

Personal Genomics:

**Handling
Exponential Data
Scaling through
Prioritizing High-
impact Variants**

Mark Gerstein
Yale

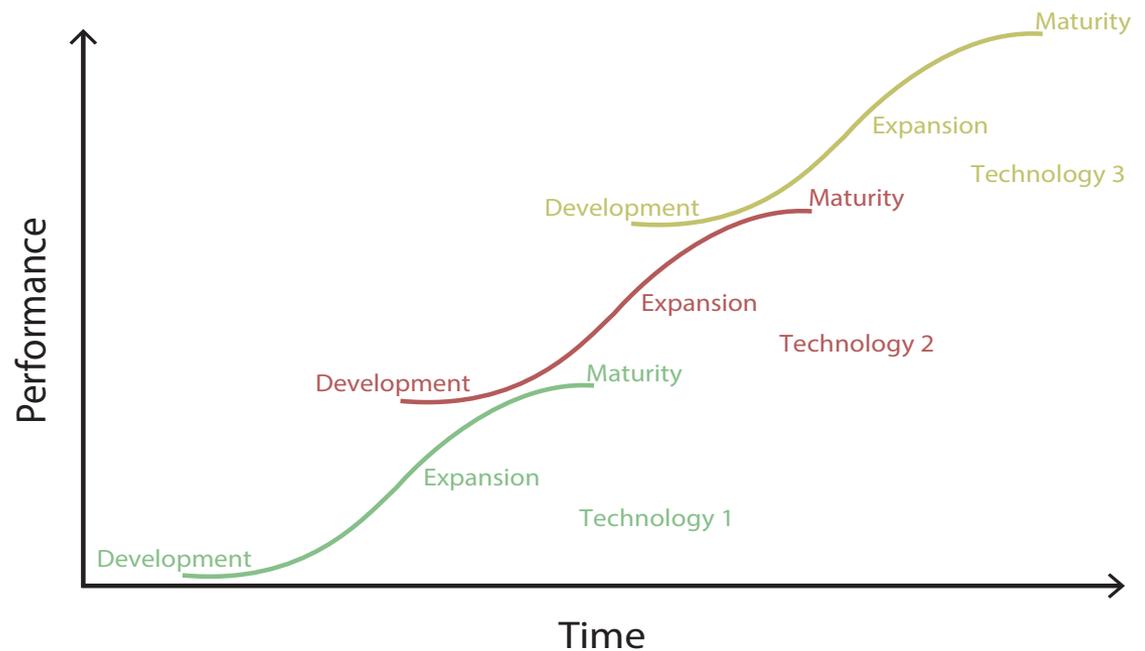
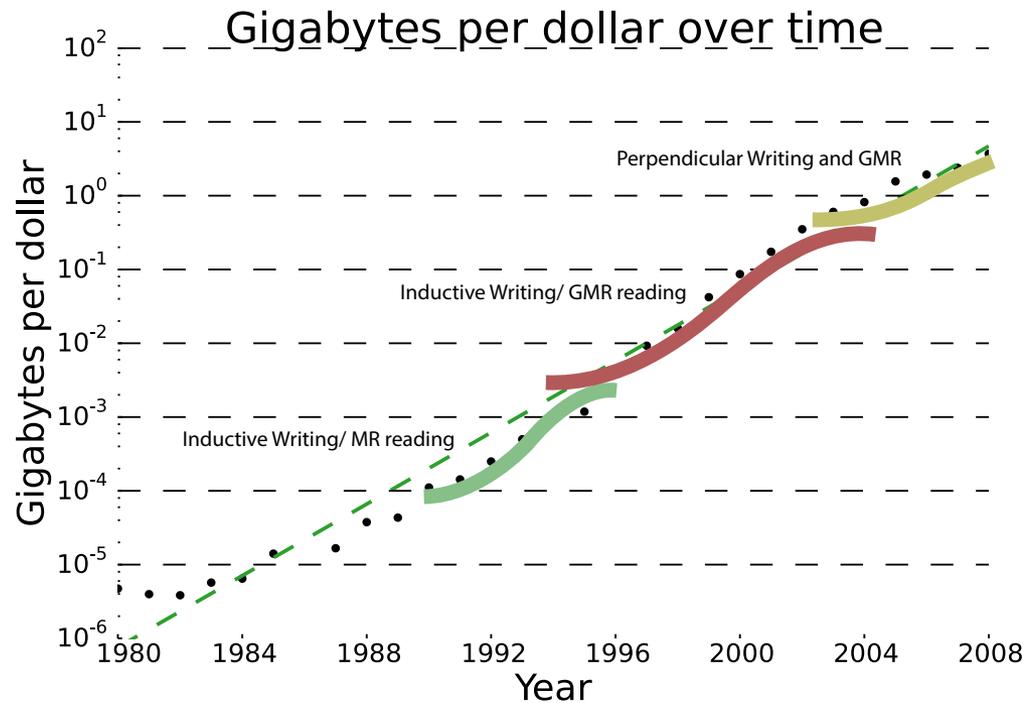
Slides freely
downloadable from
Lectures.GersteinLab.org
& “tweetable”
(via [@markgerstein](https://twitter.com/markgerstein)).
See last slide for more info.



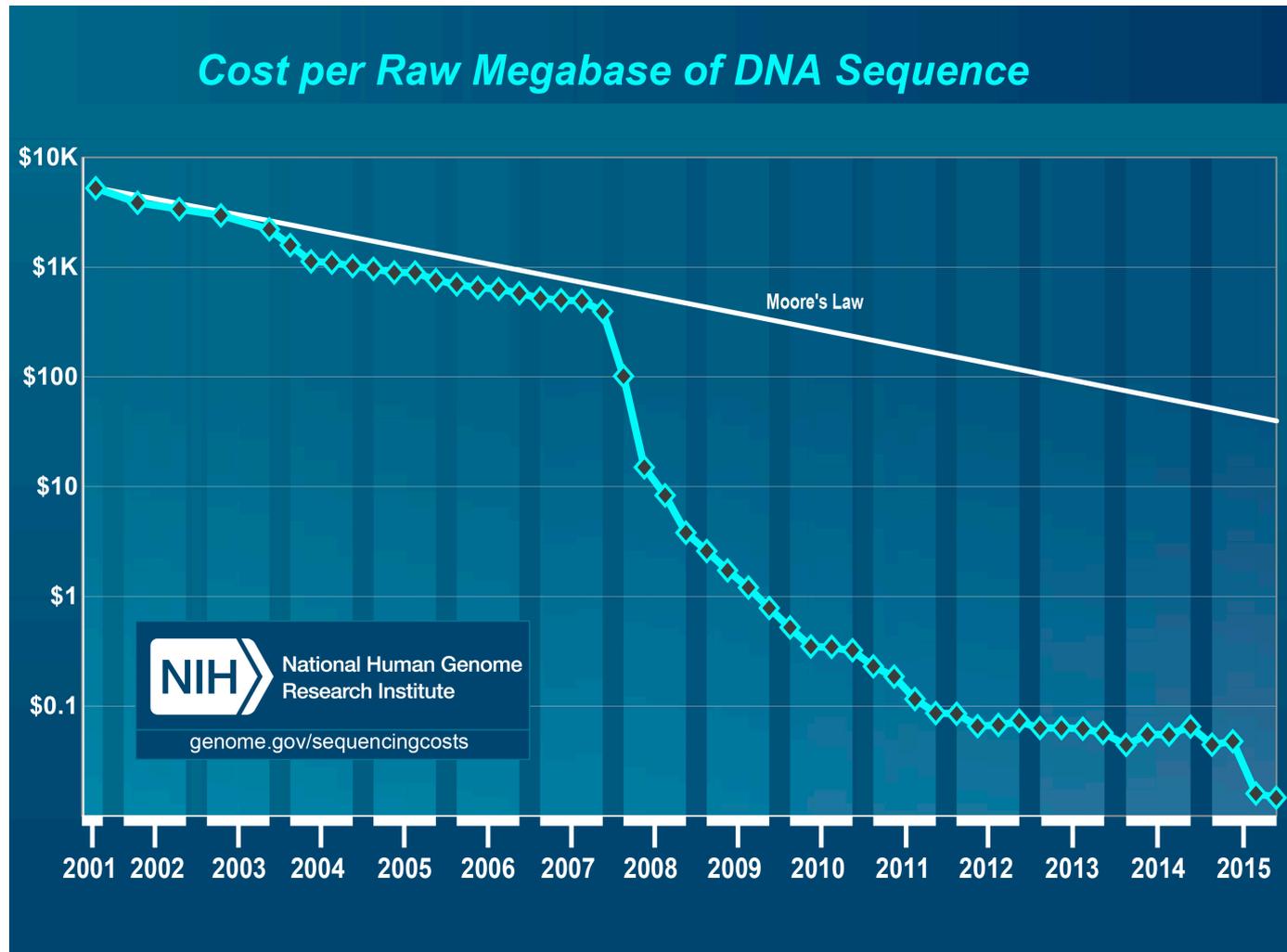
Moore's Law: Exponential Scaling of Computer Technology

– 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 Moore's law is a superposition of s-curves for different technologies



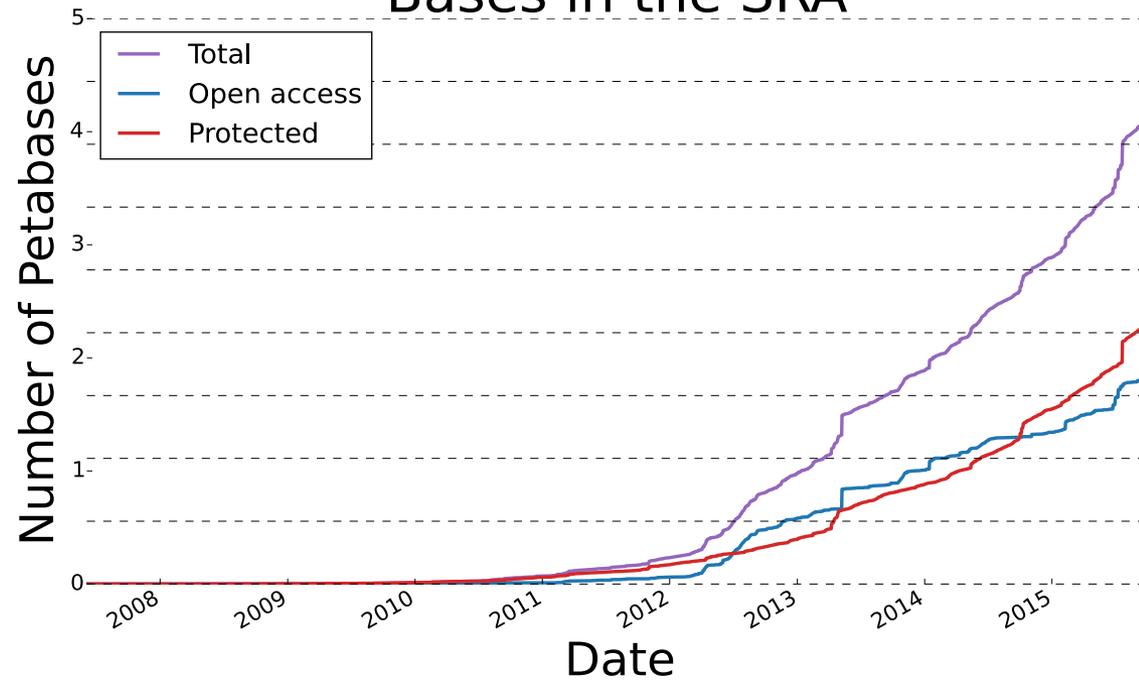
Sequencing Data Explosion: Faster than Moore's Law for a Time (or a S-curve)



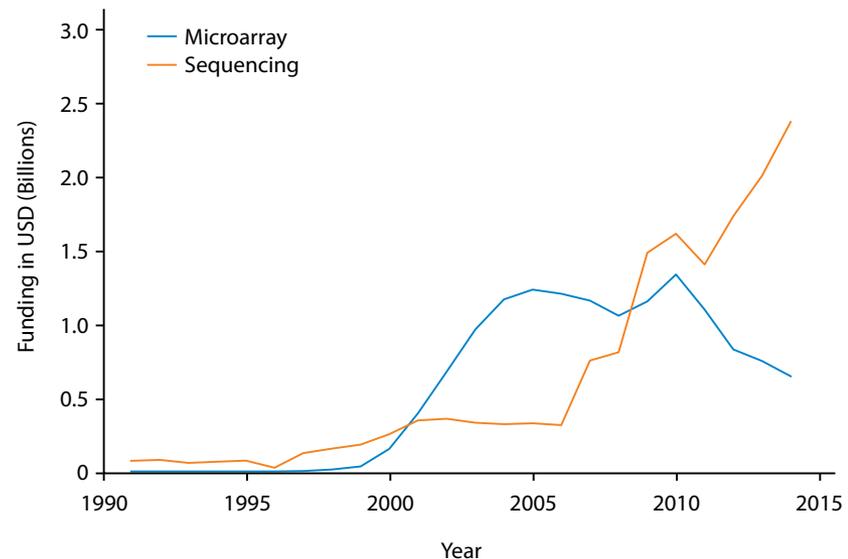
Sequencing cost reductions have resulted in an explosion of data

- 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



National Institutes of Health funding for 'microarray' and 'sequencing' projects



Seq Universe

[from Heidi Sofia, NHGRI, '14]

TCGA endpoint: ~2.5 Petabytes

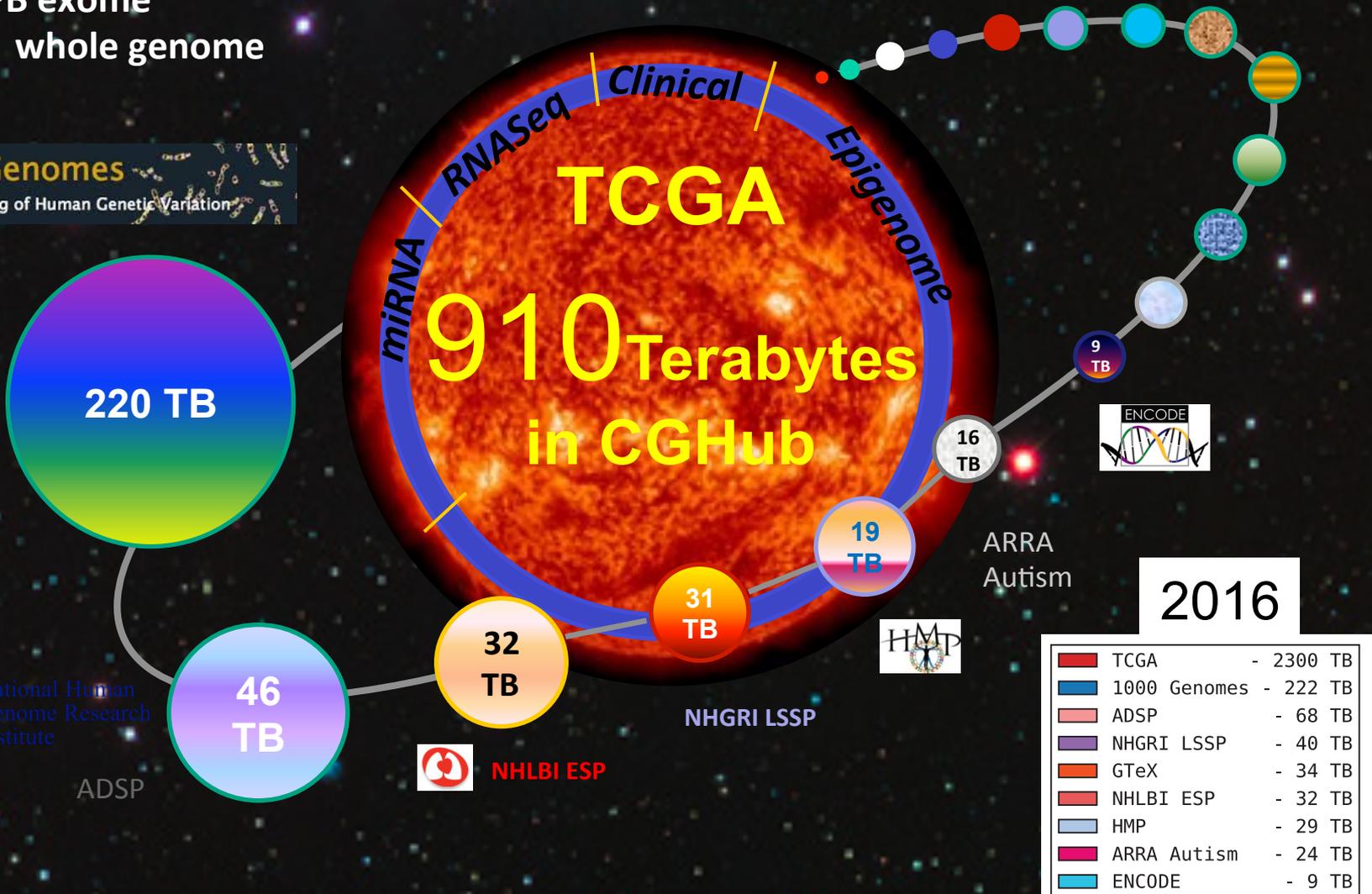
~1.5 PB exome

~1 PB whole genome

SRA >1 petabyte

1000 Genomes

A Deep Catalog of Human Genetic Variation



ADSP



NHLBI ESP

NHGRI LSSP

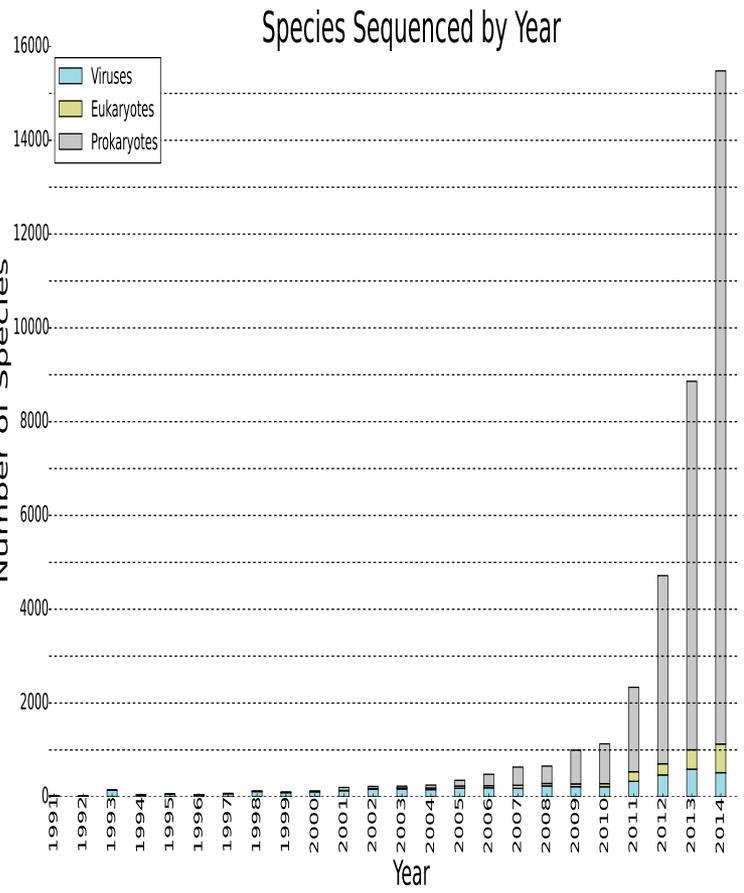
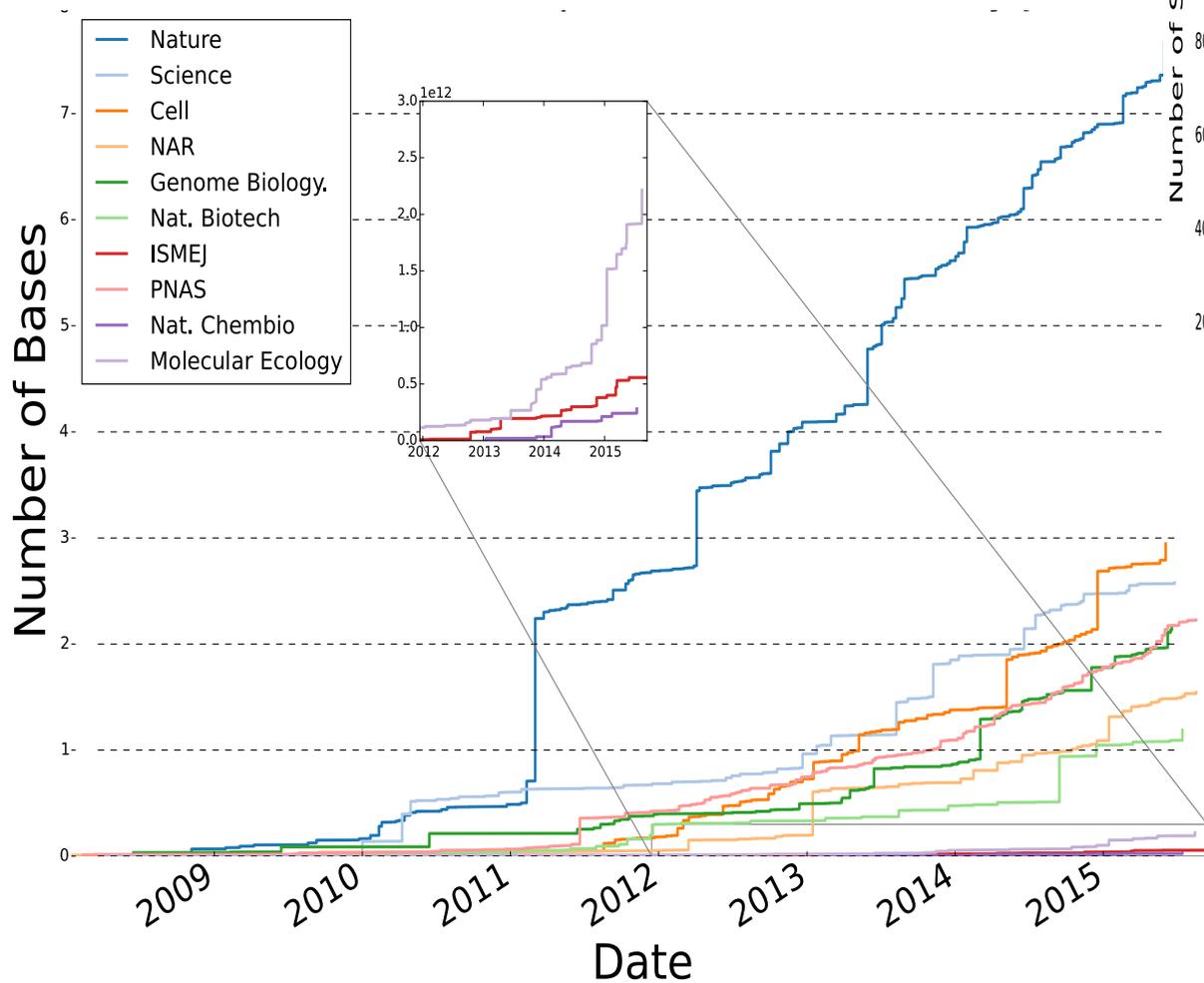


ARRA Autism

2016

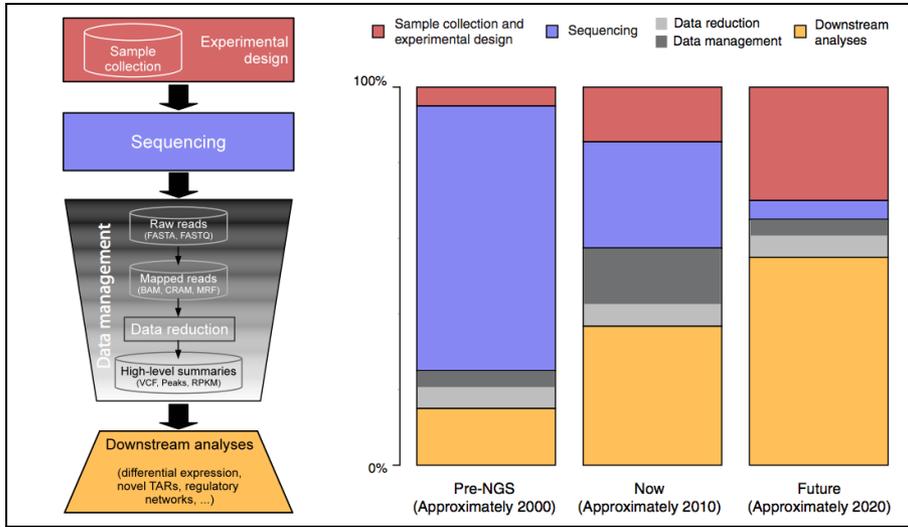
TCGA	- 2300 TB
1000 Genomes	- 222 TB
ADSP	- 68 TB
NHGRI LSSP	- 40 TB
GTeX	- 34 TB
NHLBI ESP	- 32 TB
HMP	- 29 TB
ARRA Autism	- 24 TB
ENCODE	- 9 TB

Increasing diversity in sequence data sources

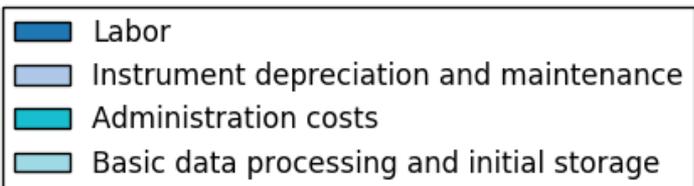
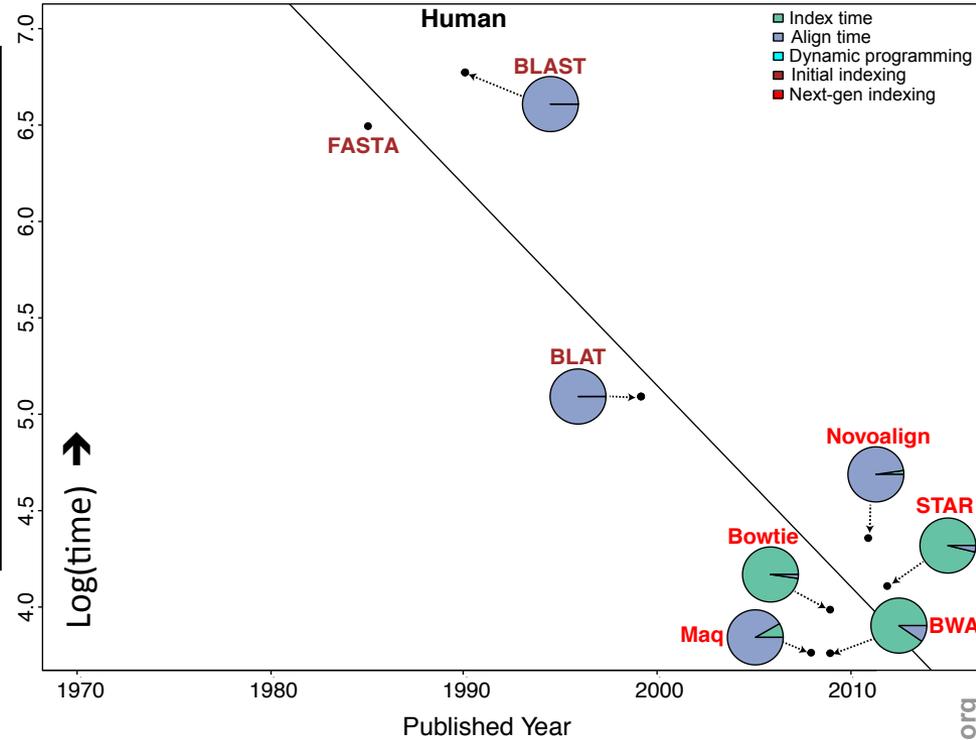


[Muir et al. ('15) GenomeBiol.]

The changing costs of a sequencing pipeline

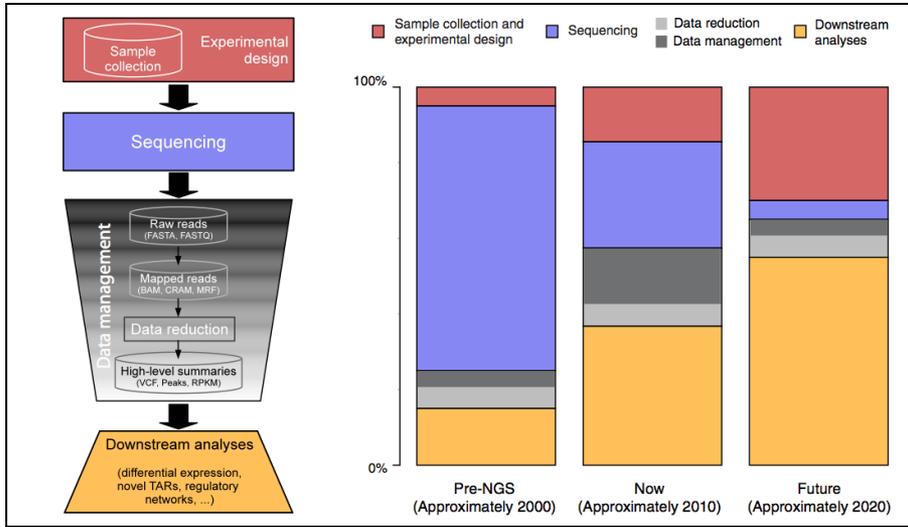


From '00 to ~' 20, cost of DNA sequencing expt. shifts from the actual seq. to sample collection & analysis

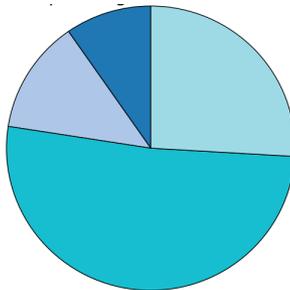
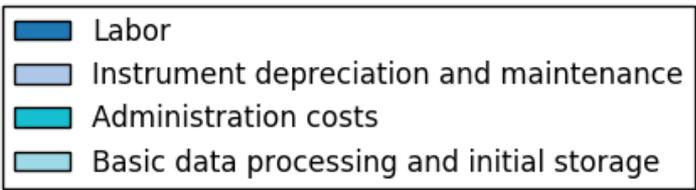
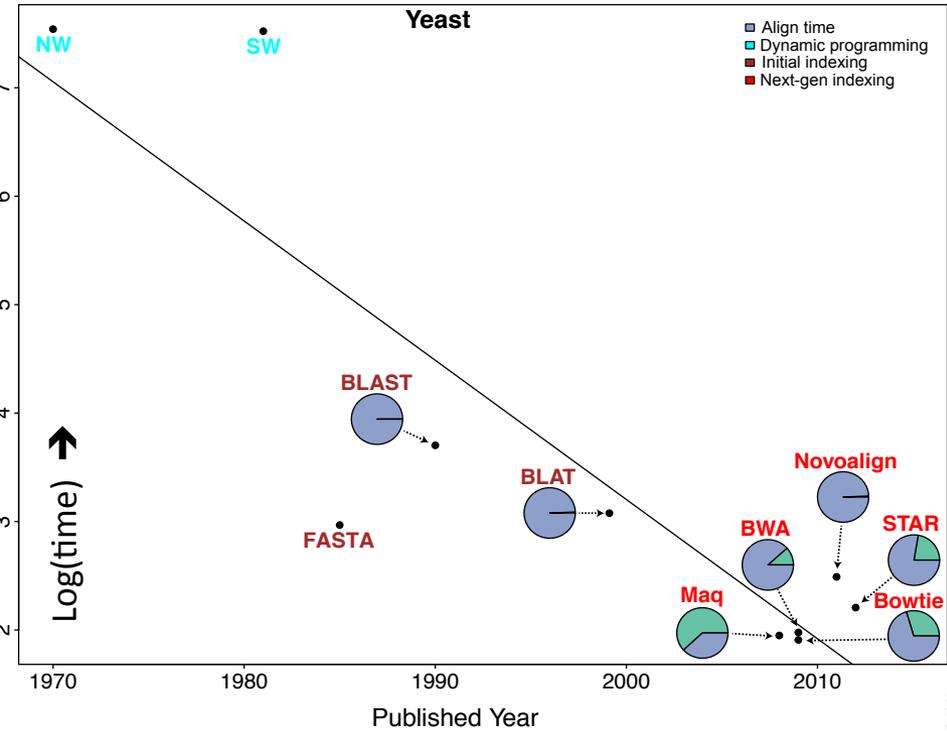


Alignment algorithms scaling to keep pace with data generation

The changing costs of a sequencing pipeline

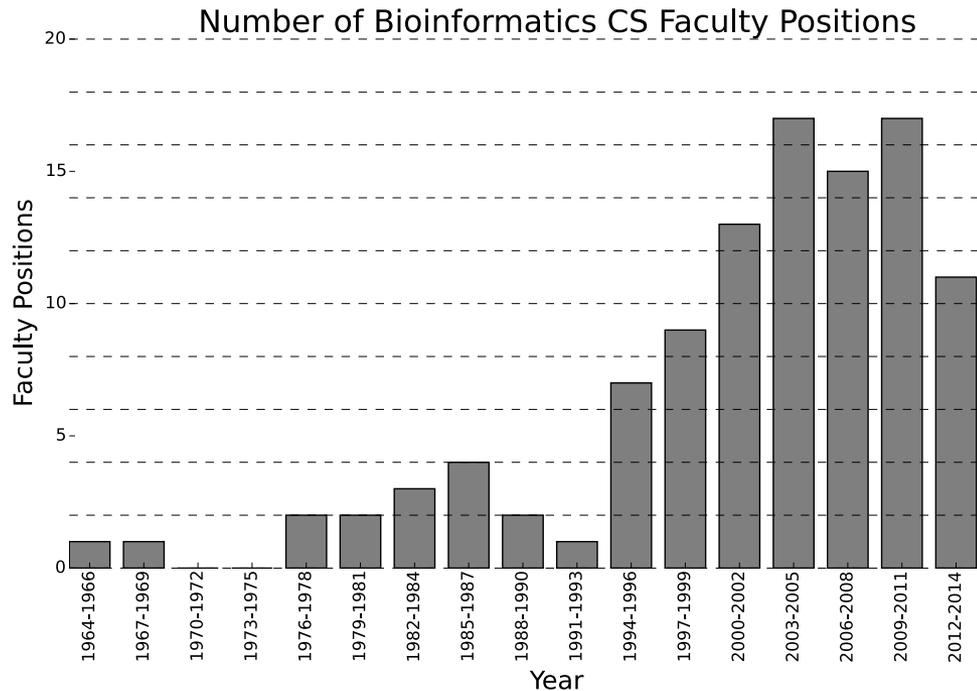


From '00 to ~' 20, cost of DNA sequencing expt. shifts from the actual seq. to sample collection & analysis



Alignment algorithms scaling to keep pace with data generation

Jobs: Bioinformatics is born!



Hires at 51 US universities in 3-year bins

http://jeffhuang.com/computer_science_professors.html

Personal Genomics: Handling Exponential Data Scaling through Prioritizing High-impact Variants

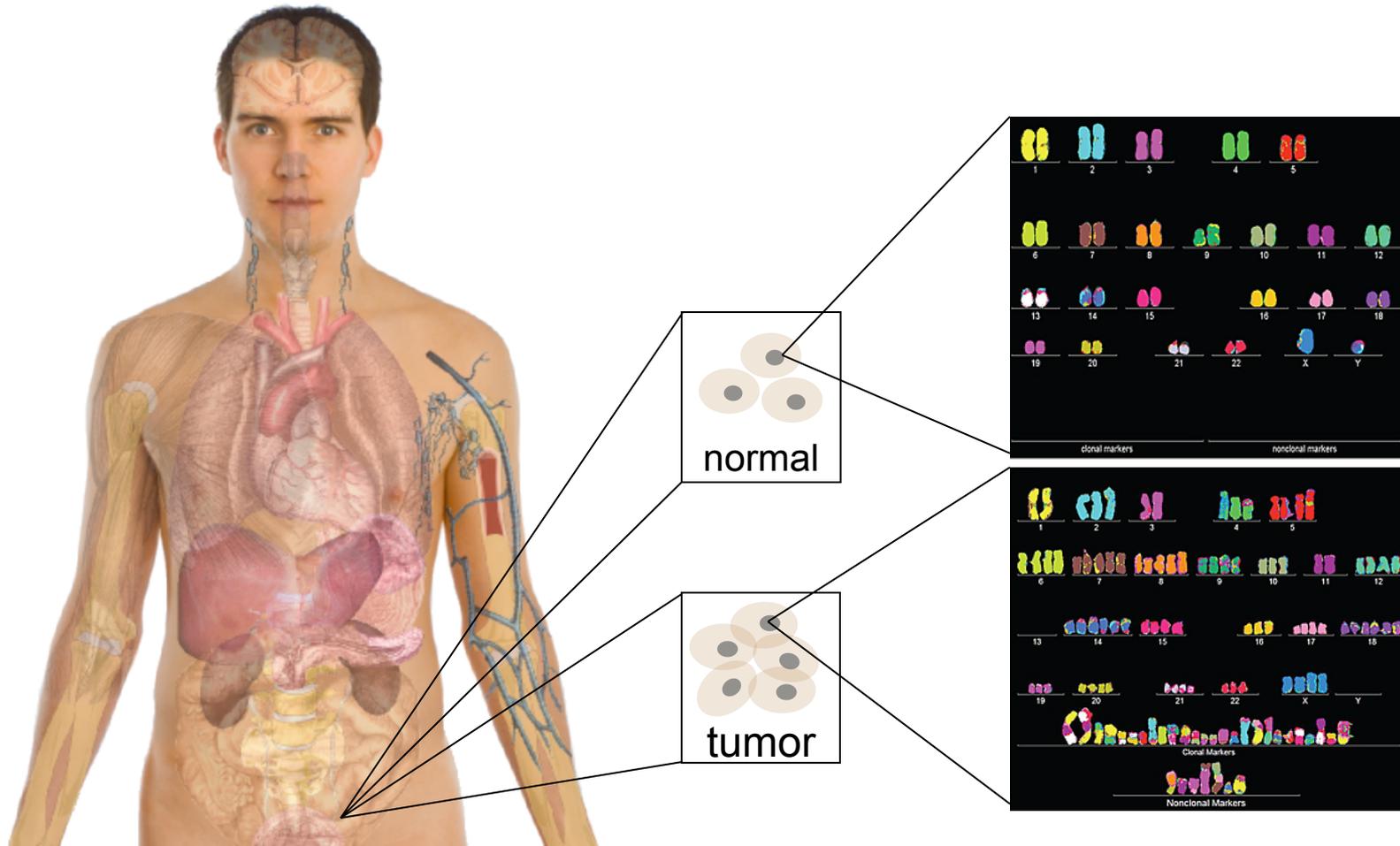
- **Introduction**
 - The exponential scaling of data generation & data processing
 - The landscape of variants in personal genomes
- **Characterizing Rare Variants in Coding Regions**
 - Identifying with **STRESS** cryptic allosteric sites
 - On surface & in interior bottlenecks
- **Non-coding Variants : Prioritizing using AlleleDB in terms of allelic elements**
 - Having observed difference in molecular activity in many contexts
- **Putting it together in workflows: Integrating evidence on non-coding variants with FunSeq**
 - Systematically weighting all the features
 - suggesting non-coding drivers
 - Prioritizing rare germline variants

Personal Genomics: Handling Exponential Data Scaling through Prioritizing High-impact Variants

- Introduction
 - The **exponential scaling** of data generation & data processing
 - The **landscape of variants** in personal genomes
- Characterizing Rare Variants in Coding Regions
 - Identifying with STRESS **cryptic allosteric sites**
 - On surface & in interior bottlenecks
- Non-coding Variants :
Prioritizing using AlleleDB in terms of **allelic elements**
 - Having observed difference in molecular activity in many contexts
- Putting it together in workflows:
Integrating evidence on non-coding variants with FunSeq
 - Systematically weighting all the features
 - suggesting non-coding drivers
 - Prioritizing rare germline variants

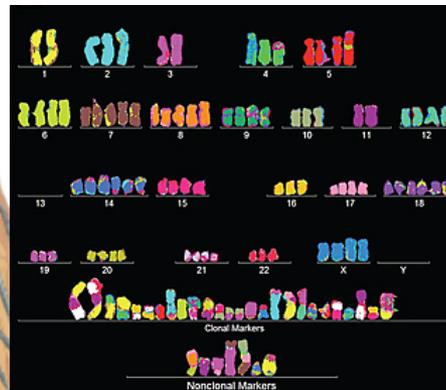
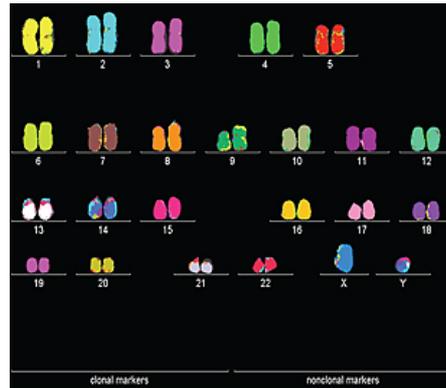
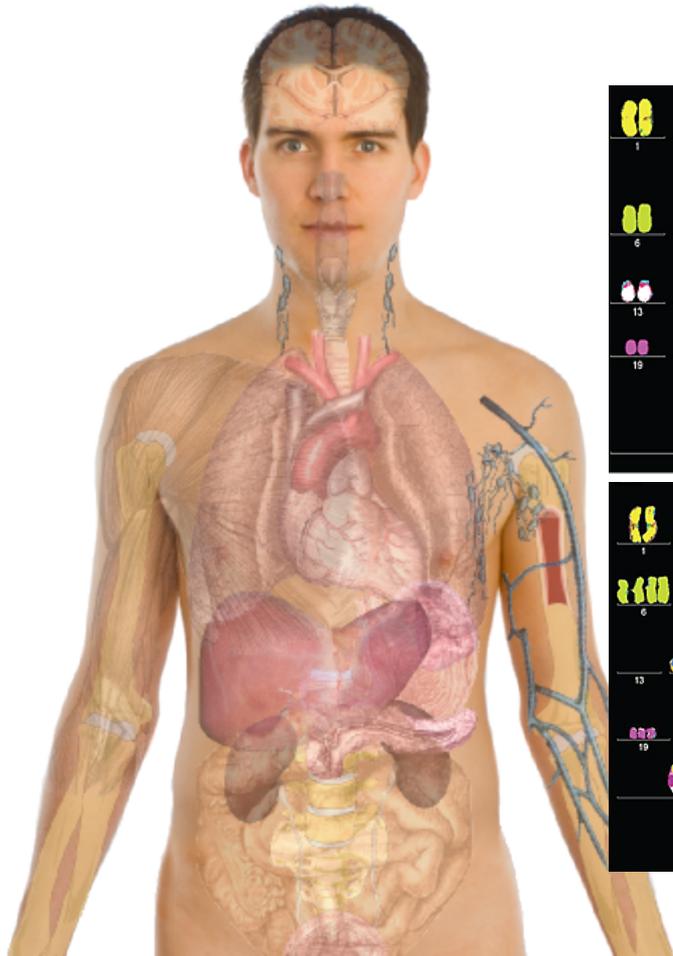
Personal Genomics as a Gateway into Biology

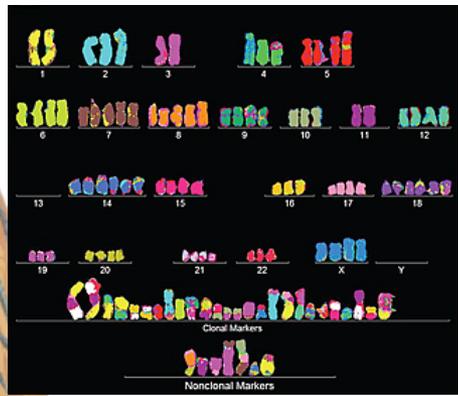
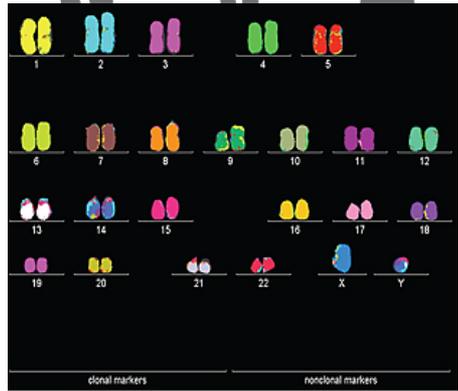
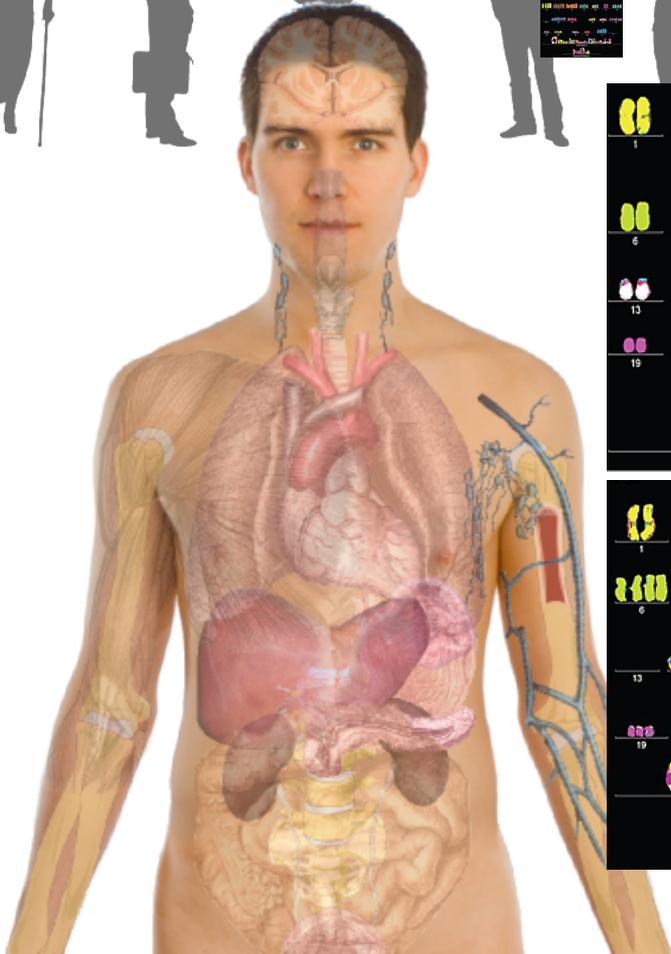
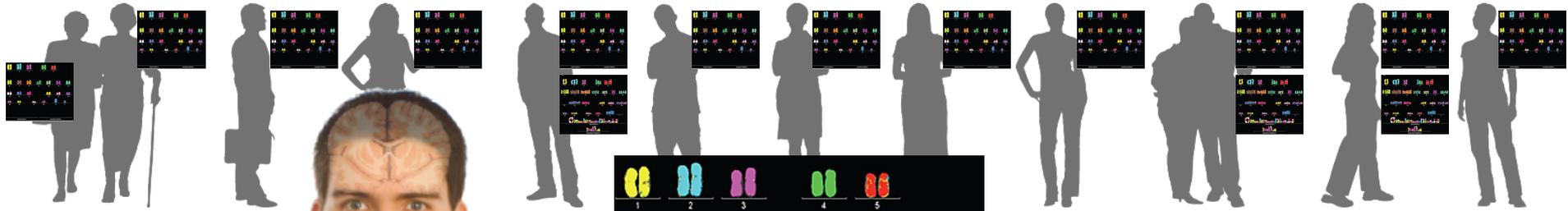
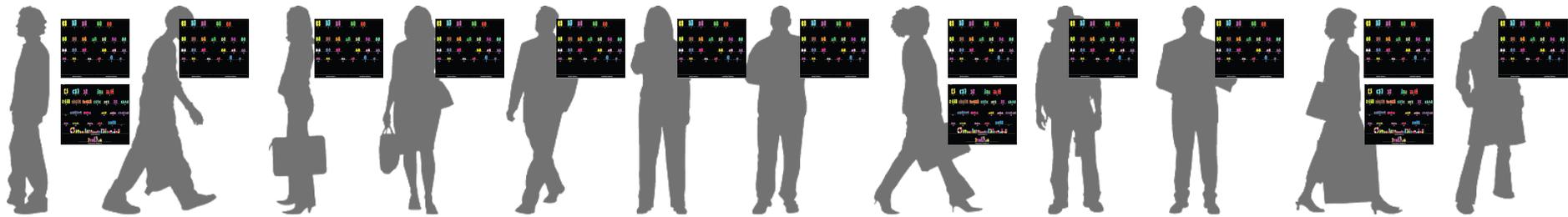
Personal genomes soon will become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.



Personal Genomics as a Gateway into Biology

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





Human Genetic Variation

A Cancer Genome



A Typical Genome

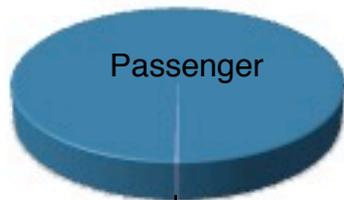


Population of 2,504 people



Origin of Variants

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



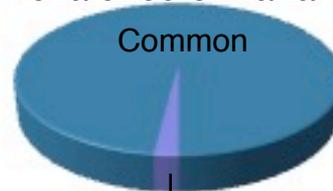
Passenger

Driver (~0.1%)

Class of Variants

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

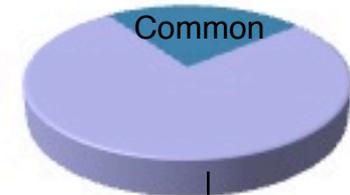
Prevalence of Variants



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.

CAN YOU FIND THE PANDA?



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.

CAN YOU FIND THE PANDA?



Finding Key Variants

Somatic

• Overall

- Often these can be conceptualized as very rare variants
- A challenge to identify somatic mutations contributing to cancer is to find driver mutations & distinguish them from passengers.

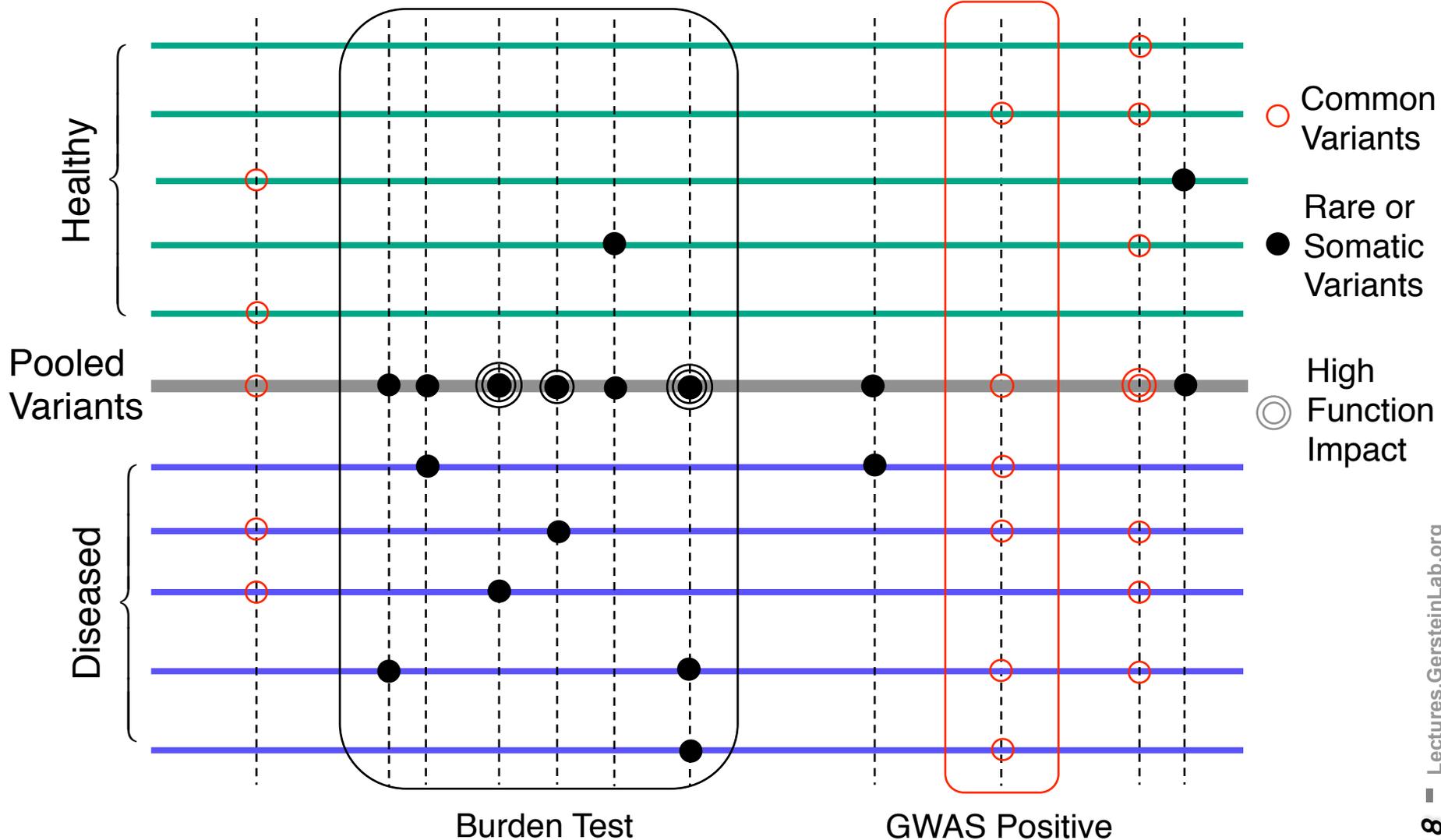
• Drivers

- Driver mutation is a mutation that directly or indirectly confers a selective growth advantage to the cell in which it occurs.
- A typical tumor contains 2-8 drivers; the remaining mutations are passengers.

• Passengers

- Conceptually, a passenger mutation has no direct or indirect effect on the selective growth advantage of the cell in which it occurred.

Association of Variants with Diseases

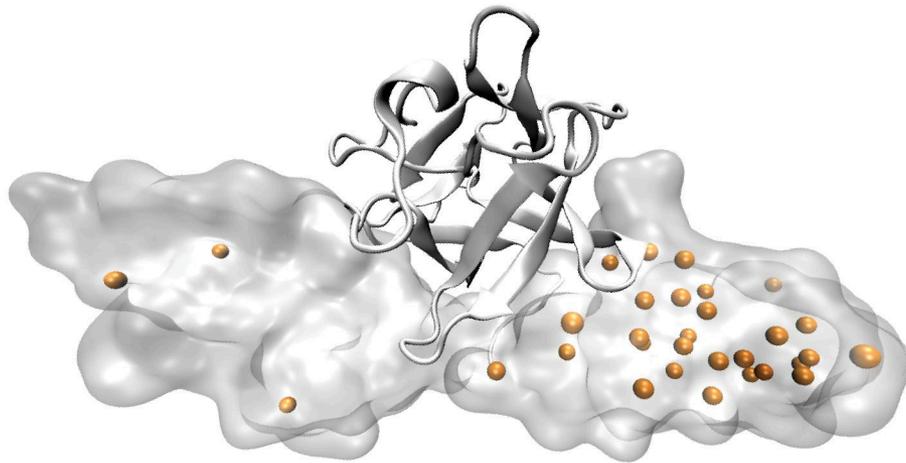


Personal Genomics: Handling Exponential Data Scaling through Prioritizing High-impact Variants

- Introduction
 - The **exponential scaling** of data generation & data processing
 - The **landscape of variants** in personal genomes
- Characterizing Rare Variants in Coding Regions
 - Identifying with STRESS **cryptic allosteric sites**
 - On surface & in interior bottlenecks
- Non-coding Variants : Prioritizing using AlleleDB in terms of **allelic elements**
 - Having observed difference in molecular activity in many contexts
- Putting it together in workflows: **Integrating evidence** on non-coding variants with FunSeq
 - Systematically weighting all the features
 - suggesting non-coding drivers
 - Prioritizing rare germline variants

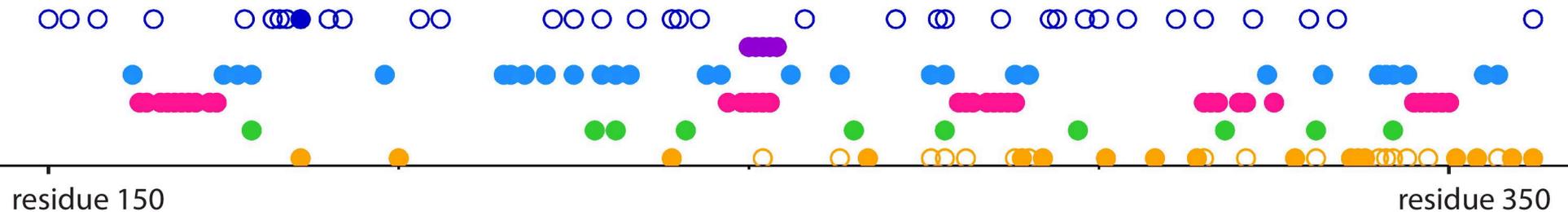
Unlike common SNVs, the statistical power with which we can evaluate rare SNVs in case-control studies is severely limited

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



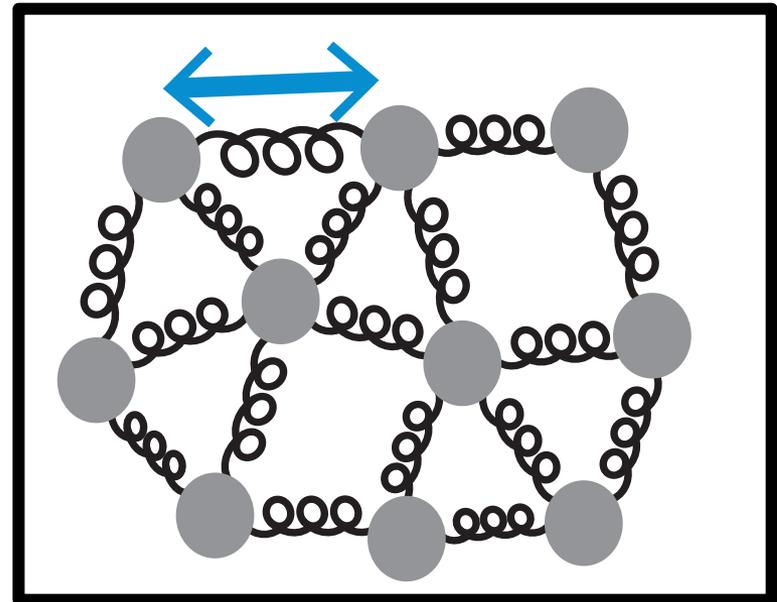
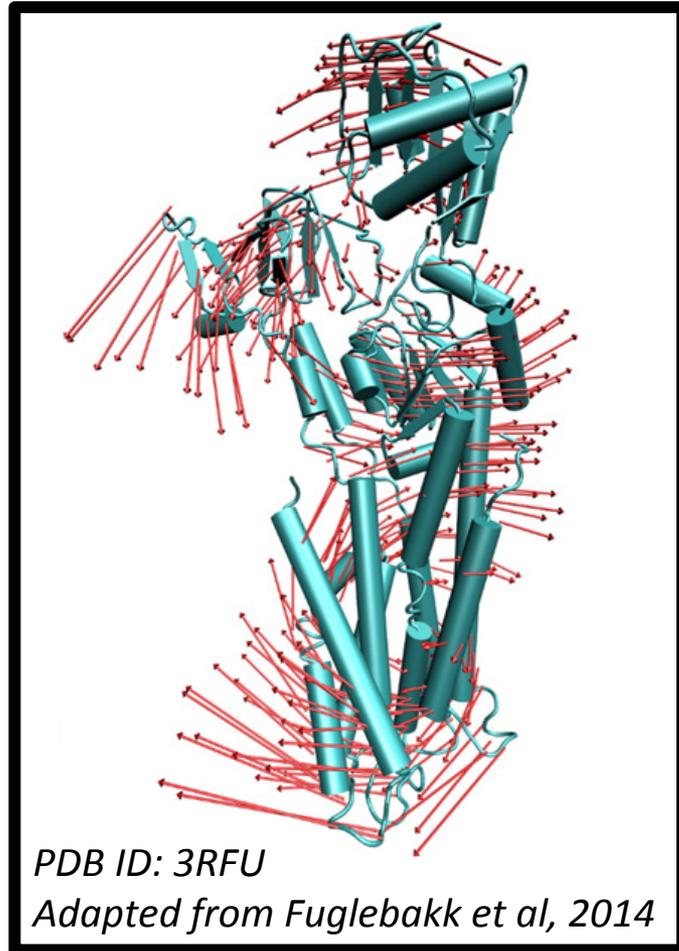
- ○ 1000G & ExAC SNVs (common | rare)
- Hinge residues
- Buried residues
- Protein-protein interaction site
- Post-translational modifications
- HGMD site (w/o annotation overlap)
- HGMD site (w/annotation overlap)

Fibroblast growth factor receptor 2 (pdb: 1IIL)



Models of Protein Conformational Change

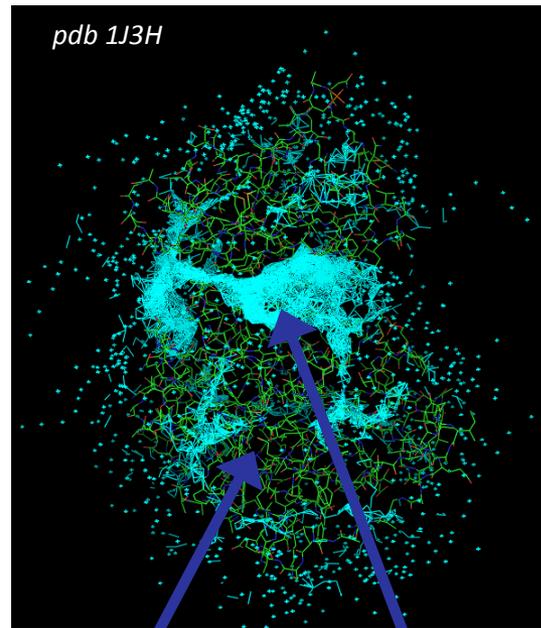
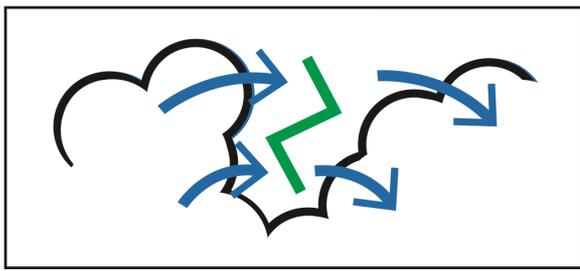
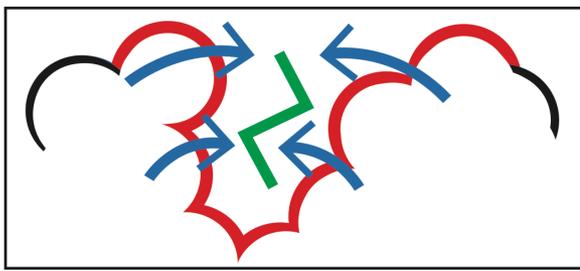
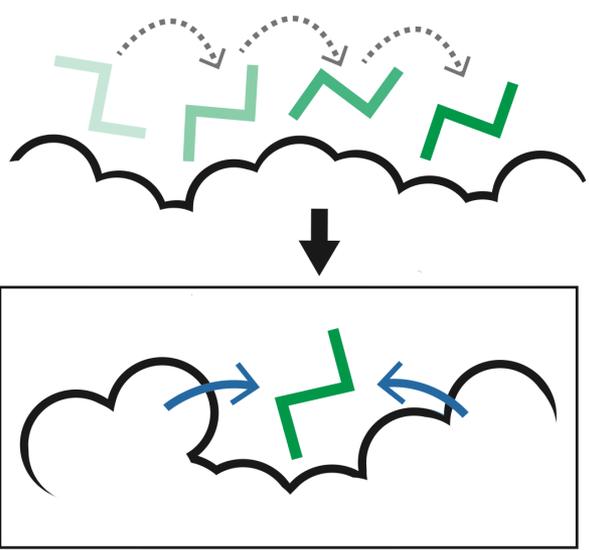
Motion Vectors from Normal Modes (ANMs)



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



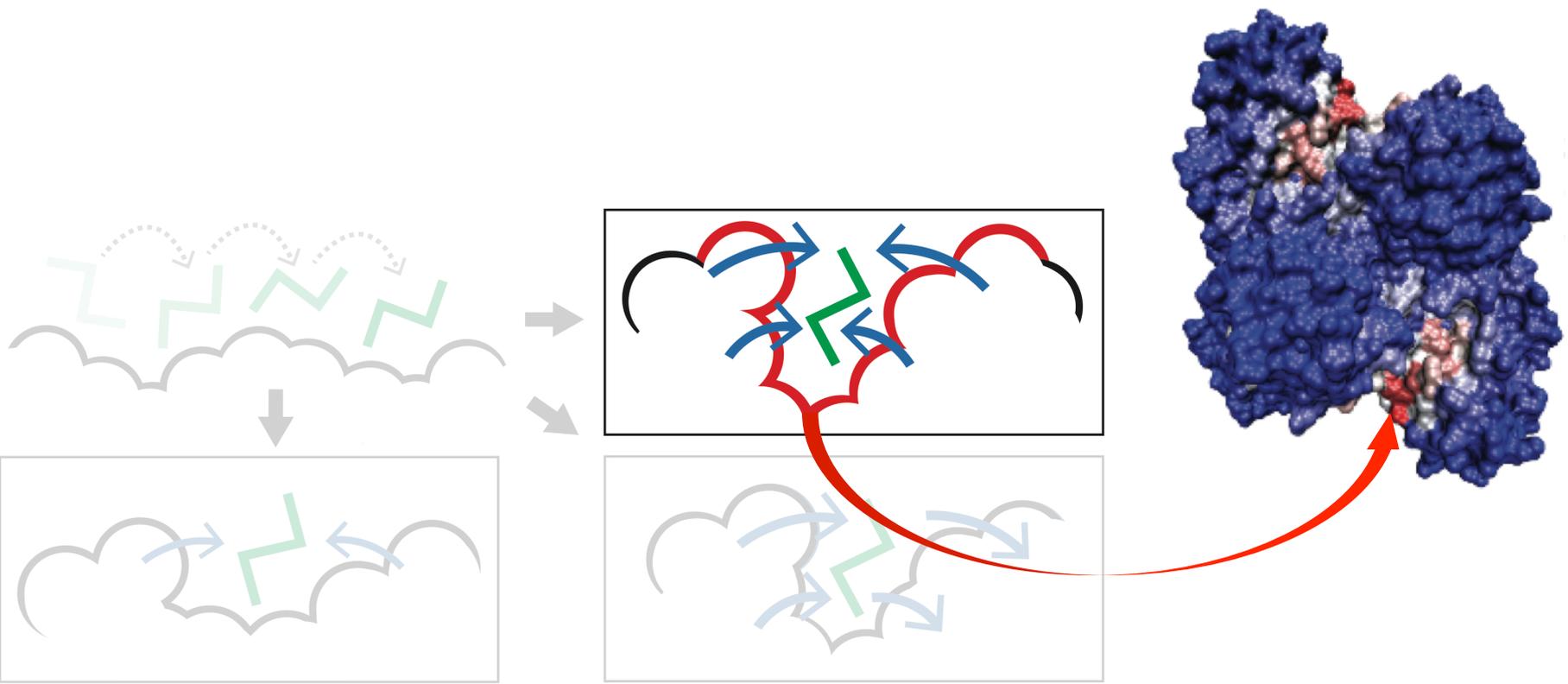
$$\text{binding leverage} = \sum_{m=1}^{10} (\sum_i \sum_j \Delta d_{ij(m)}^2)$$

Surface region with low density of candidate sites

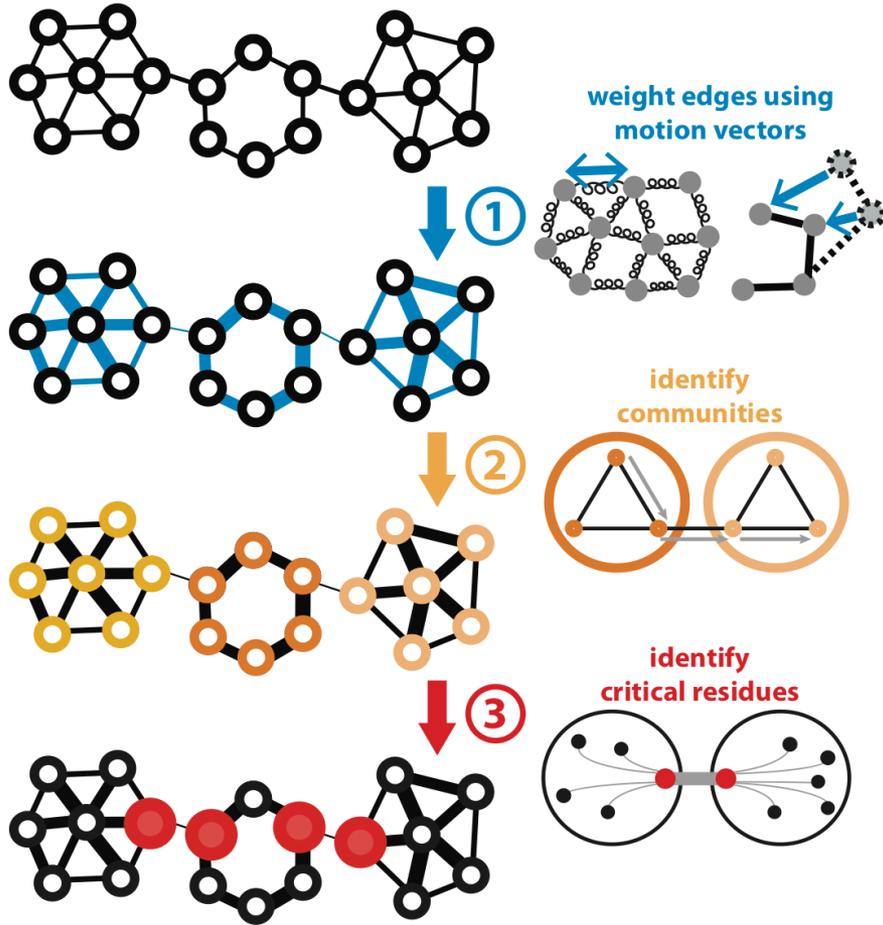
Surface region with high density of candidate sites

Predicting Allosterically-Important Residues at the Surface

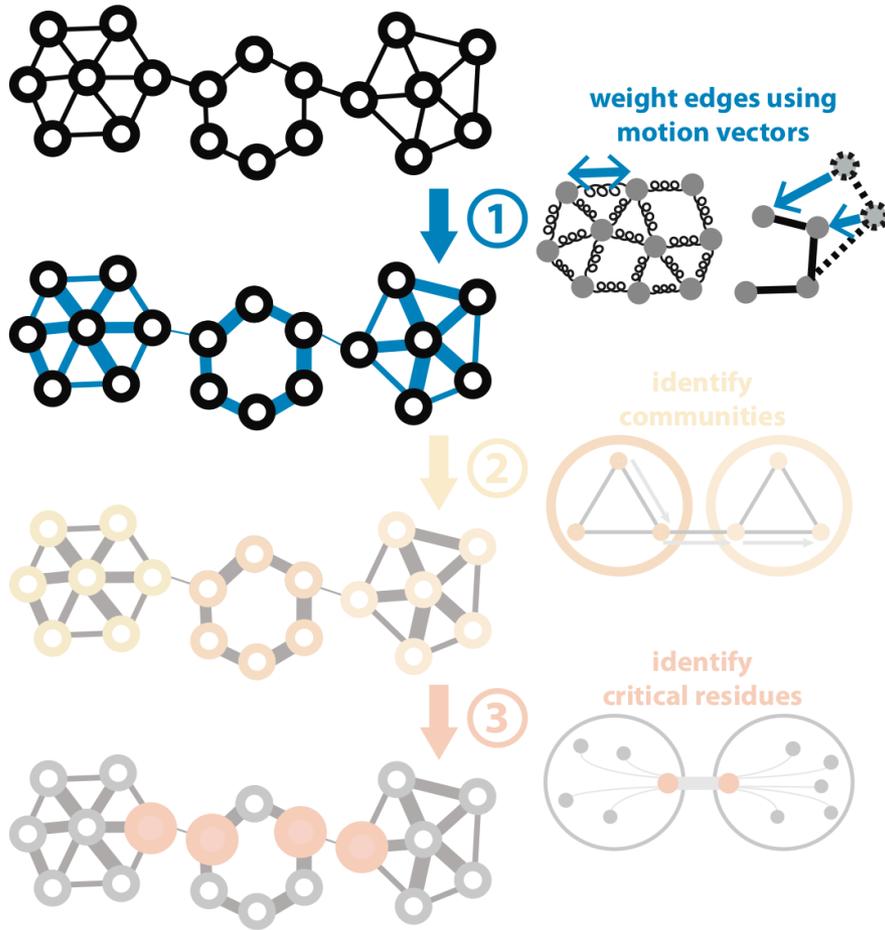
PDB: 3PFK



Predicting Allosterically-Important Residues within the Interior



Predicting Allosterically-Important Residues within the Interior

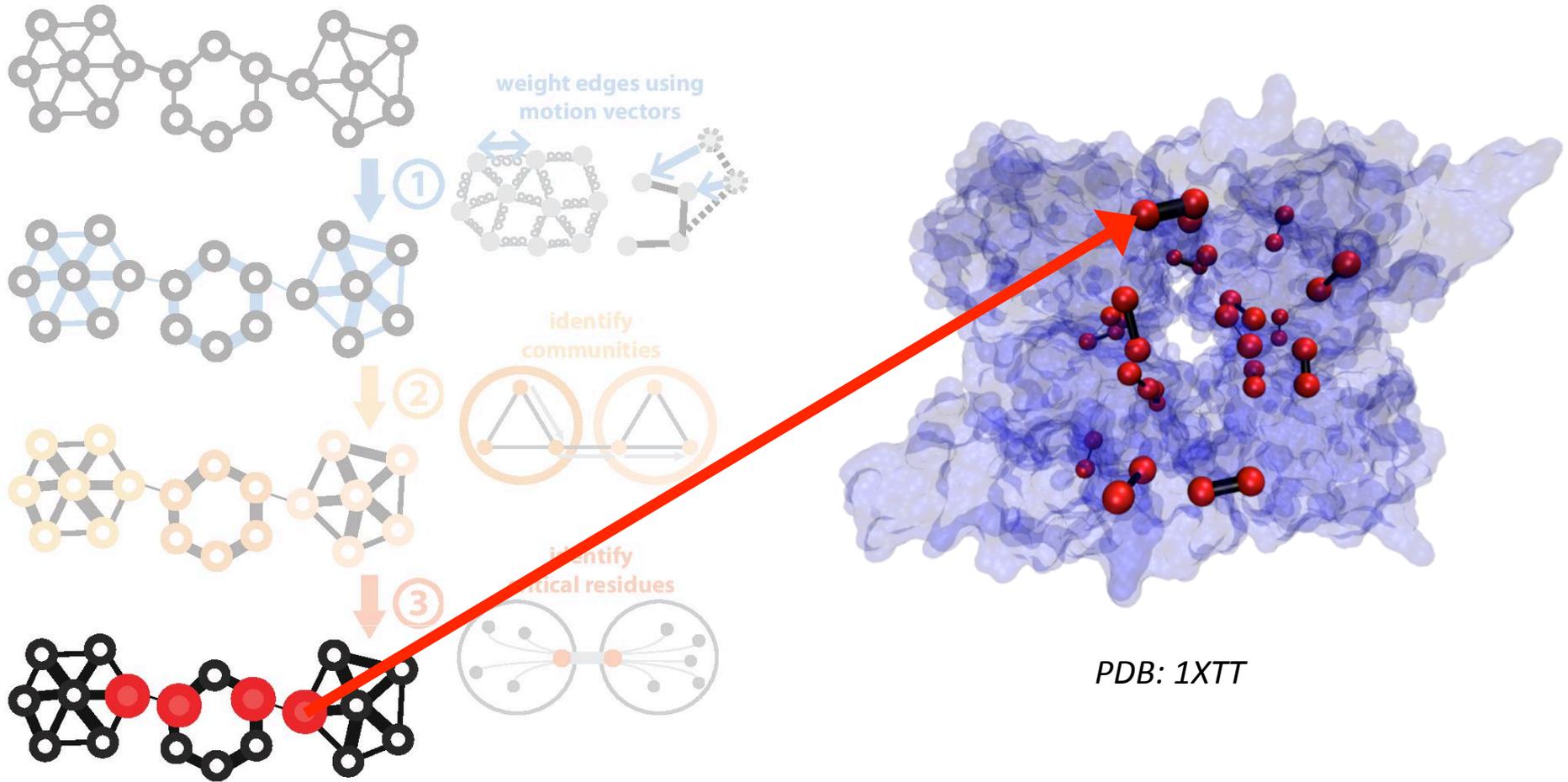


$$Cov_{ij} = \langle \mathbf{r}_i \cdot \mathbf{r}_j \rangle$$

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

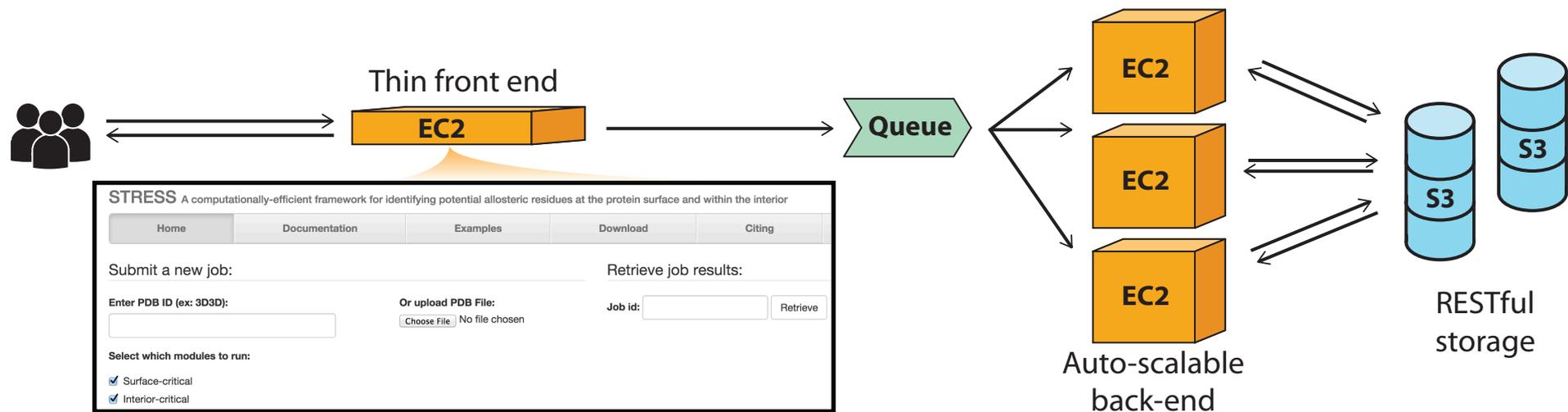
$$D_{ij} = -\log(|C_{ij}|)$$

Predicting Allosterically-Important Residues within the Interior



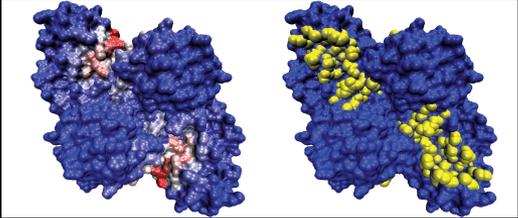
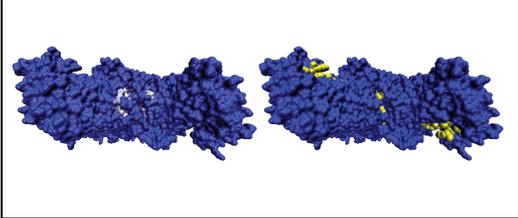
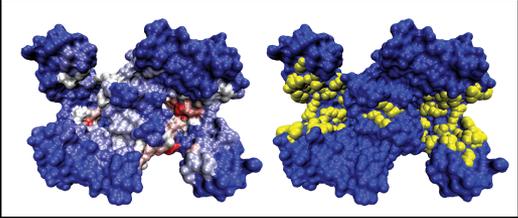
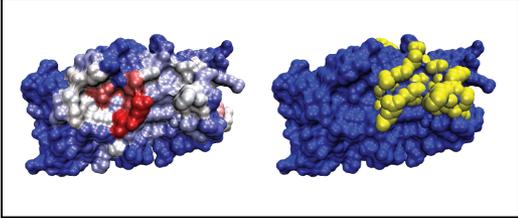
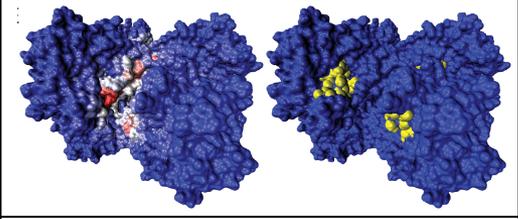
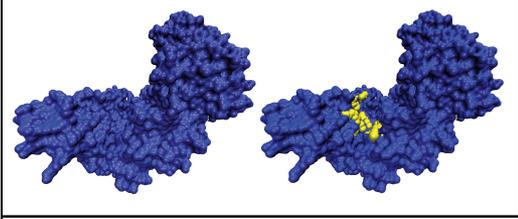
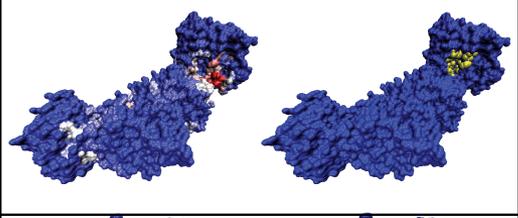
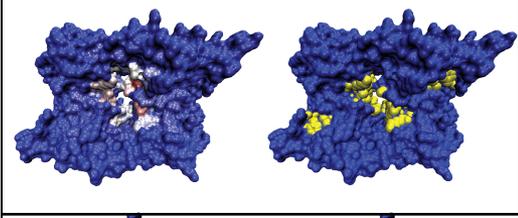
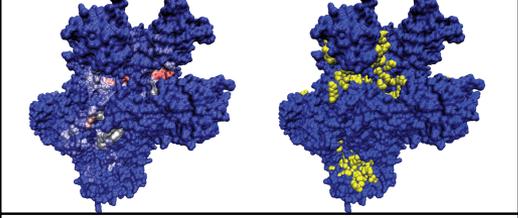
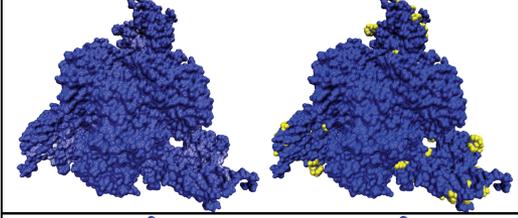
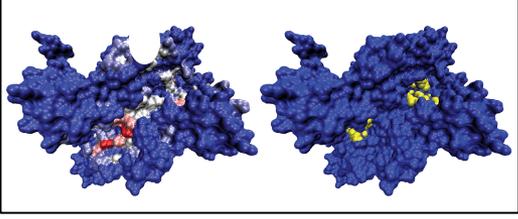
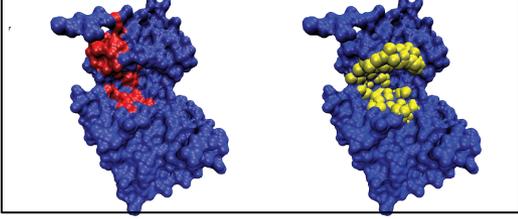
STRESS Server Architecture: Highlights

stress.molmovdb.org



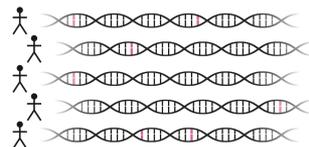
- A light front-end server handles incoming requests, and powerful back-end servers perform calculations.
- Auto Scaling adjusts the number of back-end servers as needed.
- A typical structure takes ~30 minutes on a E5-2660 v3 (2.60GHz) core.
- Input & output (i.e., predicted allosteric residues) are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

Predicting Allosteric Sites at the Surface: Predicted & Known

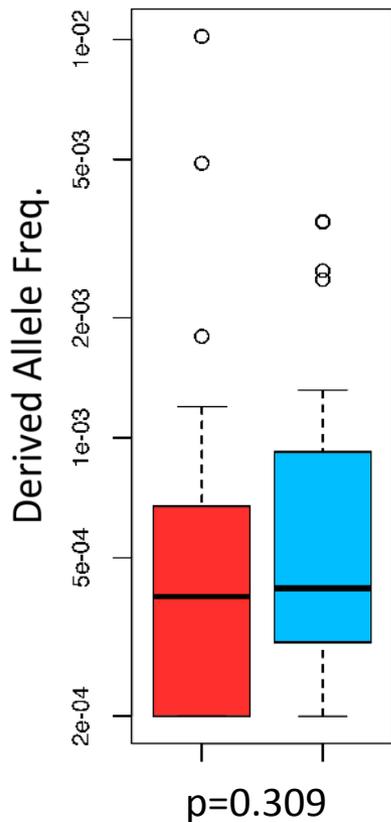
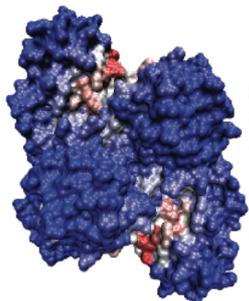
PFK (100)			Malic enz. (0)
ADK (100)			Tyr P-ase (100)
G6P deamin. (50)			Arg. Kin. (0)
Trp Synth. (25)			UPRTase (100)
Glu dehyd. (25)			ATCase (25)
Thr Synth. (67)			PKA (100)

Intra-species conservation of predicted allosteric residues

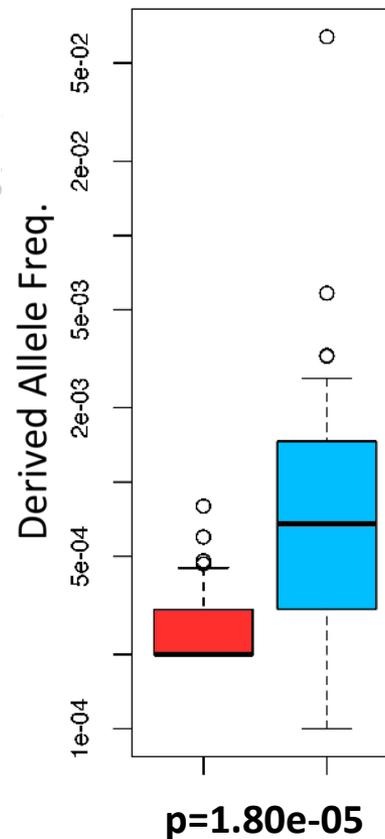
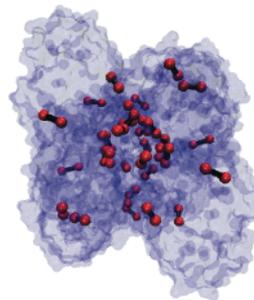
1000 Genomes



Surface



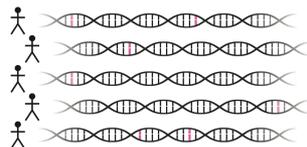
Interior



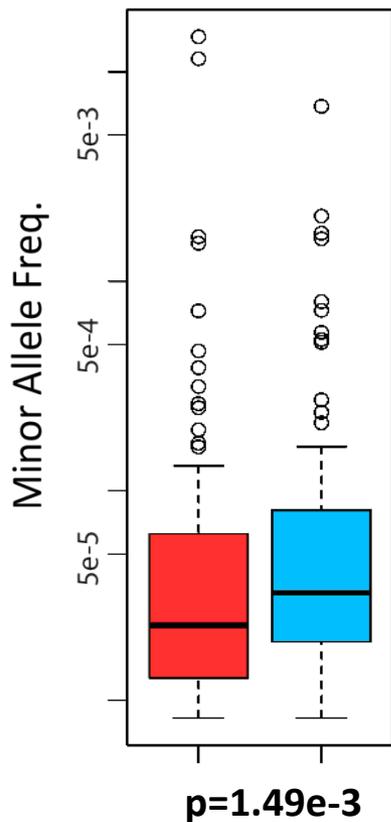
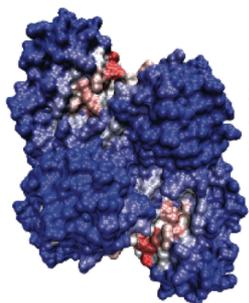
 critical
 non-critical

Intra-species conservation of predicted allosteric residues

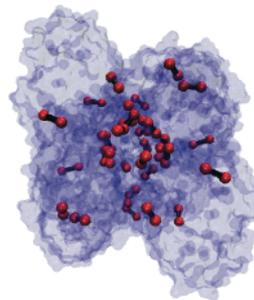
ExAC



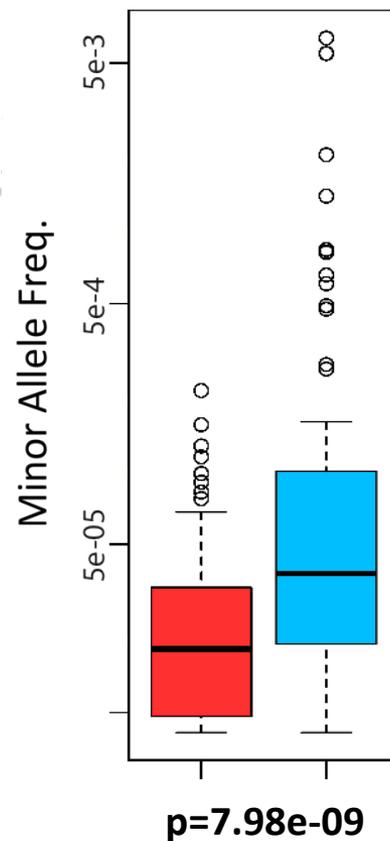
Surface



Interior

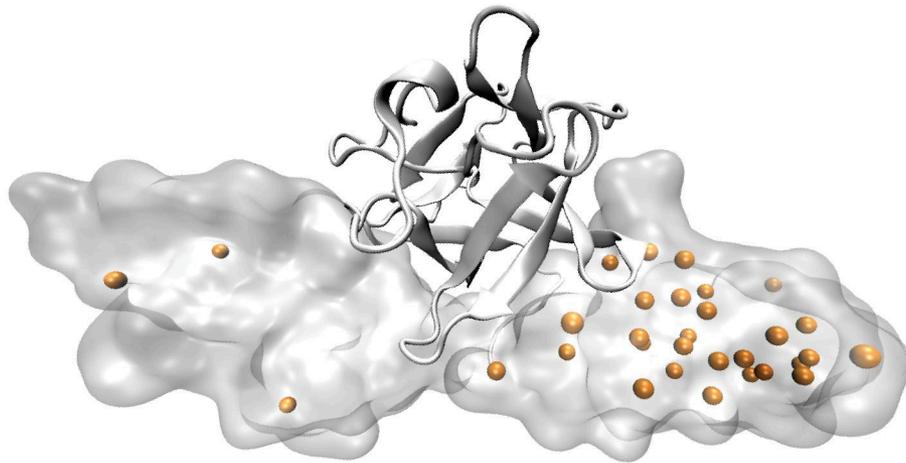


■ critical
■ non-critical



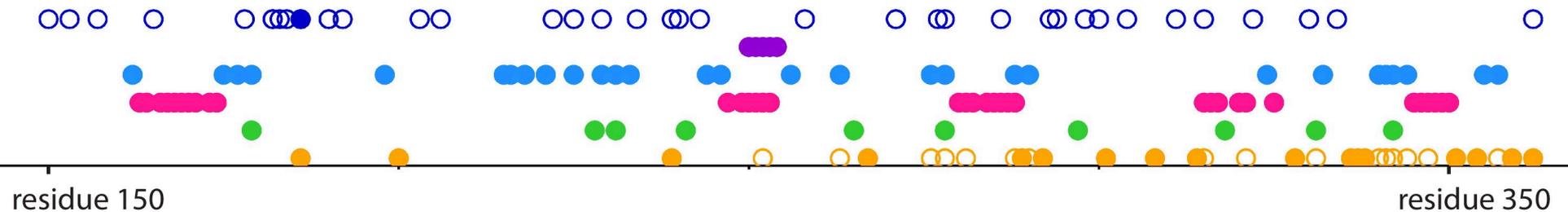
Unlike common SNVs, the statistical power with which we can evaluate rare SNVs in case-control studies is severely limited

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



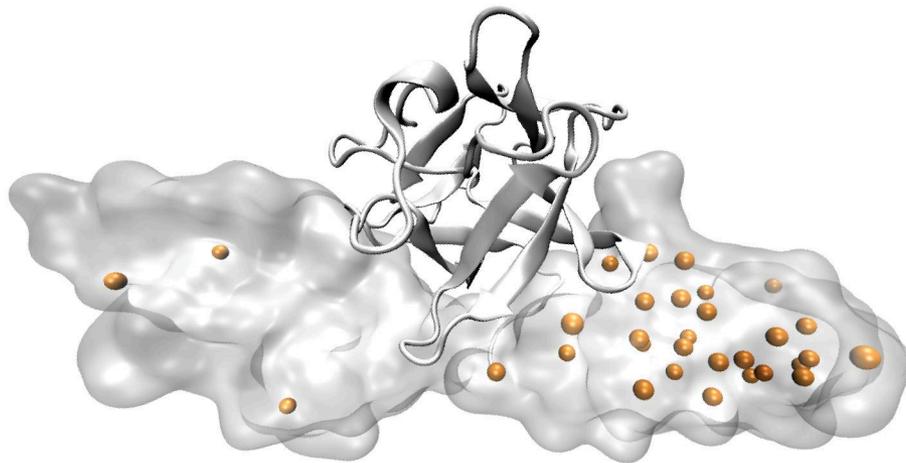
- ○ 1000G & ExAC SNVs (common | rare)
- Hinge residues
- Buried residues
- Protein-protein interaction site
- Post-translational modifications
- HGMD site (w/o annotation overlap)
- HGMD site (w/annotation overlap)

Fibroblast growth factor receptor 2 (pdb: 1IIL)



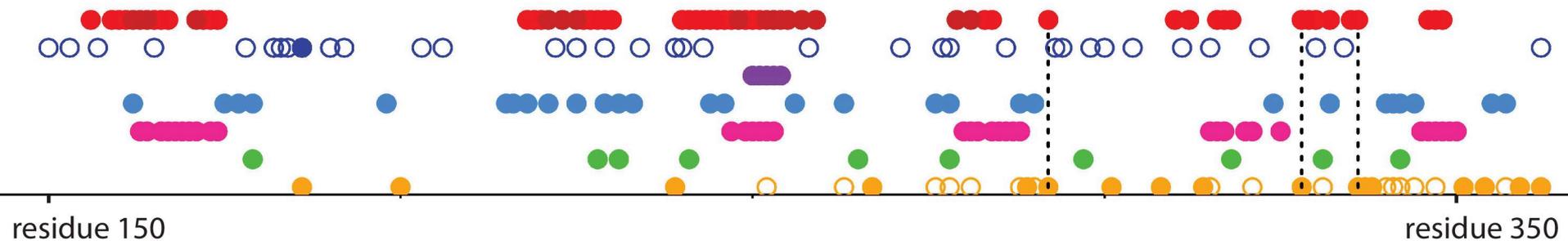
Protein structures may provide the needed alternative for evaluating rare SNVs, many of which may be disease-associated

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)

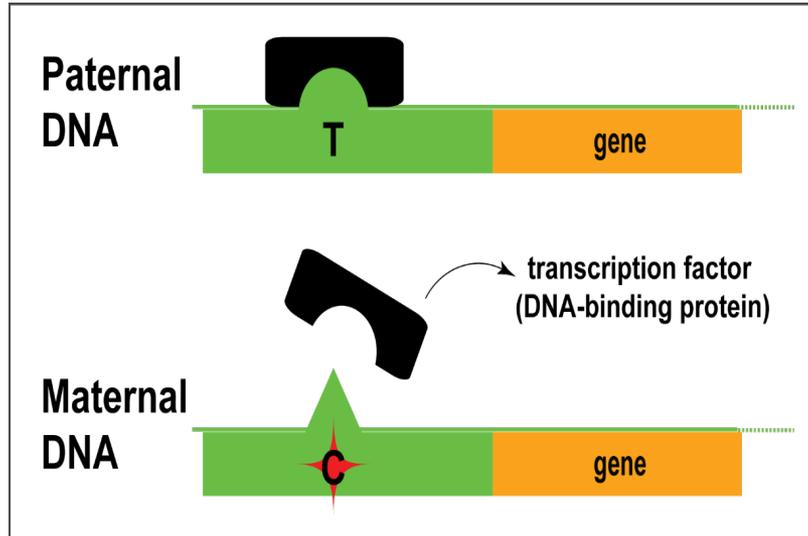
Fibroblast growth factor receptor 2 (pdb: 1IIL)



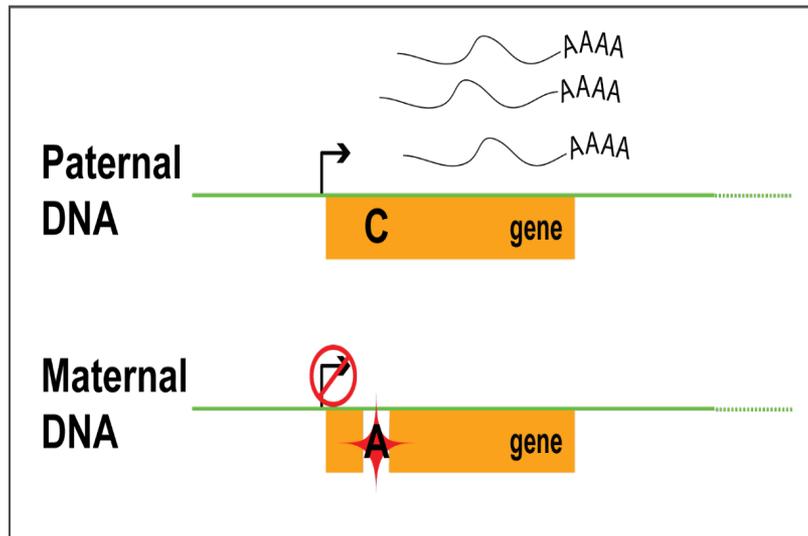
Personal Genomics: Handling Exponential Data Scaling through Prioritizing High-impact Variants

- Introduction
 - The **exponential scaling** of data generation & data processing
 - The **landscape of variants** in personal genomes
- Characterizing Rare Variants in Coding Regions
 - Identifying with STRESS **cryptic allosteric sites**
 - On surface & in interior bottlenecks
- Non-coding Variants : Prioritizing using AlleleDB in terms of **allelic elements**
 - Having observed difference in molecular activity in many contexts
- Putting it together in workflows: **Integrating evidence** on non-coding variants with FunSeq
 - Systematically weighting all the features
 - suggesting non-coding drivers
 - Prioritizing rare germline variants

Allele-specific binding and expression



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



e.g. allele-specific expression (ASE)

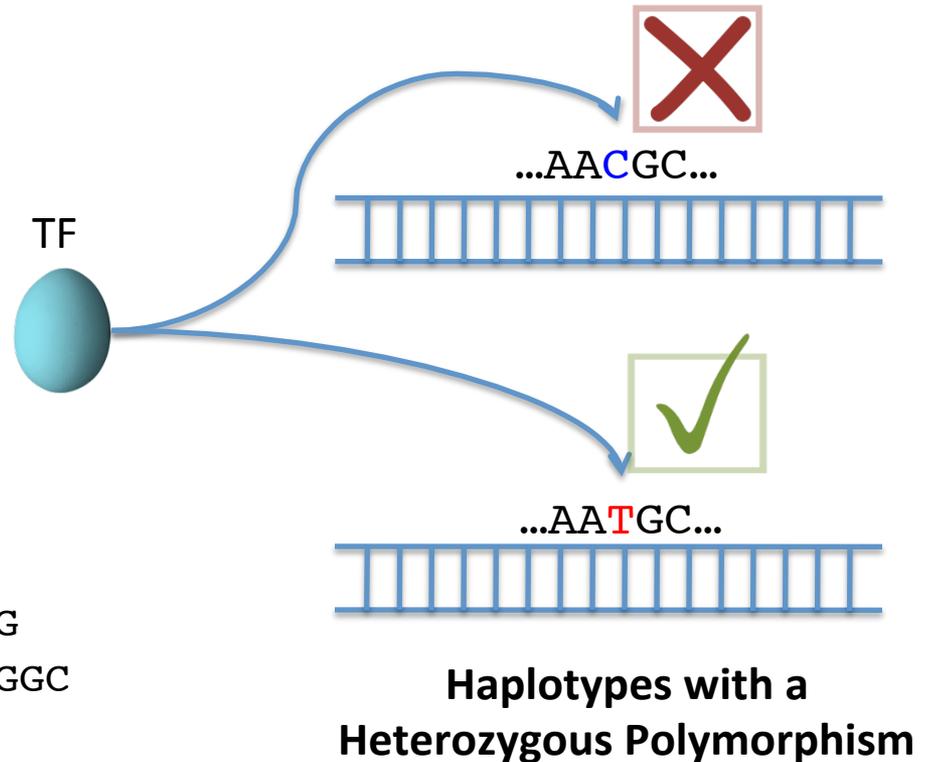
Inferring Allele Specific Binding/Expression using Sequence Reads

RNA/ChIP-Seq Reads

ACTTTGATAGCGTCAATG
 CTTTGATAGCGTCAATGC
 CTTTGATAGCGTCAACGC
 TTGACAGCGTCAATGCAC
 TGATAGCGTCAATGCACG
 ATAGCGTCAATGCACGTC
 TAGCGTCAATGCACGTCG
 CGTCAACGCACGTCGGGA
 GTCAATGCACGTCGAGAG
 CAAATGCACGTCGGGAGTT
 AAATGCACGTCGGGAGTTG
 TGCACGTTGGGAGTTGGC

10 x T

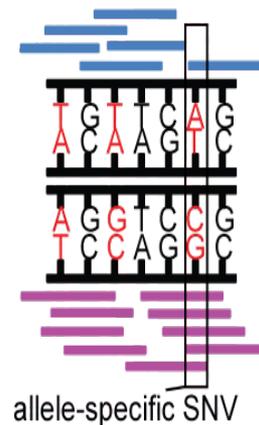
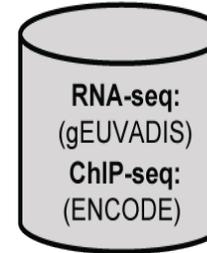
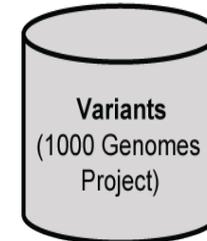
2 x C



Interplay of the annotation and individual sequence variants

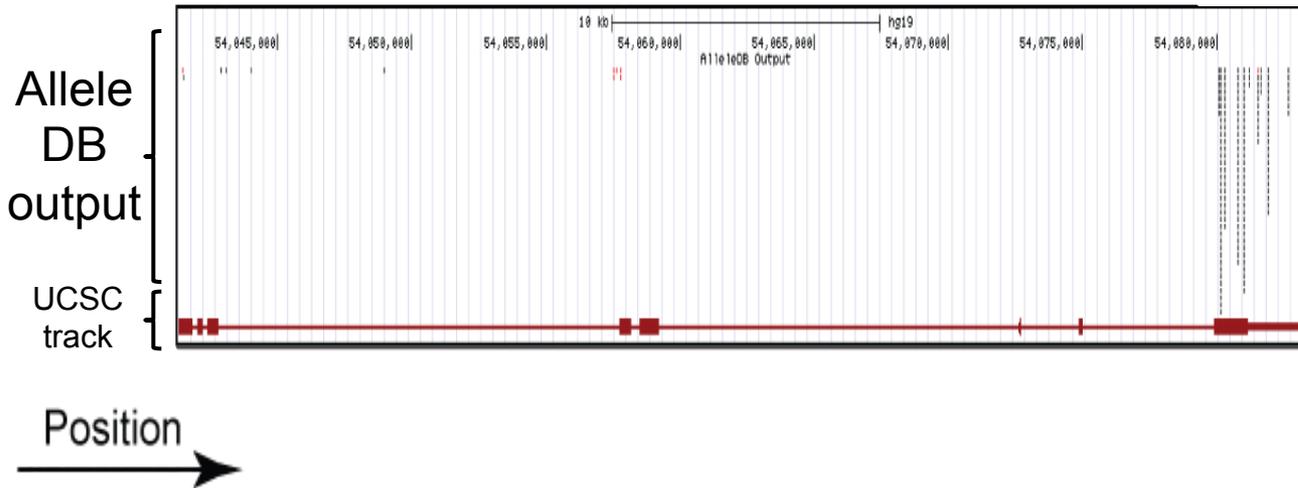
AlleleDB: Building 382 personal genomes to detect allele-specific variants on a large-scale

1. Build personal genomes
2. Align ChIP-seq & RNA-seq reads
3. Detect allele-specific variants via a series of filters and tests



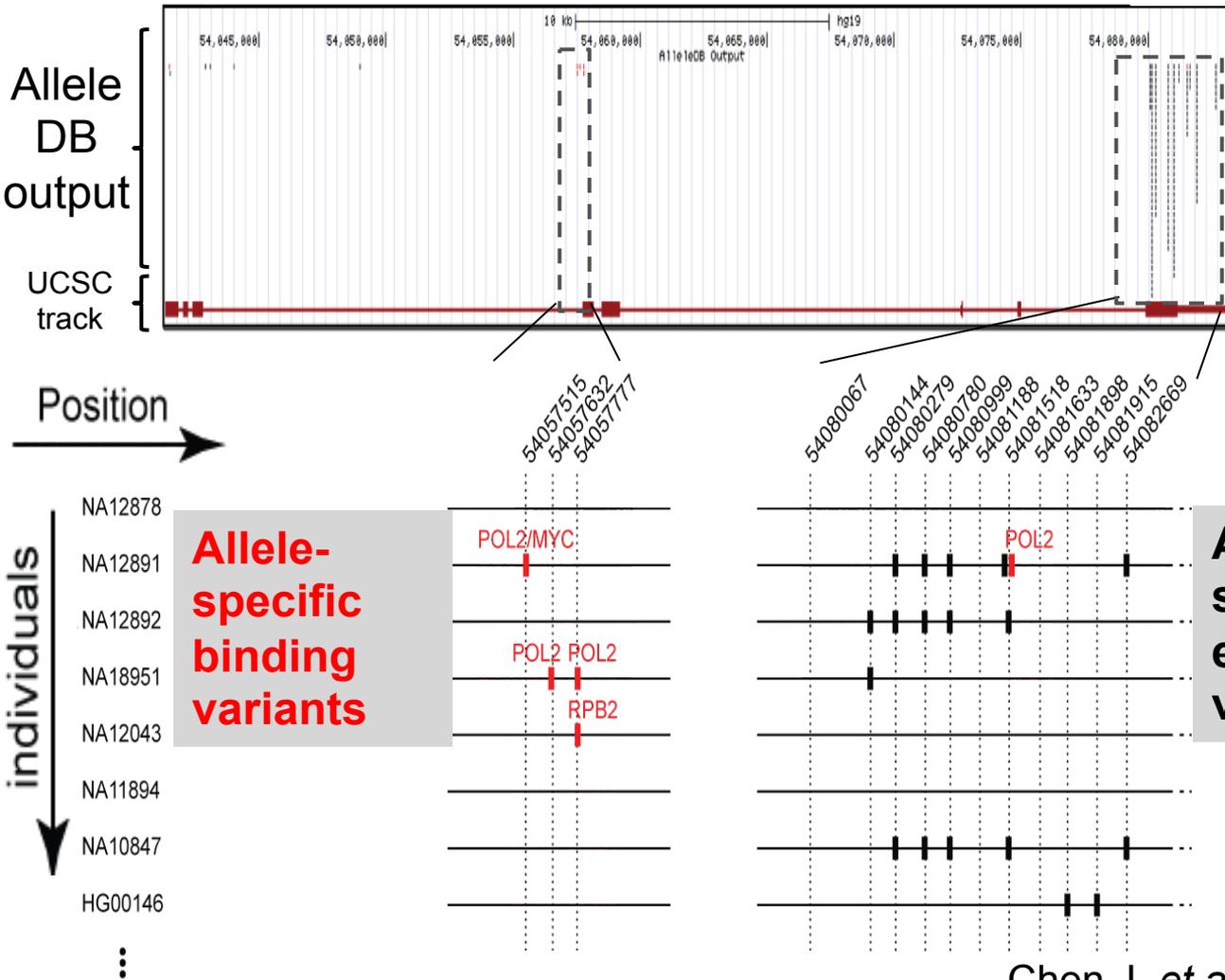
**Many Technical Issues:
Reference bias, Ambiguous
mapping bias, Over-dispersed
(non binomial null)**

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



- Interfaces with UCSC genome browser
- Showing ZNF331 gene structure

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



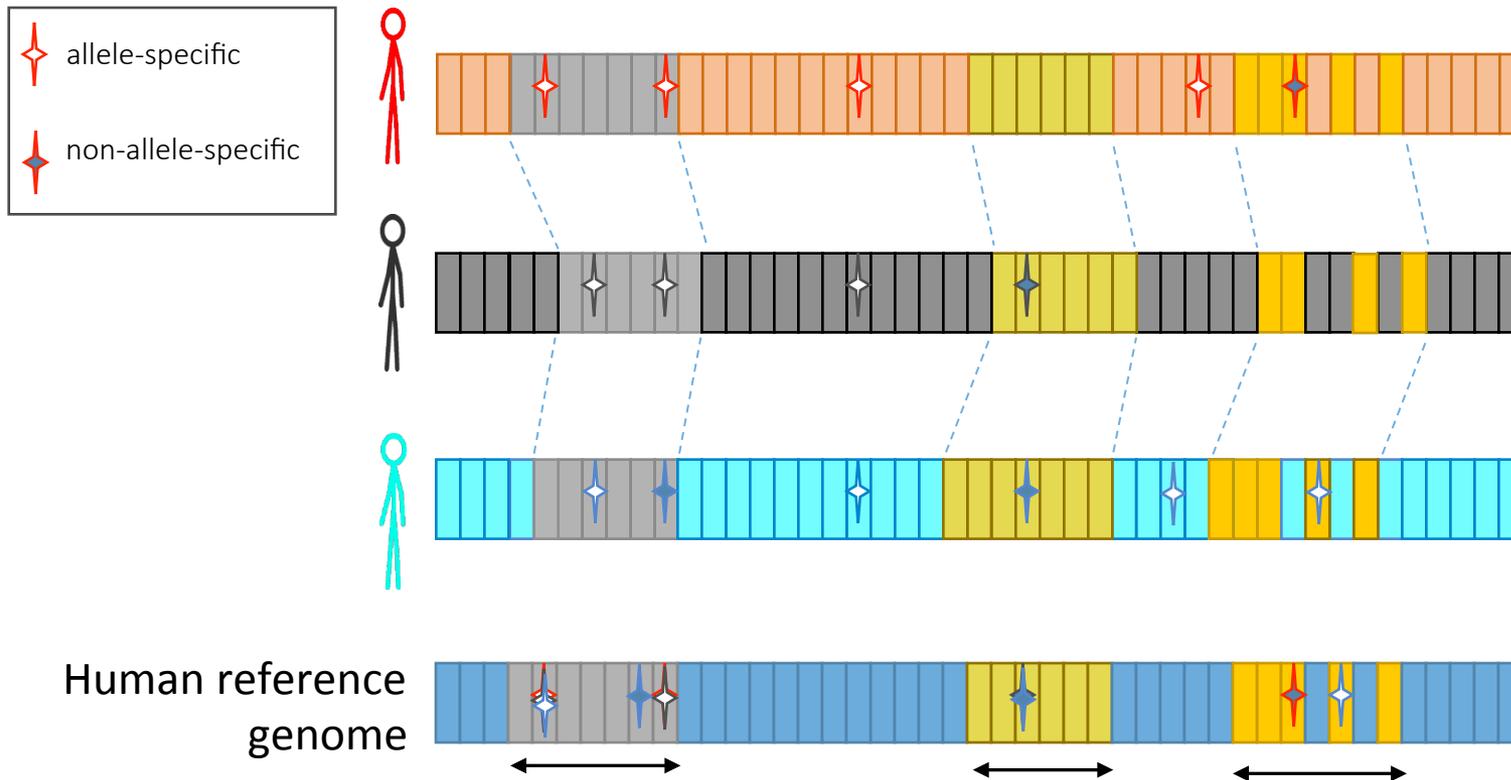
- Interfaces with UCSC genome browser
- Showing ZNF331 gene structure

Chen J. et al. (*Nature Commun*, in press)

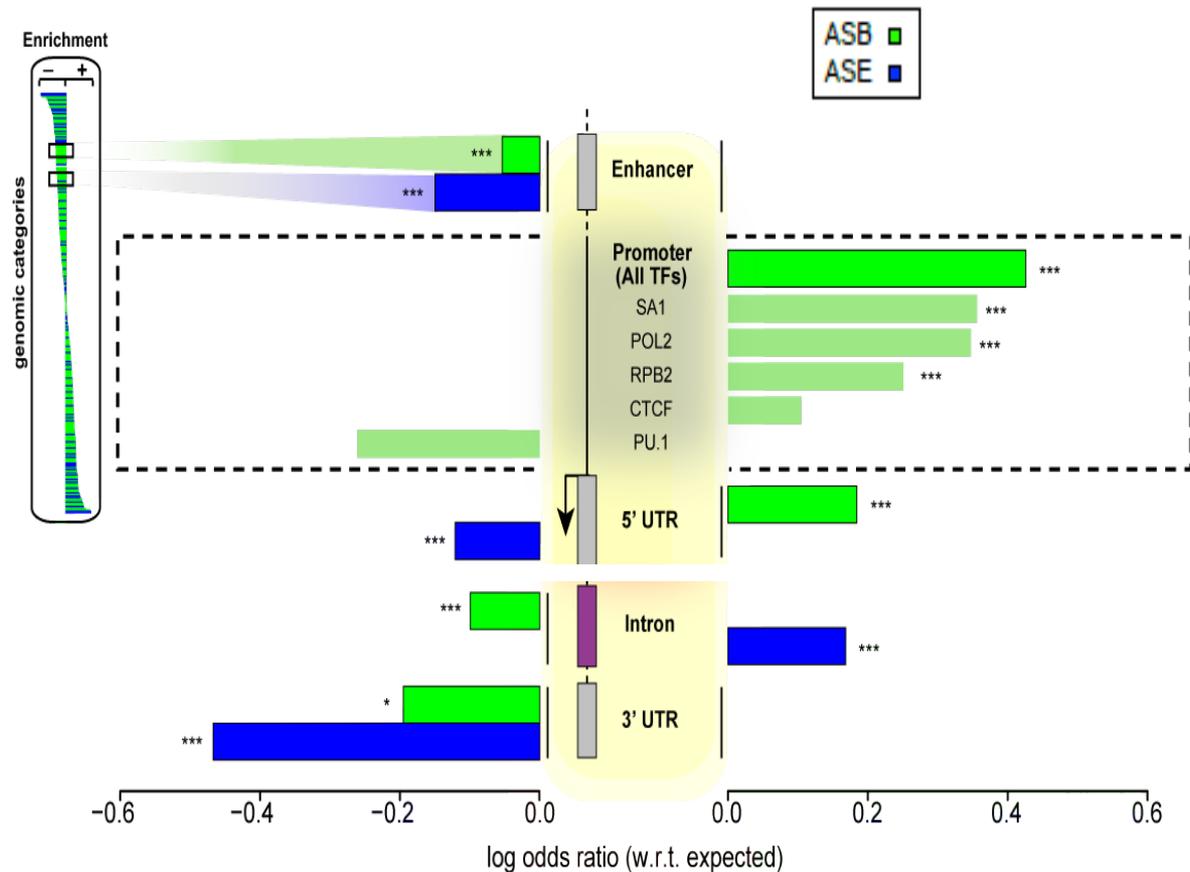
Collecting ASE/ASB variants into allele-specific genomic regions

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)



Groups of elements that are enriched or depleted in allelic activity

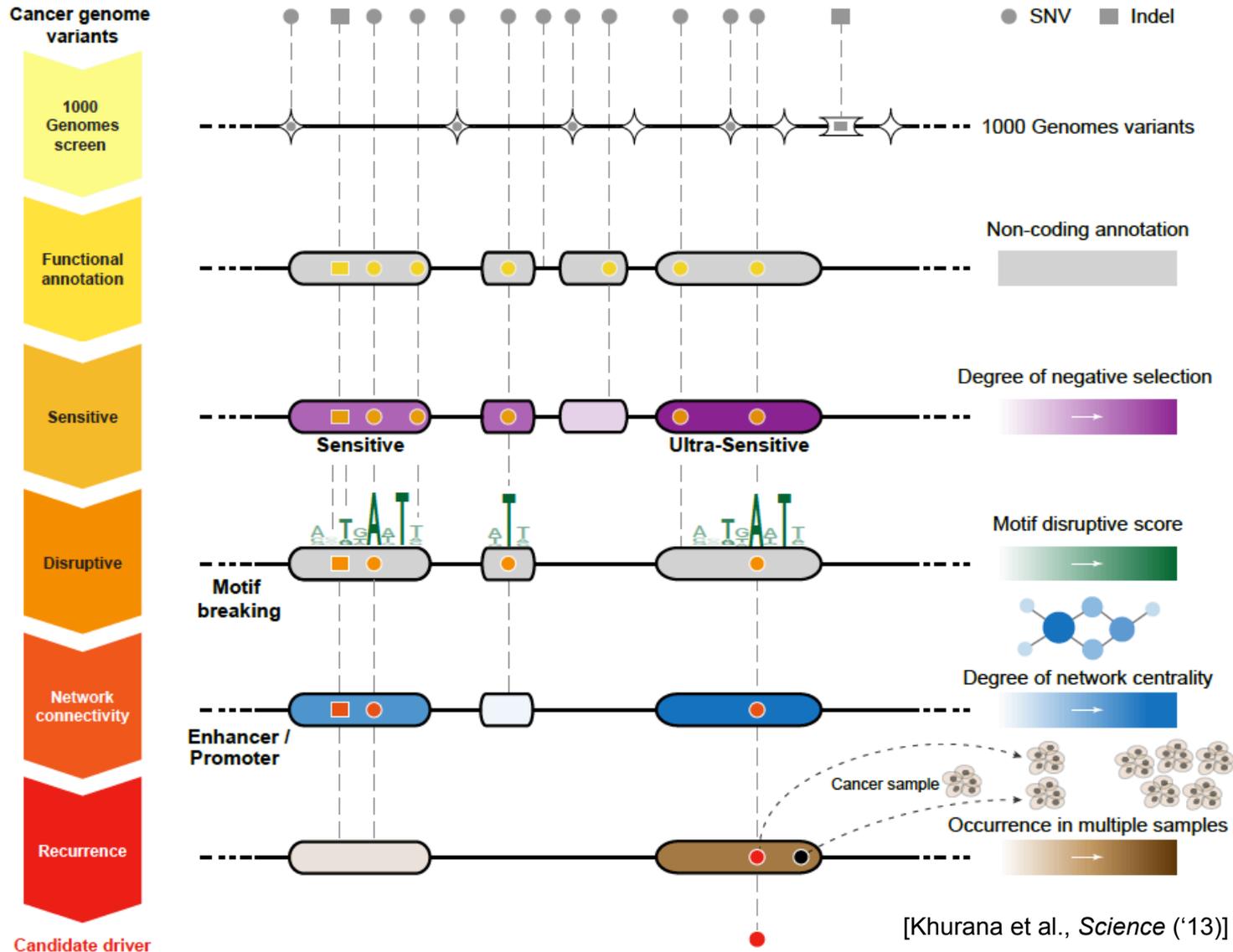


Chen J. *et al.* (*Nature Commun*, in press)

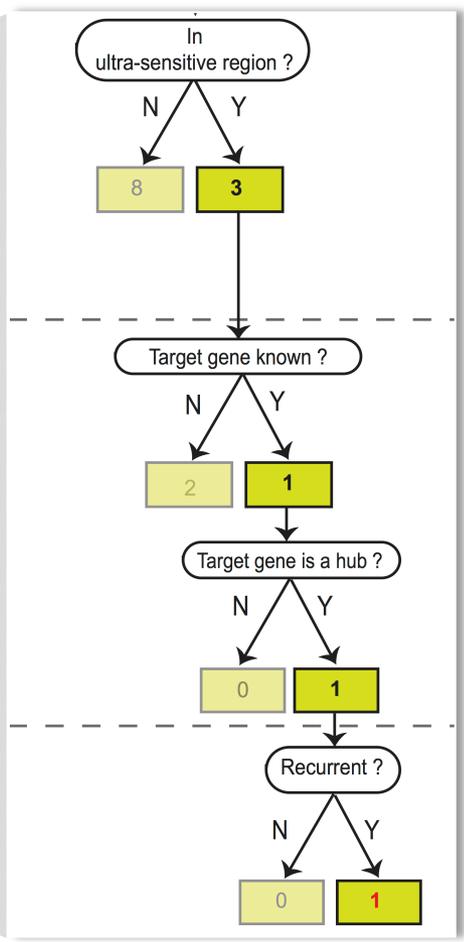
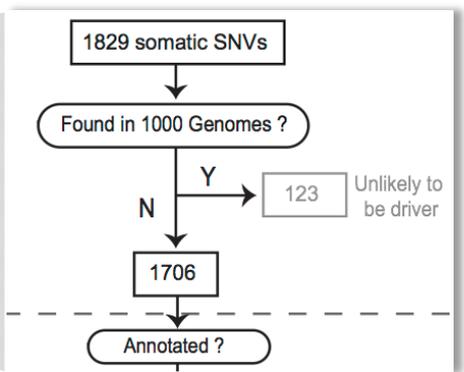
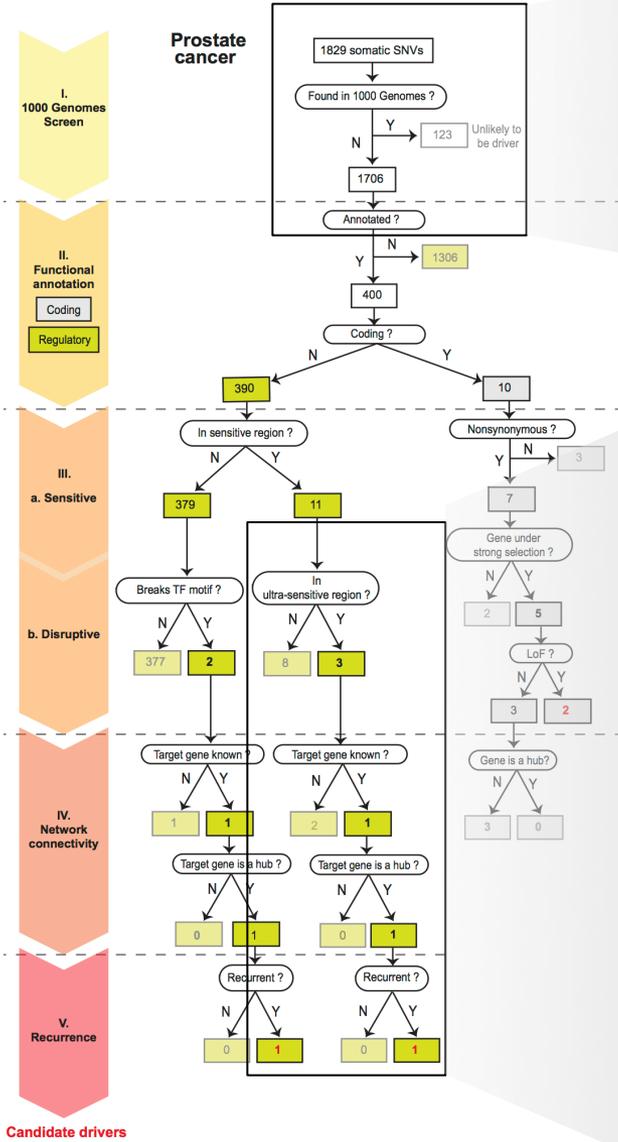
Personal Genomics: Handling Exponential Data Scaling through Prioritizing High-impact Variants

- Introduction
 - The **exponential scaling** of data generation & data processing
 - The **landscape of variants** in personal genomes
- Characterizing Rare Variants in Coding Regions
 - Identifying with STRESS **cryptic allosteric sites**
 - On surface & in interior bottlenecks
- Non-coding Variants : Prioritizing using AlleleDB in terms of **allelic elements**
 - Having observed difference in molecular activity in many contexts
- Putting it together in workflows: **Integrating evidence** on non-coding variants with FunSeq
 - Systematically weighting all the features
 - suggesting non-coding drivers
 - Prioritizing rare germline variants

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



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





Overview

This tool is specialized to prioritize somatic variants from cancer whole genome sequencing. It contains two components : 1) building data context from various resources; 2) variants prioritization. We provided downloadable scripts for users to customize the data context (found under 'Downloads'). The variants prioritization step is downloadable, and also implemented as web server (Right Panel), with pre-processed data context.

Instructions

- ♣ Input File - BED or VCF formatted. Click "green" button to add multiple files. With multiple files, the tool will do recurrent analysis. (Note: for BED format, user can put variants from multiple genomes in one file, see [Sample input file](#) .)
- ♣ Recurrence DB - User can choose particular cancer type from the database. The DB will continue be updated with newly available WGS data.
- ♣ Gene List - Option to analyze variants associated with particular set of genes. Note: Please use Gene Symbols, one row per gene.
- ♣ Differential Gene Expression Analysis - Option to detect differentially expressed genes in RNA-Seq data. Two files needed: expression file & class label file. Please refer to [Expression input files](#) for instructions to prepare those files.

♣ Note: In addition to on-site calculation, we also provide scores for all possible noncoding SNVs of GRCh37/hg19 under 'Downloads' (without annotation and recurrence analysis).

Input File: (only for hg19 SNVs)

Choose File No file chosen

BED or VCF files as input. [Sample input file](#)

Output Format:

bed

MAF:

0

Minor allele frequency threshold to filter polymorphisms from 1KG (value 0~1)

Cancer Type from Recurrence DB: [Summary table](#)

All Cancer Types

[Add a gene list](#) (Optional)

[Add differential gene expression analysis](#) (Optional)

Upload

Site integrates user variants with large-scale context

Data Context

Variant Prioritization

Weighted scoring scheme

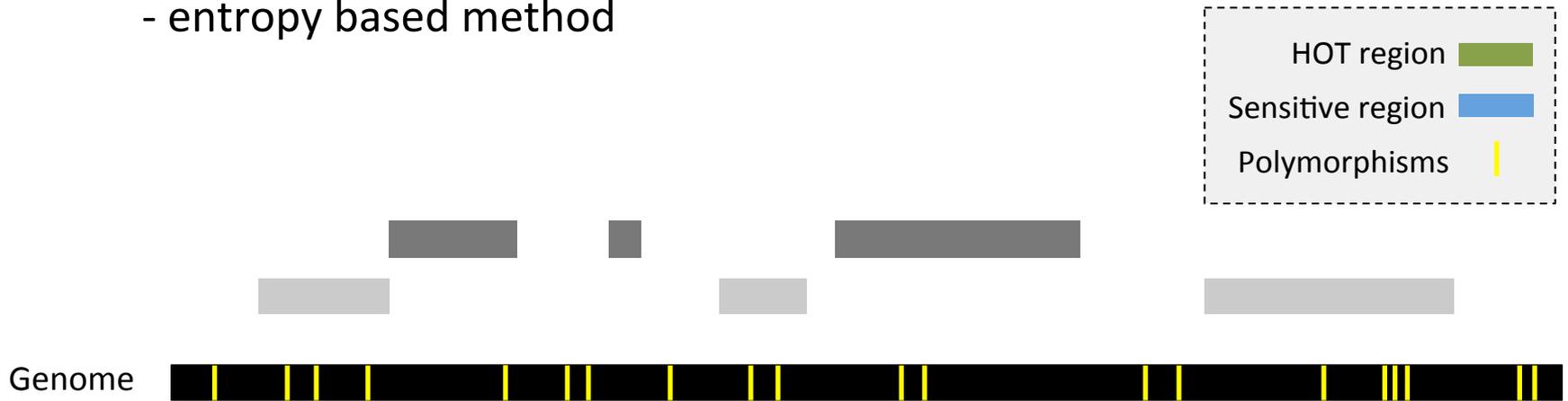
Highlighting variants

User Variants

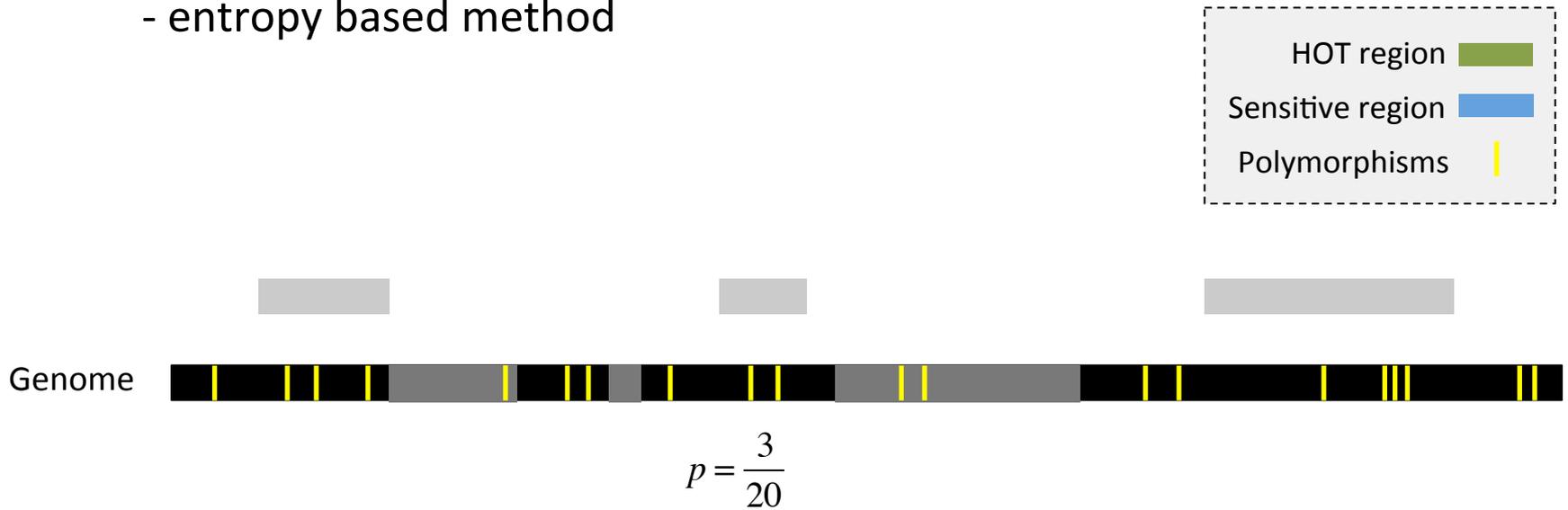
Variant Reports

FunSeq.gersteinlab.org

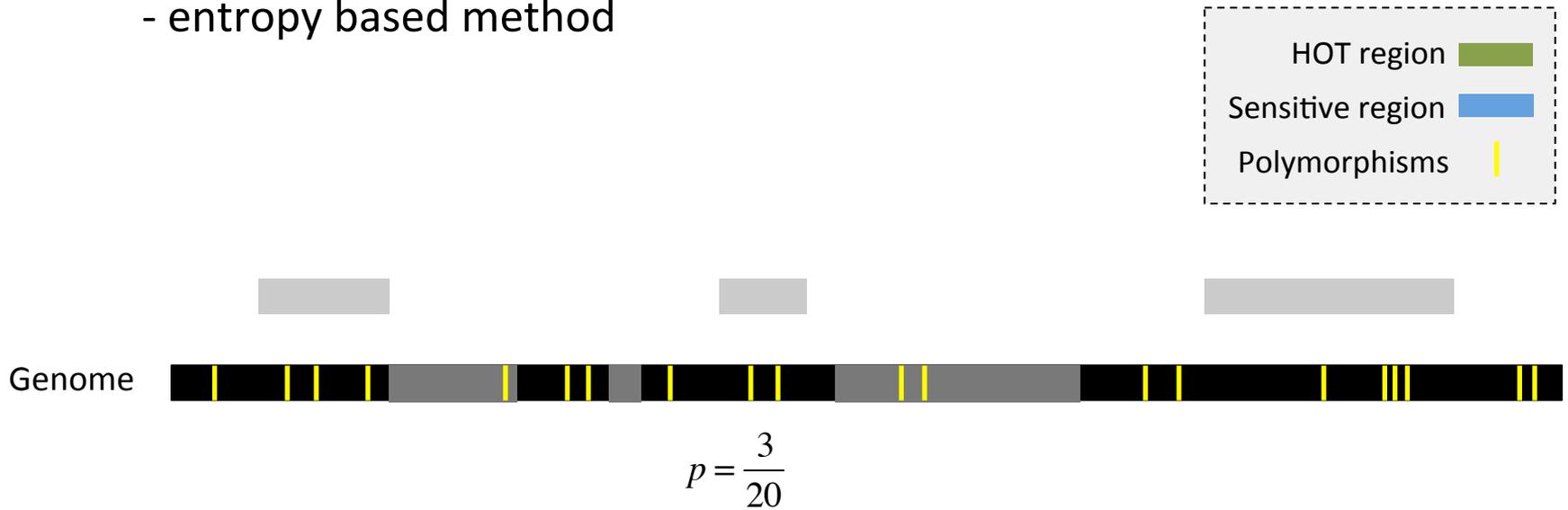
- Feature weight
 - Weighted with mutation patterns in natural polymorphisms
(features frequently observed weight less)
 - entropy based method



- Feature weight
 - Weighted with mutation patterns in natural polymorphisms
(features frequently observed weight less)
 - entropy based method



- Feature weight
 - Weighted with mutation patterns in natural polymorphisms
(features frequently observed weight less)
 - entropy based method

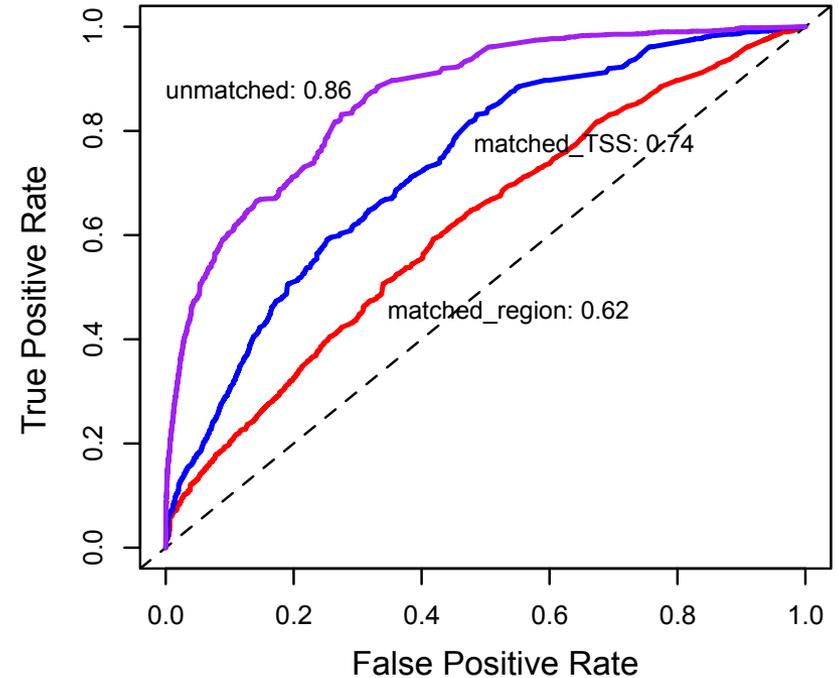
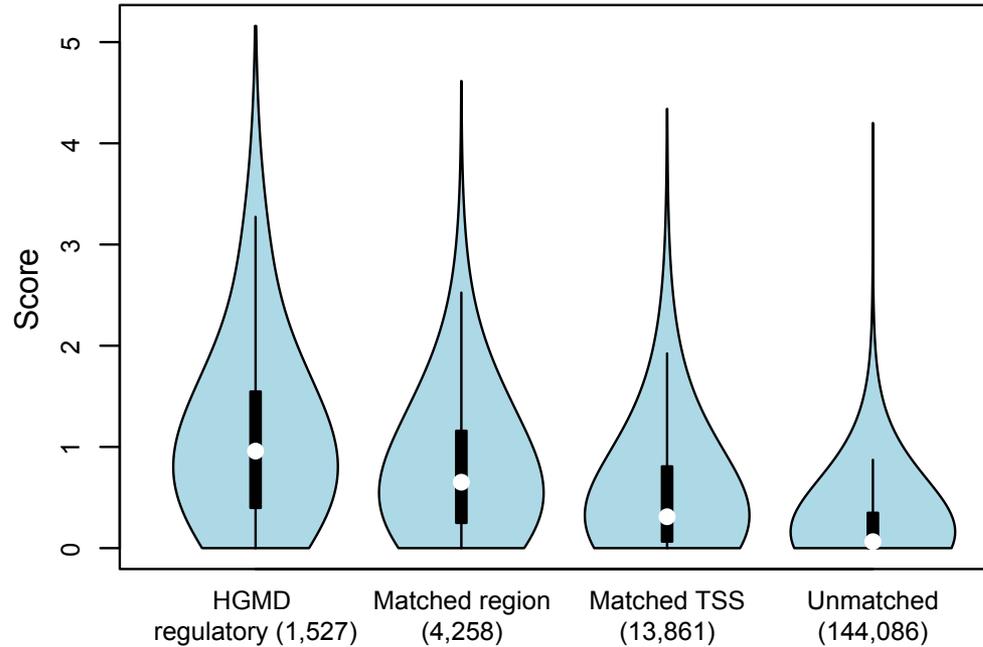


Feature weight: $w_d = 1 + p_d \log_2 p_d + (1 - p_d) \log_2 (1 - p_d)$

$p \uparrow$ $w_d \downarrow$ $p = \text{probability of the feature overlapping natural polymorphisms}$

For a variant: $\text{Score} = \sum w_d$ of observed features

Germline pathogenic variants show higher core scores than controls



3 controls with natural polymorphisms (allele frequency $\geq 1\%$)

1. Matched region: 1kb around HGMD variants

2. Matched TSS: matched for distance to TSS

3. Unmatched: randomly selected

Personal Genomics: Handling Exponential Data Scaling through Prioritizing High-impact Variants

- Introduction
 - The **exponential scaling** of data generation & data processing
 - The **landscape of variants** in personal genomes
- Characterizing Rare Variants in Coding Regions
 - Identifying with STRESS **cryptic allosteric sites**
 - On surface & in interior bottlenecks
- Non-coding Variants : Prioritizing using AlleleDB in terms of **allelic elements**
 - Having observed difference in molecular activity in many contexts
- Putting it together in workflows: **Integrating evidence** on non-coding variants with FunSeq
 - Systematically weighting all the features
 - suggesting non-coding drivers
 - Prioritizing rare germline variants

Personal Genomics: Handling Exponential Data Scaling through Prioritizing High-impact Variants

- **Introduction**
 - The exponential scaling of data generation & data processing
 - The landscape of variants in personal genomes
- **Characterizing Rare Variants in Coding Regions**
 - Identifying with **STRESS** cryptic allosteric sites
 - On surface & in interior bottlenecks
- **Non-coding Variants : Prioritizing using AlleleDB in terms of allelic elements**
 - Having observed difference in molecular activity in many contexts
- **Putting it together in workflows: Integrating evidence on non-coding variants with FunSeq**
 - Systematically weighting all the features
 - suggesting non-coding drivers
 - Prioritizing rare germline variants

Acknowledgments

Hiring Postdocs. See gersteinlab.org/jobs

AlleleDB.gersteinlab.org

J **Chen**,

J **Rozowsky**,

TR **Galeev**, A Harmanci,
R Kitchen, J Bedford,
A Abyzov, Y Kong, L Regan

CostSeq2

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

STRESS.molmovdb.org

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



FunSeq.gersteinlab.org & **FunSeq2**.gersteinlab.org

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

1000 Genomes Project Consortium, et al

Extra



Info about content in this slide pack

- General PERMISSIONS
 - This Presentation is copyright Mark Gerstein, Yale University, 2015.
 - 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>