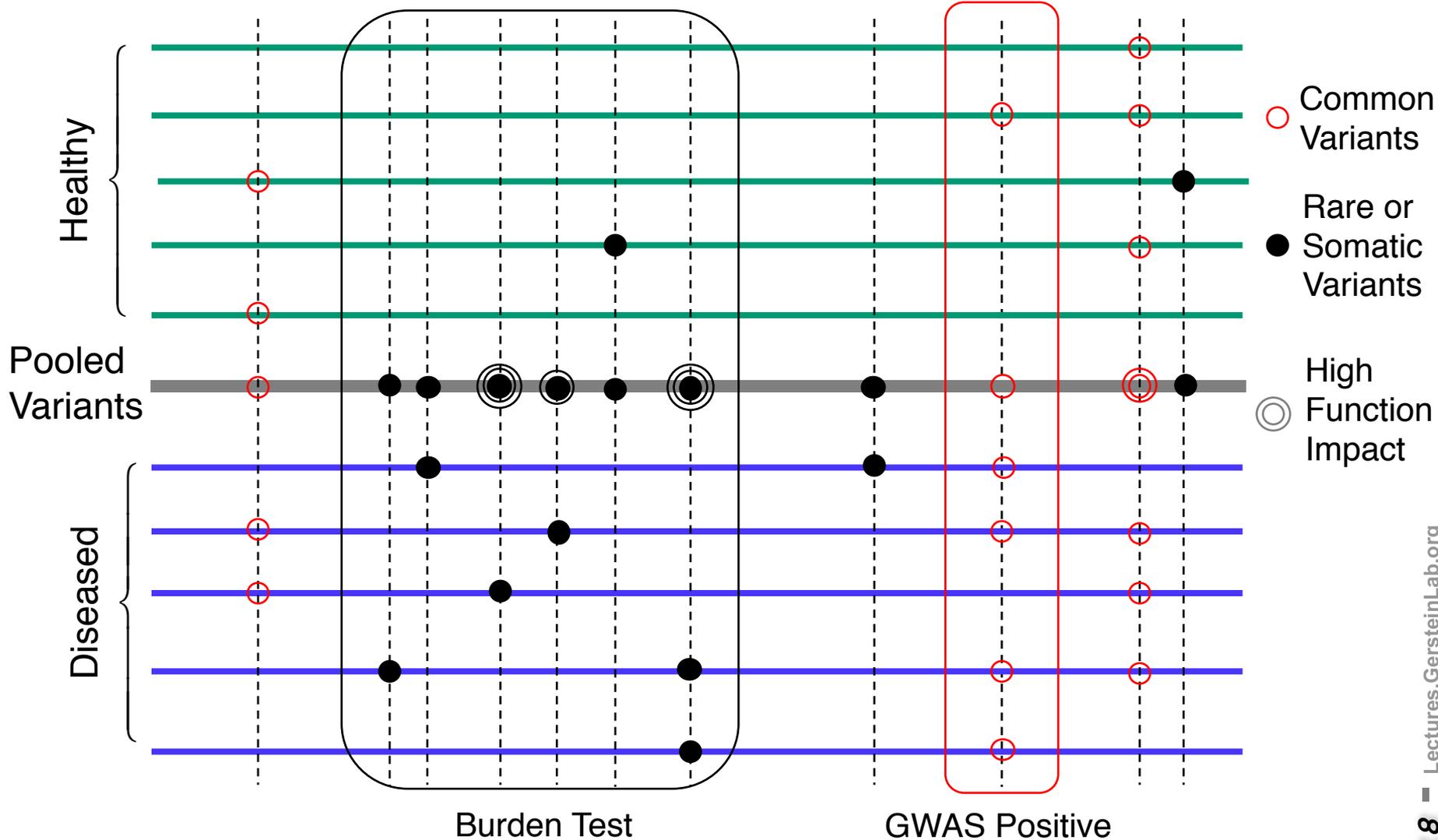  - Driver mutation is a mutation that directly or indirectly confers a selective growth advantage to the cell in which it occurs.
  - A typical tumor contains 2-8 drivers; the remaining mutations are passengers.
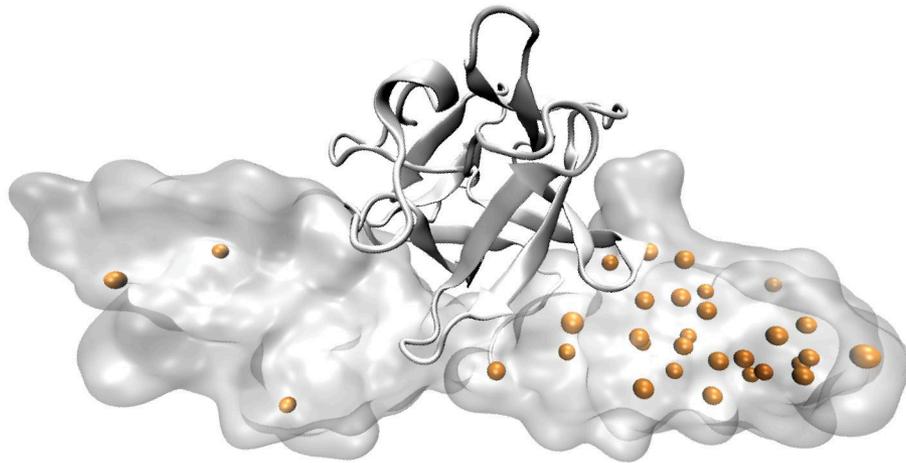- **Passengers**
  - Conceptually, a passenger mutation has no direct or indirect effect on the selective growth advantage of the cell in which it occurred.

# Association of Variants with Diseases

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
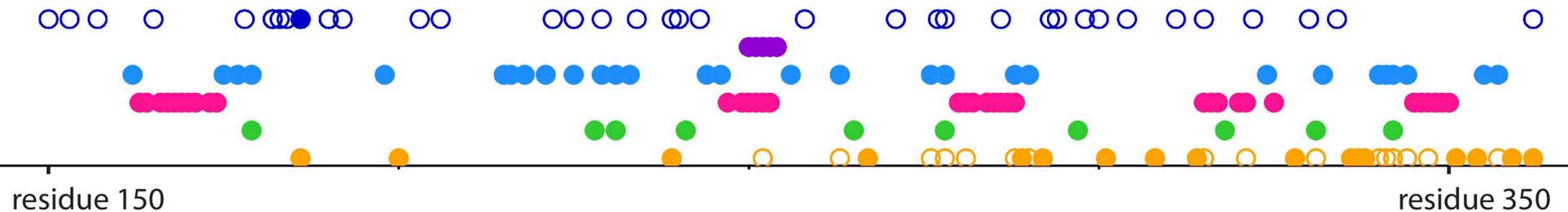    - Prioritzing rare germline variants

# Unlike common SNVs, the statistical power with which we can evaluate rare SNVs in case-control studies is severely limited

Protein structures may provide the needed alternative for evaluating rare SNVs, many of which may be disease-associated
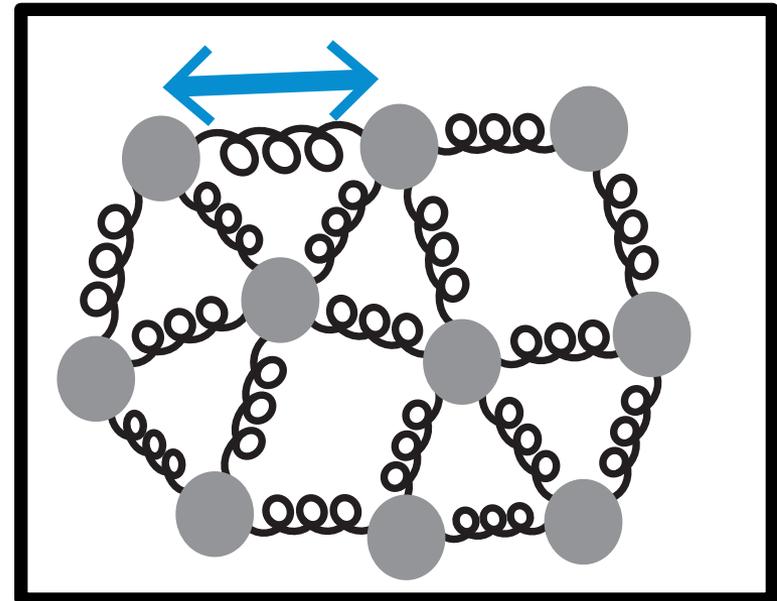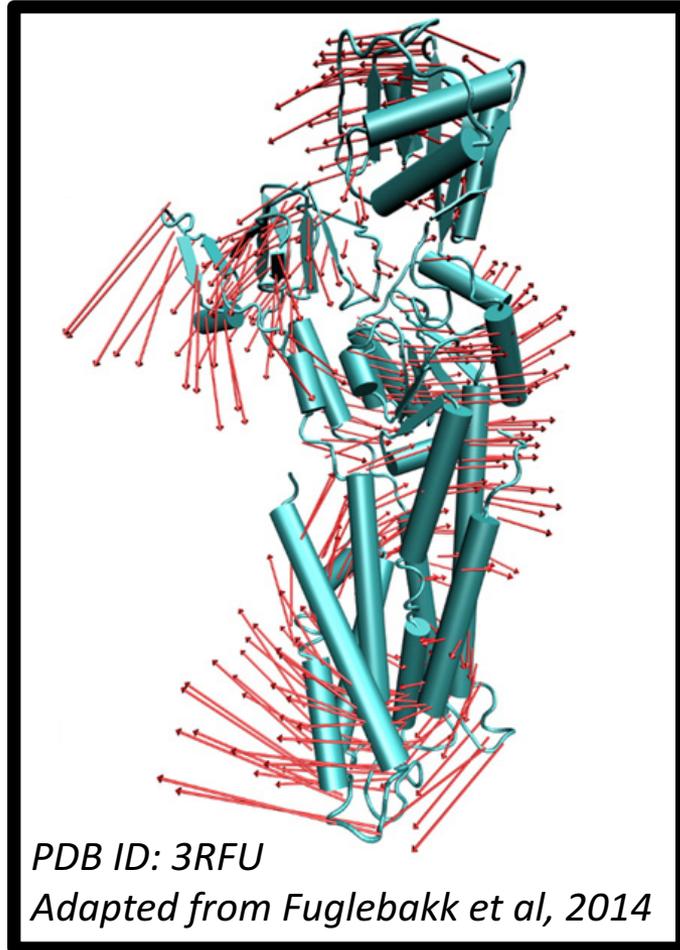


*Fibroblast growth factor receptor 2 (pdb: 1IIL)*

● ○ 1000G & ExAC SNVs (common | rare)
● Hinge residues
● Buried residues
● Protein-protein interaction site
● Post-translational modifications
● HGMD site (w/o annotation overlap)
○ HGMD site (w/annotation overlap)

residue 150

residue 350

[Sethi et al. COSB ('15)]

# Models of Protein Conformational Change

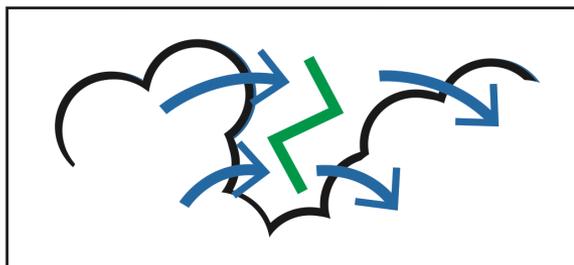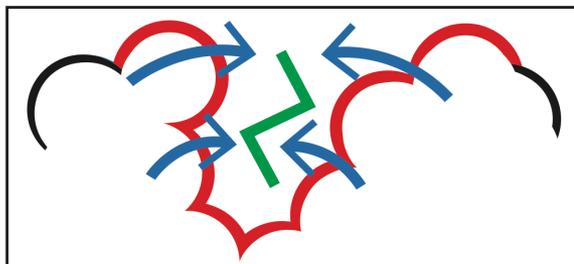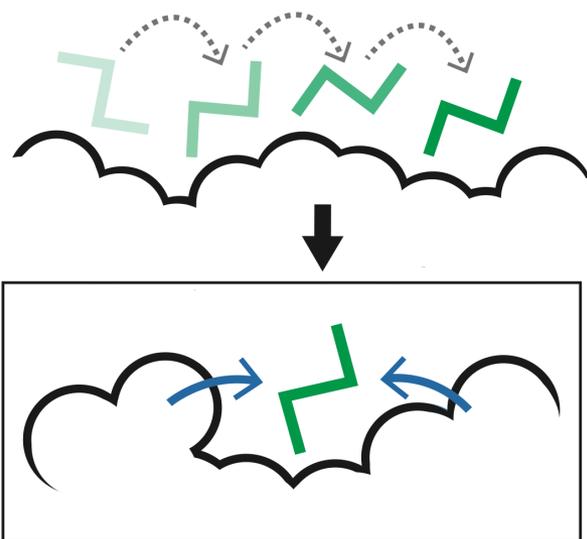## Motion Vectors from Normal Modes (ANMs)



*PDB ID: 3RFU*
*Adapted from Fuglebakk et al, 2014*



Characterizing uncharacterized variants
<= Finding Allosteric sites
<= Modeling motion

# Predicting Allosterically-Important Residues at the Surface

1. MC simulations generate a large number of candidate sites
2. Score each candidate site by the degree to which it perturbs large-scale motions
3. Prioritize & threshold the list to identify the set of high confidence-sites

$$binding\ leverage\ =\ \sum_{m=1}^{10}(\sum_i\sum_j\Delta d_{ij(m)}^2)$$

pdb 1J3H

Surface region with high density of candidate sites

Surface region with low density of candidate sites

Adapted from Clarke*, Sethi*, et al (in press)

# Predicting Allosterically-Important Residues at the Surface



*PDB: 3PFK*

Adapted from Clarke*, Sethi*, et al *(in press)*

# Predicting Allosterically-Important Residues within the Interior



weight edges using motion vectors

identify communities

identify critical residues

Adapted from Clarke*, Sethi*, et al *(in press)*

# Predicting Allosterically-Important Residues within the Interior



weight edges using
motion vectors

① identify
communities

② 

identify
critical residues

③ 

$$Cov_{ij} = \langle \mathbf{r}_i \bullet \mathbf{r}_j \rangle$$

$$C_{ij} = Cov_{ij} / \sqrt{(\langle \mathbf{r}_i^2 \rangle \langle \mathbf{r}_j^2 \rangle)}$$

$$D_{ij} = -\log(|C_{ij}|)$$

Adapted from Clarke*, Sethi*, et al *(in press)*

# Predicting Allosterically-Important Residues within the Interior



weight edges using
motion vectors

① 

identify
communities

② 

identify
critical residues

③ 

PDB: 1XTT

Adapted from Clarke*, Sethi*, et al *(in press)*

# STRESS Server Architecture: Highlights
## stress.molmovdb.org



- A light front-end server handles incoming requests, and powerful back-end servers perform calculations.

- Auto Scaling adjusts the number of back-end servers as needed.

- A typical structure takes ~30 minutes on a E5-2660 v3 (2.60GHz) core.

- Input & output (i.e., predicted allosteric residues) are stored in S3 buckets.

Adapted from Clarke*, Sethi*, et al *(in press)*

# Intra-species conservation of predicted allosteric residues
## *1000 Genomes*



### *Surface*

### *Interior*

**critical**
**non-critical**

p=0.309

**p=1.80e-05**

Adapted from Clarke*, Sethi*, et al *(in press)*

# Intra-species conservation of predicted allosteric residues
## ExAC

**Surface**

**Interior**

Minor Allele Freq.

5e-3
5e-4
5e-5

■ critical
■ non-critical

Minor Allele Freq.

5e-3
5e-4
5e-05

p=1.49e-3

p=7.98e-09

Adapted from Clarke*, Sethi*, et al (in press)

# Unlike common SNVs, the statistical power with which we can evaluate rare SNVs in case-control studies is severely limited

Protein structures may provide the needed alternative for evaluating rare SNVs, many of which may be disease-associated



*Fibroblast growth factor receptor 2 (pdb: 1IIL)*

● ○ 1000G & ExAC SNVs (common | rare)
● Hinge residues
● Buried residues
● Protein-protein interaction site
● Post-translational modifications
● HGMD site (w/o annotation overlap)
○ HGMD site (w/annotation overlap)

residue 150

residue 350

[Sethi et al. COSB ('15)]

# Protein structures may provide the needed alternative for evaluating rare SNVs, many of which may be disease-associated

## Rationalizing disease variants in the context of allosteric behavior with allostery as an added annotation



*Fibroblast growth factor receptor 2 (pdb: 1IIL)*

- ●● Predicted allosteric (surface | interior)
- ●○ 1000G & ExAC SNVs (common | rare)
- ● Hinge residues
- ● Buried residues
- ● Protein-protein interaction site
- ● Post-translational modifications
- ● HGMD site (w/o annotation overlap)
- ○ HGMD site (w/annotation overlap)

residue 150                                                    residue 350

[Sethi et al. COSB ('15)]

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

# Schematic illustration of localized frustration



ASN ⟷ ASP

more negative — more positive

favorable interaction — unfavorable interaction

[Ferreiro et al., *PNAS* ('07)]

# Workflow for evaluating localized frustration changes (ΔF)



$$F_{mut} - F_{nat} = \Delta F < 0$$

Native Structure

Mutated Structure

Model of mutated structure

$$\frac{\langle E \rangle' - E_{mut}}{\sigma'_E} = F_{mut} < 0$$

$E_{mut}$

$\langle E \rangle'$

$\langle E \rangle$

$E_{nat}$

Model of WT structure

$$\frac{\langle E \rangle - E_{nat}}{\sigma_E} = F_{nat} > 0$$

# Comparing Frustration (ΔF values) across different SNV categories



A — 1000 Genomes (Core, Surface)
B — ExAC (Core, Surface)
C — HGMD (Core, Surface)

ΔF

# ΔF distributions among rare and common SNVs

A    1000 Genomes

B    ExAC

ΔF

Core    Surf. | Core    Surf.
*common* | *rare*

# Comparison between ΔF distributions: TSGs vs. oncogenes

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

# Non-coding Annotations: Overview

Sequence features, incl. **Conservation**

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



Large-scale sequence similarity comparison

Identify large blocks of repeated and deleted sequence:

• Within the human reference genome

• Within the human population

• Between closely related mammalian genomes

Identify smaller-scale repeated blocks using statistical models

Signal processing of raw experimental data:

• Removing artefacts
• Normalization
• Window smoothing

Segmentation of processed data into active regions:

• Binding sites
• Transcriptionally active regions

Group active regions into larger annotation blocks

[Alexander et al., *Nat. Rev. Genet.* ('10)]

# 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

# Finding "Conserved" Sites in the Human Population:

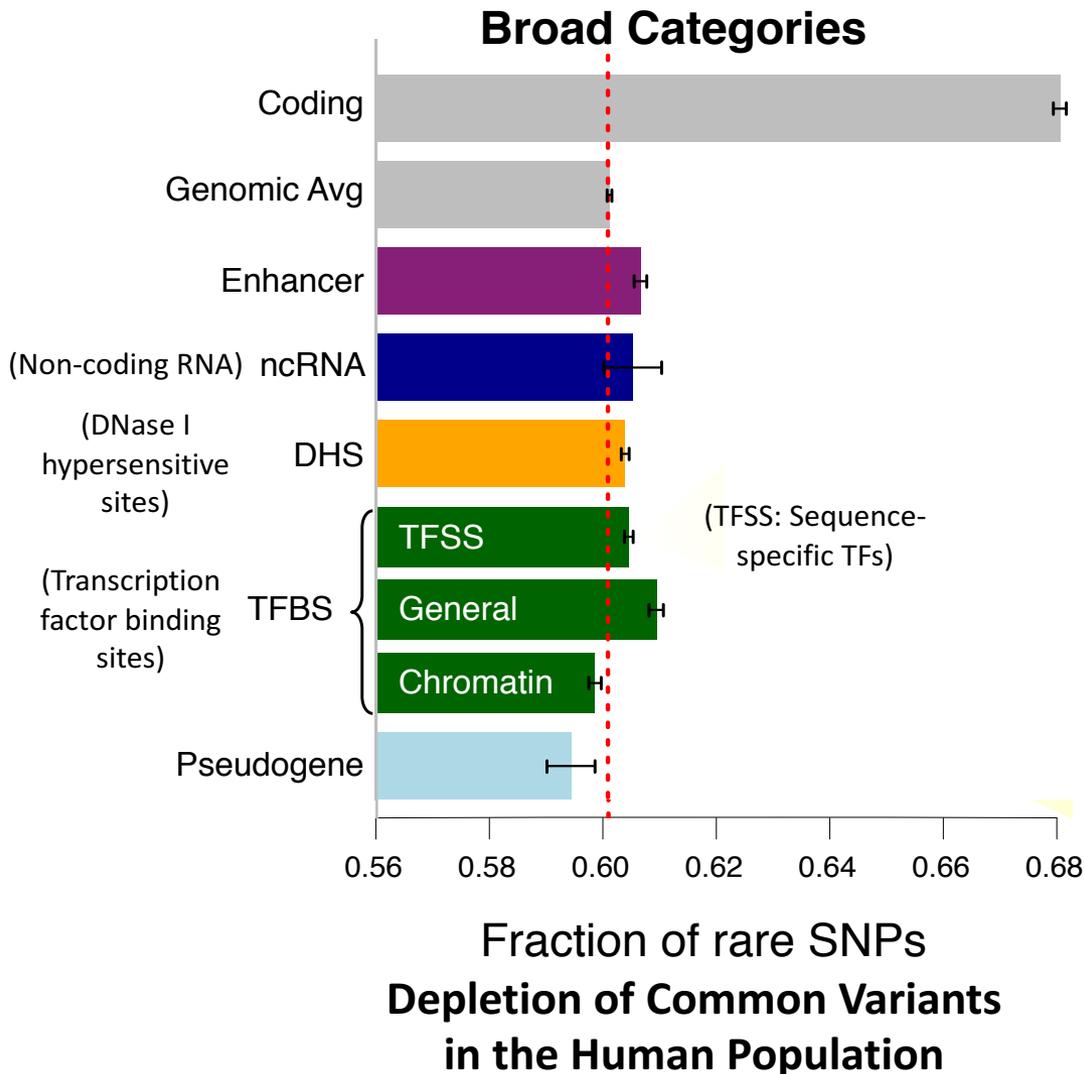## Negative selection in non-coding elements based on Production ENCODE & 1000G Phase 1



**Broad Categories**

Coding

Genomic Avg

Enhancer

(Non-coding RNA) ncRNA

(DNase I hypersensitive sites) DHS

(Transcription factor binding sites) TFBS
- TFSS
- General
- Chromatin

(TFSS: Sequence-specific TFs)

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

| | |
|---|---|
| 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

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

~0.4% genomic coverage (~ top 25)

~0.02% genomic coverage (top 5)

**A** Broad Categories

**B** Specific Categories

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

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

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
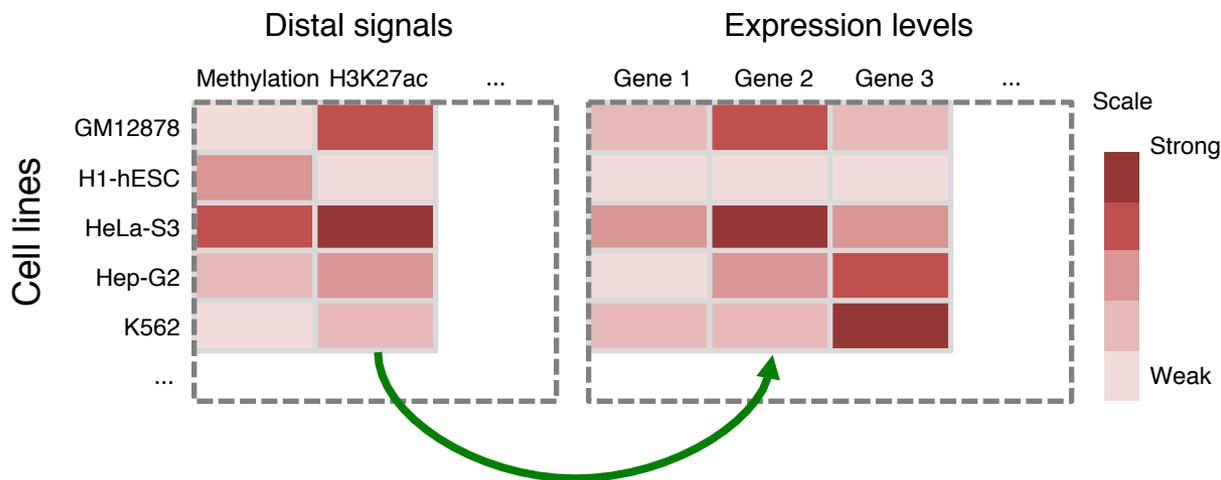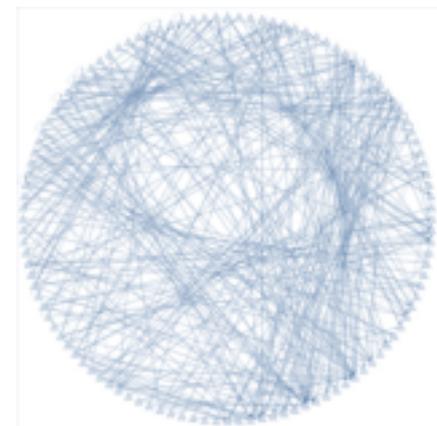    - Prioritzing rare germline variants

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

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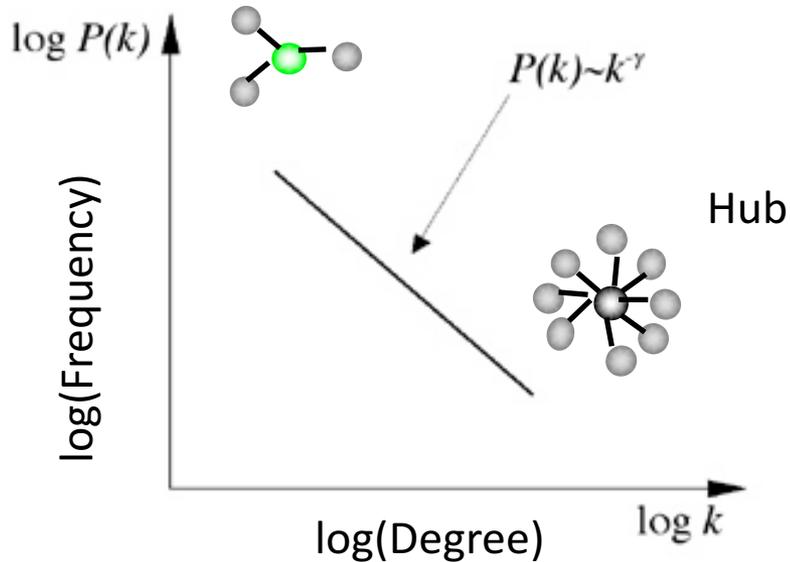Yip et al., *GenomeBiology* ('12),
Fu et al., *GenomeBiology*('14) ]

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

Expression levels

Methylation H3K27ac ...

Gene 1 Gene 2 Gene 3 ...

Scale

Cell lines

GM12878

H1-hESC

HeLa-S3

Hep-G2

K562

...

Strong

Weak

Connecting Distal Elements via **Activity Correlations**.

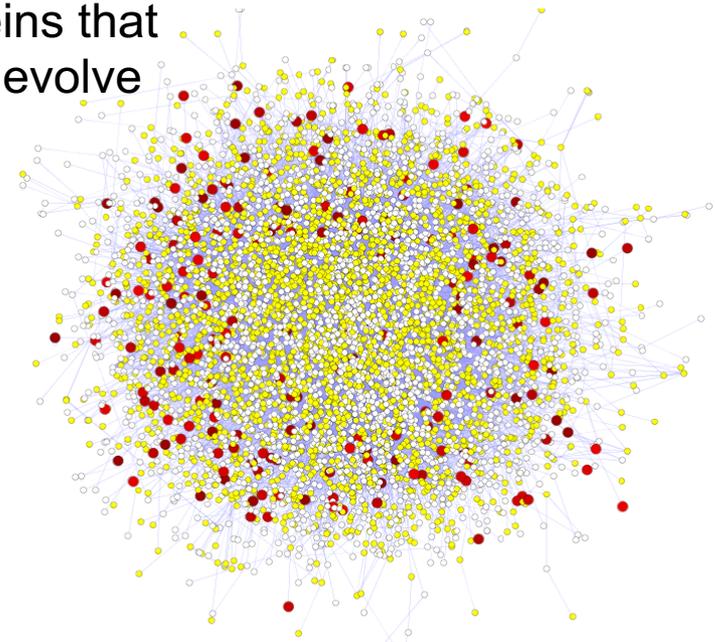Other strategies to create linkage incl. eQTL and Hi-C. Much in recent Epigenomics Roadmap.

Power-law distribution



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

log P(k)

log(Frequency)

log(Degree)

log k

Hub

High likelihood of positive selection

Not under positive selection

Lower likelihood of positive selection

No data about positive selection

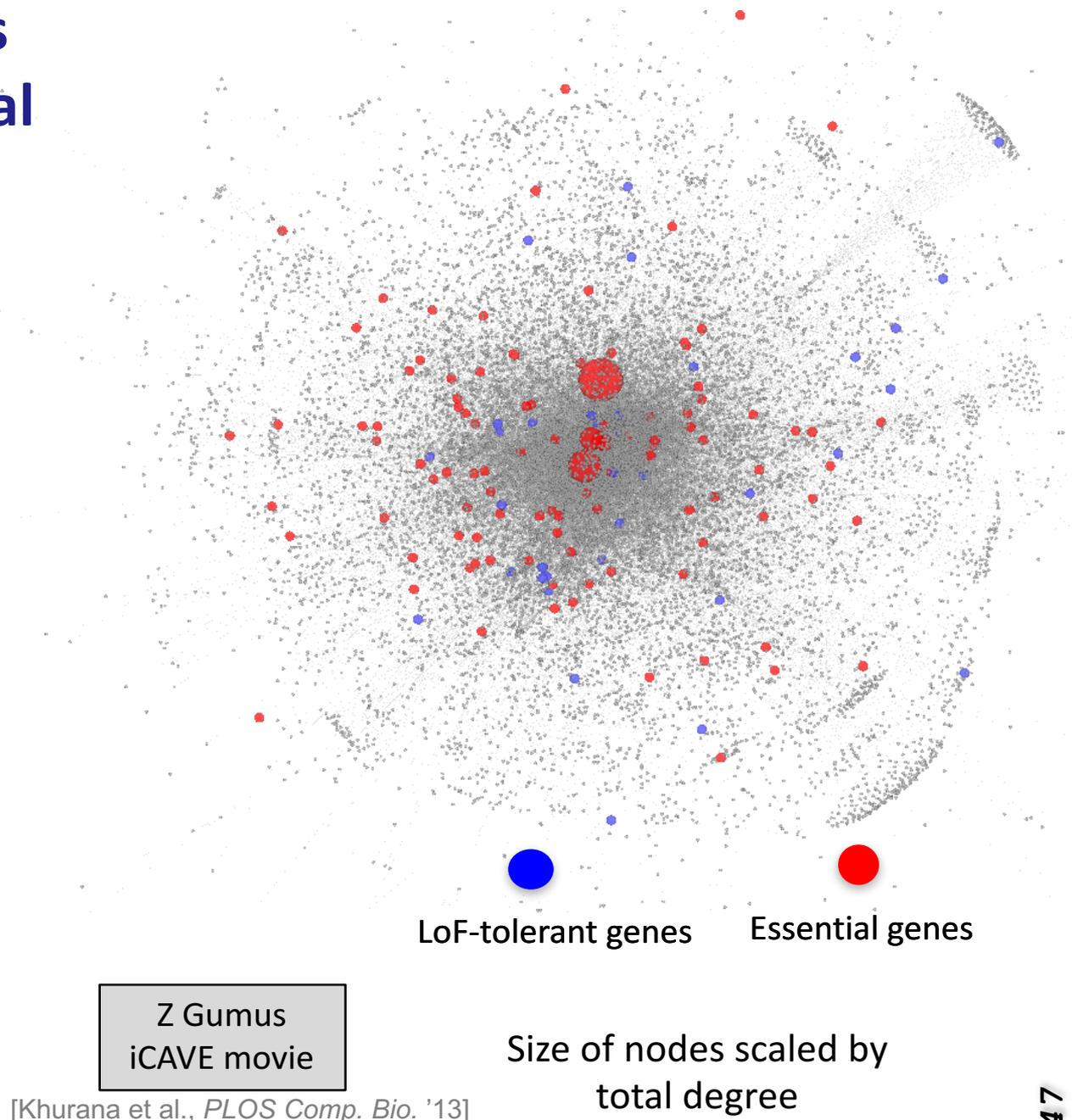[Nielsen et al. *PLoS Biol.* (2005), HPRD, Kim et al. PNAS (2007)]
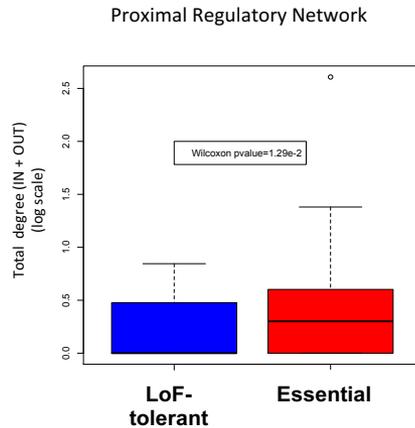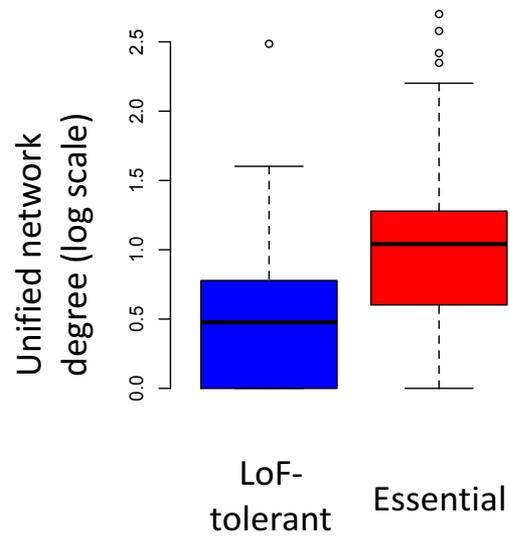
- <u>More Connectivity, More Constraint:</u> Genes & proteins that have a more central position in the network tend to evolve more slowly and are more likely to be essential.
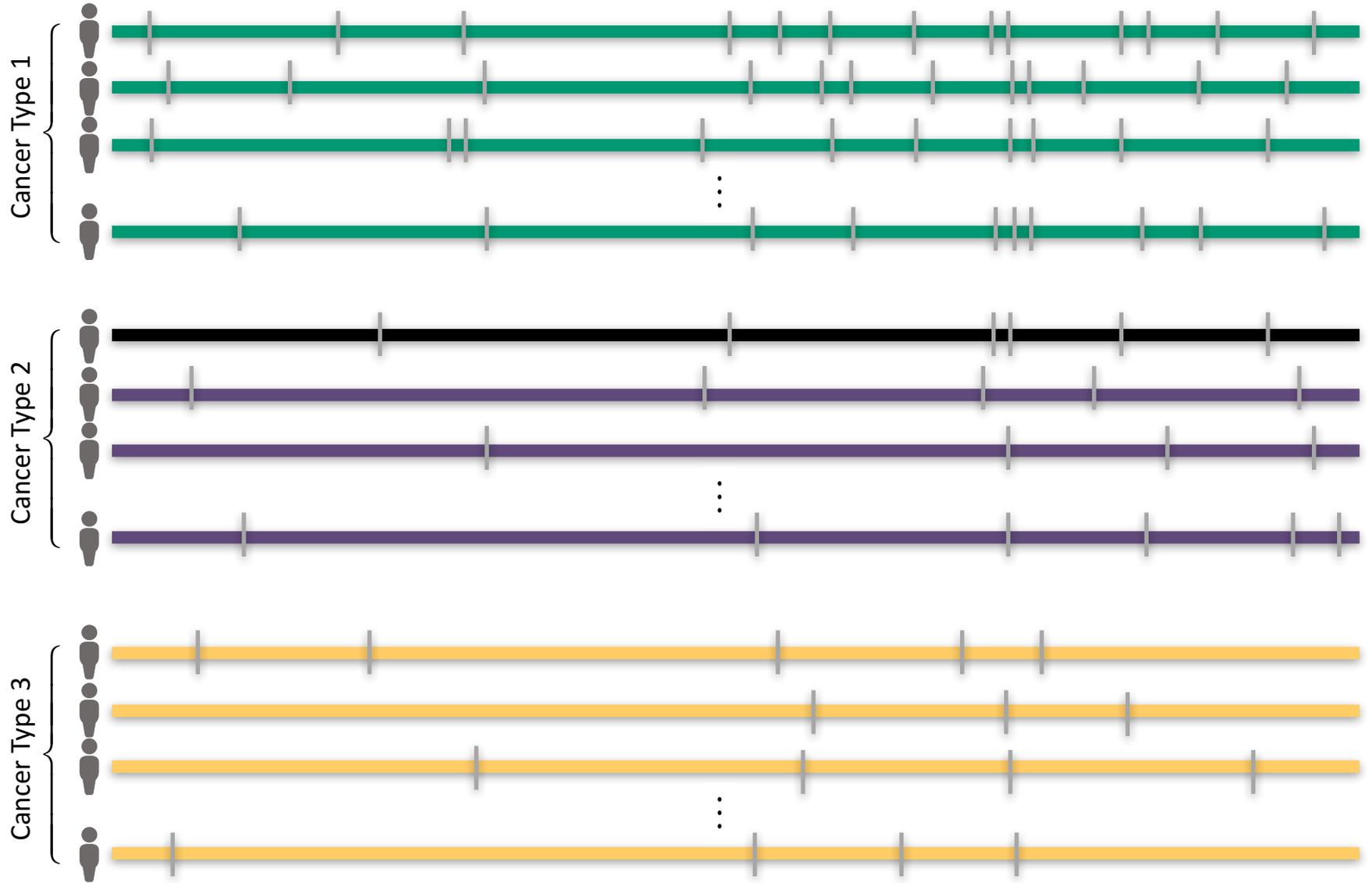
- This phenomenon is observed in **many organisms & different kinds of networks**

  - **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
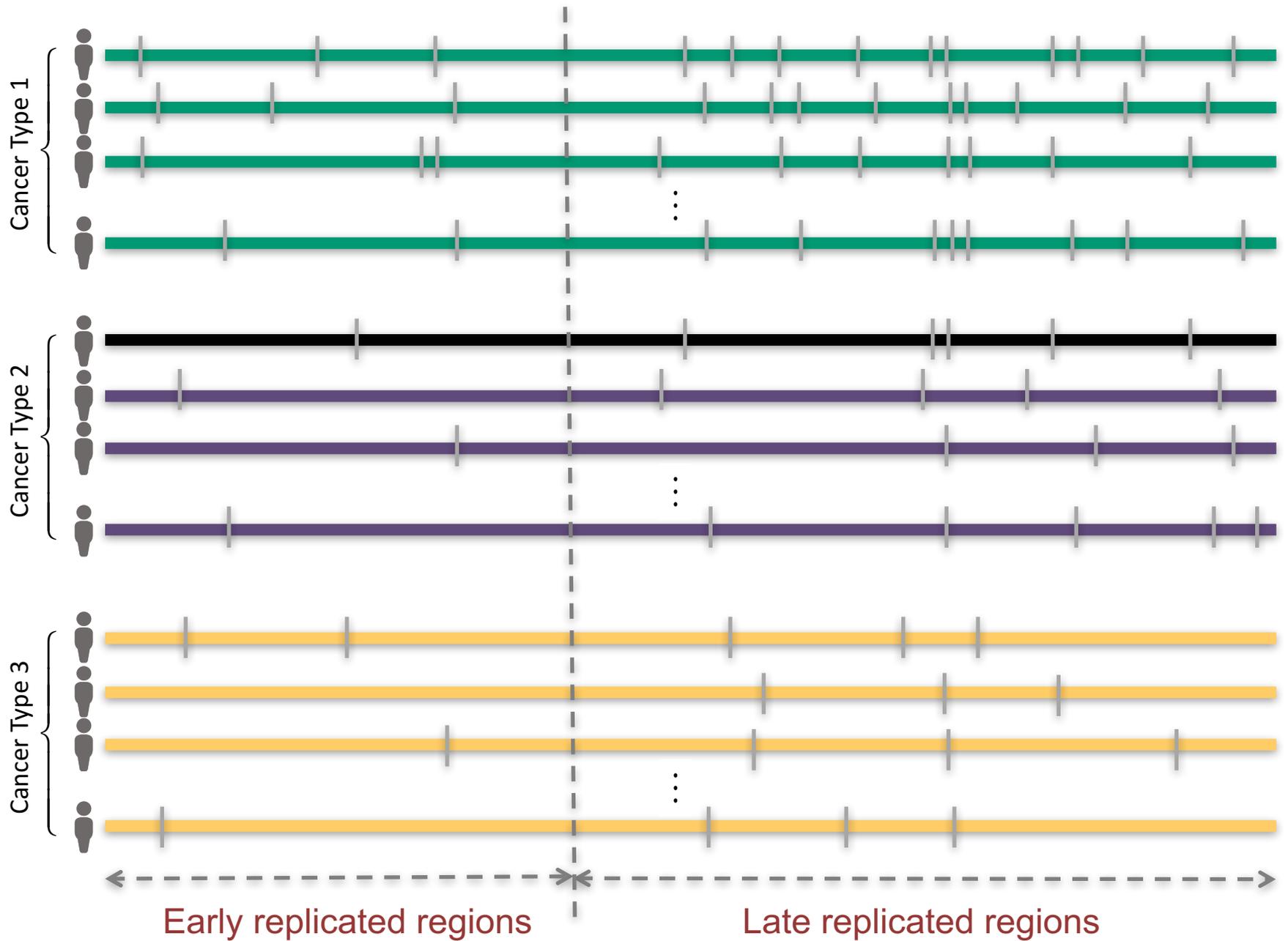
# Regulatory Hubs
# are more Essential



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

Z Gumus
iCAVE movie

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

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

# Mutation recurrence



Cancer Type 1

Cancer Type 2

Cancer Type 3

# Mutation recurrence



Cancer Type 1

Cancer Type 2

Cancer Type 3

Early replicated regions

Late replicated regions

**Noncoding annotations**

Cancer Type 1

Cancer Type 2

Cancer Type 3

Early replicated regions

Late replicated regions

**Noncoding annotations**

Cancer Type 1

Cancer Type 2

Cancer Type 3

Early replicated regions

Late replicated regions

**Cancer Somatic Mutational Heterogeneity, across cancer types, samples & regions**



(A)

y-axis: log(mutations/1Mb)

Categories: Medulloblastoma, CLL, Liver, Prostate, Lymphoma, Breast, Glial, Pancreas, Kidney, Stomach, Lung

Lower panel:
- Lung_mean
- Lung_SD

y-axis: mutations/Mbp

x-axis: 1 Mbp genome regions (locations chosen at random)

[Lochovsky et al. *NAR* ('15)]

# Cancer Somatic Mutation Modeling

- 3 models to evaluate the significance of mutation burden

- Suppose there are *k* genome elements. For element *i*, define:
  - *$n_i$*: total number of nucleotides
  - *$x_i$*: the number of mutations within the element
  - *p*: the mutation rate
  - *R*: the replication timing bin of the element

**Model 1: Constant Background Mutation Rate (Model from Previous Work)**

$$\mathbf{x_i} : Binomial(\mathbf{n_i}, \mathbf{p})$$

**Model 2: Varying Mutation Rate**

$$\mathbf{x_i}|\mathbf{p_i} : Binomial(\mathbf{n_i}, \mathbf{p_i})$$

$$\mathbf{p_i} : Beta(\mu, \sigma)$$

**Model 3: Varying Mutation Rate with Replication Timing Correction**

$$\mathbf{x_i}|\mathbf{p_i} : Binomial(\mathbf{n_i}, \mathbf{p_i})$$

$$\mathbf{p_i} : Beta(\mu|\mathbf{R}, \sigma|\mathbf{R})$$

$$\mu|\mathbf{R}, \sigma|\mathbf{R} : \text{constant within the same } \mathbf{R} \text{ bin}$$
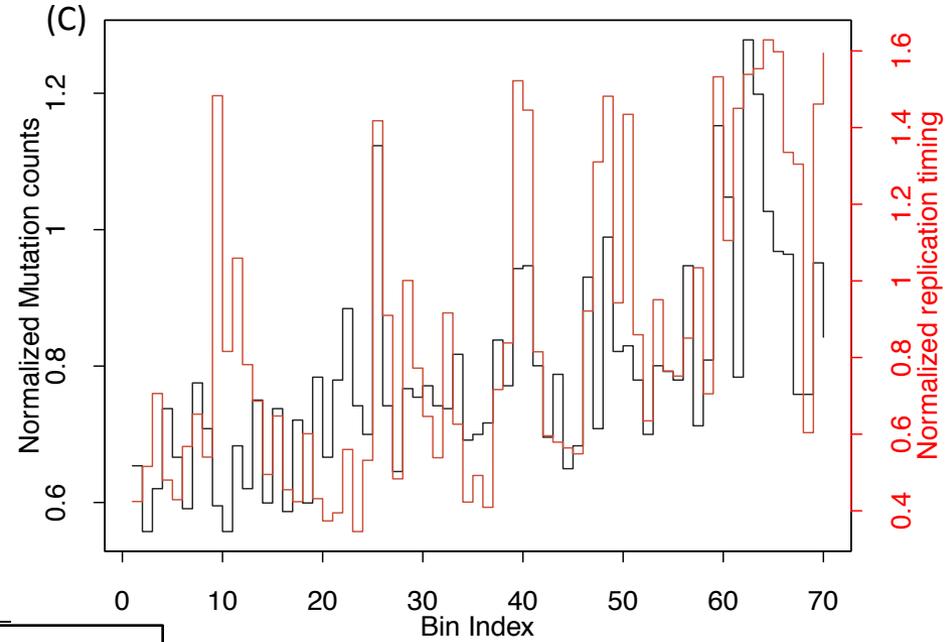
[Lochovsky et al. *NAR* ('15)]

# 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
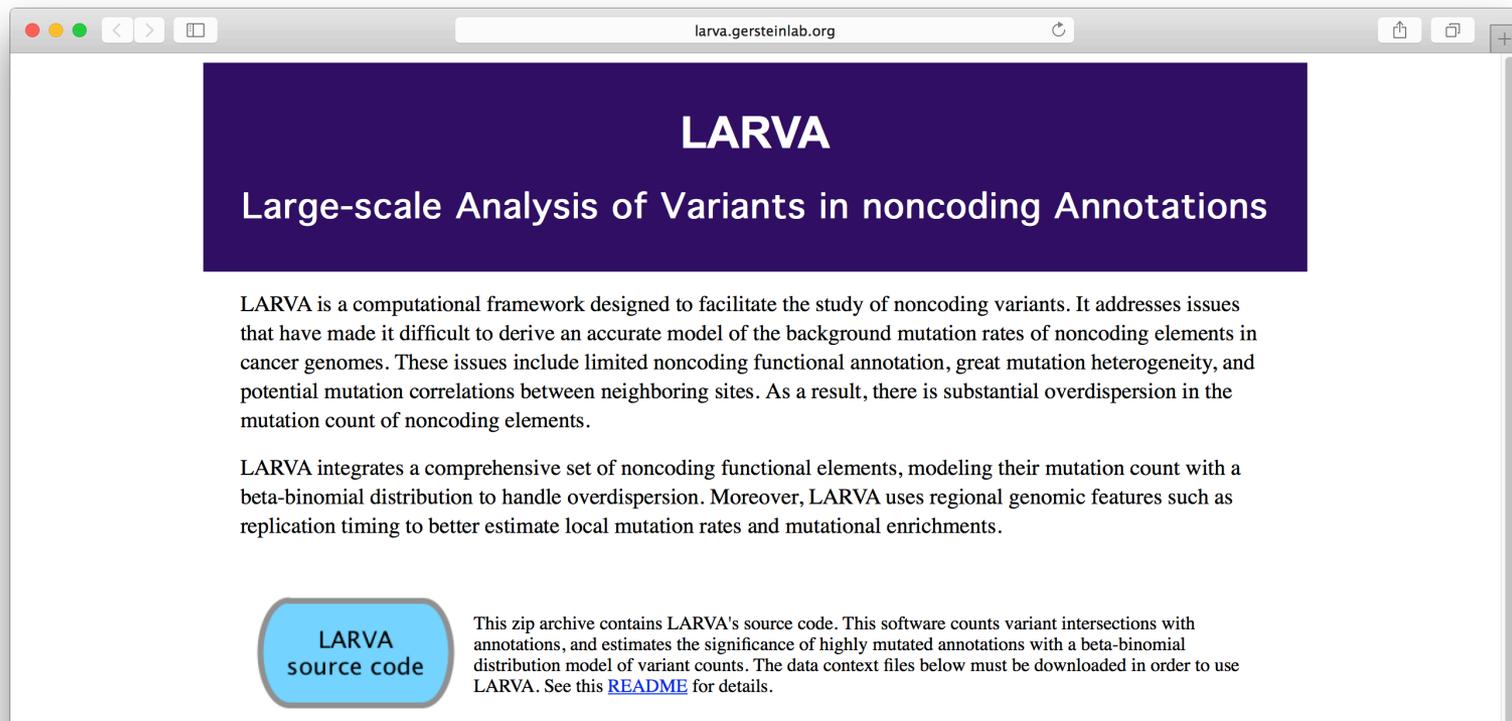
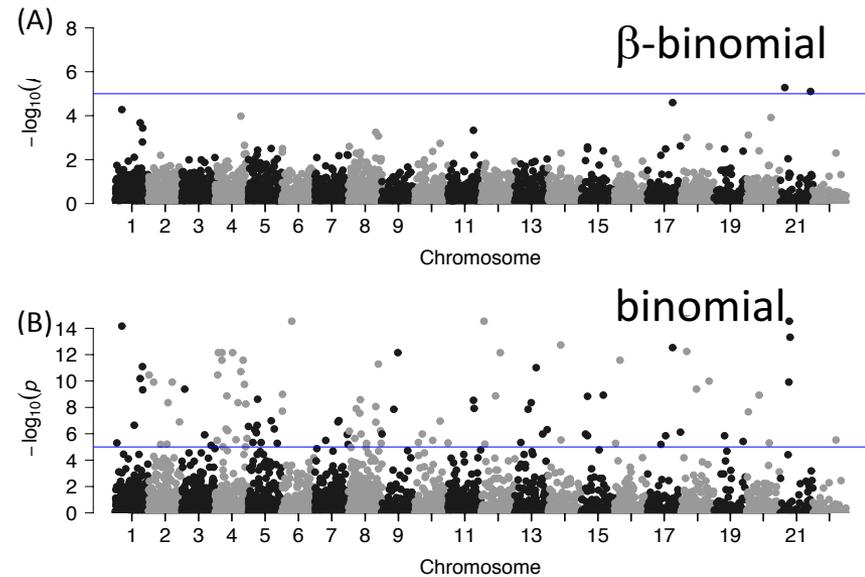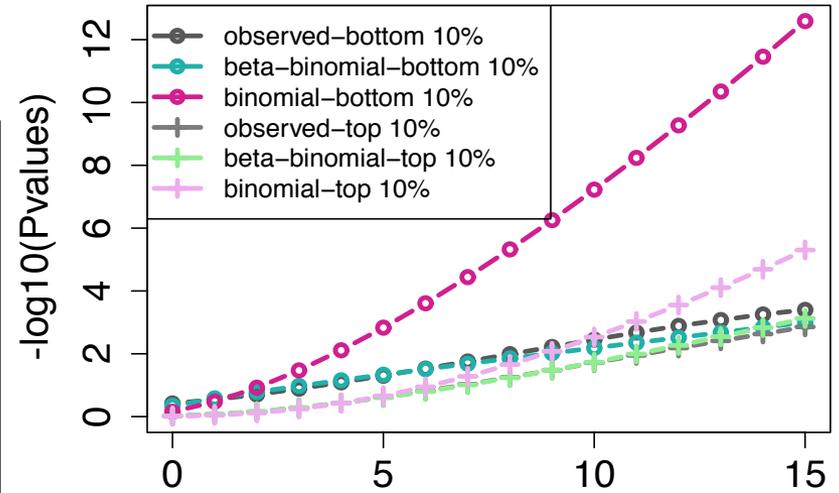# Adding DNA replication timing correction further improves the beta-binomial model



(C)

[Lochovsky et al. *NAR* ('15)]

(A)



probably

**Bottom 10% of rep. timing bins requires large correction**

**Top 10% of rep. timing bins requires little correction**

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%

somatic mutation count

# LARVA Implementation

- http://larva.gersteinlab.org/
- Freely downloadable C++ program
  - Verified compilation and correct execution on Linux
- A Docker image is also available to download
  - Runs on any operating system supported by Docker
- Running time on transcription factor binding sites (a worst case input size) is ~80 min
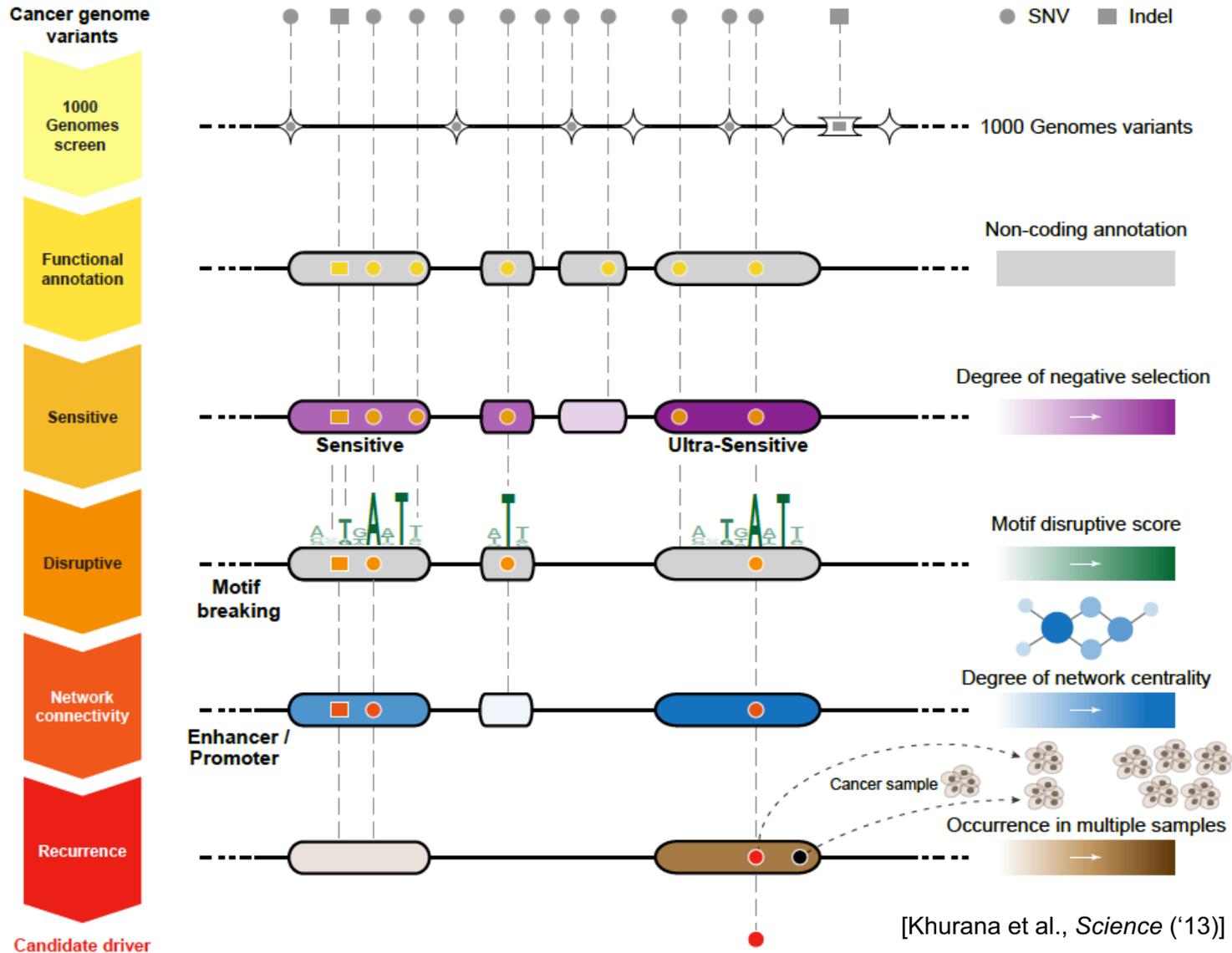  - Running time scales linearly with the number of annotations in the input



**LARVA**

**Large-scale Analysis of Variants in noncoding Annotations**

LARVA is a computational framework designed to facilitate the study of noncoding variants. It addresses issues that have made it difficult to derive an accurate model of the background mutation rates of noncoding elements in cancer genomes. These issues include limited noncoding functional annotation, great mutation heterogeneity, and potential mutation correlations between neighboring sites. As a result, there is substantial overdispersion in the mutation count of noncoding elements.

LARVA integrates a comprehensive set of noncoding functional elements, modeling their mutation count with a beta-binomial distribution to handle overdispersion. Moreover, LARVA uses regional genomic features such as replication timing to better estimate local mutation rates and mutational enrichments.

**LARVA source code**

This zip archive contains LARVA's source code. This software counts variant intersections with annotations, and estimates the significance of highly mutated annotations with a beta-binomial distribution model of variant counts. The data context files below must be downloaded in order to use LARVA. See this README for details.
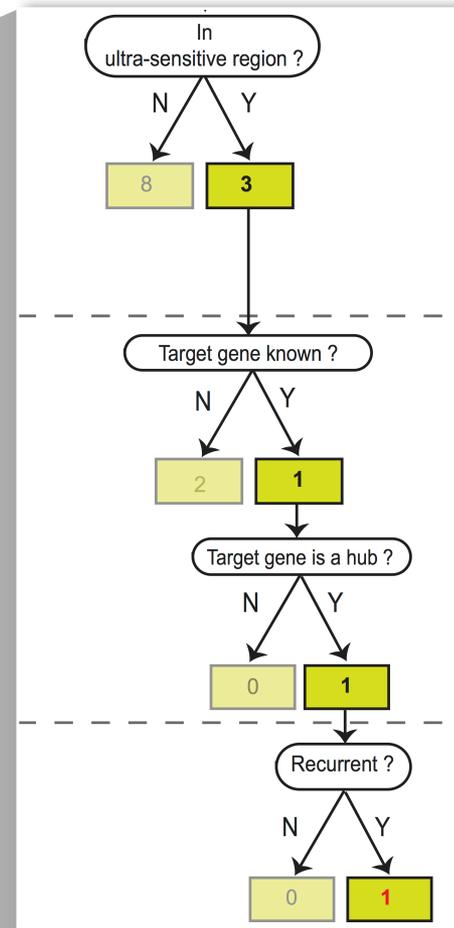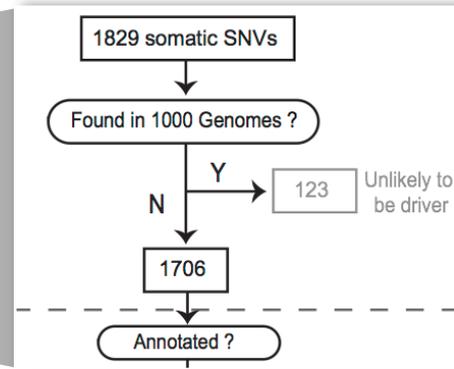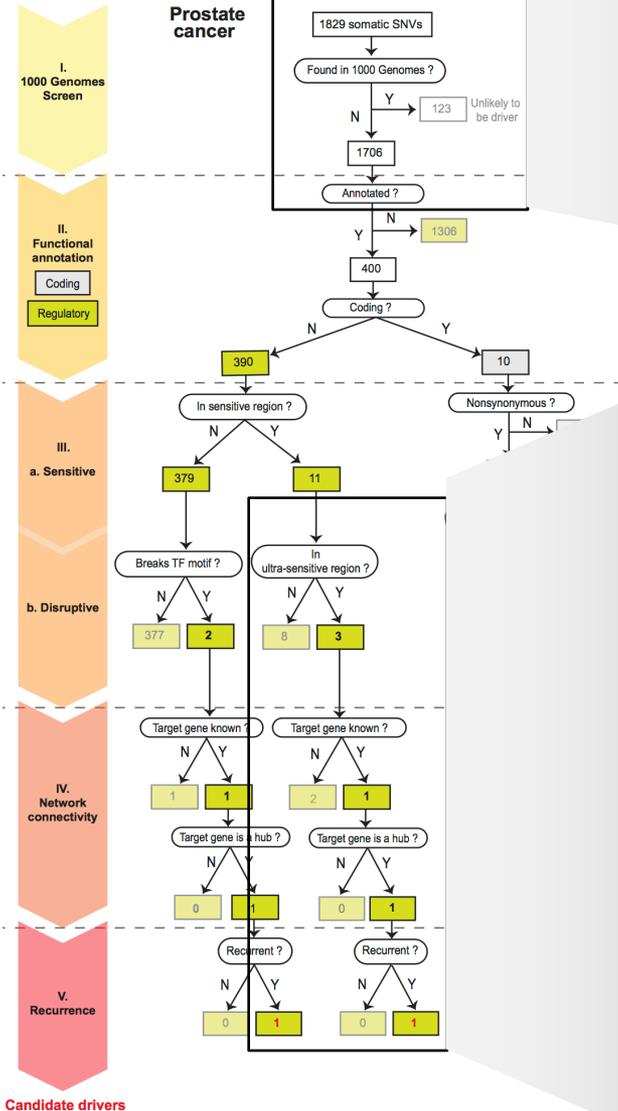
# LARVA Results



TSS LARVA results

adjusted P w/. correction

PRRC2B

TP53

LMO3

These have
literature-verified
cancer associations

AGAP5,PROZ

TERT

noncoding annotation
p-values in sorted order

-log10(Pvalues)

- observed–bottom 10%
- beta–binomial–bottom 10%
- binomial–bottom 10%
- observed–top 10%
- beta–binomial–top 10%
- binomial–top 10%

(A)    β-binomial

Chromosome

(B)    binomial

Chromosome

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones
- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)
- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

# 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)

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

Site integrates
user variants
with large-scale
context

FunSeq.gersteinlab.org

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

- Feature weight

  - Weighted with mutation patterns in natural polymorphisms

    (features frequently observed weight less)

  - entropy based method



Legend:
- HOT region
- Sensitive region
- Polymorphisms

Genome

[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

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

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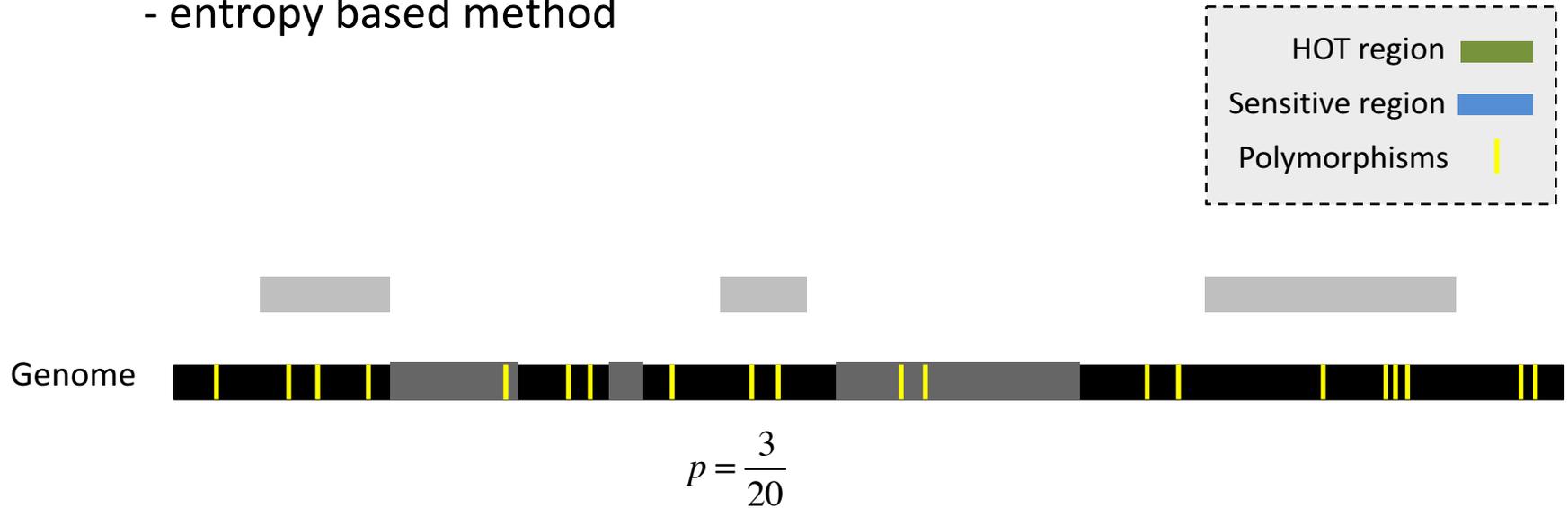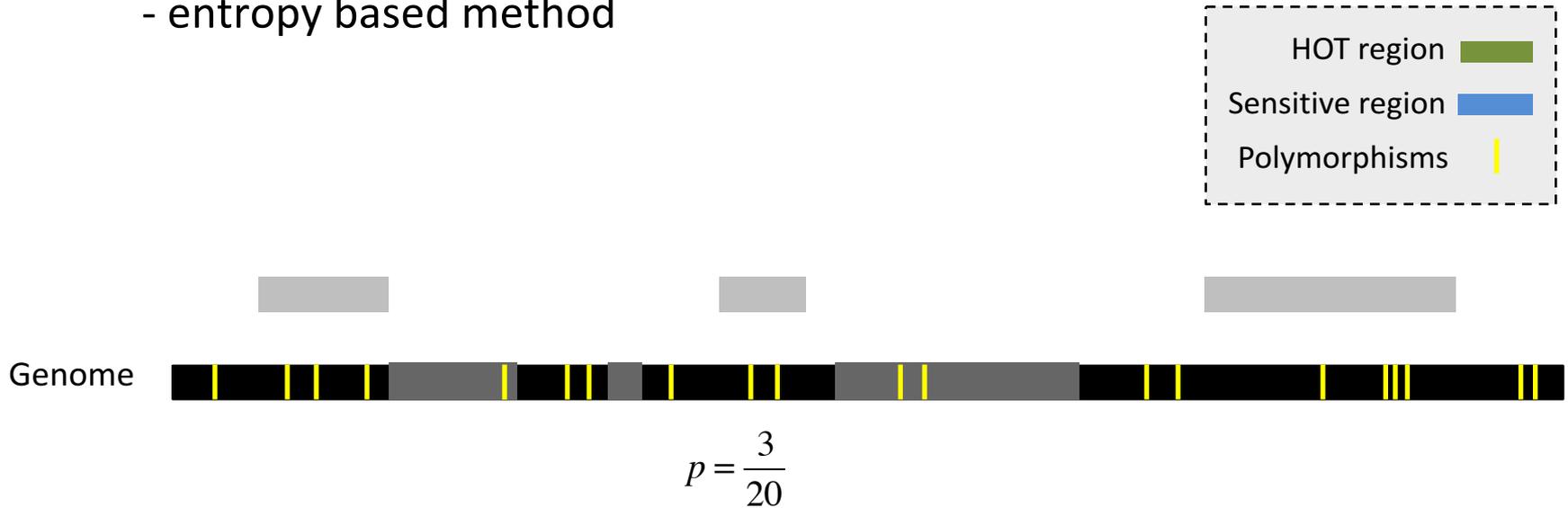
- Feature weight

    - 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: $Score = \sum w_d \quad of\ observed\ features$

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

# Germline pathogenic variants show higher core scores than controls



3 controls with natural polymorphisms (allele frequency >= 1% )

1.  Matched region:  1kb around HGMD variants

2. Matched TSS:  matched for distance to TSS

3. Unmatched: randomly selected

*Ritchie et al., Nature Methods, 2014*

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

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones

- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)

- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

**Personal Genomics:**
**Managing Rapid Data Scaling through Prioritizing High-impact Variants**

- Introduction
  - The exponential scaling of data generation & processing
  - The landscape of variants in personal genomes suggests finding a few key ones

- Characterizing Rare Variants in Coding Regions
  - Identifying with STRESS cryptic allosteric sites
    - On surface & in interior bottlenecks
  - Using changes in localized Frustration to find sites sensitive to mutations
    - Difference betw. TSGs & oncogenes

- Evaluating the Impact of Non-coding Variants with Annotation
  - Annotating non-coding regions
  - Prioritizing rare variants with "sensitive sites" (human-conserved)
  - Prioritizing in terms of network connectivity (eg hubs)

- Putting it together in Workflows
  - Using LARVA to do burden testing on non-coding annotation
    - Need to correct for over-dispersion mutation counts
    - Parameterized according to replication timing
  - Using FunSeq to integrate evidence on variants
    - Systematically weighting all the features
    - suggesting non-coding drivers
    - Prioritzing rare germline variants

## Acknowledgments

# Extra

# Info about content in this slide pack

- General PERMISSIONS

  - This Presentation is copyright Mark Gerstein, Yale University, 2016.

  - Please read permissions statement at 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.

- 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