

# Interpretable deep learning for chromatin-informed inference of transcriptional programs driven by somatic alterations across cancers

Yifeng Tao<sup>1,†</sup>, Xiaojun Ma<sup>2,3,†</sup>, Drake Palmer<sup>3</sup>, Russell Schwartz<sup>1,4</sup>, Xinghua Lu<sup>2,5</sup> and Hatice Ulku Osmanbeyoglu<sup>2,3,6,7,\*</sup>

<sup>1</sup>Computational Biology Department, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA, <sup>2</sup>Department of Biomedical Informatics, School of Medicine, University of Pittsburgh, Pittsburgh, PA, USA, <sup>3</sup>UPMC Hillman Cancer Center, University of Pittsburgh, Pittsburgh, PA, USA, <sup>4</sup>Department of Biological Sciences, Carnegie Mellon University, Pittsburgh, PA, USA, <sup>5</sup>Department of Pharmaceutical Science, School of Medicine, University of Pittsburgh, Pittsburgh, PA, USA, <sup>6</sup>Department of Bioengineering, School of Engineering, University of Pittsburgh, Pittsburgh, PA, USA and <sup>7</sup>Department of Biostatistics, School of Public Health, University of Pittsburgh, Pittsburgh, PA, USA

Received December 01, 2021; Revised September 23, 2022; Editorial Decision September 27, 2022; Accepted September 29, 2022

## ABSTRACT

Cancer is a disease of gene dysregulation, where cells acquire somatic and epigenetic alterations that drive aberrant cellular signaling. These alterations adversely impact transcriptional programs and cause profound changes in gene expression. Interpreting somatic alterations within context-specific transcriptional programs will facilitate personalized therapeutic decisions but is a monumental task. Toward this goal, we develop a partially interpretable neural network model called Chromatin-informed Inference of Transcriptional Regulators Using Self-attention mechanism (CITRUS). CITRUS models the impact of somatic alterations on transcription factors and downstream transcriptional programs. Our approach employs a self-attention mechanism to model the contextual impact of somatic alterations. Furthermore, CITRUS uses a layer of hidden nodes to explicitly represent the state of transcription factors (TFs) to learn the relationships between TFs and their target genes based on TF binding motifs in the open chromatin regions of tumor samples. We apply CITRUS to genomic, transcriptomic, and epigenomic data from 17 cancer types profiled by The Cancer Genome Atlas. CITRUS predicts patient-specific TF activities and reveals transcriptional program variations between and within tumor types. We show that CITRUS yields biological insights into delineating TFs associated with somatic alterations in indi-

vidual tumors. Thus, CITRUS is a promising tool for precision oncology.

## INTRODUCTION

The complex interplay between signaling inputs and transcriptional responses dictates important cellular functions. Dysregulation of this interplay leads to development and progression of disease, which has been most clearly delineated in the context of certain cancers. Cancer cells acquire somatic alterations that modify signaling and transcriptional programs, leading to profound changes in gene expression. We still lack a complete understanding of how somatic alterations affect cellular function in cancer. To begin to understand these effects, it is important to study somatic alterations within the specific transcriptional context in which they are found. Context- and patient-specific studies can be achieved with machine learning techniques, which are expected to facilitate personalized therapeutic decisions.

In the last decade, a monumental effort has been made to molecularly profile tumors by consortia, including The Cancer Genome Atlas (TCGA) and the International Cancer Genome Consortium (1,2). The multimodal datasets generated by these efforts include gene expression and somatic alterations, such as recurrent mutations and copy number variations (CNVs). The combination of genomic and transcriptomic information enables the integration of transcriptional states with upstream signaling pathways. Several methods have been developed to connect somatic alterations to a prior network or to gene expression (3–9). More recently, the Genomic Data Analysis Network generated assay for transposase-accessible chromatin with

\*To whom correspondence should be addressed. Tel: +1 412 623 7789; Email: osmanbeyoglu@pitt.edu

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

high-throughput sequencing (ATAC-seq) data for a subset of TCGA samples (~500 patients) (10). Although chromatin profiling helps uncover context-dependent and/or non-linear effects of transcription factors (TFs) on gene expression, it has not yet been incorporated into methods that connect somatic alterations to transcriptional programs across cancers. Incorporating DNA sequence information at promoter, intronic, and intergenic enhancers from ATAC-seq tumor profiles using TF motif analysis will improve the modeling of transcriptional regulation and delineate the impact of somatic alterations on transcriptional programs.

Deep learning is a powerful tool for capturing non-linear feature interactions that can explain the underlying biological phenomena. For example, attention mechanism is a deep learning method that has been widely used in computer vision and natural language processing. In contrast to traditional deep learning methods, the self-attention mechanism considers the contextual relationship of the input features and assigns attention weights to each input (11). In general, attention mechanisms improve the performance of deep learning models and increase the interpretability of the models. More recently, attention mechanisms have been applied to cancer genomics for cancer driver gene detection (12), drug response prediction (13) and base editing outcome prediction (14). For example, the genomic impact transformer (GIT) model utilizes a self-attention mechanism to encode the effects of somatic alterations in cancer and uses multi-layer perceptrons to predict differentially expressed genes (12). The attention mechanism enables GIT to select driver mutations that are likely to lead to downstream phenotypes. However, the GIT model lacks interpretability in the sense that it does not model intermediate TFs during modeling signaling from somatic alterations to gene expression programs.

Here, we present Chromatin-informed Inference of Transcriptional Regulators Using Self-attention mechanism (CITRUS), a partially interpretable neural network model with encoder-decoder architecture. CITRUS links somatic alterations to transcriptional programs by modeling the statistical relationships between mutations, CNVs, gene expression, and TF-target gene information derived from ATAC-seq (Figure 1). We show that CITRUS yields important biological insights into dysregulated TFs in individual tumors. Using a systematic *in silico* knock out approach, we identified key TFs associated with major somatic alterations. We believe CITRUS will assist researchers in providing actionable hypotheses for follow-up experiments and developing personalized and targeted therapeutics in a pan-cancer setting.

## MATERIALS AND METHODS

### Data pre-processing

We downloaded the batch normalized RNA-Seq expression levels quantified by RNA-Seq by Expectation Maximization (RSEM) from the Genomic Data Commons (GDC) portal (<https://gdc.cancer.gov/about-data/publications/pancanatlas>). We  $\log_2$ -transformed RSEM values and identified the 2500 most variable genes across samples within a cancer type. Then, we took the union of the

identified genes across cancer types. The final gene set included 5541 genes.

We obtained processed gene-level somatic alterations for each cancer patient from Cai *et al.* (4). Genes with non-synonymous mutations, small insert/deletion, or somatic copy number alteration (deletion or amplification) were given a value of 1, and otherwise were given a value of 0. We removed genes that were not present in at least 4% of samples for each cancer type.

We downloaded the ATAC-seq pan-cancer dataset from the GDC portal (<https://gdc.cancer.gov/about-data/publications/ATACseq-AWG>) (10).

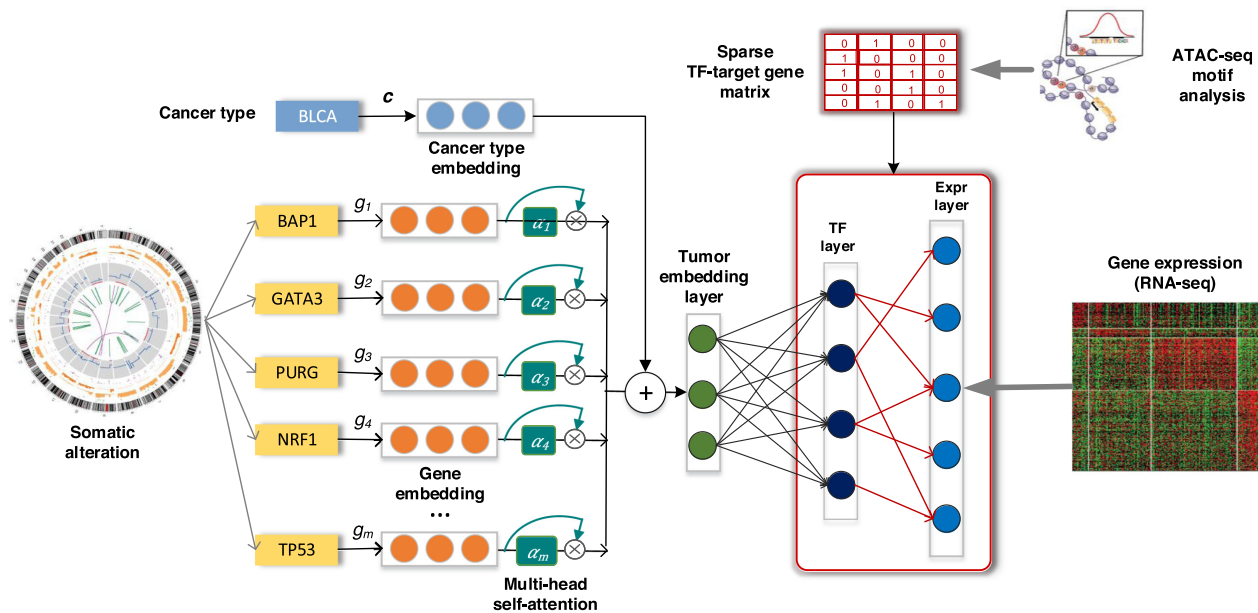
We obtained SILAC-based quantitative phosphoproteomic data set of a spontaneously immortalized non-tumorigenic breast epithelial cell line MCF10A along with two isogenic derivatives generated by knock-in of mutant alleles—one bearing the E545K mutation and the other bearing the H1047R mutation of the *PIK3CA* gene—from the originally published Supplementary Data (15). We also obtained human protein microarray-based AKT1 kinase assays from the originally published Supplementary Data (15).

### Creating TF-target gene matrix

To construct a binary TF-target gene matrix across cancer types, we started with an atlas of chromatin accessible regions derived from all tumor types (~200K peaks) based on ATAC-seq pan-cancer dataset (10). Therefore, this prior matrix was not tissue specific. We represented every gene by its feature vector of TF motif binding presence, where motif information was summarized across all promoters, intronic, and intergenic chromatin accessible sites assigned to the gene. Briefly, using the MEME (16) curated *cis*-BP (17) TF motif binding reference, we scanned the pan-cancer ATAC-seq peak atlas with FIMO (18) to find peaks likely to contain each motif ( $P < 10^{-5}$ ). We filtered TFs that were not expressed in at least 50% of samples in at least one of the seventeen tumor types. Further, similarity of predicted target peak sets was measured using the Jaccard index (size of intersection/size of union). If two TFs had a high Jaccard index ( $> 0.5$ ), we looked at the mean Jaccard index of each TF with all other TFs, and we removed the TF with the largest mean Jaccard index. The final set contained 320 TF binding motifs. Then, we associated each peak to its nearest gene in the human genome using the ChIPpeakAnno package (19). ATAC-seq peaks located in the body of the transcription unit, together with the 100 kb regions upstream of the transcription start site (TSS) and downstream of the 3' end, were assigned to the gene. We converted the assigned ATAC peaks for each gene to a feature vector of motif binding signals by assigning the maximum score of each motif across all peaks to a gene. The final TF-target gene matrix  $C \in \{0,1\}^{k \times l}$  contains a candidate set of associations between TFs and target genes.  $C_{ij} = 1$  when there is a connection from TF  $j$  to the gene  $i$  (red lines connecting the TF layer and target gene expression (Exp) layer in Figure 1).

### CITRUS model

CITRUS is a framework for modeling impact of somatic alterations on transcriptional programs through the latent



**Figure 1.** Overview of CITRUS: An attention-based model with TF-target gene priors. The input to our framework includes somatic alteration and copy number variation, assay for transposase-accessible chromatin with high-throughput sequencing (ATAC-seq), tumor expression datasets and TF recognition motifs. CITRUS takes somatic alteration and copy number variation data as input and encodes them as a tumor embedding using a self-attention mechanism. Cancer type embedding is used to stratify the confounding factor of tissue type. The middle layer further transforms the tumor embeddings into a TF layer, which represents the inferred activities of 320 TFs. Finally, gene expression levels are predicted from the TF activities through a TF-target gene priors constrained sparse layer based on ATAC-seq.

hidden status of the tumor. Figure 1 shows the model architecture with an overall encoder and decoder structure. Somatic gene alteration inputs with  $>20K$  dimensions were encoded into a compressed representation as tumor embedding. Then the tumor embedding, which represents the status of the tumor, was decoded to a large dimension data of gene expression. The encoder-decoder architecture allows the model to capture key features of the high dimension inputs and reduce the data noise as well.

In the encoder module of the CITRUS model, each mutant gene is mapped into the gene embedding, and aggregated together with cancer type embedding to form the tumor embedding through weighted sum. Like ‘word embedding’ in the natural language processing (NLP) field, our gene embedding represents each distinct somatic alteration with a particular continuous number vector such that somatic alteration that share similar biological functions are located close to each other in embedding space, which is a Euclidean space. It captures the functional similarity among somatic alterations perturbing common pathways. Our pre-trained gene embedding was derived from the gene2vec method based on the co-occurrence information of somatic alterations, which utilized the skip gram word2vec algorithm based on the context-target word pairs. If two mutant genes co-exist with a third gene in similar pattern, they tend to share similar gene embeddings. In the CITRUS model, we uploaded the pre-trained gene embedding and further updated/finetuned it under the supervision of mRNA profile. From the point of the implementation, we created embedding space only for somatic alterations of each tumor, so that the sparse high dimension vector of tumor somatic alterations ( $>10k$ ) was converted into a dense

2d array with smaller dimensions  $N \times M$ ,  $N$  is the maximal number of gene mutations/copy number events of each tumor ( $\sim 1000$ );  $M$  is the embedding size (512). This conversion also set up the foundation for self-attention mechanism which in turn to feeds into our attention weight analyses.

Cancer type embedding acts in a similar way by converting each cancer type into a vector, then merging with gene embedding to form a personalized tumor embedding. The added cancer type information provides additional guidance for the model training in a multi-cancer type data setting and further improves the prediction performance and shortens the training time for the model to converge (see the Supplementary Figure 1). The condensed tumor embedding, which comprises the attention-weighted gene embedding, and the specified cancer type embedding served as a highly informative entry point for the multi-layer perceptron decoder from which derived our transcriptional factor activities.

We design a self-attention mechanism which assigned importance weights to input features (somatic alterations) through the model training. Formally, given a specific tumor  $t$ , with the cancer type  $s$ , we have a set of somatic alterations in the tumor  $\{g_u\}_{u=1}^m$  where  $m$  is number of mutant genes. The encoder module first maps each gene  $g$  (it is  $g_u$  here, but we omit the subscript for notation simplicity) into its corresponding gene vector  $e_g$ . Then, the encoder utilizes the multi-head self-attention mechanism to calculate the weighted sum of both the gene embeddings and the cancer type embedding:

$$e_t = e_s + \alpha_1 e_1 + \alpha_2 e_2 + \alpha_3 e_3 + \dots + \alpha_m e_m$$

The self-attention mechanism takes the gene embeddings of all mutated/altered genes as an input and outputs the attention weights  $\{\alpha_u\}_{u=1}^m$  through a sub-neural network. To be more specific, in the case of single-head attention mechanism, we first calculate the unnormalized attention weights  $\{\beta_{g,j}\}_{g=1}^m$ :

$$\beta_{g,j} = \theta_j^T \tanh(W_0 e_g), \quad g = 1, 2, \dots, m.$$

where  $\theta_j$  is the single-head parameter. Then we normalize the attention weights:

$$\alpha_{1,j}, \alpha_{2,j}, \dots, \alpha_{m,j} = \text{softmax}(\beta_{1,j}, \beta_{2,j}, \dots, \beta_{m,j}).$$

In the case of multi-head attention mechanism, we have parameters for multiple heads  $\{\theta_j\}_{j=1}^h$ . The final attention weights are the sum of single head weights:

$$\alpha_g = \alpha_{g,1} + \alpha_{g,2} + \dots + \alpha_{g,h}, \quad g = 1, 2, \dots, m$$

The attention mechanism captures the context of co-existing somatic alterations and their complex interactions, which is lost in simpler models. Interested readers can find the mathematical details of self-attention mechanisms in the cited reference (12).

In the decoder module of CITRUS, we first infer the TF activities from the encoded tumor embedding  $e_t$ :

$$e_f = \tanh(W_f e_t + b_f).$$

We used tanh activation instead of ReLU operation, which is more widely used in deep learning, because it has similar performance to that of ReLU in our model and generates more biologically meaningful results (e.g., distribution of TFs  $e_f$ ). Finally, CITRUS predicts cancer type-specific mRNA expression from TF activities:

$$\hat{y} = W e_f + b_r$$

where  $W$  corresponds to the sparse TF–target gene matrix constrained by the prior  $C \in \{0, 1\}^{k \times l}$ . More specifically, to integrate priors into our model,  $W$  shares the same shape with prior  $C$ , and  $W_{i,j}$  is allowed to be nonzero only when  $C_{i,j} = 1$ , and  $W_{i,j}$  is constrained to be non-negative value. We use mean square loss function as:  $MSE(y, \hat{y})$ . The TF–target gene matrix contains the binary constraints that element-wisely define the sparse connections between TFs and potential target genes. During the training of model, only non-zero part of the parameters is allowed to be back-propagated and updated. The rest part of the weight parameters are always zero values throughout the training and inference.

One might use other common approaches to integrate prior  $C$  into the  $W$ , i.e. by applying a Gaussian prior to  $W$ , which is equivalent to adding an additional penalty to the loss function  $\sum_{i,j:C_{i,j}=0} (W)_{i,j}^2$ . However, this ‘soft’ constraint

tends to generate less stable TF layers across different runs of training compared to the ‘hard’ constraints shown in our model.

To prevent overfitting and to increase robustness to noise, we introduced additional dropout operations with a dropout rate of 0.2 after the input layer, activated tumor embedding layer, and activated TF layer.

## Training and evaluation

We implemented CITRUS through the PyTorch package (<https://pytorch.org/>), and training was performed using the Adam optimizer with default parameters except for the learning rate (15) and weight decay. We set the learning rate to  $1 \times 10^{-3}$  and the weight decay to  $1 \times 10^{-5}$ . We used early stopping with patience of 30 steps to stop training.

For statistical evaluation, we computed the mean Spearman correlation ( $\rho$ ) between predicted and measured gene expression profiles for each tumor. We held out 20% of samples as the testing set with stratified splitting, preserving the percentage of samples for each cancer type (Supplementary Table 1). The remaining 80% of samples were used for training and validation to determine the optimal hyperparameters of the model. For hyperparameter setting, we did 5-fold cross validation by stratified data splitting into 5-fold. Each fold was used for validation after training on the other 4-fold. We utilized mean of the performance over all runs to decide the optimal values of hyperparameters such as the learning rate and batch size. Then, we applied the trained model with selected hyperparameters to the testing set for performance evaluation. To increase the stability of inferred TF activity analysis, we assembled multiple CITRUS models trained with different random initialization state and integrated the TF layer based on the average of 10 trials.

**Parameter selection:** CITRUS includes >10 hyperparameters that are described in the following paragraphs. These hyperparameters were tuned for optimal performance in the validation set. Ideally, hyperparameter optimization is performed using a grid search of all parameters. However, this is not practical due to the tremendous computational cost. For example, three options for each parameter leads to  $3^{10}$  possible combinations for just 10 parameters. In addition, we guide the performance metric by  $k$ -fold cross-validation, and the total experiments necessary would be  $5 \times 3^{10}$  ( $k = 5$ ). Therefore, our hyperparameter tuning strategy combined automatic and manual tuning. First, we created empirical settings for each parameter and randomly selected a set of parameters from 100 combinations. We utilized the best-performing settings to narrow down the preliminary decisions and correlation among parameters. Then, we tuned parameters independently or in sub-groups manually or by grid search.

**Model robustness:** The learning rate is perhaps the most important hyperparameter in neural network training. We first tested the learning rate in a range of settings [ $10^{-5}$ ,  $10^{-4}$ ,  $10^{-3}$ ,  $10^{-2}$  ...], starting with the lowest setting and progressing to larger values until validation loss started to diverge. We found that if the learning rate was too small, overfitting occurred and picked up input noise. Additionally, overfitting reduced the number of driver genes that were covered in downstream attention weight analyses. If the learning rate was too big, however, the model could not converge to an optima and yielded higher validation loss. Ultimately, we selected learning rates of  $10^{-3}$  and  $10^{-4}$  and applied a weight penalty (weight decay) to find an optimal combination of settings. We set the weight decay range from  $10^{-6}$  to  $10^{-4}$  and performed a grid search. The optimal settings for learning rate and weight decay were determined to be  $10^{-3}$  and  $10^{-5}$ , respectively. Although large batch sizes

can accelerate learning rates and training, our experiments indicated that a learning rate of  $10^{-3}$  was the largest value that maintained validation accuracy when tested on increasing batch sizes (16, 64, 100 and 300, which is the maximum value that could run in GPU). We found that larger batch sizes tended to have slightly higher gene-wise correlation at the cost of longer training time. To balance execution time, we selected a batch size of 100. The early stopping patience setting is also related to the learning rate and batch size. Specifically, higher learning rates and larger batch sizes require smaller patience to stop training. Higher patience settings may otherwise cause overfitting. Using our selected learning rate and batch size settings, a patience of 30 was generally sufficient to maintain training without stopping too early (underfitting) due to fluctuation and without halting too far from the optima (overfitting). We validated a patience setting of 30 by comparing it with a case of overfitting. We selected the lowest loss point in the overfit training and measured how far it was from the model with early stopping. During early stages of training, the model showed an initial drop in validation performance followed by a rise. To avoid this inconsistency, we did not apply early stopping for the first 180 test steps. To test the attention mechanism, we created a mesh grid for two attention sizes (256, 128) and four attention head settings (32, 16, 8, 4). We then performed an exhaustive grid search within these settings. Based on prediction performance, we selected 256 and eight as the optimal values for attention size and attention head, respectively.

Finally, we fine-tuned our model by adjusting the dropout rate. Because we used weight decay for regularization, dropout is considered a secondary regularization for our model. In addition to hidden layer dropout, we also applied dropout to our input to reduce input noise and network redundancy and to generate a more stable hidden TF layer. We tested a sequence of five dropout rates (0.1, 0.2, 0.3, 0.4, 0.5). All dropout rate settings yielded performances above 0.9 for average sample correlation in the testing set. We determined the dropout rate optimal value (0.2) primarily based on driver gene coverage in self-attention analysis.

As we used an early stopping mechanism, we set the maximum iteration parameter to 1000. This setting ensures that the training process stops either once the patience setting is satisfied or once the maximum iterations is reached. Code testing and quick runs were performed with a maximum iteration of one.

We tested two activation functions: 'ReLU' and 'tanh'. Although both activation functions performed similarly, 'tanh' generated more biologically meaningful results and was selected. We also tested l2, minimax and standard normalization (scale) to normalize gene expression and found that scale normalization generated the best prediction accuracy for our model settings.

### Training the affinity regression (AR) models

AR is an algorithm for efficiently solving a regularized bilinear regression problem (20–23) and was defined in our model as follows. For a data set of  $M$  tumor samples profiled using RNA-seq with  $N$  genes, we let  $\mathbf{Y} \in \mathbb{R}^{N \times M}$  be the  $\log_{10}$

gene expression profiles of tumor samples. Each column of  $\mathbf{Y}$  corresponds to an RNA-seq experiment for a cancer type. We define the TF attributes of each gene in a matrix  $\mathbf{D} \in \mathbb{R}^{N \times Q}$ , where each row represents a gene, and each column represents a TF vector. The TF vector indicates whether there is a binding site for the TF on each gene based on ATAC-seq data. We define the somatic alteration attributes of tumor samples as a matrix  $\mathbf{P} \in \mathbb{R}^{M \times S}$  where each row represents a tumor sample, and each column represents the somatic alteration status for the tumor sample. We set up a bilinear regression problem to learn the weight matrix  $\mathbf{W} \in \mathbb{R}^{Q \times S}$  on paired TF and somatic alteration features:

$$\mathbf{DWP}^T \sim \mathbf{Y}$$

We can transform the system to an equivalent system of equations by reformulating the matrix products as Kronecker products:

$$\mathbf{DWP}^T \approx \mathbf{Y} \Leftrightarrow (\mathbf{P} \otimes \mathbf{D})\text{vec}(\mathbf{W}) \approx \text{vec}(\mathbf{Y})$$

where  $\otimes$  is a Kronecker product, and  $\text{vec}$  is a vectorizing operator that stacks a matrix and produces a vector. The result of this system is a standard (if large-scale) regression problem. Full details and a derivation of the reduced optimization problem are provided elsewhere (21). For statistical evaluation, we separated datasets by cancer type and conducted 5-fold cross-validation to tune hyperparameters in the training and validation sets as we performed in CITRUS setting.

### In silico knockout analysis

We implemented an *in silico* knock out approach that removes a specific somatic mutation (or copy number variation)  $g$  from all the tumor samples that carry it. The new somatic alteration profiles and the CITRUS-inferred TF activities generate a 'wild type' corpus that does not contain the alteration  $g$ . In contrast, the original samples containing the alteration  $g$  serve as the 'mutant/altered' group. We then conducted t-tests between the mutant and wild type groups to evaluate the impact of mutation  $g$ . This approach captures the contextual effects of mutations through the non-linear attention module of CITRUS and provides a controlled experimental environment that holds all mutations constant except for mutation  $g$ . For complex genotypes, the model explains TF activity across tumors. We then corrected for multiple hypotheses across models, treating inferred TF activities as separate groups of tests.

### Association score between TF activity subtypes and frequent somatic alterations

For each somatic mutation or copy number variation, we calculated the P-value of its frequency in a cancer subtype compared to other subtypes using Fisher's exact test. The P-value was further adjusted through FDR across subtypes. To identify the relative frequency of a somatic alteration in a subtype, we defined an association score, which is the product of the relative frequency direction and  $-\log_{10}\text{FDR}$ .

## Statistical analysis

Statistical tests were performed with the R statistical environment (4.0.2) and Python. For population comparisons of inferred TF activities, we performed Student's t-tests and determined the direction of shifts by comparing the mean of the two populations. We corrected raw *P*-values for multiple hypothesis testing based on two methods: Bonferroni and FDR (BH method).

## RESULTS

### Pan-cancer modeling of transcriptional programs

To systematically interpret somatic alterations within context-specific transcriptional programs and to identify disrupted TFs that drive tumor-specific gene expression patterns across multiple cancer types, we developed CITRUS (Figure 1, see Materials and Methods for details). CITRUS traces biological signaling from somatic alterations to signaling pathways, to TFs, and finally to target gene expression (mRNA levels). To enable this tracing, CITRUS employs an encoder-decoder architecture. The encoder module compresses input somatic alterations into a latent vector variable called a tumor embedding. The decoder predicts TF activities from the tumor embedding and then predicts target gene expression. More intuitively, the model learns the flow of information from somatic alterations to altered activity of TFs to their transcriptional changes in target genes.

We applied this approach to 17 tumors from TCGA and identified key TFs associated with somatic alterations. Our dataset included samples from 17 different tumor types for which mRNA, somatic mutation and copy number variation data were available. ATAC-seq data were available for subset of patients: bladder urothelial carcinoma (BLCA,  $n = 371$ ,  $n_{\text{ATAC-seq}} = 10$ ), breast cancer (BRCA,  $n = 719$ ,  $n_{\text{ATAC-seq}} = 75$ ), cervical squamous cell carcinoma and endocervical adenocarcinoma (CESC,  $n = 267$ ,  $n_{\text{ATAC-seq}} = 4$ ), colorectal adenocarcinoma (COAD,  $n = 271$ ,  $n_{\text{ATAC-seq}} = 41$ ), esophageal carcinoma (ESCA,  $n = 170$ ,  $n_{\text{ATAC-seq}} = 18$ ), glioblastoma multiforme (GBM,  $n = 143$ ,  $n_{\text{ATAC-seq}} = 9$ ), head and neck squamous carcinoma (HNSC,  $n = 475$ ,  $n_{\text{ATAC-seq}} = 9$ ), kidney renal cell-clear carcinoma (KIRC,  $n = 357$ ,  $n_{\text{ATAC-seq}} = 16$ ), kidney renal papillary cell carcinoma (KIRP,  $n = 272$ ,  $n_{\text{ATAC-seq}} = 34$ ), liver hepatocellular carcinoma (LIHC,  $n = 336$ ,  $n_{\text{ATAC-seq}} = 17$ ), lung adenocarcinoma (LUAD,  $n = 459$ ,  $n_{\text{ATAC-seq}} = 22$ ), lung squamous cell carcinoma (LUSC,  $n = 430$ ,  $n_{\text{ATAC-seq}} = 16$ ), pheochromocytoma and paraganglioma (PCPG,  $n = 109$ ,  $n_{\text{ATAC-seq}} = 9$ ), prostate cancer (PRAD,  $n = 449$ ,  $n_{\text{ATAC-seq}} = 26$ ), stomach adenocarcinoma (STAD,  $n = 373$ ,  $n_{\text{ATAC-seq}} = 21$ ), thyroid carcinoma (THCA,  $n = 216$ ,  $n_{\text{ATAC-seq}} = 14$ ), and uterine corpus endometrial carcinoma (UCEC,  $n = 361$ ,  $n_{\text{ATAC-seq}} = 13$ ).

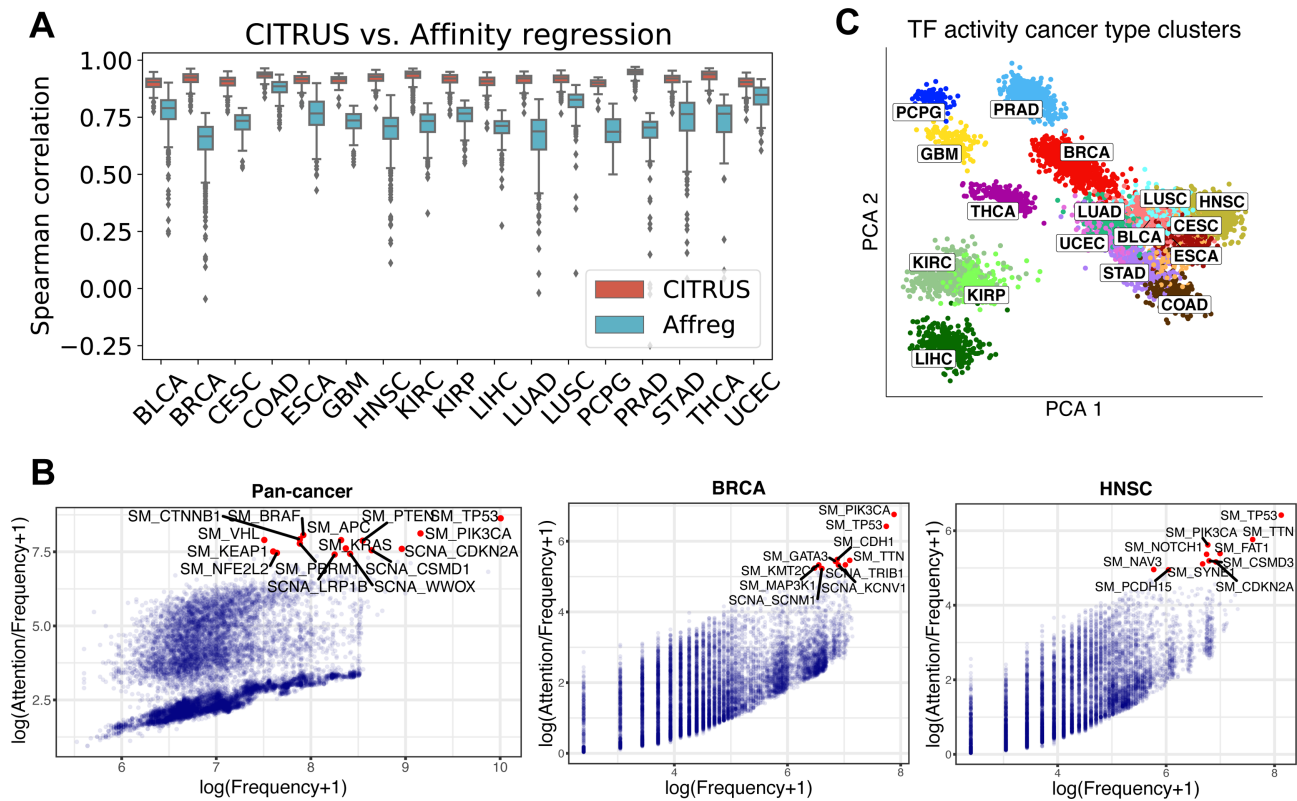
For statistical evaluation, we computed the mean Spearman correlation and mean squared error (MSE) between predicted and measured gene expression profiles on the testing set (see Materials and Methods). CITRUS achieved significantly better performance than a regularized bilinear regression algorithm called affinity regression (AR) (21–24) that was trained independently for each cancer type and ex-

plain gene expression across tumors in terms of somatic alteration status and presence of TF motif binding sites based on a pan-cancer ATAC-seq atlas (Figure 2A and Supplementary Table 2). We also observed better prediction performance with models built using cancer type embedding compared to without cancer type embedding (Supplementary Figure 1). Moreover, models built without cancer type took  $\sim 1.7\times$  time longer to coverage than the cancer type model.

To identify somatic alterations that influenced gene expression programs, we compared the relationship of overall attention weights (inferred by CITRUS) and the frequencies of somatic alterations (used as the control group) across all cancer types and within each cancer type (Figure 2B and Supplementary Figure 2). In general, attention weights were positively correlated with the frequency of somatic alteration. For example, the top altered genes *TP53* and *PIK3CA* had high attention weights. However, our self-attention mechanism assigned low attention weights to many frequently altered genes. We also observed a few infrequently altered genes with high attention weights. For example, the H3K4 methyltransferase *KMT2C* had a high attention weight in BRCA but was infrequently altered. Indeed, *KMT2C* is a key regulator of ER $\alpha$  activity and anti-estrogen response in breast cancer (24,25).

We found genes with high attention weights were enriched for known cancer drivers using the IntOGen<sup>9</sup> database. We first grouped all the genes into two parts with the threshold of 2 ( $\log(\text{attention} + 1) \geq 2$  as the more attended group, and  $\log(\text{attention} + 1) < 2$  as the less attended group). Using Fisher's exact test, we verified that known cancer driver genes were enriched in the highly attention group ( $P = 4.48 \times 10^{-41}$ ) in the pan-cancer analysis. We further examined driver enrichment using the frequently mutated genes. We found the frequently mutated genes with lower attention weight were not significantly enriched in cancer drivers ( $P > 0.05$ , Fisher's exact test) unlike frequently mutated genes with higher attention weight ( $P < 0.05$ , Fisher's exact test). Briefly, we selected the top 100 frequently mutated genes and divided them into two groups (high attention and low attention group) by several different separation thresholds (Supplementary Table 3). For each threshold split, we conducted driver enrichment analyses on these two groups respectively with Fisher's Exact test. It consistently shows that the high attention group has a significant driver enrichment ( $P < 0.05$ ) while the low attention group does not ( $P > 0.05$ ).

We used CITRUS to infer patient-specific TF activities across tumor types. Clustering tumors by these inferred TF activities largely recovered the distinction between major tumor types (Figure 2C). Interestingly, samples with squamous morphology components (BLCA, CESC, ESCA, HNSC and LUSC) clustered together. Tumors with tissue or organ similarities or proximity were also clustered together. These included neuroendocrine and glioma tumors (GBM and PCPG), clear cell and papillary renal carcinomas (KIRC and KIRP), a gastrointestinal group (COAD, and STAD), and breast and endometrial cancer (BRCA and UCEC). We also observed similar clustering of the tumor embeddings (Supplementary Fig. 3).



**Figure 2.** CITRUS models the impact of somatic alterations on gene expression programs. (A) Performance of CITRUS in each cancer type compared to the regularized bilinear regression method Affinity regression (Affreg). Boxplots show the mean Spearman correlations between predicted and actual gene expression based on CITRUS (orange) and Affreg (light blue) in TCGA datasets for each cancer type. Both CITRUS and Affreg were tuned on the same training and validation sets and evaluated on the same testing set. (B) Somatic alteration frequencies and CITRUS-inferred attention weights of genes. Cumulative pan-cancer results are shown on the left, and individual BRCA and HNSC results are shown in the middle and on the right, respectively. See Supplementary Fig. 1 for complete results from each cancer type. (C) Principal component analysis (PCA) of TF activity colored by cancer type. Standard TCGA tumor symbols are used to indicate tumor type.

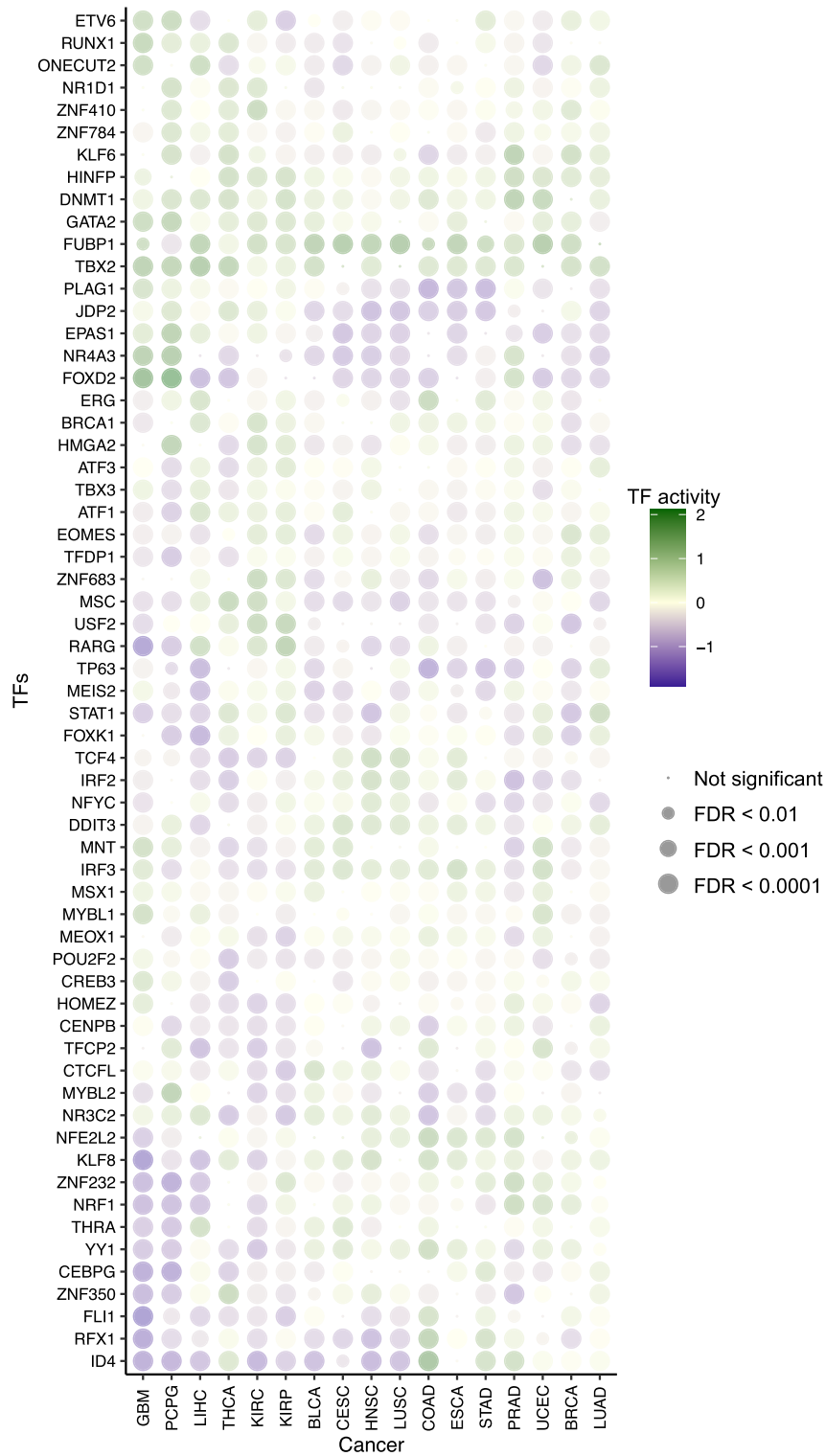
Next, we assessed TF–tumor type associations by *t*-test and compared inferred TF activities between samples in each tumor type versus those in all other tumor types. We corrected for false discovery rate (FDR) across TFs and identified significant shared and cancer-specific TFs, which are listed in Supplementary Data 1. The average TF activity and significance of the four most significant TFs in each cancer are shown in Figure 3. Our results highlight both known and novel cancer-specific TF regulators. For example, FUBP1, which regulates *c-Myc* gene transcription, had significantly higher inferred activity in many cancer types, including LIHC, HNSC, BLCA, ESCA, CESC, LUSC, PRAD, BRCA and UCEC. Consistent with previous reports, IRF3 activity was significantly higher in GBM (26). KLF8 had decreased activity in GBM, LIHC and KIRC, which is consistent with its role in suppressing cell apoptosis during tumor progression (27). Additionally, YY1, which regulates various developmental processes (28), had increased activity in CESC and COAD.

#### Cancer subtype identification from CITRUS-inferred TF activity and somatic alterations

Next, we asked whether CITRUS could identify cancer subtypes based on the TF activity associated with somatic alterations. We conducted *k*-means clustering of inferred TF

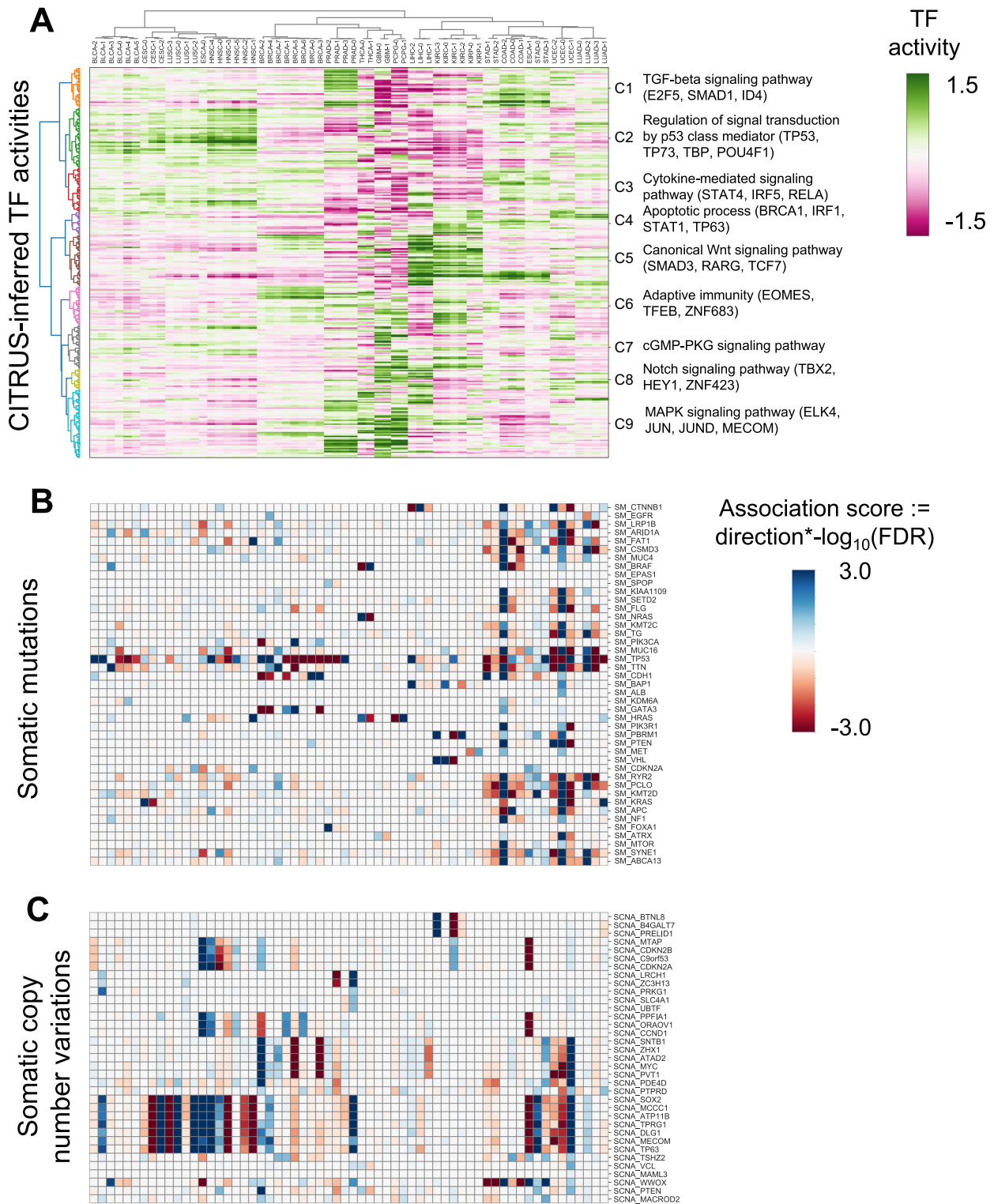
activities for each cancer type to define subtypes, and then we conducted hierarchical clustering of both the cancer subtypes and TF activities. Figure 4 shows the clustering of subtypes by CITRUS-inferred mean TF activities and corresponding somatic alteration associations (see Methods). We observed major differences in mean TF activities across cancer types and minor but significant differences within cancer types. Variations within a cancer type may arise from distinct mutation or CNV profiles of subgroups. For example, clustering by TF activities revealed subclasses of CESC enriched with *KRAS* mutations; KIRC enriched with *VHL*, *BAP1*, *PBRM1* and *TP53* mutations; LIHC enriched with *CTNNB1*, *BAP1* and *TP53* mutations; THCA enriched with *NRAS*, *HRAS* and *BRAF* mutations; and PCPG enriched with *HRAS* mutations (multiple hypothesis corrected Fisher's exact test *P*-value < 0.05).

As our goal was to decipher cancer-specific downstream effects of targeted therapies and to discover secondary targets for combination drug strategies, we developed a systematic statistical approach for modeling the impact of somatic alterations on TF activity. We implemented an *in silico* knock out approach that removes a specific somatic mutation (or CNV) *g* from all carrier tumor samples in each TCGA cancer study and then predicts altered TF activity (see Methods). Using this approach, we were able to identify TFs whose inferred activity was significantly dysregu-



**Figure 3.** CITRUS identifies regulatory features of tumor types. Dot plot shows the mean inferred TF activity differences between samples in a given tumor type versus those in all other tumor types by *t*-test. We corrected for FDR across TFs for each pairwise comparison and identified significant TFs. The complete results are included in Supplementary Data 1. The dot size indicates  $-\log_{10}(\text{FDR})$ . For clarity, the union of the top four significant TFs in each cancer type is shown.





**Figure 4.** Landscape of somatic alterations and inferred TF activities. (A) Heatmap shows tumor subtypes clustered by mean inferred TF activity. The color scale is proportional to TF activity. (B, C) Heatmaps of association scores for (B) mutations and (C) copy number variations. Association scores were calculated by multiplying the  $-\log_{10}$  FDR by the direction derived from Fisher's exact test.

lated by somatic alterations in known cancer driver genes. Figure 5A demonstrates TF activities that were associated with somatic alterations in UCEC. CITRUS identified mutations in *PIK3CA*, *PTEN*, *KRAS*, *TP53* and *CTNNB1* that were significantly associated with various TF activities across UCEC tumors (~66% of tumors have *PTEN* inactivating mutations, ~50% have *PIK3CA* activating mutations, ~38% have *TP53* mutations, ~26% have *CTNNB1* mutations, and ~20% have *KRAS* mutations). UCEC samples with *PTEN* mutations were mutually exclusive with *TP53*, *CTNNB1*, and *KRAS* mutations and showed distinct TF activity patterns. Mutations in *PTEN* that inactivate its phosphatase activity result in increased PI3K signaling. Consistent with this effect, TFs associated with *PTEN* mutations were involved in cell cycle and differentiation, including E2F5, TP63, ELF3, DBP, ZKSCAN3, LHX2, HOXB6, SOX9, DBP, MYLB1 and GLIS1. TFs associated with *CTNNB1* mutant status were involved in WNT and TGF-beta signaling including TCF7, TCF7L2, TCF7L1, FOXH1, EMX1 and MYBL1.

Similarly, CITRUS identified TF activities that were associated with somatic alterations in BRCA (Figure 5B). Mutations in *PIK3CA*, *PTEN*, *MAP2K4*, *GATA3*, *TP53* and *CDH1* were significantly associated with various TF activities. In BRCA, ~36% of tumors have *PIK3CA* activating mutations, ~35% have *TP53* mutations, ~15% have *GATA3* mutations, ~15% have *CDH1* mutations, ~10% have *PTEN* mutations, and ~7% have *MAP2K4* mutations. Activating mutations in *PIK3CA* often occur in one of three hotspot locations (E545K, E542K and H1047R) and promote constitutive signaling through the pathway. TFs associated with *PIK3CA* mutations were involved in WNT signaling, epithelial-mesenchymal transition, and cancer stem cell transition, including ELF3, TFEC, STAT4, STAT5B, NFATC1, GLIS1, CDC5L and AR. We also examined protein microarray-based AKT1 kinase assay and SILAC-based phosphoproteomic data from isogenic knock-in breast cell lines harboring mutations of *PIK3CA* (15). Of 20 TFs represented in the phosphoproteomic data associated with mutant *PIK3CA*, 16 of them associated with *PIK3CA* mutation in our analysis (Supplementary Table 4). Moreover, of the seven TFs identified as AKT substrates, six of them were associated with *PIK3CA* mutation in our analysis (FDR-adjusted  $P$ -value < 0.05) (Supplementary Table 5).

BRCA samples with *PIK3CA* and *TP53* mutations were mutually exclusive, and our *in silico* knock out analysis associated distinct TFs with these mutations. *TP53* mutant tumors were associated with increased activity of TFs that have roles in tumor growth, such as ETS2 and FOSB, growth modulation, such as THAP1, CREB3L1 and CEBPZ, and development, such as MEF2C/D, MEOX1 and MSX1. We performed similar analyses for other cancer types (Supplementary Figure 4).

Although the TFs affected by some somatic alterations differed between cancer types, mutation of *TP53* was associated with similar TFs across cancer types (Supplementary Figure 5). *TP53* is one of the most frequently inactivated tumor suppressor genes that suffers from missense mutations in human cancer. These missense mutations result in the expression of a mutant form of p53 protein. We observed that

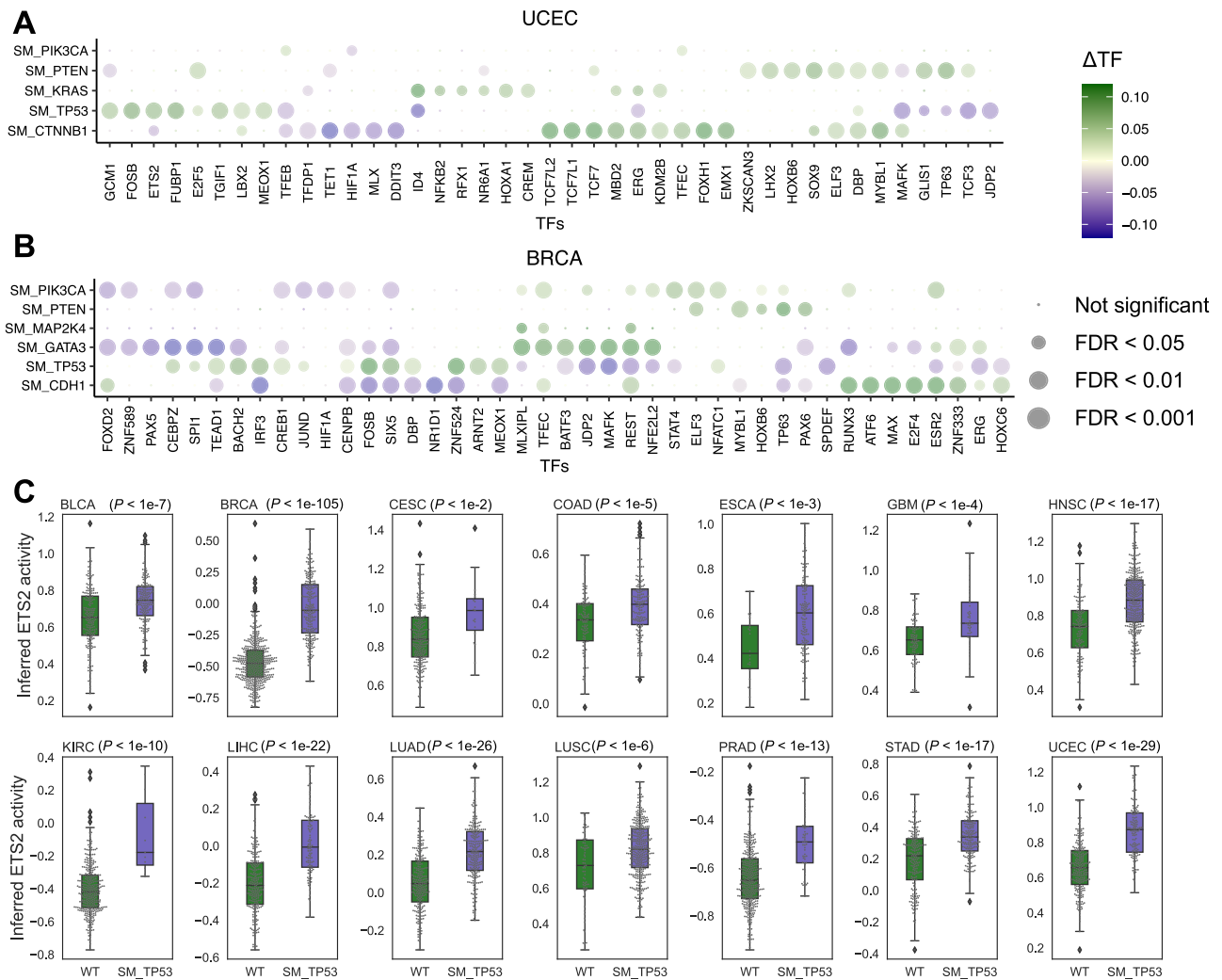
inferred *TP53* activity was lower when *TP53* was mutated compared to wild type for most cancer types (Supplementary Figure 6). Mutant p53 protein can disable other tumor suppressors (e.g. p63 and p73) or enable oncogenes, such as ETS2 (29). Indeed, the inferred TF activity of ETS2 was higher in mutant versus WT *TP53* tumors across cancers (Figure 5C); however, these differences were not as significant at the gene expression level (Supplementary Figure 7). We also observed significant upregulation of motif-based targets of ETS2 based on ATAC-seq relative to all genes in the TCGA tumor data excluding GBM ( $P$ -value <  $1e-6$ , Kolmogorov-Smirnov test, Supplementary Figure 8) consistent with our analysis.

## DISCUSSION

Analysis of the regulatory network in tumor datasets is challenging due to the complexity of the cancer genome (e.g. aneuploidy, CNVs, structural variation, and mutations). CITRUS provides a systematic framework for integrating regulatory genomics with tumor expression and somatic alterations to better understand how expression programs are affected by somatic alterations in cancers and to infer patient-specific TF activities. Our method uses a deep learning framework called a self-attention mechanism to capture the complex contextual interactions between somatic alterations. For a more accurate representation of TF-target gene relationships, we leveraged ATAC-seq tumor data from TCGA patients. CITRUS is designed to capture the flow of information from altered genes (e.g. signaling proteins) to TFs to target genes, and our *in silico* knock out analysis predicts the causal impact of somatic alterations. Joint modeling across different tumor types also revealed patient subgroups associated with somatic alterations.

The key advantage of the CITRUS model is that it enables the detection and representation of the functional impact of somatic alterations on transcription factors, which in turn enables detection of common mechanisms of tumors even if tumors host different somatic alteration. Further, a specific mutated gene usually contains a specific type of mutations, so in a lot of cases, the mutated gene is already very informative of the mutation types. In cases where a somatic alteration is associated with the activity of a targetable TF or their upstream/downstream component the knowledge of putative downstream positive and negative effectors can aid in the identification of combination therapies. For example, one could perform shRNA or CRISPR/Cas screening of downstream TFs to identify those whose knockdown/deletion leads to reduced/enhanced proliferation. This mechanistic information could inform the development of combination therapies, e.g. by selecting agents known to alter activity of these specific TFs.

One limitation of the TF binding motif approach utilized by CITRUS is that TFs of the same family often share a similar motif and thus are difficult to disambiguate. Therefore, TF motifs may encompass the activities of multiple TFs. Moreover, co-binding TF binding patterns (e.g. AP-1-IRF complexes) can be biologically meaningful for gene expression and are not currently represented in our model. Future models will work to represent these composite elements as features. Another limitation is that we do not represent



**Figure 5.** Somatic alterations are associated with dysregulated TF activity. Impact of somatic alterations on individual TFs based on *in silico* knock out experiments in (A) UCEC and (B) BRCA datasets from TCGA. The dot plot shows mean TF activity, and dot size indicates  $-\log_{10}(\text{FDR})$ . Vertical axis are mutations and horizontal axis are TFs associated with at least on the mutations. The experiment and control groups are whether the mutation present or not in the patient. To evaluate the impact of mutations on TF activities, we conducted Student's *t*-test on the two groups. See Supplementary Figure 4 for the full list of cancer types. (C) Inferred ETS2 activity in TCGA studies and impact of *TP53* mutations. Tumors with mutant *TP53* have significantly higher ETS2 activity than WT tumors ( $P < 0.01$ , *t*-test). This association is not significant using mRNA levels of *ETS2* (Supplementary Figure 6). Box edges represent the upper and lower quantile with median value shown as a bold line in the middle of the box. Whiskers extend to 1.5 times the quantile.

directionality in the TF–target gene priors (i.e. whether a gene is activated or repressed by a TF). Prior knowledge of whether the TF is acting as an activator or as a repressor would add meaningful interpretation to inferred TF activities. These limitations may confound the interpretation of the activity of TFs with context-specific activator and repressor roles. Further, regulatory network analysis of tumor datasets is also complicated by the presence of stromal/immune cells within the tumor and the heterogeneity of the cancer cells themselves. However, our framework can be extended to model single-cell RNA-seq or deconvoluted RNA-seq via computational methods.

Despite these limitations, modeling the impact of somatic alterations on transcriptional programs may ultimately enable the development of individualized therapies, aid in understanding mechanisms of drug resistance, and allow the identification of biomarkers of response. We anticipate

that computational modeling of transcriptional regulation across different tumor types will emerge as an important tool in precision oncology, aiding in the eventual goal of selecting the best therapeutic option for individual patients.

#### DATA AVAILABILITY

ATAC-seq data are available in the public repository Genomic Data Commons (<https://gdc.cancer.gov/about-data/publications/ATACseq-AWG>). RNA-seq gene expression, somatic mutation, copy number variation, and clinical data are available in a public repository from TCGA's Firehose data run (<https://confluence.broadinstitute.org/display/GDAC/Dashboard-Stddata>). Only the samples 'whitelisted' by TCGA for the Pan-Cancer Analysis Working Group were used in the study. For our analysis, we only used samples with parallel RNA-seq, somatic mutation,

and GISTIC copy number data. Processed input and output files have been made available at the supplementary website for the paper: <https://sites.pitt.edu/~xim33/CITRUS>.

## CODE AVAILABILITY

The software for CITRUS is available at <https://github.com/osmanbeyoglulab/CITRUS>.

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## ACKNOWLEDGEMENTS

The results published here are based on data generated by The Cancer Genome Atlas project established by the NCI and NHGRI (accession number: phs000178.v7p6). Information about TCGA and the investigators and institutions that constitute TCGA research network can be found at <http://cancergenome.nih.gov/>. We thank Jacob Stewart-Ornstein, Eneda Toska, Jing Hong Wang and Harinder Singh for helpful discussions. The content of this manuscript is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health or other funding agencies. The Pennsylvania Department of Health specifically disclaims responsibility for any analyses, interpretations, or conclusions.

## FUNDING

National Institutes of Health [R00 CA207871 to H.U.O.]; National Institutes of Health [R01 LM012011 to X.L.]; Fellowship in Digital Health from the Center for Machine Learning and Health at Carnegie Mellon University (to Y.T.); UPMC-ITTC fund (to R.S.); National Institutes of Health [R01HG010589, R21CA216452 to R.S.]; Pennsylvania Department of Health [FP00003273 to R.S.]; Mario Lemieux Foundation (to R.S.); AWS Machine Learning Research Award (to R.S.). Funding for open access charge: National Institutes of Health.

*Conflict of interest statement.* None declared.

## REFERENCES

1. ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium (2020) Pan-cancer analysis of whole genomes. *Nature*, **578**, 82–93.
2. Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.R., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C. and Stuart, J.M. (2013) The cancer genome atlas pan-cancer analysis project. *Nat. Genet.*, **45**, 1113–1120.
3. Wang, Z., Ng, K.S., Chen, T., Kim, T.B., Wang, F., Shaw, K., Scott, K.L., Meric-Bernstam, F., Mills, G.B. and Chen, K. (2018) Cancer driver mutation prediction through bayesian integration of multi-omic data. *PLoS One*, **13**, e0196939.
4. Cai, C., Cooper, G.F., Lu, K.N., Ma, X., Xu, S., Zhao, Z., Chen, X., Xue, Y., Lee, A.V., Clark, N. *et al.* (2019) Systematic discovery of the functional impact of somatic genome alterations in individual tumors through tumor-specific causal inference. *PLoS Comput. Biol.*, **15**, e1007088.
5. Hofree, M., Shen, J.P., Carter, H., Gross, A. and Ideker, T. (2013) Network-based stratification of tumor mutations. *Nat. Methods*, **10**, 1108–1115.
6. Paull, E.O., Carlin, D.E., Niepel, M., Sorger, P.K., Haussler, D. and Stuart, J.M. (2013) Discovering causal pathways linking genomic events to transcriptional states using tied diffusion through interacting events (TieDIE). *Bioinformatics*, **29**, 2757–2764.
7. Basha, O., Mauer, O., Simonovsky, E., Shpringer, R. and Yeager-Lotem, E. (2019) ResponseNet v.3: revealing signaling and regulatory pathways connecting your proteins and genes across human tissues. *Nucleic Acids Res.*, **47**, W242–W247.
8. Bashashati, A., Haffari, G., Ding, J., Ha, G., Lui, K., Rosner, J., Huntsman, D.G., Caldas, C., Aparicio, S.A. and Shah, S.P. (2012) DriverNet: uncovering the impact of somatic driver mutations on transcriptional networks in cancer. *Genome Biol.*, **13**, R124.
9. Ng, S., Collisson, E.A., Sokolov, A., Goldstein, T., Gonzalez-Perez, A., Lopez-Bigas, N., Benz, C., Haussler, D. and Stuart, J.M. (2012) PARADIGM-SHIFT predicts the function of mutations in multiple cancers using pathway impact analysis. *Bioinformatics*, **28**, i640–i646.
10. Corces, M.R., Granja, J.M., Shams, S., Louie, B.H., Seoane, J.A., Zhou, W., Silva, T.C., Groeneveld, C., Wong, C.K., Cho, S.W. *et al.* (2018) The chromatin accessibility landscape of primary human cancers. *Science*, **362**, eaav1898.
11. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, L. and Polosukhin, I. (2017) *Adv. Neural Inf. Process. Syst.*, 5998–6008.
12. Tao, Y., Cai, C., Cohen, W.W. and Lu, X. (2020) From genome to phenotype: predicting multiple cancer phenotypes based on somatic genomic alterations via the genomic impact transformer. *Pac. Symp. Biocomput.*, **25**, 79–90.
13. Cadow, J., Born, J., Manica, M., Oskooei, A. and Rodriguez Martinez, M. (2020) PaccMann: a web service for interpretable anticancer compound sensitivity prediction. *Nucleic Acids Res.*, **48**, W502–W508.
14. Marquart, K.F., Allam, A., Janjuha, S., Sintsova, A., Villiger, L., Frey, N., Krauthammer, M. and Schwank, G. (2021) Predicting base editing outcomes with an attention-based deep learning algorithm trained on high-throughput target library screens. *Nat. Commun.*, **12**, 5114.
15. Wu, X., Renuse, S., Sahasrabudhe, N.A., Zahari, M.S., Chaerkady, R., Kim, M.S., Nirujogi, R.S., Mohseni, M., Kumar, P., Raju, R. *et al.* (2014) Activation of diverse signalling pathways by oncogenic PIK3CA mutations. *Nat. Commun.*, **5**, 4961.
16. Bailey, T.L., Boden, M., Buske, F.A., Frith, M., Grant, C.E., Clementi, L., Ren, J., Li, W.W. and Noble, W.S. (2009) MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.*, **37**, W202–W208.
17. Weirauch, M.T., Yang, A., Albu, M., Cote, A.G., Montenegro-Montero, A., Drewe, P., Najafabadi, H.S., Lambert, S.A., Mann, I., Cook, K. *et al.* (2014) Determination and inference of eukaryotic transcription factor sequence specificity. *Cell*, **158**, 1431–1443.
18. Grant, C.E., Bailey, T.L. and Noble, W.S. (2011) FIMO: scanning for occurrences of a given motif. *Bioinformatics*, **27**, 1017–1018.
19. Zhu, L.J., Gazin, C., Lawson, N.D., Pages, H., Lin, S.M., Lapointe, D.S. and Green, M.R. (2010) ChIPpeakAnno: a bioconductor package to annotate chip-seq and chip-chip data. *BMC Bioinformatics*, **11**, 237.
20. Osmanbeyoglu, H.U., Pelossof, R., Bromberg, J.F. and Leslie, C.S. (2014) Linking signaling pathways to transcriptional programs in breast cancer. *Genome Res.*, **24**, 1869–1880.
21. Pelossof, R., Singh, I., Yang, J.L., Weirauch, M.T., Hughes, T.R. and Leslie, C.S. (2015) Affinity regression predicts the recognition code of nucleic acid-binding proteins. *Nat. Biotechnol.*, **33**, 1242–1249.
22. Osmanbeyoglu, H.U., Toska, E., Chan, C., Baselga, J. and Leslie, C.S. (2017) Pancancer modelling predicts the context-specific impact of somatic mutations on transcriptional programs. *Nat. Commun.*, **8**, 14249.
23. Ma, X., Somasundaram, A., Qi, Z., Hartman, D.J., Singh, H. and Osmanbeyoglu, H.U. (2021) SPaRTAN, a computational framework for linking cell-surface receptors to transcriptional regulators. *Nucleic Acids Res.*, **49**, 9633–9647.
24. Gala, K., Li, Q., Sinha, A., Razavi, P., Dorso, M., Sanchez-Vega, F., Chung, Y.R., Hendrickson, R., Hsieh, J.J., Berger, M. *et al.* (2018) KMT2C mediates the estrogen dependence of breast cancer through regulation of ERalpha enhancer function. *Oncogene*, **37**, 4692–4710.
25. Jozwik, K.M., Chernukhin, I., Serandour, A.A., Nagarajan, S. and Carroll, J.S. (2016) FOXA1 directs H3K4 monomethylation at

- enhancers via recruitment of the methyltransferase MLL3. *Cell Rep.*, **17**, 2715–2723.
26. Tarassishin, L. and Lee, S.C. (2013) Interferon regulatory factor 3 alters glioma inflammatory and invasive properties. *J. Neurooncol.*, **113**, 185–194.
27. Wang, M.D., Xing, H., Li, C., Liang, L., Wu, H., Xu, X.F., Sun, L.Y., Wu, M.C., Shen, F. and Yang, T. (2020) A novel role of Kruppel-like factor 8 as an apoptosis repressor in hepatocellular carcinoma. *Cancer Cell Int.*, **20**, 422.
28. Zhang, Q., Stovall, D.B., Inoue, K. and Sui, G. (2011) The oncogenic role of yin yang 1. *Crit. Rev. Oncog.*, **16**, 163–197.
29. Martinez, L.A. (2016) Mutant p53 and ETS2, a tale of reciprocity. *Front. Oncol.*, **6**, 35.