

Using ENCODE Data for Cancer Genomics



Slides freely downloadable from [Lectures.GersteinLab.org](https://lectures.gersteinlab.org) & “tweetable” (via [@MarkGerstein](https://twitter.com/MarkGerstein)).
No Conflicts for this Talk. See last slide for more info.

BIOSAMPLE →

ENCODEC

86 Cancerous (40 Cancer Types) + 143 Composite Normal (inc. Roadmap)

		K562	HepG2	A549	MCF-7	HeLa-S3	H1-hESC	Caco-2	HCT116	Panc1	LNCaP	PC-3	PC-9	SK-N-MC	DND-41	SK-N-SH	...						
		CML	LIHC	LUAD	BRCA	Cervix	ESC	COAD+READ	PAAD	PRAD	LUAD	SARC	LAML	NB	...								
Chromatin Accessibility	DS	DNase-seq		◆	◆	◆	◆	◆	◆	◆	◆	◆	◆	◆	◆	◆	◆						
Histone Modification	HM	Histone ChIP-seq		19	14	85	16	14	53	3	16	7	1	11	11	8	11	19					
Transcription	TX	RNA-seq		◆	◆	◆	◆	◆	◆	▼	◆	◆	▼	◆	▼	▼	▼	◆					
		RAMPAGE		◆																			
RNA-binding Proteins	RP	eCLIP		191	164																		
RNAi/CRISPR Knockdown	KD	shRNA/siRNA KD		326	257	2																	
		CRISPR KD/KO		108	19																		
3D Chromatin Structure	3D	ChIA-PET		9	2		5	1															
		Hi-C		▼	◆	◆	▼	◆	▼														
Enhancers	SS	STARR-seq		◆	◆		◆																
Methylation	ME	WGBS		◆	◆	◆	▼	◆	◆														
		RRBS		◆	◆	◆	◆	◆	◆														
Replication Timing	RT	Repli-chip						◆	◆														
		Repli-seq		◆	◆	◆	◆	◆															
Transcription Factors	TF	TF ChIP-seq		558	300	240	149	78	89														
Cell Line WGS	WG	SNV		▼		▼	▼	▼															
		SV		▼		▼	▼	▼															

528 ENCODE Cell Types	→	229 Deduplicated & Selected Human Biosamples
-----------------------	---	--

[Zhang et al. (‘20), biorxiv + Nat. Comm. (in press)]

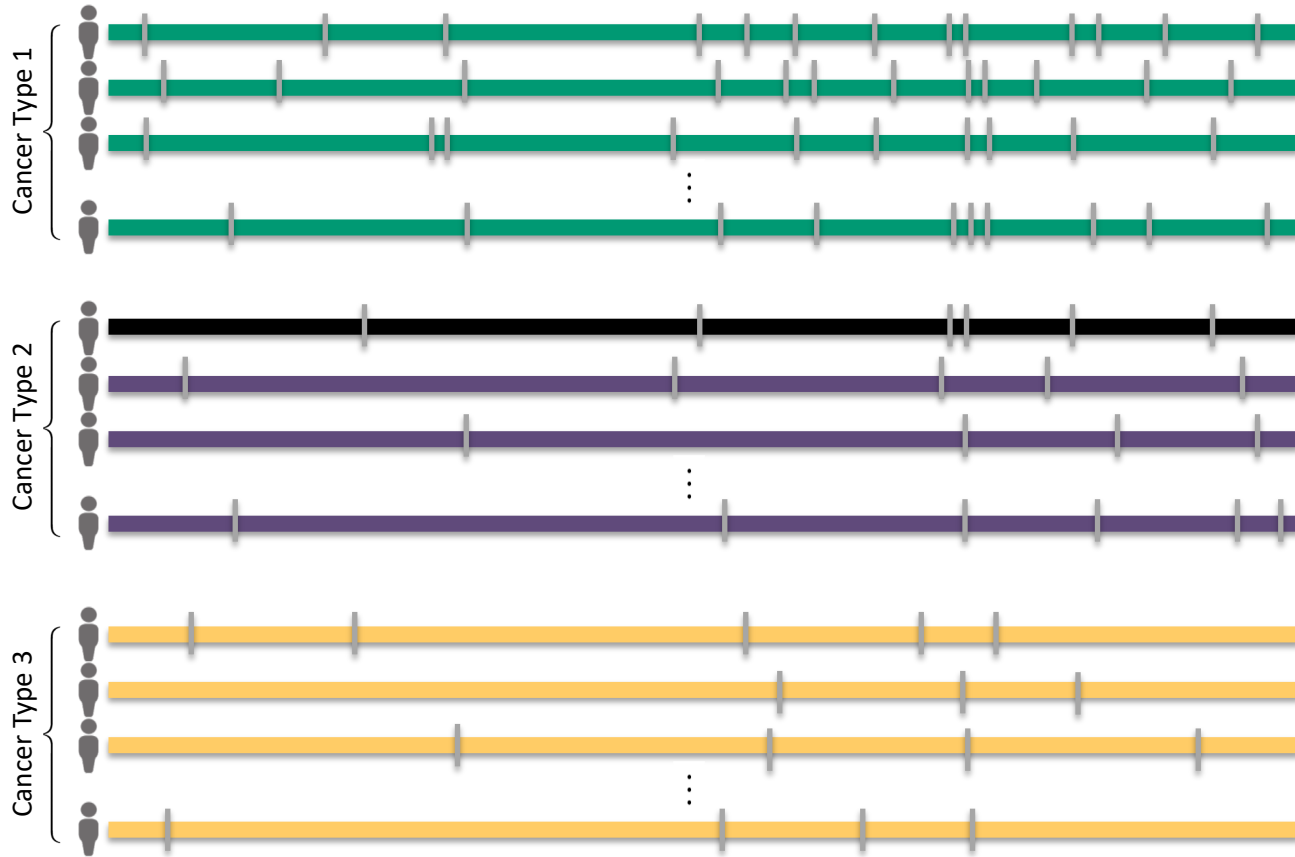
Using ENCODE Data for Cancer Genomics

- **BMR Correction:**
LARVA/MOAT/NIMBUS
 - Parametric models explicitly modeling genomic covariates
 - Many ENCODE covariates useful in accurately estimating background mutation rate
- **Network Rewiring in Cancer**
 - Large-scale ENCODE chip-seq data in certain cell lines highlights TFs changing targets greatly in oncogenesis. (Focus on CML)
 - TopicNet LDA approach (from text-mining) finds regulators that greatly change their gene communities
- **RADAR Variant Prioritization**
 - Prioritizes germline & somatic variants based on post-transcriptional regulome using ENCODE eCLIP
 - Incorporates new features related to RNA sec. struc & tissue specific effects
- **Regulatory Drivers of Differential Expression**
 - Highlighting regulators in terms of their power to drive differential expression.
 - Relationship of this to network hierarchy & RBP-TF cross talk
 - Example of MYC & SUB1

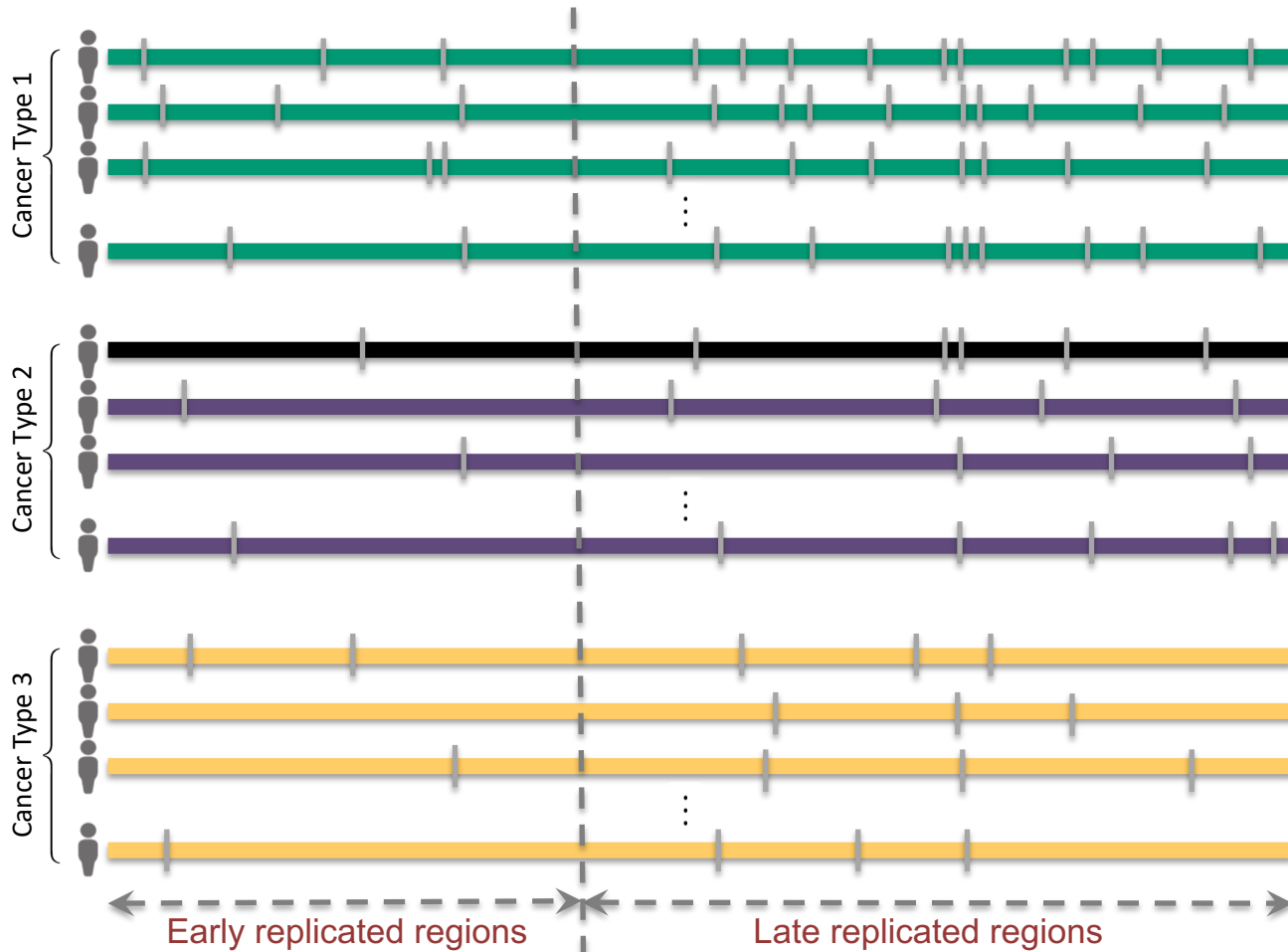
Using ENCODE Data for Cancer Genomics

- **BMR Correction:**
LARVA/MOAT/NIMBUS
 - Parametric models explicitly modeling genomic covariates
 - Many ENCODE covariates useful in accurately estimating background mutation rate
- **Network Rewiring in Cancer**
 - Large-scale ENCODE chip-seq data in certain cell lines highlights TFs changing targets greatly in oncogenesis. (Focus on CML)
 - TopicNet LDA approach (from text-mining) finds regulators that greatly change their gene communities
- **RADAR Variant Prioritization**
 - Prioritizes germline & somatic variants based on post-transcriptional regulome using ENCODE eCLIP
 - Incorporates new features related to RNA sec. struc & tissue specific effects
- **Regulatory Drivers of Differential Expression**
 - Highlighting regulators in terms of their power to drive differential expression.
 - Relationship of this to network hierarchy & RBP-TF cross talk
 - Example of MYC & SUB1

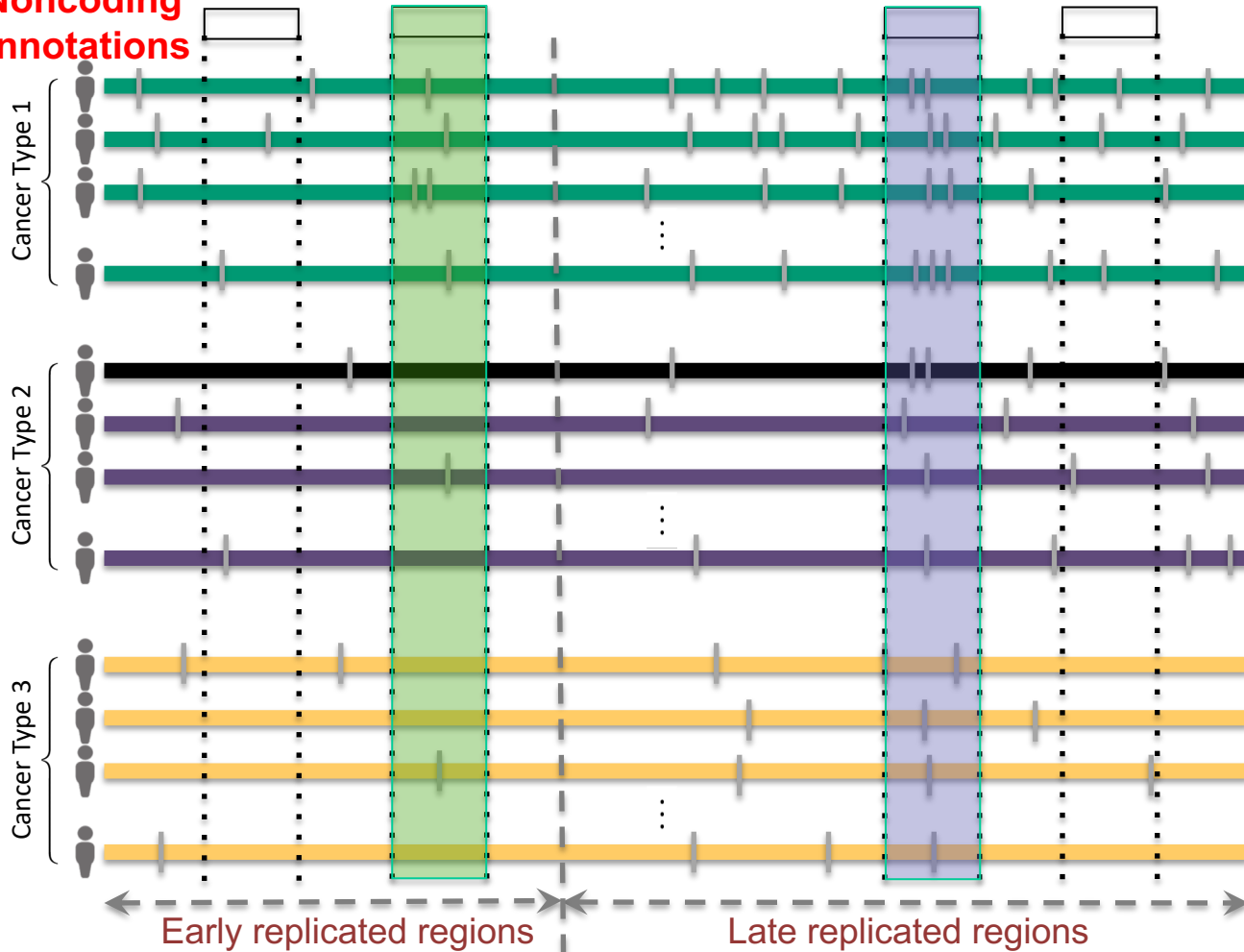
Mutation recurrence



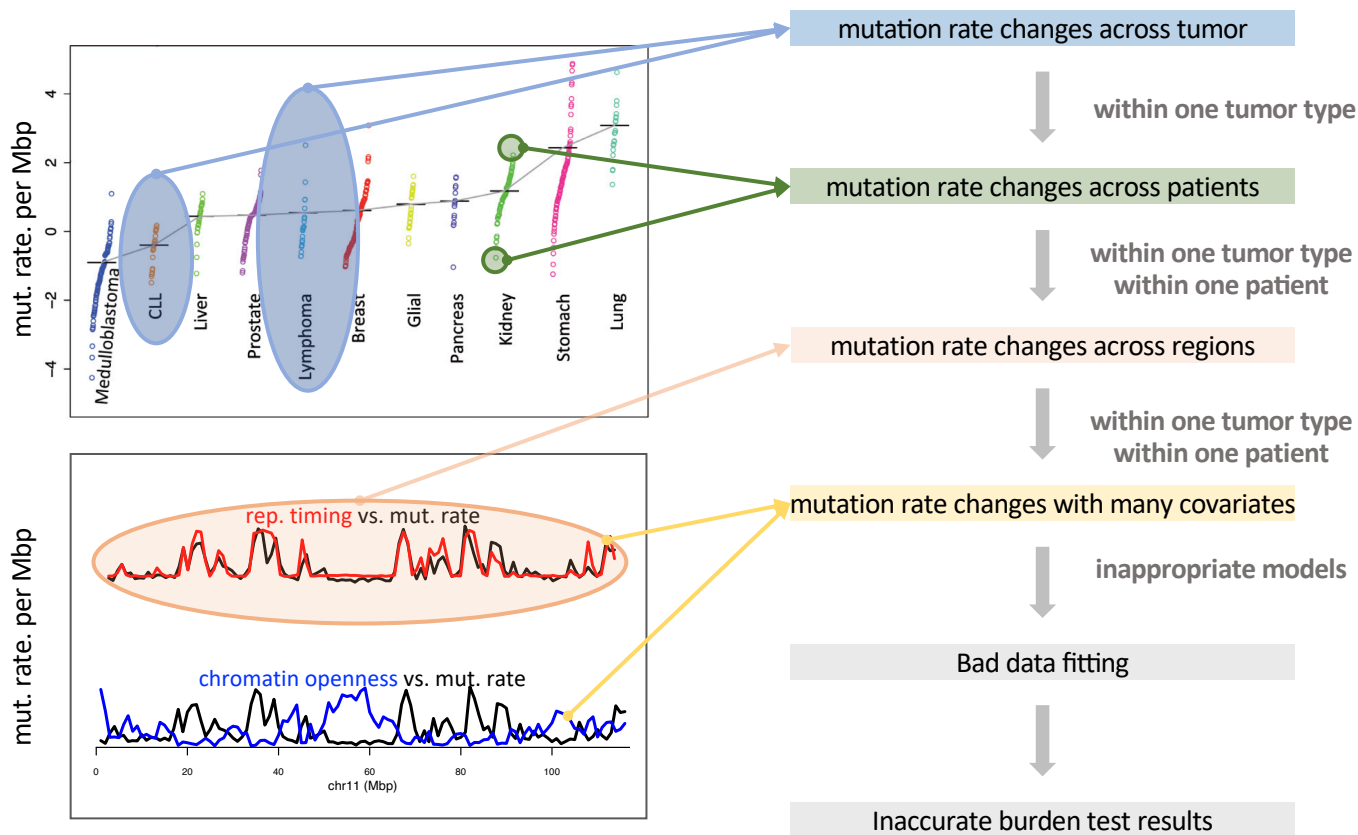
Mutation recurrence

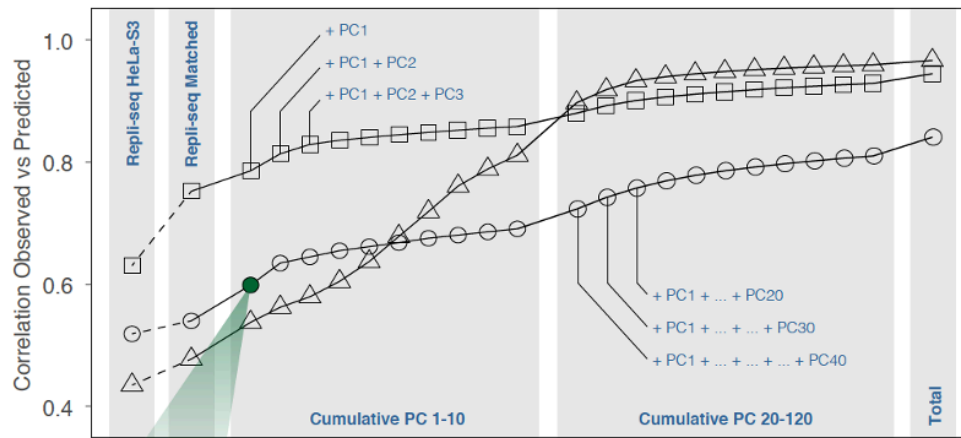


Noncoding annotations

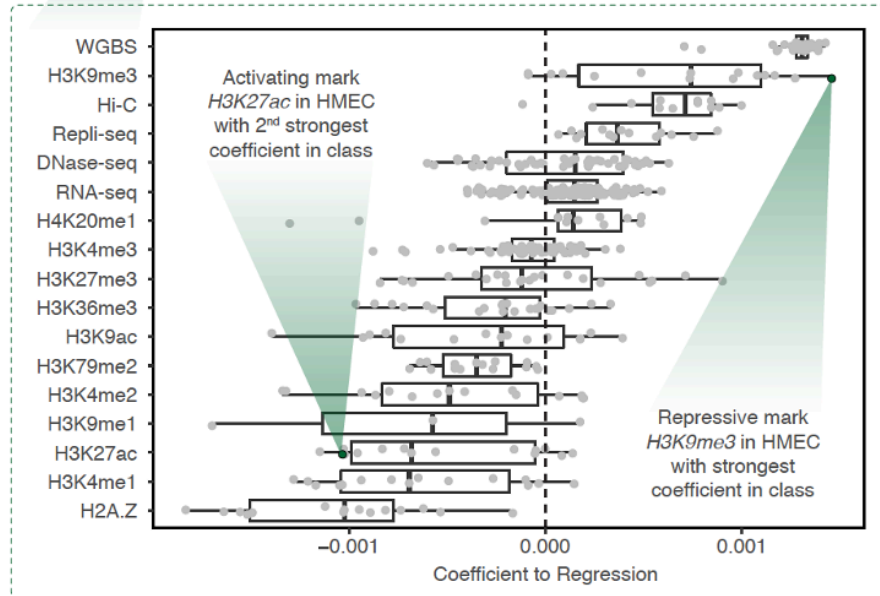


violation of the constant mutation rate assumption





Accurately modeling background mutation rate with full spectrum of ENCODE data



[Zhang et al. *Nat. Comm.* ('20);
Zhang et al. *bioRxiv* + *BMC Bioinfo* ('20), in press]

Cancer Somatic Mutation Modeling

PARAMETRIC MODELS (LARVA/NIMBUS)

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

$$x_i : \text{Binomial}(n_i, p)$$

Model 2: Varying Mutation Rate with Covariate Correction (Beta Binom.)

$$x_i : \text{Binomial}(n_i, p_i)$$

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

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

Model 3: Varying Mutation Rate with Covariate Correction (Neg. Binom.)

$$x_i | p_i \sim \text{Pois}(p_i)$$

$$p_i \sim \text{gamma}(\mu_i, \theta_i)$$

$$\log(\mu_i) \sim \beta_0 + \beta_1 v_1 + \dots + \beta_k v_k$$

- Suppose there are L genome elements. For element i , define:
 - n_i : total number of nucleotides
 - x_i : the number of mutations within the element
 - p : the mutation rate
 - R_i & v_k : covariates
- Non-parametric model is useful when covariate data is missing for the studied annotations
 - Also sidesteps issue of properly identifying and modeling every relevant covariate (possibly hundreds)

NON-PARAMETRIC MODELS (MOAT)

Assume constant background mutation rate in local regions.

Model 3a: Random Permutation of Input Annotations

Shuffle annotations within local region to assess background mutation rate.

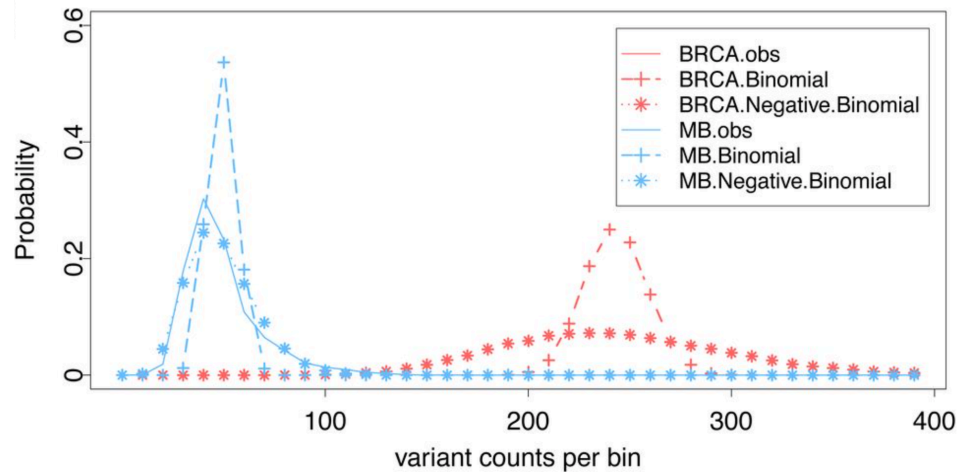
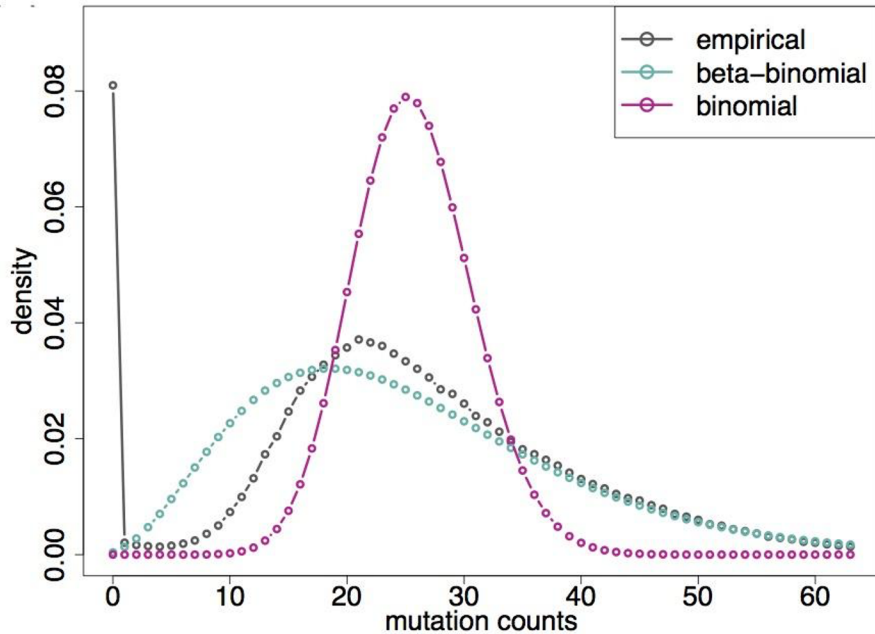
Model 3b: Random Permutation of Input Variants

Shuffle variants within local region to assess background mutation rate.

[Lochovsky et al. *Bioinformatics* ('17)]

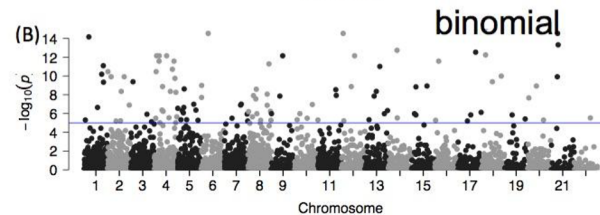
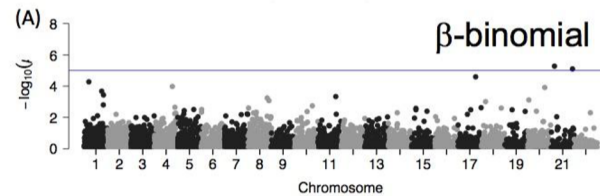
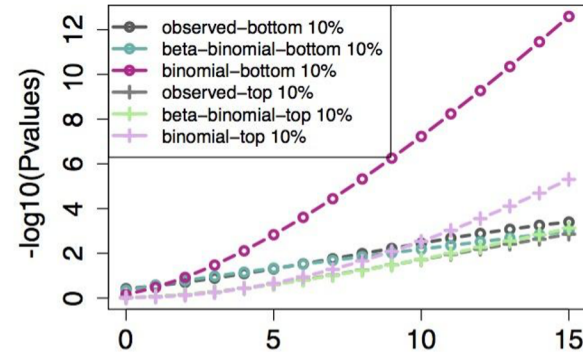
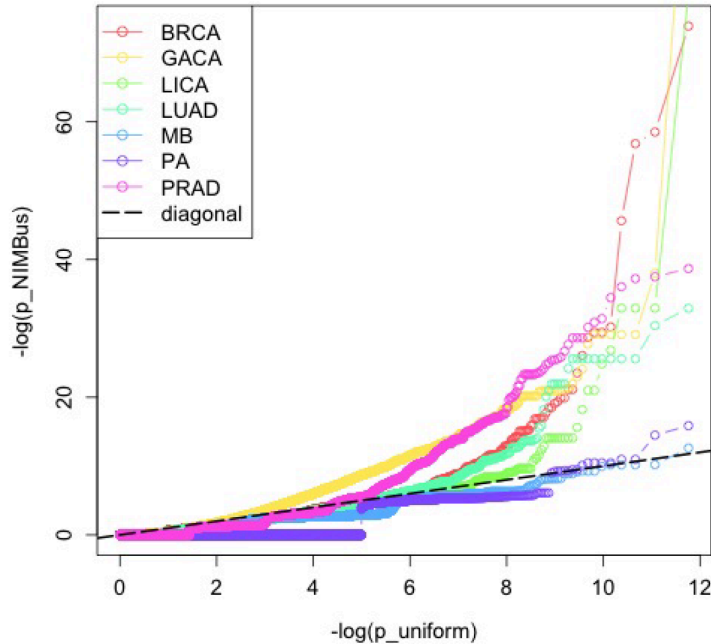
LARVA/NIMBUS Model Comparison

- Comparison of mutation count frequency implied by the binomial model (model 1) and the beta-binomial model (model 2) relative to the empirical distribution
- The beta-binomial/negative binomial distribution is significantly better, especially for accurately modeling the over-dispersion of the empirical distribution



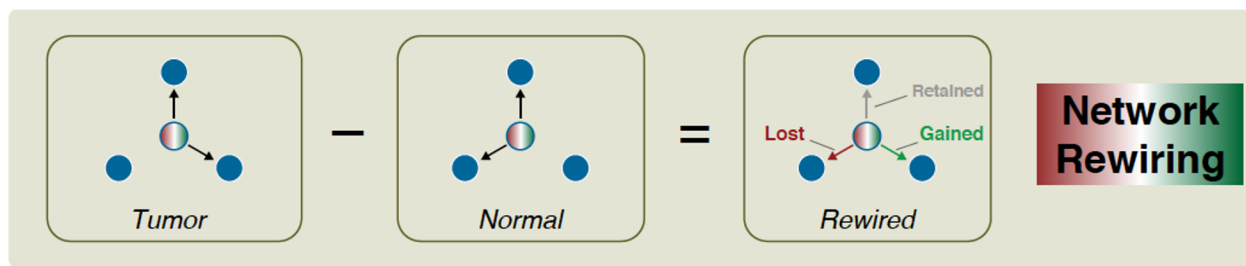
LARVA/NIMBUS Results: Reducing P-value inflation

(B)

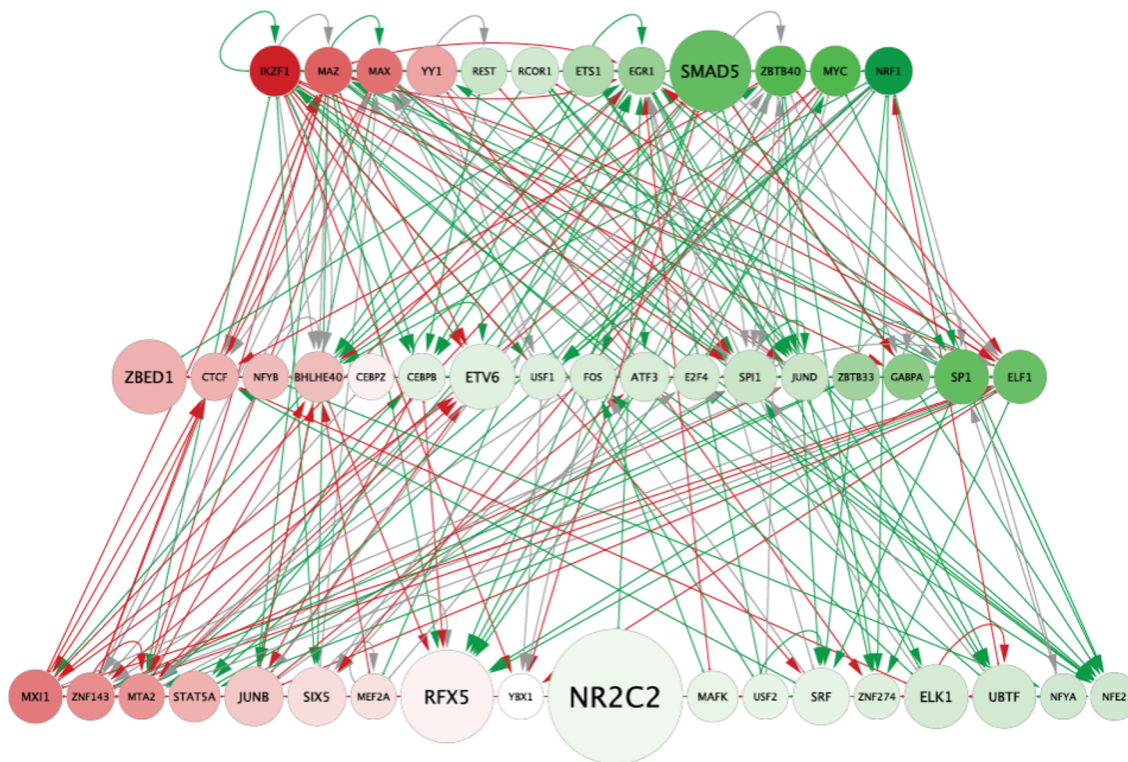


Using ENCODE Data for Cancer Genomics

- **BMR Correction:**
LARVA/MOAT/NIMBUS
 - Parametric models explicitly modeling genomic covariates
 - Many ENCODE covariates useful in accurately estimating background mutation rate
- **Network Rewiring in Cancer**
 - Large-scale ENCODE chip-seq data in certain cell lines highlights TFs changing targets greatly in oncogenesis. (Focus on CML)
 - TopicNet LDA approach (from text-mining) finds regulators that greatly change their gene communities
- **RADAR Variant Prioritization**
 - Prioritizes germline & somatic variants based on post-transcriptional regulome using ENCODE eCLIP
 - Incorporates new features related to RNA sec. struc & tissue specific effects
- **Regulatory Drivers of Differential Expression**
 - Highlighting regulators in terms of their power to drive differential expression.
 - Relationship of this to network hierarchy & RBP-TF cross talk
 - Example of MYC & SUB1

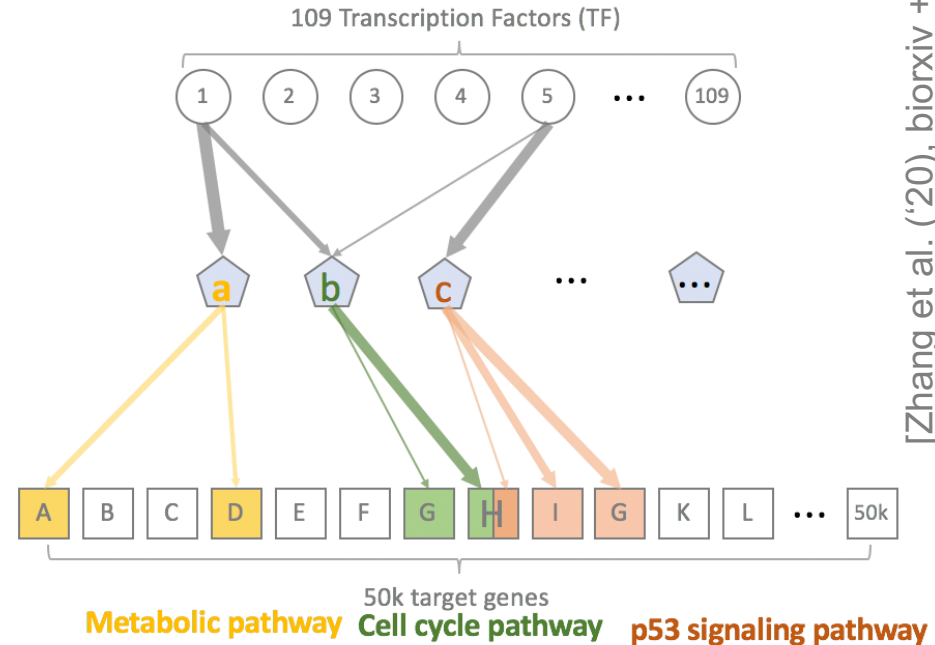
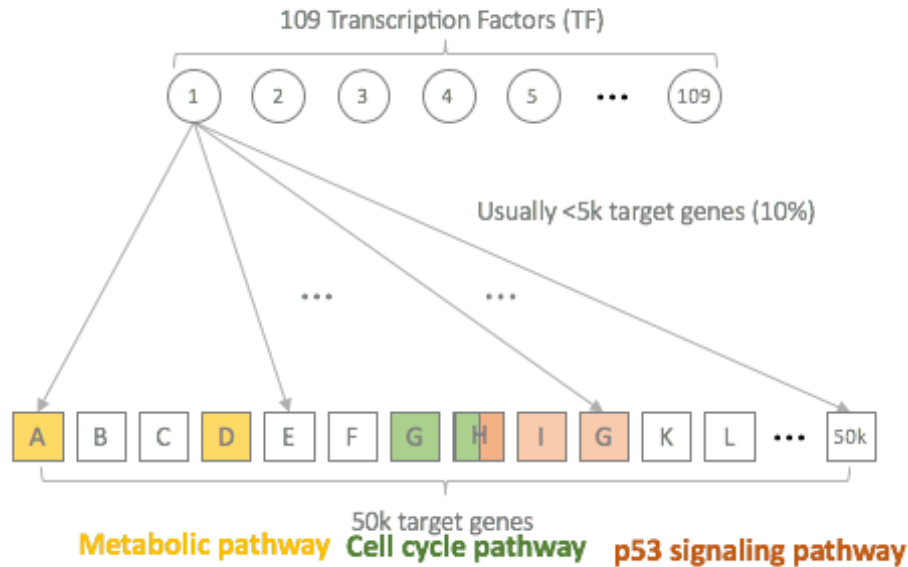


Rewired edges in comparison of GM12878 to K562 109 node TF-TF network (approx. CML)

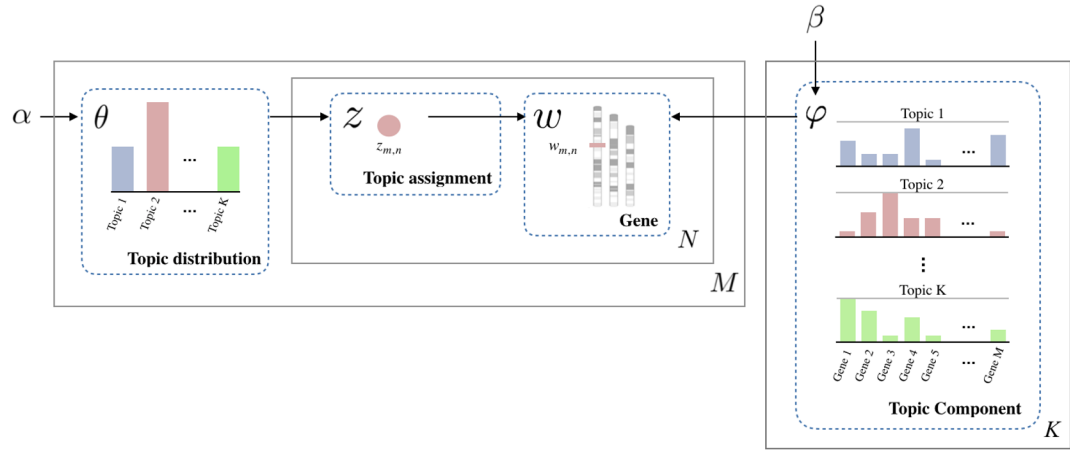
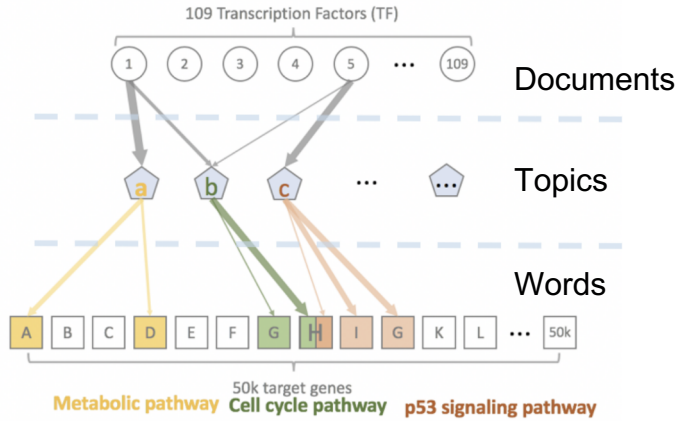


Simplifying Network Rewiring

From $TF \rightarrow gene$ ($109 \times 50,000$)
to $TF \rightarrow pathway$ (109×50)



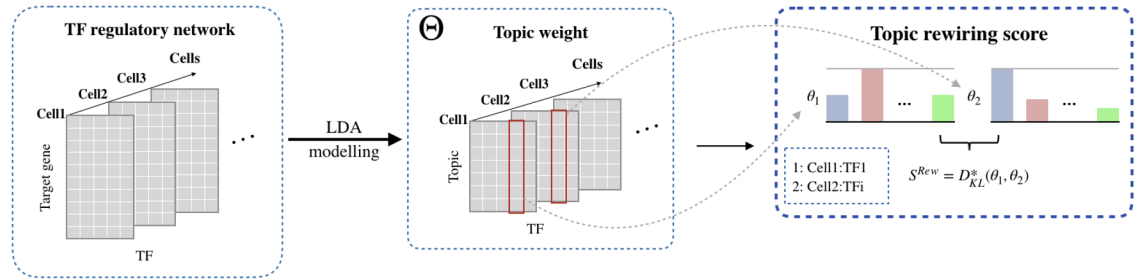
TopicNet: Measuring transcriptional regulatory network change using LDA



α Prior info β Prior info

θ : topic distribution per document

φ : word distribution per topic

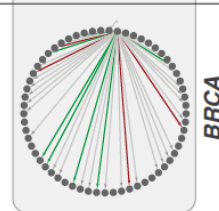
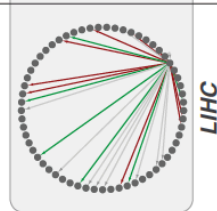
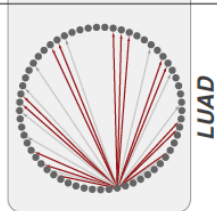
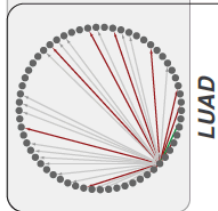
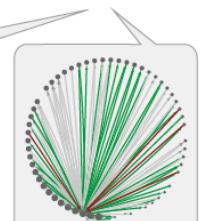
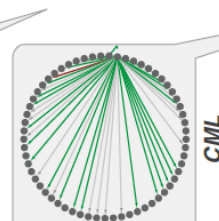
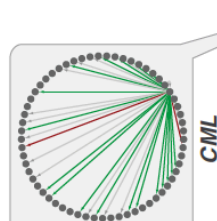
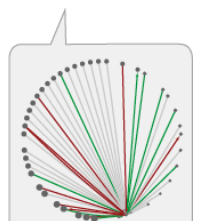
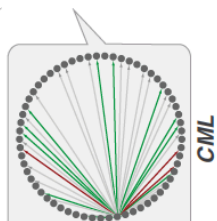
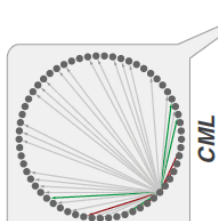
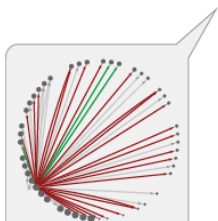
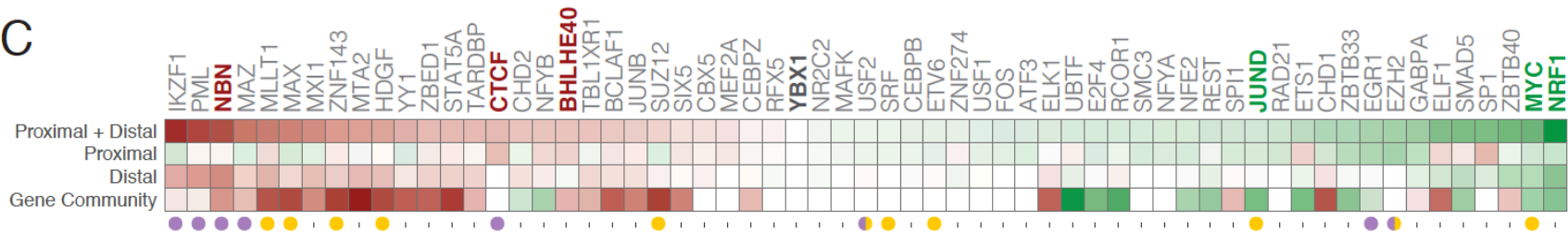


Loser

TF-Gene Network Rewiring

Gainer

C



- ➔ Gained Edge
- ➔ Retained Edge
- ➔ Lost Edge
- High Rewiring
- Low Rewiring
- TSG
- Oncogene/TSG
- Oncogene

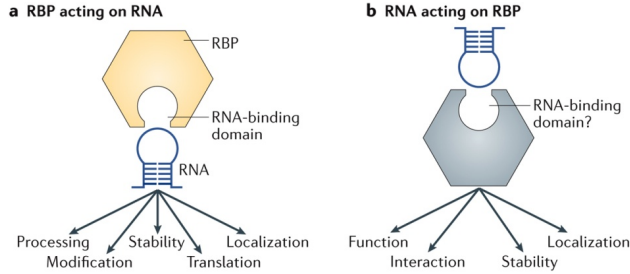
CML

Other Cancers

Using ENCODE Data for Cancer Genomics

- **BMR Correction:**
LARVA/MOAT/NIMBUS
 - Parametric models explicitly modeling genomic covariates
 - Many ENCODE covariates useful in accurately estimating background mutation rate
- **Network Rewiring in Cancer**
 - Large-scale ENCODE chip-seq data in certain cell lines highlights TFs changing targets greatly in oncogenesis. (Focus on CML)
 - TopicNet LDA approach (from text-mining) finds regulators that greatly change their gene communities
- **RADAR Variant Prioritization**
 - Prioritizes germline & somatic variants based on post-transcriptional regulome using ENCODE eCLIP
 - Incorporates new features related to RNA sec. struc & tissue specific effects
- **Regulatory Drivers of Differential Expression**
 - Highlighting regulators in terms of their power to drive differential expression.
 - Relationship of this to network hierarchy & RBP-TF cross talk
 - Example of MYC & SUB1

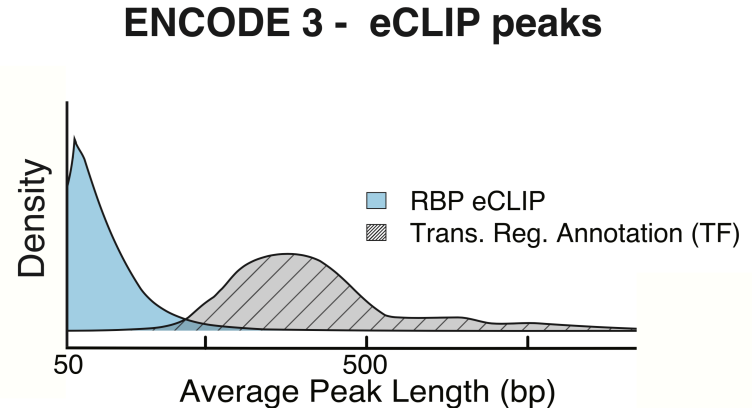
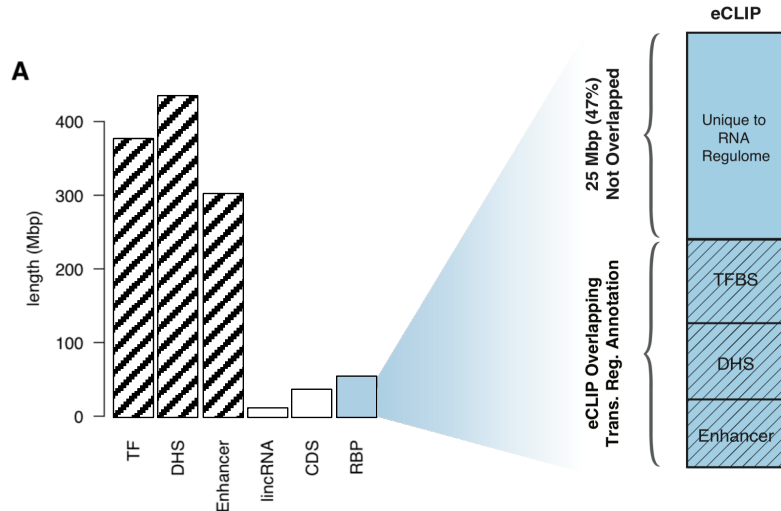
RNA Binding Proteins (RBPs)



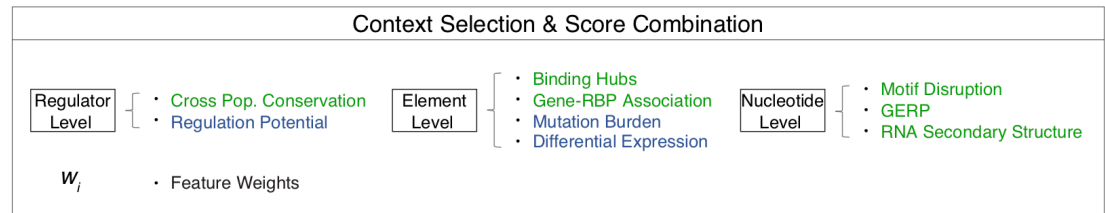
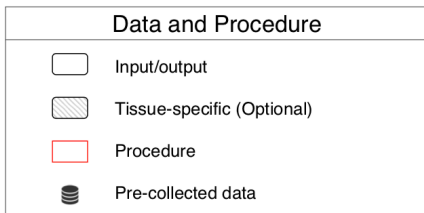
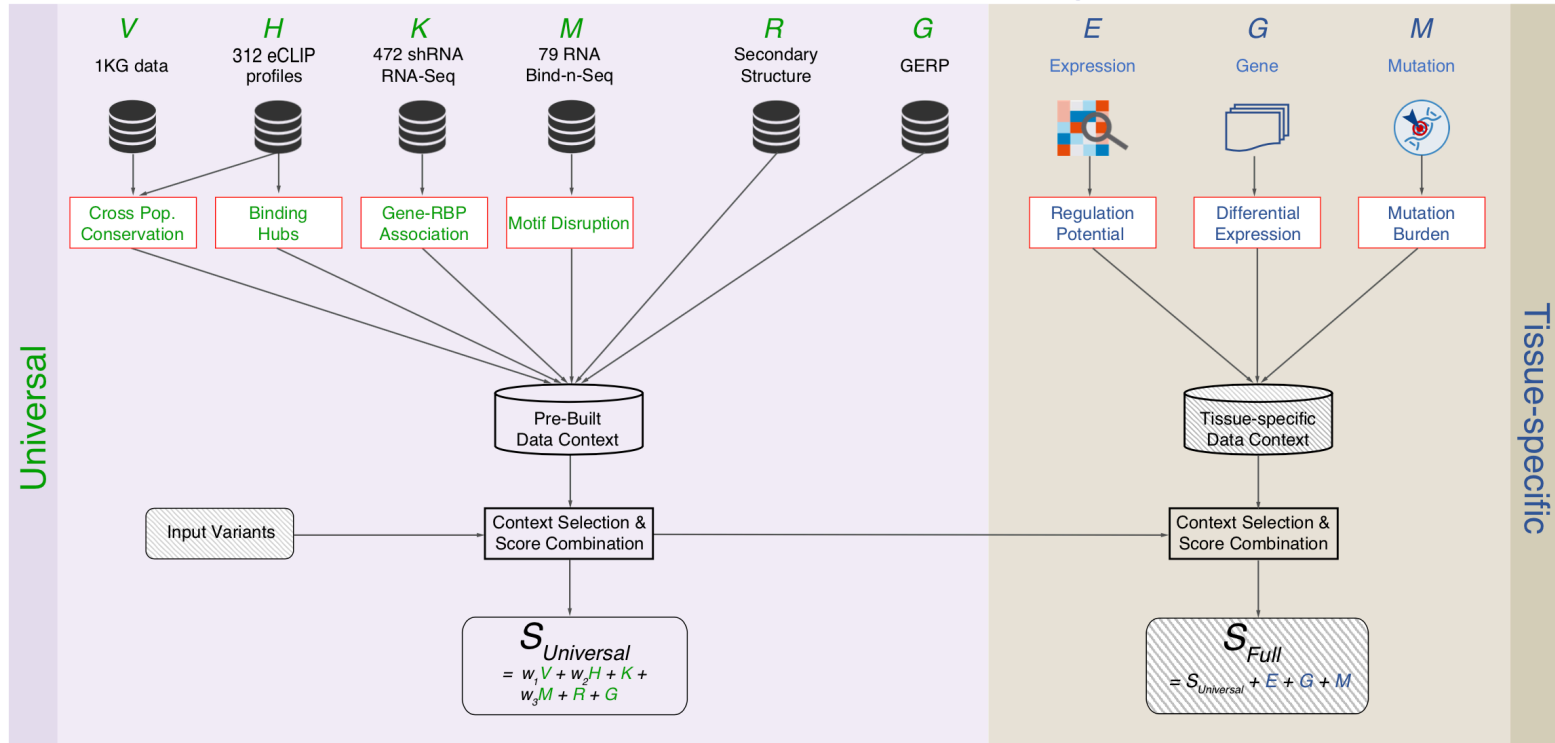
Nature Reviews | Molecular Cell Biology

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

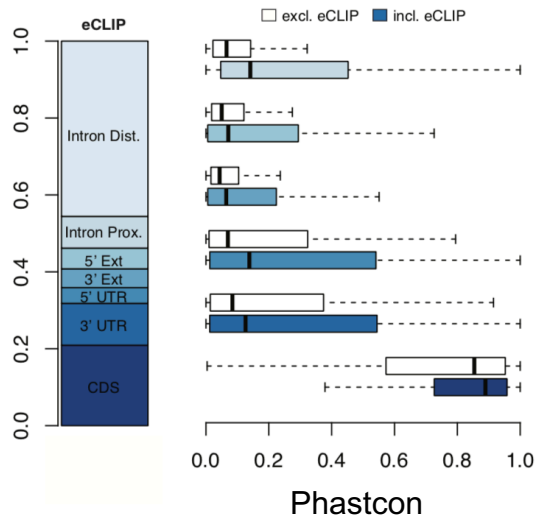
- **Before ENCODE3: >150 expt.** in many different cell types
- **ENCODE3 did ~350 focused eCLIP expt.** for >110 RBPs on HepG2 & K562 (Van Nostrand...Yeo. Nat. Meth. '16; Van Nostrand...Graveley, Yeo (submitted in relation to ENCODE3))



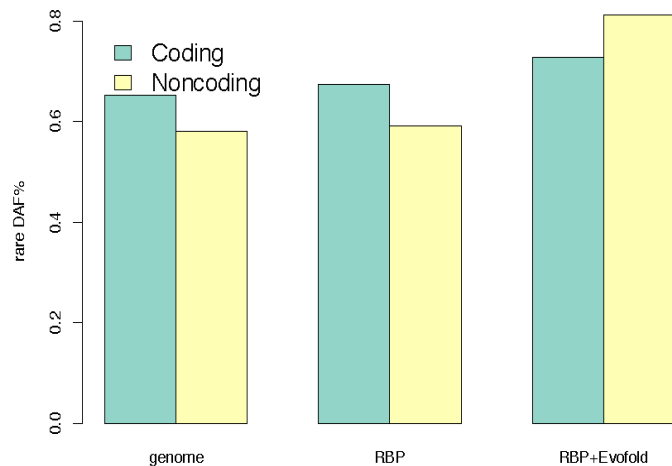
Schematic of RADAR Scoring



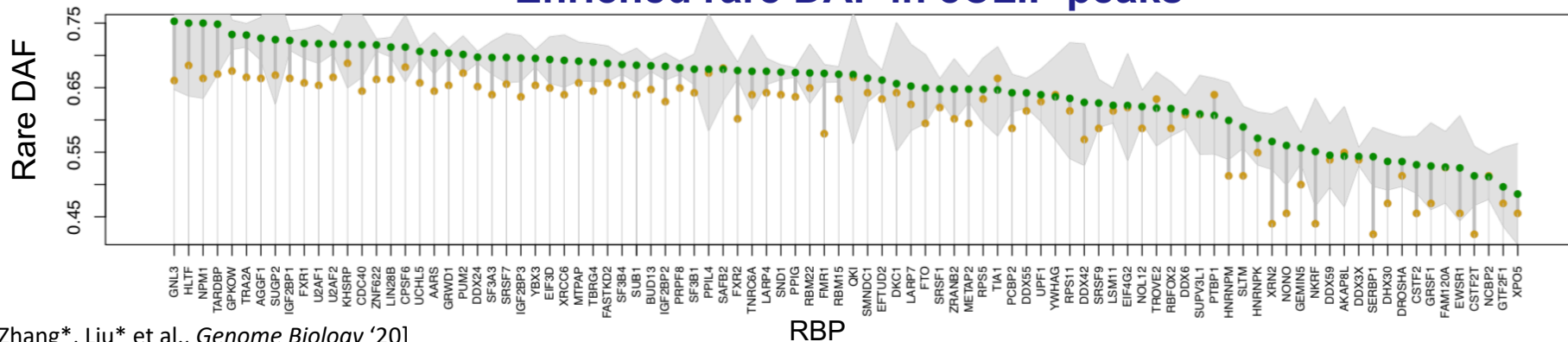
High Phastcon in RBP-overlapped annotations



RNA Structure Cons. from EvoFold

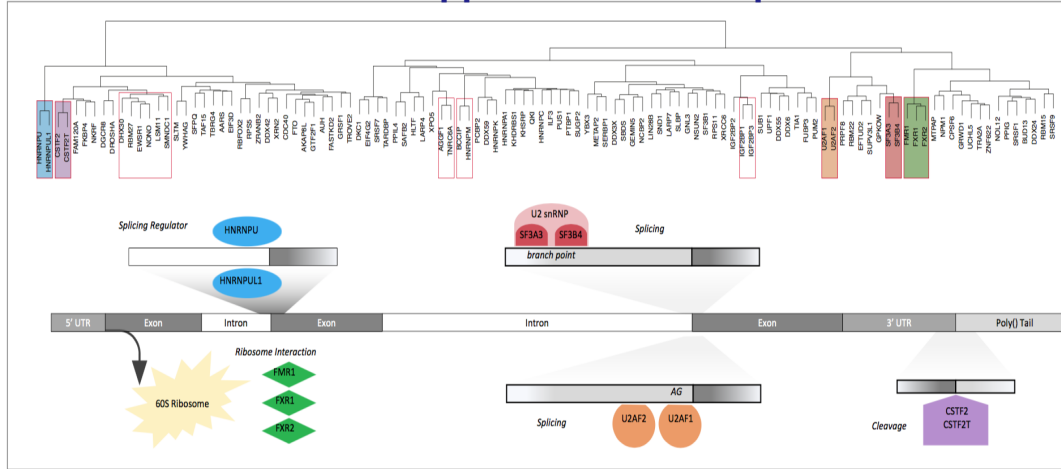


Enriched rare DAF in eCLIP peaks

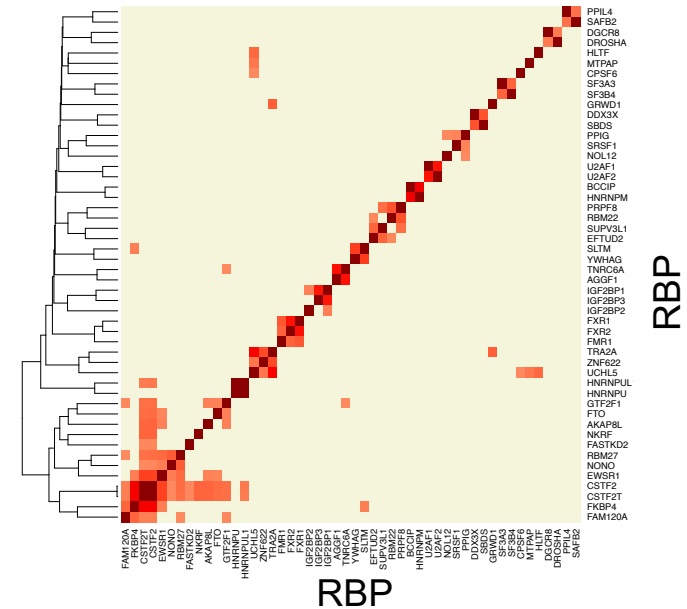


Co-binding of RBPs form biologically relevant complexes

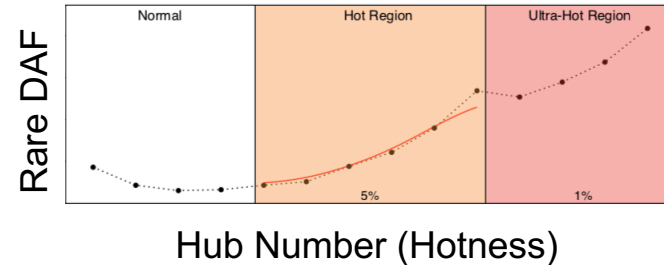
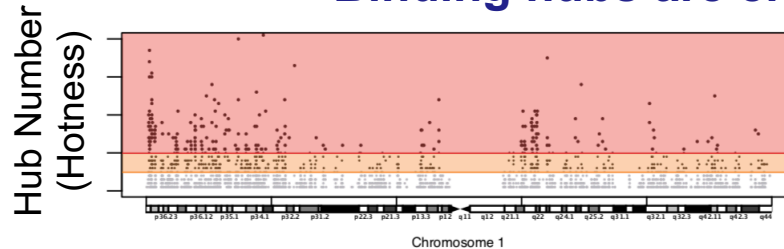
Literature supported RBP complexes



Unique co-binding patterns of RBPs

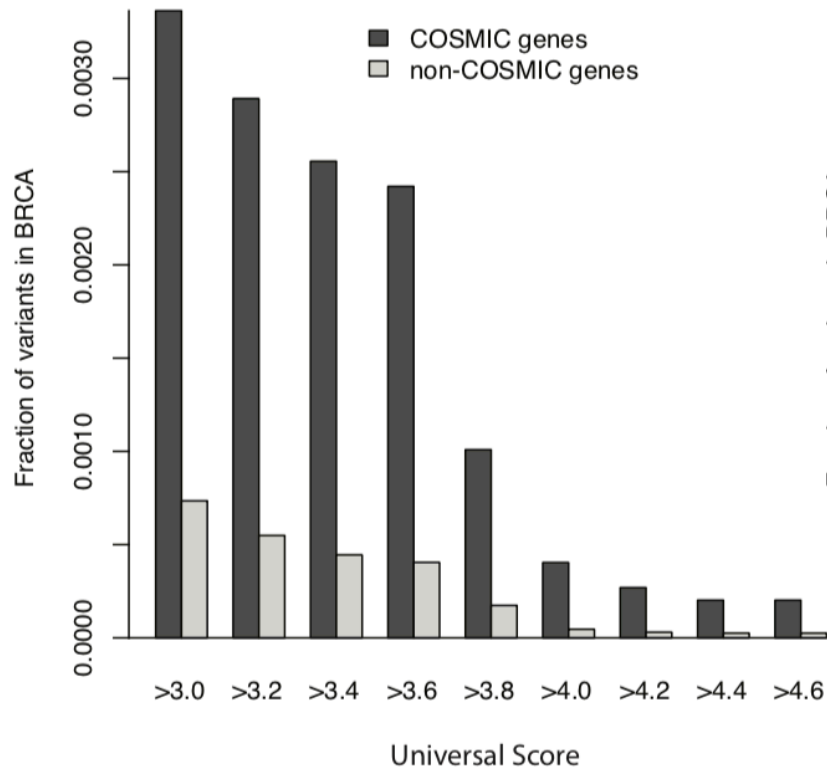


Binding hubs are enriched for rare variants

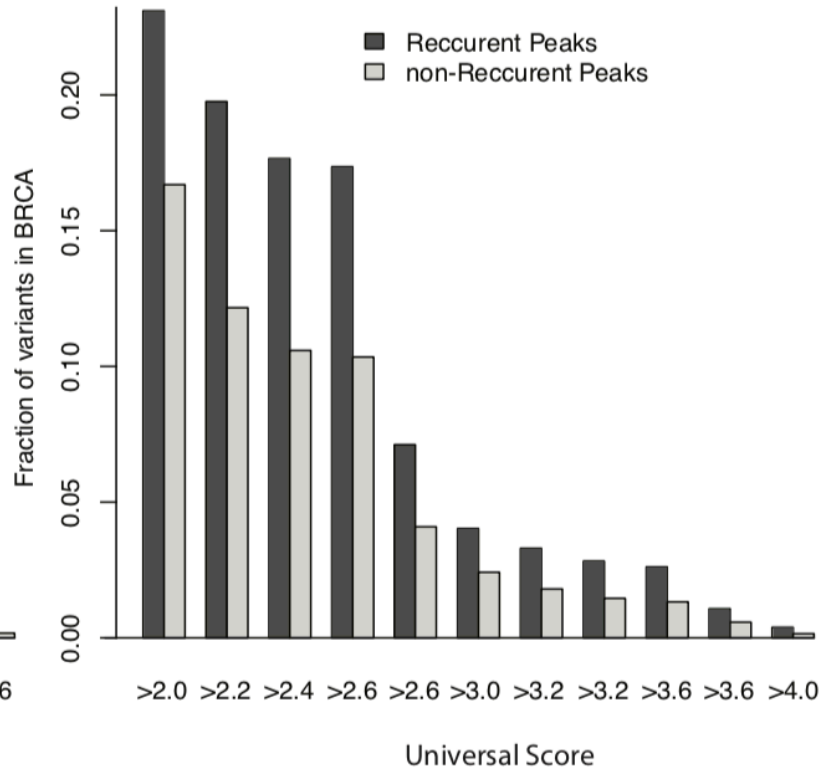


RADAR Scores enriched in COSMIC genes and recurrently mutated regions

A

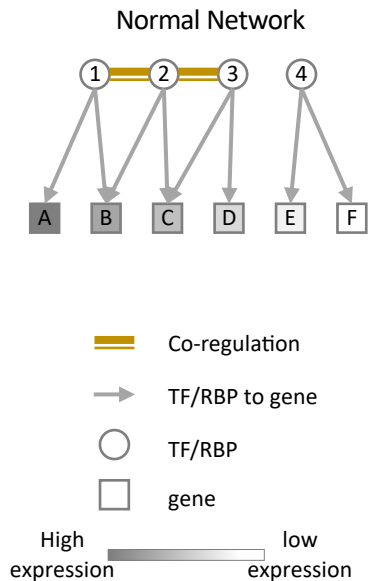


B

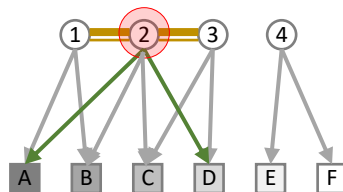


Using ENCODE Data for Cancer Genomics

- **BMR Correction:**
LARVA/MOAT/NIMBUS
 - Parametric models explicitly modeling genomic covariates
 - Many ENCODE covariates useful in accurately estimating background mutation rate
- **Network Rewiring in Cancer**
 - Large-scale ENCODE chip-seq data in certain cell lines highlights TFs changing targets greatly in oncogenesis. (Focus on CML)
 - TopicNet LDA approach (from text-mining) finds regulators that greatly change their gene communities
- **RADAR Variant Prioritization**
 - Prioritizes germline & somatic variants based on post-transcriptional regulome using ENCODE eCLIP
 - Incorporates new features related to RNA sec. struc & tissue specific effects
- **Regulatory Drivers of Differential Expression**
 - Highlighting regulators in terms of their power to drive differential expression.
 - Relationship of this to network hierarchy & RBP-TF cross talk
 - Example of MYC & SUB1

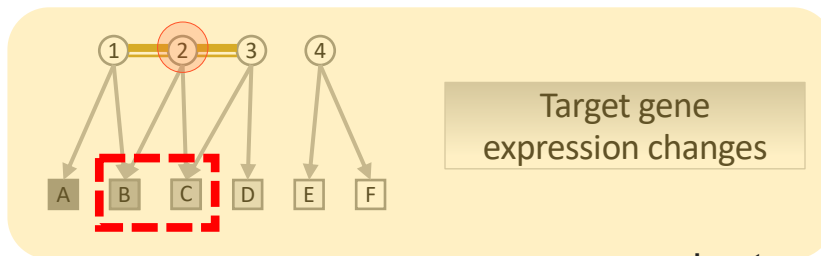


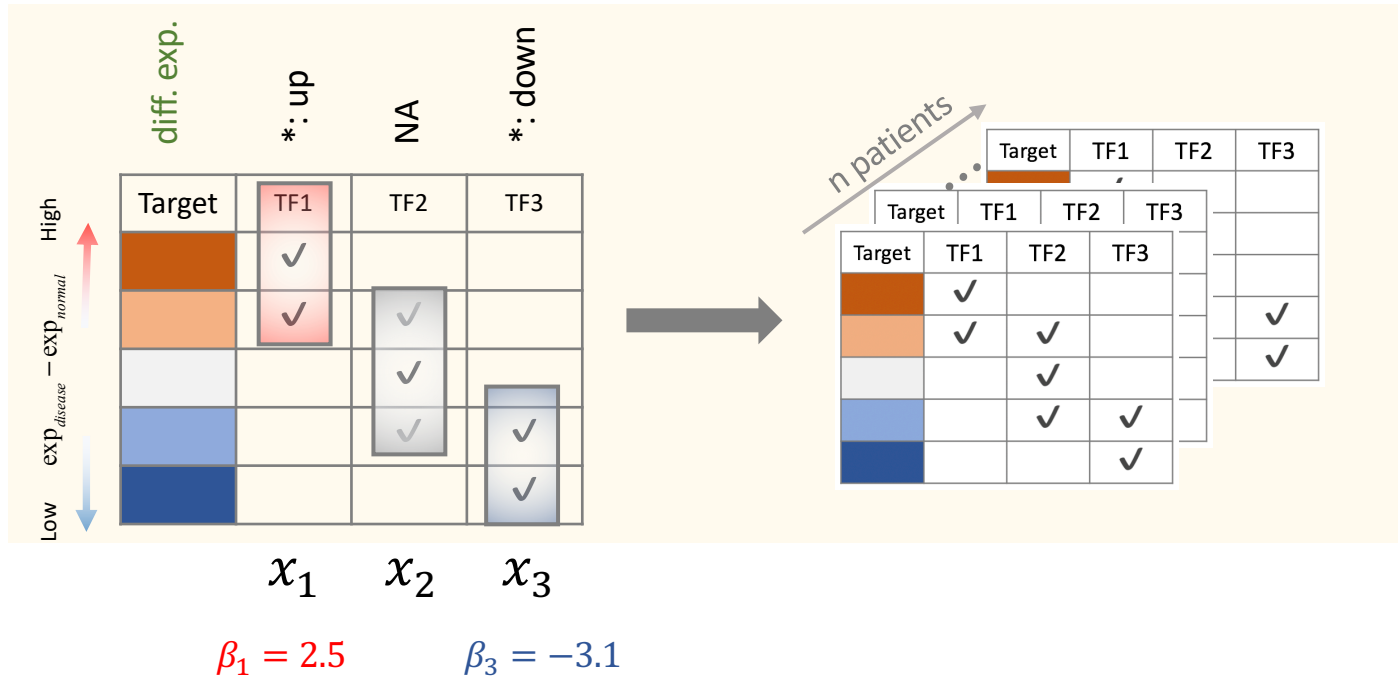
Disease Network :
dotted line = lost edge



Principles

Direct target
gain/loss

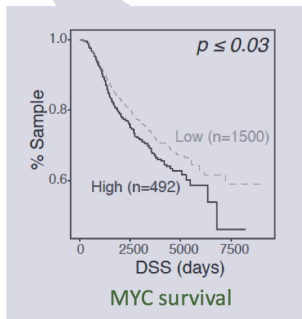
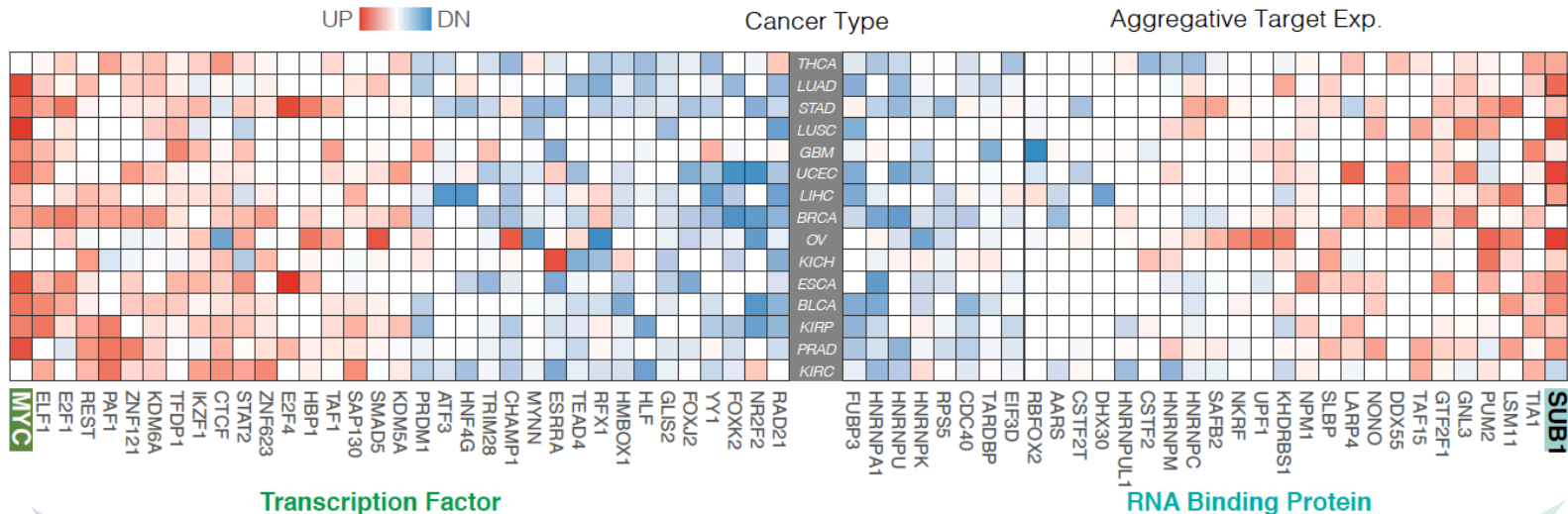




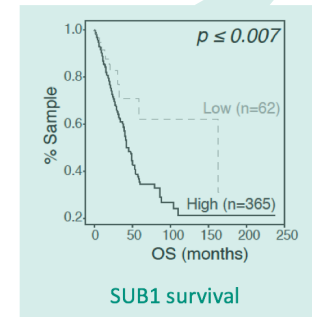
2198 ChIP-seq
459 eCLIP

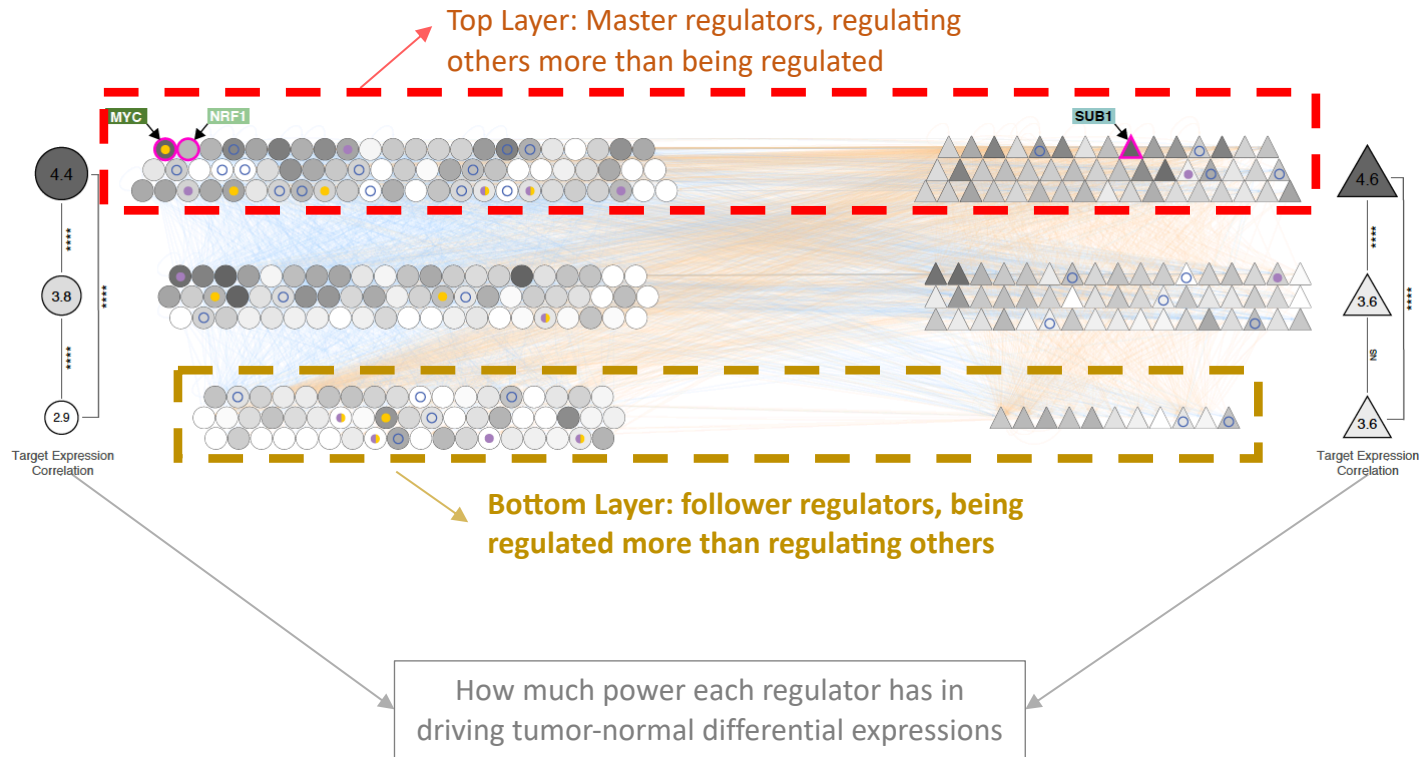
$$y = (\text{exp}_{\text{disease}} - \text{exp}_{\text{normal}}) \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

differential expression Network for Regulator 1 to k



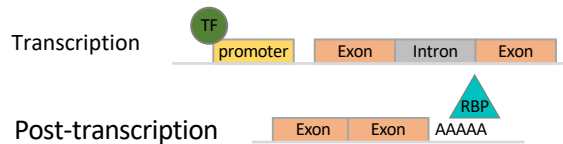
Aggregated t-statistic in regression over TCGA samples





TF-RBP crosstalk

TF-RBP regulate the same gene at different levels



Using ENCODE Data for Cancer Genomics

- **BMR Correction:**
LARVA/MOAT/NIMBUS
 - Parametric models explicitly modeling genomic covariates
 - Many ENCODE covariates useful in accurately estimating background mutation rate
- **Network Rewiring in Cancer**
 - Large-scale ENCODE chip-seq data in certain cell lines highlights TFs changing targets greatly in oncogenesis. (Focus on CML)
 - TopicNet LDA approach (from text-mining) finds regulators that greatly change their gene communities
- **RADAR Variant Prioritization**
 - Prioritizes germline & somatic variants based on post-transcriptional regulome using ENCODE eCLIP
 - Incorporates new features related to RNA sec. struc & tissue specific effects
- **Regulatory Drivers of Differential Expression**
 - Highlighting regulators in terms of their power to drive differential expression.
 - Relationship of this to network hierarchy & RBP-TF cross talk
 - Example of MYC & SUB1

Using ENCODE Data for Cancer Genomics

- **BMR Correction:**
LARVA/MOAT/NIMBUS
 - Parametric models explicitly modeling genomic covariates
 - Many ENCODE covariates useful in accurately estimating background mutation rate
- **Network Rewiring in Cancer**
 - Large-scale ENCODE chip-seq data in certain cell lines highlights TFs changing targets greatly in oncogenesis. (Focus on CML)
 - TopicNet LDA approach (from text-mining) finds regulators that greatly change their gene communities
- **RADAR Variant Prioritization**
 - Prioritizes germline & somatic variants based on post-transcriptional regulome using ENCODE eCLIP
 - Incorporates new features related to RNA sec. struc & tissue specific effects
- **Regulatory Drivers of Differential Expression**
 - Highlighting regulators in terms of their power to drive differential expression.
 - Relationship of this to network hierarchy & RBP-TF cross talk
 - Example of MYC & SUB1



ENCODEC.gersteinlab.org

J **Zhang**, D **Lee**, V **Dhiman**, P **Jiang**, J **Xu**,
P McGillivray, H Yang.... S Liu, K White

NIMBus.gersteinlab.org

J **Zhang**, J **Liu**, P McGillivray, C Yi, L Lochovsky, D Lee

RADAR.gersteinlab.org

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

{LARVA,MOAT}.gersteinlab.org

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

github.com/gersteinlab/**TopicNet**

S **Lou**, T **Li**, X **Kong**, J Zhang, J Liu, D Lee



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)