# Analysis of Personal Genomes:
# Evaluating the impact of variants in exomes using protein structure & allelic activity

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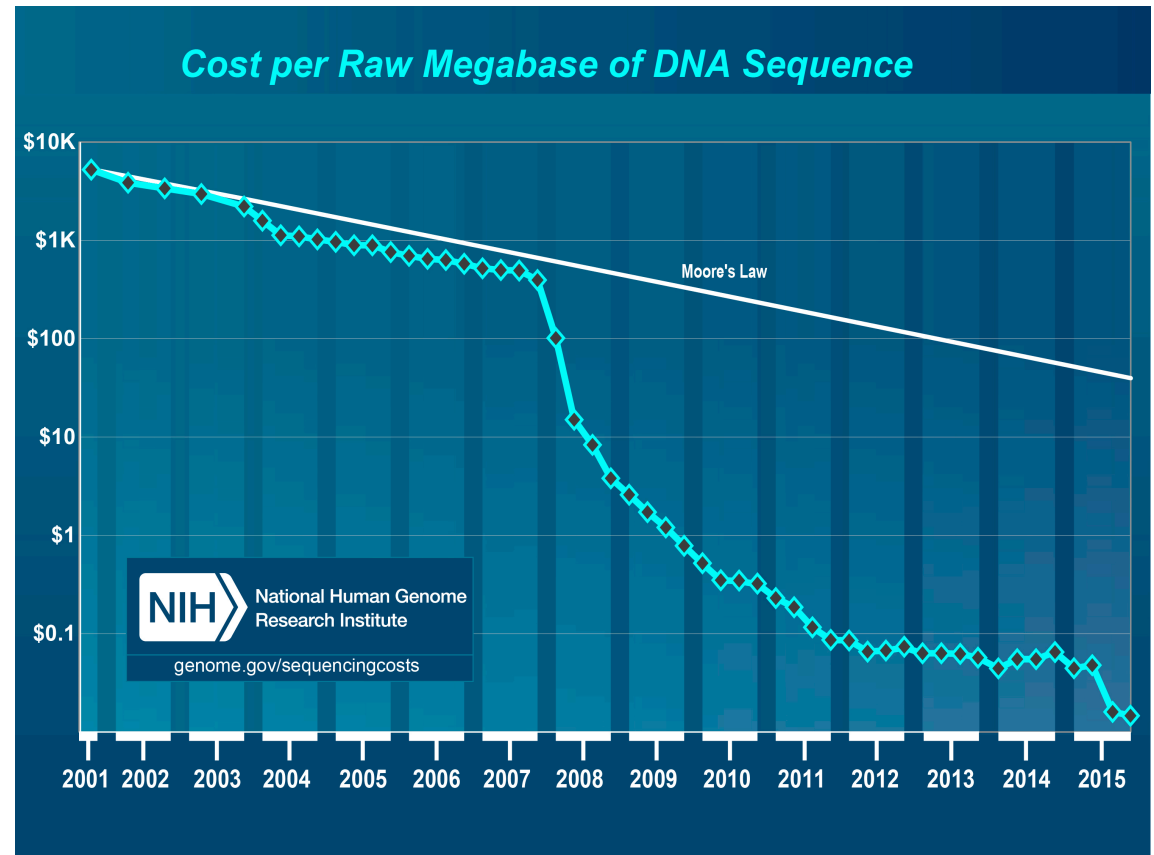
Slides freely downloadable from Lectures.GersteinLab.org

& "tweetable" (via @markgerstein). See last slide for more info.

# 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
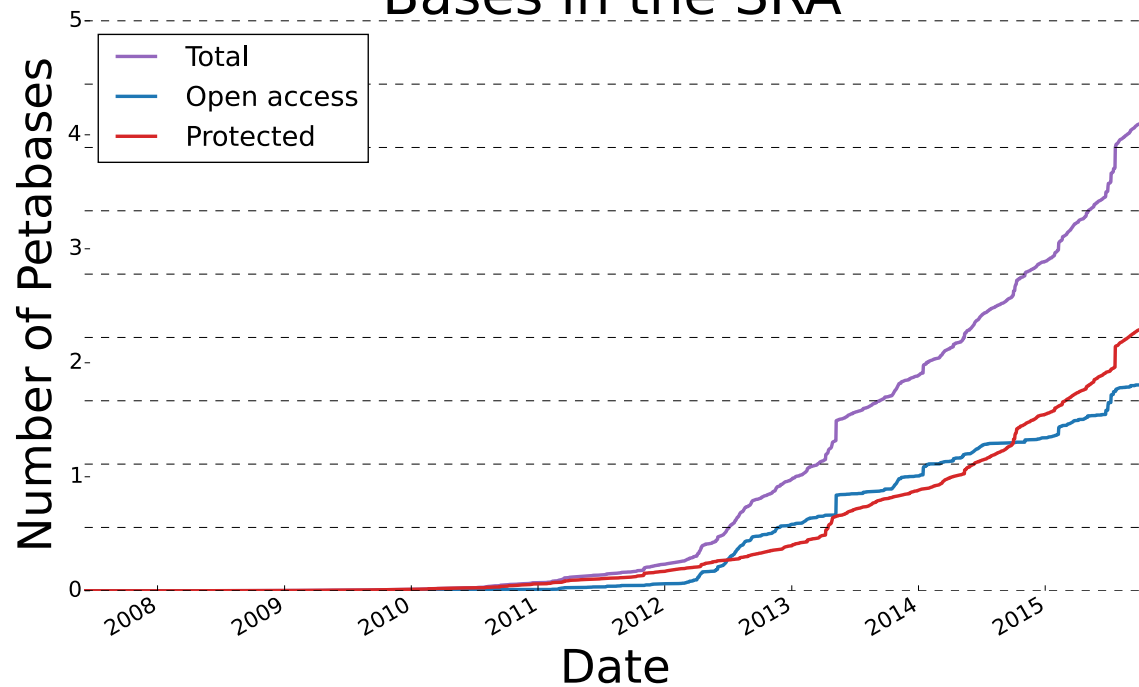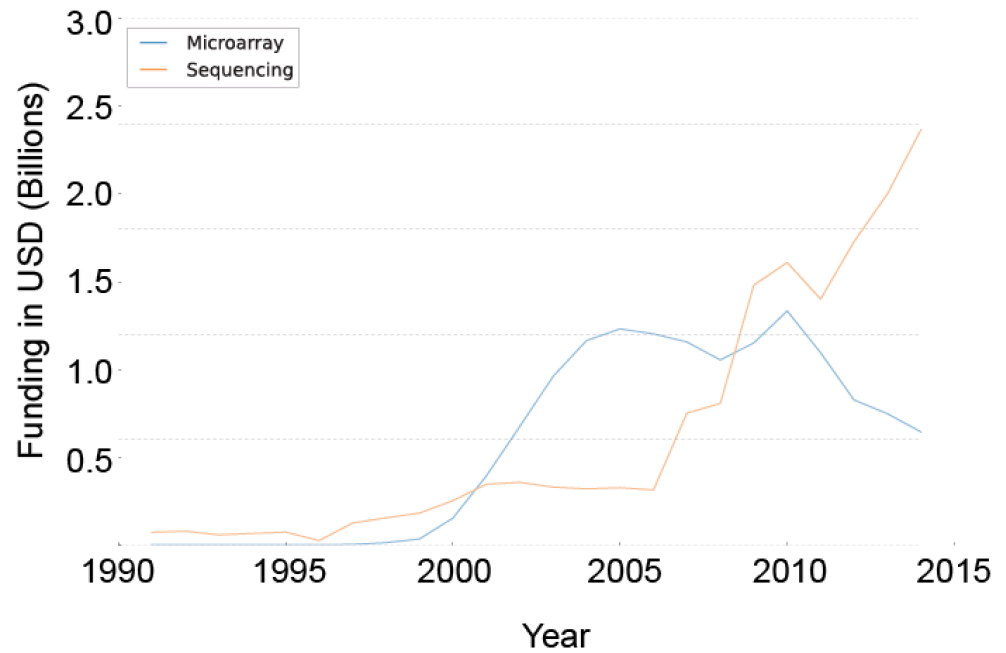
genome.gov/sequencingcosts

# 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

Y-axis: Number of Petabases
X-axis: Date (2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015)

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



Legend: Microarray, Sequencing

Y-axis: Funding in USD (Billions)
X-axis: Year (1990, 1995, 2000, 2005, 2010, 2015)

# 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 |

### Prevalence of Variants



Passenger

Driver (~0.1%)

Common

Rare* (1-4%)

Common

Rare (~75%)

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

4

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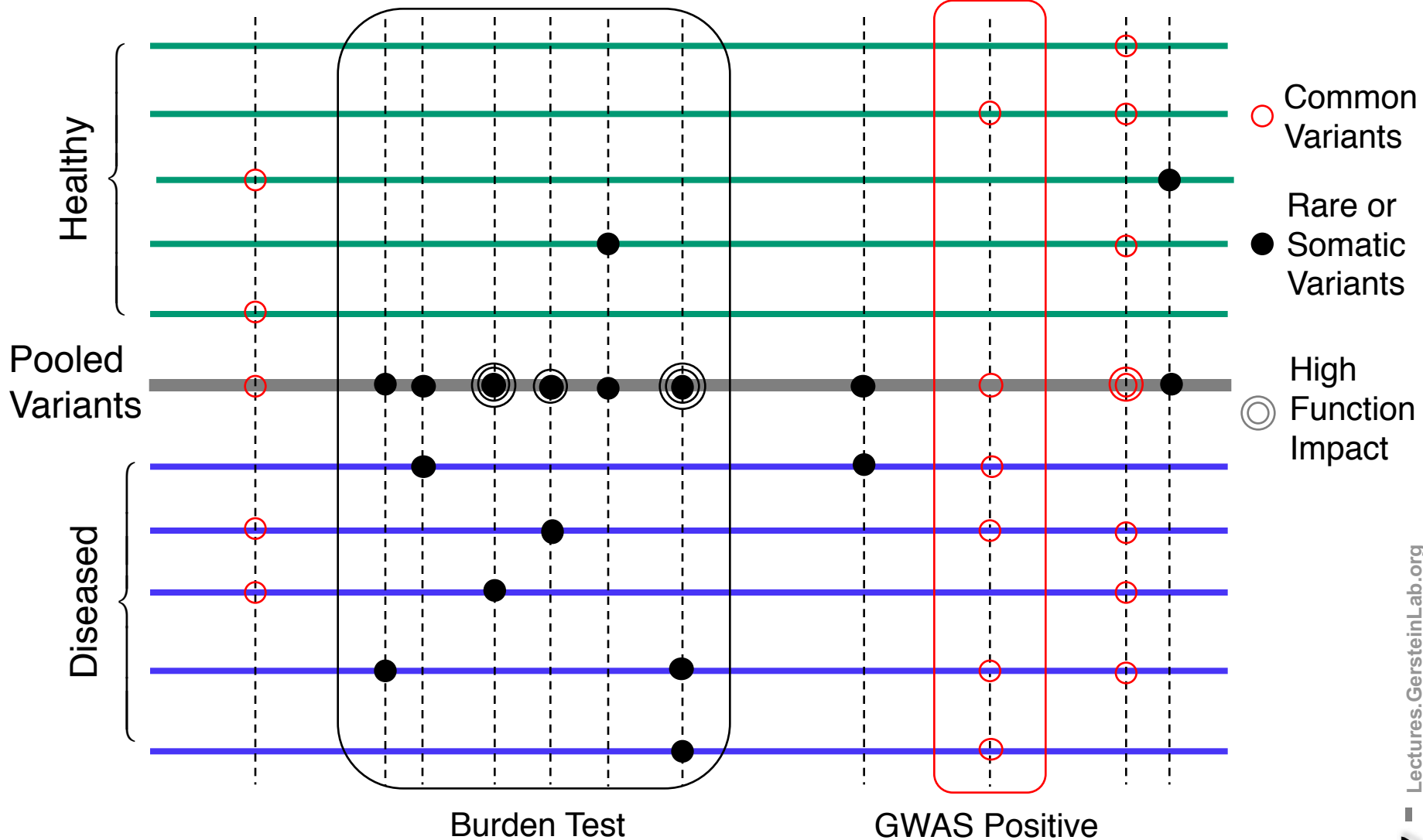- ## Drivers
  - Driver mutation is a mutation that directly or indirectly confers a selective growth advantage to the cell in which it occurs.
  - A typical tumor contains 2-8 drivers; the remaining mutations are passengers.
- ## Passengers
  - Conceptually, a passenger mutation has no direct or indirect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

Vogelstein B. Science 2013. 339(6127):1546-1558

# Association of Variants with Diseases



Healthy

Pooled Variants

Diseased

Burden Test

GWAS Positive

○ Common Variants

● Rare or Somatic Variants

◎ High Function Impact

# Rare variant analysis particularly applicable at the moment to Exomes



- CMG rare disease variants & TCGA somatic variants

  - Main NIH disease genomic project
  - Both of these focus on "rare" variant for which GWAS is not meaningful
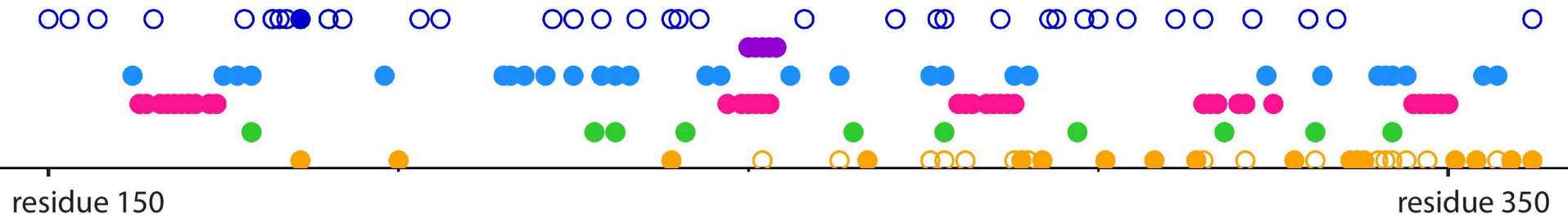  - Larger numbers of individual exomes more important than WGS

- Exomes have the current potential for great scale with the better impact interpretability of coding variants, often in a region of known protein structure

  - Scale of EXAC, >60K exomes
    [Lek et al. '16]

# 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)]

## Developing Tools for evaluating the impact of rare variants in coding regions

- New tools to wring everything out of protein structure
    - Stress for finding cryptic sites
    - Frustration for rapidly evaluating packing changes
    - (MotifVar) Intensification for using the amplifying power of protein structural motifs (eg TPR)
- Another approach – looking for allelic variants

**Analysis of Personal Genomes:
Evaluating the impact of variants in exomes
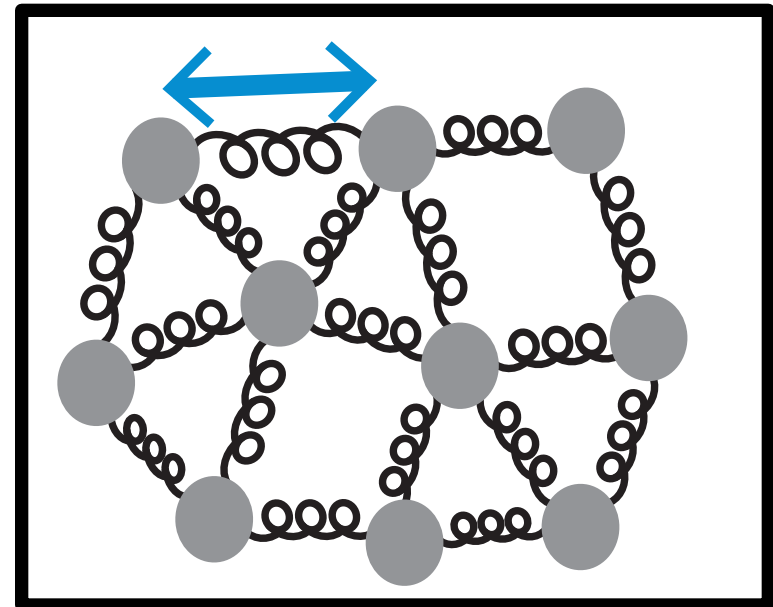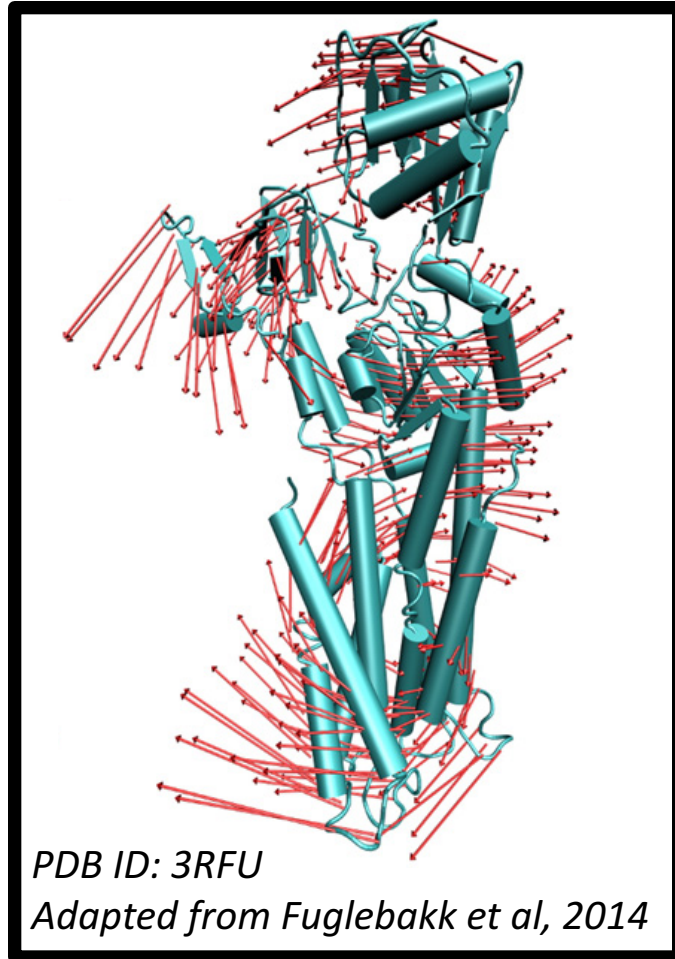using protein structure & allelic activity**

- Introduction
  - Rare v common variants
  - The importance of interpreting rare coding variants in the context of disease genomics (CMG,TCGA)
- Identifying cryptic allosteric sites with STRESS
  - On surface & in interior bottlenecks
- Using changes in localized frustration to find further sites sensitive to mutations
  - Difference betw. TSGs & oncogenes

- Using structural motifs (eg TPR) for intensification of weak population genetic signals
  - For both negative and positive selection
- Prioritizing allelic genes using AlleleDB
  - Having observed difference in molecular activity in many contexts

# Analysis of Personal Genomes:
# Evaluating the impact of variants in exomes
# using protein structure & allelic activity

- Introduction
  - Rare v common variants
  - The importance of interpreting rare coding variants in the context of disease genomics (CMG,TCGA)
- Identifying cryptic allosteric sites with STRESS
  - On surface & in interior bottlenecks
- Using changes in localized frustration to find further sites sensitive to mutations
  - Difference betw. TSGs & oncogenes

- Using structural motifs (eg TPR) for intensification of weak population genetic signals
  - For both negative and positive selection
- Prioritizing allelic genes using AlleleDB
  - Having observed difference in molecular activity in many contexts
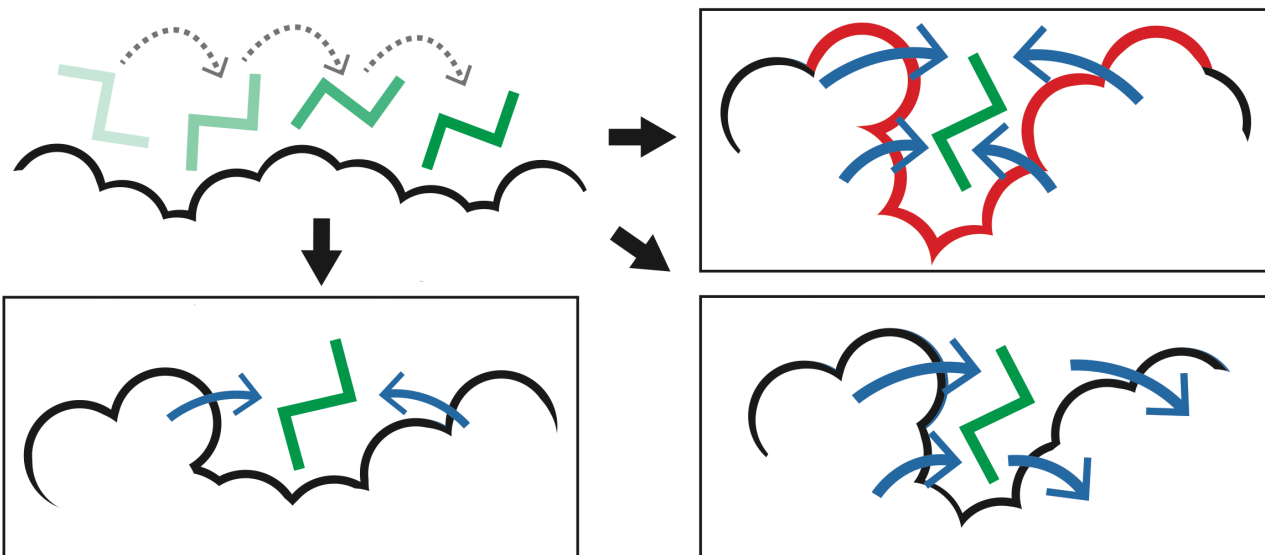
# Models of Protein Conformational Change

## Motion Vectors from Normal Modes (ANMs)



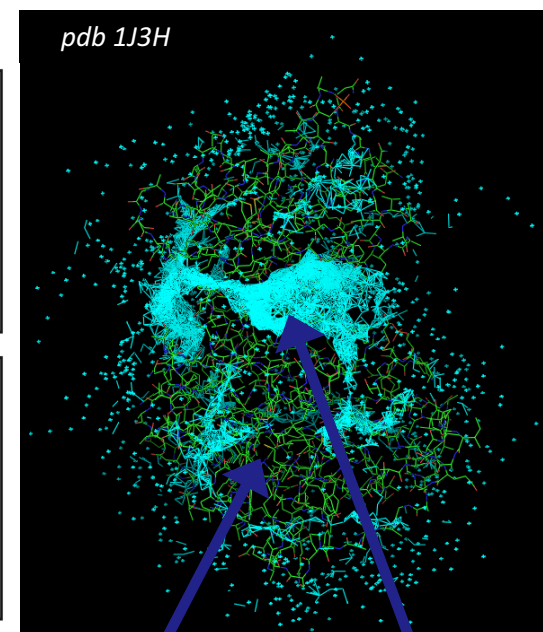*PDB ID: 3RFU*
*Adapted from Fuglebakk et al, 2014*



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



pdb 1J3H

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

Surface region with high density of candidate sites
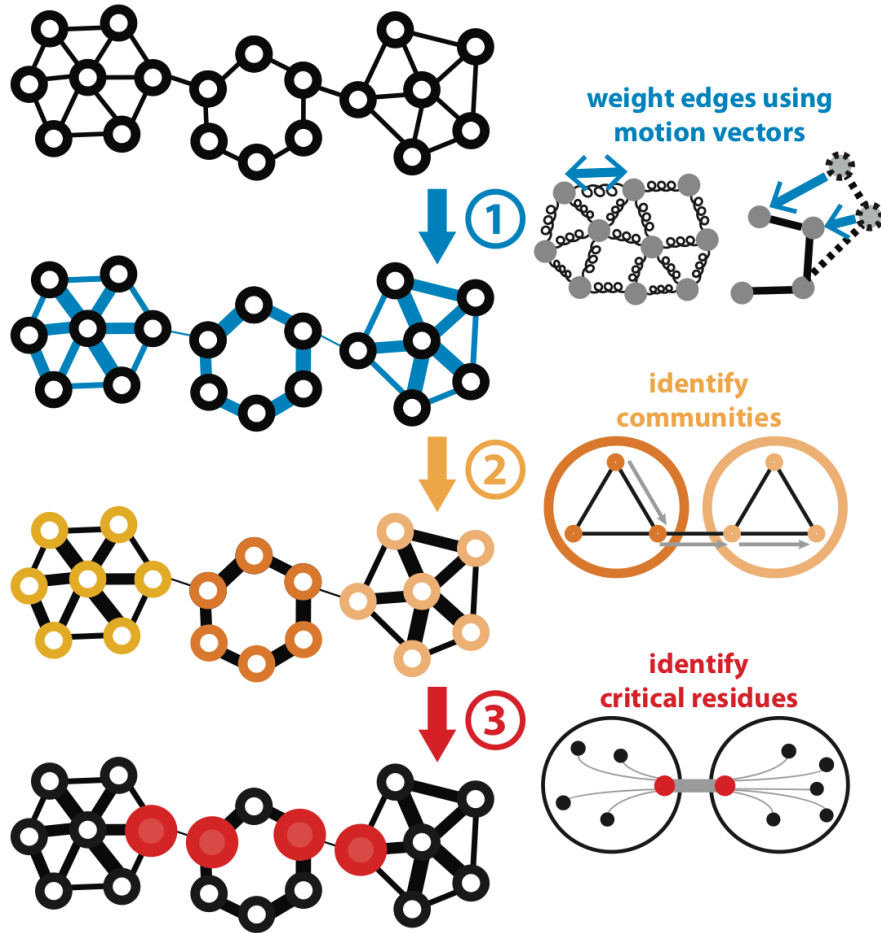
Surface region with low density of candidate sites

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

① 

identify communities

② 

identify critical residues

③

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

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

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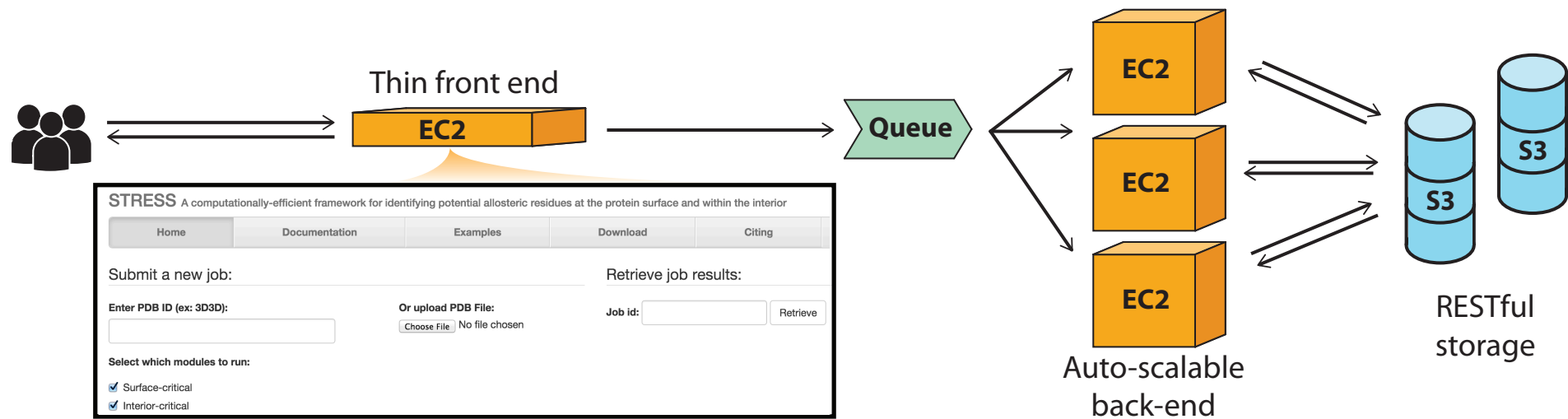
$$D_{ij} = -\log(|C_{ij}|)$$

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



Thin front end

EC2

Queue

EC2
EC2
EC2

Auto-scalable back-end

S3
S3

RESTful storage

**STRESS** A computationally-efficient framework for identifying potential allosteric residues at the protein surface and within the interior

| Home | Documentation | Examples | Download | Citing |
|---|---|---|---|---|

Submit a new job:                                    Retrieve job results:

Enter PDB ID (ex: 3D3D):        Or upload PDB File:        Job id: [_____] [Retrieve]
[_____]                  [Choose File] No file chosen

Select which modules to run:
☑ Surface-critical
☑ Interior-critical

- 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*



p=0.309

### *Interior*



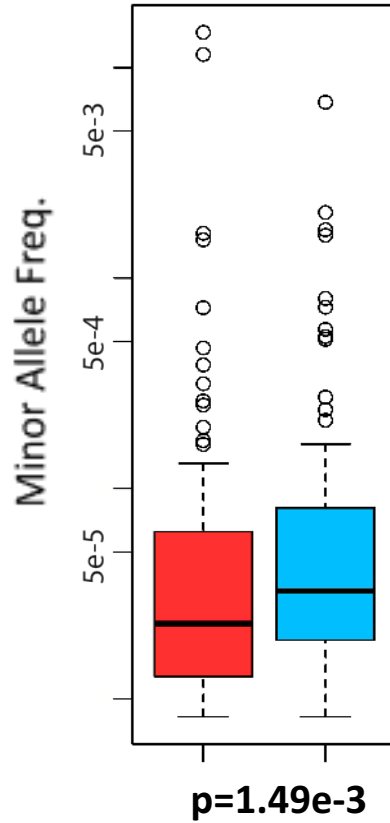**p=1.80e-05**

- critical (red)
- non-critical (blue)

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

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

### *Surface*



**p=1.49e-3**

### *Interior*



**p=7.98e-09**


critical
non-critical

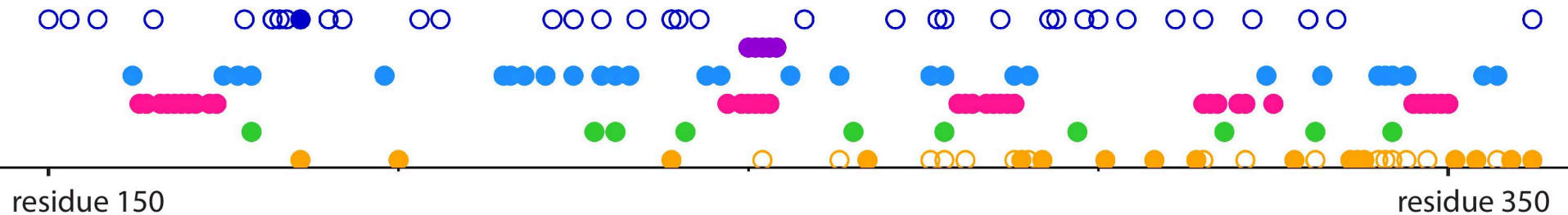Adapted from Clarke*, Sethi*, et al *(in press)*

# 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
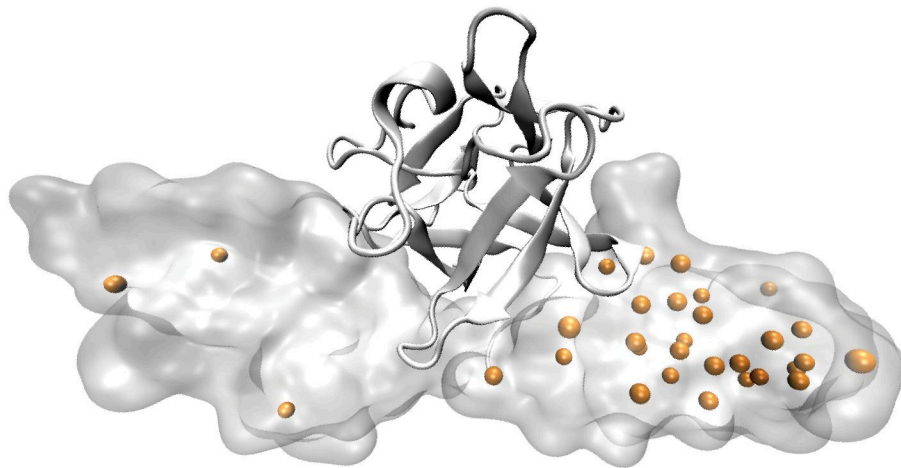


*Fibroblast growth factor receptor 2 (pdb: 1IIL)*

● ○ 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
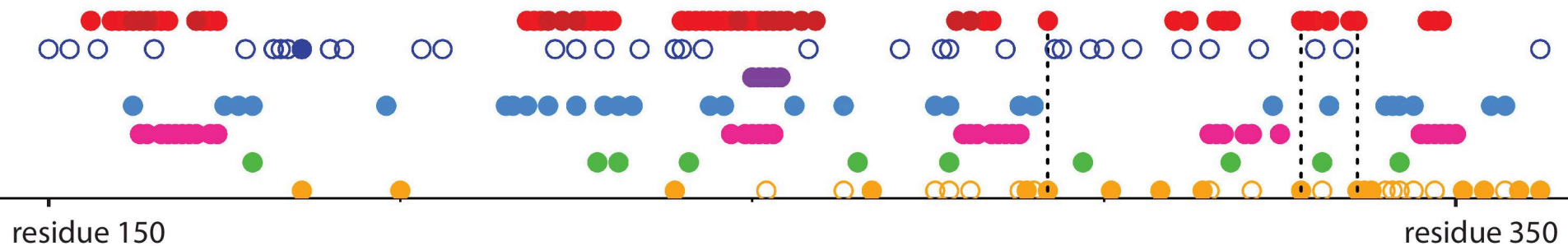
[Sethi et al. COSB ('15)]

# 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
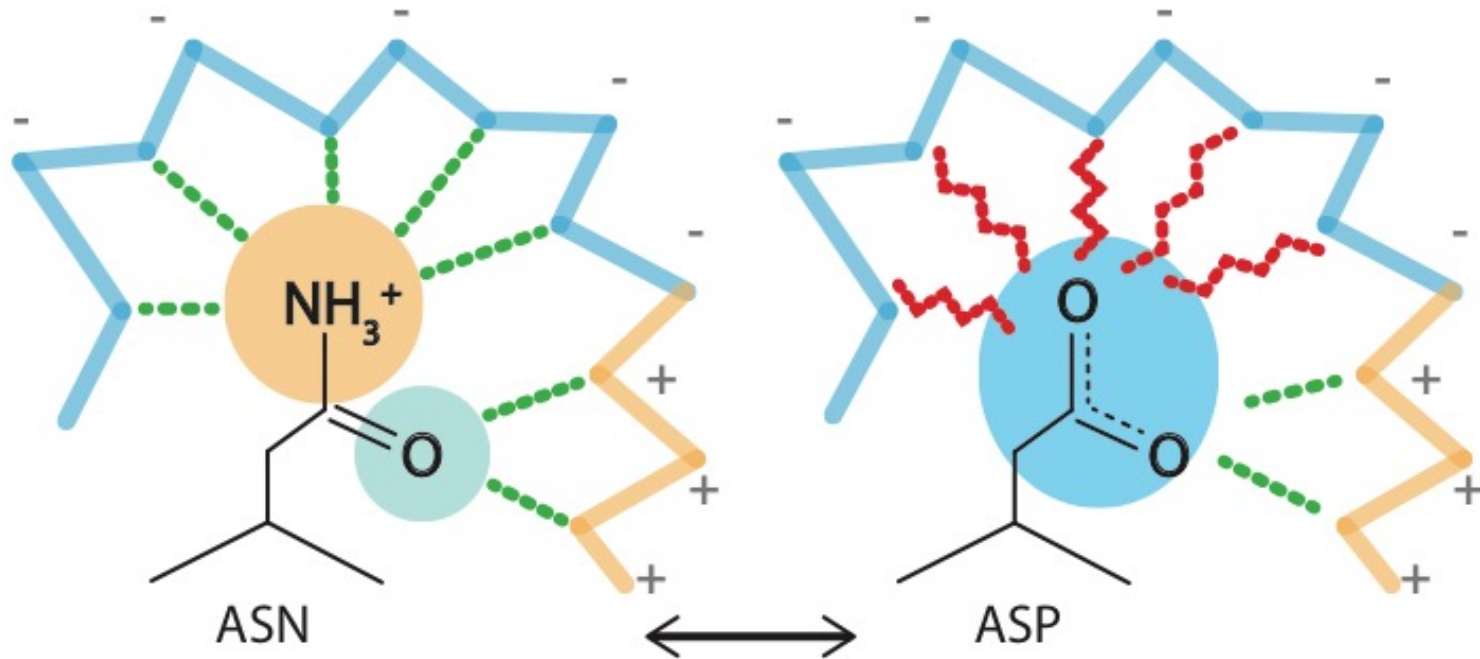- ● HGMD site (w/o annotation overlap)
- ○ HGMD site (w/annotation overlap)



residue 150                                                          residue 350

[Sethi et al. COSB ('15)]

# Analysis of Personal Genomes:
## Evaluating the impact of variants in exomes using protein structure & allelic activity

- Introduction
  - Rare v common variants
  - The importance of interpreting rare coding variants in the context of disease genomics (CMG,TCGA)
- Identifying cryptic allosteric sites with STRESS
  - On surface & in interior bottlenecks
- Using changes in localized frustration to find further sites sensitive to mutations
  - Difference betw. TSGs & oncogenes

- Using structural motifs (eg TPR) for intensification of weak population genetic signals
  - For both negative and positive selection
- Prioritizing allelic genes using AlleleDB
  - Having observed difference in molecular activity in many contexts

# Schematic illustration of localized frustration



ASN ⟷ ASP

Legend:
more negative — more positive
favorable interaction (green dashed)
unfavorable interaction (red zigzag)

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

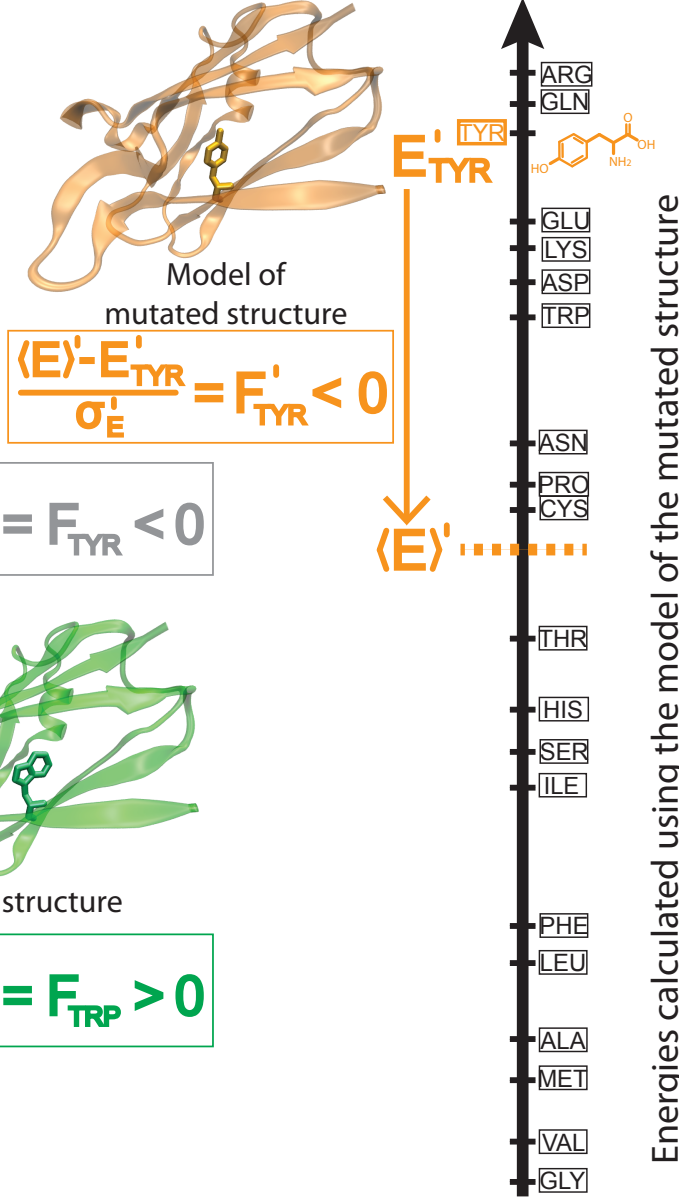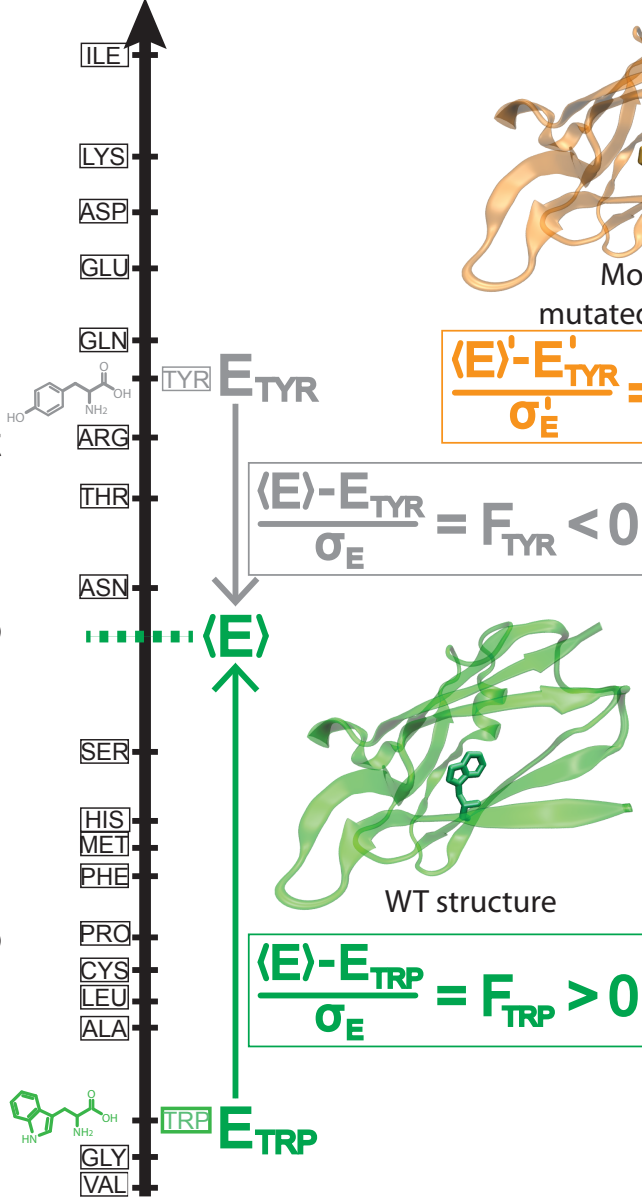Workflow for evaluating localized frustration changes (ΔF)

Measuring perturbation with naive calculation

$$F_{TYR} - F_{TRP} = \widetilde{\Delta F} < 0$$

Measuring perturbation with secondary calculation
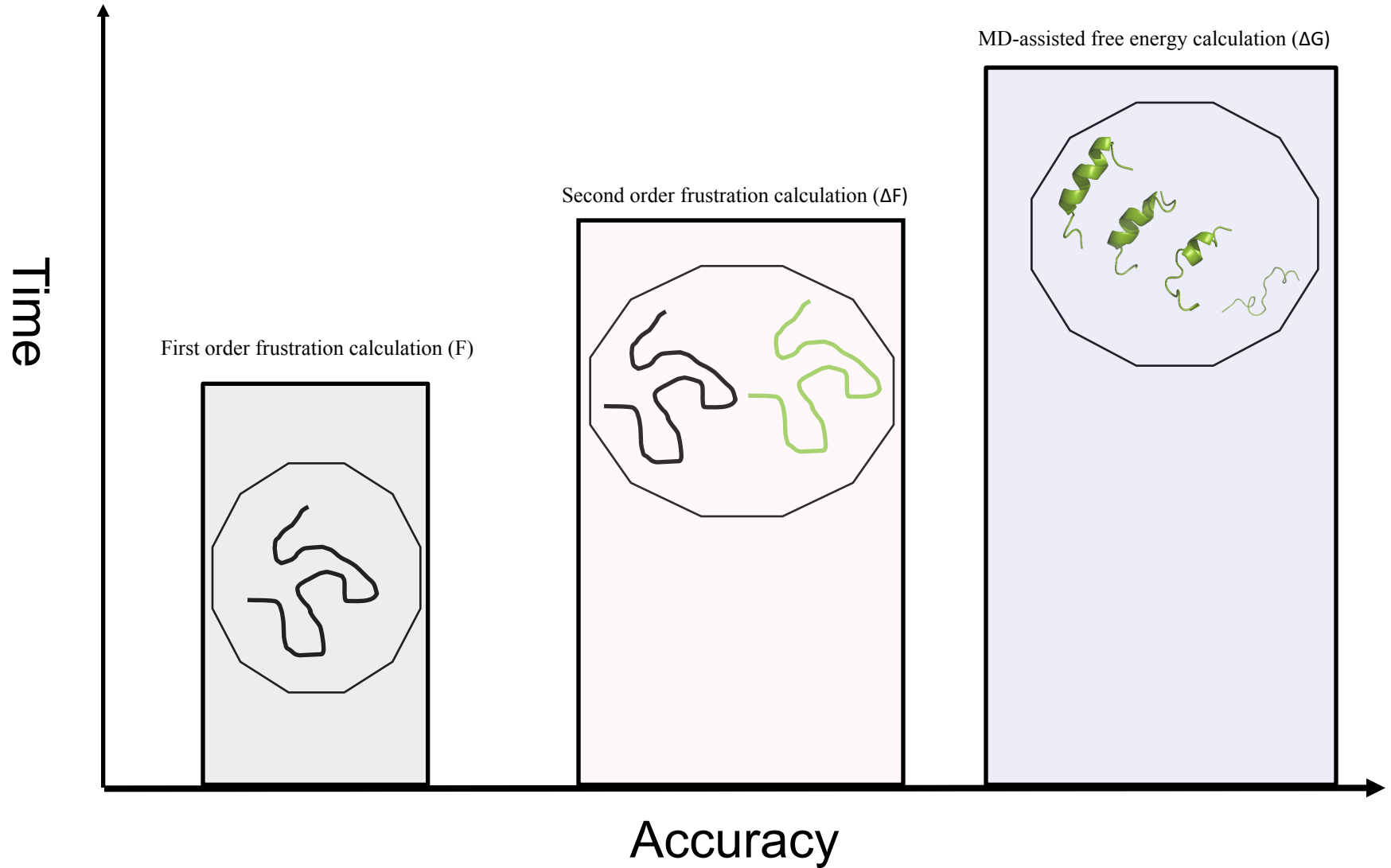
$$F'_{TYR} - F_{TRP} = \Delta F < 0$$

Energies calculated using the wild-type structure

ILE
LYS
ASP
GLU
GLN
TYR   $E_{TYR}$
ARG
THR

$$\frac{\langle E \rangle - E_{TYR}}{\sigma_E} = F_{TYR} < 0$$

ASN

⟨E⟩

SER
HIS
MET
PHE
PRO
CYS
LEU
ALA

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

WT structure

TRP   $E_{TRP}$
GLY
VAL

Model of mutated structure

$$\frac{\langle E \rangle' - E'_{TYR}}{\sigma'_E} = F'_{TYR} < 0$$

$E'_{TYR}$   TYR

⟨E⟩'

Energies calculated using the model of the mutated structure

ARG
GLN
GLU
LYS
ASP
TRP
ASN
PRO
CYS
THR
HIS
SER
ILE
PHE
LEU
ALA
MET
VAL
GLY

# Striking a balance:
# the complexity of the second order frustration calculation



MD-assisted free energy calculation (ΔG)

Second order frustration calculation (ΔF)

First order frustration calculation (F)

Time

Accuracy

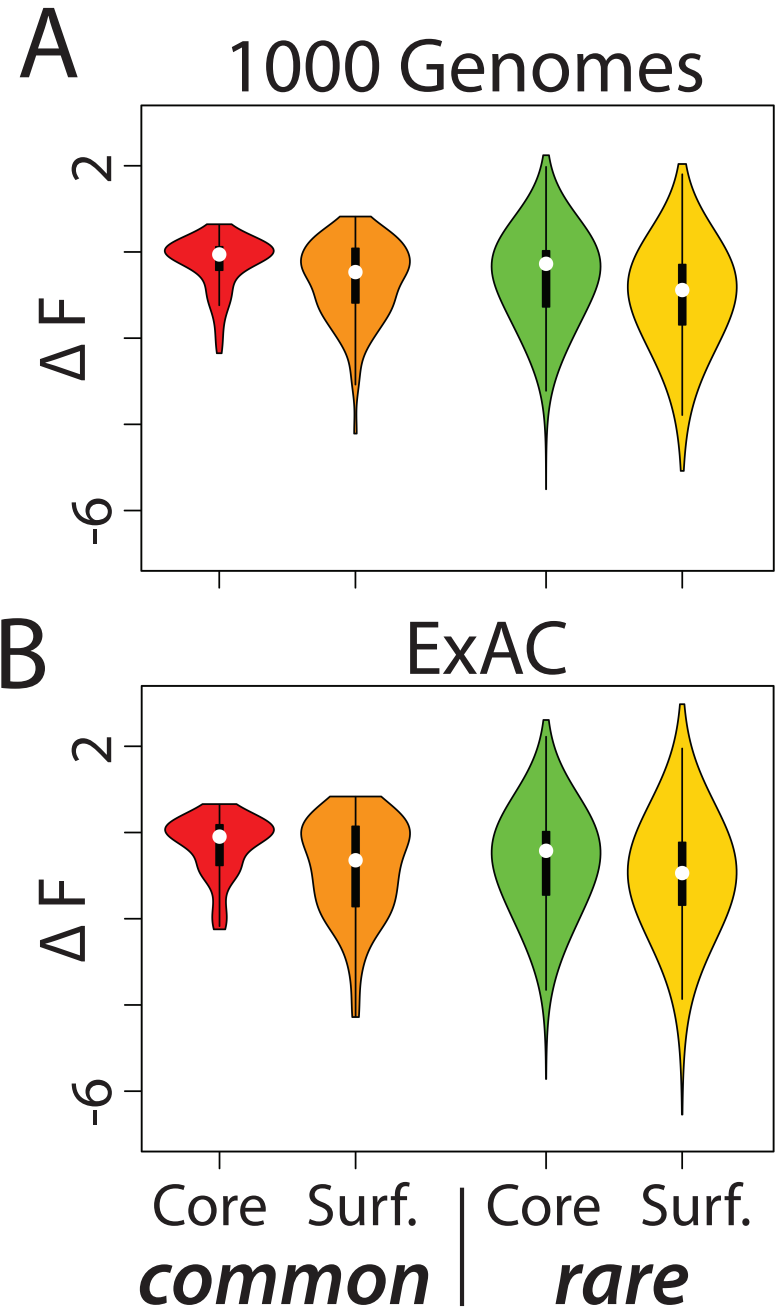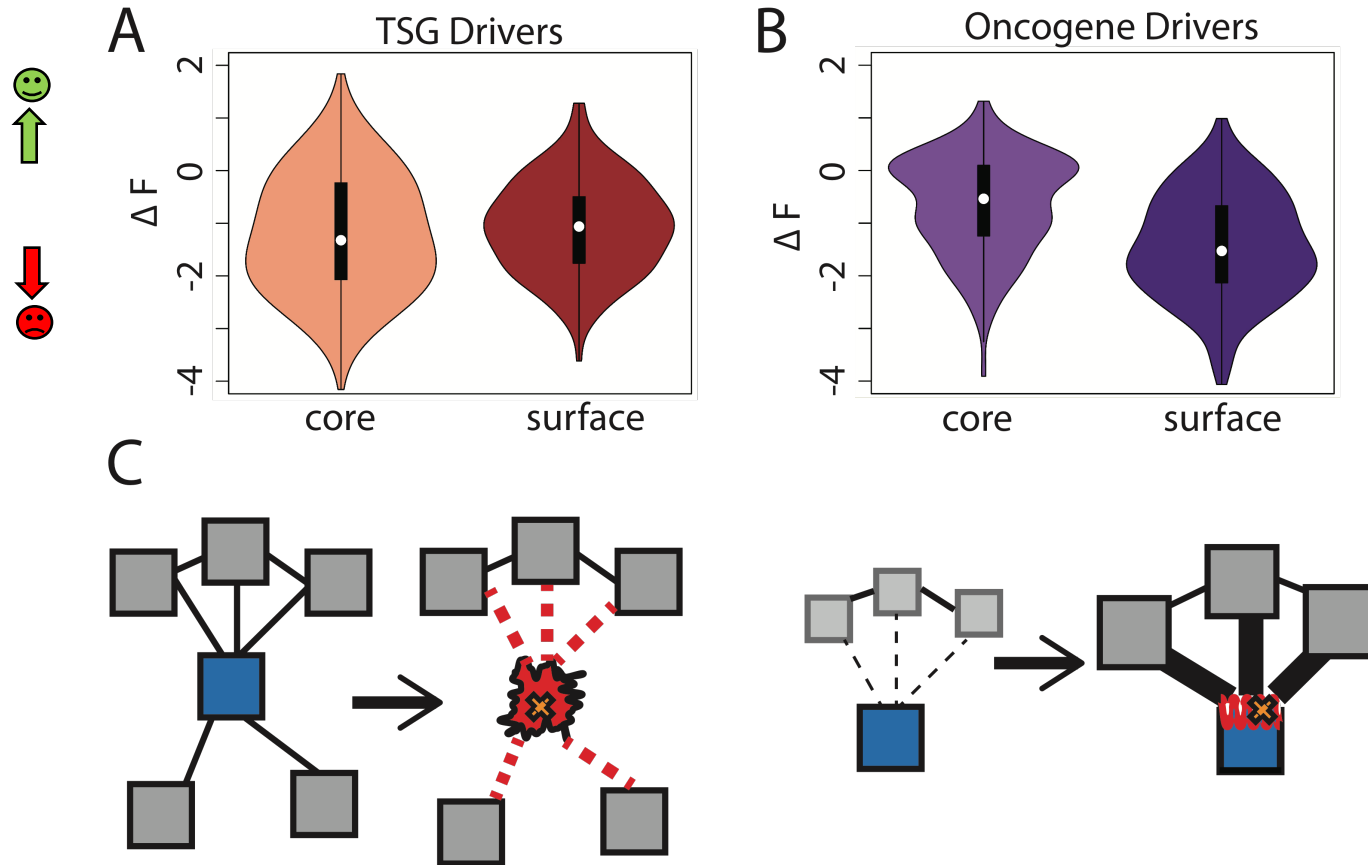# Comparing ΔF values across different SNV categories: Normal v disease



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

# ΔF distributions among rare v. common SNVs

Rare mutations cause more unfavorable frustration change than common ones



A — 1000 Genomes

B — ExAC

Core  Surf.  Core  Surf.

*common* | *rare*

# Comparison between ΔF distributions: TSGs v. oncogenes



A — TSG Drivers — core, surface

B — Oncogene Drivers — core, surface

C

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.

**Analysis of Personal Genomes:
Evaluating the impact of variants in exomes
using protein structure & allelic activity**

- Introduction
  - Rare v common variants
  - The importance of interpreting rare coding variants in the context of disease genomics (CMG,TCGA)
- Identifying cryptic allosteric sites with STRESS
  - On surface & in interior bottlenecks
- Using changes in localized frustration to find further sites sensitive to mutations
  - Difference betw. TSGs & oncogenes

- Using structural motifs (eg TPR) for intensification of weak population genetic signals
  - For both negative and positive selection
- Prioritizing allelic genes using AlleleDB
  - Having observed difference in molecular activity in many contexts

# Intensification amplifies signals from motif-based MSAs

# Intensification amplifies signals from motif-based MSAs

1. **Find motifs**

1. **Generate motif-MSA**
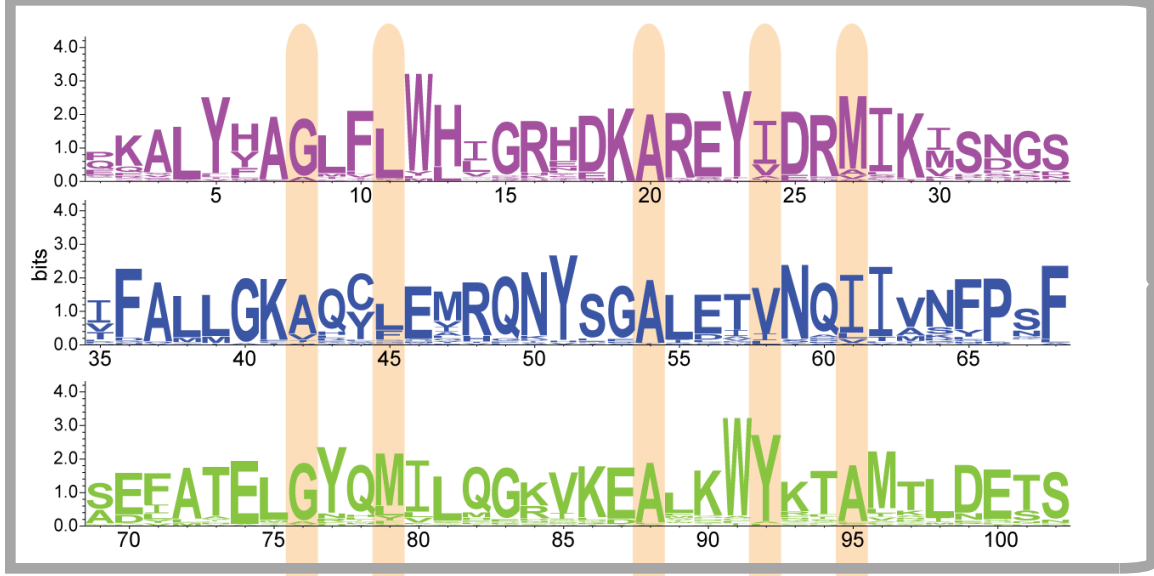
1. **Map SNVs to motif-MSA**

1. **Evaluate SNV profiles**

1. **Store in database**
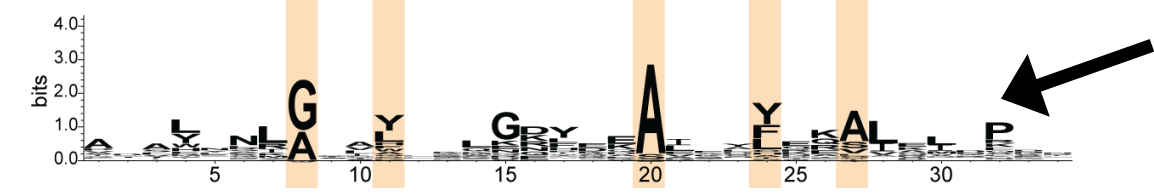


Motif-MSA and SNV profiles for:
a) amino acid freq
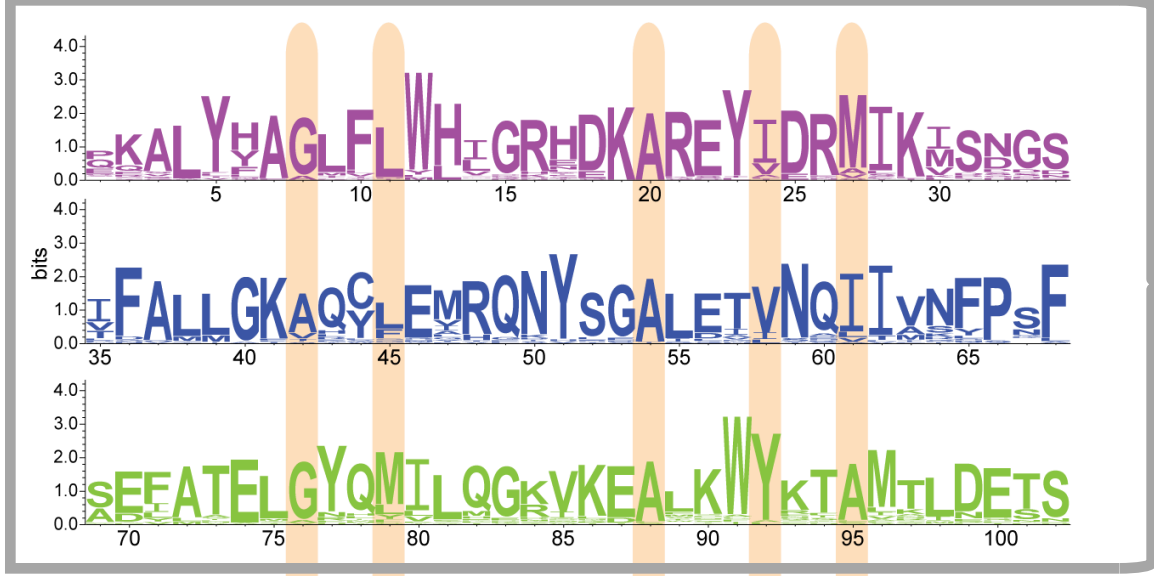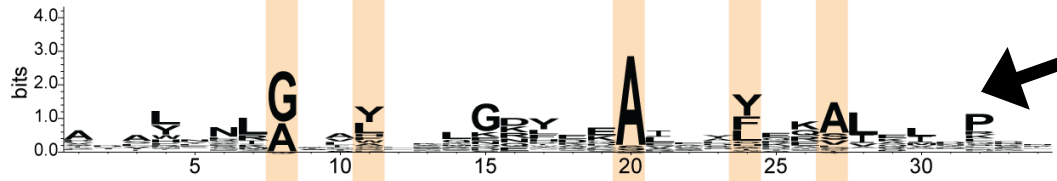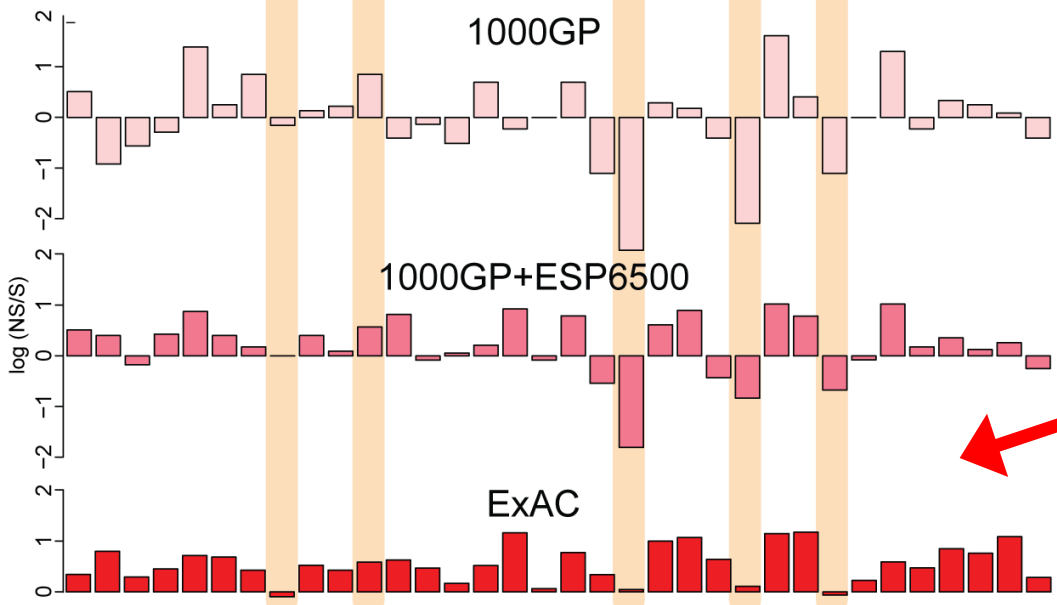b) SIFT scores
c) R/C
d) NS/S
e) ΔDAF (pop)

Species MSAs

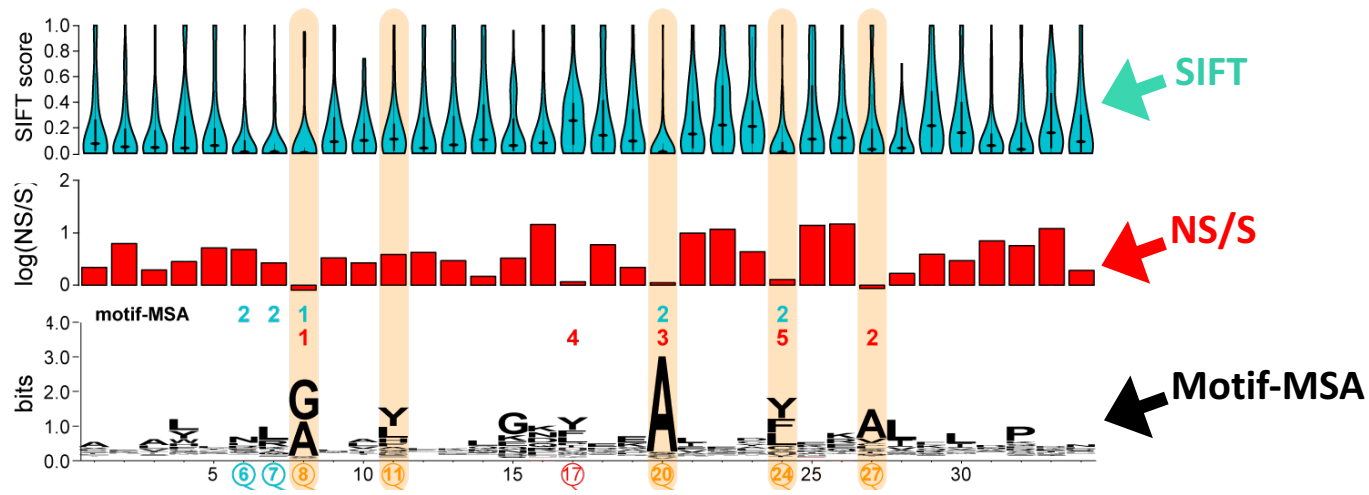**Motif-MSA uncovers important positions missed by species-MSA**

Species MSAs

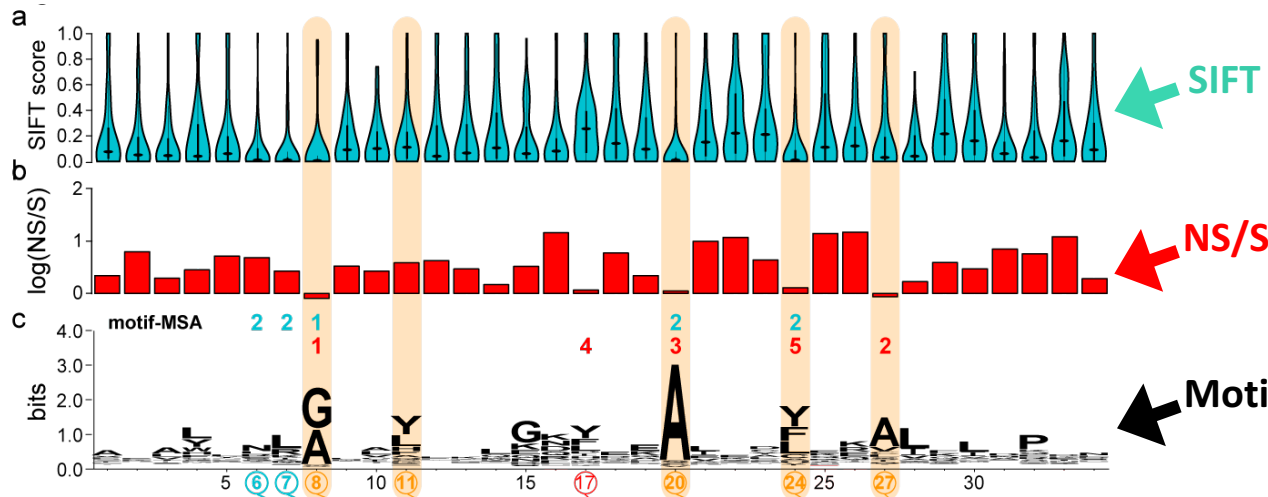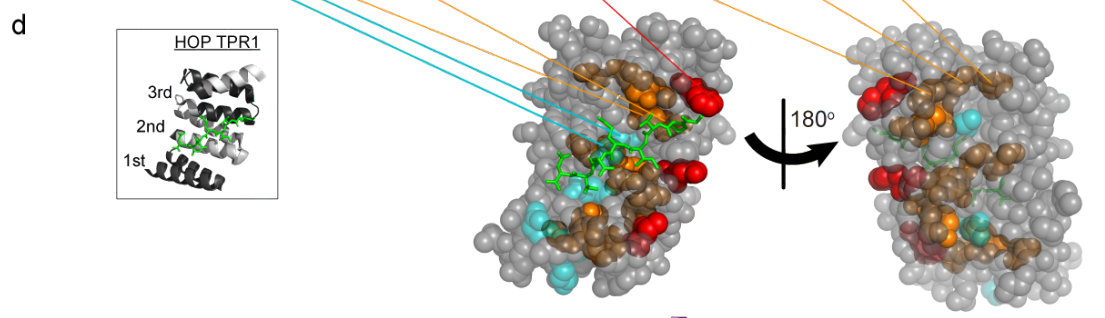Motif-MSA uncovers important
positions missed by species-MSA

1000GP

1000GP+ESP6500

ExAC

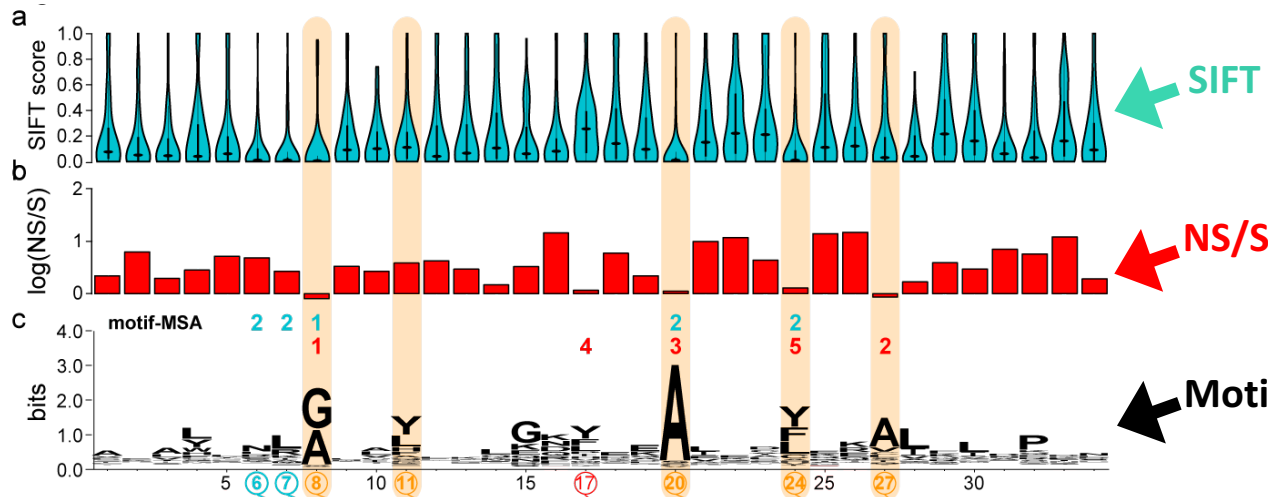**Signal-to-noise is the best in ExAC**

**Selection in PPI motifs**

**Selection
in PPI motifs**

SIFT

NS/S

Motif-MSA

How to check possible significance:
Burial within structure

**Selection in PPI motifs**

a — SIFT score — SIFT

b — log(NS/S) — NS/S

c — motif-MSA — Motif-MSA

d — HOP TPR1 — 3rd, 2nd, 1st — 180°

e — number of SNVs — ClinVar, HGMD, combined

**How to check possible significance:**

**-> Burial within structure**

**-> more SNVs implicated in diseases in ClinVar and HGMD**
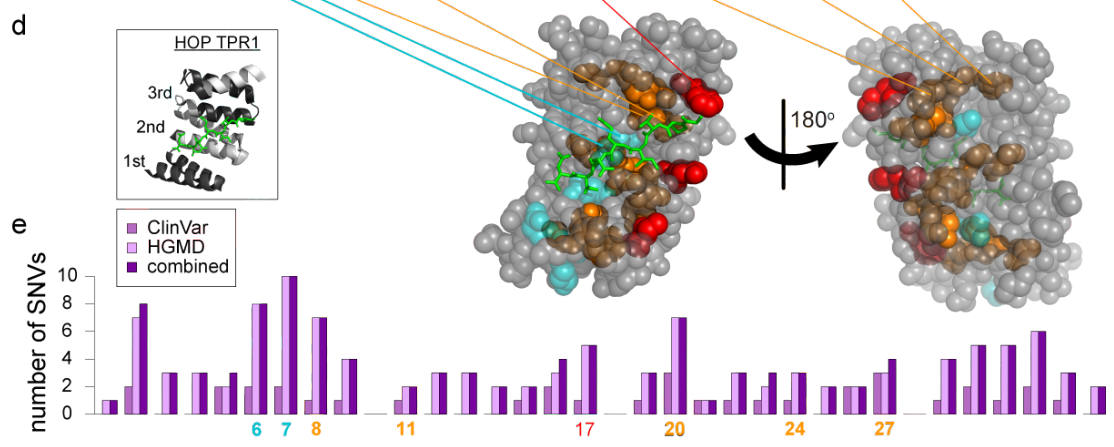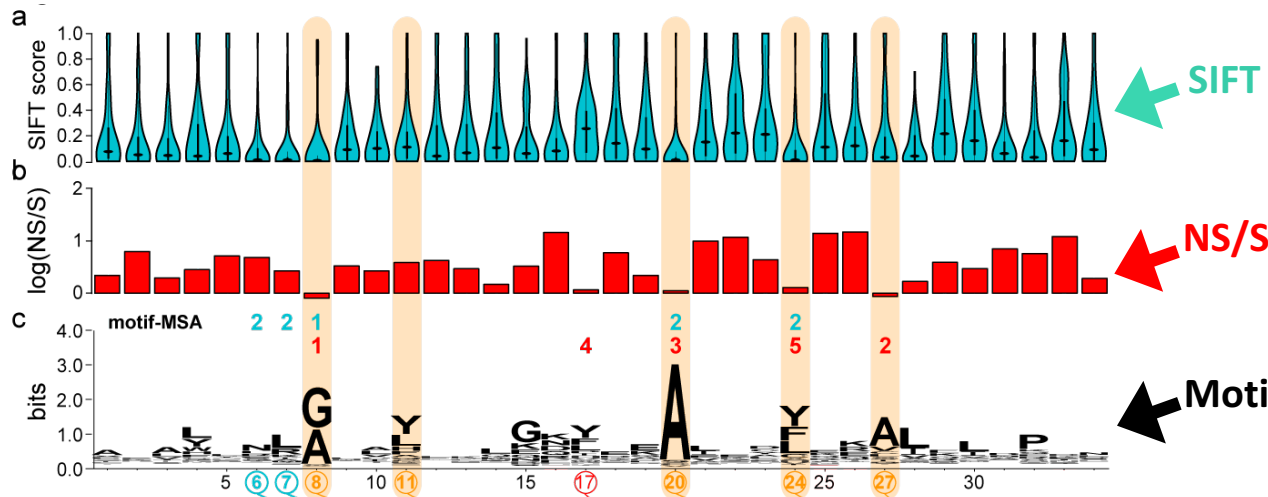
**Selection in PPI motifs**

SIFT

NS/S

Motif-MSA

How to check possible significance:

-> burial within structure
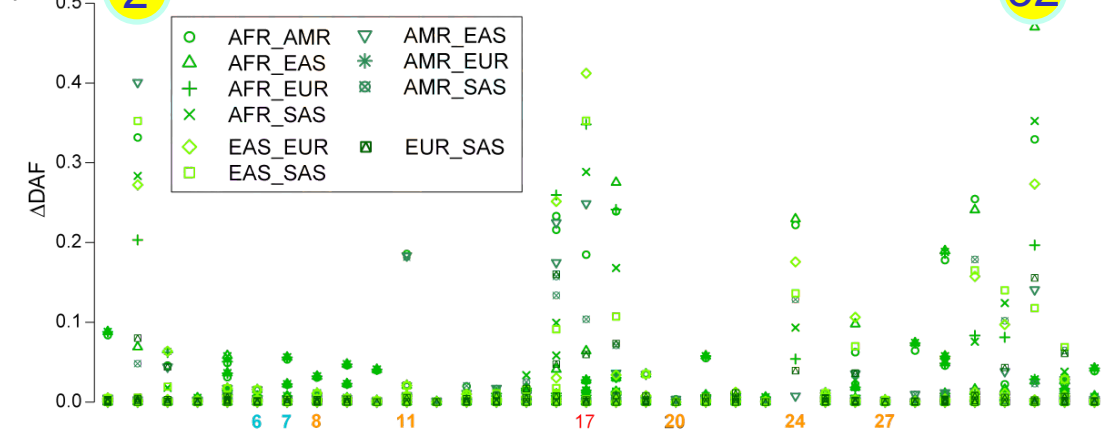
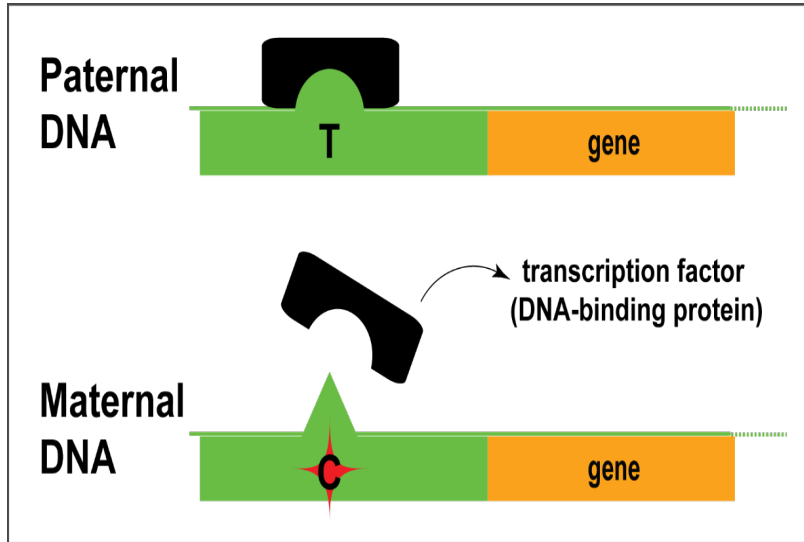-> more SNVs implicated in diseases in ClinVar and HGMD

-> sites with increased human pop. differentiation might indicate important position
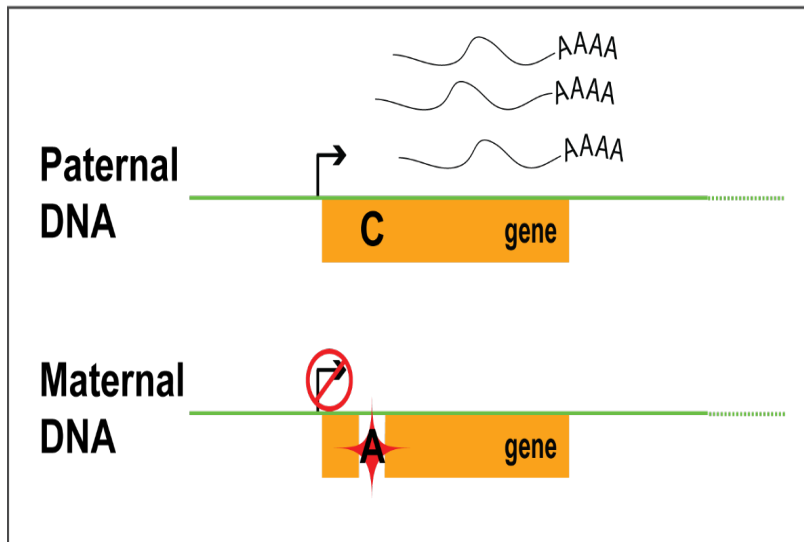
# Analysis of Personal Genomes:
## Evaluating the impact of variants in exomes using protein structure & allelic activity

- Introduction
  - Rare v common variants
  - The importance of interpreting rare coding variants in the context of disease genomics (CMG,TCGA)
- Identifying cryptic allosteric sites with STRESS
  - On surface & in interior bottlenecks
- Using changes in localized frustration to find further sites sensitive to mutations
  - Difference betw. TSGs & oncogenes

- Using structural motifs (eg TPR) for intensification of weak population genetic signals
  - For both negative and positive selection
- Prioritizing allelic genes using AlleleDB
  - Having observed difference in molecular activity in many contexts

# Allele-specific binding and expression



Genomic variants affecting allele-specific behavior e.g. allele-specific binding (ASB)

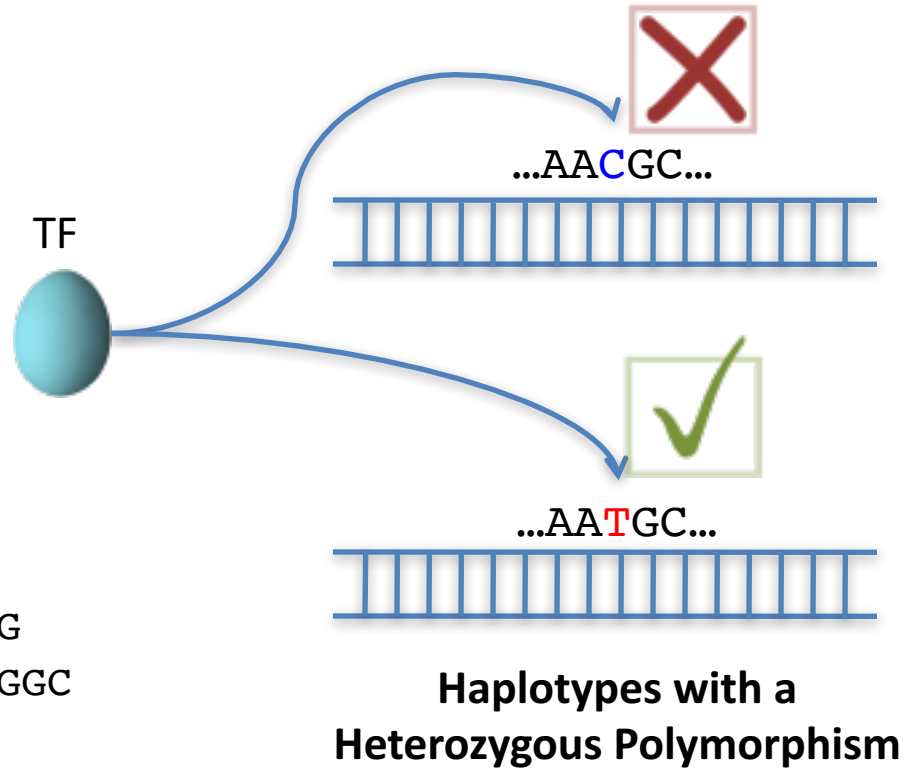e.g. allele-specific expression (ASE)

# Inferring Allele Specific Binding/Expression using Sequence Reads

**RNA/ChIP–Seq Reads**

```
ACTTTGATAGCGTCAATG
 CTTTGATAGCGTCAATGC
 CTTTGATAGCGTCAACGC
    TTGACAGCGTCAATGCAC
     TGATAGCGTCAATGCACG
      ATAGCGTCAATGCACGTC
       TAGCGTCAATGCACGTCG
        CGTCAACGCACGTCGGGA
         GTCAATGCACGTCGAGAG
          CAATGCACGTCGGGAGTT
           AATGCACGTCGGGAGTTG
            TGCACGTTGGGAGTTGGC
```
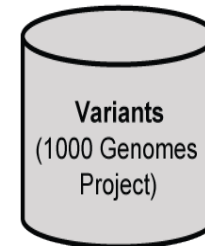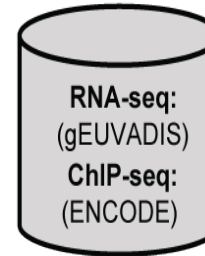
```
10 x T
 2 x C
```

TF

...AACGC...

...AATGC...

**Haplotypes with a
Heterozygous Polymorphism**

# AlleleDB: Building 382 personal genomes to detect allele-specific variants on a large-scale

1. Build personal genomes

x 382

Variants
(1000 Genomes Project)

RNA-seq:
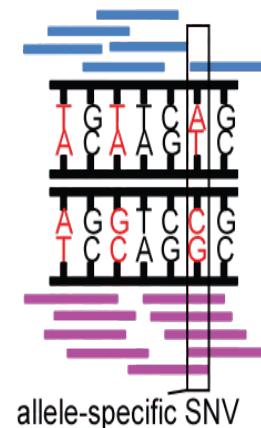(gEUVADIS)
ChIP-seq:
(ENCODE)

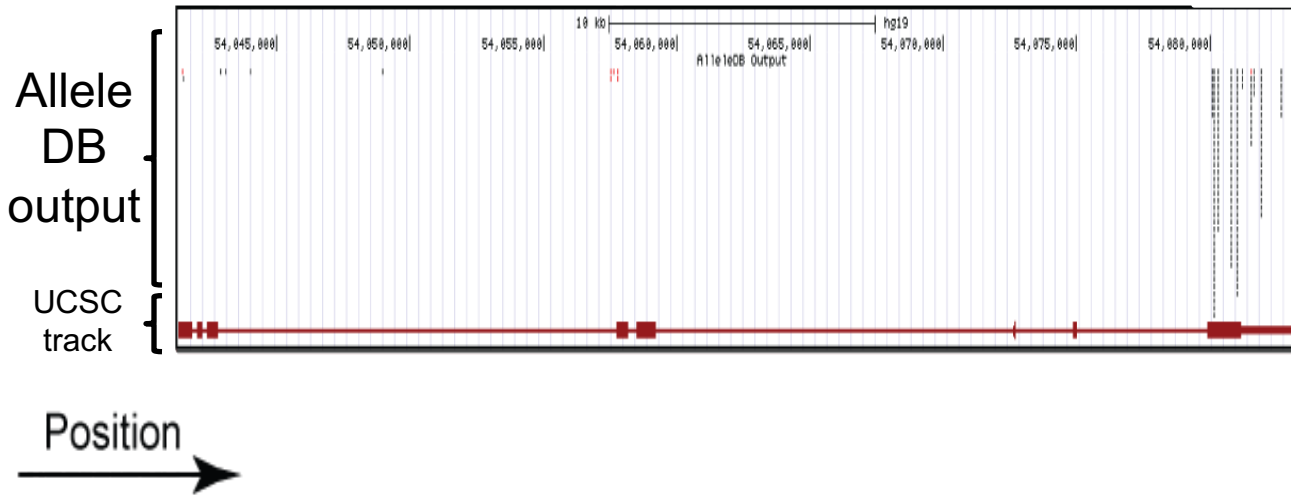2. Align ChIP-seq & RNA-seq reads

1. Detect allele-specific variants via a series of filters and tests

**Many Technical Issues: Reference bias, Ambiguous mapping bias, Over-dispersed (non binomial null)**

allele-specific SNV

alleledb.gersteinlab.org

# AlleleDB: Annotating rare & common allele-specific variants over a population



- Interfaces with UCSC genome browser
- Showing ZNF331 gene structure

[Chen *et al.* ('16) *Nat. Comm.*]
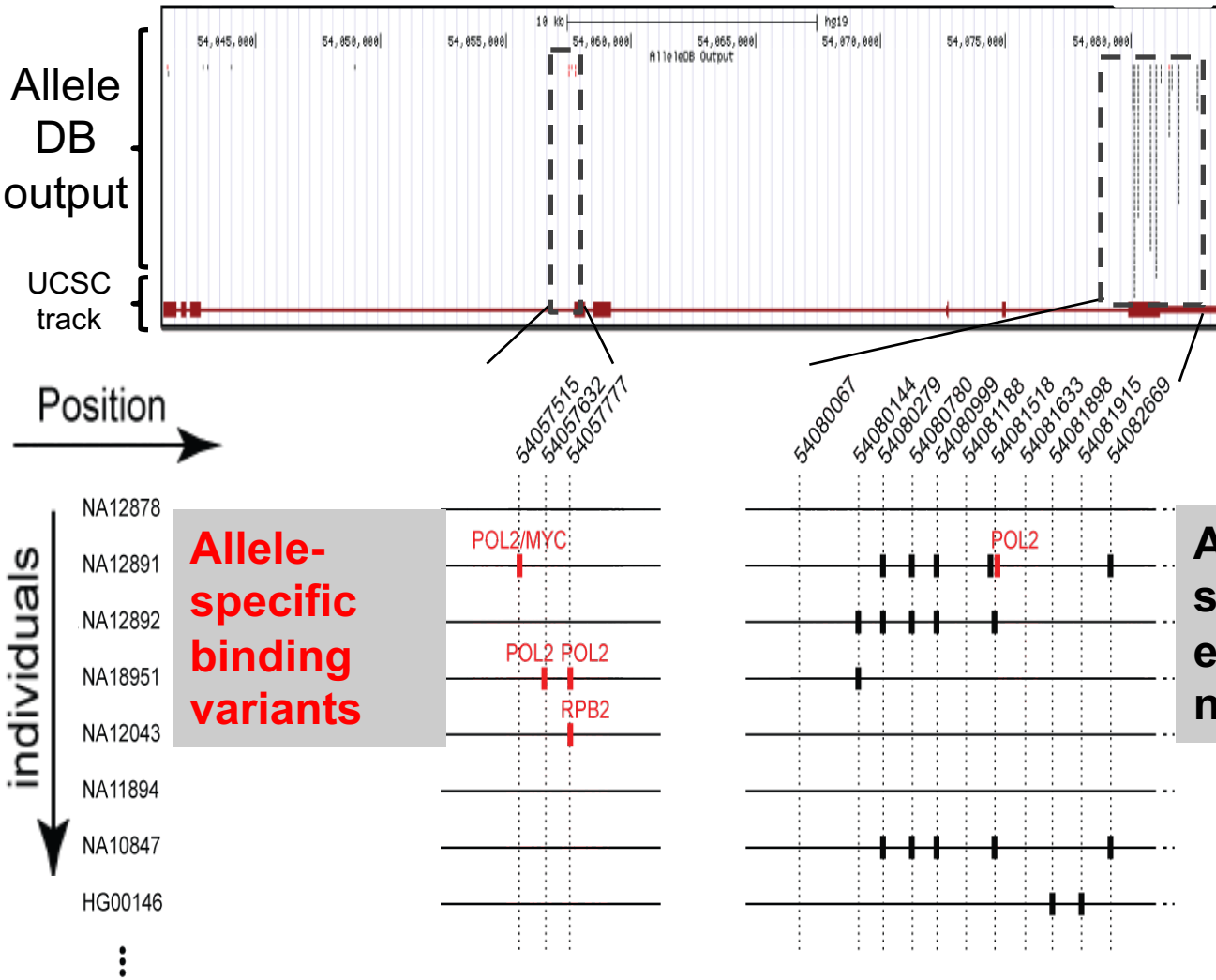
# AlleleDB: Annotating rare & common allele-specific variants over a population



- Interfaces with UCSC genome browser
- Showing ZNF331 gene structure

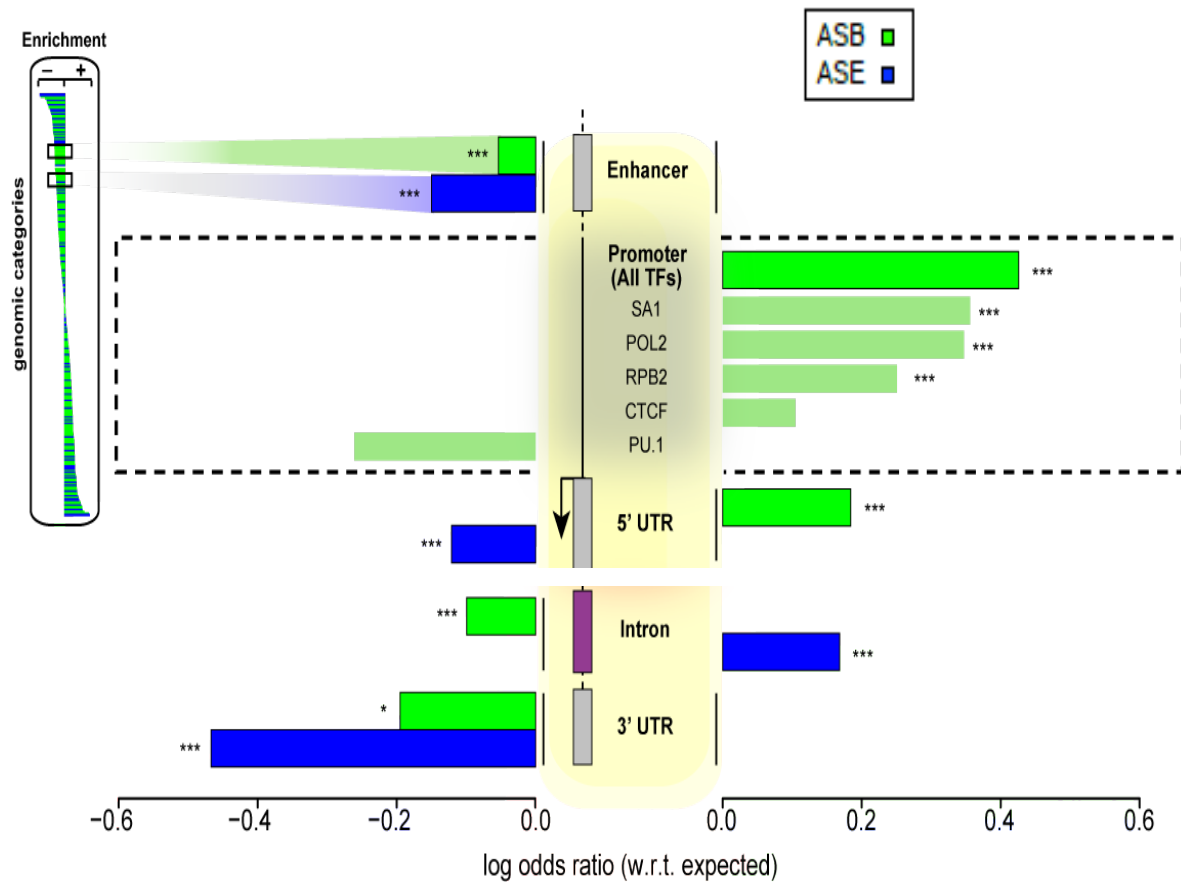[Chen *et al.* ('16) *Nat. Comm.*]

# Collecting ASE/ASB variants
# into allele-specific genomic regions

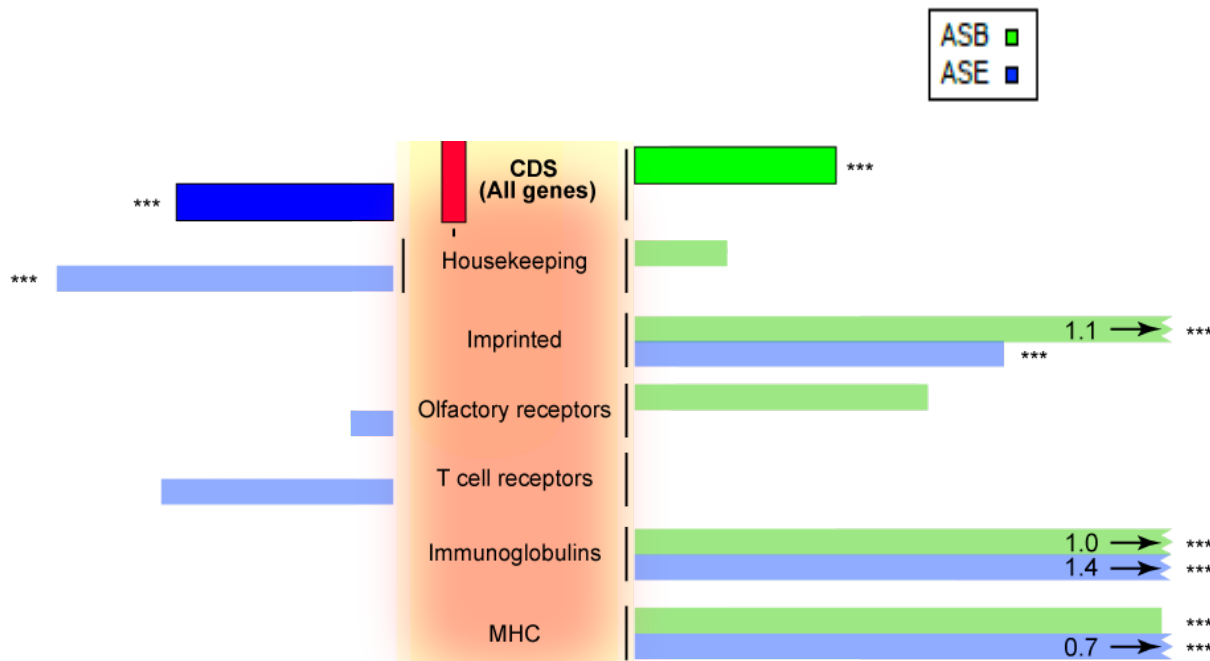Does a particular genomic element have a higher tendency to be allele-specific?
Fisher's exact test, for the **enrichment** of allele-specific variants in the element (with respect to non-allele-specific variants that could potentially be called as allelic)



[Chen *et al.* ('16) *Nat. Comm.*]

# Groups of elements that are enriched or depleted in allelic activity

[Chen *et al.* ('16) *Nat. Comm.*]

# Groups of elements that are enriched or depleted in allelic activity



ASB ■ (green)
ASE ■ (blue)

CDS (All genes) — ASB ***, ASE ***
Housekeeping
Imprinted — 1.1 → ***, ***
Olfactory receptors
T cell receptors
Immunoglobulins — 1.0 → ***, 1.4 → ***
MHC — ***, 0.7 → ***

- <u>SNURF</u> imprinted and implicated in Prader-Willi/Angelman Syndrome
- <u>KCNQ1</u> is an imprinted gene

# Analysis of Personal Genomes:
## Evaluating the impact of variants in exomes using protein structure & allelic activity

- Introduction
  - Rare v common variants
  - The importance of interpreting rare coding variants in the context of disease genomics (CMG,TCGA)
- Identifying cryptic allosteric sites with STRESS
  - On surface & in interior bottlenecks
- Using changes in localized frustration to find further sites sensitive to mutations
  - Difference betw. TSGs & oncogenes

- Using structural motifs (eg TPR) for intensification of weak population genetic signals
  - For both negative and positive selection
- Prioritizing allelic genes using AlleleDB
  - Having observed difference in molecular activity in many contexts

Analysis of Personal Genomes:
Evaluating the impact of variants in exomes
using protein structure & allelic activity

- Introduction
  - Rare v common variants
  - The importance of interpreting rare coding variants in the context of disease genomics (CMG,TCGA)
- Identifying cryptic allosteric sites with STRESS
  - On surface & in interior bottlenecks
- Using changes in localized frustration to find further sites sensitive to mutations
  - Difference betw. TSGs & oncogenes

- Using structural motifs (eg TPR) for intensification of weak population genetic signals
  - For both negative and positive selection
- Prioritizing allelic genes using AlleleDB
  - Having observed difference in molecular activity in many contexts

# 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