

# Developing Informatics Technologies to Model Cancer Gene Regulation

Xiaole Shirley Liu

*Dept of Biostatistics & Computational Biology  
Center for Functional Cancer Epigenetics  
Dana-Farber Cancer Institute  
Harvard School of Public Health  
Broad Institute of Harvard and MIT*

# Public Data

- ChIP-seq and related data in human / mouse:
  - Factor ChIP-seq: > 11K, > 700 different TFs
  - Histone Mark ChIP-seq: > 10K
  - Chromatin accessibility: > 1300
- TCGA: gene expression profiles in ~ 12K tumors and 31 cancer types
- Model transcriptional and epigenetic gene regulation

# Cistrome: cis-elements bound by trans-factor

## Cistrome Analysis Pipeline: <http://cistrome.org/ap/>

The screenshot shows the Cistrome Analysis Pipeline interface. On the left, a sidebar lists various tools: Import Data, Peak Calling (with options for Upload File from your computer, MA2C Peak Calling for ChIP-chip (Nimblegen), MACS Peak Calling for ChIP-Seq, and Call Peaks from Wiggle file Peak Calling from Wiggle file), Peak Calling - MAT, Correlation, Annotate and Visualize, Expression, Motif, Liftover/Others, Galaxy:Get Data, Galaxy:Send Data, Galaxy:Lift-Over, Galaxy:Text Manipulation, Galaxy:Filter and Sort, Galaxy:Join, Subtract and Group, Galaxy:Convert Formats, Galaxy:Extract Features, Galaxy:Fetch Sequences, Galaxy:Get Genomic Scores, Galaxy:Operate on Genomic Intervals. The main area is titled "MACS" and contains the following form fields:

- Treatment file: 150: PU1-all-clean.bed
- Input file: 154: input-clean.bed
- Format: Bed
- Effective Genome Size: 2170000000
- Tag size: 25
- P-Value: 1e-8
- Use Model?: True
- Model fold: 20
- Advanced Options: No
- Execute button

A note at the bottom states: "This tool performs peak calling for ChIP-Seq data. MACS is developed in Xiaole Shirley Liu's lab, by Tao Liu and Yong Zhang, and published on Genome Biology (pubmed: 18798982). The version deployed here is 1.3.7.1." A tip message says: "TIP: Please first upload your treatment and control files using the Upload File from your computer tool."

The right side features a "History" panel listing previous jobs:

- refresh | collapse all
- Unnamed history
- 188: Venn Diagram on data 181 and data 185
- 187: MACS xls on CEBPB-treat-clean.bed
- 186: MACS job log on CEBPB-treat-clean.bed
- 185: MACS peaks on CEBPB-treat-clean.bed
- 184: MACS wiggle on CEBPB-treat-clean.bed
- 183: MACS xls on PU1-all-clean.bed
- 182: MACS job log on PU1-all-clean.bed
- 181: MACS peaks on PU1-all-clean.bed
- 167: MACS xls on PU1-all-clean.bed
- 166: MACS job log on PU1-all-clean.bed
- 165: MACS peaks on PU1-all-clean.bed
- 164: MACS wiggle on PU1-all-clean.bed
- 155: H3K4me1-clean.bed
- 154: input-clean.bed

# Cistrome DB

## <http://cistrome.org/db/>

- All the ChIP/DNase-seq data human (hg38) / mouse (mm10) data as of Dec 2015 processed (peaks, signals, QC)

Cistrome Data Browser    Home    About    LOGOUT wuqiu ➔

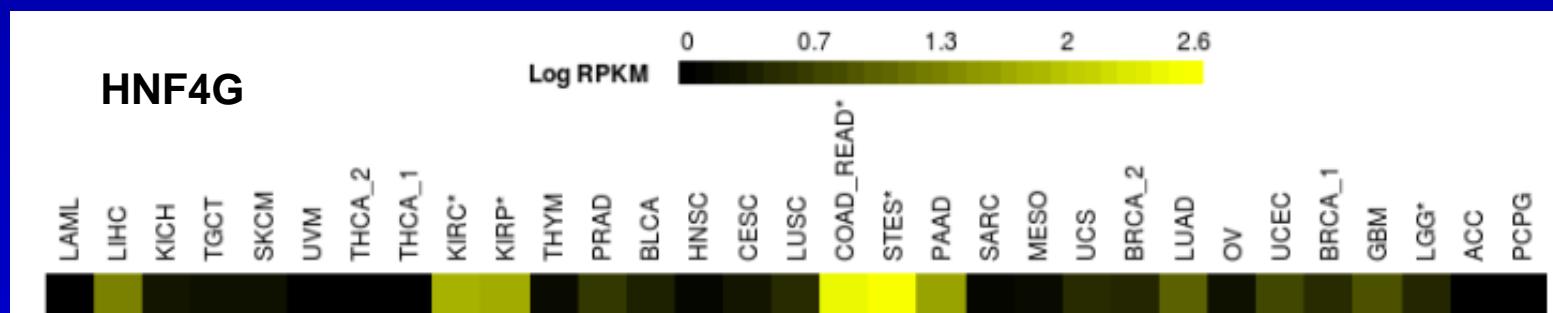
## Dataset Browser

Containing word(s):   Options ▾

Species	Biological Sources	Factors
All Homo sapiens Mus musculus	All 7438 BT-474 DBDmut DL23 DLD-1	All EP300 ESR1 FOXA1 GATA3 H3K27ac

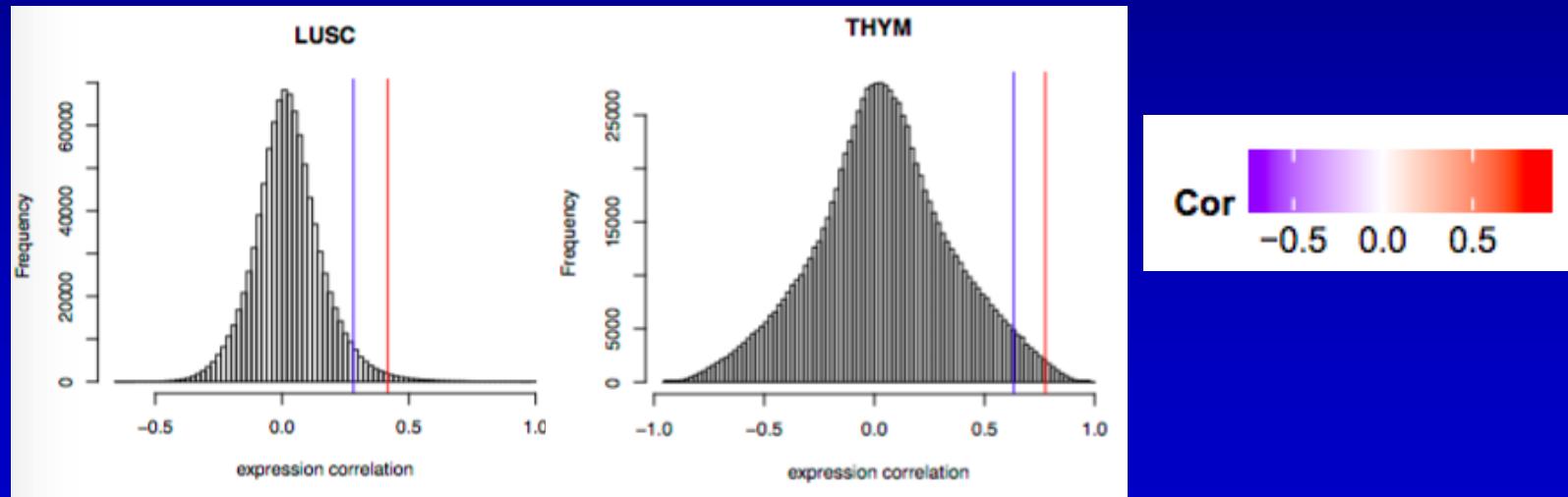
# Factors Influencing Cancer-Specific Transcription Factor Targets

- Re-cluster all TCGA gene expression
  - Separate strong subtypes (luminal vs basal BRCA)
  - Combine similar cancer types (READ + COAD)
- TF expression level in each cancer type
  - Eliminate TF if median TF expression RPKM < 1



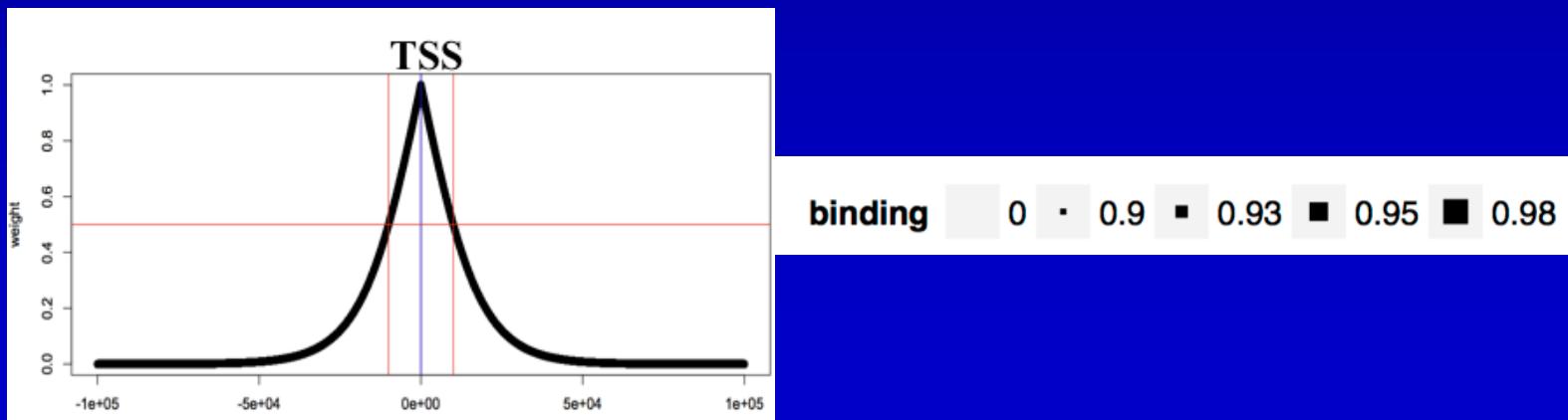
# Factors Influencing Cancer-Specific Transcription Factor Targets

- Gene expression correlation between TF and putative targets
- Cutoff (5%) based on cancer-specific null distribution

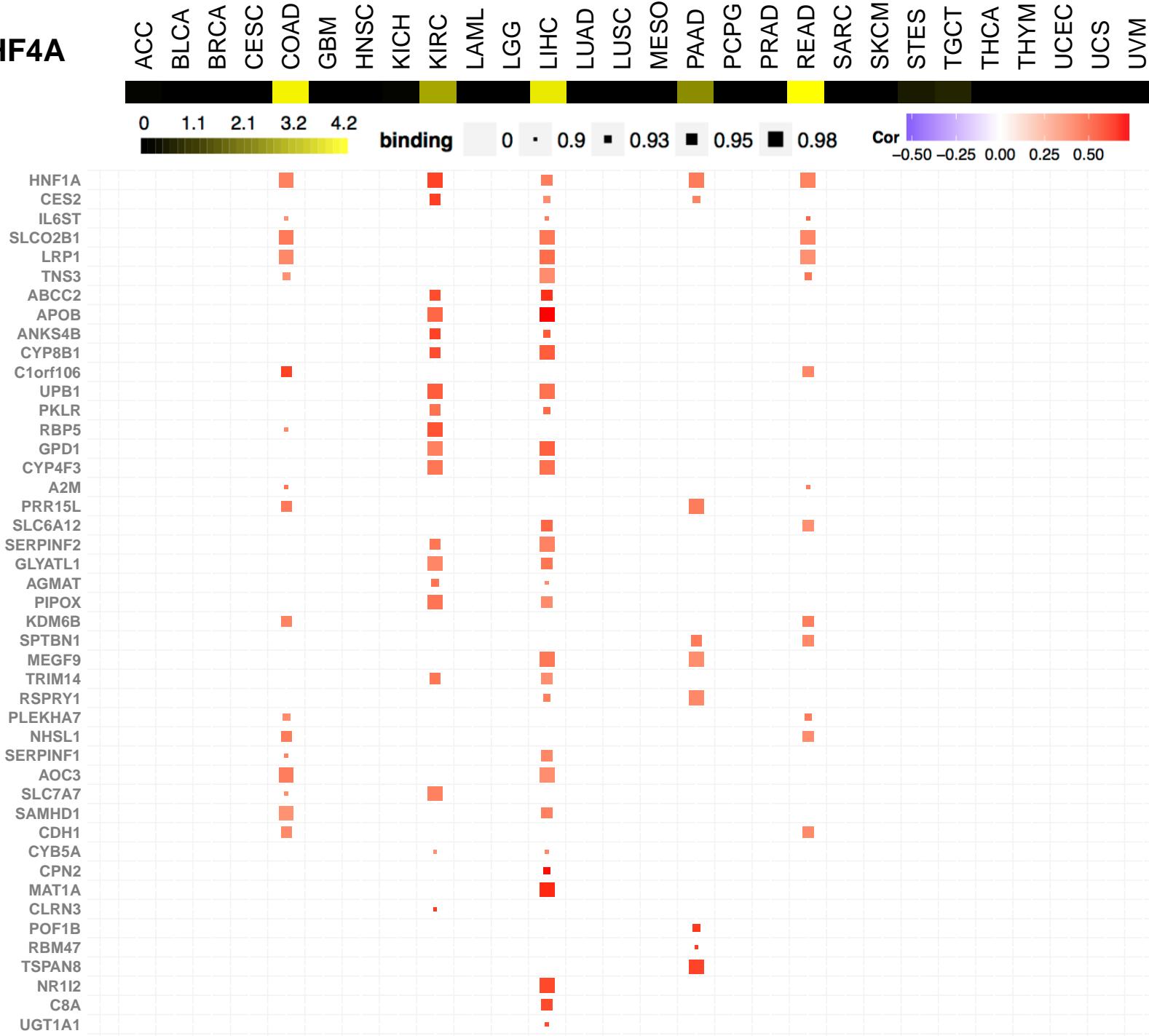


# Factors Influencing Cancer-Specific Transcription Factor Targets

- For each TF ChIP-seq, calculate **regulatory potential**, the sum of binding sites weighted by distance to TSS with exponential decay
- For each cancer, use logistic regression to select TF ChIP-seq data that best explain the top correlated genes with TF

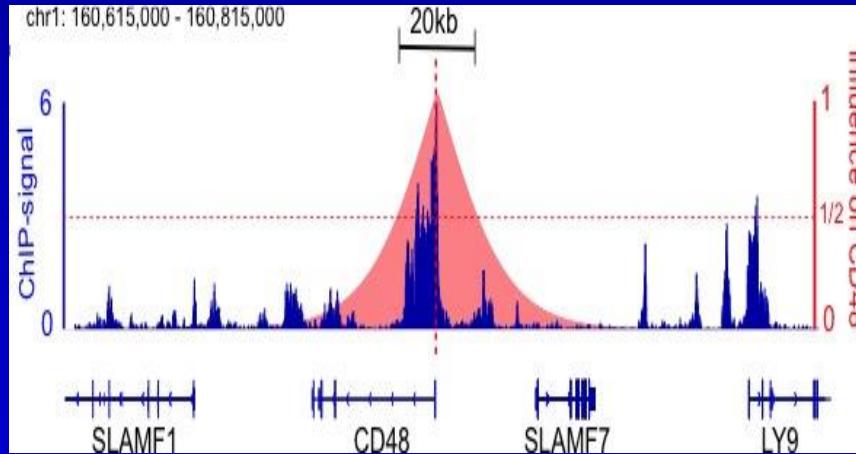


# NHF4A



# Reuse Public H3K27ac

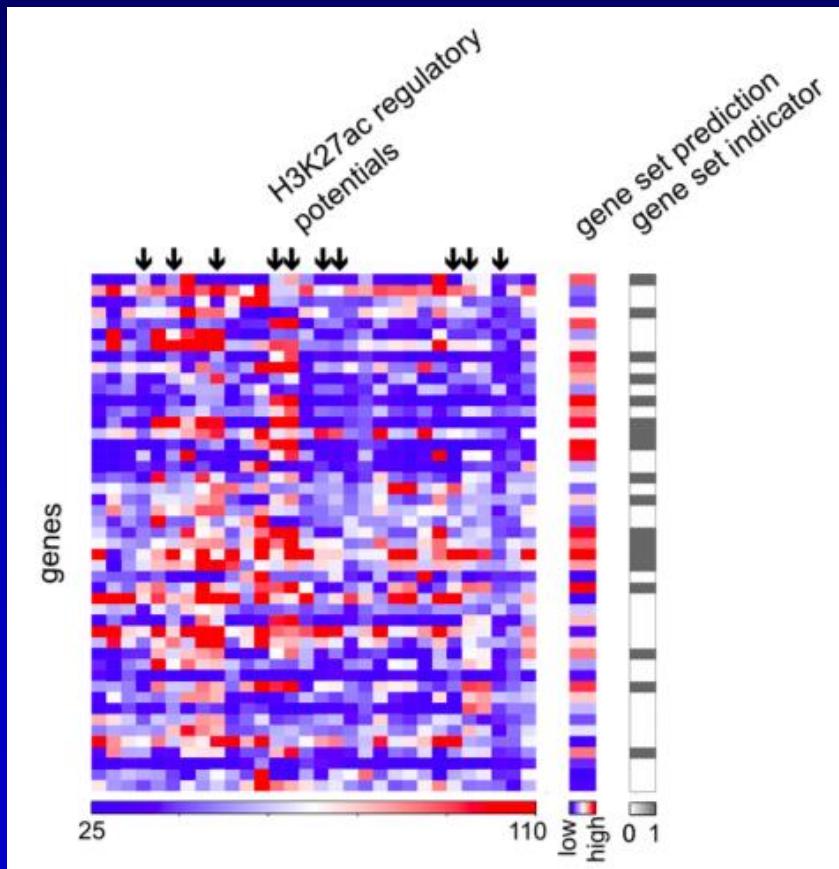
- Many laboratories have differential gene expression without matched H3K27ac ChIP-seq
- Regulatory Potential from H3K27ac
  - No peak calls, distance weight all H3K27ac reads by their distance to TSS with 10kb exponential decay



# MARGE

## Model-based Analysis of Regulation of Gene Expression

- Find relevant public H3K27ac ChIP-seq data that best represent the gene list



Logistic regression

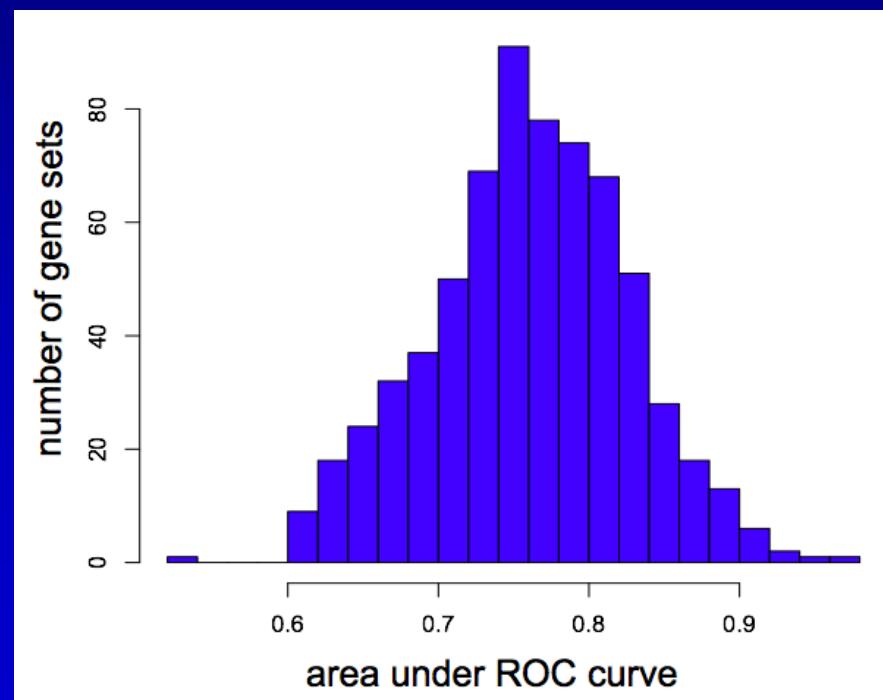
$$y_i \sim \alpha_0 + \sum_j a_j \sqrt{RP_{ij}}$$



# The Power of Public H3K27ac Data

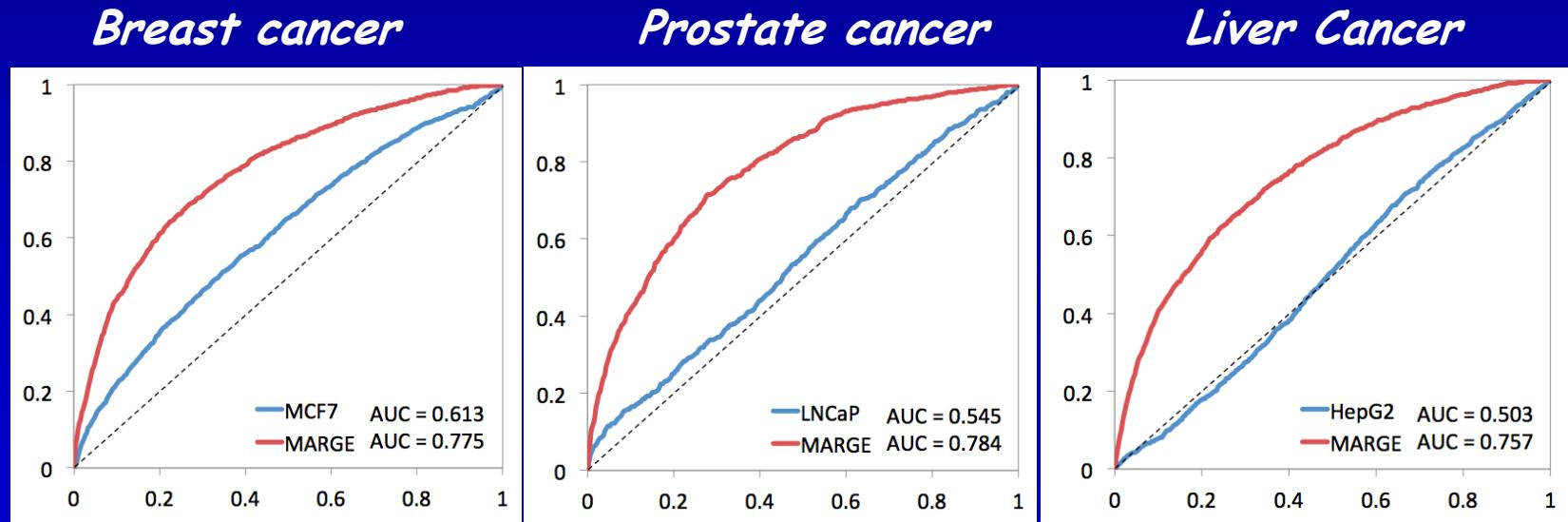
- Given a gene list (e.g. up-regulated upon perturbation), MARGE uses regression to retrieve informative H3K27ac ChIP-seq profiles in the public to explain the genes

Train on genes in odd chromosomes and test on genes in even chromosomes



# H3K27ac ChIP-seq for TCGA

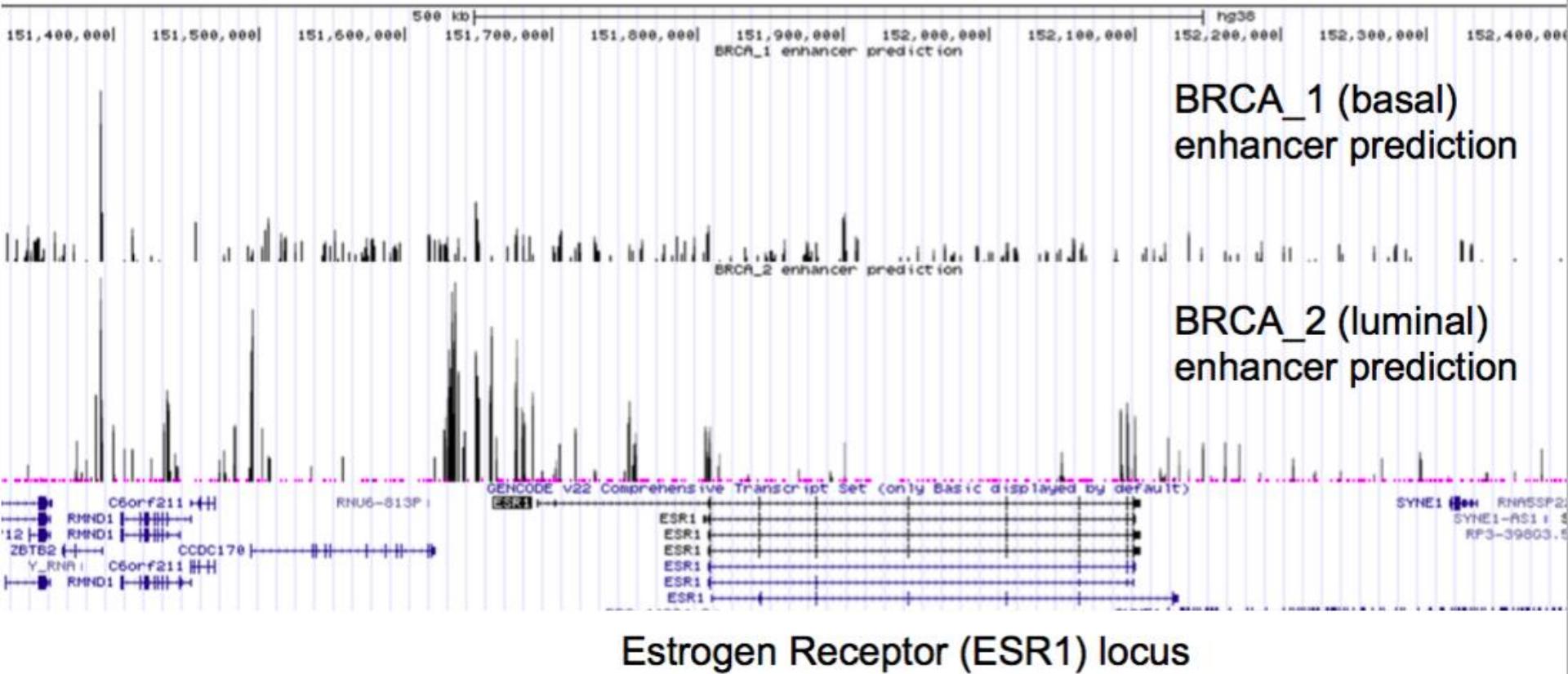
- No histone mark ChIP-seq data for TCGA
- Given genes up-regulated in one cancer type, could we find relevant public H3K27ac ChIP-seq collections?
- Use DNase-seq peaks to improve H3K27ac resolution



# TCGA Cancer Enhancer Prediction (Beta Version)

No	ID	Description	# Cancer samples	# Normal samples	# Cancer genes	Enhancer Prediction	
1	BLCA	Bladder urothelial carcinoma	242	11	<a href="#">2724</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
2	BRCA_1	Breast invasive carcinoma 1 (basal)	199	17	<a href="#">1860</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
3	BRCA_2	Breast invasive carcinoma 2 (luminal)	844	94	<a href="#">1760</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
4	COAD_READ*	Colon & colorectal adenocarcinoma	349	30	<a href="#">2031</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
5	HNSC	Head and Neck squamous cell carcinoma	335	33	<a href="#">2165</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
6	KICH	Kidney chromophobe	64	24	<a href="#">1504</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
7	KIRC*	Kidney renal clear cell carcinoma	501	69	<a href="#">2727</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
8	KIRP*	Kidney renal papillary cell carcinoma	287	32	<a href="#">1691</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
9	LIHC	Liver hepatocellular carcinoma	354	47	<a href="#">4220</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
10	LUAD	Lung adenocarcinoma	467	51	<a href="#">2161</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
11	LUSC	Lung squamous cell carcinoma	356	34	<a href="#">3744</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
12	PRAD	Prostate adenocarcinoma	490	52	<a href="#">967</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
13	STES*	Stomach & Esophageal carcinoma	362	33	<a href="#">4162</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>
14	THCA_1	Thyroid carcinoma 1	357	42	<a href="#">1723</a>	<a href="#">download</a>	<a href="#">Visualize WashU Browser</a> <a href="#">Visualize UCSC Browser</a>

# Cancer-Specific Enhancer Profile



# Summary

- Use regulatory potential (exponential decay with distance) to model binding site and H3K27ac effect on gene expression
- MARGE uses logistic regression to retrieve public H3K27ac ChIP-seq to explain differential expression gene set of interest
- Can predict enhancer profiles for up-regulated genes in each TCGA cancer type
- <http://cistrome.org/CistromeCancer/>

# Acknowledgement

## Shirley Liu Laboratory

- Shenglin Mei
- Chongzhi Zang
- Cliff Meyer
- Su Wang
- Qian Qin
- Qiu Wu
- Jingyu Fan
- Hanfei Sun
- Rongbin Zheng
- Jian Ma
- Len Taing
- Muyuan Zhu
- Jiaxin Wu
- Anya Zhang

## Collaborators

- Myles Brown
- Jun Liu
- Yong Zhang
- Tao Liu
- Fugen Li
- Henry Long
- ITCR and colleagues



YouTube Video

[https://www.youtube.com/watch?  
v=pQ02wg8MI6o](https://www.youtube.com/watch?v=pQ02wg8MI6o)