

# Computational analysis of variants: coding versus non-coding

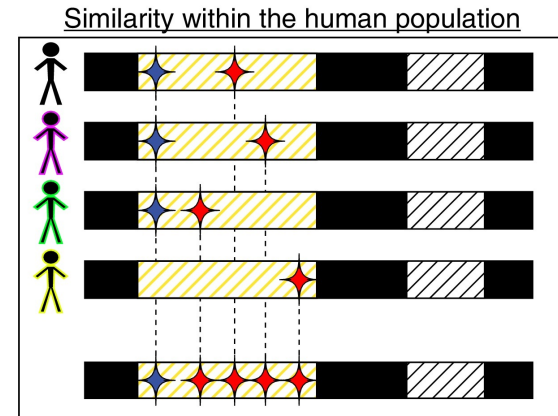
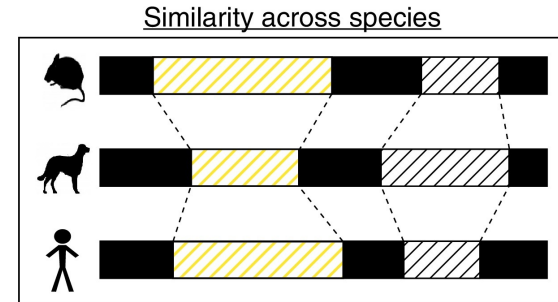
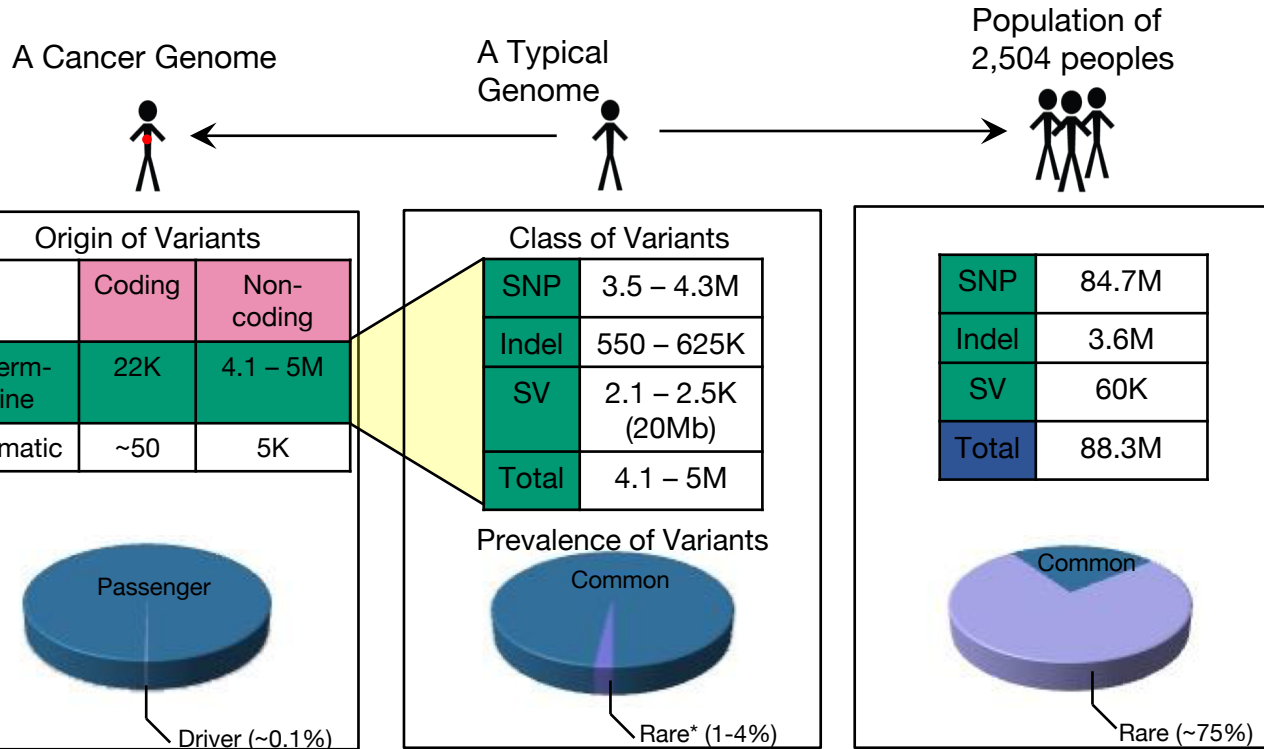


Mark Gerstein, Yale

Slides freely downloadable from [Lectures.GersteinLab.org](http://Lectures.GersteinLab.org) & “tweetable” (via [@MarkGerstein](https://twitter.com/MarkGerstein)).

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

# Human Genetic Variation: the prevalence of rare variants in population studies



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



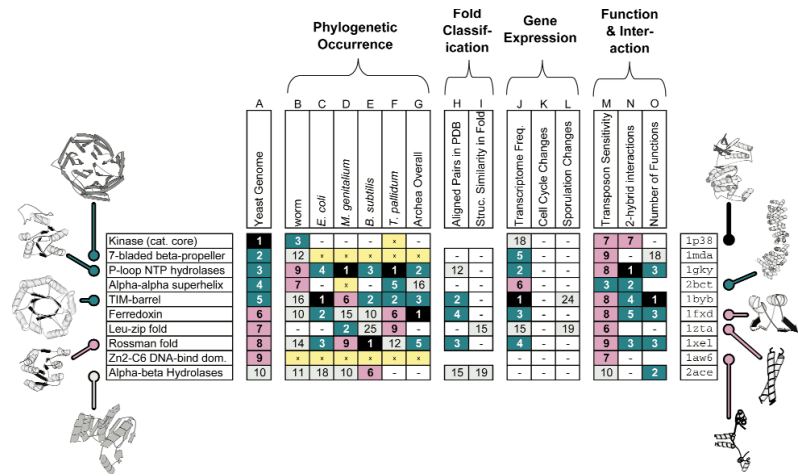
# Coding v Non-coding

- Coding
  - Easily interpretable, particularly related to structure
  - Available in large quantities
  - Exomes have the current potential for great scale (Scale of EXAC, >60K exomes [Lek et al. '16])
- Non-coding
  - Not as interpretable & hard to connect to genes
- “Near coding”
  - Bits of non-coding, close to genes & readily linked to them
  - EX: Splice sites, promoters, uORF

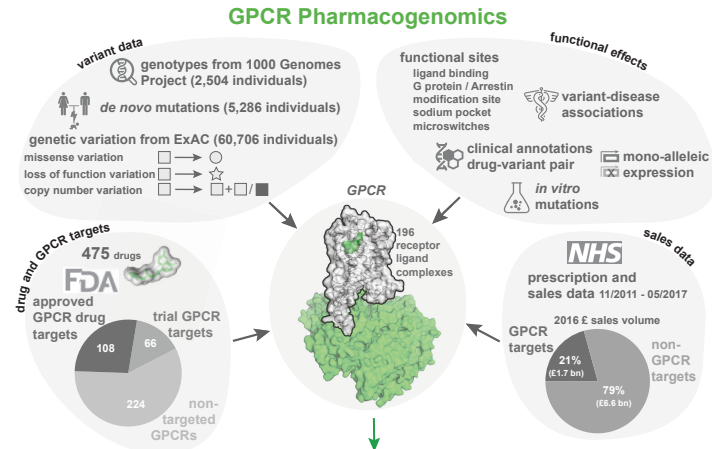
# Structure & genomics

Structure particularly useful for interpreting the impact of the many rare variants whose effect can not be found via GWAS

Also, integration of structure data with genomic variants, EHR & drug data will be key for realizing the goal of precision medicine.



Gerstein et al., Nat. Struct. Bio. 2000

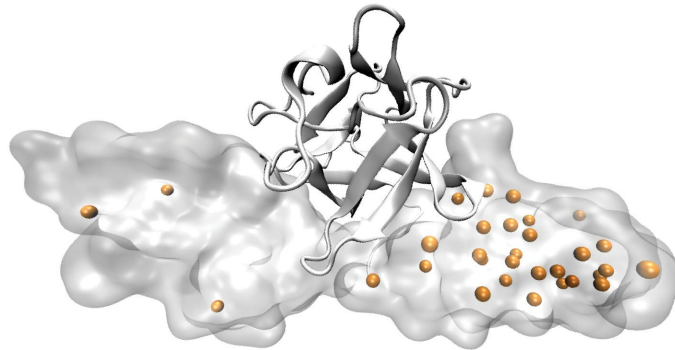


Assessing the spectrum, prevalence and functional impact of genetic variation for alteration in drug response

Hauser et al., Cell 2018

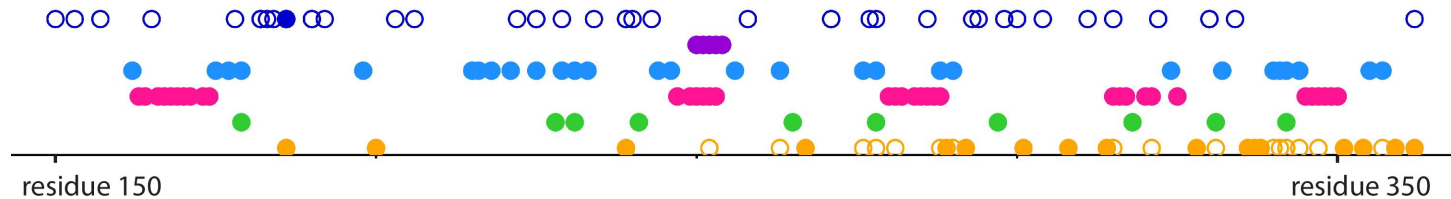
Unlike common SNVs, the statistical power with which we can evaluate rare SNVs in case-control studies is severely limited

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)



[Sethi et al. COSB ('15)]

# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**

- Rare v common,  
somatic v germline, coding v noncoding

- **Identifying cryptic allosteric sites with STRESS**

- On surface & in interior bottlenecks

- **Frustration as a localized metric of SNV impact**

- Differential profiles for oncogenes v. TSGs

- **ALoFT: Annotation of LoF Transcripts**

- **Using dynamics to help identify mutation clusters (Hotcommics)**

- Find dynamic sub-communities & determine aggregated mutational burden within these

- **RADAR Prioritization for RBP sites**

- Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
- Incorporates new features related to RNA sec. struc & tissue specific effects

- **uORF Prioritization**

- Feature integration to find small subset of upstream mutations that potentially alter translation

- **GRAM to assess the molecular effect of (promotor) mutations**

- Universal score + cell type specific score

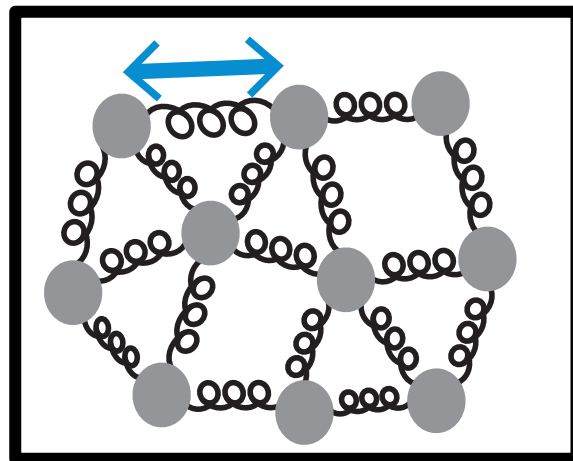
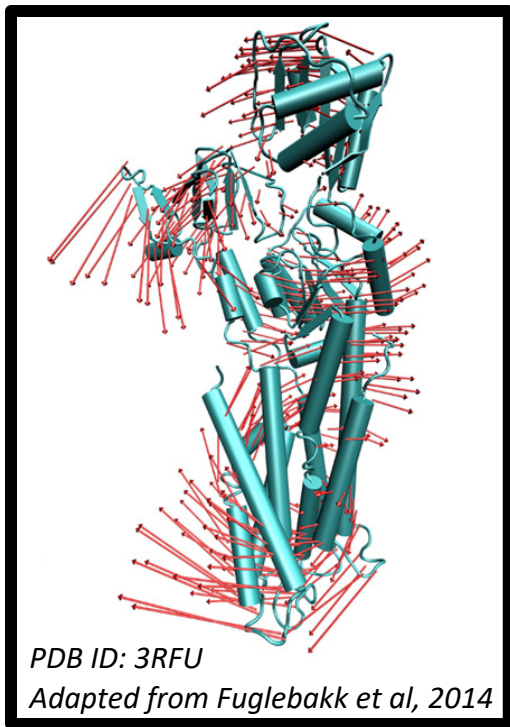
# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with STRESS**
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (Hotcommics)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score



# Models of Protein Conformational Change

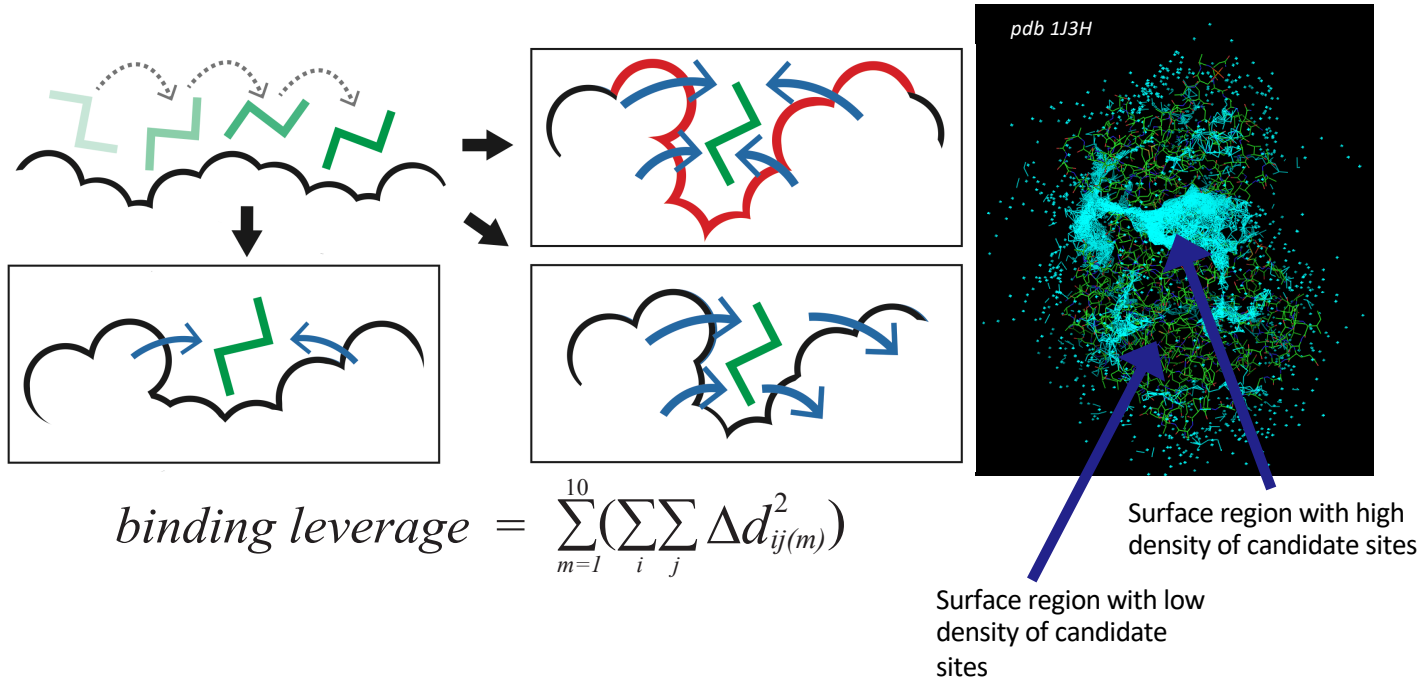
## Motion Vectors from Normal Modes (ANMs)



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

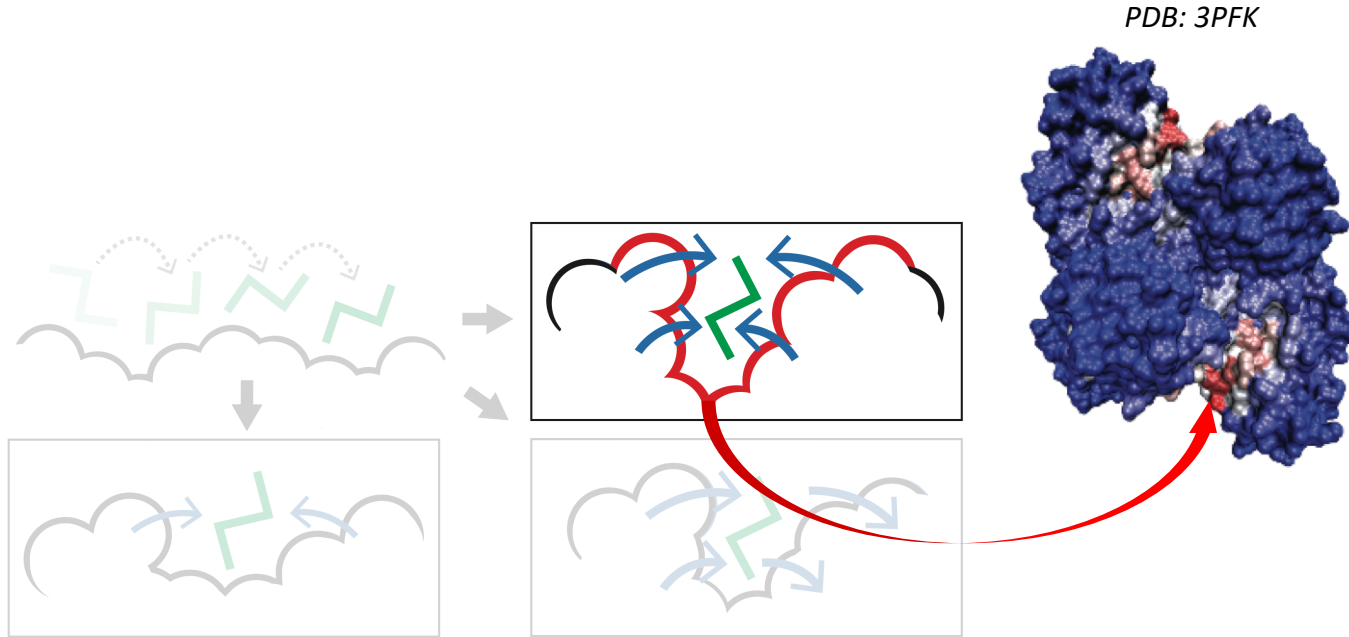
# Predicting Allosterically-Important Residues at the Surface

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



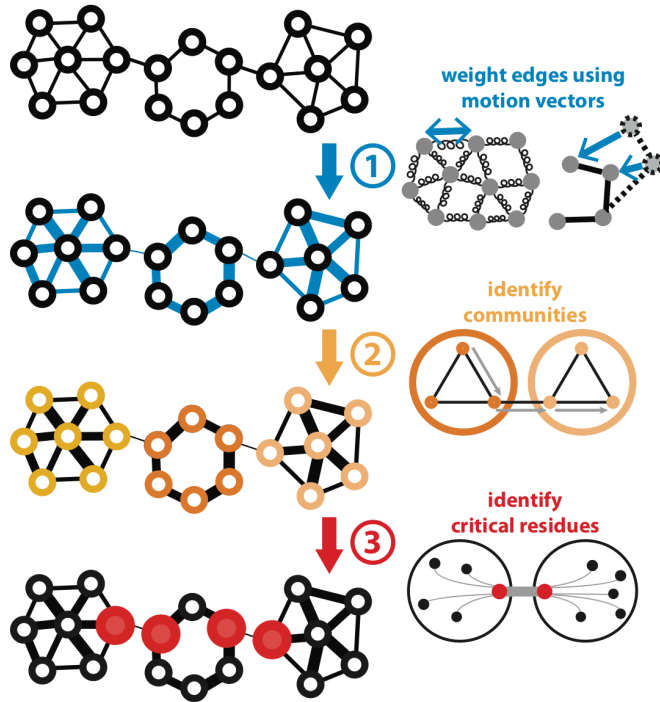
Adapted from Clarke\*, Sethi\*, et al ('16)

# Predicting Allosterically-Important Residues at the Surface



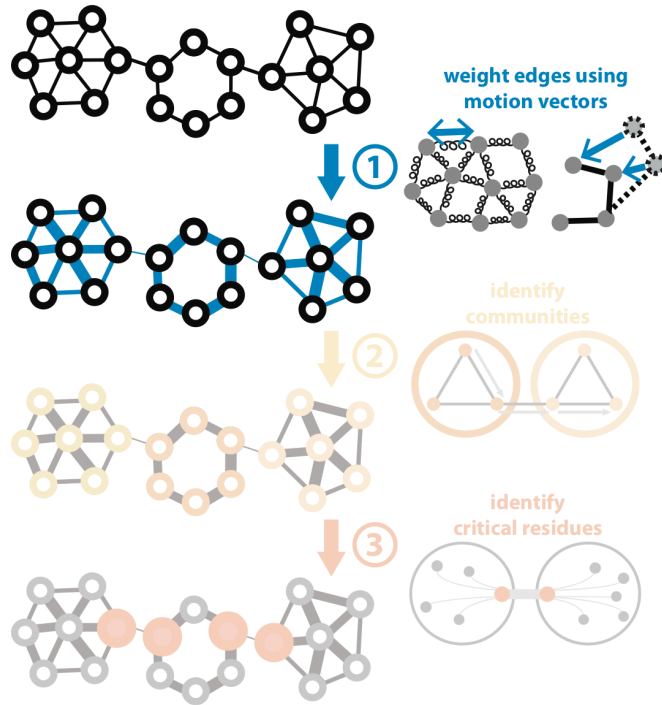
Adapted from Clarke\*, Sethi\*, et al ('16)

# Predicting Allosterically-Important Residues within the Interior



Adapted from Clarke\*, Sethi\*, et al ('16)

# Predicting Allosterically-Important Residues within the Interior

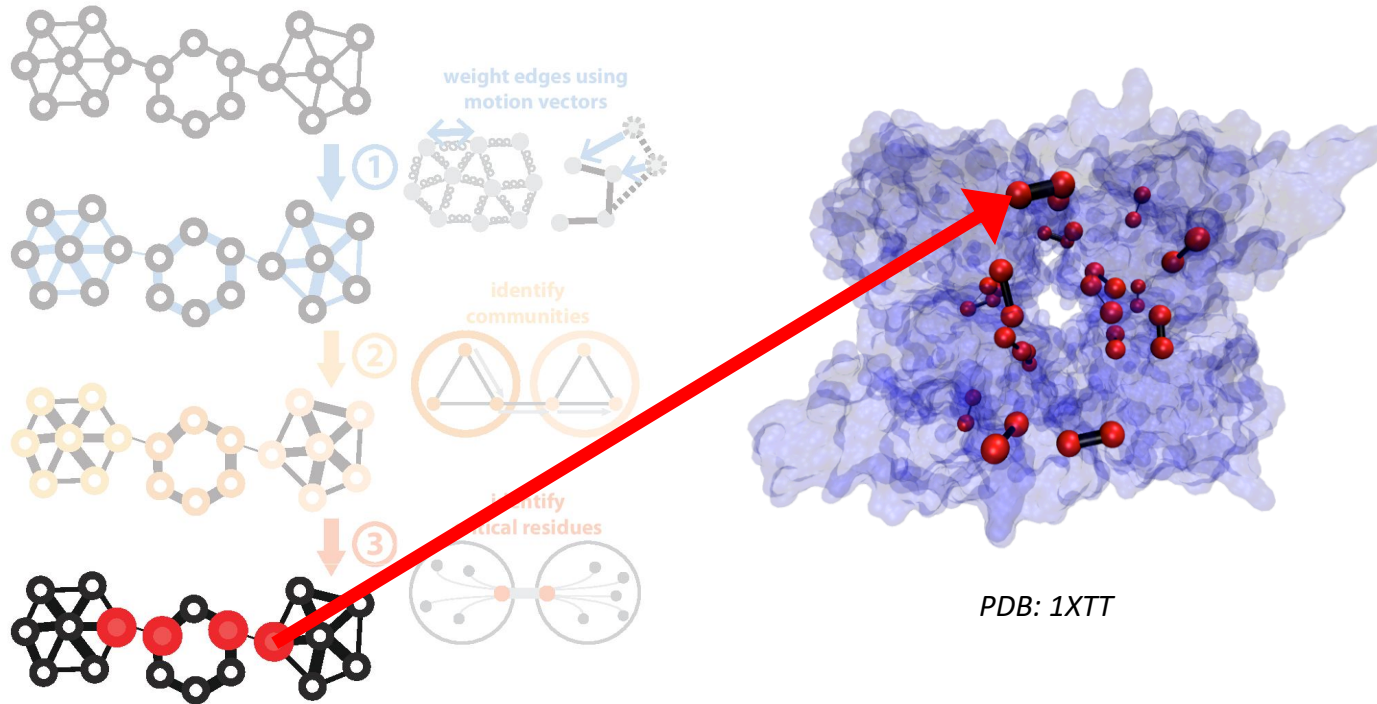


$$Cov_{ij} = \langle \mathbf{r}_i \cdot \mathbf{r}_j \rangle$$

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

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

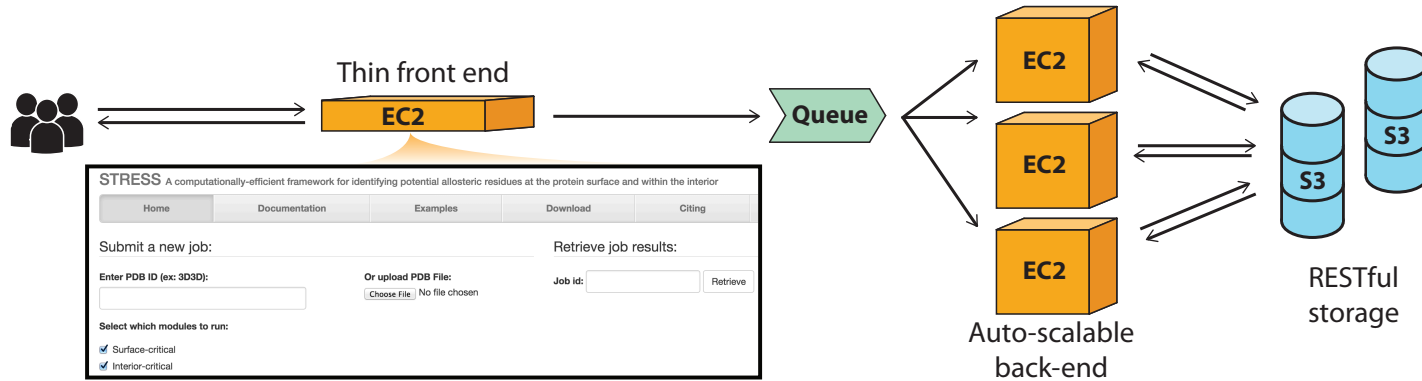
# Predicting Allosterically-Important Residues within the Interior



Adapted from Clarke\*, Sethi\*, et al ('16)

# STRESS Server Architecture: Highlights

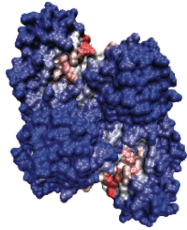
stress.molmovdb.org



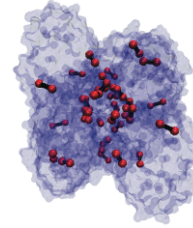
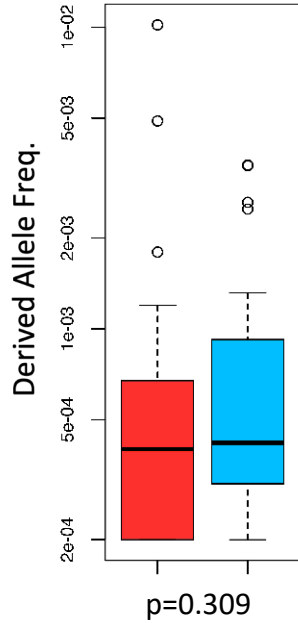
- A light front-end server handles incoming requests, and powerful back-end servers perform calculations.
- Auto Scaling adjusts the number of back-end servers as needed.
- A typical structure takes ~30 minutes on a E5-2660 v3 (2.60GHz) core.
- Input & output (i.e., predicted allosteric residues) are stored in S3 buckets.

# Intra-species conservation of predicted allosteric residues

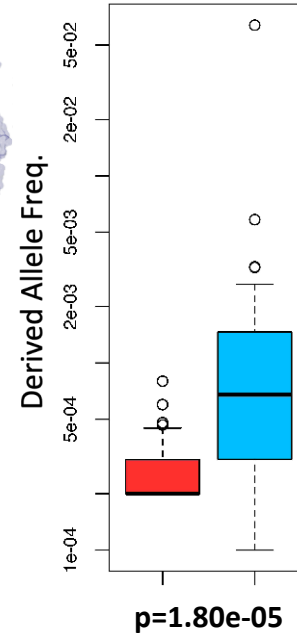
## 1000 Genomes



### Surface



### Interior



 critical  
 non-critical

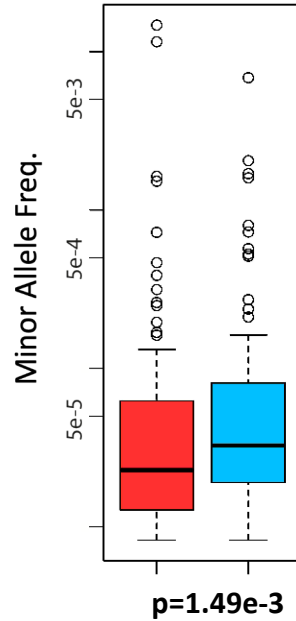
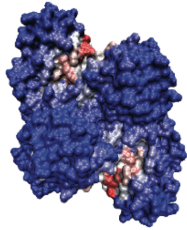


# Intra-species conservation of predicted allosteric residues

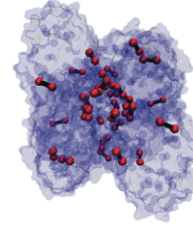
*ExAC*



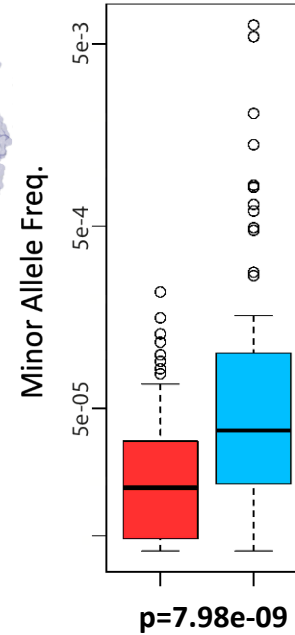
Surface



Interior

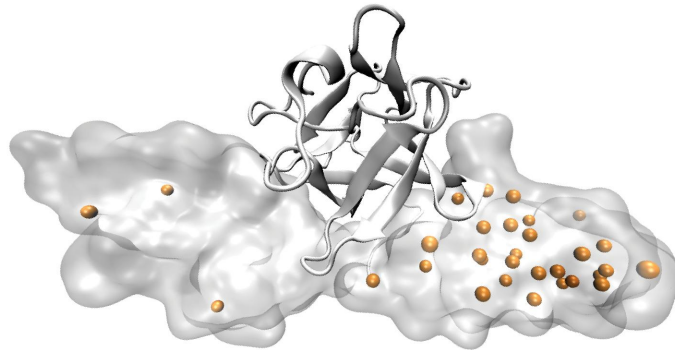


**critical**  
**non-critical**



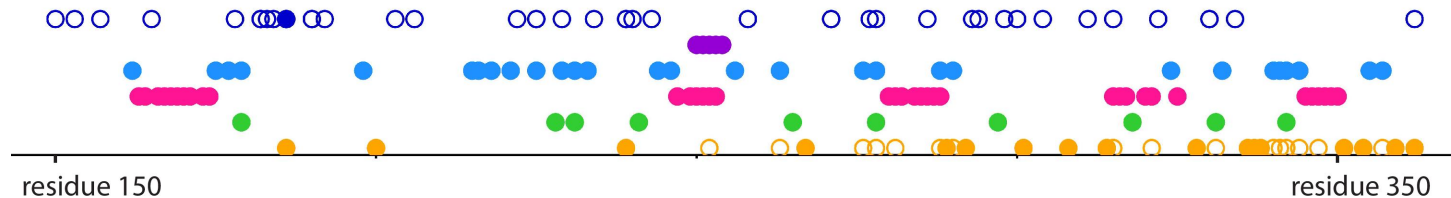
Unlike common SNVs, the statistical power with which we can evaluate rare SNVs in case-control studies is severely limited

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)*

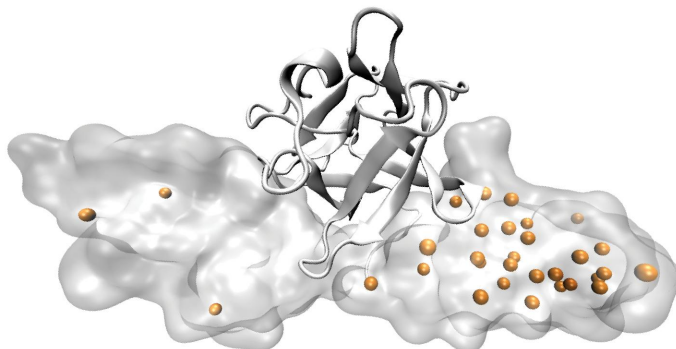
- 1000G & ExAC SNVs (common | rare)
- Hinge residues
- Buried residues
- Protein-protein interaction site
- Post-translational modifications
- HGMD site (w/o annotation overlap)
- HGMD site (w/annotation overlap)



[Sethi et al. COSB ('15)]

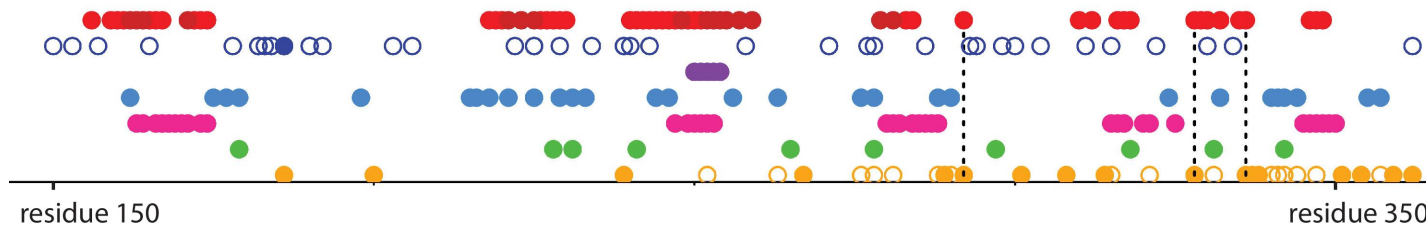
# Protein structures may provide the needed alternative for evaluating rare SNVs, many of which may be disease-associated

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



- 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)

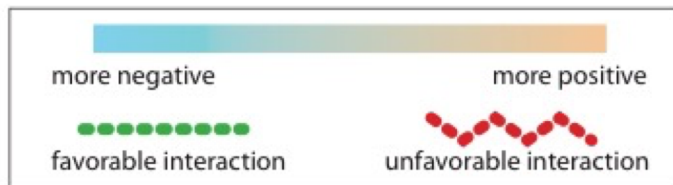
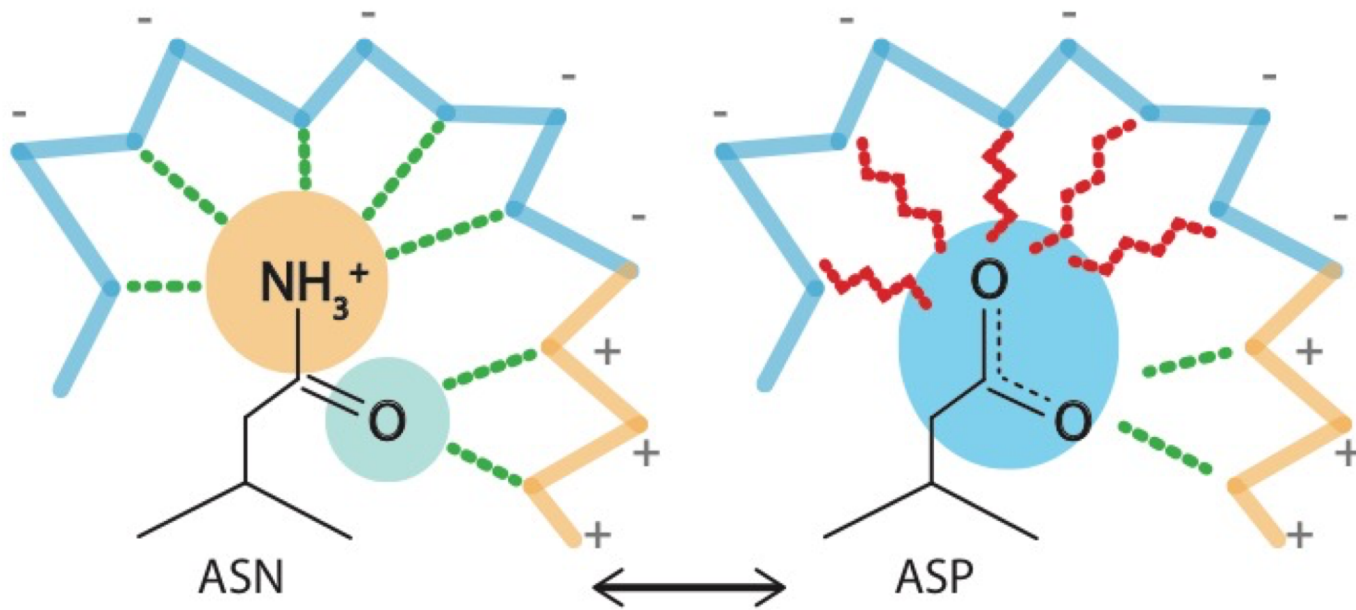
*Fibroblast growth factor receptor 2 (pdb: 1IIL)*



[Sethi et al. COSB ('15)]

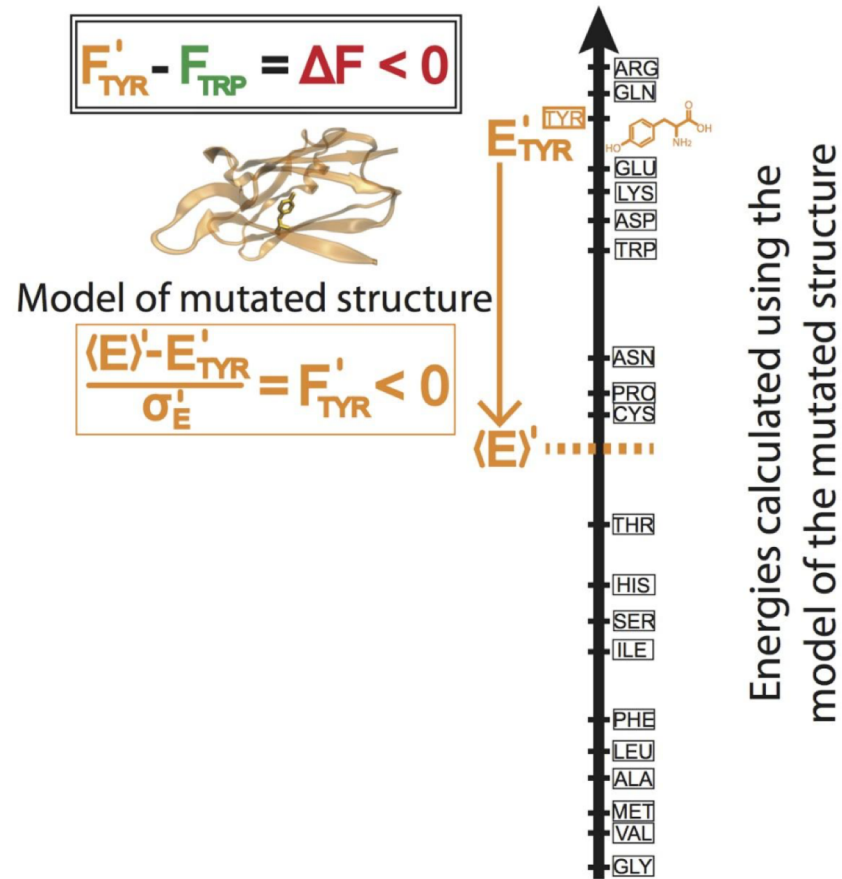
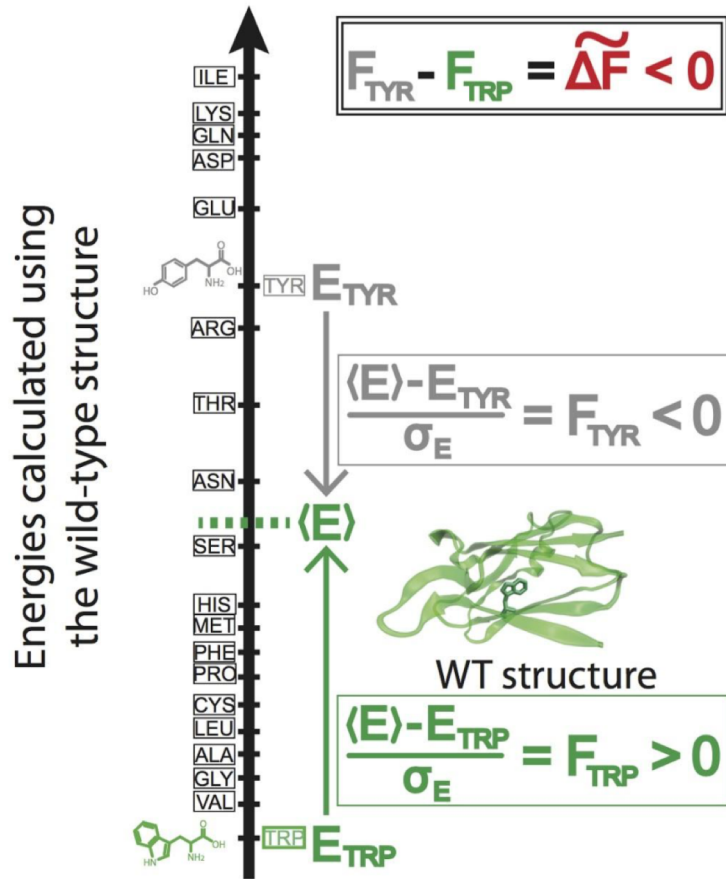
# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with STRESS**
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (Hotcommics)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score

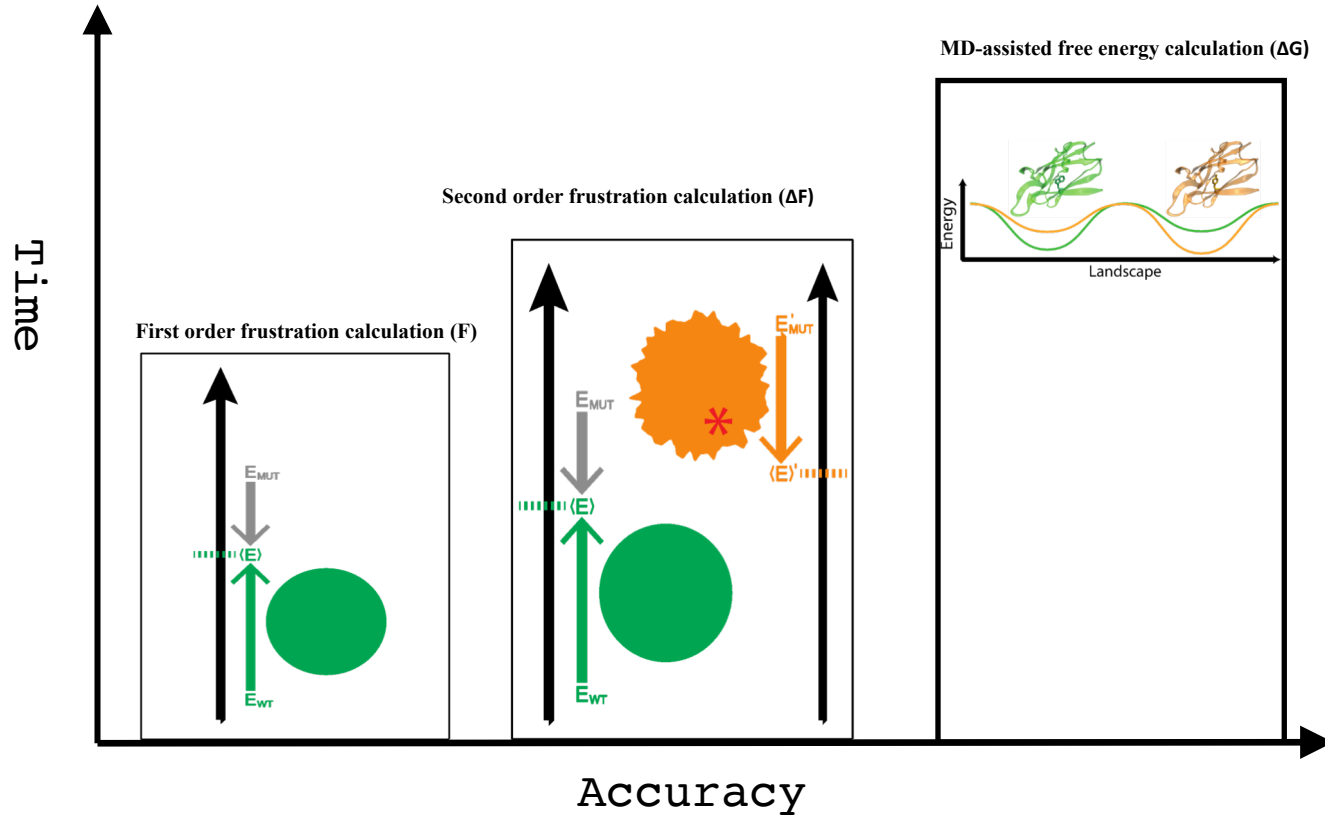


What is  
localized  
frustration  
?

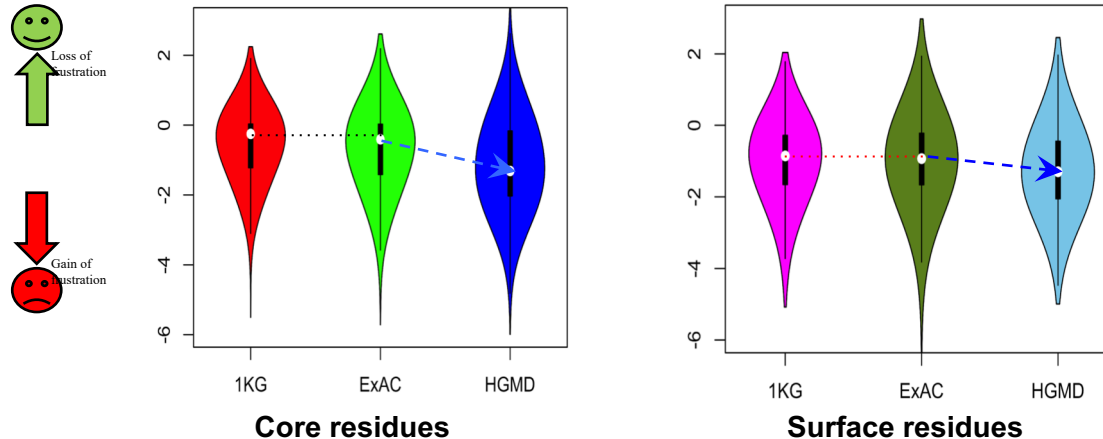
# Workflow for evaluating localized frustration changes ( $\Delta F$ )



# Complexity of the second order frustration calculation



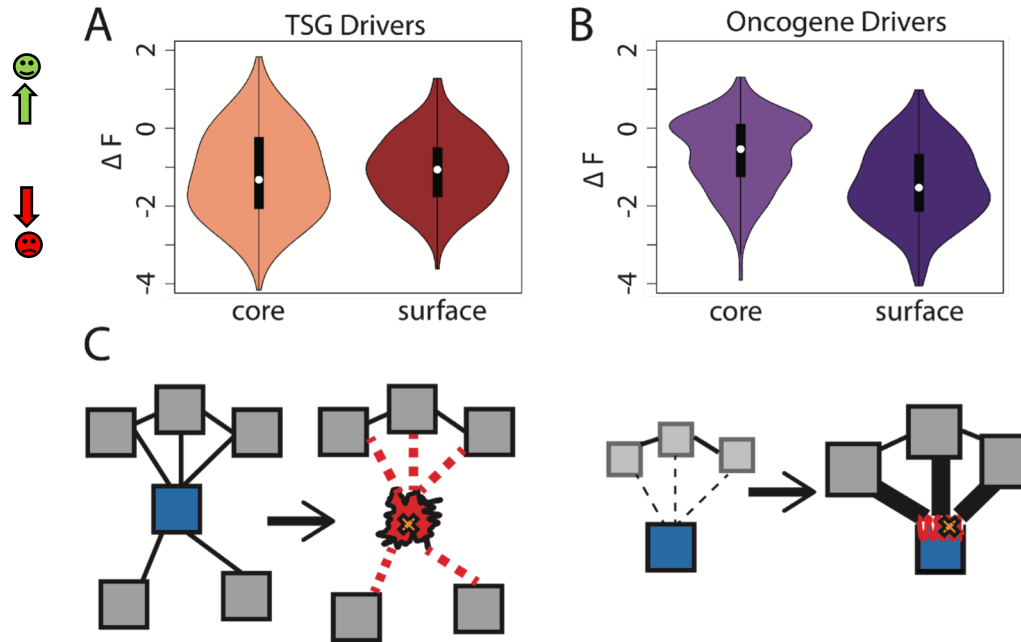
# Comparing $\Delta F$ values across different SNV categories: disease v normal



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



# Comparison between $\Delta F$ 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.

# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with STRESS**
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (Hotcommics)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score

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

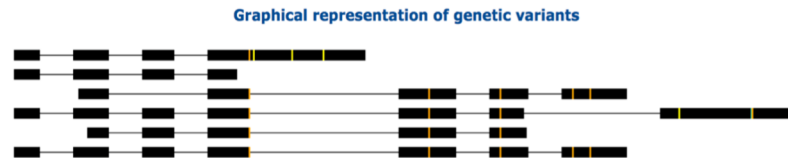
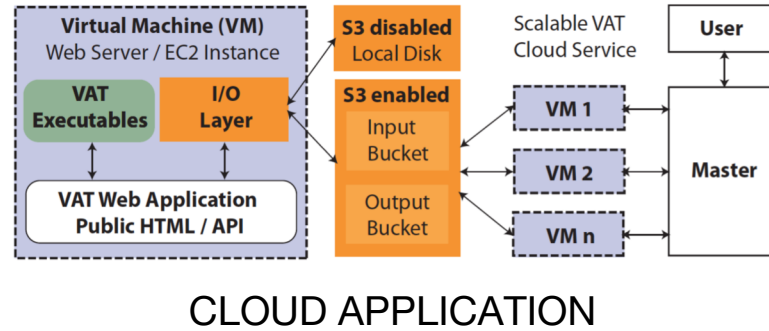
VCF Input

Output:

- Annotated VCFs
- Graphical representations of functional impact on transcripts

Access:

- Webserver
- AWS cloud instance
- Source freely available



[vat.gersteinlab.org](http://vat.gersteinlab.org)

Habegger L. \*, Balasubramanian S. \*, et al. *Bioinformatics*, 2012

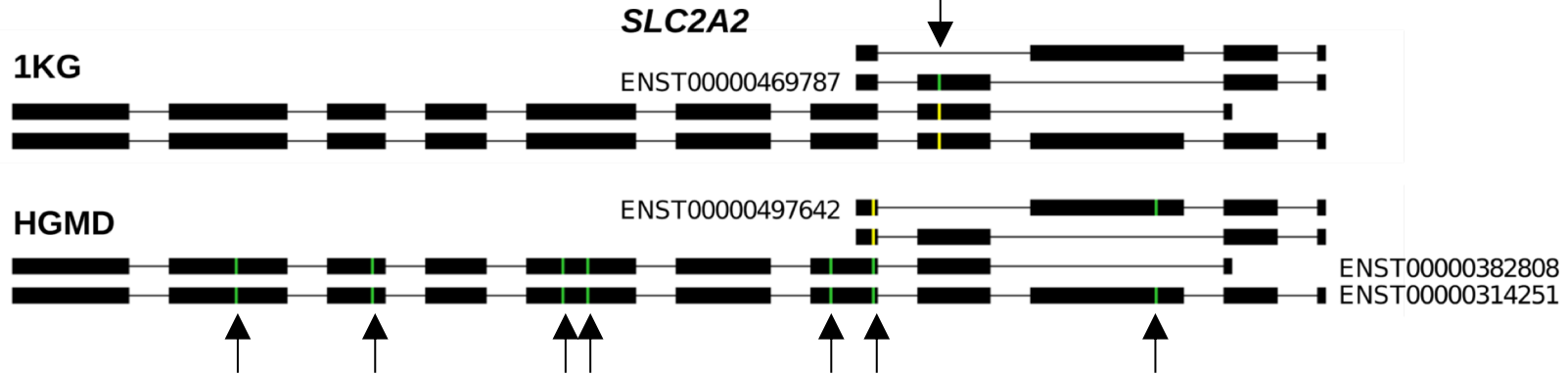
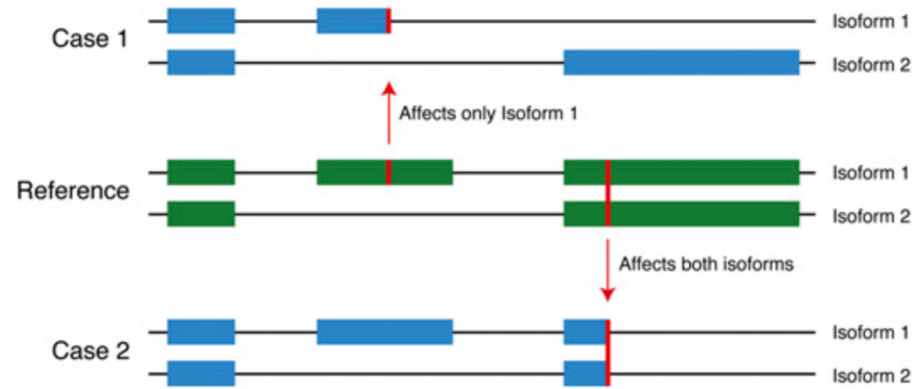
# Complexities in LOF annotation

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



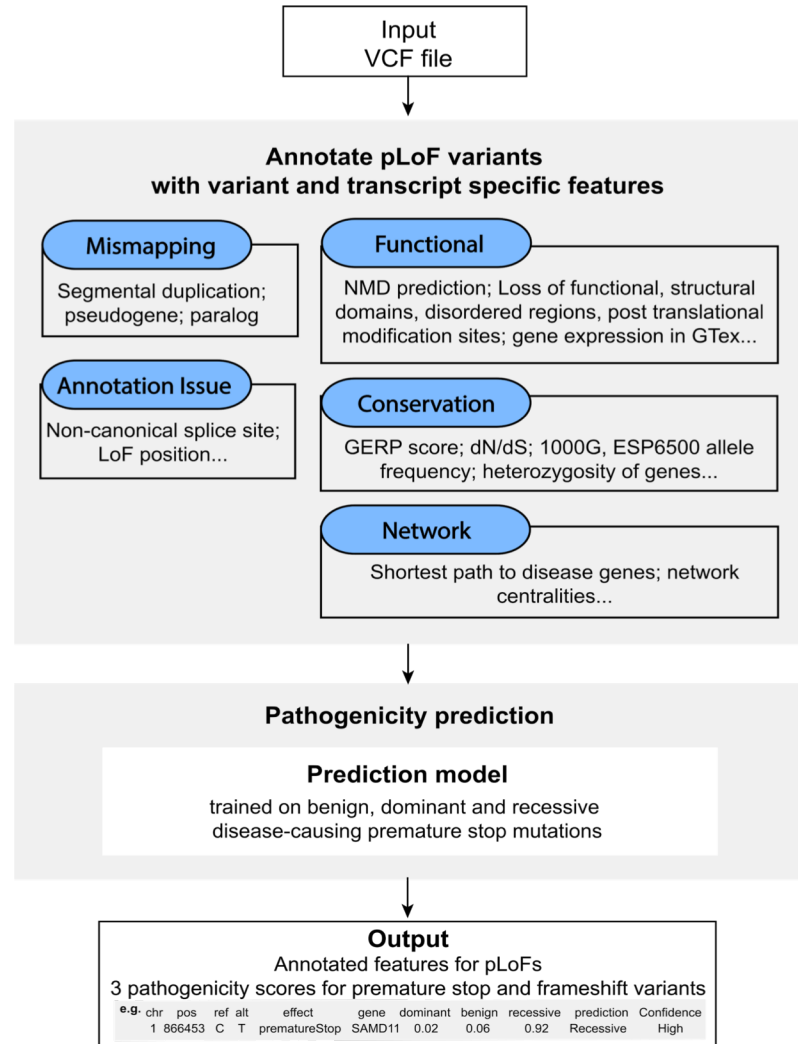
# Annotation of Loss-of-Function Transcripts (ALoFT)

Runs on top of VAT

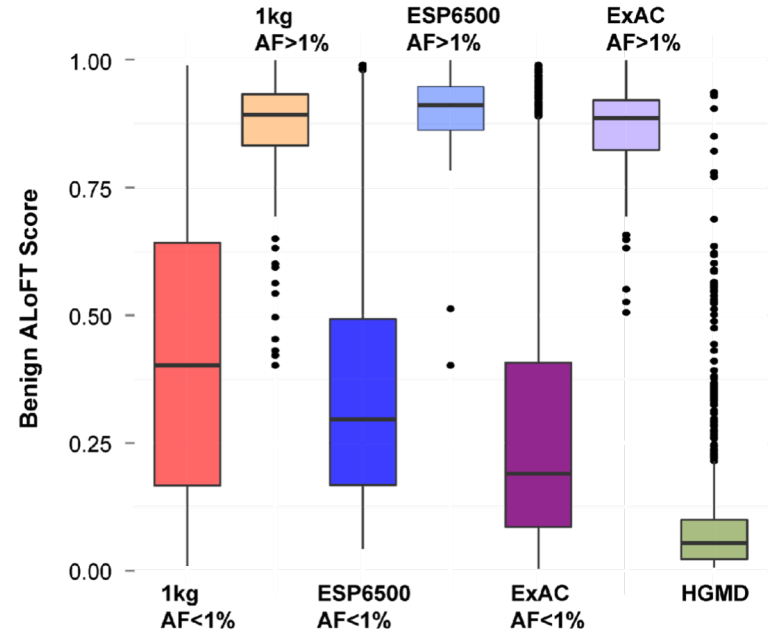
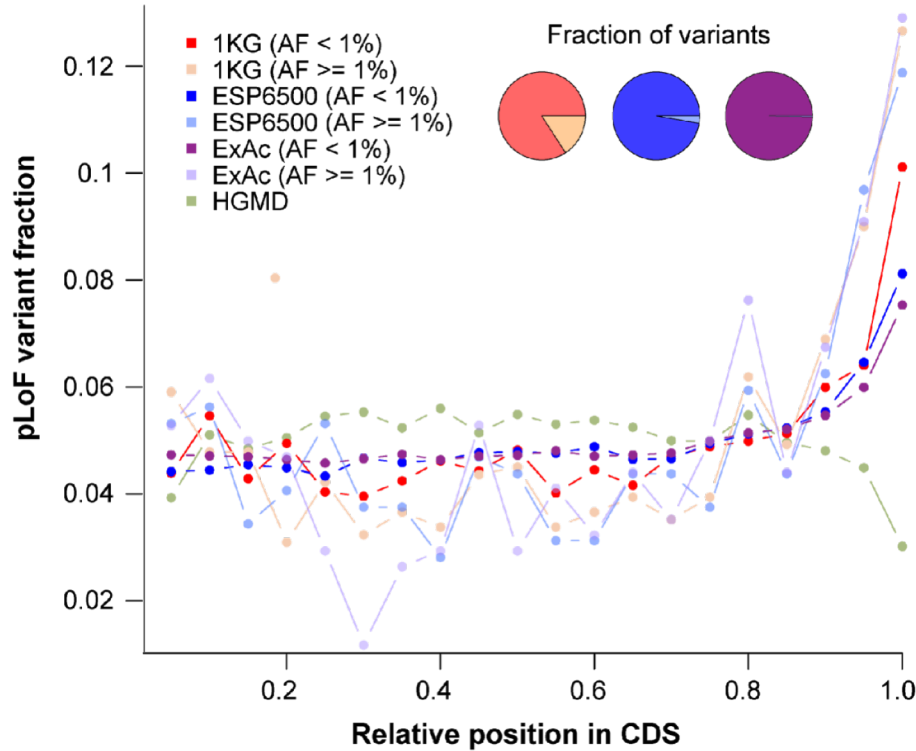
Output:

- 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

# ALoFT identifies deleterious somatic LoF variants

## Cancer genes:

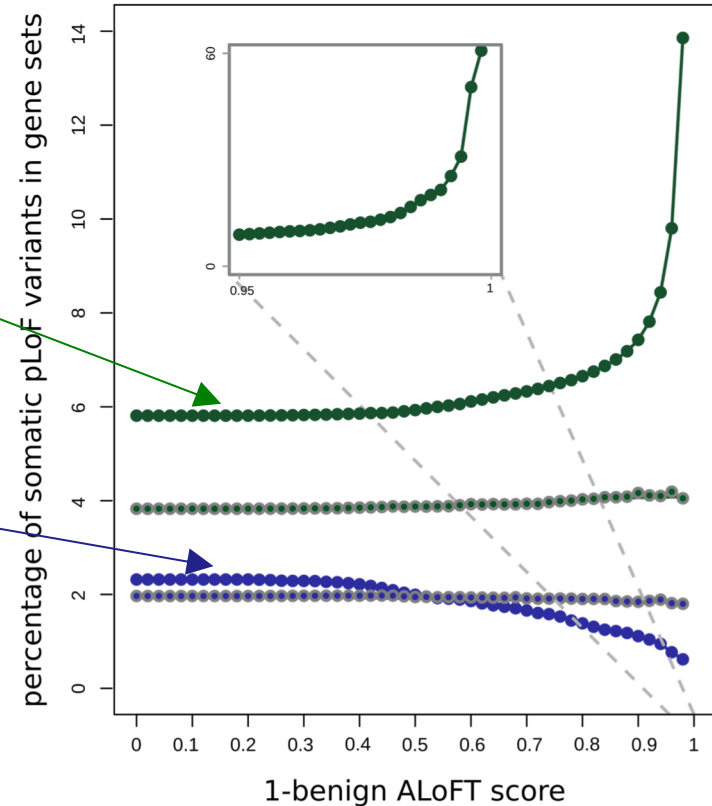
- COSMIC consensus.
- *Enriched in deleterious LoFs.*

## LoF tolerant genes:

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

## cancer genes vs. LoF tolerant genes

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

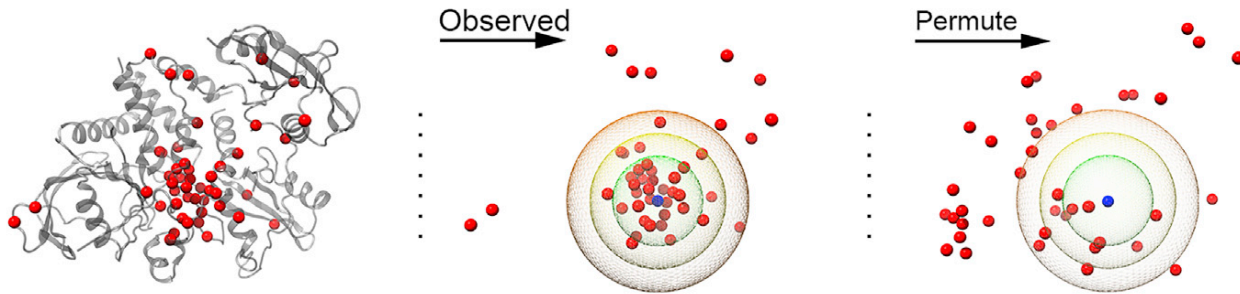


# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with **STRESS****
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (**Hotcommics**)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score



# Structures have been used successfully to “aggregate” the burden of mutations

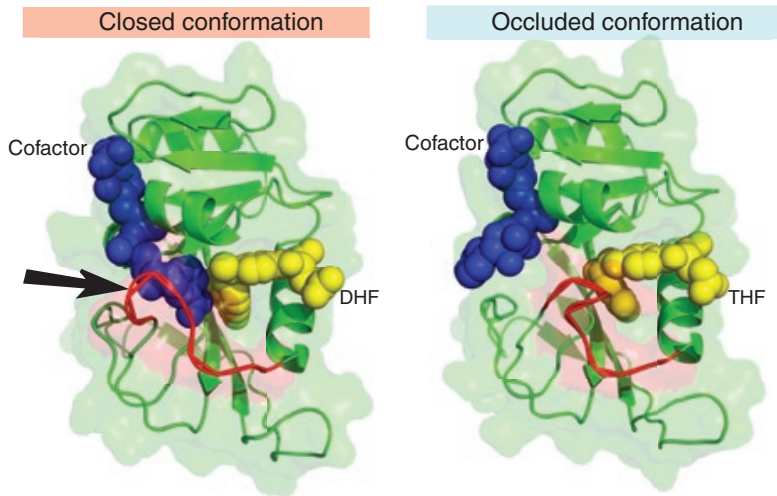


These approaches search for mutational clusters on protein structure using distance cutoff.

Permutation is performed to identify statistically significant mutational clusters on static protein structure.

Both rare germline & somatic

# Protein dynamics is important for protein function



Proteins are inherently dynamic bio-molecules and sample large ensembles of conformations.

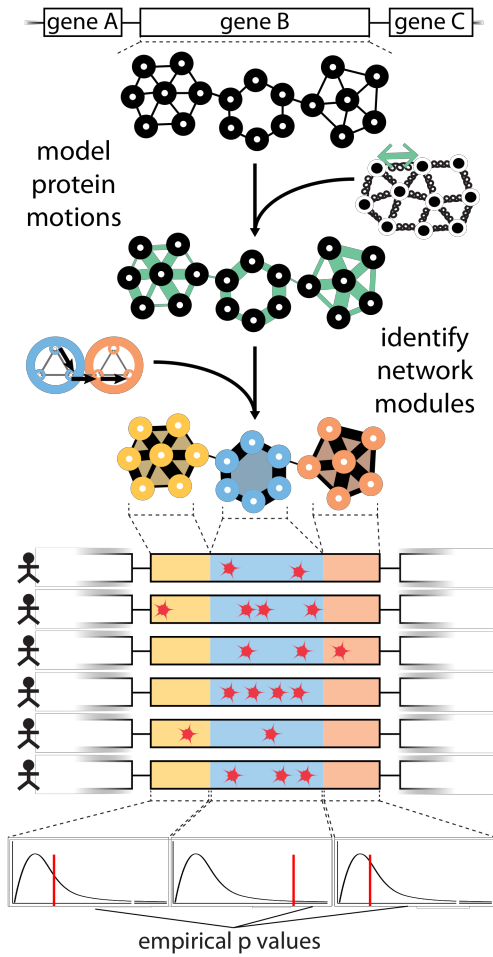
prior structure-based methods are potentially **less sensitive to identify functional residues** through the mutation clustering approach.

Potentially **miss many critical mutational clusters**

We leverage protein dynamics to identify mutation clusters

Focus on data from from TCGA pancan atlas data

# Workflow to aggregate mutations taking into account dynamics

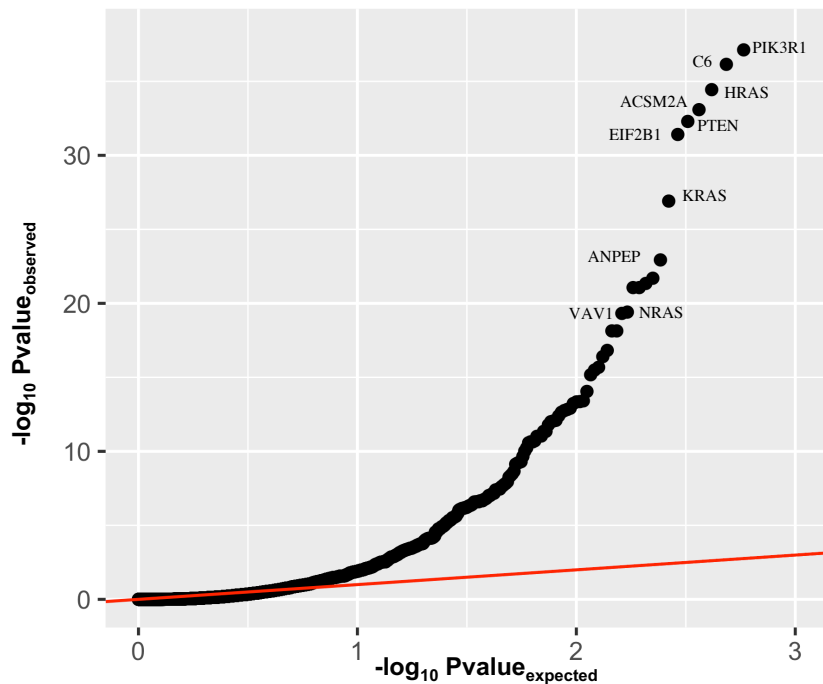


Our framework leverages large-scale conformational changes of a protein to identify dynamic sub-regions of proteins (or “communities”).

we mapped missense mutations onto three-dimensional protein structures.

For each community with mapped mutations, we performed a Fisher exact test to determine whether variants fall within a given community is more frequently observed than what would be expected by chance.

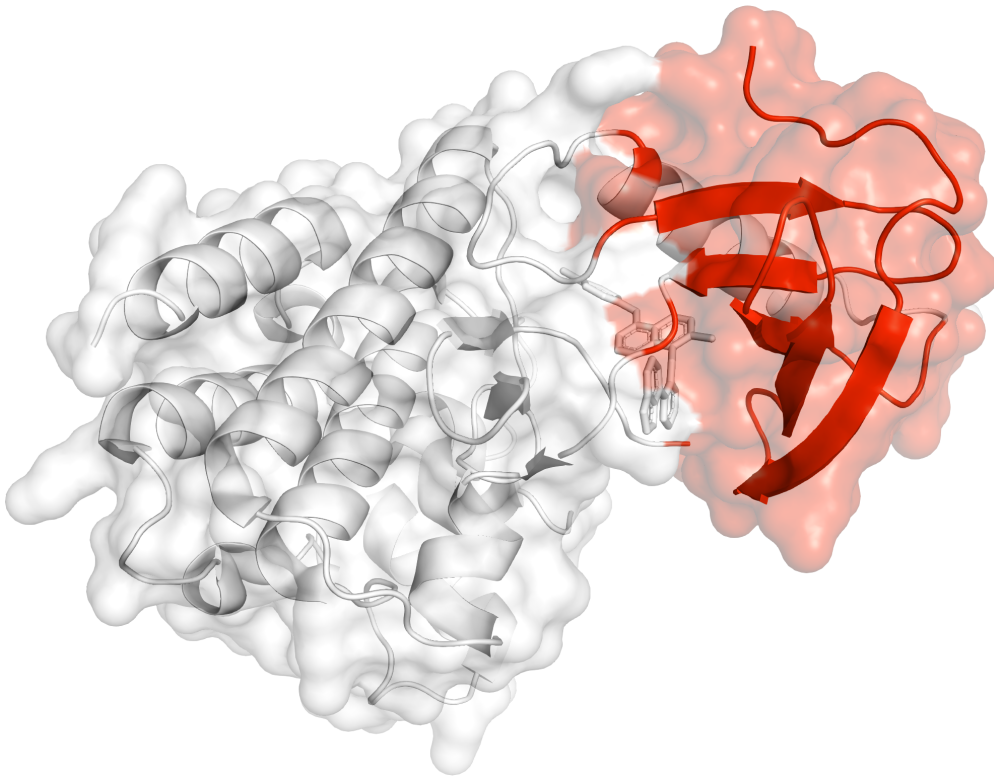
# Pancancer Q-Q plot for genes with hotspot communities



Our pan-cancer analysis identifies hotspot communities present on protein structures of 434 putative driver genes.

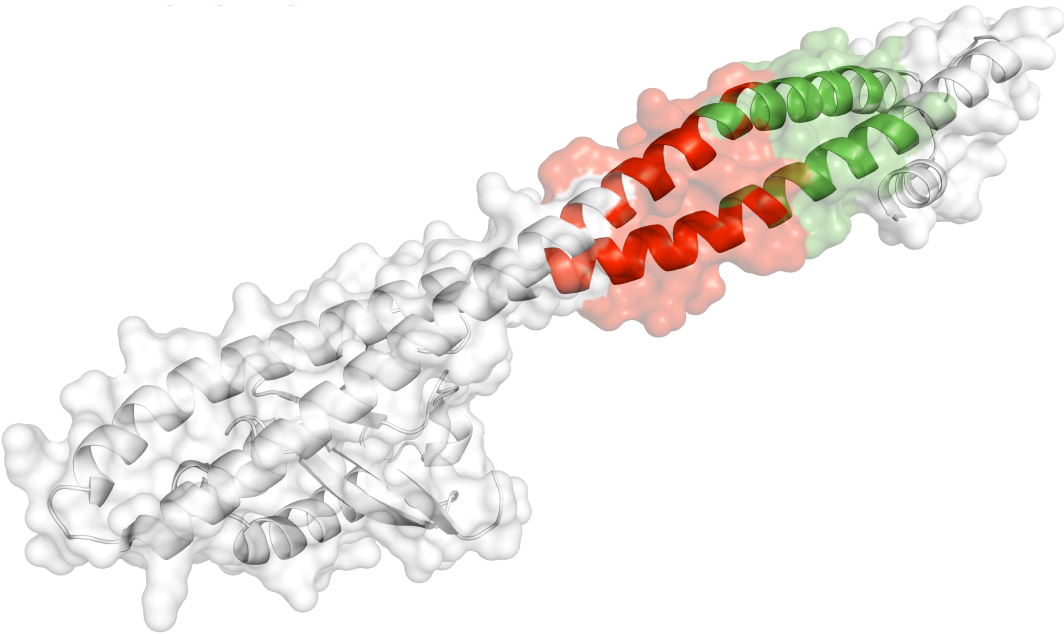
Our workflow identifies well known driver genes as well as novel putative driver genes.

# Example of oncogene with hotspot communities



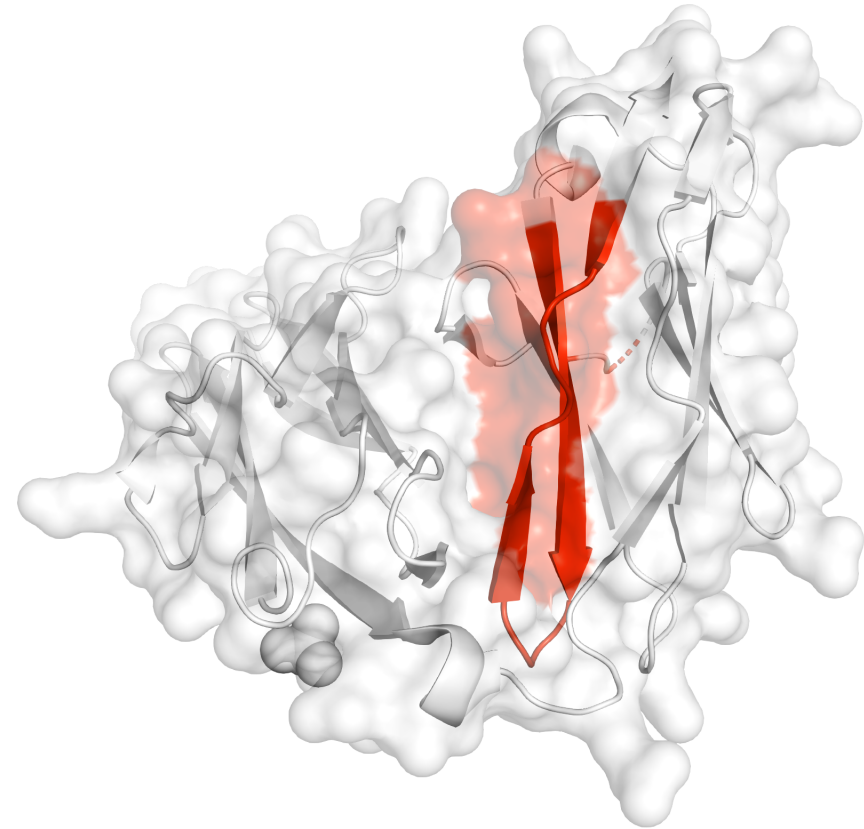
BRAF: We identify **one hotspot community comprising of 52 residues** on the co-crystal structure of the BRAfV600E kinase domain .

# Example of TSG with hotspot communities



We identify **two hotspot communities** adjacent to each other on the co-crystal structure of the PIK3R1 gene.

# Example of novel drivers with hotspot communities



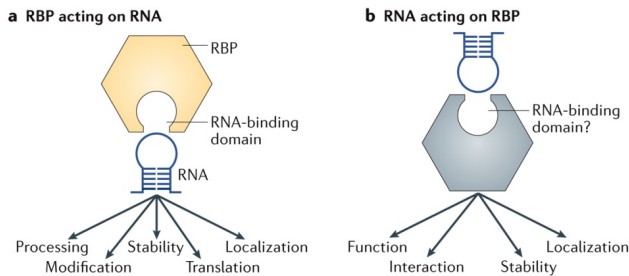
Our workflow predicts **one hotspot community** that comprise of 47 residues in the crystal structure of PTPRD gene.

# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with STRESS**
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (Hotcommics)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score



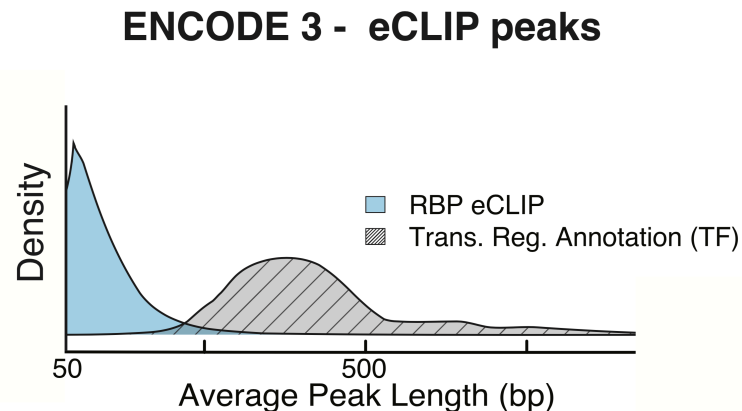
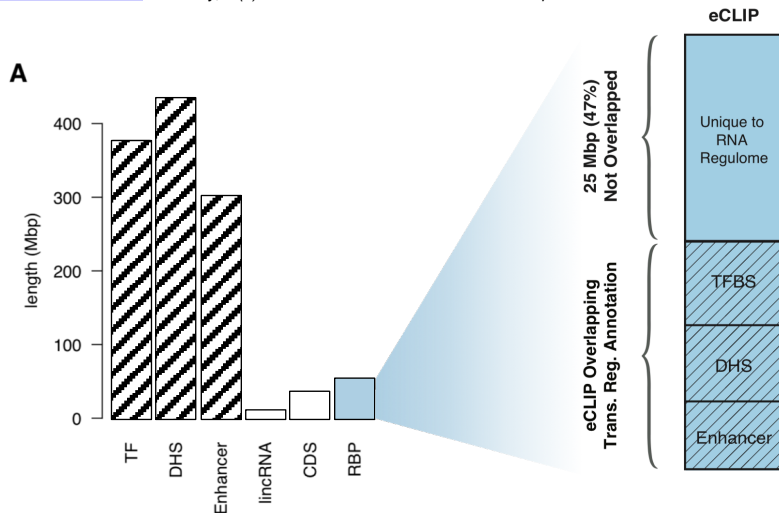
# RNA Binding Proteins (RBPs)



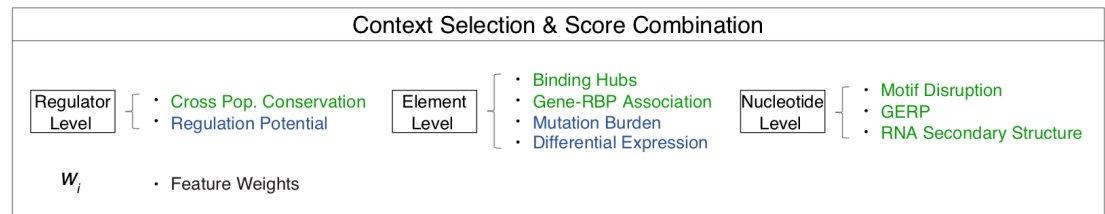
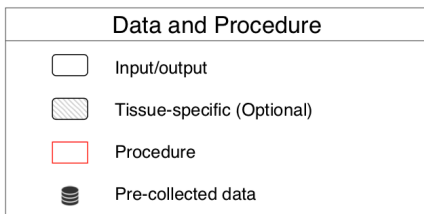
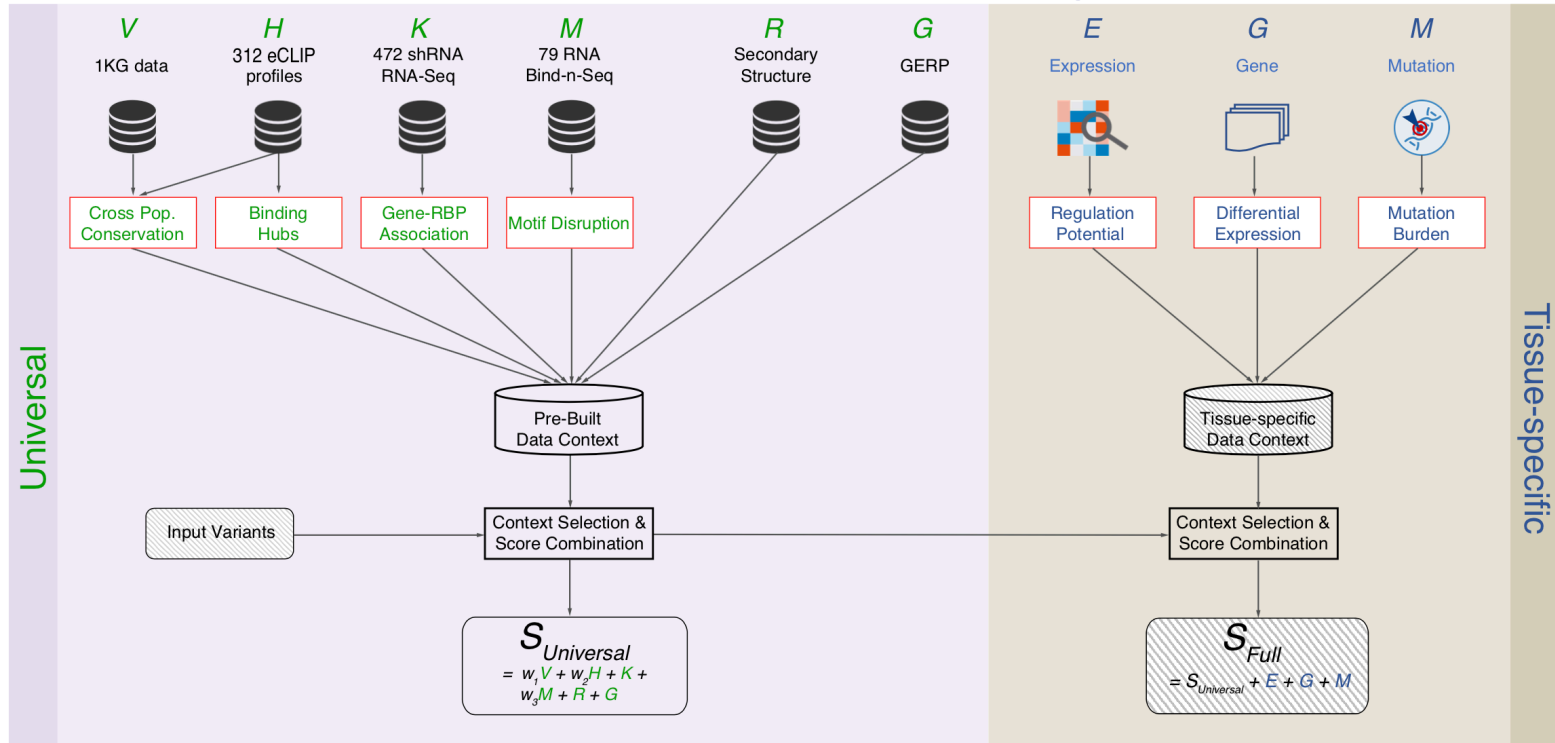
Nature Reviews | Molecular Cell Biology

[Nat Rev Mol Cell Biol.](https://doi.org/10.1038/nrm.2017.130) 2018 May;19(5):327-341. doi: 10.1038/nrm.2017.130. Epub 2018 Jan 17.

- **Before ENCODE3: >150 expt.** in many different cell types
- **ENCODE3 did ~350 focused eCLIP expt.** for >110 RBPs on HepG2 & K562 (Van Nostrand...Yeo. Nat. Meth. '16; Van Nostrand...Graveley, Yeo (submitted in relation to ENCODE3))

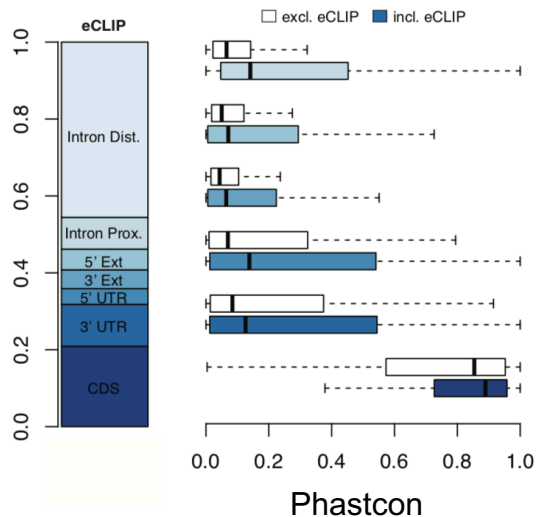


# Schematic of RADAR Scoring

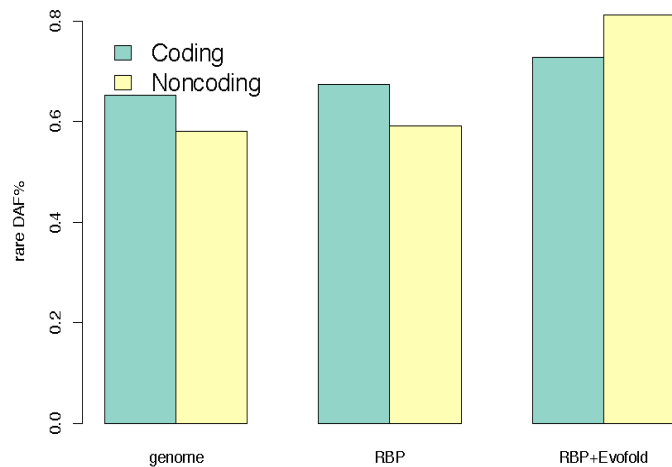




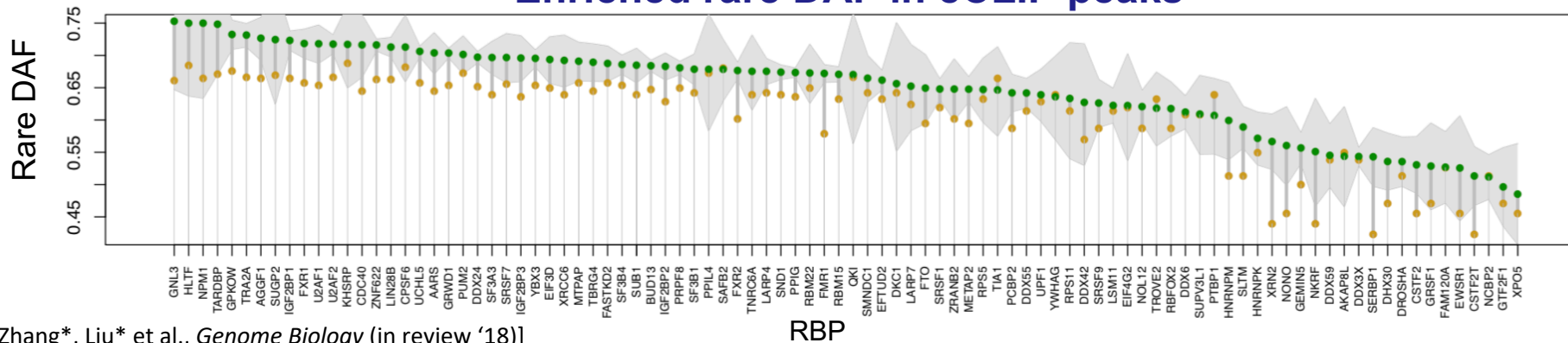
# High Phastcon in RBP-overlapped annotations



# RNA Structure Cons. from EvoFold



# Enriched rare DAF in eCLIP peaks

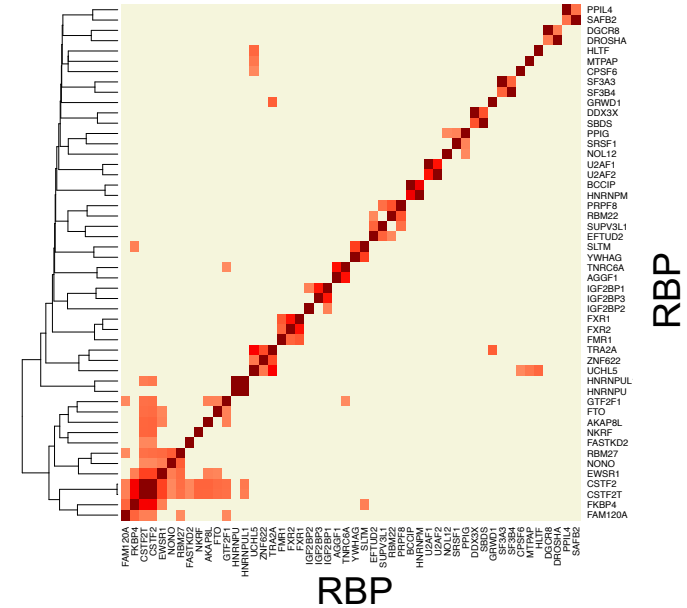
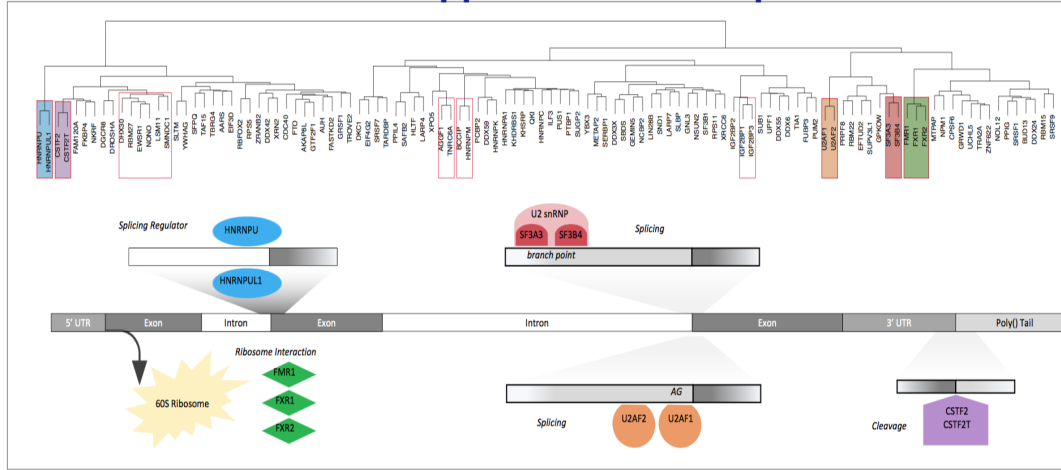


[Zhang\*, Liu\* et al., *Genome Biology* (in review '18)]

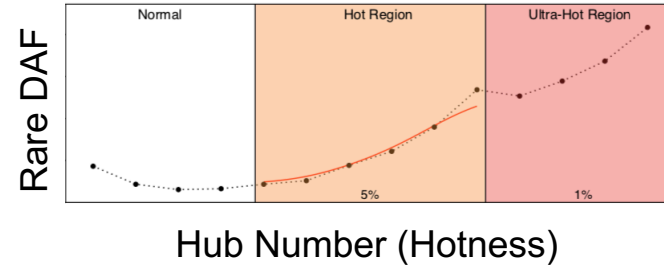
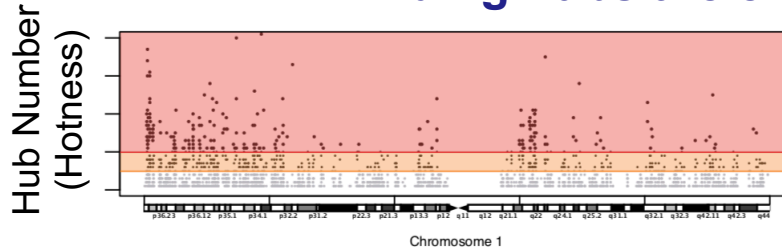
# Co-binding of RBPs form biologically relevant complexes

## Unique co-binding patterns of RBPs

### Literature supported RBP complexes



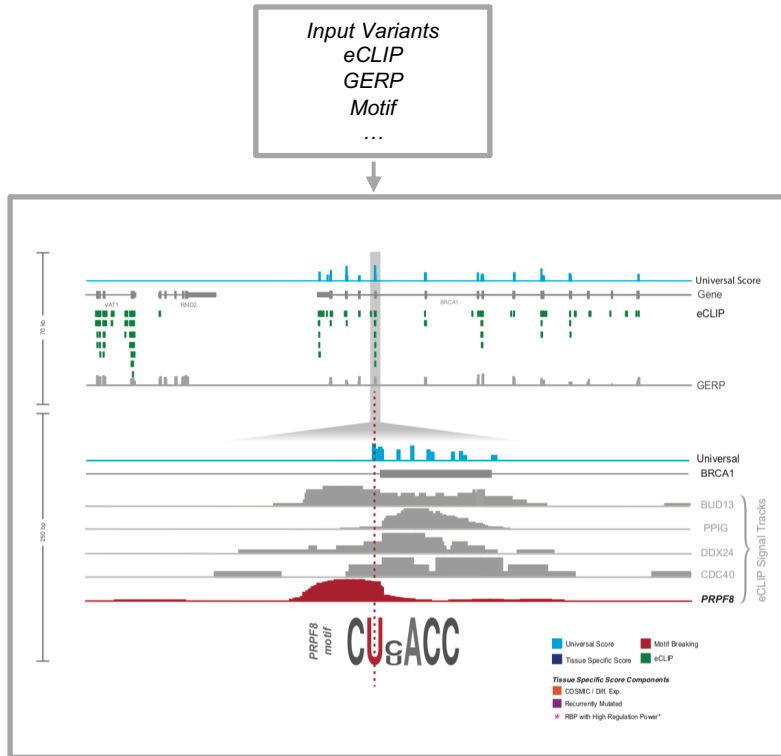
## Binding hubs are enriched for rare variants



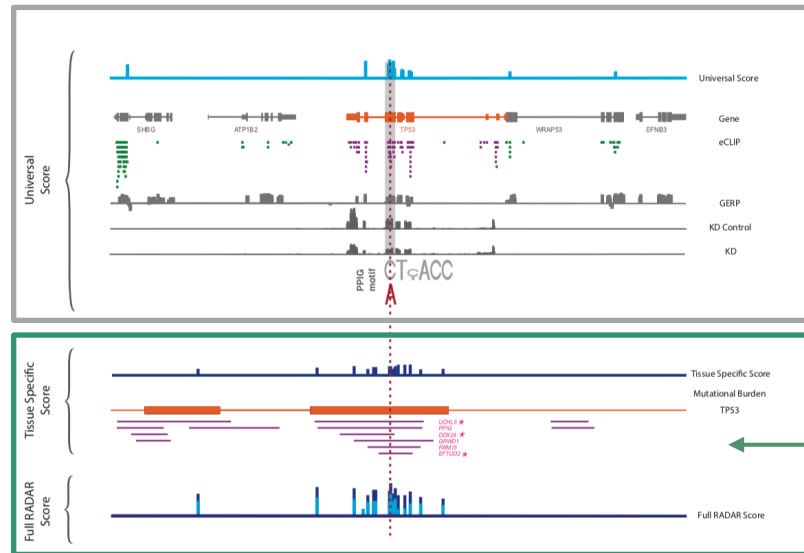
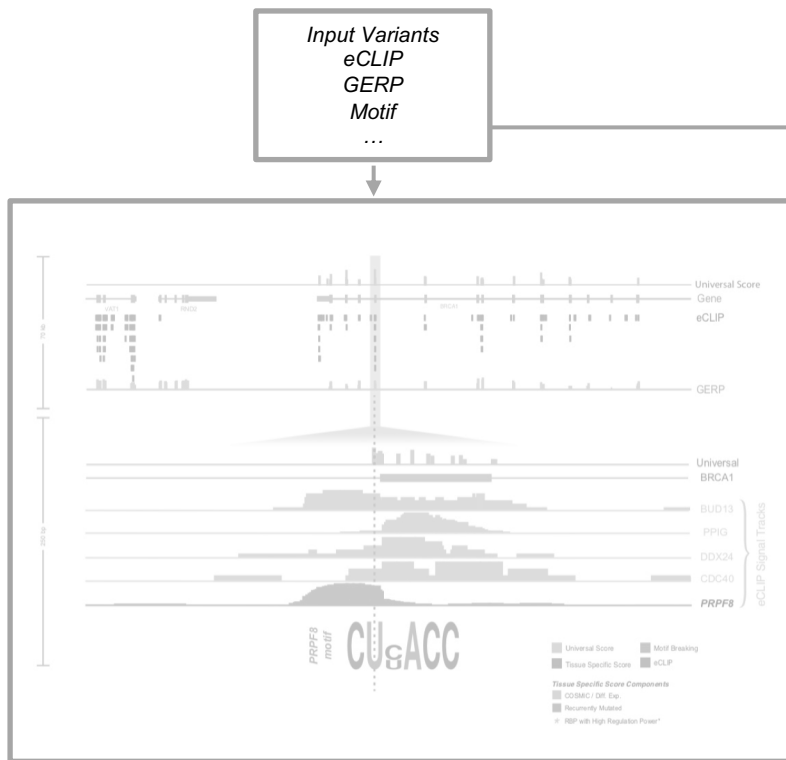


# Visualization of RADAR Features and Scoring

## Germline Variants are Score Using a Universal Scoring Scheme



# Visualization of RADAR Features and Scoring



**Tissue Specific:**  
Variants  
Expression  
Regulatory Potential

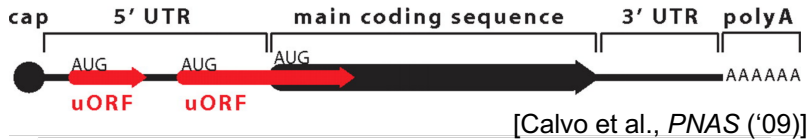
**Somatic Variant Scored with Universal + Tissue specific context score**



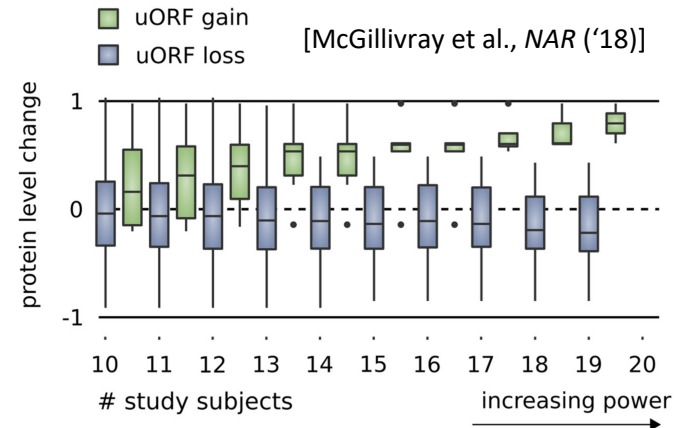
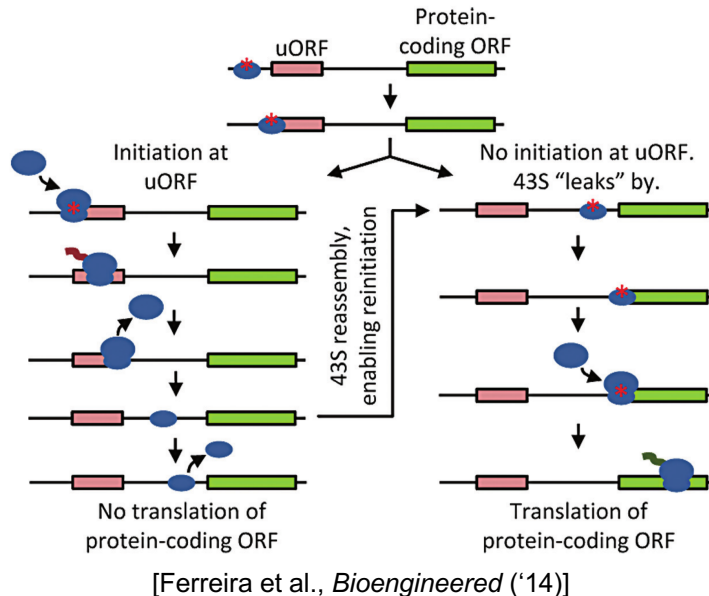
# Computational analysis of variants: coding versus non-coding

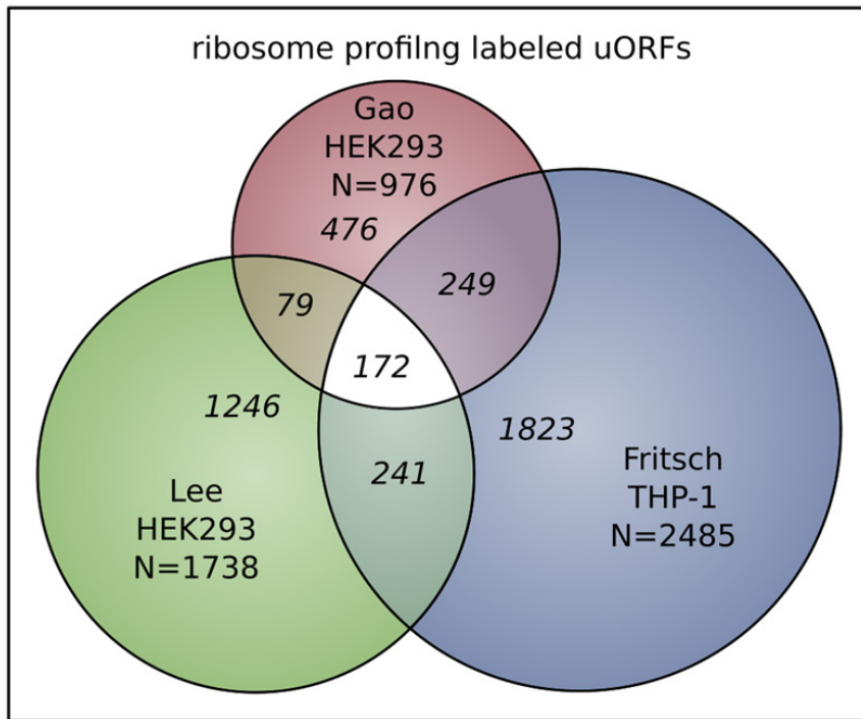
- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with STRESS**
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (Hotcommics)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score

# Upstream open reading frames (uORFs) regulate translation are affected by mutation



- uORFs regulate the translation of downstream coding regions.
- In Battle et al. 2014 data uORF gain & loss assoc. protein level change.

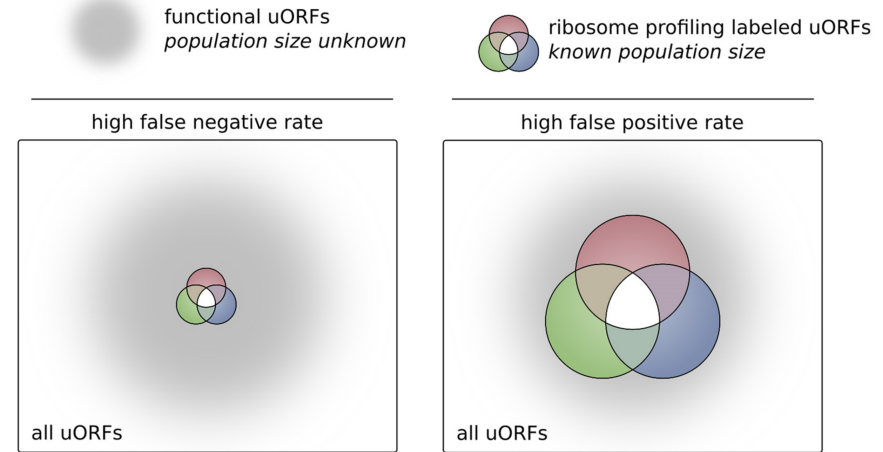




**From a “Universe” of  
1.3 M pot. uORFs**

## The population of functional uORFs may be significant

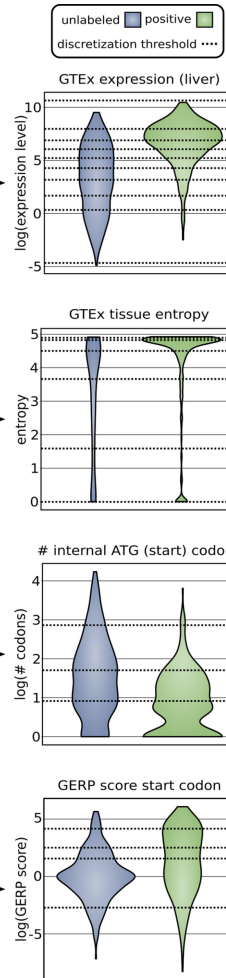
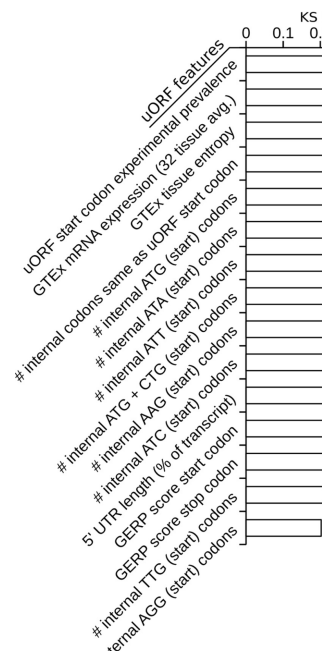
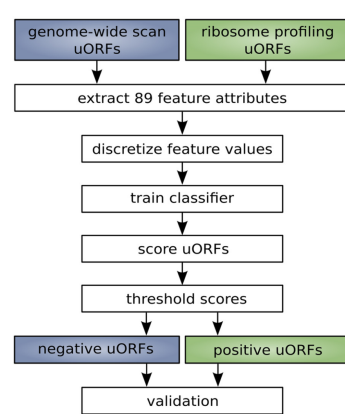
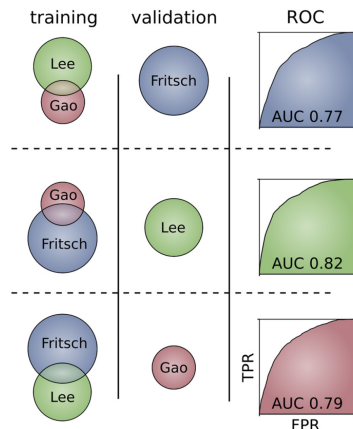
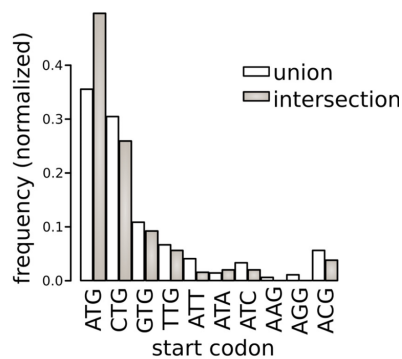
**C**



- Ribosome profiling experiments have low overlap in identified uORFs.
- This suggests high false-negative rate, and more functional uORFs than currently known.

# Prediction & validation of functional uORFs using 89 features

- All near-cognate start codons predicted.
- Cross-validation on independent ribosome profiling datasets and validation using in vivo protein levels and ribosome occupancy in humans (Battle et al. 2014).



Expr. Level

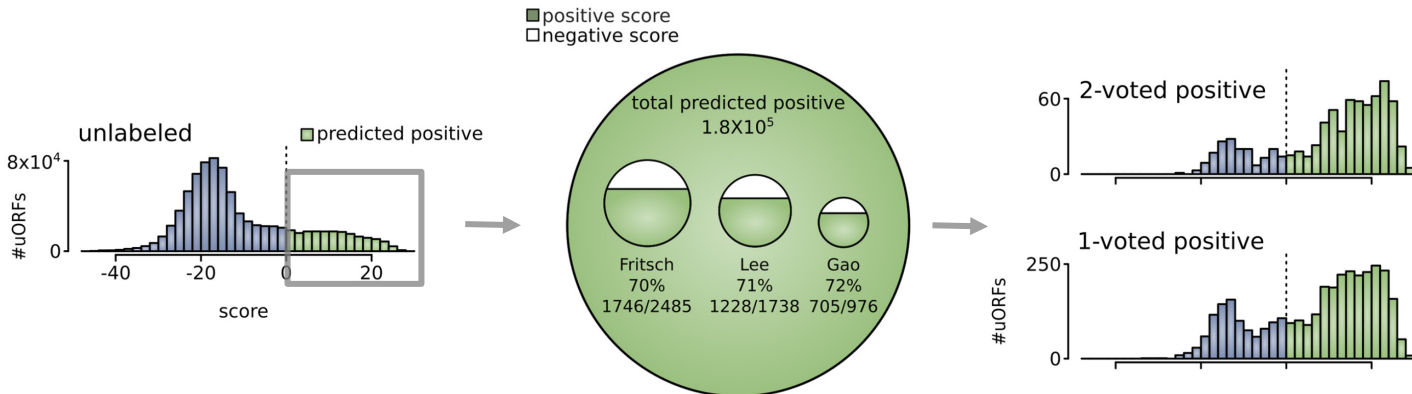
Tissue Dist.

Int. ATG Start

Conser-  
vation

# A comprehensive catalog of functional uORFs

Universe of **1.3M**  
uORFs scored via  
Simple Bayes algo.



- Predicted functional uORFs may be intersected with disease associated variants.

- **180K**: Large predicted positive set likely to affect translation
- Calibration on gold standards, suggests getting **~70%** of known

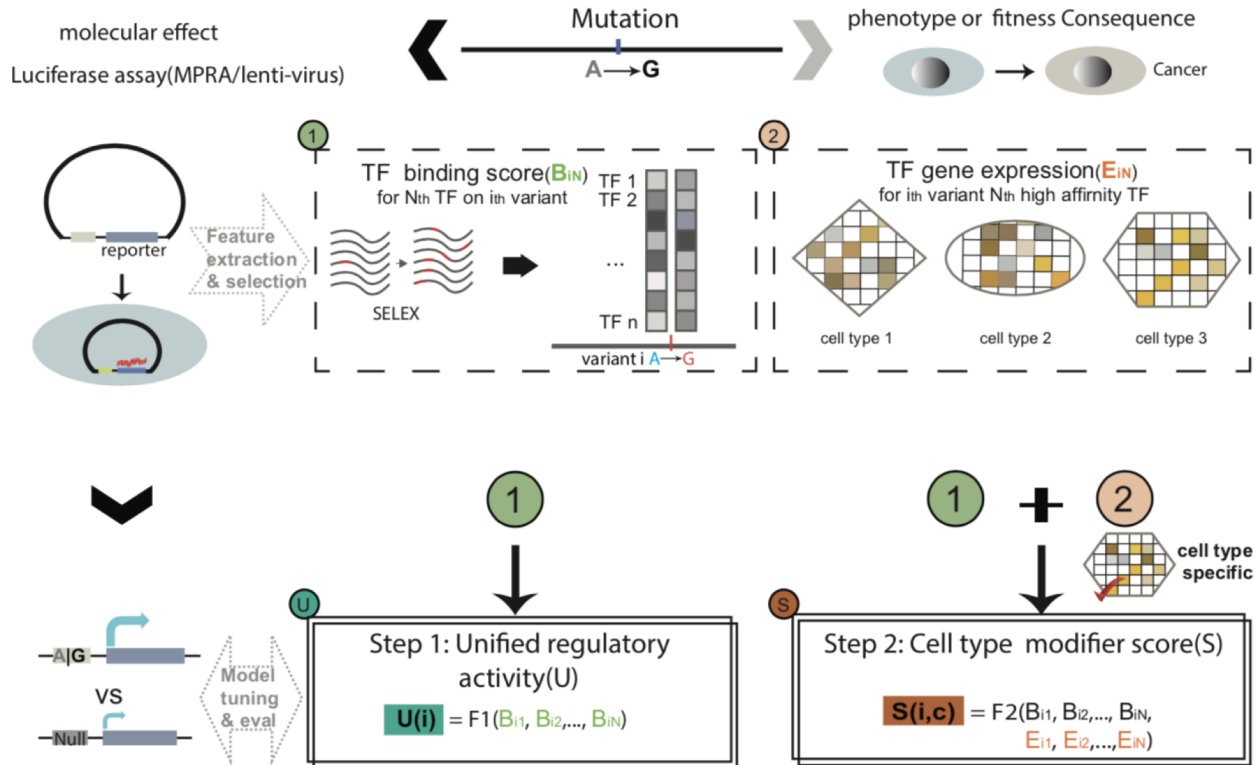
# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with STRESS**
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (Hotcommics)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score

# Promotor Mutations

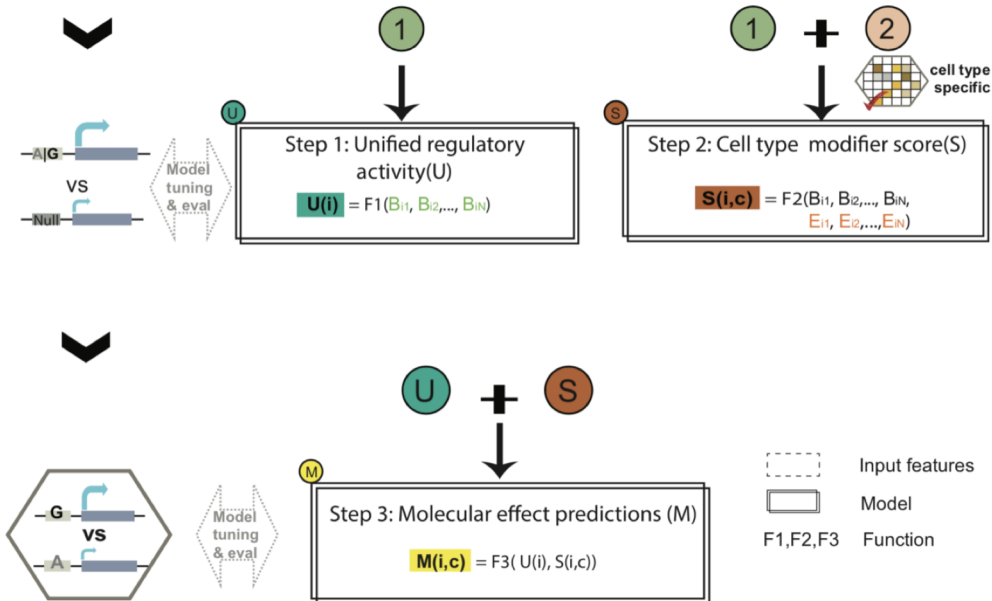
- How do we assess their effect?
  - What's the readout? Expression?
- Molecular v Organismic phenotype
- Importance of specific cellular context

# GRAM approach for assessing molecular impact

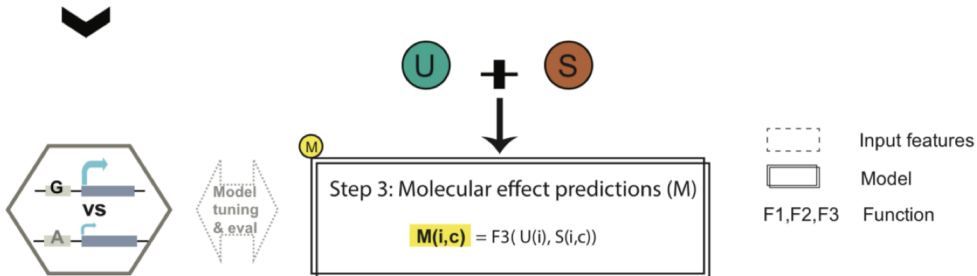
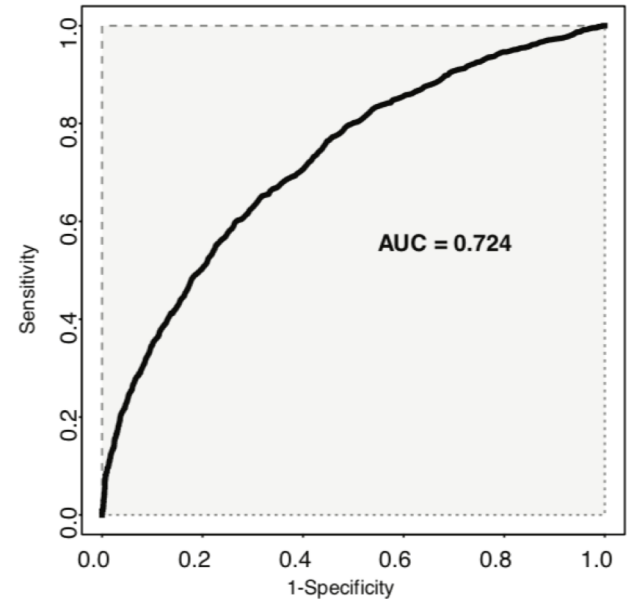




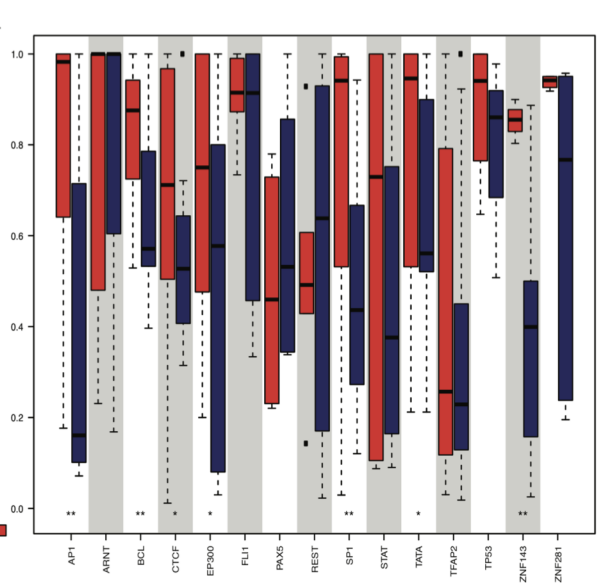
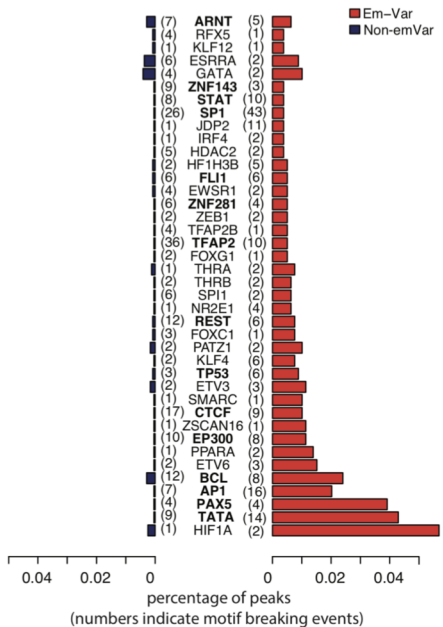
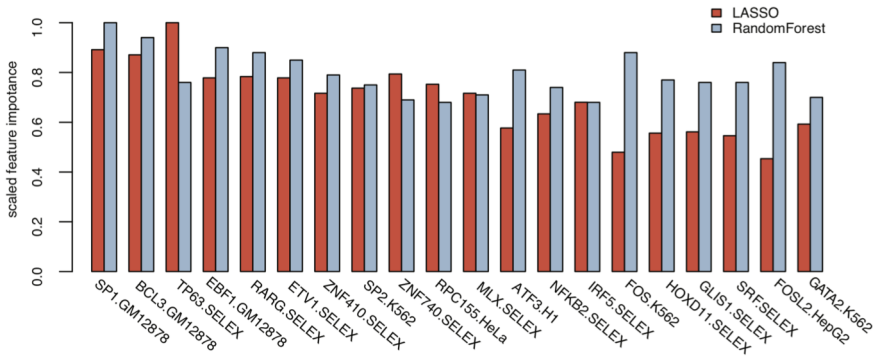
# GRAM performance



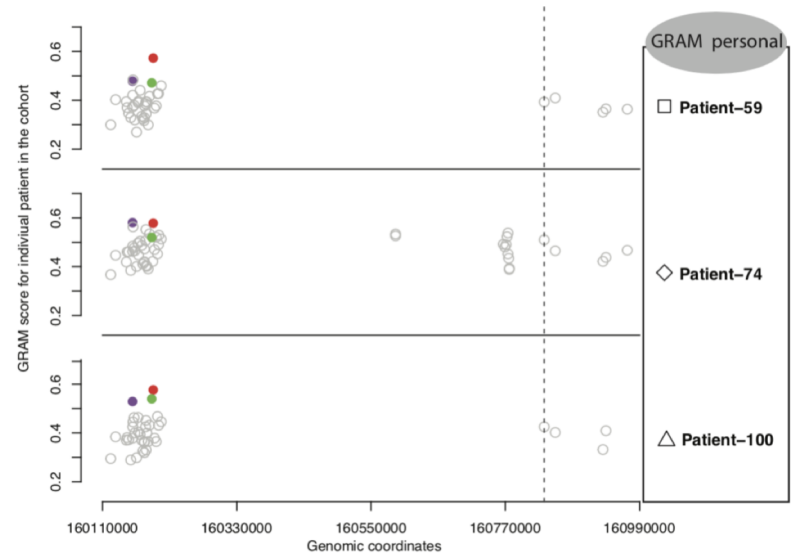
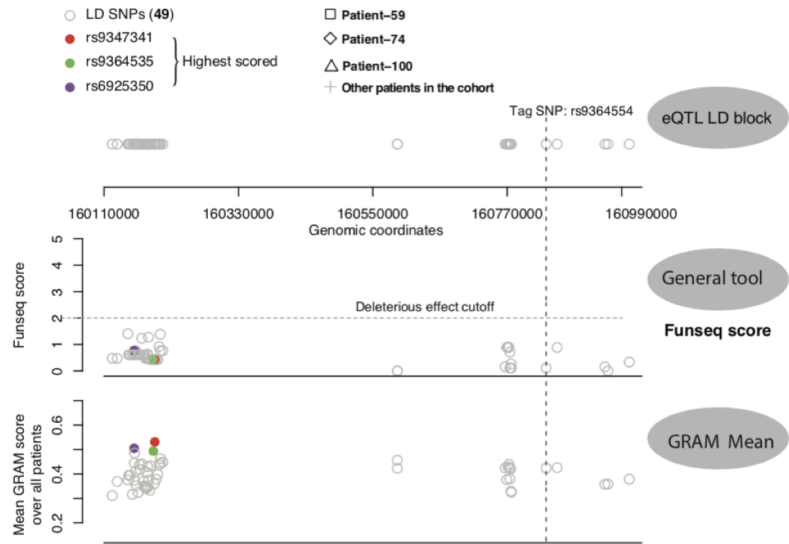
Model trained by GM12878 ChIP-Seq + SELEX dataset in step3



# GRAM Important Features



# GRAM application to find molecular cause within group of eQTL variants



# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**
  - Rare v common, somatic v germline, coding v noncoding
- **Identifying cryptic allosteric sites with STRESS**
  - On surface & in interior bottlenecks
- **Frustration as a localized metric of SNV impact**
  - Differential profiles for oncogenes v. TSGs
- **ALoFT: Annotation of LoF Transcripts**
- **Using dynamics to help identify mutation clusters (Hotcommics)**
  - Find dynamic sub-communities & determine aggregated mutational burden within these
- **RADAR Prioritization for RBP sites**
  - Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
  - Incorporates new features related to RNA sec. struc & tissue specific effects
- **uORF Prioritization**
  - Feature integration to find small subset of upstream mutations that potentially alter translation
- **GRAM to assess the molecular effect of (promotor) mutations**
  - Universal score + cell type specific score

# Computational analysis of variants: coding versus non-coding

- **Intro: types of variants**

- Rare v common,  
somatic v germline, coding v noncoding

- **Identifying cryptic allosteric sites with STRESS**

- On surface & in interior bottlenecks

- **Frustration as a localized metric of SNV impact**

- Differential profiles for oncogenes v. TSGs

- **ALoFT: Annotation of LoF Transcripts**

- **Using dynamics to help identify mutation clusters (Hotcommics)**

- Find dynamic sub-communities & determine aggregated mutational burden within these

- **RADAR Prioritization for RBP sites**

- Prioritizes variants based on post-transcriptional regulome using ENCODE eCLIP
- Incorporates new features related to RNA sec. struc & tissue specific effects

- **uORF Prioritization**

- Feature integration to find small subset of upstream mutations that potentially alter translation

- **GRAM to assess the molecular effect of (promotor) mutations**

- Universal score + cell type specific score

**RADAR**.gersteinlab.org

J **Zhang**, J **Liu**, D Lee, J-J Feng,  
L Lochovsky, S Lou, M Rutenberg-Schoenberg

github.gersteinlab.org/**uORFs**

P **McGillivray**, R Ault, M Pawashe,  
R Kitchen, S Balasubramanian

github.com/gersteinlab/**Frustration**

S **Kumar**, D Clarke

**STRESS**.molmovdb.org

D **Clarke**, A **Sethi**, S Li, S Kumar,  
R Chang, J Chen

github.com/gersteinlab/**gram**

S **Lou**, KA Cotter, T Li, J Liang, H Mohsen, J Liu,  
J Zhang, S Cohen, J Xu, H Yu, MA Rubin

github.com/gersteinlab/**hotcommics**

S **Kumar**, D Clarke

**ALoFT**.gersteinlab.org

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





# Info about this talk

## No Conflicts

Unless explicitly listed here. There are no conflicts of interest relevant to the material in this talk

## General PERMISSIONS

- This Presentation is copyright Mark Gerstein, Yale University, 2019.
- Please read permissions statement at  
**[sites.gersteinlab.org/Permissions](https://sites.gersteinlab.org/Permissions)**
- Basically, feel free to use slides & images in the talk with PROPER acknowledgement (via citation to relevant papers or website link). 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 [streams.gerstein.info](https://streams.gerstein.info) . In particular, many of the images have particular EXIF tags, such as `kwpotppt` , that can be easily queried from flickr, viz: [flickr.com/photos/mbgmbg/tags/kwpotppt](https://www.flickr.com/photos/mbgmbg/tags/kwpotppt)