# Interpretable Deep Learning for Chromatin-Informed Inference of Transcriptional Programs Driven by Somatic Alterations Across Cancers

Yifeng Tao[1,†], Xiaojun Ma[2,†], Drake Palmer[3], Russell Schwartz[1,4], Xinghua Lu[2,5], Hatice Ulku Osmanbeyoglu[2,6,*]

[1]Computational Biology Department, School of Computer Science, Carnegie Mellon University

[2]Department of Biomedical Informatics, School of Medicine, University of Pittsburgh

[3]Department of Biological Sciences, University of Pittsburgh School of Arts & Sciences

[4]Department of Biological Sciences, Carnegie Mellon University

[5]Department of Pharmaceutical Science, School of Medicine, University of Pittsburgh

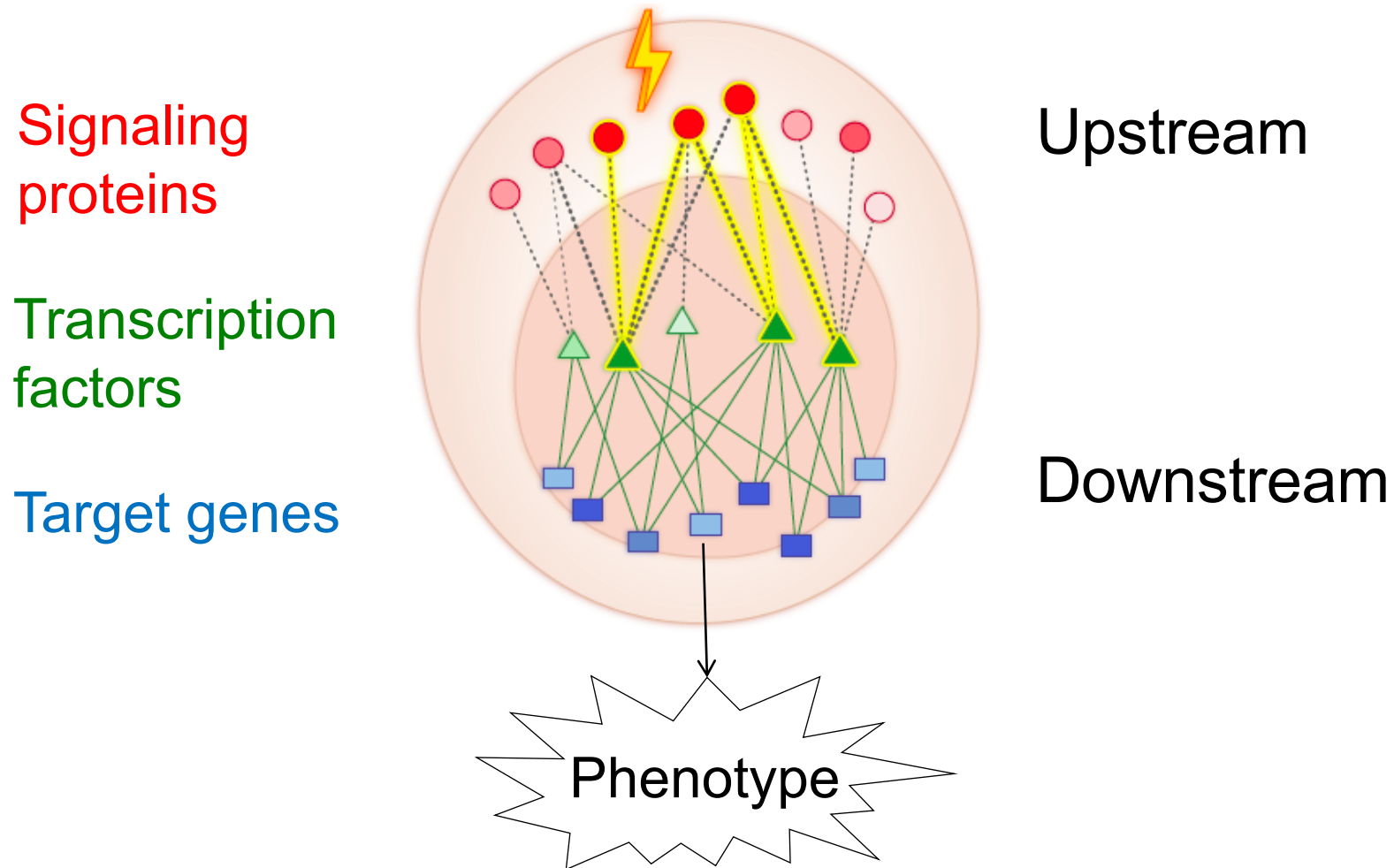[6]Department of Bioengineering, School of Engineering, University of Pittsburgh
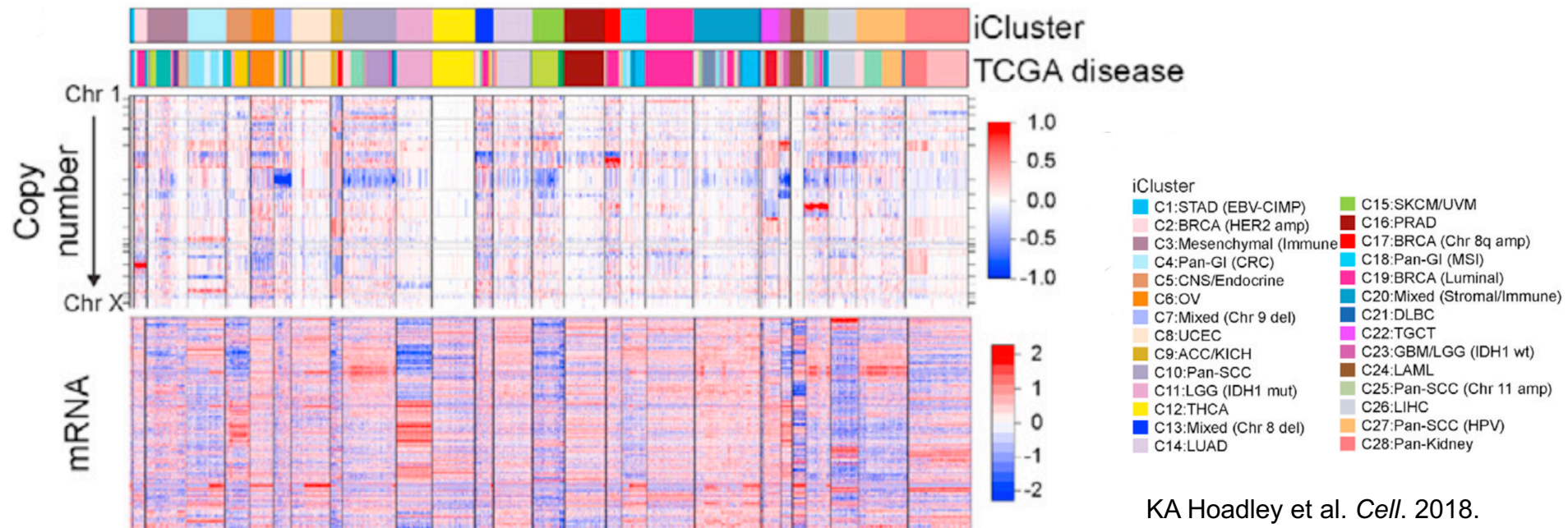
[†]Contributed equally: Y.T., X.M.

# Signaling and transcriptional response

- Cancers are caused by the perturbations of multiple pathways and transcriptional regulatory programs
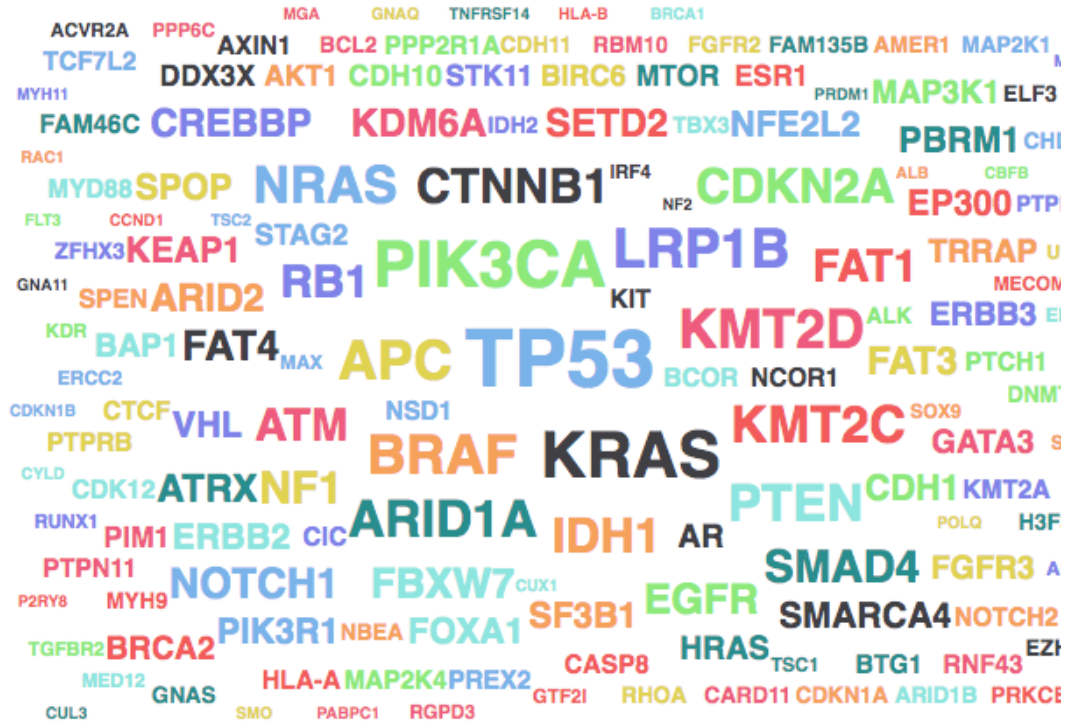


Signaling proteins

Transcription factors

Target genes

Upstream

Downstream

Phenotype

# Pan-cancer modeling of regulatory programs



KA Hoadley et al. *Cell*. 2018.

- Similar TFs may be dysregulated across cancers

- Similarities between cancer types can inform new therapies

- Extensive training data from more common tumor types also compensates for smaller sample sizes in similar but rarer cancers (e.g. pheochromocytoma and paraganglioma; PCPG)

# Modeling non-linear relationships



https://www.intogen.org

- Effects of upstream alterations not equal, e.g., cancer drivers vs. passengers
- Complex interactions between genes, e.g., mutual exclusivity
- Role of genomic alterations is context specific

- Attention mechanism!

# Attention mechanism

- A deep learning method to assign importance weights to input features
  - Widely used in Computer Vision/Natural Language Processing
  - Computed in a contextual manner

Toilet tissue    Eskimo dog    Snow leopard

S Woo et al. *ECCV*. 2018.

The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
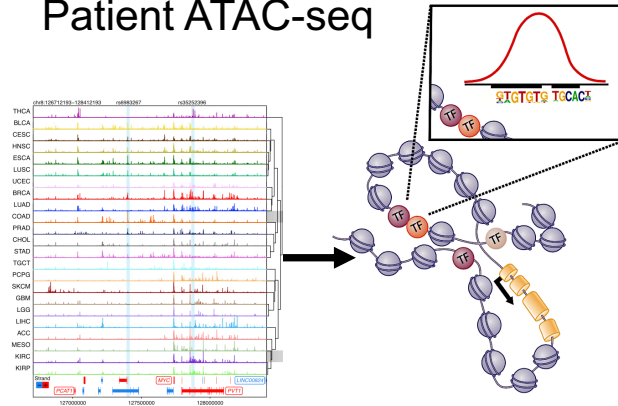The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .

J Cheng et al. *EMNLP*. 2016.

5

# Datasets/Approach: Modeling impact of somatic alterations on gene expression programs



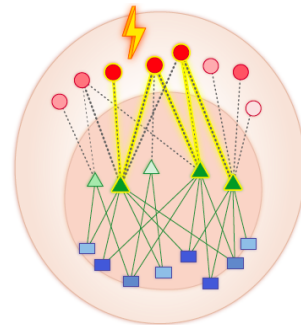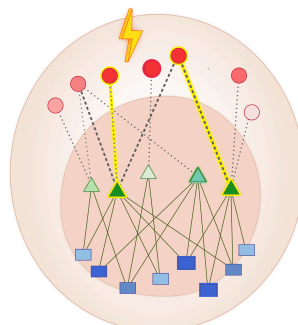Patient somatic alterations

Patient ATAC-seq

Patient RNA-seq

Mutation

Mutation

Mutation

Cancer type 1

Cancer type 2

...

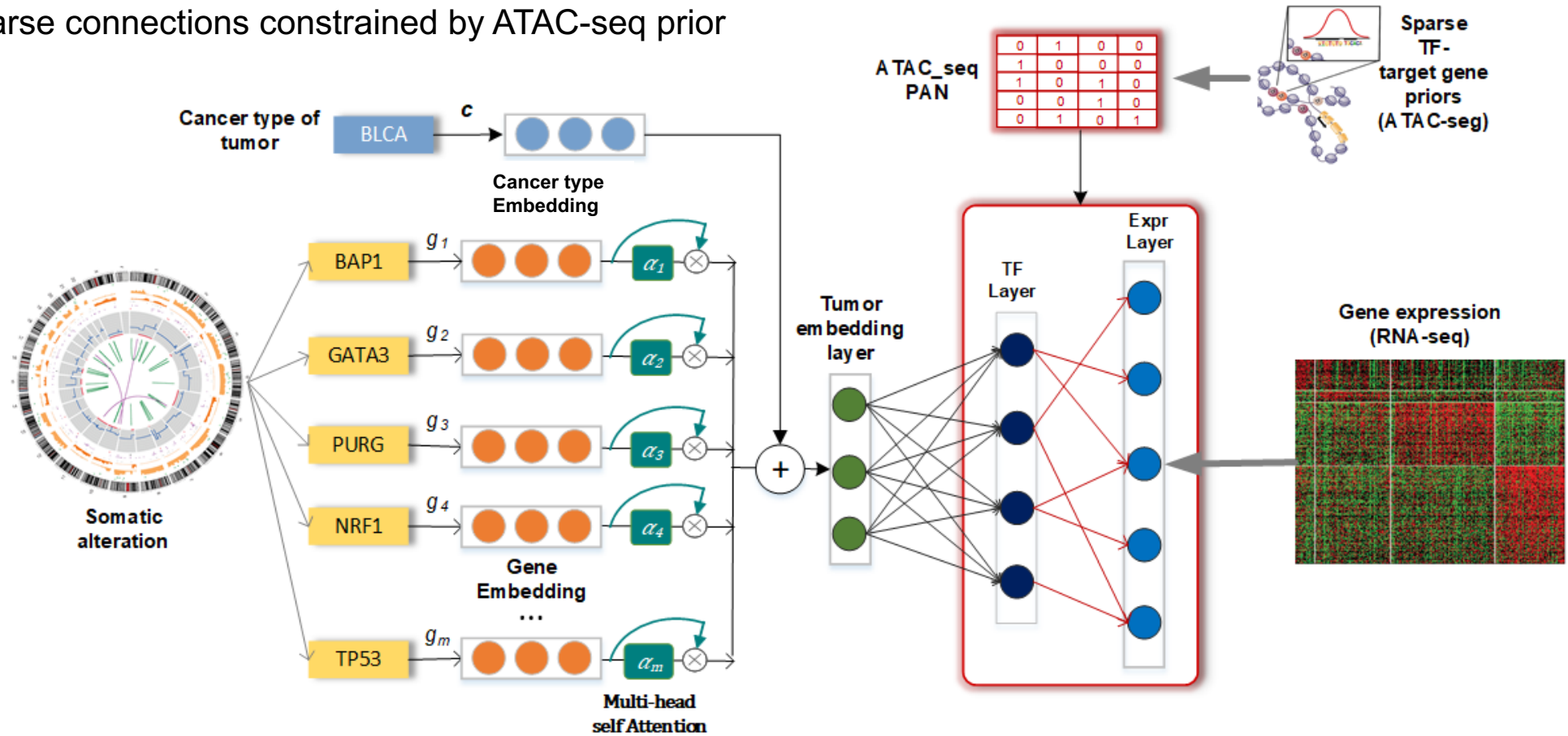Cancer type 17

Patients-specific regulatory networks
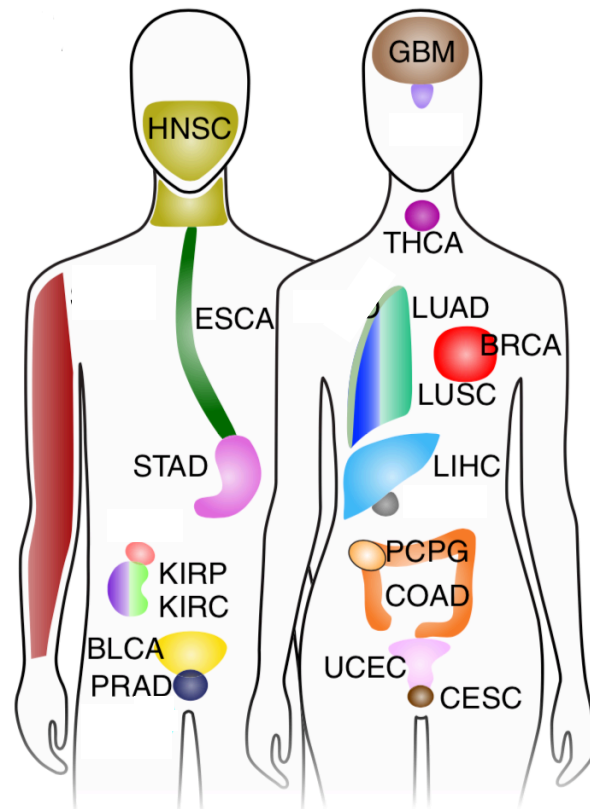
- Transcription factor
- Target genes

# Approach: interpretable deep learning

- CITRUS
  - Chromatin-informed Inference of Transcriptional Regulators Using Self-attention mechanism
  - Self-attention mechanism
  - Sparse connections constrained by ATAC-seq prior

# Pan-cancer data sources



MR Corces et al. *Science*. 2018.

| Datasets | Summary |
|---|---|
| ATAC-seq | 410 tumors |
| Bladder (BLCA) | 371 tumors |
| Breast (BRCA) | 719 tumors |
| Cervical and endocervical (CESC) | 267 tumors |
| Colon (COAD) | 271 tumors |
| Esophageal (ESCA) | 170 tumors |
| Glioblastoma (GBM) | 143 tumors |
| Head and Neck (HNSC) | 475 tumors |
| Kidney renal clear cell (KIRC) | 357 tumors |
| Kidney renal papillary cell (KIRP) | 272 tumors |
| Liver hepatocellular (LIHC) | 336 tumors |
| Lung adenocarcinoma (LUAD) | 459 tumors |
| Lung squamous (LUSC) | 430 tumors |
| Pheochromocytoma and Paraganglioma (PCPG) | 109 tumors |
| Prostate (PRAD) | 449 tumors |
| Stomach (STAD) | 373 tumors |
| Thyroid (THCA) | 216 tumors |
| Uterine corpus endometrial (UCEC) | 361 tumors |

The Cancer Genome Atlas Research Network (TCGA)

# ATAC-seq identifies shared and unique epigenetic landscape across cancers



MR Corces et al. *Science*. 2018.

TF motif prediction in ATAC-seq peak regions

# CITRUS better predicts gene expression in held-out tumors compared to bilinear models

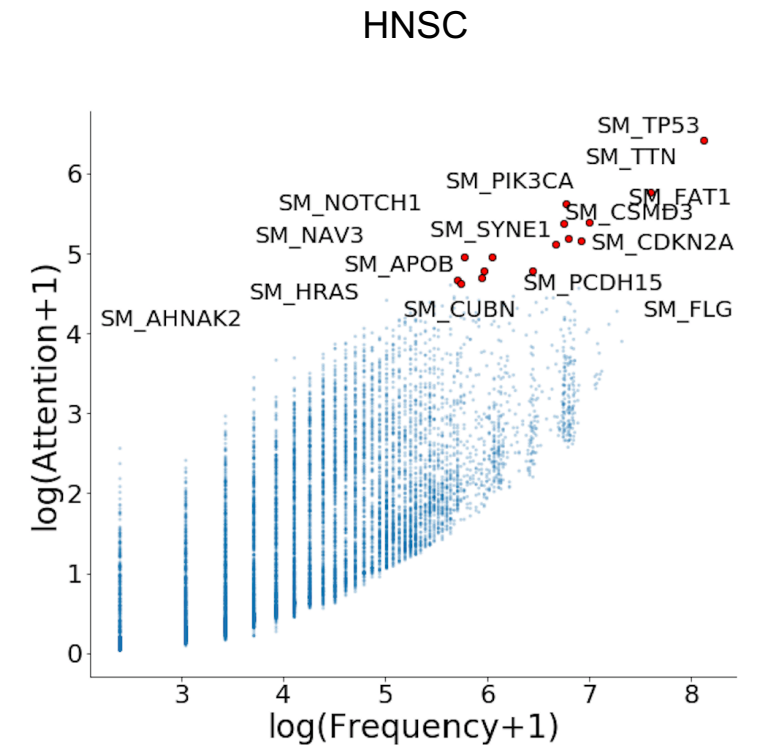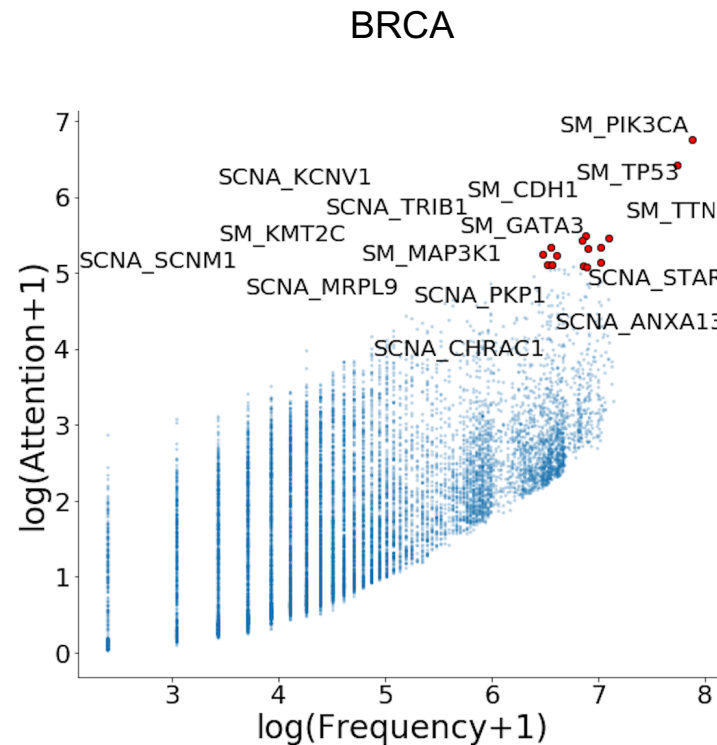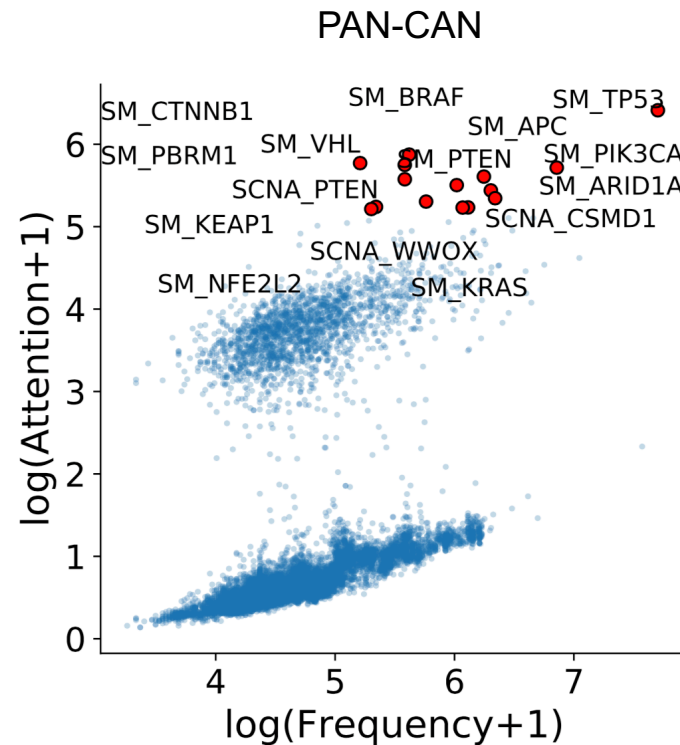- Affinity regression (bilinear) vs. CITRUS (deep learning)



CITRUS vs. Affinity regression

R Pelossof et al. *Nature Biotech*. 2015.
HU Osmanbeyoglu et al. *Nature Comm*. 2017.

# Overall attention weights

- Impacts of somatic alterations

# Clustering based on inferred TF activity largely recovered the distinction between the major tumor types
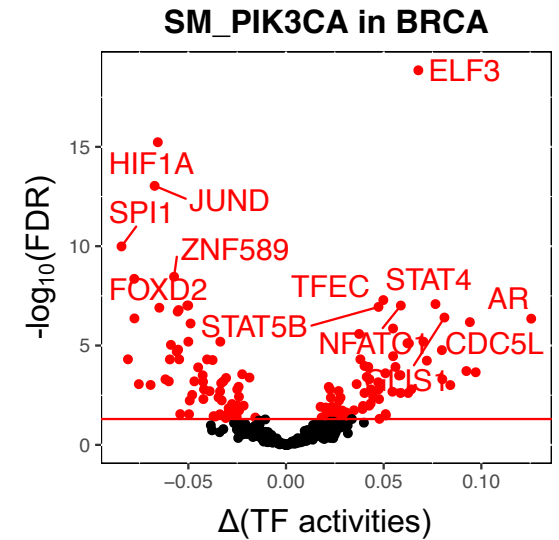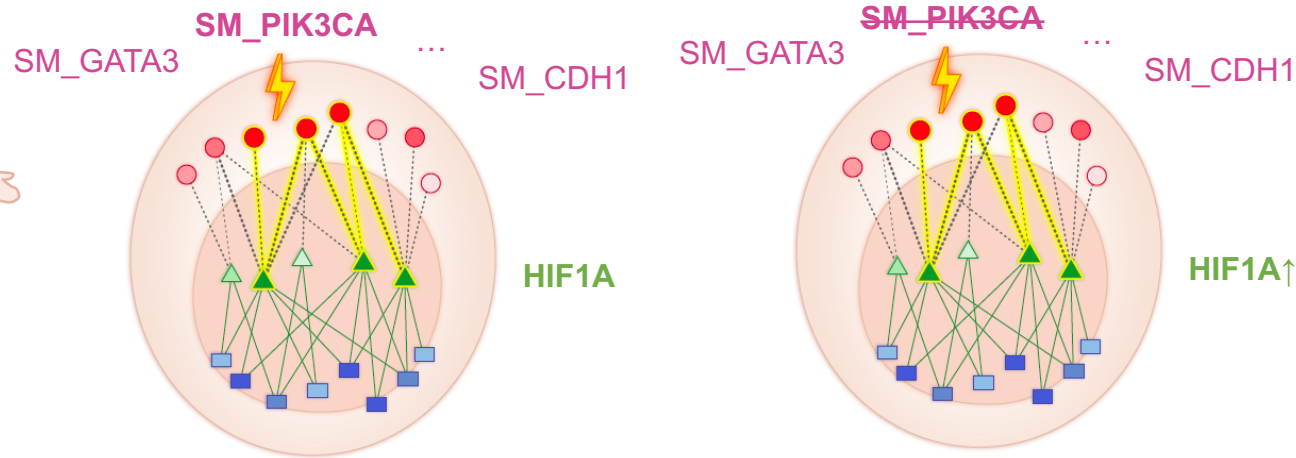


TF activity cancer type clusters

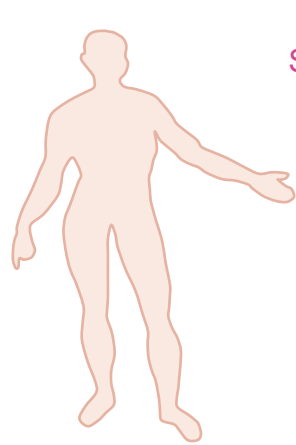# Landscape of mutations and inferred TF activities



CITRUS-inferred TF activities

Somatic mutations

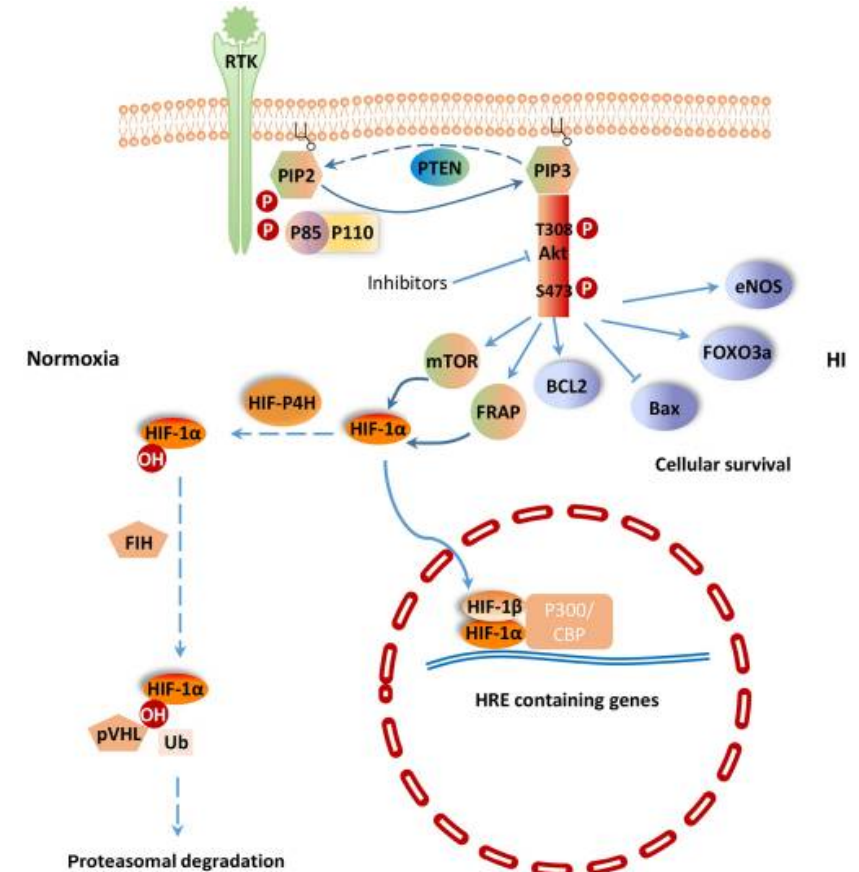Somatic copy number alterations
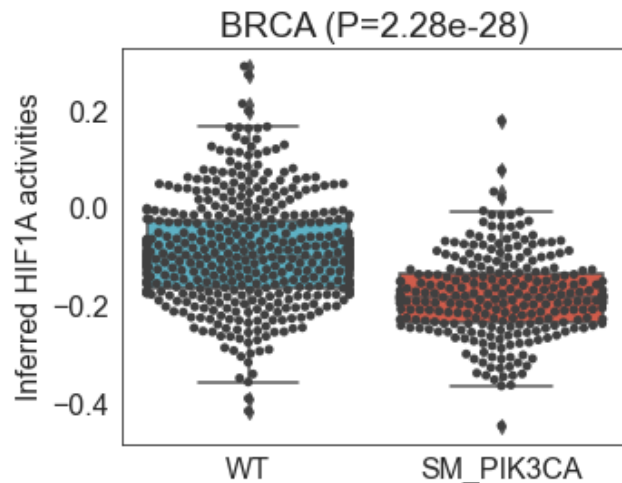
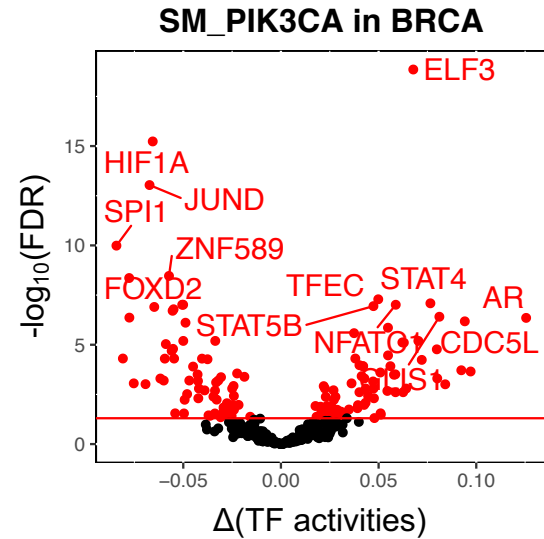Association score := direction*-log₁₀(FDR)

# Impact of mutations on TFs in breast cancer

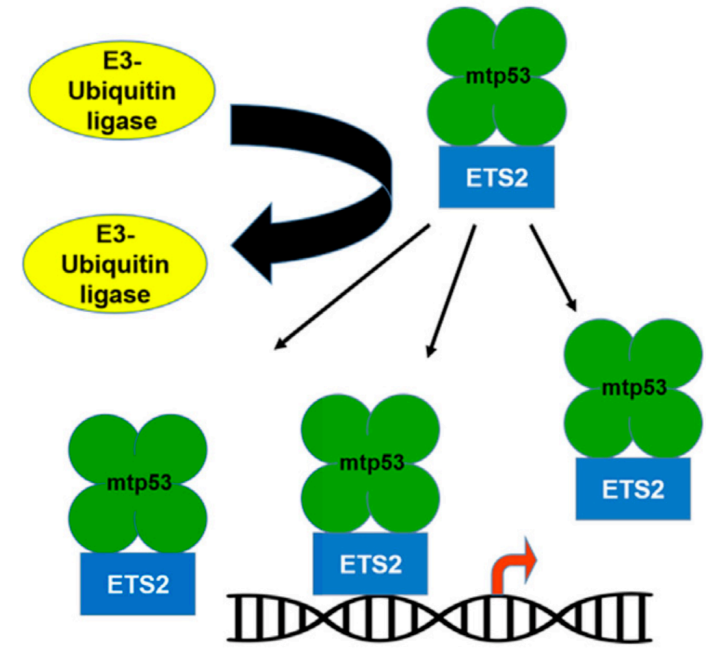- Knock out *in silico*: different from t-test, simulates the knockout of mutations
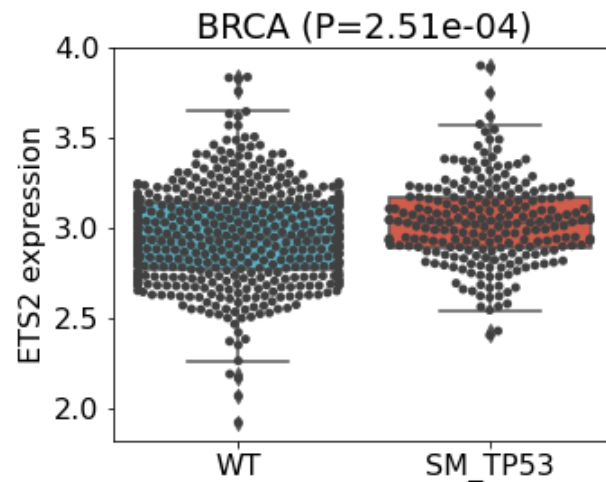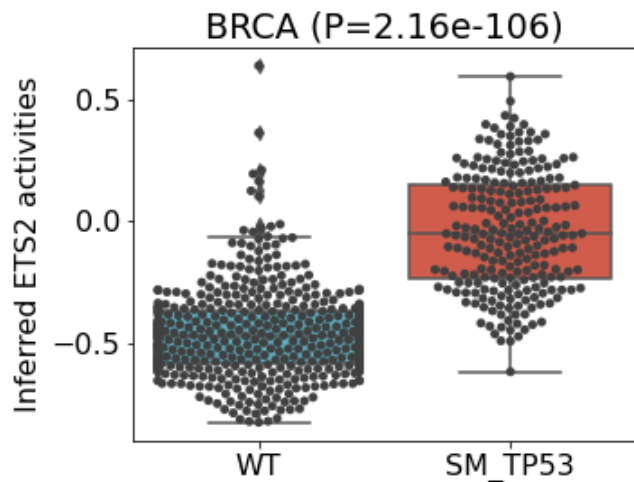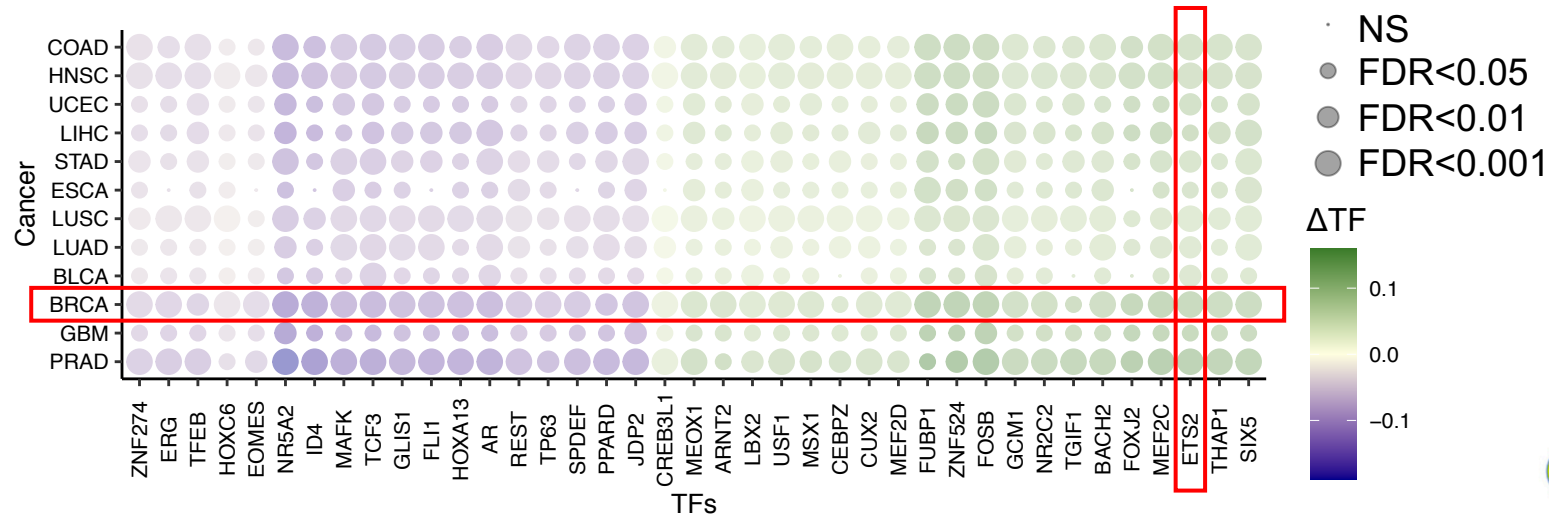
# Impact of PIK3CA mutation on TFs in breast cancer



Z Zhang et al. *Mol Med Rep*. 2018.

# Impact of TP53 mutation across cancers



LA Martinez. *Front Oncol*. 2016.

# Conclusion and future work

- CITRUS: deep learning approach modeling transcriptional programs in pan-cancer

- Utilize self-attention mechanism to capture non-linear effects of mutations

- Integrate ATAC-seq as knowledge base

- Further explore potential clinical relevance
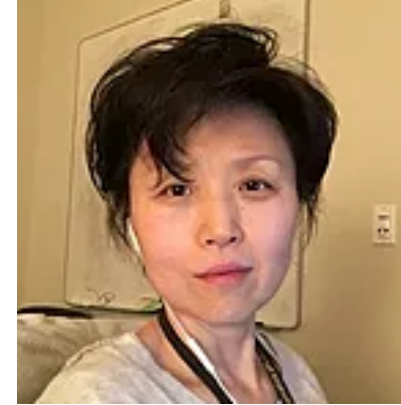
# Acknowledgments



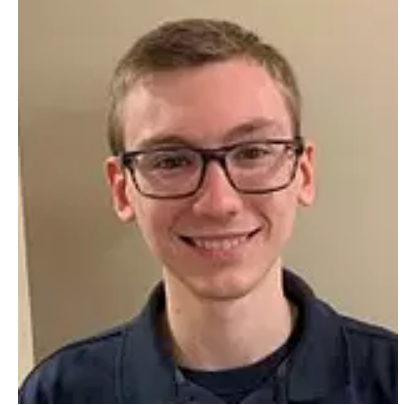**Dr. Hatice Ulku Osmanbeyoglu**
University of Pittsburgh

**Dr. Xinghua Lu**
University of Pittsburgh

**Dr. Russell Schwartz**
Carnegie Mellon University

**Xiaojun Ma**
University of Pittsburgh

**Drake Palmer**
University of Pittsburgh

Looking for students and postdocs!
Please reach out at
osmanbeyogluhu@pitt.edu