

Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

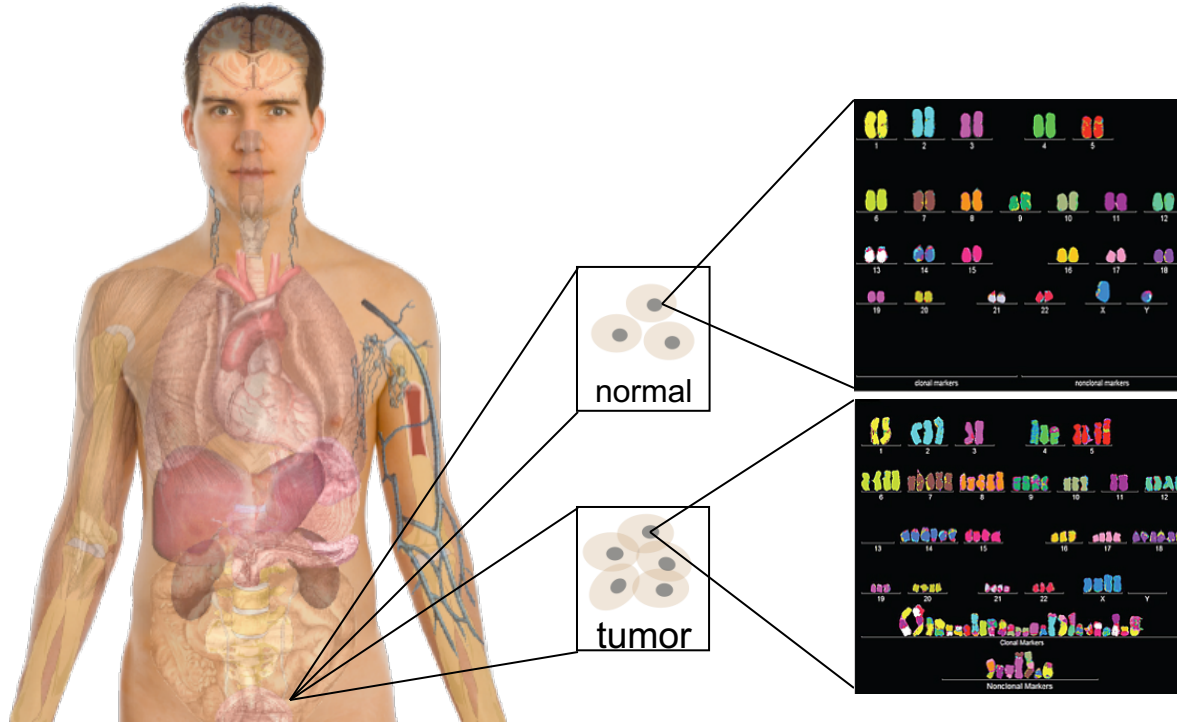
Mark Gerstein
Yale

Slides freely
downloadable from
Lectures.GersteinLab.org
& “tweetable” (via **@MarkGerstein**).

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

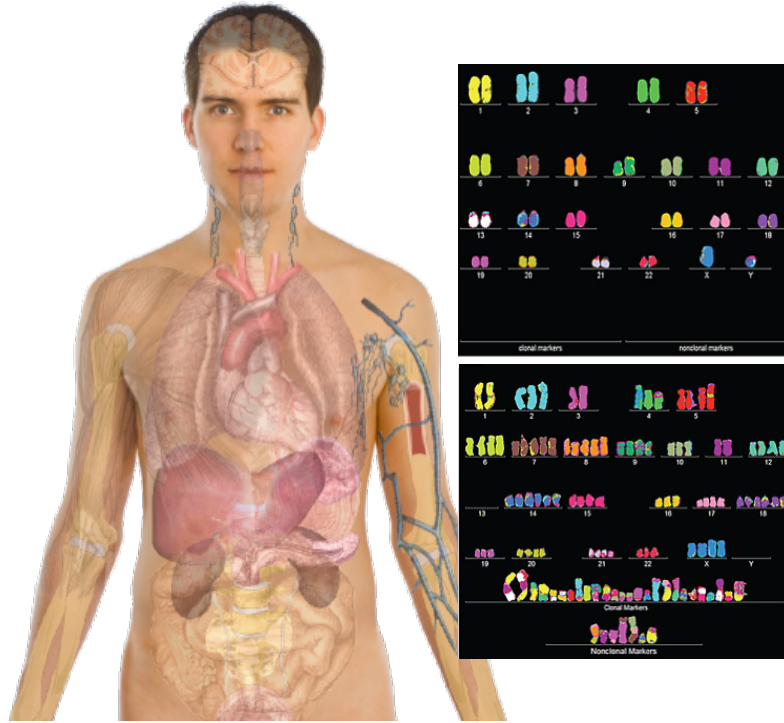
Personal Genomics as a Gateway into Biology

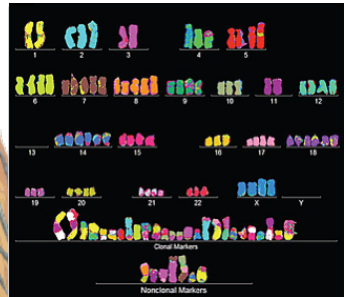
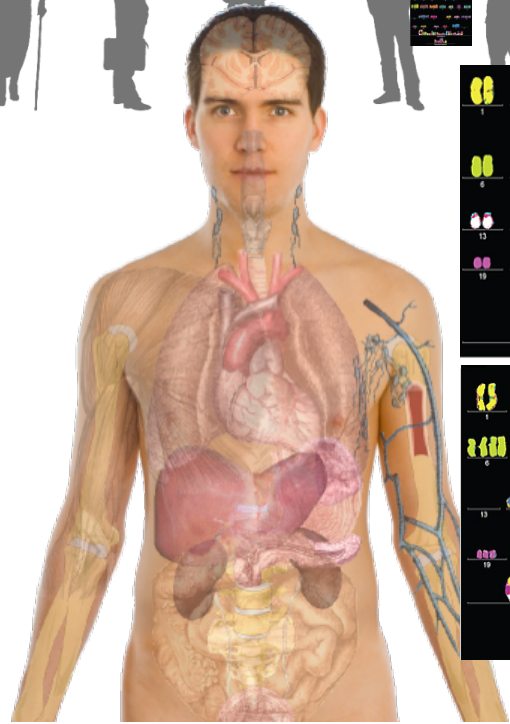
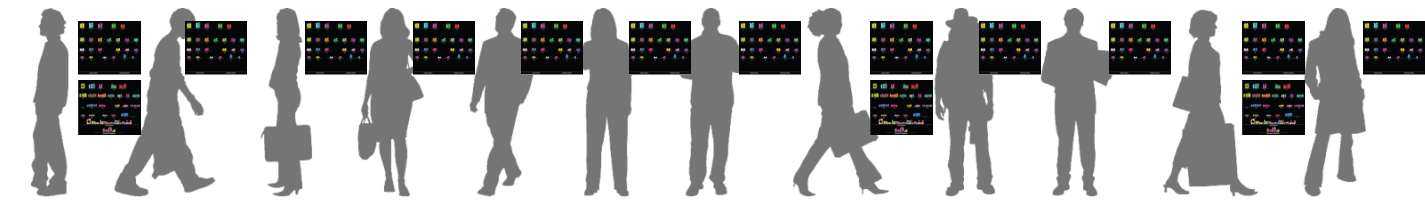
Personal genomes will soon become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.



Personal Genomics as a Gateway into Biology

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





Keys to genome interpretation

Relating individuals' variants to **DBs**

Scaling DBs to the **population**

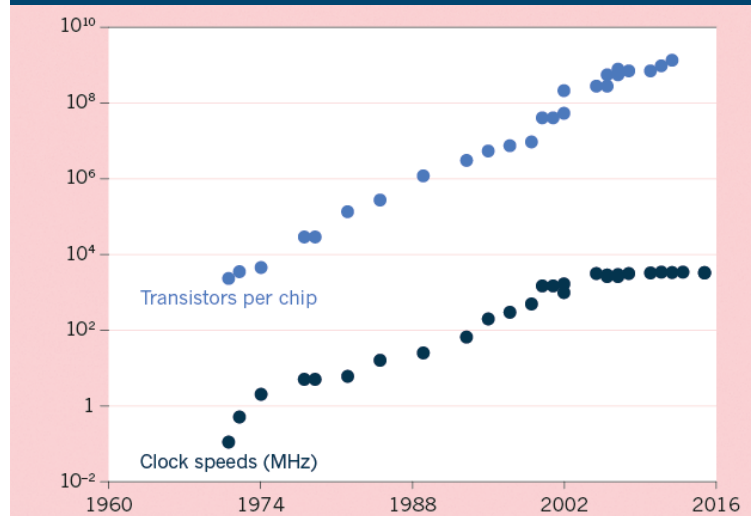
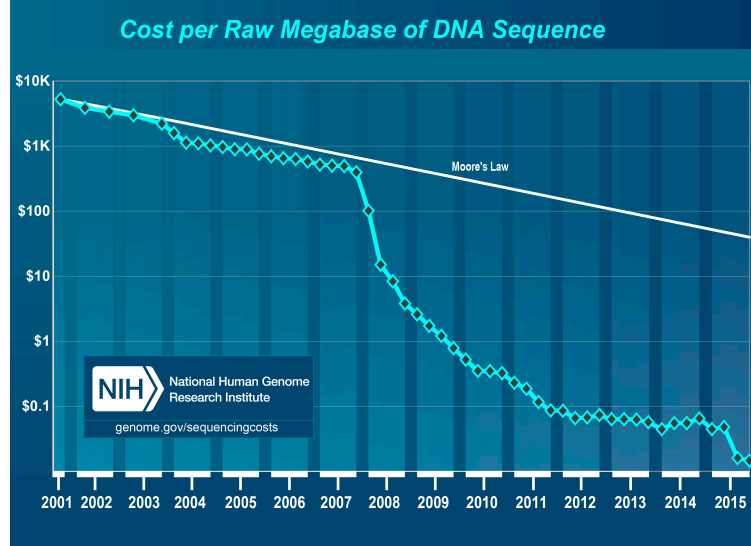
Identifying **key variants** -
separating into rare, recurrent,
common, &c

The **Scaling** of Genomic Data Science:

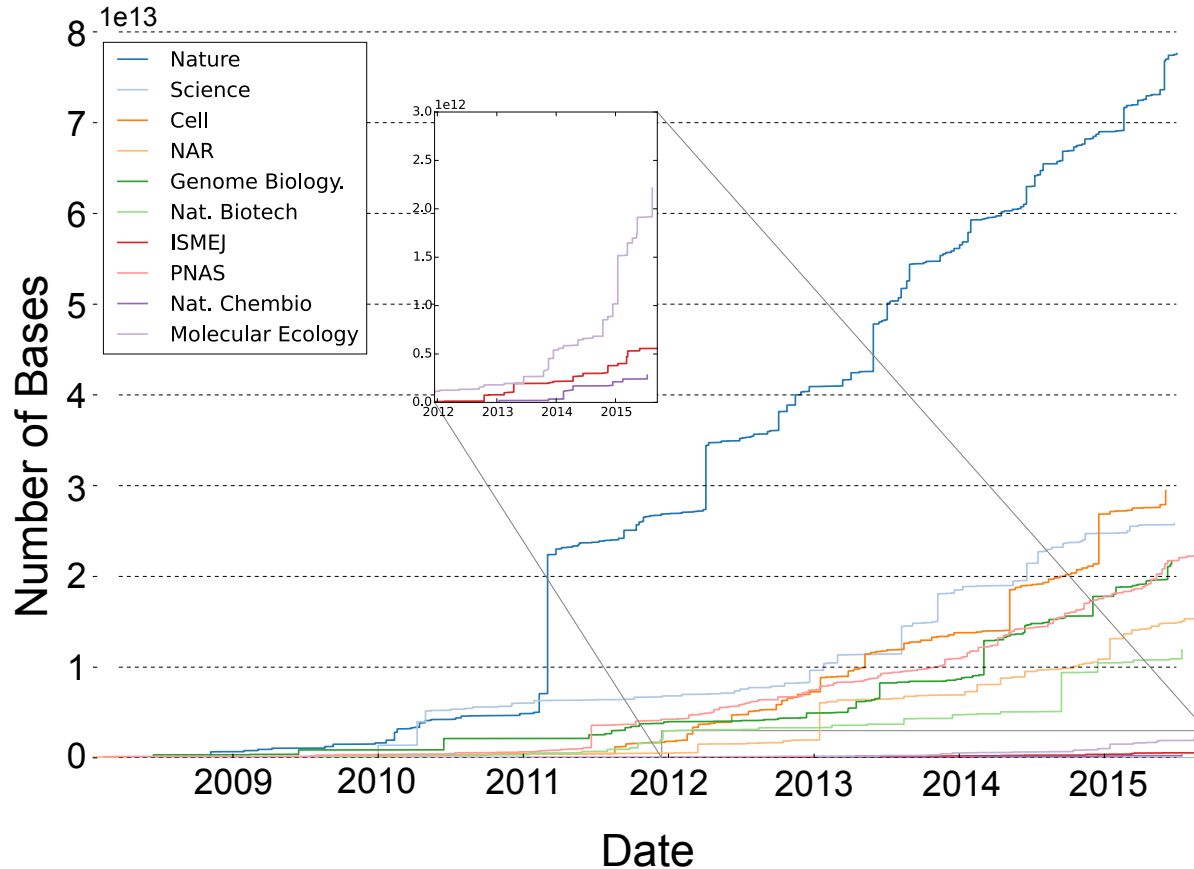
Powered by exponential increases in data & computing

(**Moore's Law**)

[NHGRI website + Waldrop ('15) Nature]



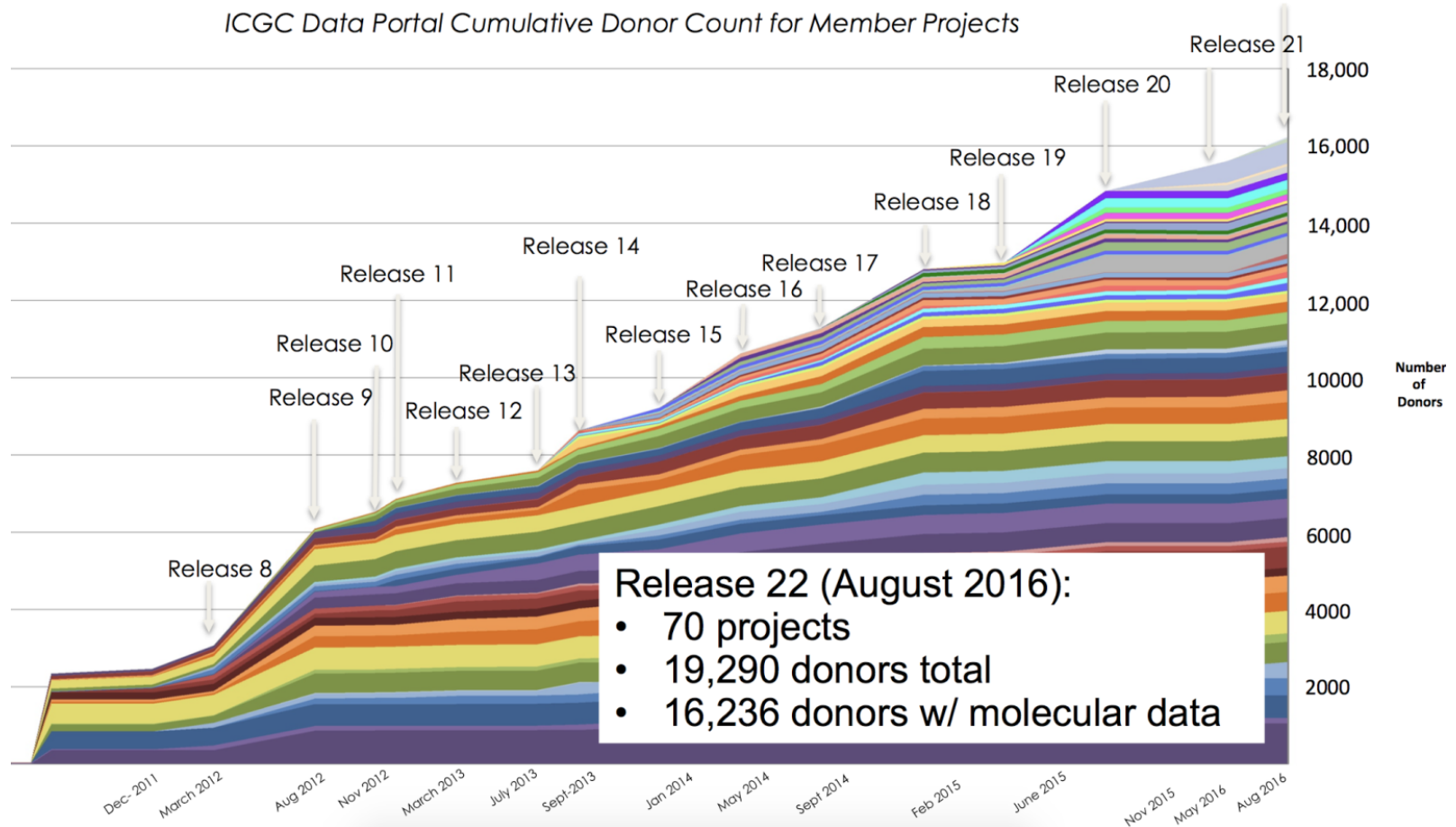
Exponential **Scaling** Changes Fields Using Genomic Data



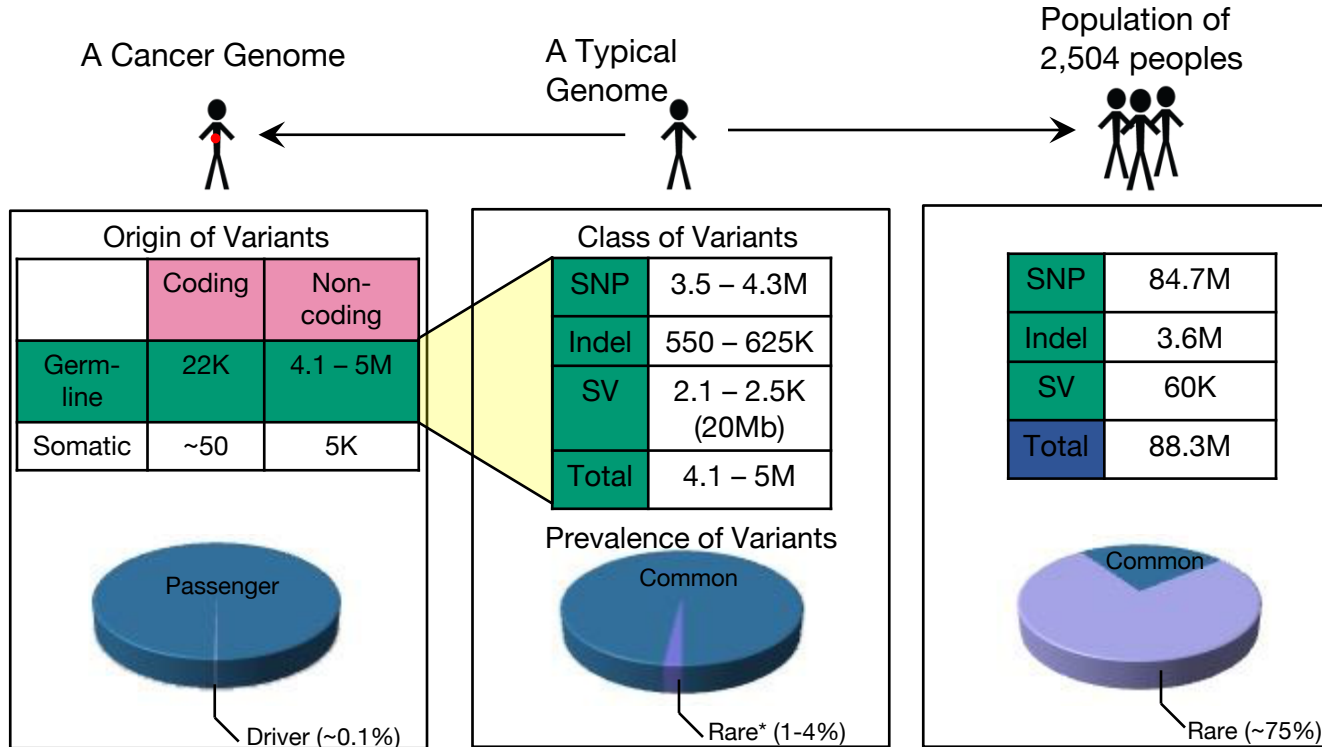
Growth of ICGC datasets

Release 22
70 ICGC
projects

ICGC Data Portal Cumulative Donor Count for Member Projects



Human Genetic Variation



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

Finding Key Variants

Germline

CAN YOU FIND THE PANDA?

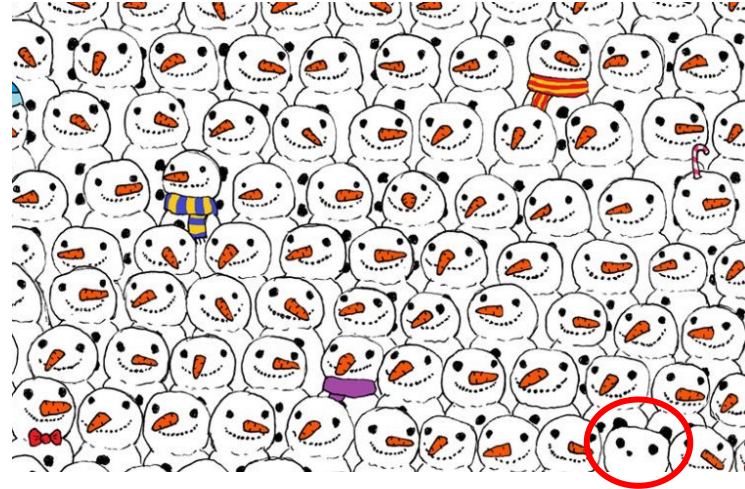


- **Common variants**
 - Can be most readily associated with phenotype (ie disease) via GWAS
 - Usually their functional effect is weaker
 - Many are non-coding
 - Issue of LD in identifying the actual causal variant.
- **Rare variants**
 - Associations are usually underpowered due to low frequencies but often have larger functional impact
 - Can be collapsed in the same element to gain statistical power (burden tests).

CAN YOU FIND THE PANDA?

Finding Key Variants

Somatic



- **Overall**

- Often these can be thought of as very rare variants

- **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**

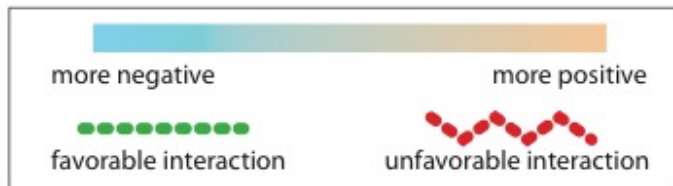
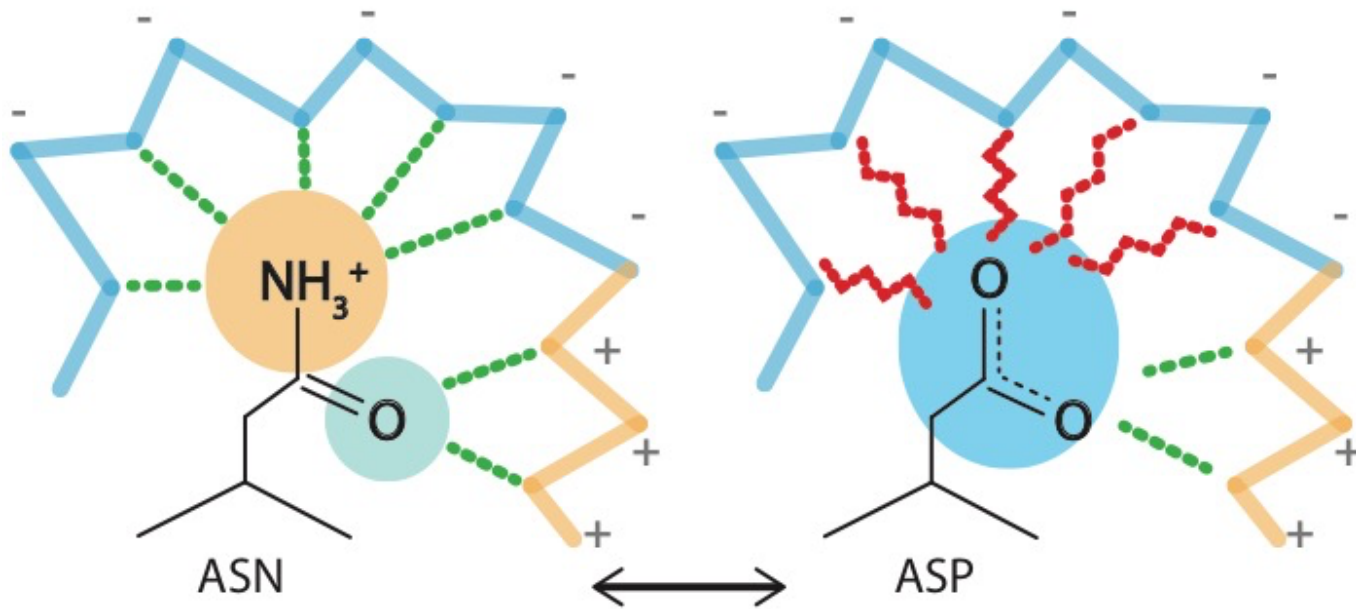
- Conceptually, a passenger mutation has no direct or indirect effect on the selective growth advantage of the cell in which it occurred.

Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - **The exponential scaling** of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - **FunSeq** integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - **RADAR**: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - **uORFs**: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

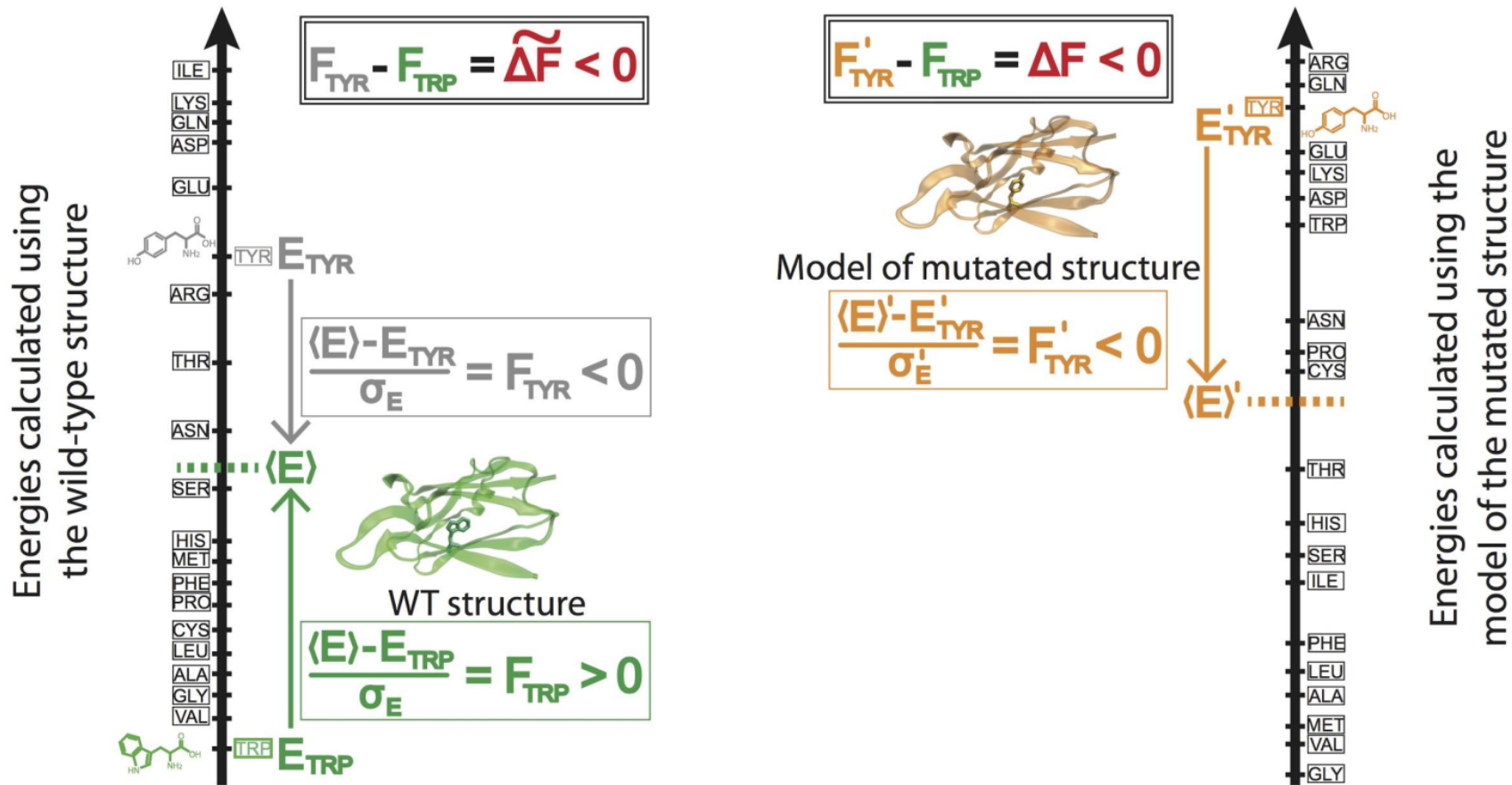
Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - The exponential scaling of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - Frustration as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - FunSeq integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - RADAR: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - uORFs: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

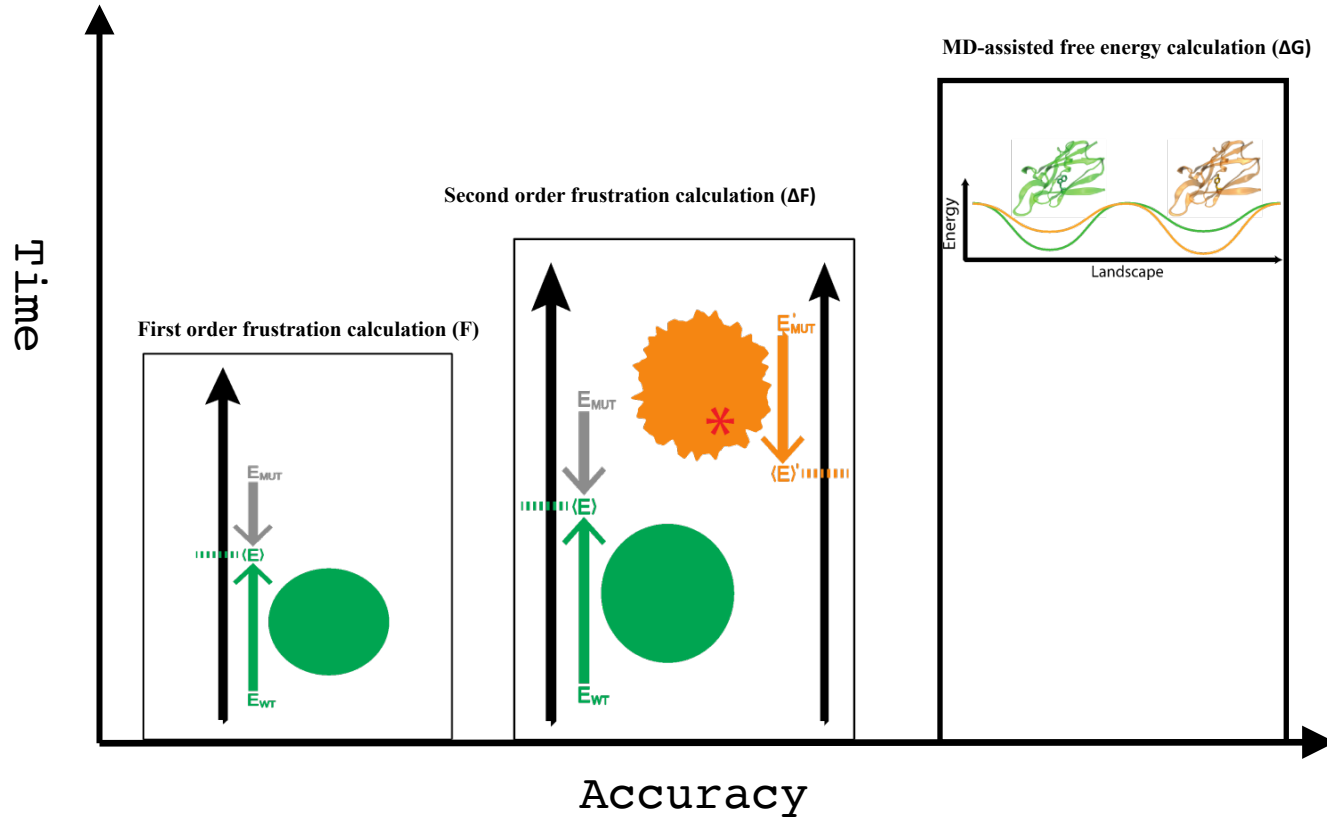


What is
localized
frustration
?

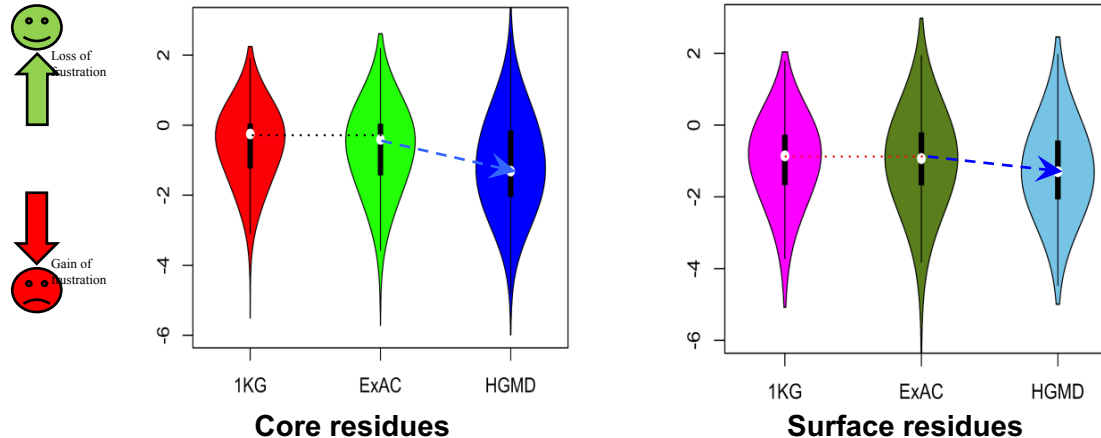
Workflow for evaluating localized frustration changes (ΔF)



Complexity of the second order frustration calculation

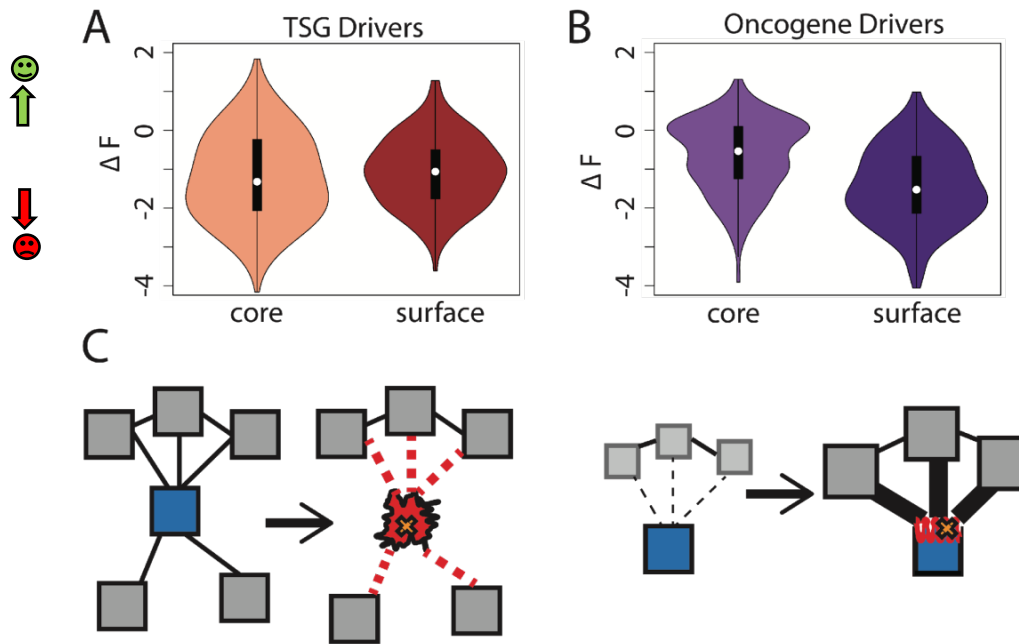


Comparing ΔF values across different SNV categories: disease v normal



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

Comparison between ΔF distributions: TSGs v. oncogenes



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

Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - The exponential scaling of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - Frustration as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - FunSeq integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - RADAR: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - uORFs: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

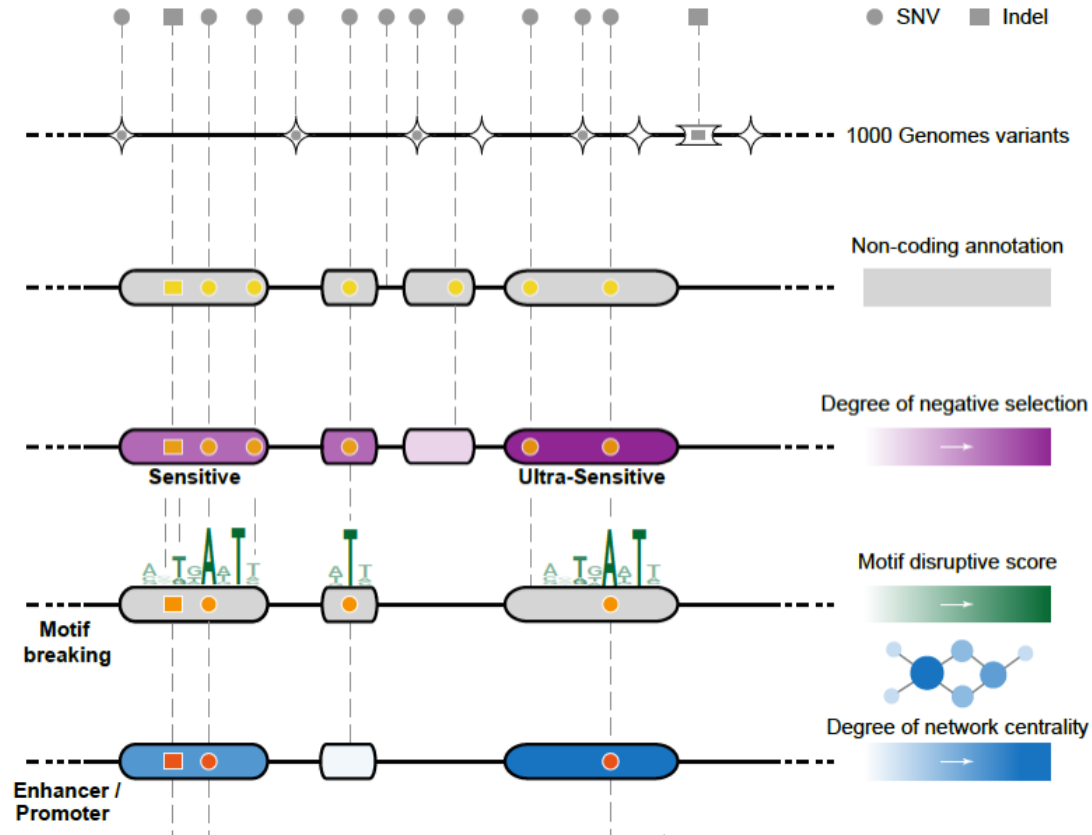
Funseq: a flexible framework to determine functional impact & use this to prioritize variants

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

Conservation (GERP, allele freq.)

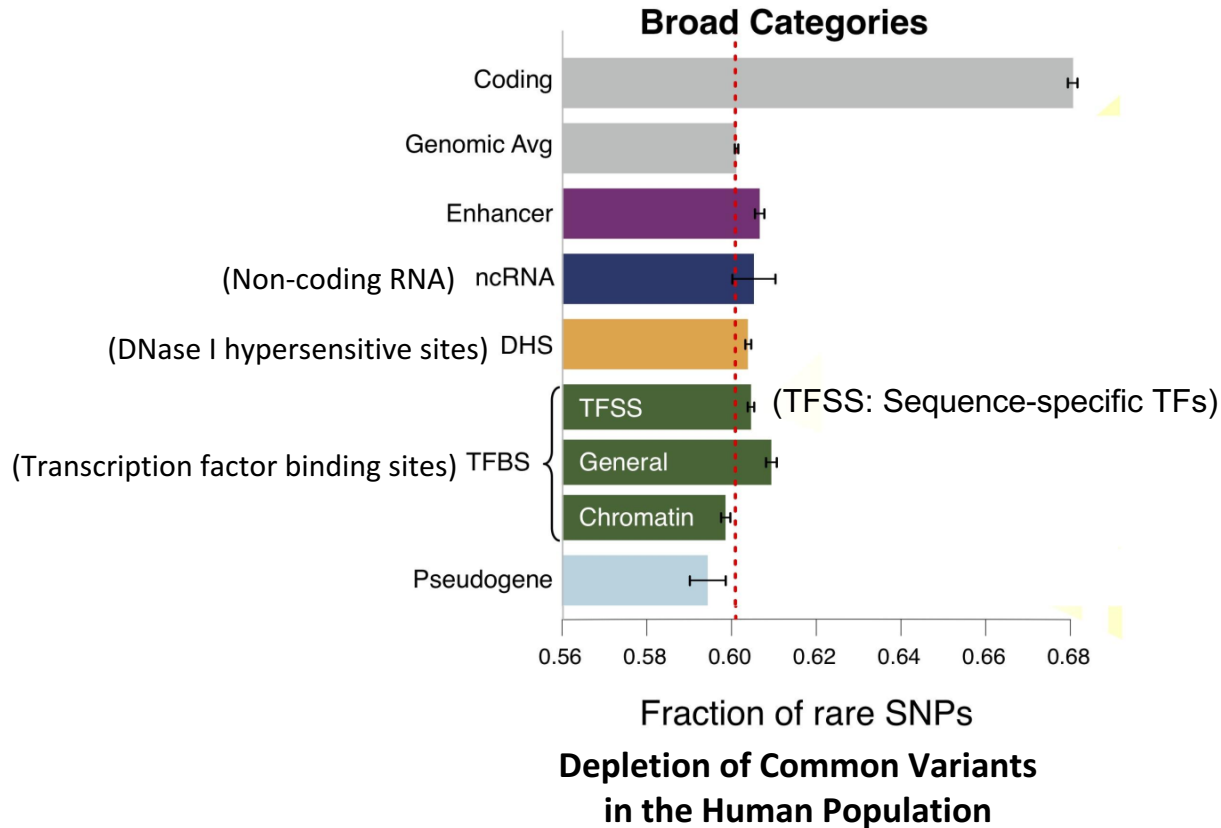
Mutational impact (motif breaking, Lof)

Network (centrality position)



Finding "Conserved" Sites in the Human Population:

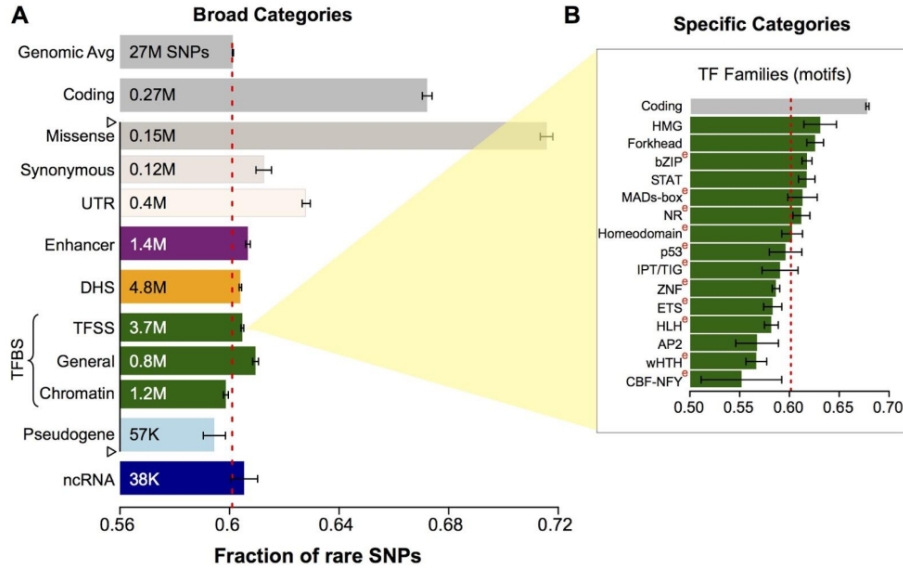
Negative selection in non-coding elements based on
Production ENCODE & 1000G Phase 1



Broad categories of
regulatory regions under
negative selection
Related to:

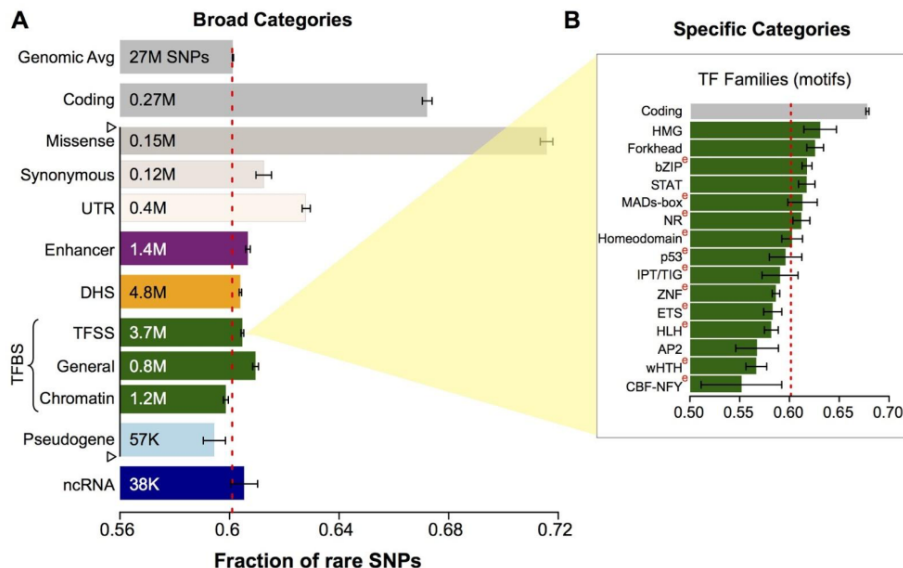
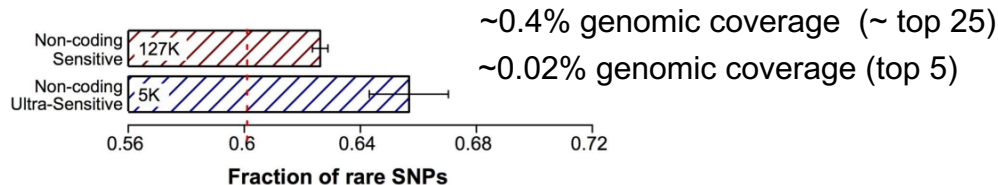
ENCODE, *Nature*, 2012
Ward & Kellis, *Science*, 2012
Mu et al, *NAR*, 2011

Differential selective constraints among specific sub-categories



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

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

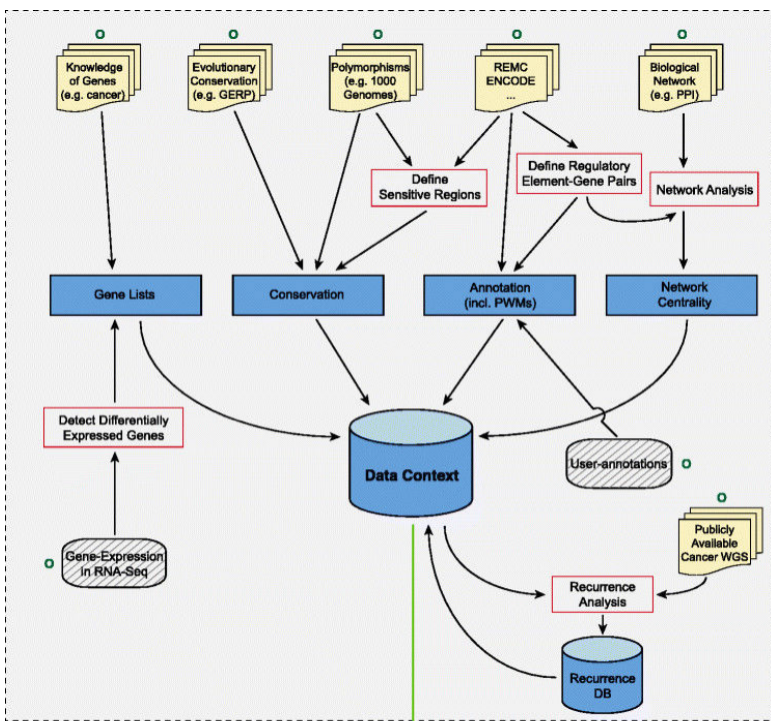
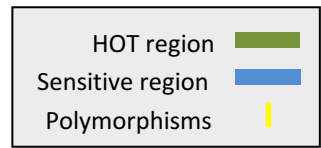


Sub-categorization possible because of better statistics from 1000G phase 1 v pilot

Defining Sensitive non-coding Regions

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

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



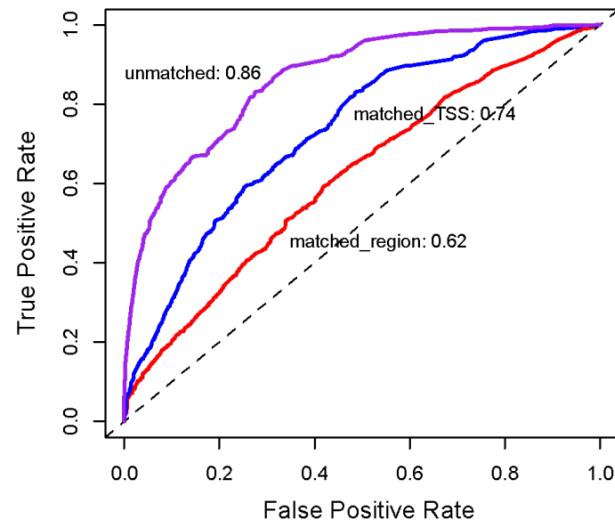
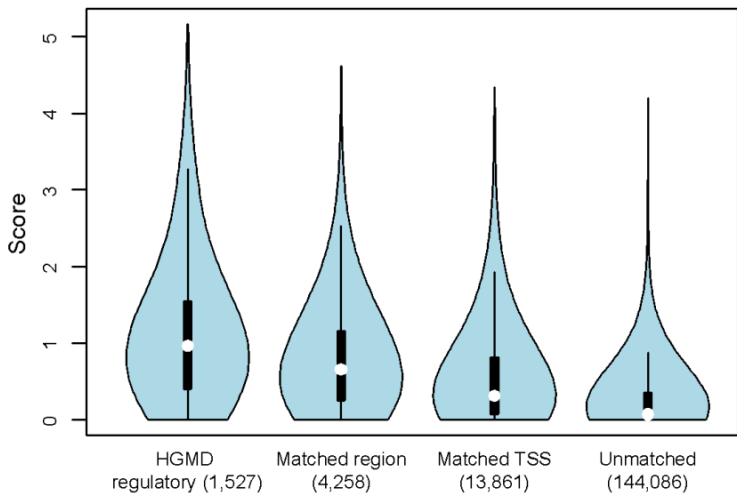
Genome



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

- Info. theory based method (ie annotation “surprisal”) for weighting consistently many genomic features
- Practical web server
- Submission of variants & pre-computed large data context from uniformly processing large-scale datasets

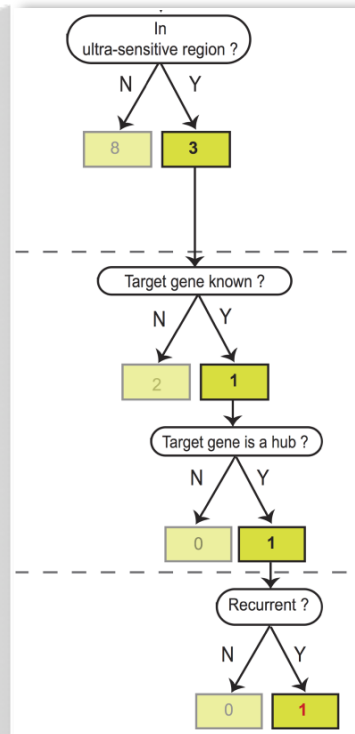
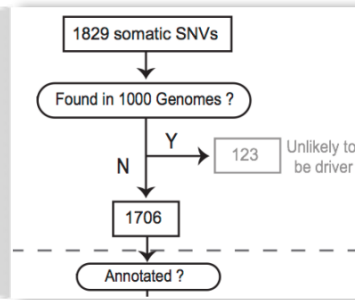
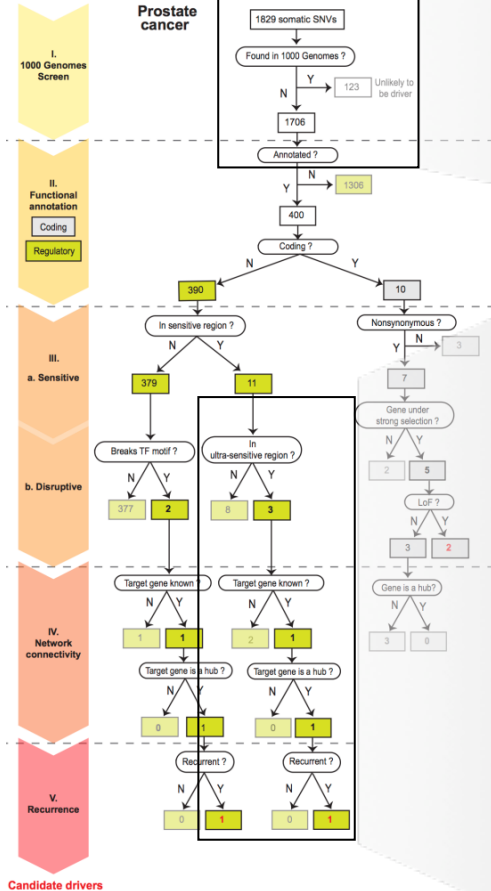
Germline pathogenic variants show higher core scores than controls



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

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

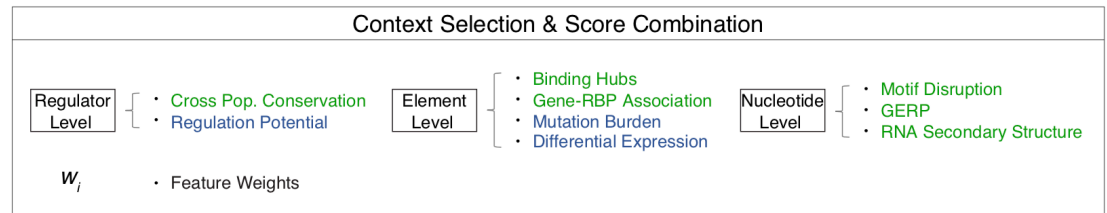
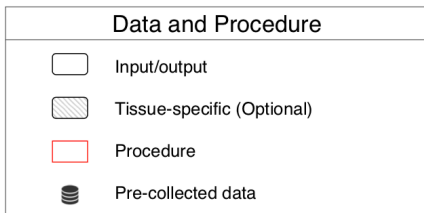
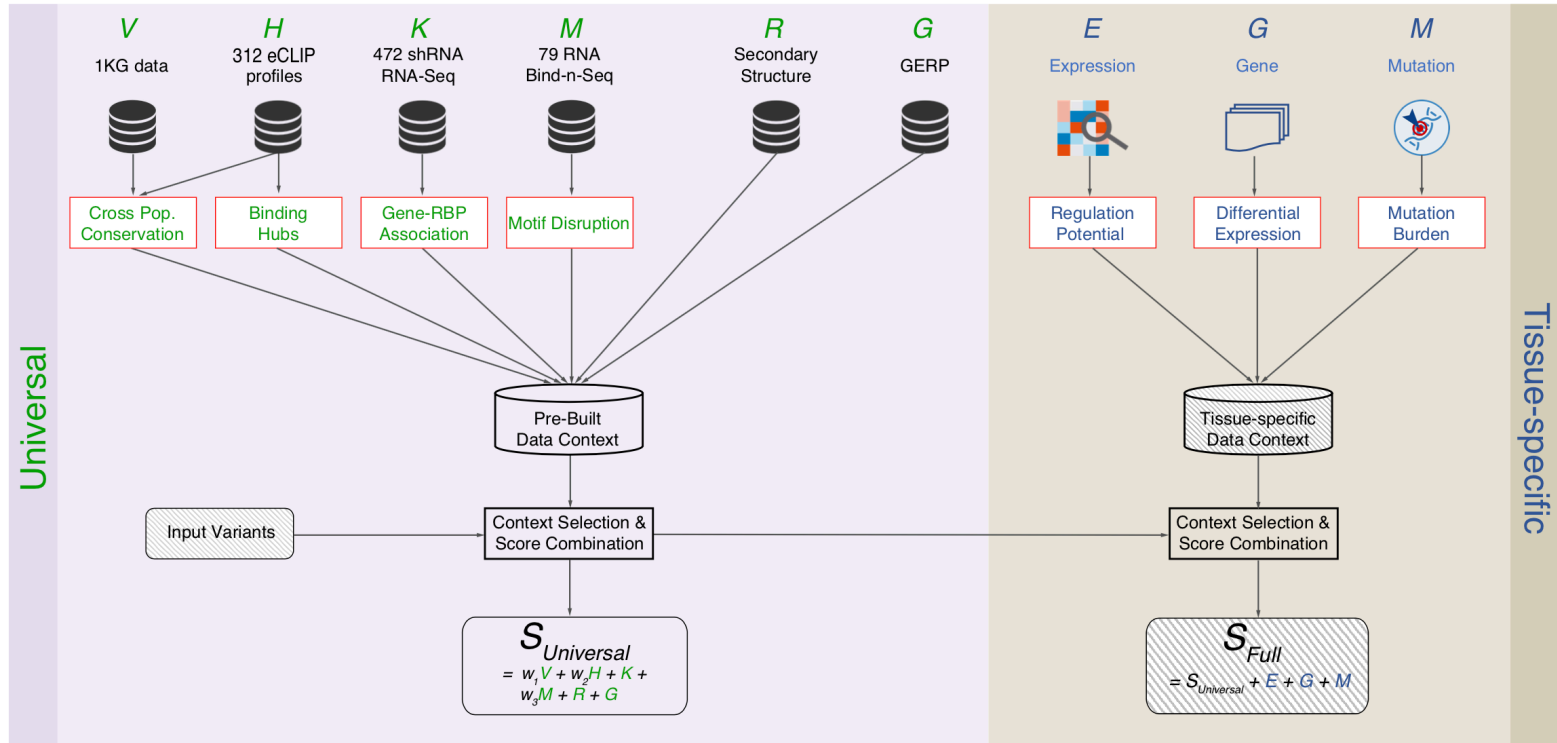


[Khurana et al., Science ('13)]

Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

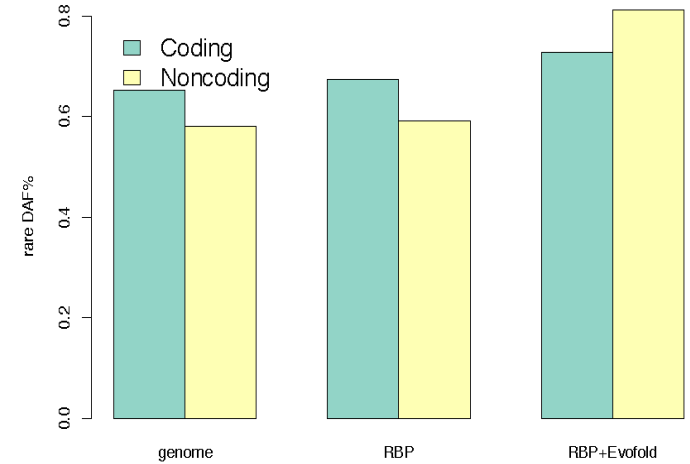
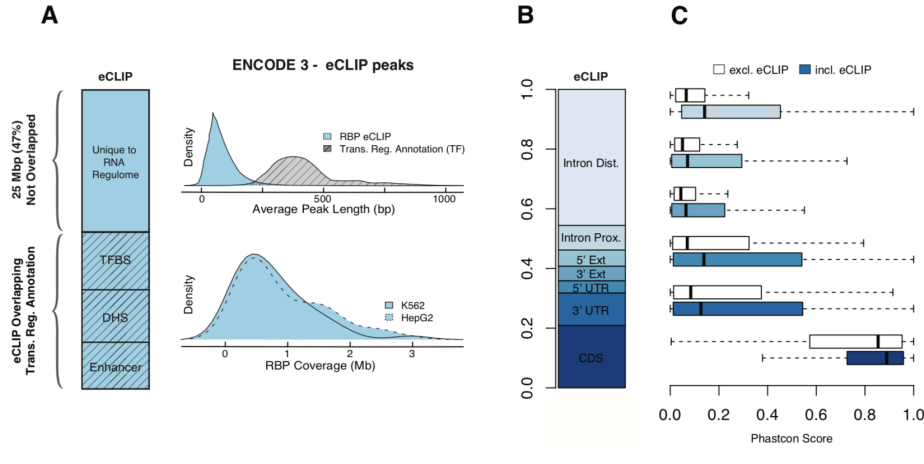
- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - The exponential scaling of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - Frustration as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - FunSeq integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - RADAR: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - uORFs: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

Schematic of RADAR Scoring

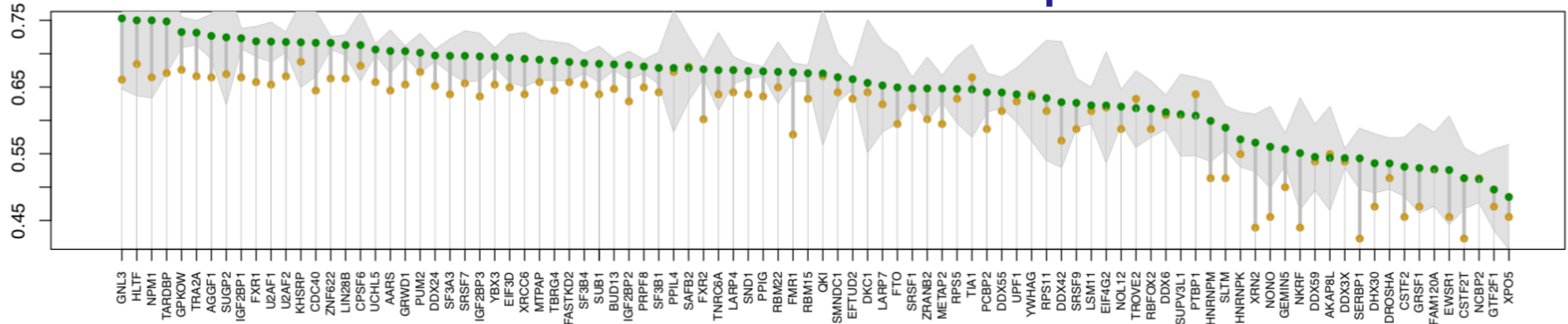


Summary of eCLIP and Phastcon

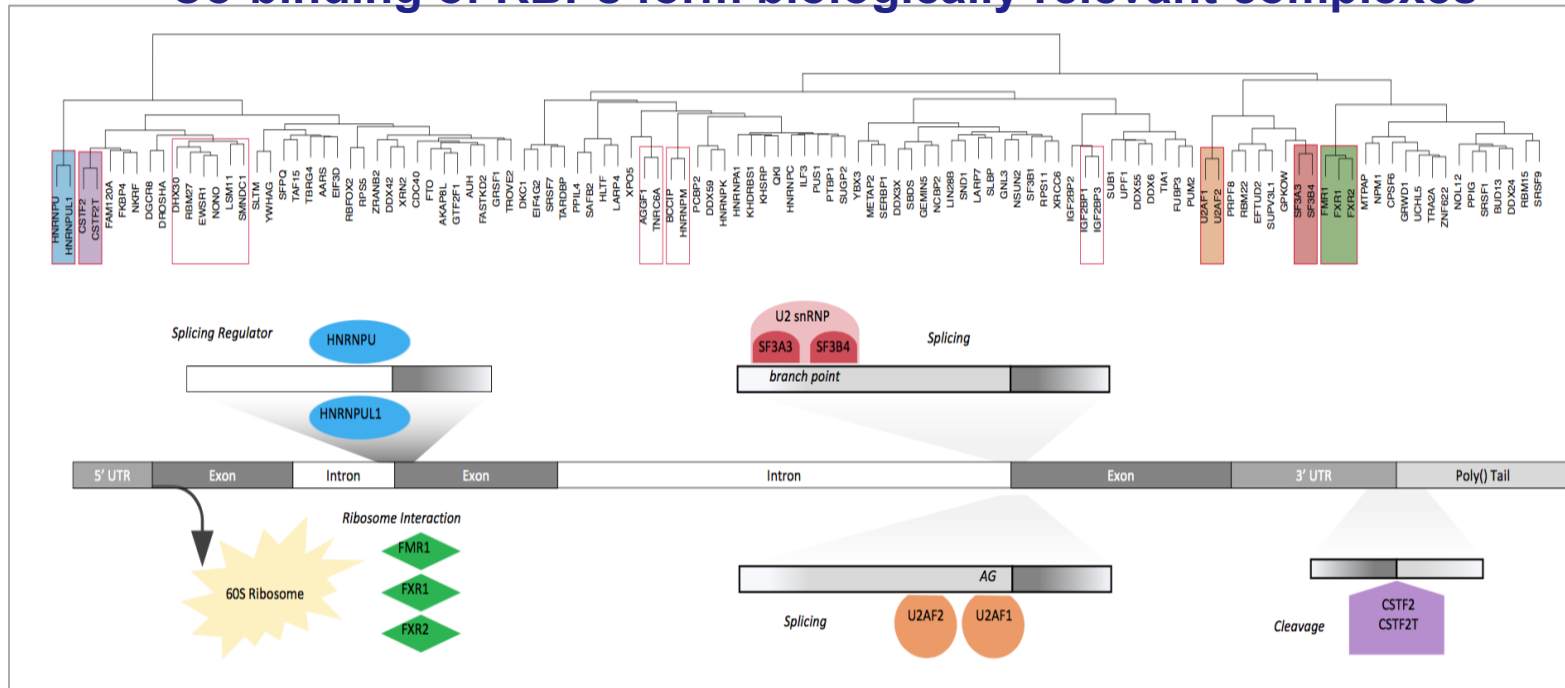
RNA Structure Cons. from Evofold



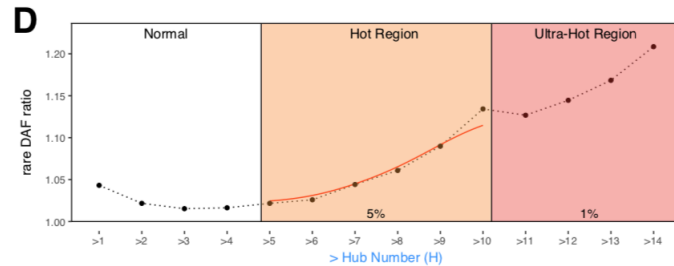
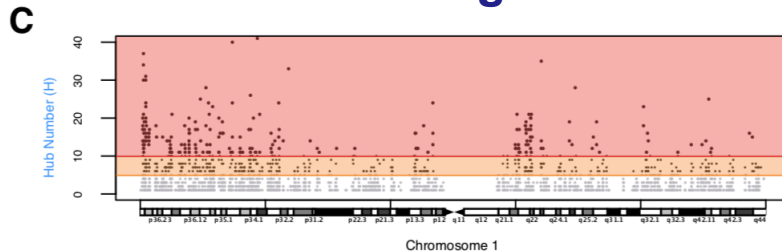
Enriched rare DAF in eCLIP peaks



Co-binding of RBPs form biologically relevant complexes

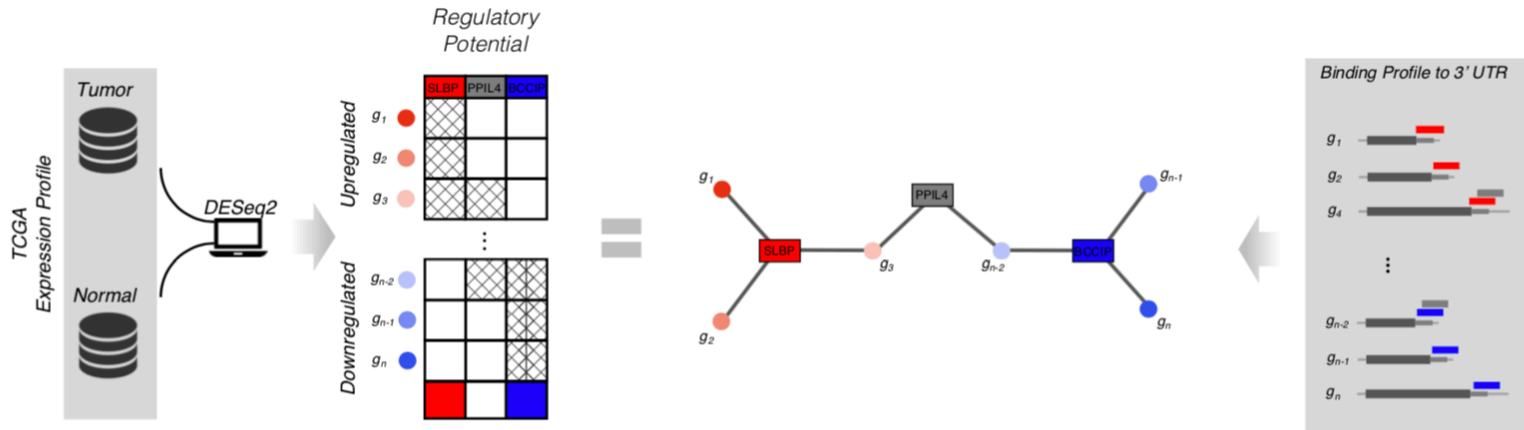


Binding hubs are enriched for rare variants

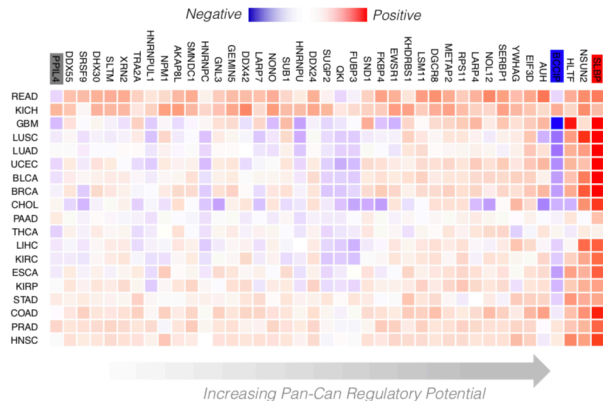


Regulatory Potential of RBPs derived from regression between gene network and expression levels

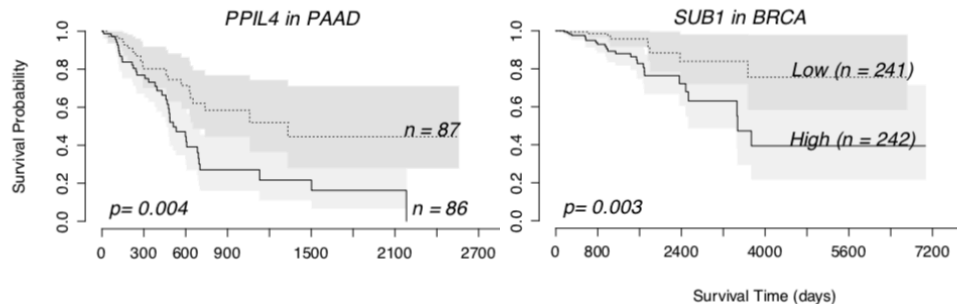
A



B

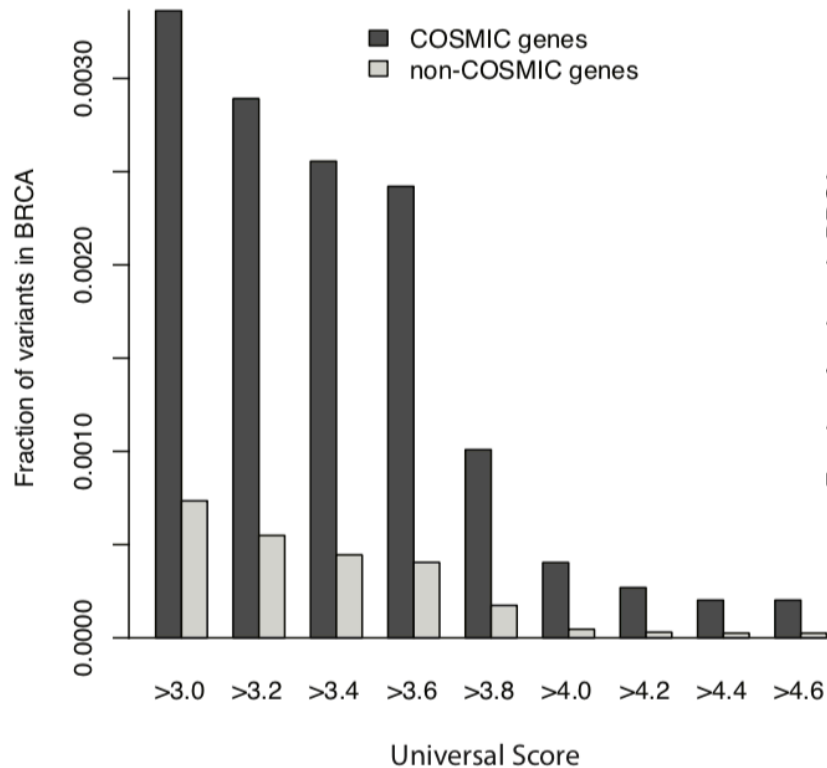


C

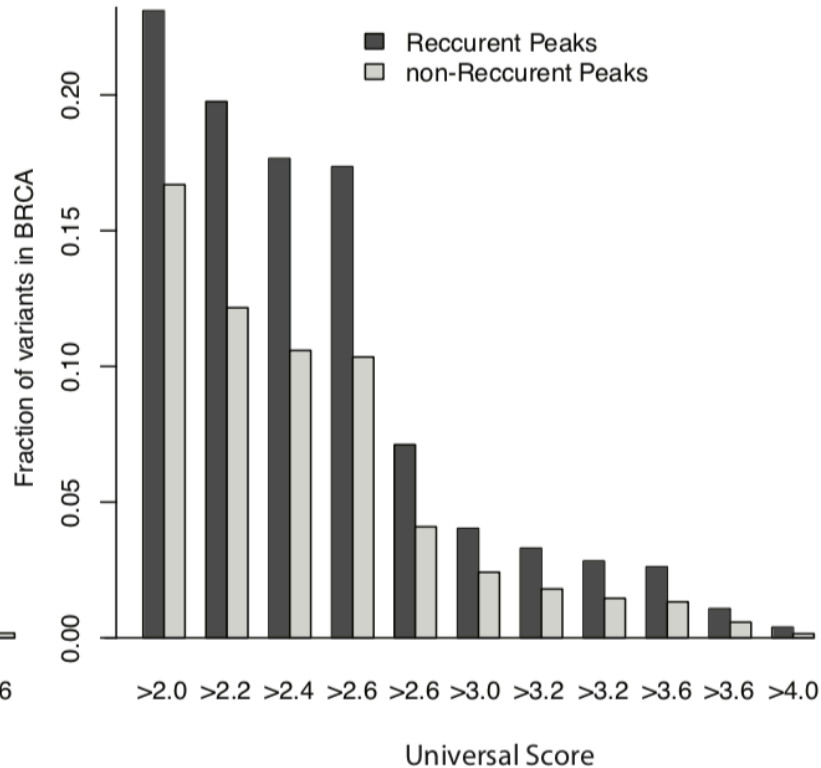


RADAR Scores enriched in COSMIC genes and recurrently mutated regions

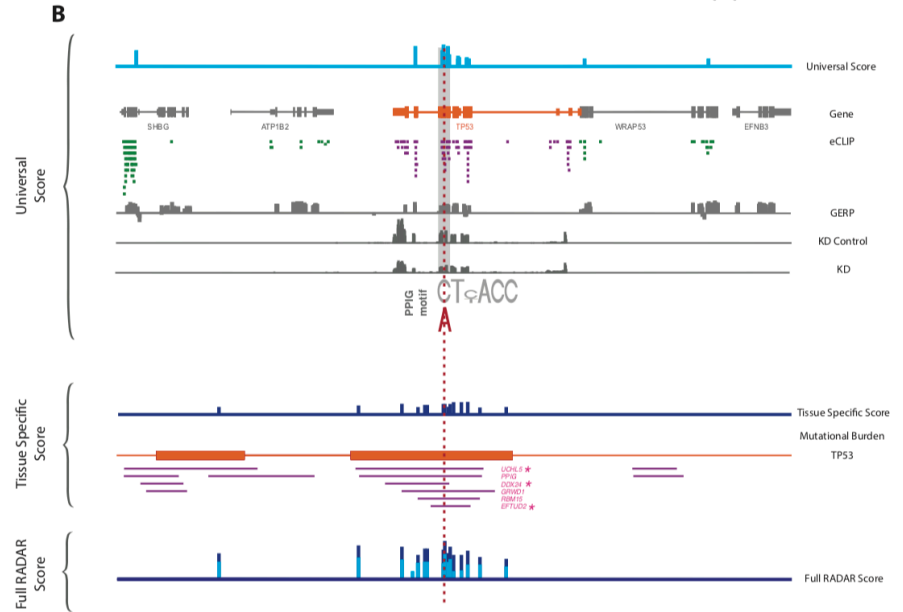
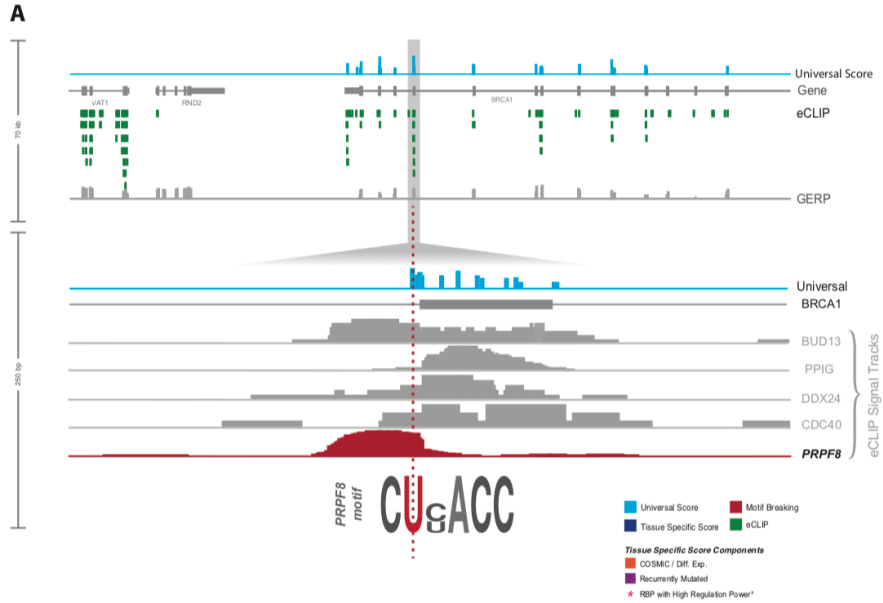
A

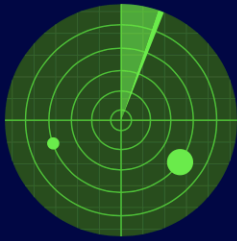


B



Visualization of RADAR Features and Scoring



[HOME](#)[DOWNLOADS](#)[DOCS](#)[EXAMPLE](#)

RADAR can be run on the command line by following the instructions on the [Docs page](#) or through the web here. Running RADAR through the website will print the results after several moments. You can try running RADAR through the web form with a [sample file with one variant](#). Alternatively, you may also input a list of variants into the form as text. If variants are provided in both file and text formats, the variant file will be scored and the text field will be ignored.

More details on the RADAR inputs can be found on the [Docs page](#).

- Variants: a list of variants
 - BED file: a BED file containing the variants
 - Text format: type variants directly into a text box, lines may be tab- or space-delimited
- Cancer type: a TCGA cancer type, only needed if any tissue-specific scores are to be included.
- Tissue-specific scores: which tissue-specific scores should be included along with the universal scores for each variant.

Variants:

no file selected

E.g. chr1 13506 13507 G A

Cancer type:

Select a cancer

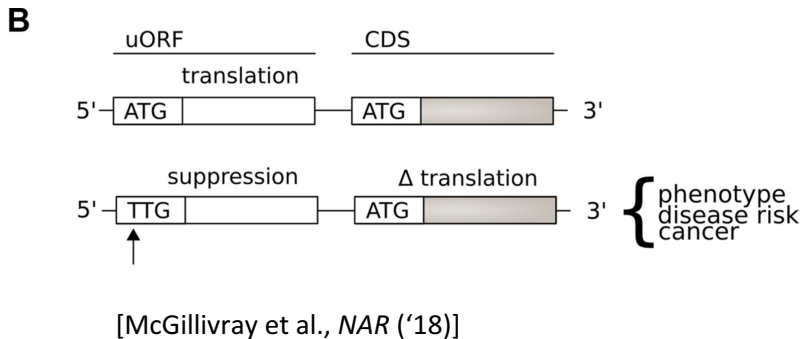
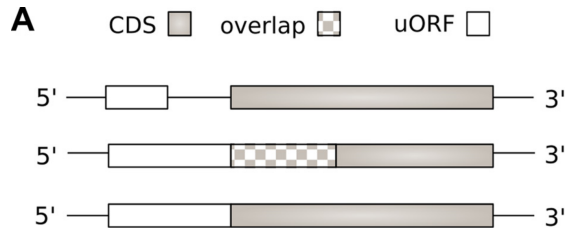
Tissue-specific scores:

- Key genes
- Mutation recurrence
- RBP-regulation power

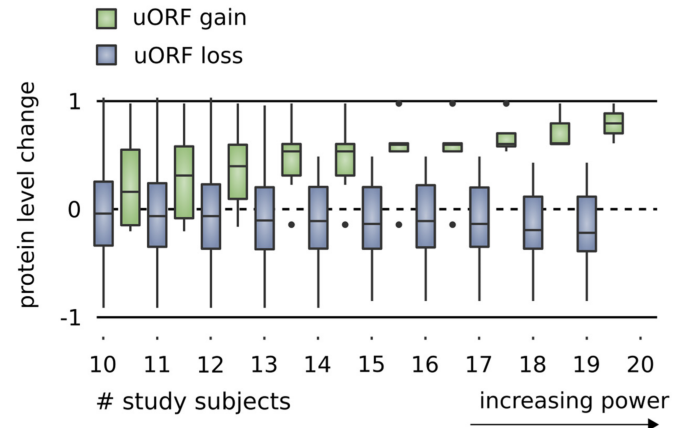
Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

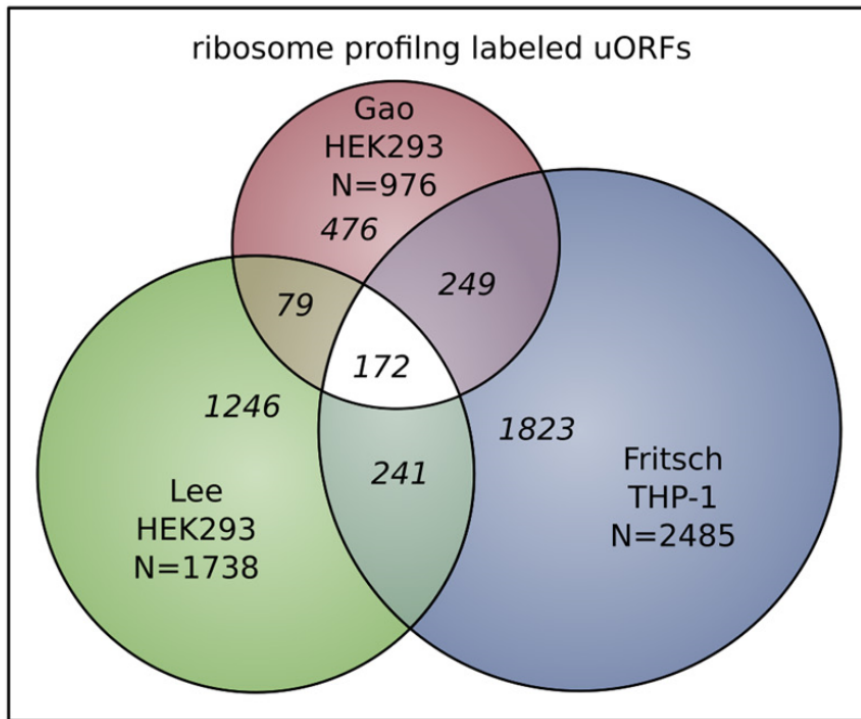
- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - The exponential scaling of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - Frustration as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - FunSeq integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - RADAR: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - uORFs: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

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



- uORFs regulate the translation of downstream coding regions.
- This regulation may be altered by somatic mutation in cancer.
- In Battle et al. 2014 data uORF gain & loss assoc. protein level change.

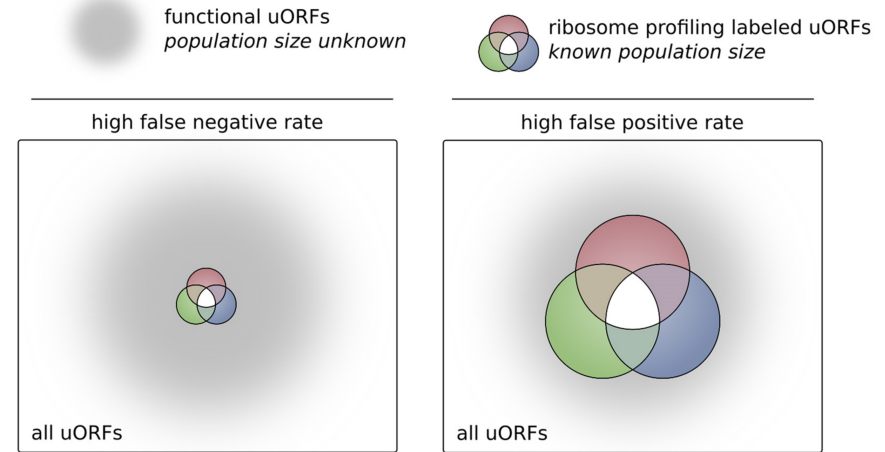




From a “Universe” of
1.3 M pot. uORFs

The population of functional uORFs may be significant

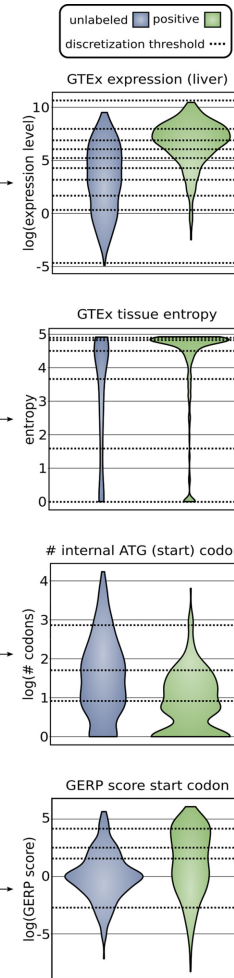
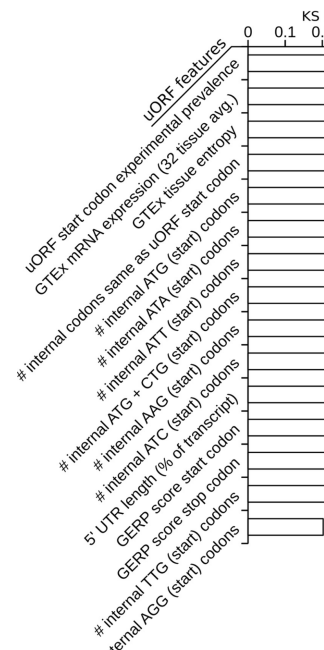
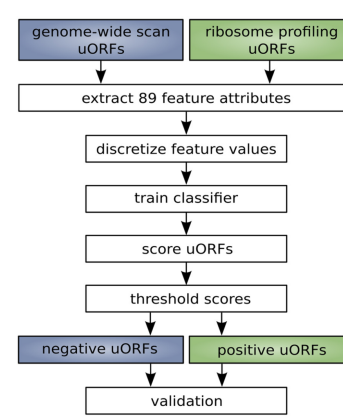
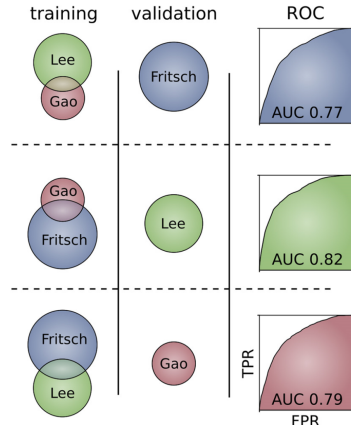
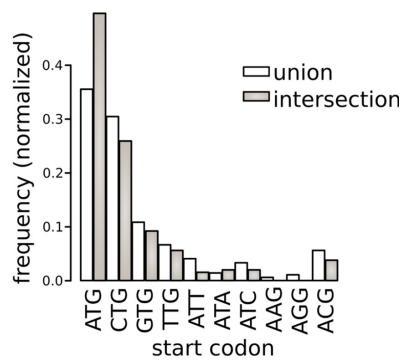
C



- 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

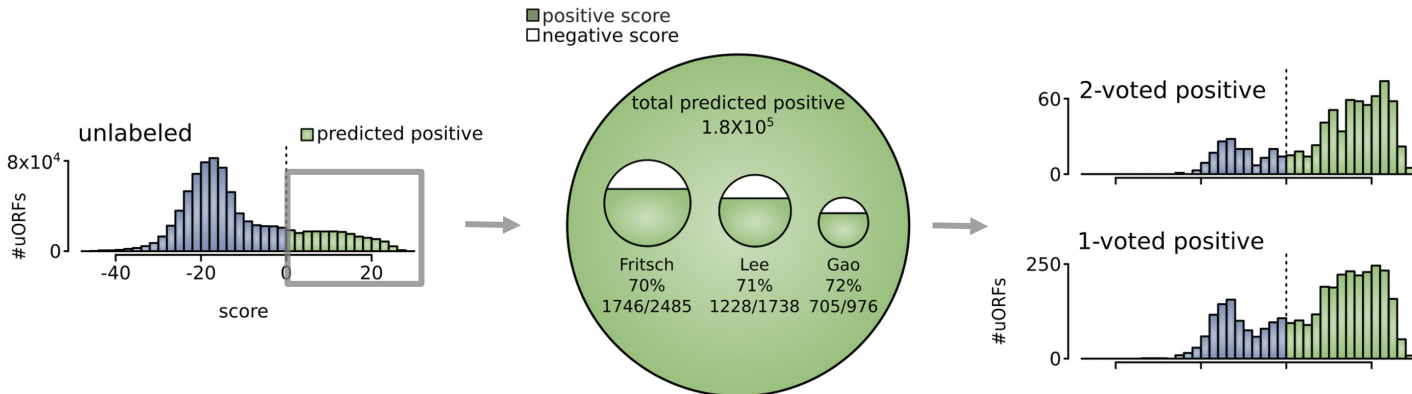
Tissue Dist.

Int. ATG Start

Conservation

A comprehensive catalog of functional uORFs

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

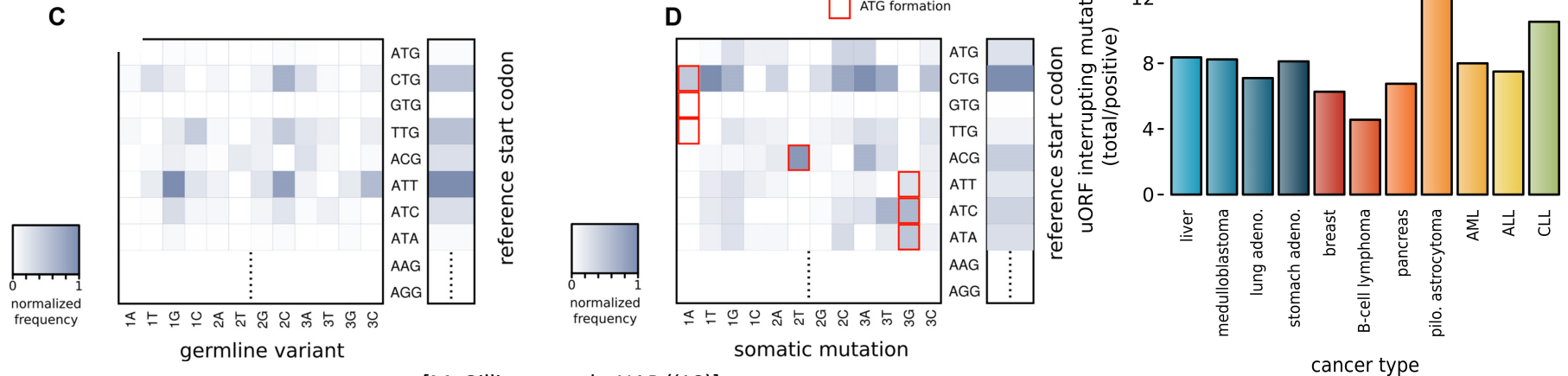


- Predicted functional uORFs may be intersected with disease associated variants.

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

Somatic alteration of uORFs disproportionately affects certain cancers and molecular pathways

- uORF gain and loss occurs in cancer (incl. in cancer associated genes, e.g., MYC, BCL2, etc.).
- Alteration of translation may contribute to cancer.
- These changes are concentrated in certain cancers and pathways.
- Mutations leading to uORFs diff in somatic vs. germline.



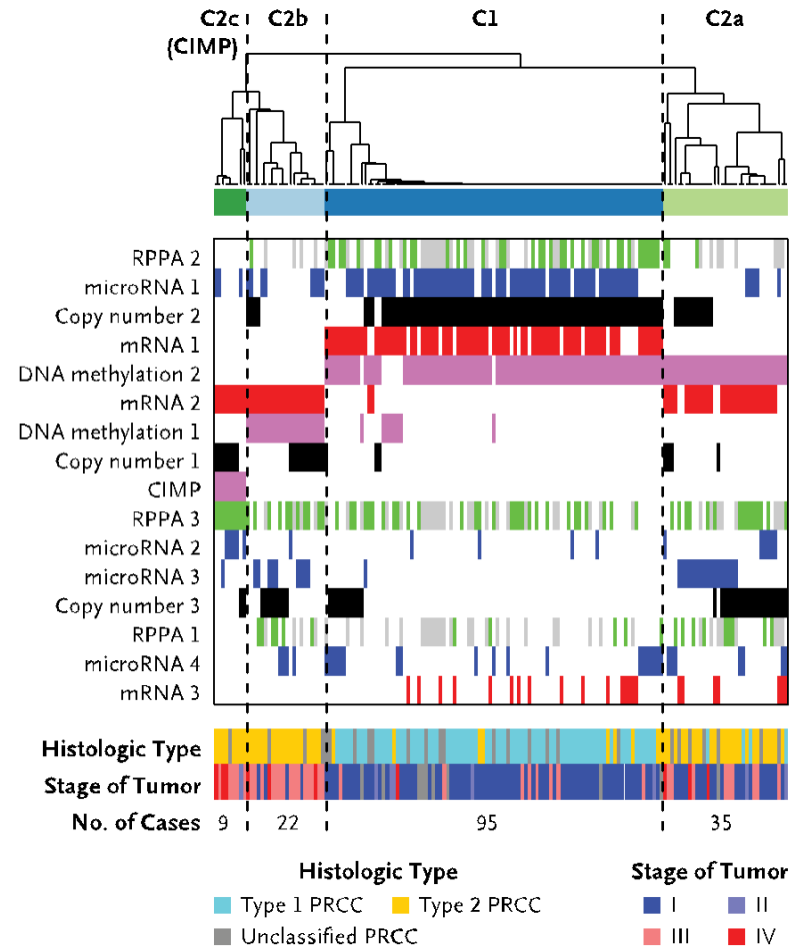
[McGillivray et al., *NAR* ('18)]

Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - The exponential scaling of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - Frustration as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - FunSeq integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - RADAR: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - uORFs: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

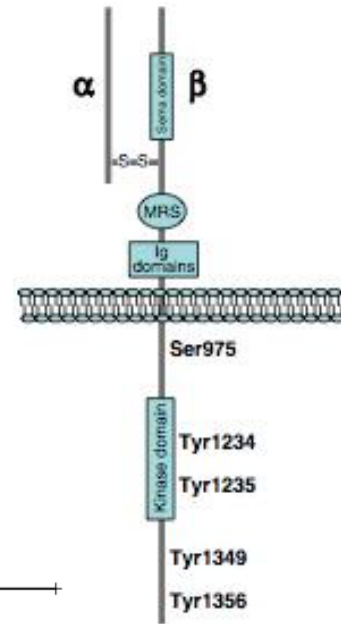
An (underpowered) case study: pRCC

- Kidney cancer lifetime risk of 1.6% & the papillary type (pRCC) counts for ~10% of all cases
- TCGA project sequenced 161 pRCC exomes & classified them into subtypes
 - Yet, cannot pin down the cause for a significant portion of cases....
- 35 WGS of TN pairs, perhaps useful? But not that definitive from a recurrence perspective

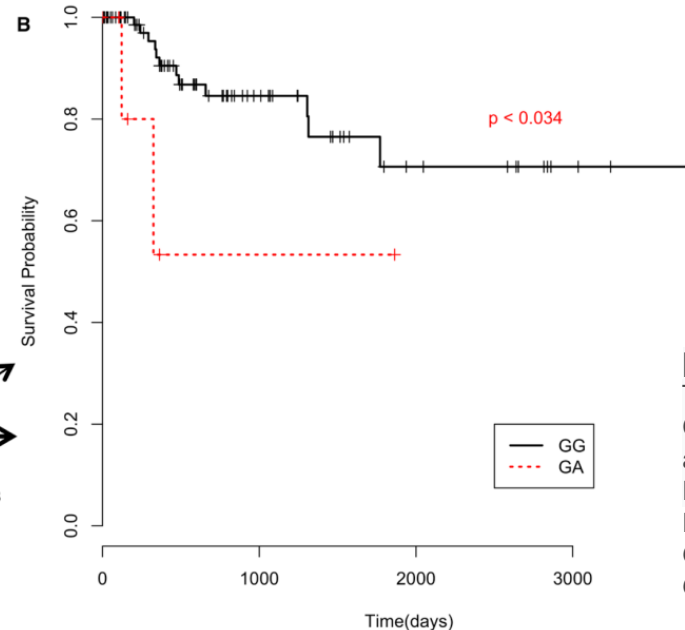
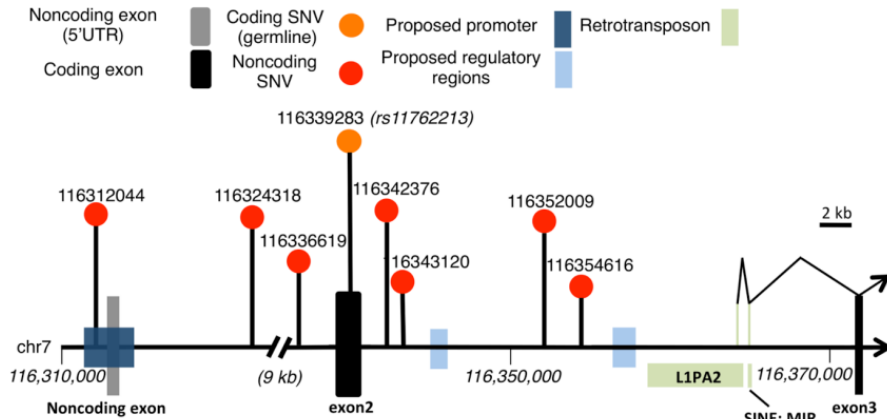


- MET is long known pRCC driver
- In MET, TCGA found somatic SNVs, duplications & an alt. splicing event as drivers (43/161).
- In addition, from 35 WGS we found
 - A noncoding hotspot associated with *MET*
 - Lack of SVs & breakpoints disrupting *MET*
 - Germline SNP (rs11762213) predicts survival in type 2 patients

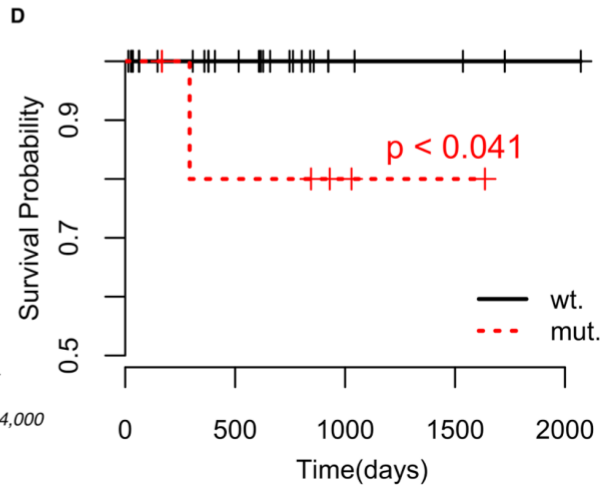
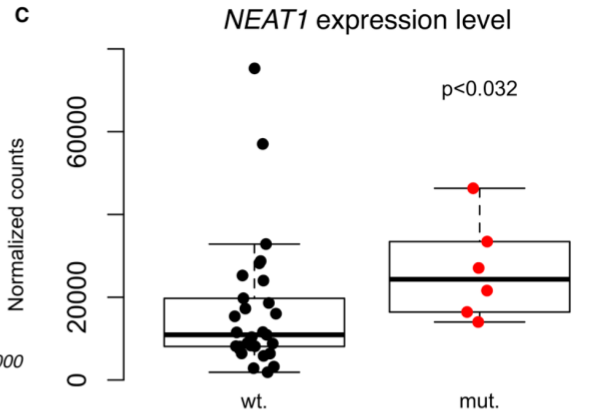
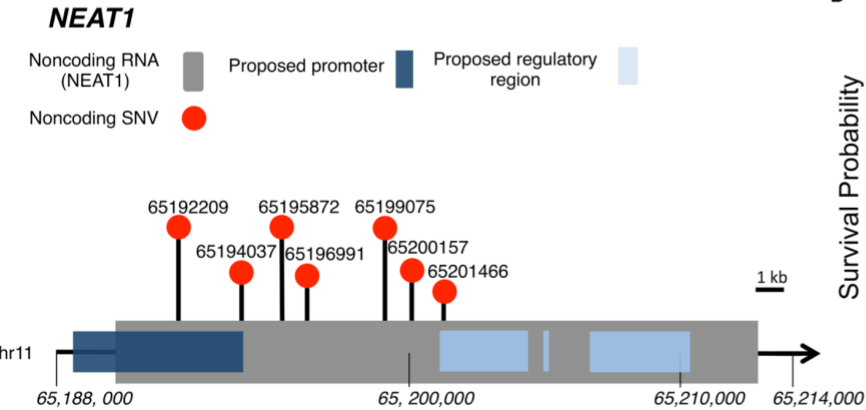
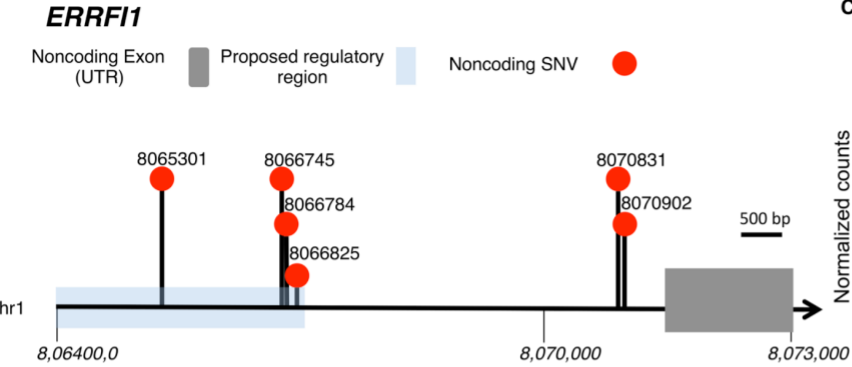
Tyr-kinase MET: Known Facts & New Results



A *MET*



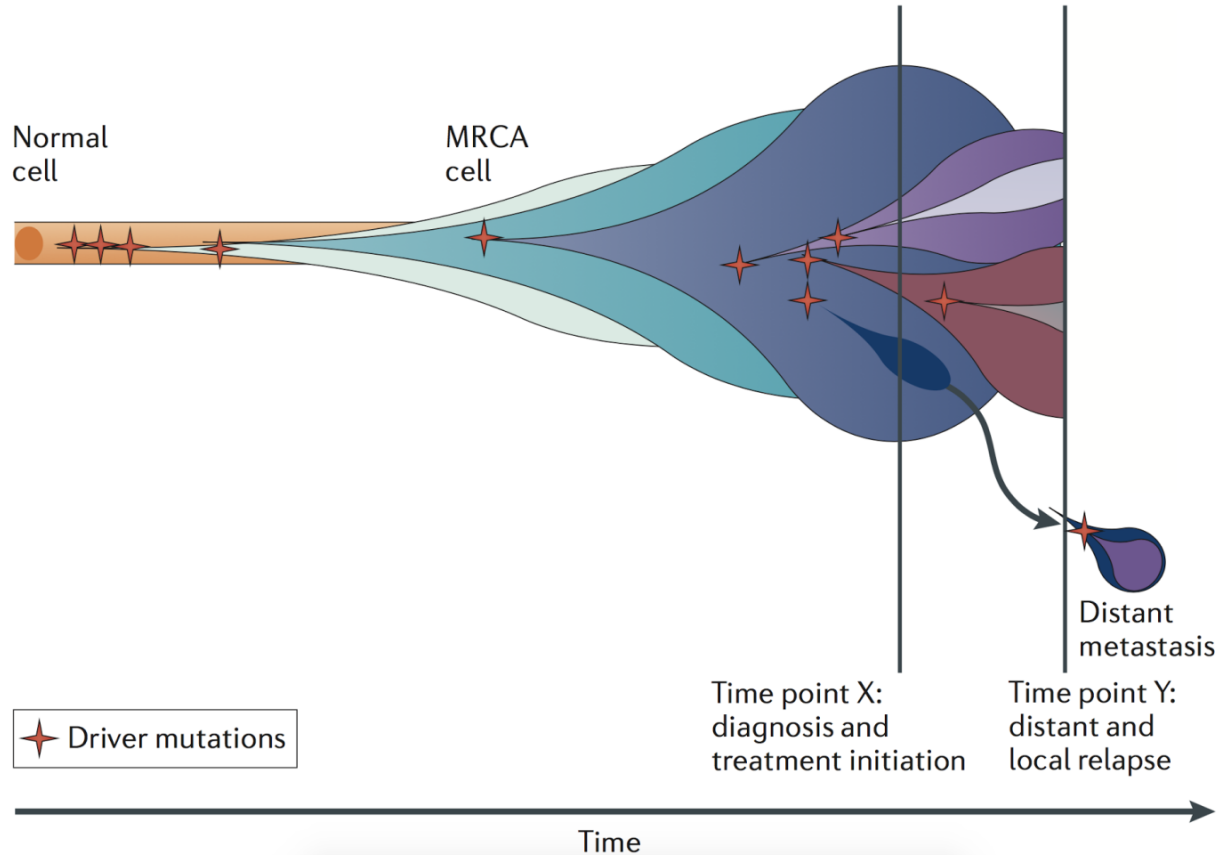
[A. Gentile, L. Trusolino and PM. Comoglio, Cancer and Metastasis Reviews ('08); S. Li, B. Shuch and M. Gerstein PLOS Genetics ('17)]



**Beyond
MET: 2
non-coding
hotspots in
NEAT &
ERRFI1,**

**supported by expr.
changes &
survival
analysis**

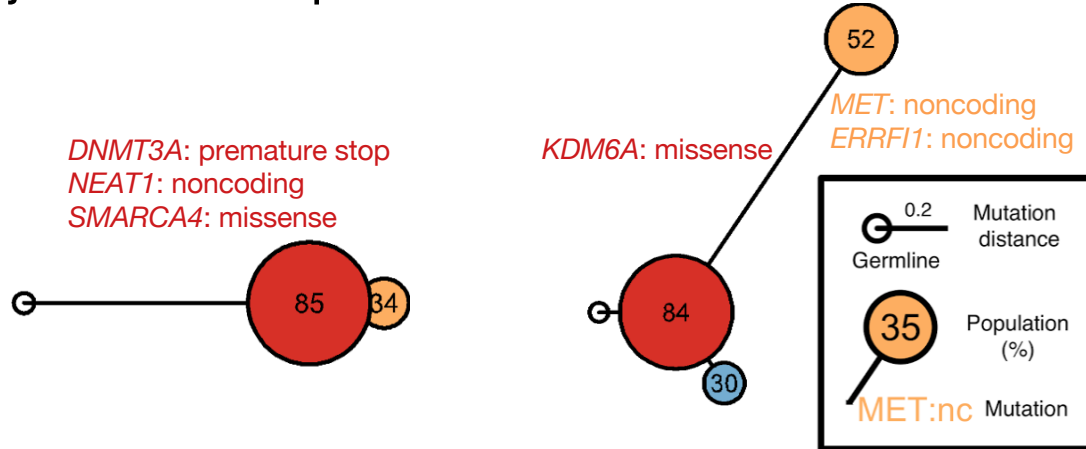
Tumor Evolution: Highlight the Ordering of Key Mutations



Yates et al, NRG (2012)

Construct evolutionary trees in pRCC

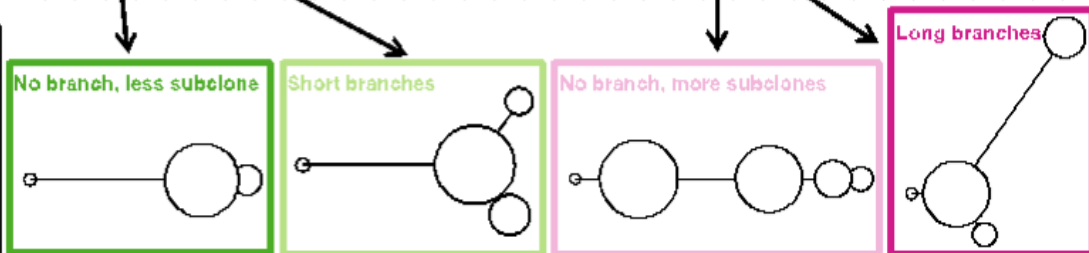
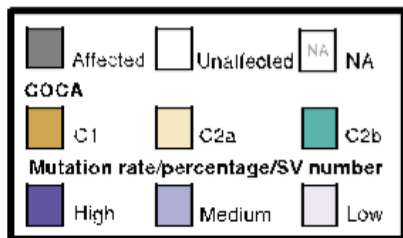
- Infer mutation order and tree structure based on mutation abundance (PhyloWGS, Deshwar et al., 2015)
- Some of the key mutations occur in all the clones while others are just in some parts of the tree



[S. Li, B. Shuch and M. Gerstein PLOS Genetics ('17)]

Tree topology correlates with molecular subtypes

		Type 1										Type 2								Unclassified																
Histological type/Patient ID		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35
COCA		Affected															NA	Affected																		
Coding	MET	Affected										Unaffected										Affected														
	OTHs	Unaffected										Affected										Unaffected														
	MET	Affected										Unaffected										Affected														
	OTHs	Unaffected										Affected										Unaffected														
Noncoding	MET	Affected										Unaffected										Affected														
	OTHs	Unaffected										Affected										Unaffected														
	MET	Affected										Unaffected										Affected														
Mutation Processes	OTHs	Affected										Unaffected										Affected														
	MET	Affected										Unaffected										Affected														
	OTHs	Unaffected										Affected										Unaffected														
	MET	Unaffected										Affected										Unaffected														
Whole genome mutation rate		High										Medium										Low														
DHS mutation percentage		High										Medium										Low														
SV number		High										Medium										Low														
Evolution tree topology		No branch, less subclone										Short branches										No branch, more subclones					Long branches									



Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

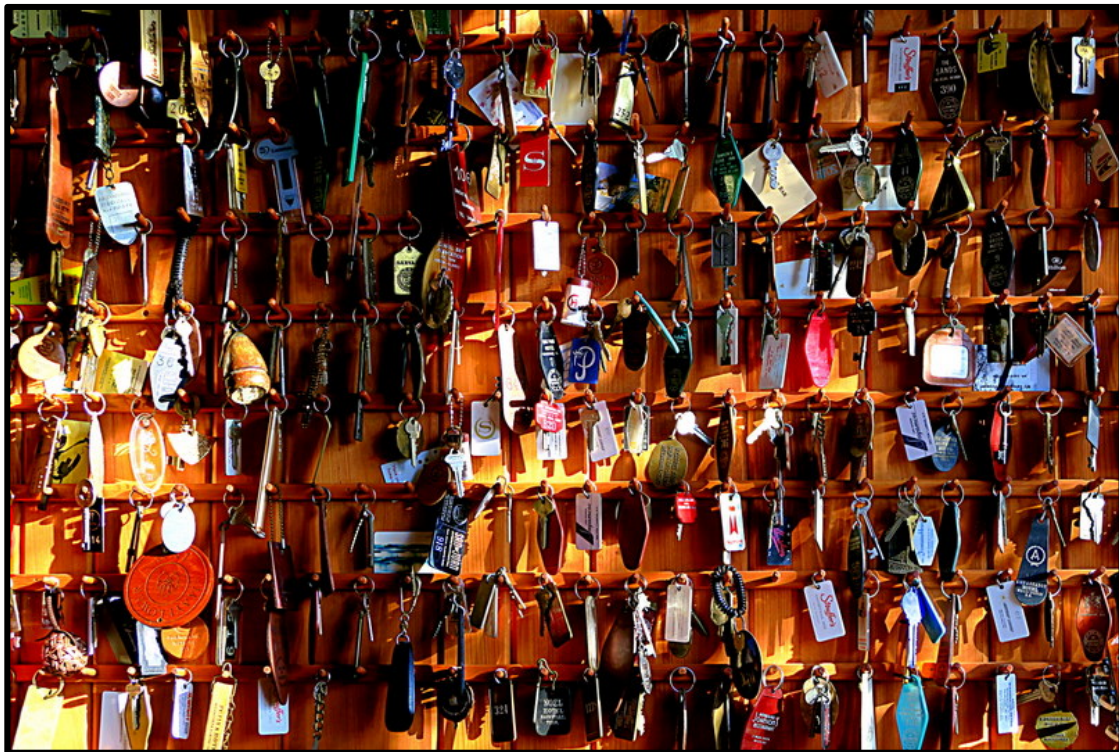
- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - The exponential scaling of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - Frustration as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - FunSeq integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - RADAR: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - uORFs: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

Prioritizing Variants in Personal Genomes: Using functional impact, with particular application to cancer

- Introduction
 - An individual's disease variants as the public's gateway into genomics & biology
 - **The exponential scaling** of data gen. & processing
 - Big-data mining to prioritize key variants as drivers
- Functional impact #1: Coding
 - **Frustration** as a localized metric of SNV impact. Differential profiles for oncogenes v. TSGs
- Functional impact #2: Non-coding
 - **FunSeq** integrates evidence, with a “surprisal” based weighting scheme. Prioritizing rare variants with “sensitive sites” (human conserved)
 - **RADAR**: prioritize variants based on post-transcriptional regulome using ENCODE eCLIP
 - **uORFs**: Feature integration to find small subset of upstream mutations that potentially alter translation
- (Low-power) application to **pRCC**
 - WGS finds additional facts on the canonical driver, MET. Other suggestive non-coding hotspots.
 - Analysis of signatures & tumor evolution helps identify key mutations in different ways 15

github.com/gersteinlab/**Frustration** - S **Kumar**, D Clarke

pRCC - S **Li**, B Shuch



CostSeq2 - P **Muir**, S Li, S Lou, D Wang, DJ Spakowicz, L Salichos, J Zhang, GM Weinstock, F Isaacs, J Rozowsky

github.gersteinlab.org/**uORFs**

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

FunSeq.gersteinlab.org

Y **Fu**, E **Khurana**, Z Liu, S Lou, J Bedford, X Mu, K Yip

RADAR.gersteinlab.org

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



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, 2017.
- Please read permissions statement at
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 . 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)