# Genomics & Data Science:
# Approaches to identifying key variants through functional impact & recurrence



Mark Gerstein
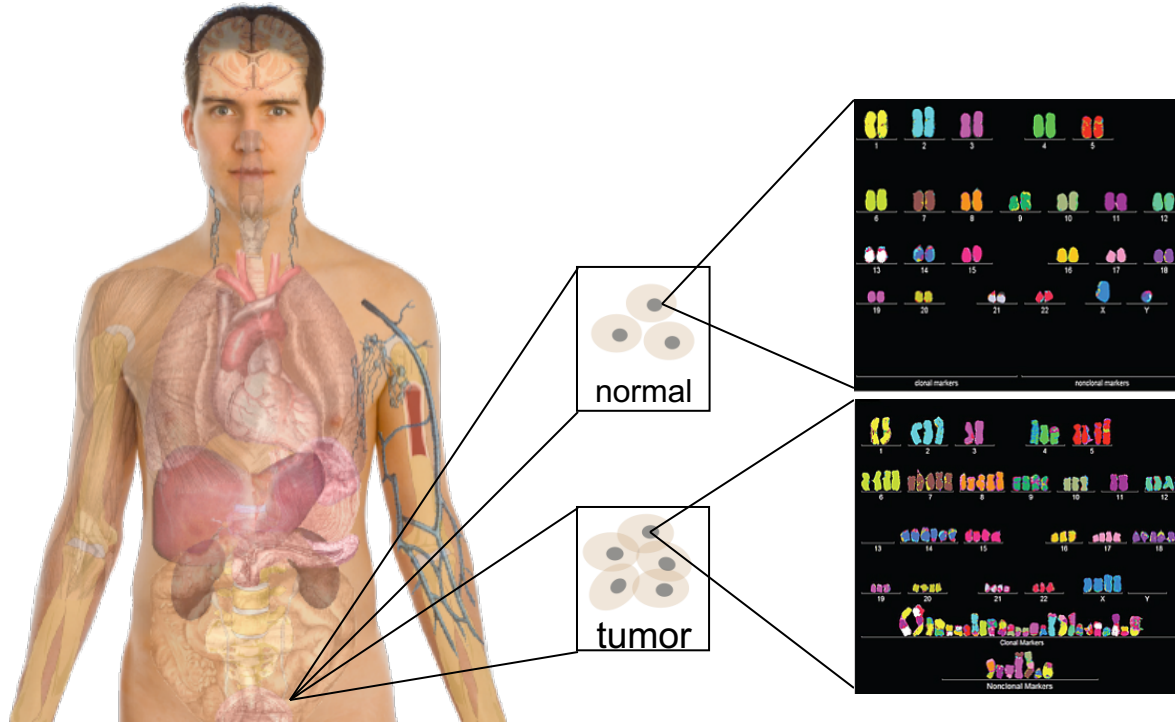Yale

Slides freely
downloadable from
Lectures.GersteinLab.org
& "tweetable"
(via @MarkGerstein).

No Conflicts for this Talk
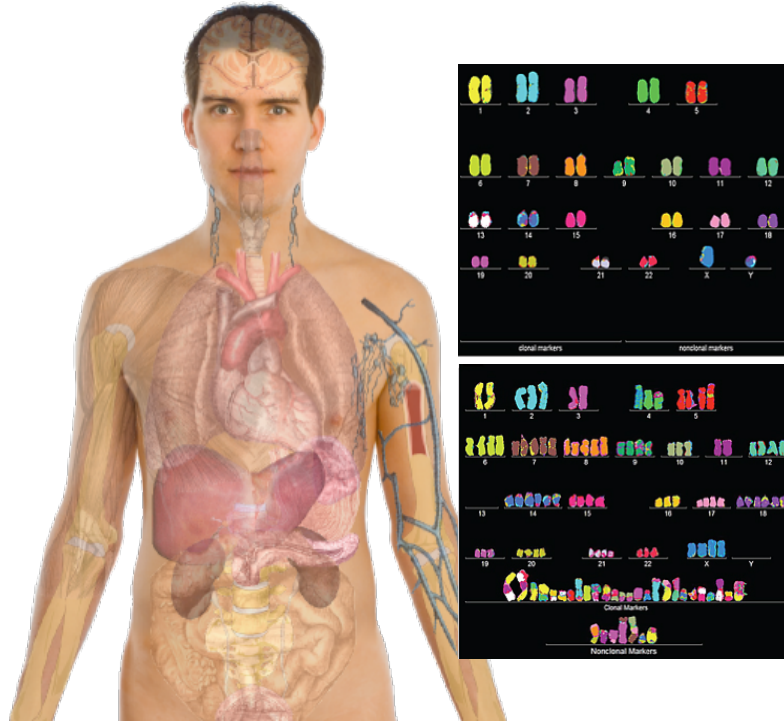
See last slide for more info.

# Personal Genomics as a Gateway into Biology

Personal genomes will soon 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.



normal

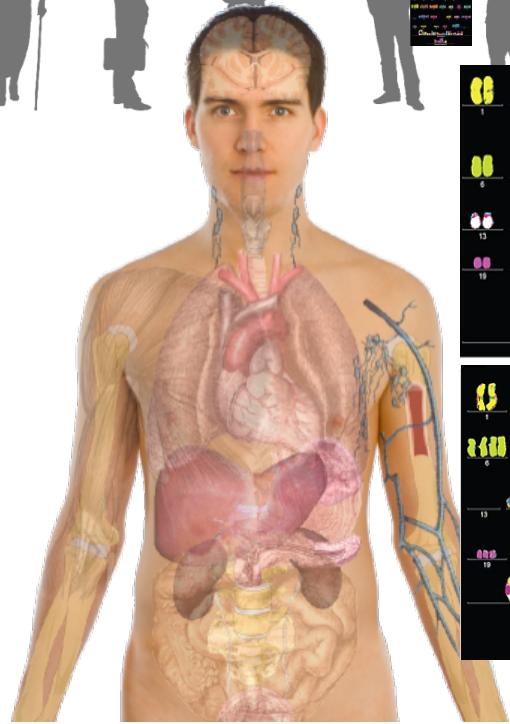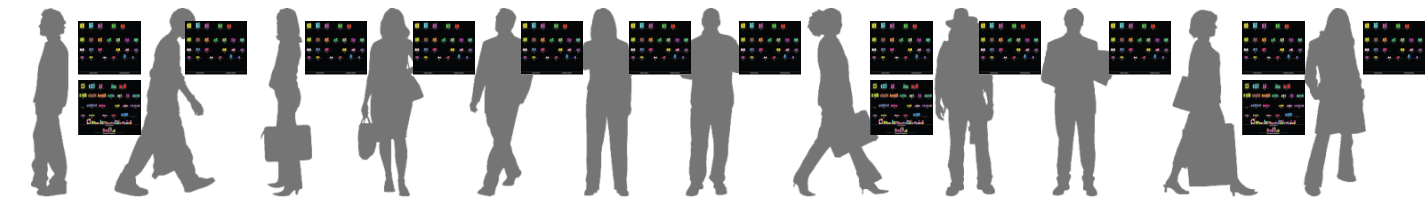tumor

# Personal Genomics
# as a Gateway into Biology

Personal genomes will soon 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.

Keys to genome interpretation
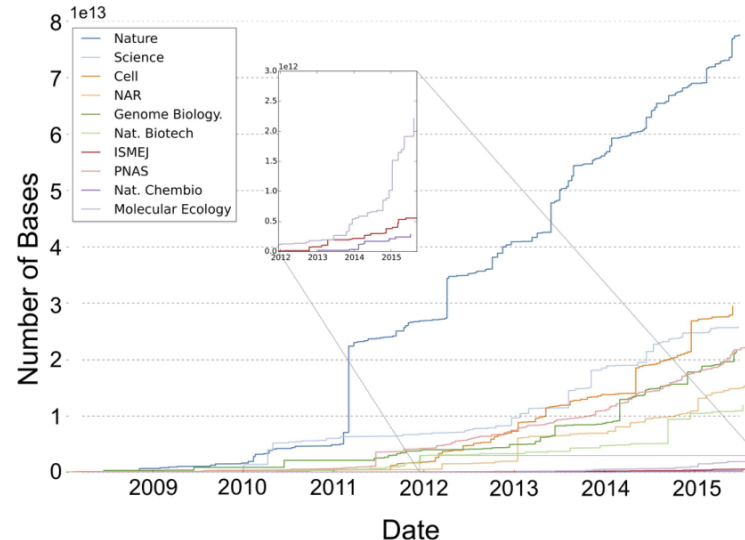
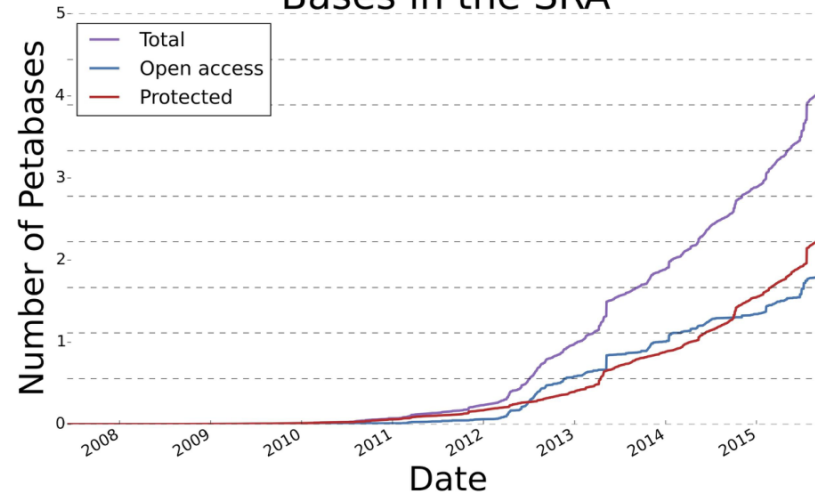Relating individuals' variants to **DBs**

**Scaling** DBs to the **population**

Identifying **key variants** - separating into rare, recurrent, common, &c

## DB Growth: explosion of data scale & a diversity of uses



Bases in the SRA

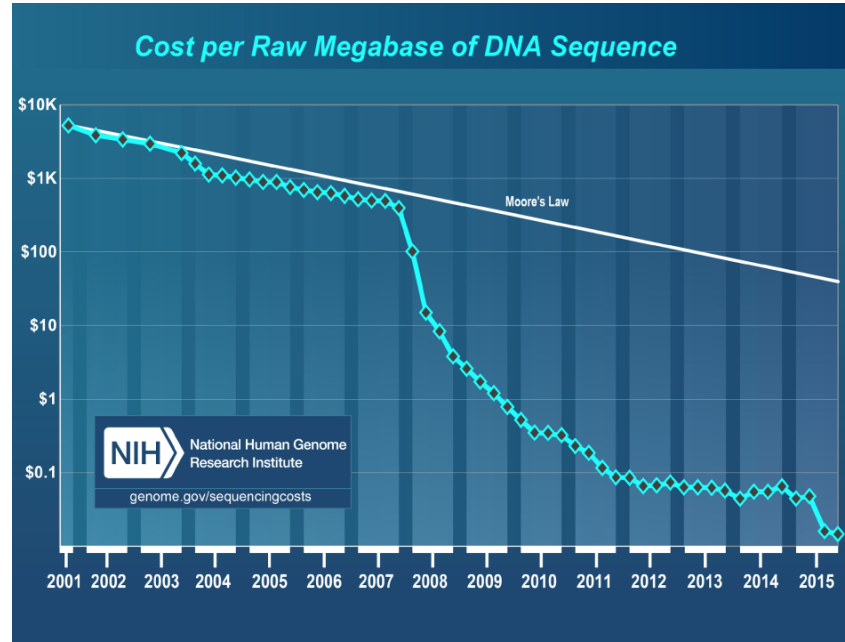- The type of sequence data deposited has changed as well.
  - Protected data represents an increasing fraction of all submitted sequences.



[Muir et al. ('15) GenomeBiol.]

# Sequencing Data Explosion: Faster than Moore's Law?

- 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).

## 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.



[Waldrop ('15) Nature]

# Kryder's Law and S-curves underlying exponential growth
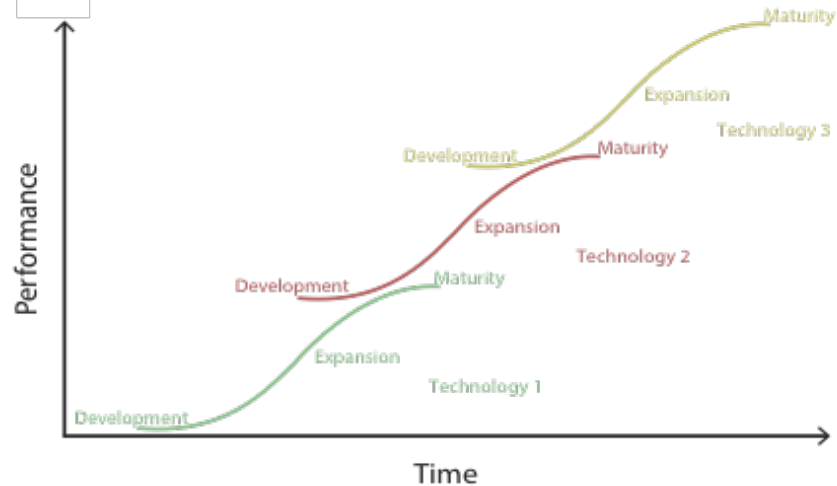
- Moore's & Kryder's Laws
  - As important as the increase in computer speed has been, the ability to store large amounts of information on computers is even more crucial
- Exponential increase seen in Kryder's law is a superposition of S-curves for different technologies

## Gigabytes per dollar over time

# 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

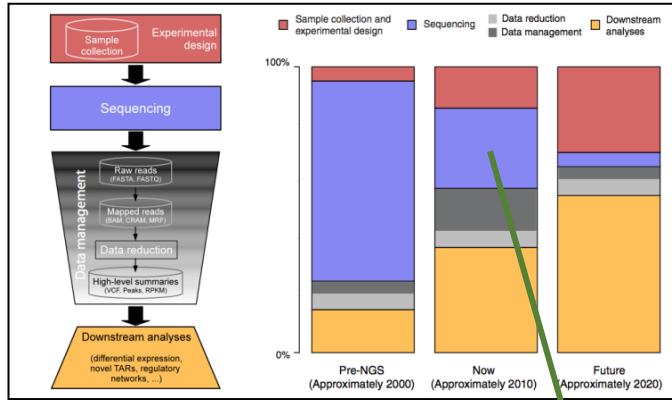[Sboner et al. ('11), Muir et al. ('15) Genome Biology]

# 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

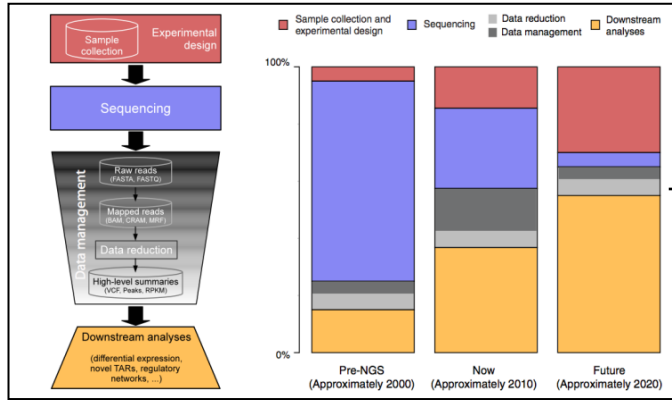[Sboner et al. ('11), Muir et al. ('15) Genome Biology]

# The changing costs of a sequencing pipeline



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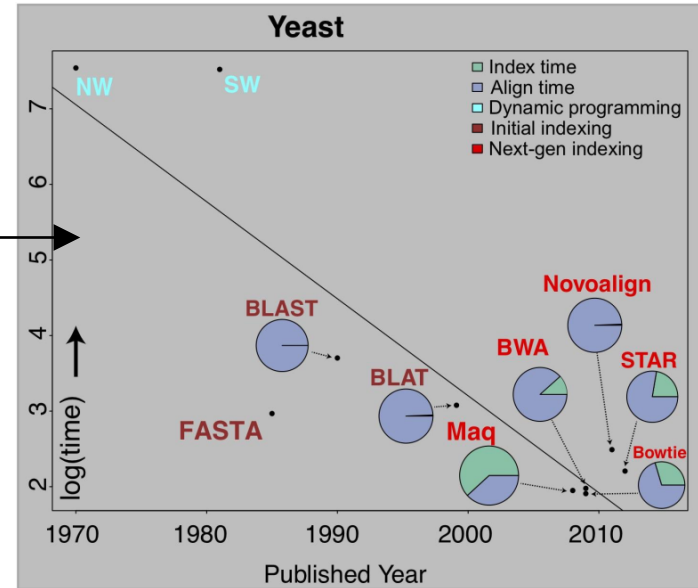from the actual seq. to sample
collection & analysis

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

# 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 |


Passenger
Driver (~0.1%)

### Class of Variants

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

### Prevalence of Variants


Common
Rare* (1-4%)

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


Common
Rare (~75%)

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

# Finding Key Variants

# Germline

**CAN YOU FIND THE PANDA?**



- **Common variants**
  - Can be most readily associated with phenotype (ie disease) via GWAS
  - 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 but often have larger functional impact
  - Can be collapsed in the same element to gain statistical power (burden tests).
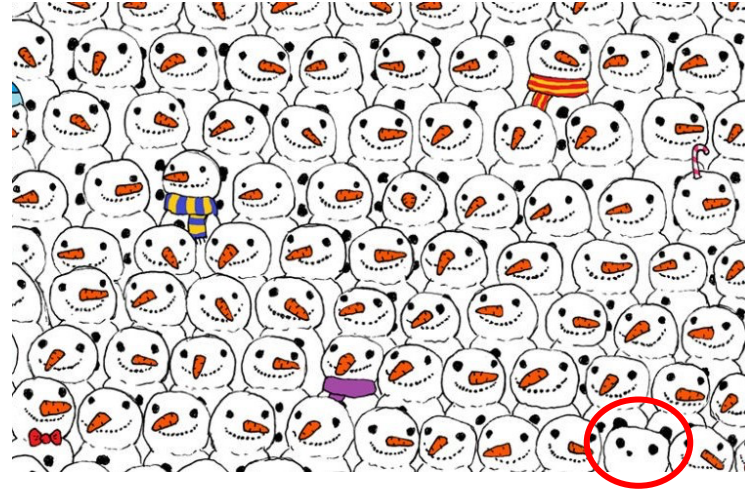
McCarthy, M. et al. Nat. Rev. Genet. 2008. 9, 356-369, Zuk, O. et al. PNSA. 2014. Vol. 11, no. 4, MacArthur DG et al. Nature 2014. 508:469-476

# Finding Key Variants

## Somatic



**CAN YOU FIND THE PANDA?**

- **Overall**
  - Often these can be thought of as <u>very rare variants</u>
- **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.

# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence:
  Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs

# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence:
  Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs

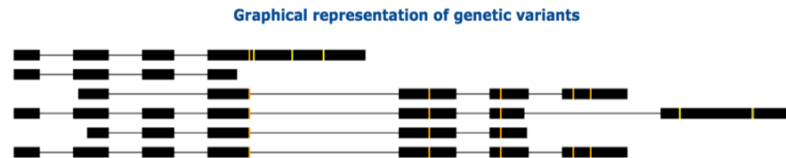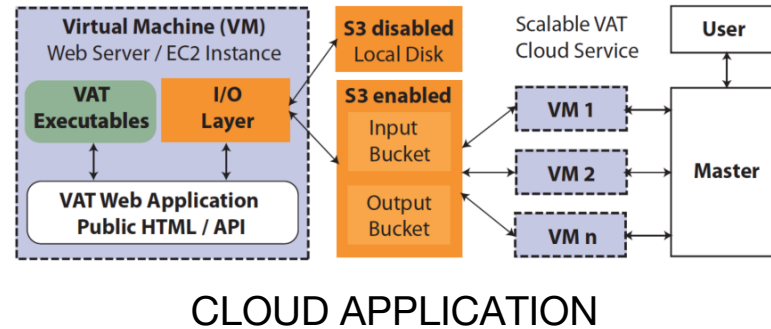# Variant Annotation Tool (VAT), developed for 1000G FIG

VCF Input

Output:
- Annotated VCFs
- Graphical representations of functional impact on transcripts

Access:
- Webserver
- AWS cloud instance
- Source freely available



CLOUD APPLICATION



Graphical representation of genetic variants

## vat.gersteinlab.org

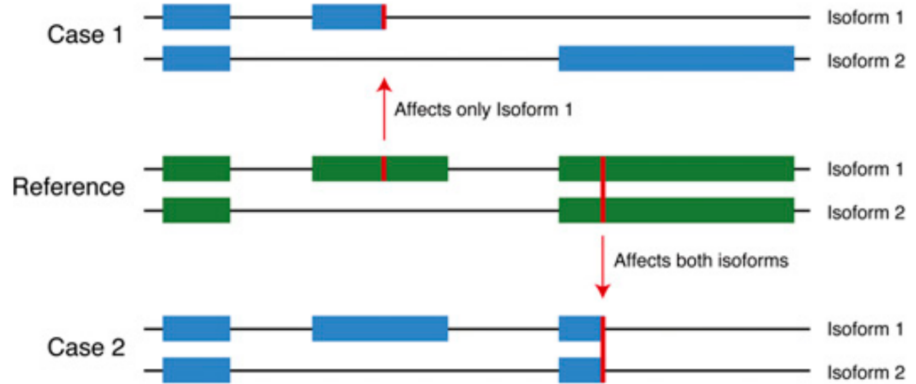*Habegger L.[*], Balasubramanian S.[*], et al. Bioinformatics, 2012*

# Complexities in LOF annotation

Transcript isoforms,
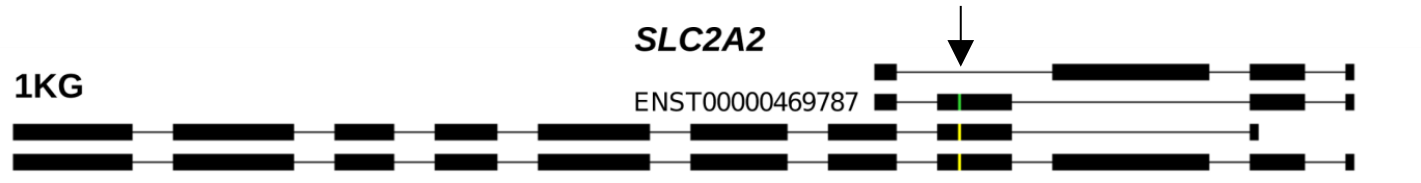distance to stop,
functional domains,
protein folding,
etc.

Balasubramanian S. et al., *Genes Dev.,* '11
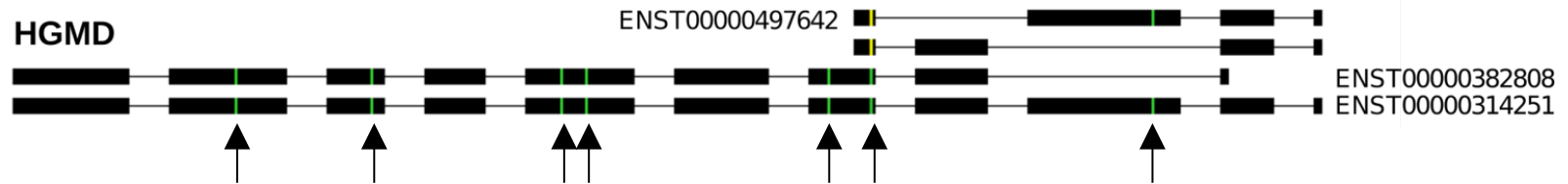Balasubramanian S.*, Fu Y.* et al., *NComms.*, '17



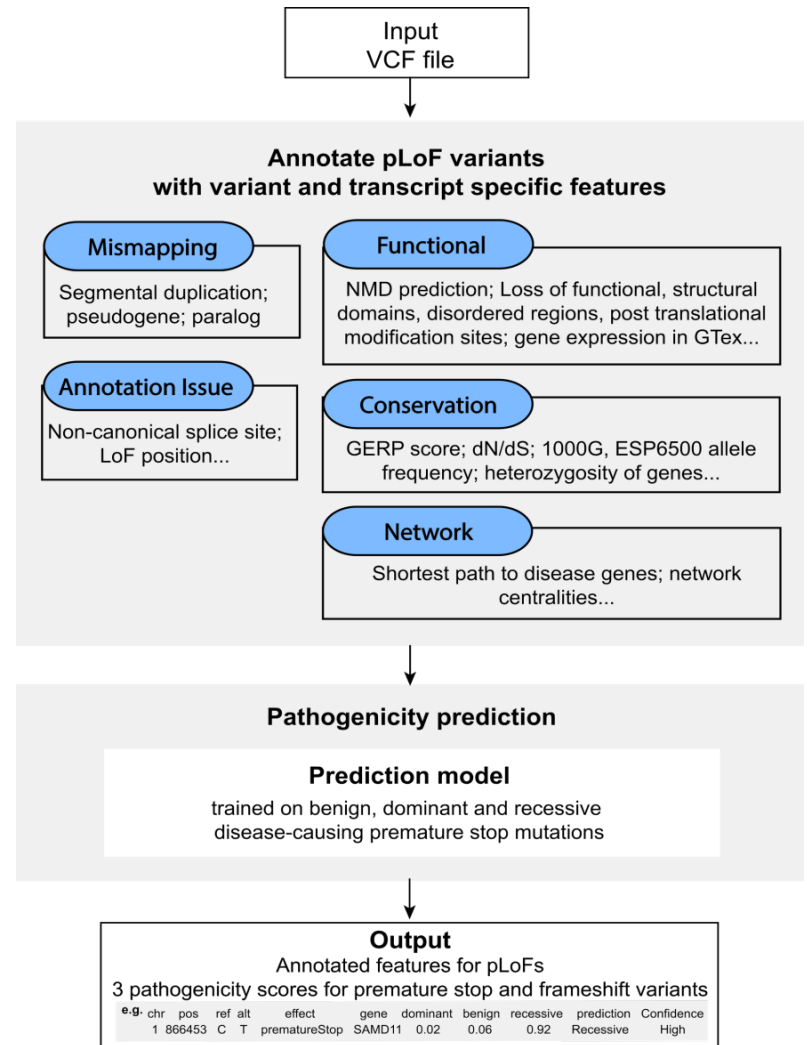Impact of a SNP on alternate splice forms

# Annotation of Loss-of-Function Transcripts (ALoFT)

Runs on top of VAT
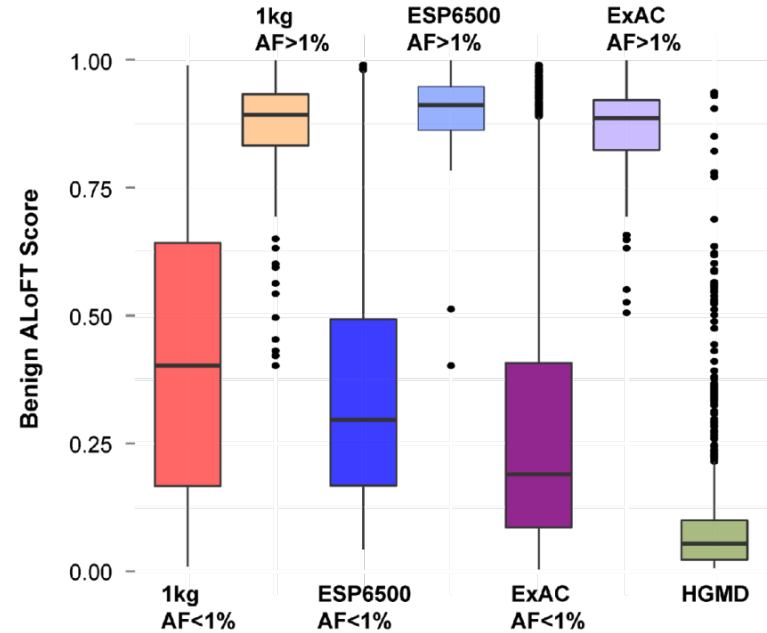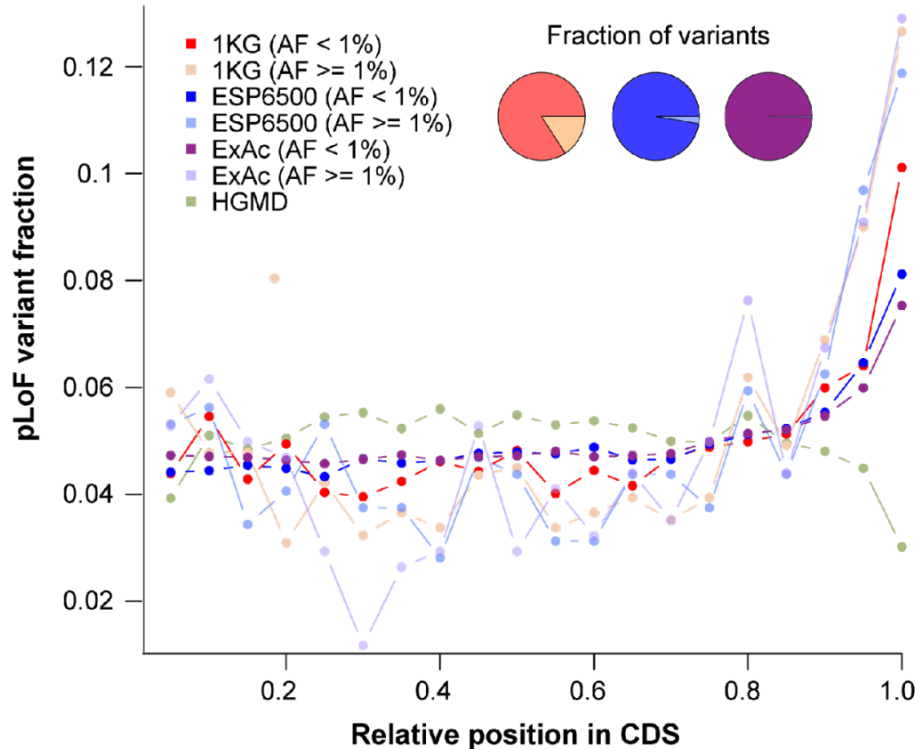
Output:

- Impact score: benign or deleterious.
- Decorated VCF.
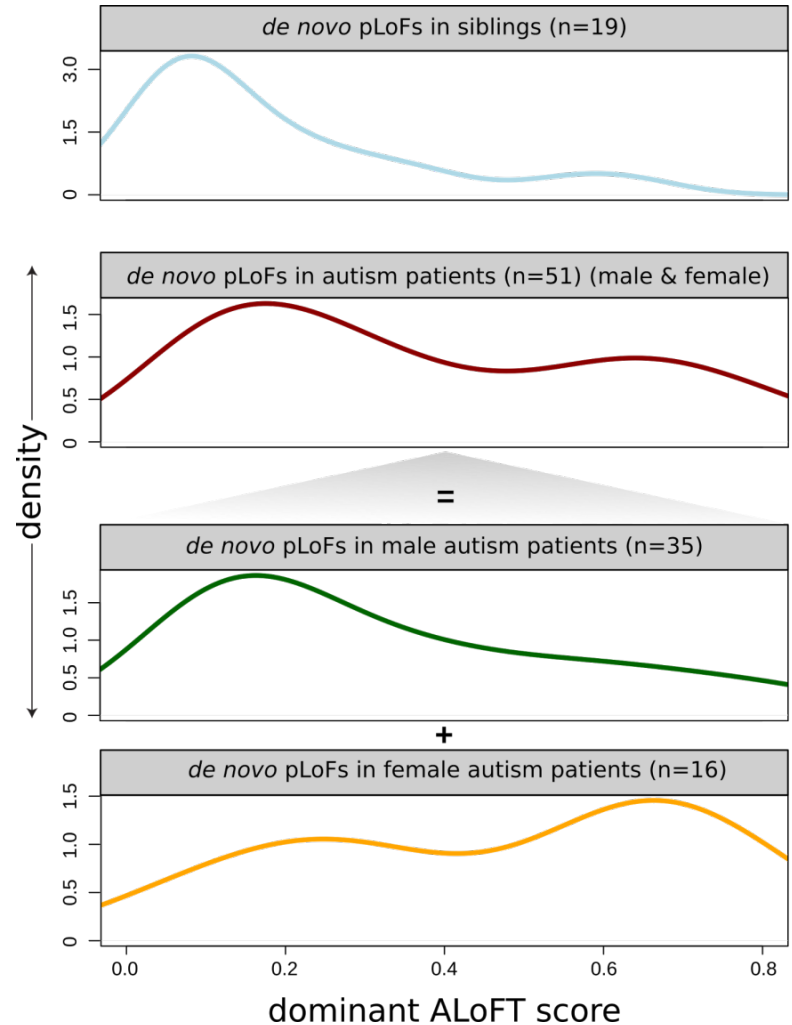


Balasubramanian S.*, Fu Y.*  et al., *NComms.,* '17

# LoF distribution varies as expected by mutation set (from healthy people v from disease)



Balasubramanian S.*, Fu Y.* et al., *NComms.,* '17

# Application to LoF mutations in autism spectrum disorder



*Balasubramanian S.\*, Fu Y.\*  et al., NComms., '17*

# ALoFT identifies deleterious somatic LoF variants

**Cancer genes:**

- COSMIC consensus.

- *Enriched in deleterious LoFs.*

**LoF tolerant genes:**
- LoF in the 1KG cohort.
- *Depleted in deleterious LoFs.*

Balasubramanian S.*, Fu Y.*  et al., *NComms.*, '17



cancer genes vs. LoF tolerant genes

- 504 cancer genes
- 387 LoF-tolerant genes
- 504 random genes
- 387 random genes

y-axis: percentage of somatic pLoF variants in gene sets

x-axis: 1-benign ALoFT score

# ALoFT refines cancer mutation characterization

## 20/20 rule ALoFT stratification



*Vogelstein et al. '13:* if >20% of mutations in gene inactivating → tumor suppressor gene (TSG).

ALoFT further refines 20/20 rule predictions.

## deleterious LoFs / total mutations



## deleterious LoFs / total LoFs



Balasubramanian S.*, Fu Y.* et al., *NComms.,* '17

# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
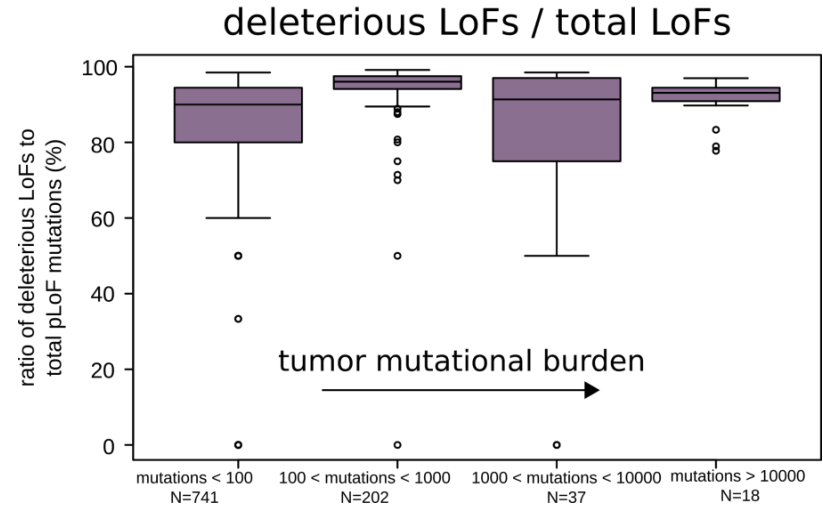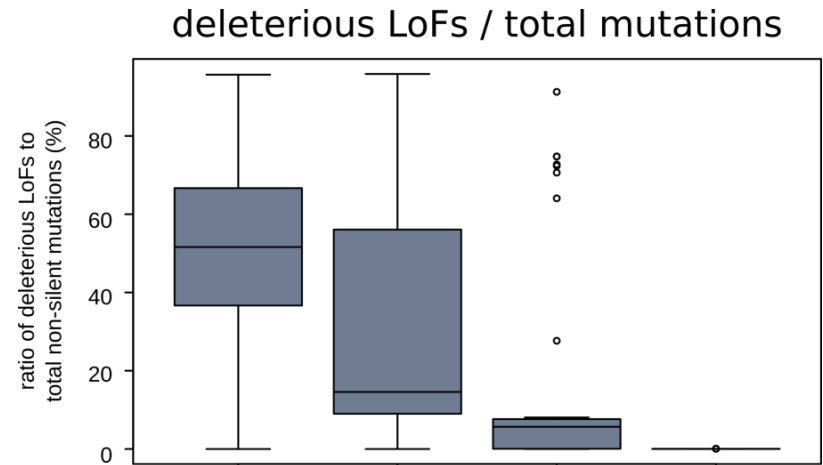  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence:
  Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs

**What is localized frustration?**

more negative — more positive

favorable interaction

unfavorable interaction

ASN ⟷ ASP

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

# Workflow for evaluating localized frustration changes (ΔF)

# Complexity of the second order frustration calculation

# Comparing ΔF values across different SNV categories: disease v normal



Normal mutations (1000G) tend to unfavorably frustrate (less frustrated) surface more than core, but for disease mutations (HGMD) no trend & greater changes

[Kumar et al, *NAR* (2016)]

# Comparison between ΔF distributions: TSGs v. oncogenes



SNVs in TSGs change frustration more in core than the surface, whereas those associated with oncogenes manifest the opposite pattern. This is consistent with differences in LOF v GOF mechanisms.

# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence:
  Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs
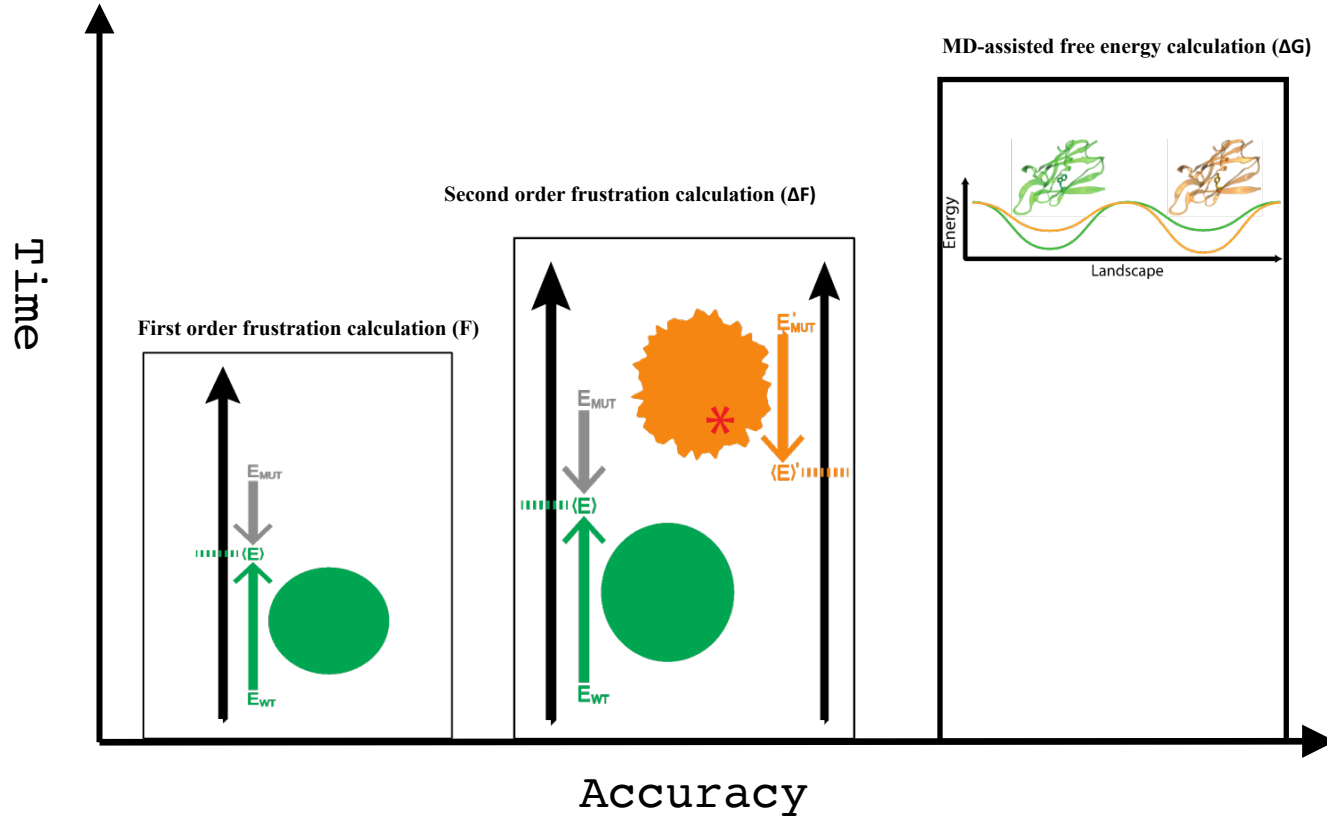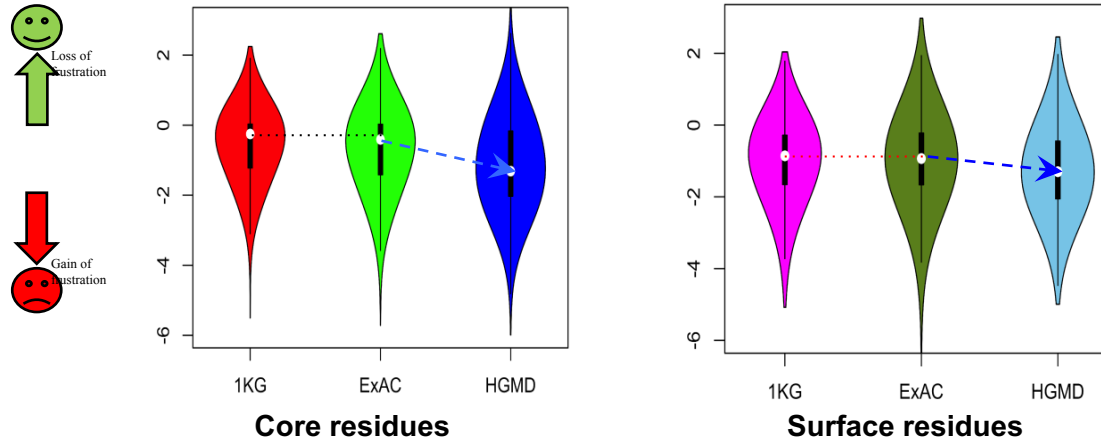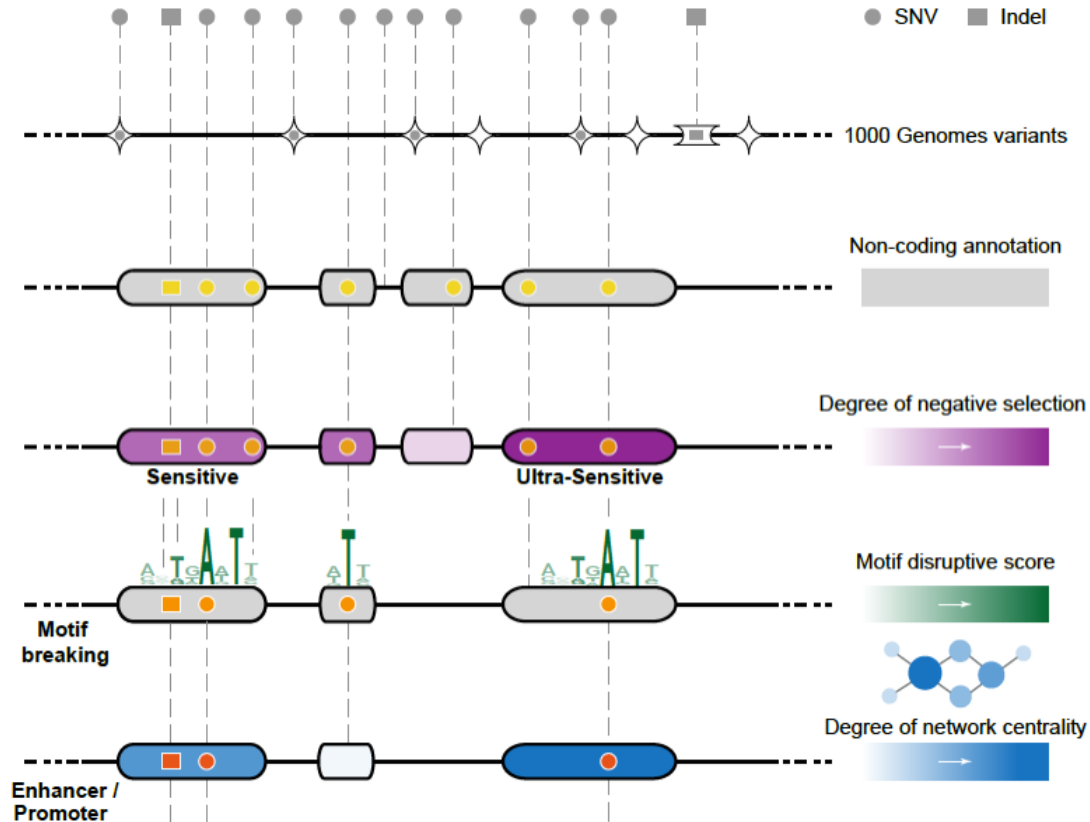
# Funseq: a flexible framework to determine functional impact & use this to prioritize variants



**Annotation (tf binding sites open chromatin, ncRNAs) & Chromatin Dynamics**

**Conservation (GERP, allele freq.)**

**Mutational impact (motif breaking, Lof)**

**Network (centrality position)**

# Finding "Conserved" Sites in the Human Population:

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

**Broad Categories**



(Non-coding RNA) ncRNA

(DNase I hypersensitive sites) DHS

(Transcription factor binding sites) TFBS

(TFSS: Sequence-specific TFs)

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

**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

# Differential selective constraints among specific sub-categories

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

~0.4% genomic coverage (~ top 25)

~0.02% genomic coverage (top 5)

# Defining Sensitive non-coding Regions

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)]

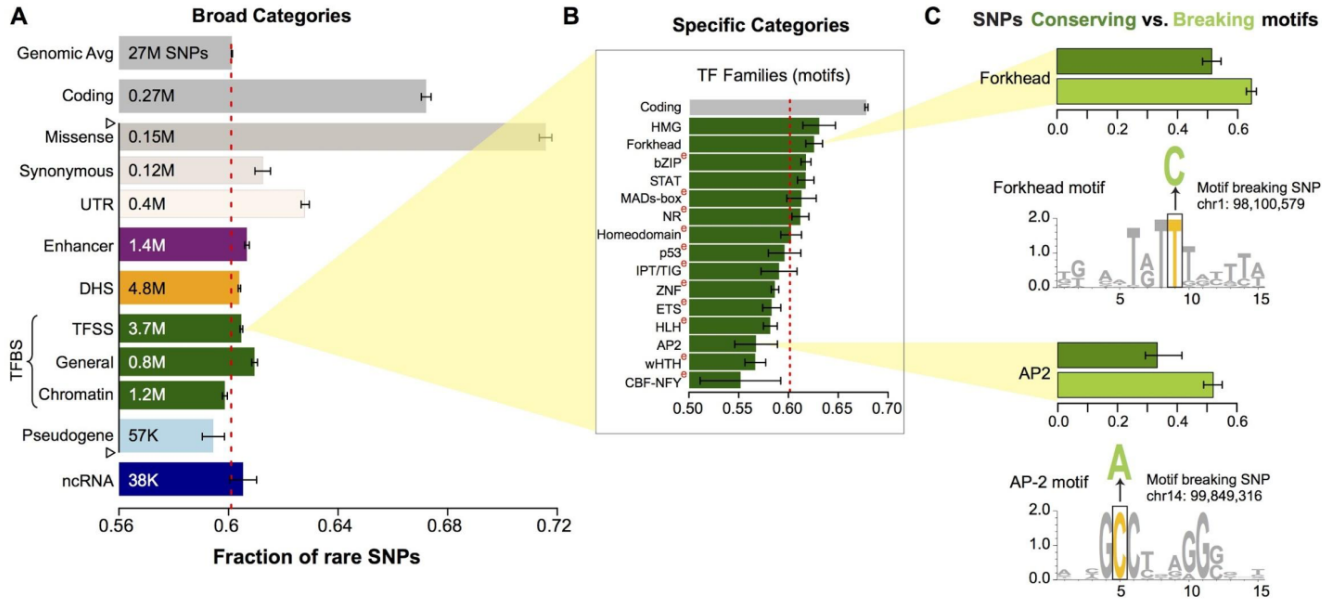# SNPs which break TF motifs are under stronger selection



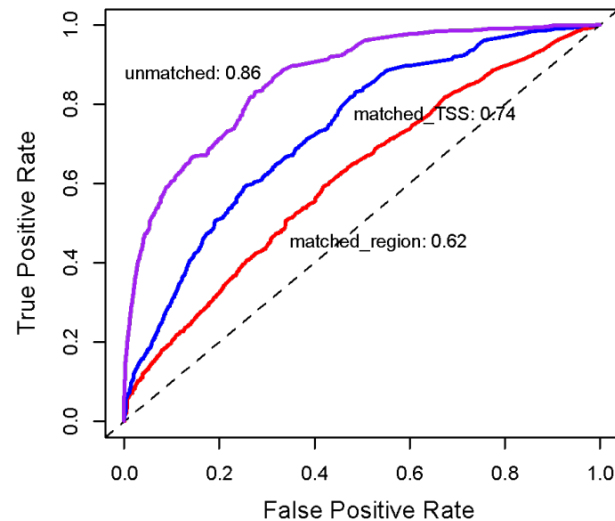[Khurana et al., *Science* ('13)]

# FunSeq.gersteinlab.org

$$w_d = 1 + p_d log_2 p_d + (1 - p_d) log_2 (1 - p_d)$$

- Entropy based method for weighting consistently many genomic features

- Practical web server
- Submission of variants & pre-computed large data context from uniformly processing large-scale datasets

[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)]

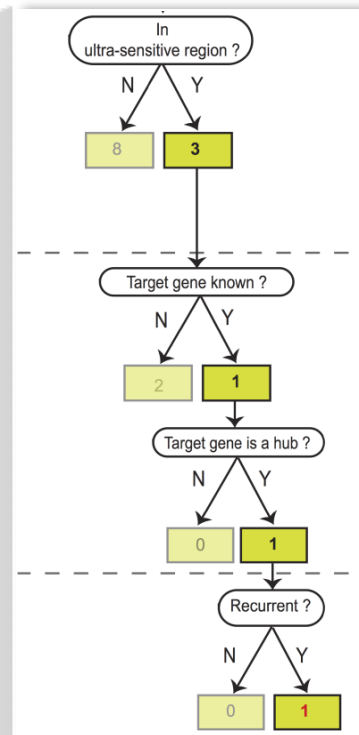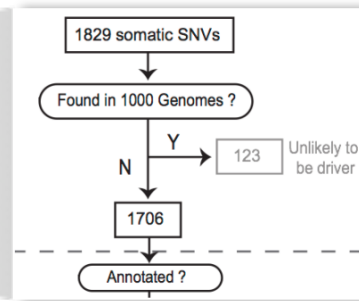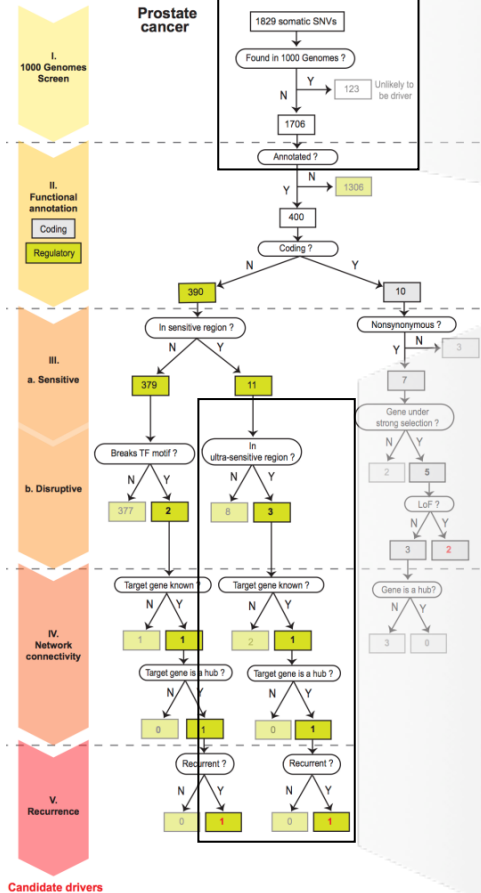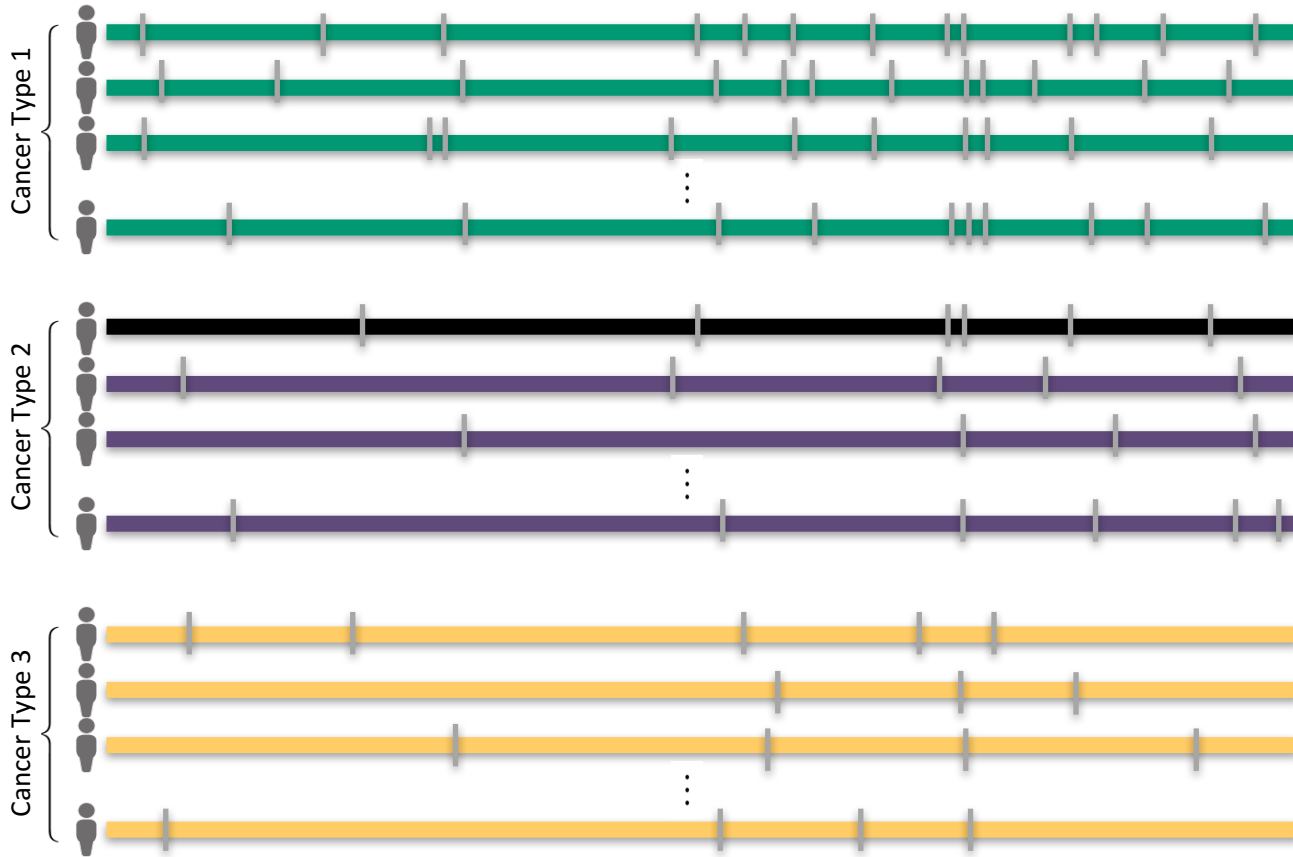Flowchart for 1 Prostate Cancer Genome (from Berger et al. '11)
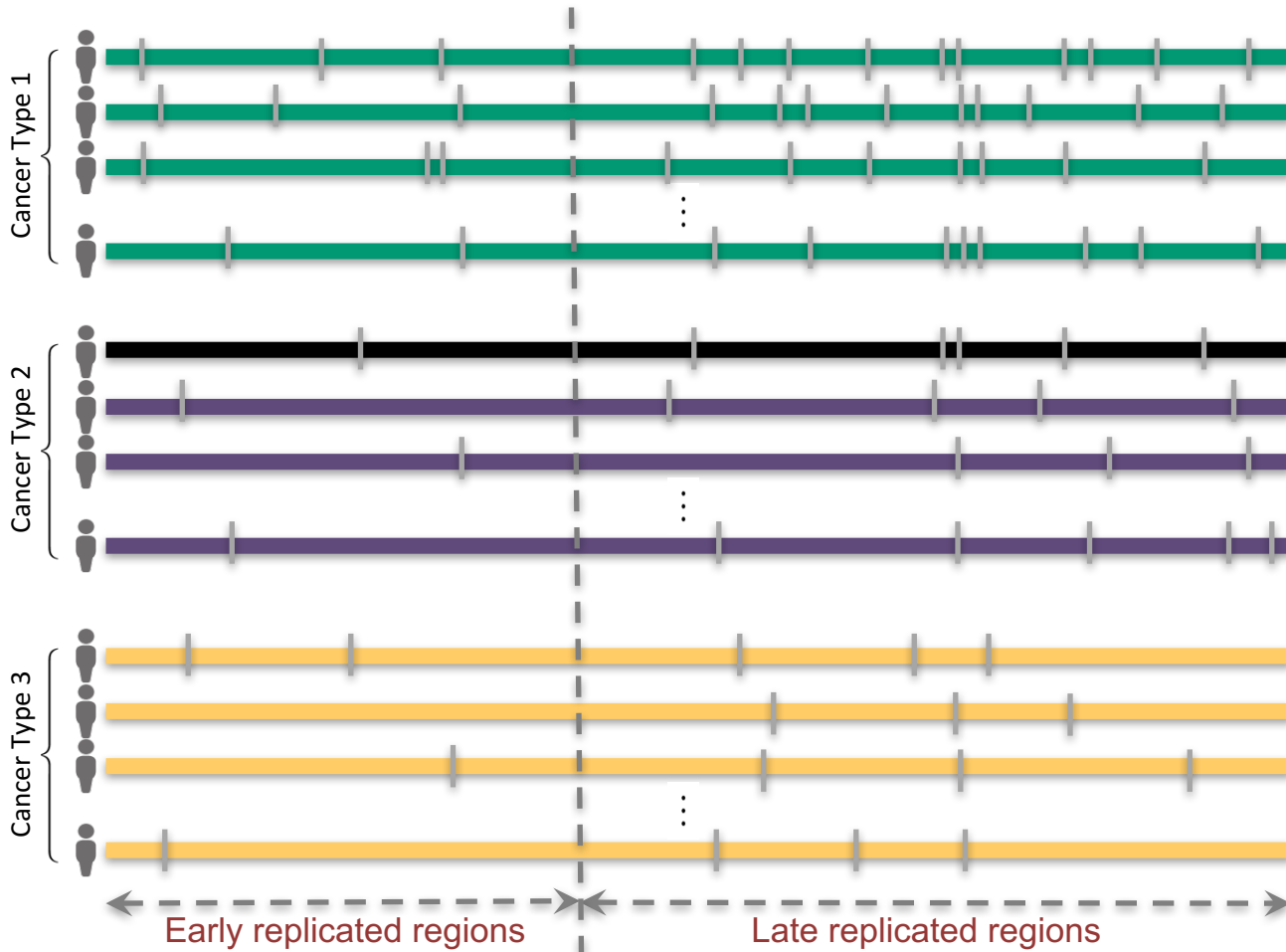
# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence: Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs

# 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**



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

46

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

**Chromatin remodeling failure leads to more mutations in early-replicating regions**

**Variation in somatic mutations is closely associated with chromatin structure (TADs) & replication timing**

genomic distance **from the TAD boundary**

mutation load (standardized)

[Yan et al., *PLOS Comp. Bio. ('17)*;  S. Li et al., PLOS Genetics ('17)] ]

# mrTADFinder:
## Identifying TADs at multiple resolutions by maximizing modularity vs appropriate null



input: contact map W

null model E

Choose a particular resolution γ
Optimize Q over all possible partitions

$$Q = \frac{1}{2N} \sum_{ij} (W_{ij} - \gamma E_{ij}) \delta_{\sigma_i \sigma_j}$$   γ: resolution parameter

Multiple runs to define boundary scores for all pairs of adjacent bins

consensus boundaries based on the boundary scores

consensus TADs   output

γ=2

22.0   24.0   26.0   28.0   30.0   32.0   34.0   36.0

γ=2.5

22.0   24.0   26.0   28.0   30.0   32.0   34.0   36.0

γ=3

22.0   24.0   26.0   28.0   30.0   32.0   34.0   36.0

[Yan et al., *PLOS Comp. Bio.* ('17)]

# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence: Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs

# Cancer Somatic Mutation Modeling

## PARAMETRIC MODELS

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

$x_i : Binomial(n_i, p)$

**Model 2a: Varying Mutation Rate with Single Covariate Correction**

$x_i : Binomial(n_i, p_i)$

$p_i : Beta(\mu|R_i, \sigma|R_i)$

$\mu|R_i, \sigma|R_i$ : constant within the same covariate rank

**Model 2b: Varying Mutation Rate with Multiple Covariate Correction**

$x_i : Binomial(n_i, p_i)$

$p_i : Beta(\mu|\boldsymbol{R_i}, \sigma|\boldsymbol{R_i})$

$\mu|\boldsymbol{R_i}, \sigma|\boldsymbol{R_i}$ : constant within the same covariate rank

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

- 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
  - $R_i$: the covariate rank of the element

- Non-parametric model is useful when covariate data is missing for the studied annotations
  - Also sidesteps issue of properly identifying and modeling every relevant covariate (possibly hundreds)

## NON-PARAMETRIC MODELS

Assume constant background mutation rate in local regions.

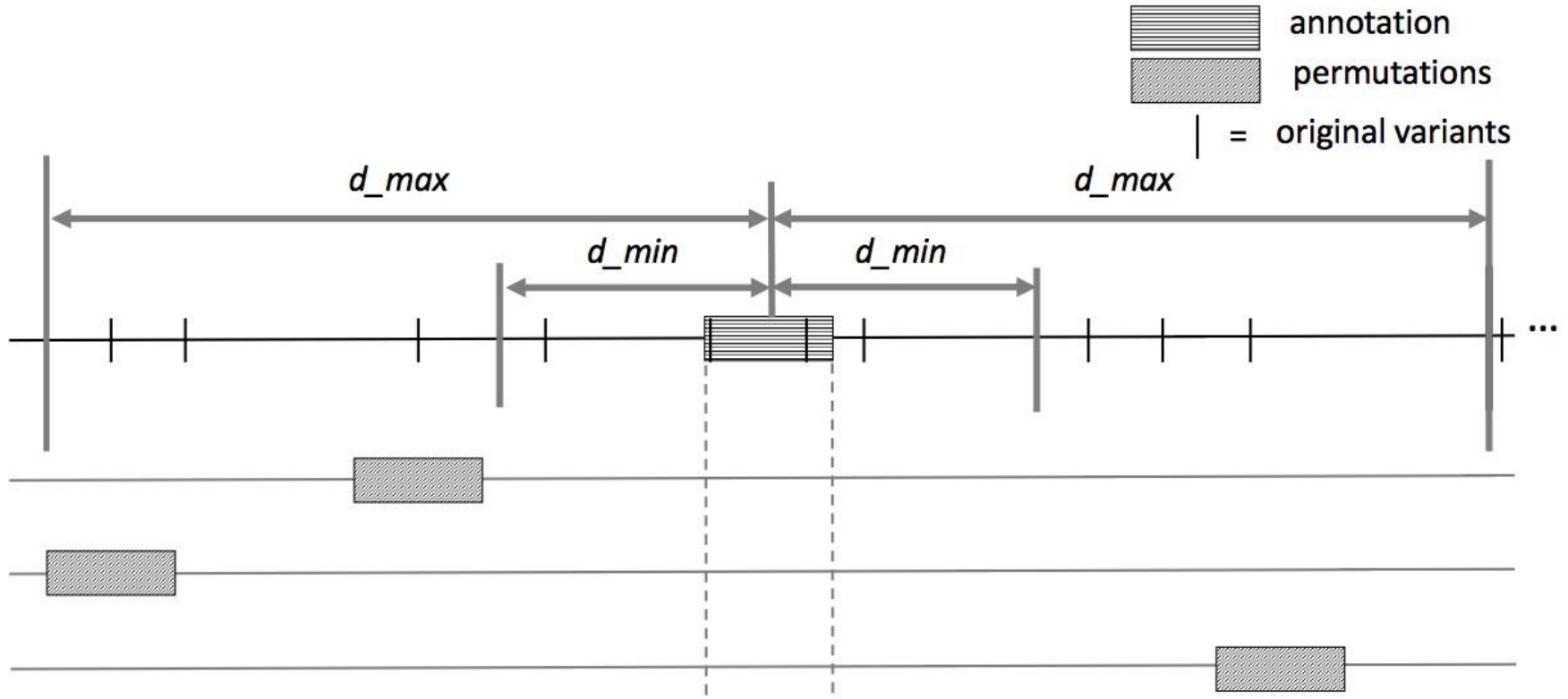**Model 3a: Random Permutation of Input Annotations**

Shuffle annotations within local region to assess background mutation rate.

**Model 3b: Random Permutation of Input Variants**

Shuffle variants within local region to assess background mutation rate.

[Lochovsky et al. *Bioinformatics* in press]

# MOAT-a: Annotation-based permutation

[Lochovsky et al. *Bioinformatics* in

# MOAT-v: Variant-based Permutation

Can preserve tri-nt context in shuffle



annotation

| = original variants

⋮ = permuted variants

bin width $W$

$W \approx 2 \cdot d\_max$

# MOAT-s: a variant on MOAT-v

- A somatic variant simulator
  - Given a set of input variants, shuffle to new locations, taking genome structure into account

# 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

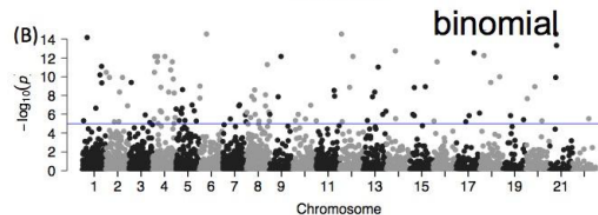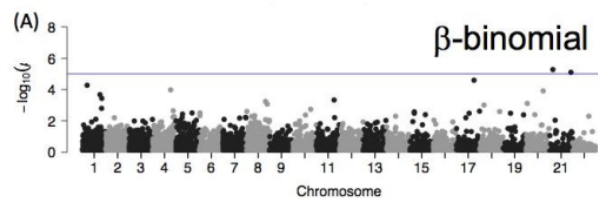# LARVA Results



TSS LARVA results

adjusted P w/. correction

PRRC2B
TP53
LMO3
AGAP5,PROZ
TERT

These have literature-verified cancer associations

noncoding annotation p-values in sorted order

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

$-\log10(\text{Pvalues})$

(A) β-binomial

(B) binomial

# MOAT: recapitulates LARVA with GPU-driven runtime scalability

| Gene Name | Documented role with cancer | Pubmed ID |
|---|---|---|
| SLC3A1 | Cysteine transporter SLC3A1 promotes breast cancer tumorigenesis | 28382174 |
| ADRA2B | reduce cancer cell proliferation, invasion, and migration | 25026350 |
| SIL1 | subtype-specific proteins in breast cancer | 23386393 |
| TCF24 | NA | NA |
| AGAP5 | significant mutation hotspots in cancer | 25261935 |
| TMPRSS13 | Type II transmembrane serine proteases in cancer and viral infections | 19581128 |
| ERO1L | Overexpression of ERO1L is Associated with Poor Prognosis of Gastric Cancer | 26987398 |

.
.
.

MOAT's high mutation burden elements recapitulate LARVA's results & published noncoding cancer-associated elements.

Computational efficiency of MOAT's NVIDIA™ CUDA™ version, with respect to the number of permutations, is dramatically enhanced compared to CPU version.

| Number of permutations | Fold speedup of CUDA version |
|---|---|
| 1k | 14x |
| 10k | 100x |
| 100k | 256x |

# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence: Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs

# Genomics & Data Science:
## Approaches to identifying key variants through functional impact & recurrence

- Introduction
  - An individual's disease variants as the public's gateway into genomics & biology
  - **The exponential scaling** of data generation & processing
  - Mining the data to prioritize variants for key drivers
- Functional impact #1: Coding
  - **ALoFT**: Annotation of Loss-of-Function Transcripts.
  - LoF annotation as a complex problem + finding deleterious LoFs
  - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs

- Functional impact #2: Non-coding
  - **FunSeq** integrates evidence, with an entropy based weighting scheme.
  - Prioritizing rare variants with "sensitive sites" (human conserved)
- Recurrence:
  Statistics for driver identification
  - **Background mutation rate** significantly varies & is correlated with replication timing & TADs
  - Developed a variety of parametric & non-parametric methods taking this into account
  - **LARVA** uses parametric beta-binomial model, explicitly modeling covariates
  - **MOAT** does a variety of non-parm. shuffles (annotation, variants, &c). Useful when explicit covariates not available. Slower than but speeded up w/ GPUs

github.com/gersteinlab/**Frustration**

S **Kumar**, D Clarke

github.com/gersteinlab/**MrTADfinder**

KK **Yan**, S Lou

**VAT**.gersteinlab.org

L **Habegger**, S Balasubramanian,
DZ Chen, E Khurana, A Sboner,
A Harmanci, J Rozowsky, D Clarke, M Snyder

**ALoFT.**gersteinlab.org

S **Balasubramanian**,
Y **Fu**, M Pawashe, P McGillivray,
M Jin, J Liu, K Karczewski, D MacArthur

**FunSeq**.gersteinlab.org

Y **Fu**, E **Khurana**, Z Liu,
S Lou, J Bedford, XJ Mu, KY Yip

**CostSeq2**

P **Muir**, S Li, S Lou, D Wang, DJ Spakowicz, L
Salichos, J Zhang, GM Weinstock, F Isaacs, J
Rozowsky

**LARVA**.gersteinlab.org

L **Lochovsky**, J **Zhang**,
Y Fu, E Khurana

**MOAT**.gersteinlab.org

L **Lochovsky**, J **Zhang**

# Info about this talk

## General PERMISSIONS

## PHOTOS & IMAGES