

1 **Interpretable deep learning for chromatin-informed inference of transcriptional** 2 **programs driven by somatic alterations across cancers**

3 Yifeng Tao¹⁺, Xiaojun Ma^{2,3+}, Drake Palmer³, Russell Schwartz^{1,4}, Xinghua Lu^{2,5}, Hatice Ulku
4 Osmanbeyoglu^{2,3,6*}

5
6 ¹Computational Biology Department, School of Computer Science, Carnegie Mellon University,
7 Pittsburgh, PA, USA

8 ²Department of Biomedical Informatics, School of Medicine, University of Pittsburgh, Pittsburgh,
9 PA, USA

10 ³UPMC Hillman Cancer Center, University of Pittsburgh, Pittsburgh, PA, USA

11 ⁴Department of Biological Sciences, Carnegie Mellon University, Pittsburgh, PA, USA

12 ⁵Department of Pharmaceutical Science, School of Medicine, University of Pittsburgh,
13 Pittsburgh, PA, USA

14 ⁶Department of Bioengineering, School of Engineering, University of Pittsburgh, Pittsburgh, PA,
15 USA

16
17 ⁺These authors contributed equally: Y.T., X.M.

18
19 ^{*} Correspondence to: Hatice Ulku Osmanbeyoglu (osmanbeyoglu@pitt.edu)
20 ORCID ID: 0000-0002-4972-4347

21 **Abstract**

22
23 Cancer is a disease of gene dysregulation, where cells acquire somatic and epigenetic alterations
24 that drive aberrant cellular signaling. These alterations adversely impact transcriptional programs
25 and cause profound changes in gene expression. Interpreting somatic alterations within context-
26 specific transcriptional programs will facilitate personalized therapeutic decisions but is a
27 monumental task. Toward this goal, we develop a partially interpretable neural network model
28 called **Chromatin-informed Inference of Transcriptional Regulators Using Self-attention**
29 **mechanism (CITRUS)**. CITRUS models the impact of somatic alterations on transcription factors
30 and downstream transcriptional programs. Our approach employs a self-attention mechanism to
31 model the contextual impact of somatic alterations. Furthermore, CITRUS uses a layer of hidden
32 nodes to explicitly represent the state of transcription factors (TFs) to learn the relationships
33 between TFs and their target genes based on TF binding motifs in the open chromatin regions of
34 tumor samples. We apply CITRUS to genomic, transcriptomic, and epigenomic data from 17
35 cancer types profiled by The Cancer Genome Atlas. CITRUS predicts patient-specific TF activities
36 and reveals transcriptional program variations between and within tumor types. We show that
37 CITRUS yields biological insights into delineating TFs associated with somatic alterations in
38 individual tumors. Thus, CITRUS is a promising tool for precision oncology.

39 **Introduction**

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

49
50 In the last decade, a monumental effort has been made to molecularly profile tumors by consortia,
51 including The Cancer Genome Atlas (TCGA) and the International Cancer Genome Consortium
52 (1,2). The multimodal datasets generated by these efforts include gene expression and somatic
53 alterations, such as recurrent mutations and copy number variations (CNVs). The combination of
54 genomic and transcriptomic information enables the integration of transcriptional states with
55 upstream signaling pathways. Several methods have been developed to connect somatic
56 alterations to a prior network or to gene expression (3-9). More recently, the Genomic Data
57 Analysis Network generated assay for transposase-accessible chromatin with high-throughput
58 sequencing (ATAC-seq) data for a subset of TCGA samples (~500 patients) (10). Although
59 chromatin profiling helps uncover context-dependent and/or non-linear effects of transcription
60 factors (TFs) on gene expression, it has not yet been incorporated into methods that connect
61 somatic alterations to transcriptional programs across cancers. Incorporating DNA sequence
62 information at promoter, intronic, and intergenic enhancers from ATAC-seq tumor profiles using
63 TF motif analysis will improve the modeling of transcriptional regulation and delineate the impact
64 of somatic alterations on transcriptional programs.

65
66 Deep learning is a powerful tool for capturing non-linear feature interactions that can explain the
67 underlying biological phenomena. For example, attention mechanism is a deep learning method
68 that has been widely used in computer vision and natural language processing. In contrast to
69 traditional deep learning methods, the self-attention mechanism considers the contextual
70 relationship of the input features and assigns attention weights to each input (11). In general,
71 attention mechanisms improve the performance of deep learning models and increase the
72 interpretability of the models. More recently, attention mechanisms have been applied to cancer
73 genomics for cancer driver gene detection (12), drug response prediction (13), and base editing
74 outcome prediction (14). For example, the genomic impact transformer (GIT) model utilizes a self-
75 attention mechanism to encode the effects of somatic alterations in cancer and uses multi-layer
76 perceptrons to predict differentially expressed genes (12). The attention mechanism enables GIT
77 to select driver mutations that are likely to lead to downstream phenotypes. However, the GIT
78 model lacks interpretability in the sense that it does not model intermediate TFs during modeling
79 signaling from somatic alterations to gene expression programs.

80
81 Here, we present **Chromatin-informed Inference of Transcriptional Regulators Using Self-**
82 **attention mechanism (CITRUS)**, a partially interpretable neural network model with encoder-
83 decoder architecture. CITRUS links somatic alterations to transcriptional programs by modeling
84 the statistical relationships between mutations, CNVs, gene expression, and TF-target gene
85 information derived from ATAC-seq (**Fig. 1**). We show that CITRUS yields important biological
86 insights into dysregulated TFs in individual tumors. Using a systematic *in silico* knock out
87 approach, we identified key TFs associated with major somatic alterations. We believe CITRUS
88 will assist researchers in providing actionable hypotheses for follow-up experiments and
89 developing personalized and targeted therapeutics in a pan-cancer setting.

90 **Material & Methods**

91 **Data pre-processing**

92 We downloaded the batch normalized RNA-Seq expression levels quantified by RNA-Seq by
93 Expectation Maximization (RSEM) from the Genomic Data Commons (GDC) portal
94 (<https://gdc.cancer.gov/about-data/publications/pancanatlas>). We log₂-transformed RSEM values
95 and identified the 2,500 most variable genes across samples within a cancer type. Then, we took
96 the union of the identified genes across cancer types. The final gene set included 5541 genes.

97
98 We obtained processed gene-level somatic alterations for each cancer patient from Cai et al. (4).
99 Genes with non-synonymous mutations, small insert/deletion, or somatic copy number alteration
100 (deletion or amplification) were given a value of 1, and otherwise were given a value of 0. We
101 removed genes that were not present in at least 4% of samples for each cancer type.

102
103 We downloaded the ATAC-seq pan-cancer dataset from the GDC portal
104 (<https://gdc.cancer.gov/about-data/publications/ATACseq-AWG>) (10). Using the MEME (15)
105 curated Cis-BP (16) TF-binding motif reference, we scanned the pan-cancer ATAC-seq peak atlas
106 with FIMO (17) to find peaks likely to contain each motif ($P < 10^{-5}$). The final set contained 320
107 motifs. We associated each peak with its nearest gene in the human genome using the
108 ChIPpeakAnno package (18). ATAC-seq peaks located in the body of the transcription unit,
109 100 kb upstream of the transcription start site (TSS), and 100 kb downstream of the 3' end were
110 assigned to the associated gene. TF-binding site identification was used to convert the assigned
111 ATAC peaks for each gene to a feature vector of binding signals by assigning the maximum score
112 of each motif across all peaks to a gene. Then, we created a matrix $\mathbf{C} \in \{0,1\}^{k \times l}$ that defines a
113 candidate set of associations between TFs and target genes. $C_{i,j} = 1$ when there is a connection
114 from TF j to the gene/RNA i (red lines connecting the TF layer and target gene expression (Exp)
115 layer in **Fig. 1**).

116 **CITRUS model**

117
118 CITRUS is a framework for modeling impact of somatic alterations on transcriptional programs.
119 **Fig. 1** shows the model architecture with an overall encoder and decoder structure. Somatic gene
120 alteration inputs with more than 20K dimensions were encoded into a compressed representation
121 as tumor embedding and then decode to a large dimension data of gene expression. This allows
122 the model to capture key features of the high dimension inputs and reduce the data noise as well.

123
124 We designed a self-attention mechanism which assigned importance weights to input features
125 (somatic alterations) through the model training. Formally, given a specific tumor t , with the
126 cancer type s , we have a set of somatic alterations in the tumor $\{g_u\}_{u=1}^m$ where m is number of
127 mutant genes. The encoder module first maps each gene g (it is g_u here, but we omit the
128 subscript for notation simplicity) into its corresponding gene vector e_g . Then, the encoder utilizes
129 the multi-head self-attention mechanism to calculate the weighted sum of both the gene
130 embeddings and the cancer type embedding:

$$131 \quad \quad \quad 132 \quad \quad \quad 133 \quad \quad \quad e_t = e_s + \alpha_1 e_1 + \alpha_2 e_2 + \alpha_3 e_3 + \dots + \alpha_m e_m$$

134 The self-attention mechanism takes the gene embeddings of all mutated/changed genes as an
135 input and outputs the attention weights $\{\alpha_u\}_{u=1}^m$ through a sub-neural network. The attention
136 mechanism captures the context of co-existing somatic alterations and their complex interactions,
137 which is lost in simpler models. Interested readers can find the mathematical details of self-
138 attention mechanisms in the cited reference (12).

139

140 The decoder first infers the TF activities from the encoded tumor embedding e_t :

$$141 \quad e_f = \tanh(W_f e_t + b_f).$$

142
143 We used tanh activation instead of ReLU operation, which is more widely used in deep learning,
144 because it has similar performance to that of ReLU in our model and generates more biologically
145 meaningful results (e.g., distribution of TFs e_f). Finally, CITRUS predicts cancer type-specific
146 mRNA expression from TF activities:

$$147 \quad \hat{y} = W e_f + b_r$$

148
149 where W corresponds to the sparse TF-target gene matrix constrained by the prior $C \in \{0,1\}^{k \times l}$.
150 More specifically, to integrate priors into our model, W shares the same shape with prior C , and
151 $W_{i,j}$ is allowed to be nonzero only when $C_{i,j} = 1$, and $W_{i,j}$ is constrained to be non-negative value.
152 We use mean square loss function as: $MSE(y, \hat{y})$.

153
154 One might use other common approaches to integrate prior C into the W , i.e., by applying a
155 Gaussian prior to W , which is equivalent to adding an additional penalty to the loss function
156 $\sum_{i,j:C_{i,j}=0} (W)_{i,j}^2$. However, this “soft” constraint tends to generate less stable TF layers across
157 different runs of training compared to the “hard” constraints shown in our model.

158
159 To prevent overfitting and to increase robustness to noise, we introduced additional dropout
160 operations with a dropout rate of 0.2 after the input layer, activated tumor embedding layer, and
161 activated TF layer.

162 **Training and evaluation**

163 We implemented CITRUS through the PyTorch package (<https://pytorch.org/>), and training was
164 performed using the Adam optimizer with default parameters except for the learning rate¹⁵ and
165 weight decay. We set the learning rate to 1×10^{-3} and the weight decay to 1×10^{-5} . We used
166 early stopping with patience of 30 steps to stop training.

167
168 For statistical evaluation, we computed the mean Spearman correlation (ρ) between predicted
169 and measured gene expression profiles for each tumor type. We split datasets into training (40%),
170 validation (20%), and testing (20%) sets. For CITRUS, we utilized the training and validation sets
171 to tune hyperparameters, such as the learning rate and training steps, and then evaluated these
172 parameters on the testing set. For affinity regression (see below), we separated datasets by
173 cancer type and conducted 5-fold cross-validation to tune hyperparameters in the training and
174 validation sets. Then, we applied the trained model with selected hyperparameters to the testing
175 set for performance evaluation. To increase the stability of inferred TF activity analysis, we
176 assembled multiple CITRUS models trained with different random initialization state and
177 integrated the TF layer based on the average of 10 trials.

178
179 **Parameter selection:** CITRUS includes more than 10 hyperparameters that are described in the
180 following paragraphs. These hyperparameters were tuned for optimal performance in the
181 validation set. Ideally, hyperparameter optimization is performed using a grid search of all
182 parameters. However, this is not practical due to the tremendous computational cost. For
183 example, three options for each parameter leads to 3^{10} possible combinations for just 10
184 parameters. In addition, we guide the performance metric by k-fold cross-validation, and the total
185 experiments necessary would be 5×3^{10} ($k=5$). Therefore, our hyperparameter tuning strategy
186 combined automatic and manual tuning. First, we created empirical settings for each parameter
187 and randomly selected a set of parameters from 100 combinations. We utilized the best-

189 performing settings to narrow down the preliminary decisions and correlation among parameters.
190 Then, we tuned parameters independently or in sub-groups manually or by grid search.

191
192 **Model robustness:** The learning rate is perhaps the most important hyperparameter in neural
193 network training. We first tested the learning rate in a range of settings [10^{-5} , 10^{-4} , 10^{-3} , 10^{-2} ...],
194 starting with the lowest setting and progressing to larger values until validation loss started to
195 diverge. We found that if the learning rate was too small, overfitting occurred and picked up input
196 noise. Additionally, overfitting reduced the number of driver genes that were covered in
197 downstream attention weight analyses. If the learning rate was too big, however, the model could
198 not converge to an optima and yielded higher validation loss. Ultimately, we selected learning
199 rates of 10^{-3} and 10^{-4} and applied a weight penalty (weight decay) to find an optimal combination
200 of settings. We set the weight decay range from 10^{-6} to 10^{-4} and performed a grid search. The
201 optimal settings for learning rate and weight decay were determined to be 10^{-3} and 10^{-5} ,
202 respectively. Although large batch sizes can accelerate learning rates and training, our
203 experiments indicated that a learning rate of 10^{-3} was the largest value that maintained validation
204 accuracy when tested on increasing batch sizes (16, 64, 100, and 300, which is the maximum
205 value that could run in GPU). We found that larger batch sizes tended to have slightly higher
206 gene-wise correlation at the cost of longer training time. To balance execution time, we selected
207 a batch size of 100. The early stopping patience setting is also related to the learning rate and
208 batch size. Specifically, higher learning rates and larger batch sizes require smaller patience to
209 stop training. Higher patience settings may otherwise cause overfitting. Using our selected
210 learning rate and batch size settings, a patience of 30 was generally sufficient to maintain training
211 without stopping too early (underfitting) due to fluctuation and without halting too far from the
212 optima (overfitting). We validated a patience setting of 30 by comparing it with a case of overfitting.
213 We selected the lowest loss point in the overfit training and measured how far it was from the
214 model with early stopping. During early stages of training, the model showed an initial drop in
215 validation performance followed by a rise. To avoid this inconsistency, we did not apply early
216 stopping for the first 180 test steps. To test the attention mechanism, we created a mesh grid for
217 two attention sizes (256, 128) and four attention head settings (32, 16, 8, 4). We then performed
218 an exhaustive grid search within these settings. Based on prediction performance, we selected
219 256 and eight as the optimal values for attention size and attention head, respectively.

220
221 Finally, we fine-tuned our model by adjusting the dropout rate. Because we used weight decay
222 for regularization, dropout is considered a secondary regularization for our model. In addition to
223 hidden layer dropout, we also applied dropout to our input to reduce input noise and network
224 redundancy and to generate a more stable hidden TF layer. We tested a sequence of five dropout
225 rates (0.1, 0.2, 0.3, 0.4, 0.5). All dropout rate settings yielded performances above 0.9 for average
226 sample correlation in the testing set. We determined the dropout rate optimal value (0.2) primarily
227 based on driver gene coverage in self-attention analyses.

228
229 As we used an early stopping mechanism, we set the maximum iteration parameter to 1000. This
230 setting ensures that the training process stops either once the patience setting is satisfied or once
231 the maximum iterations is reached. Code testing and quick runs were performed with a maximum
232 iteration of one.

233
234 We tested two activation functions: 'ReLU' and 'tanh'. Although both activation functions
235 performed similarly, 'tanh' generated more biologically meaningful results and was selected. We
236 also tested l2, minimax, and standard normalization (scale) to normalize gene expression and
237 found that scale normalization generated the best prediction accuracy for our model settings.

238
239

240 **Training the affinity regression (AR) models**

241 AR is an algorithm for efficiently solving a regularized bilinear regression problem (19,20) and
242 was defined in our model as follows. For a data set of M tumor samples profiled using RNA-seq
243 with N genes, we let $\mathbf{Y} \in \mathbb{R}^{N \times M}$ be the log10 gene expression profiles of tumor samples. Each
244 column of \mathbf{Y} corresponds to an RNA-seq experiment for a cancer type. We define the TF attributes
245 of each gene in a matrix $\mathbf{D} \in \mathbb{R}^{N \times Q}$, where each row represents a gene, and each column
246 represents a TF vector. The TF vector indicates whether there is a binding site for the TF on each
247 gene based on ATAC-seq data. We define the somatic alteration attributes of tumor samples as
248 a matrix $\mathbf{P} \in \mathbb{R}^{M \times S}$ where each row represents a tumor sample, and each column represents the
249 somatic alteration status for the tumor sample. We set up a bilinear regression problem to learn
250 the weight matrix $\mathbf{W} \in \mathbb{R}^{Q \times S}$ on paired TF and somatic alteration features:

$$251 \quad \mathbf{DWP}^T \sim \mathbf{Y}$$

252
253 We can transform the system to an equivalent system of equations by reformulating the matrix
254 products as Kronecker products:

$$255 \quad \mathbf{DWP}^T \approx \mathbf{Y} \Leftrightarrow (\mathbf{P} \otimes \mathbf{D}) \text{vec}(\mathbf{W}) \approx \text{vec}(\mathbf{Y})$$

256 where \otimes is a Kronecker product, and vec is a vectorizing operator that stacks a matrix and
257 produces a vector. The result of this system is a standard (if large-scale) regression problem. Full
258 details and a derivation of the reduced optimization problem are provided elsewhere (20).

259 ***In silico* knockout analysis**

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

270 **Statistical analysis**

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

275
276
277 Association score between TF activity subtypes and frequent somatic alterations. For each
278 somatic mutation or copy number variation, we calculated the P -value of its frequency in a cancer
279 subtype compared to other subtypes using Fisher's exact test. The P -value was further adjusted
280 through FDR across subtypes. To identify the relative frequency of a somatic alteration in a
281 subtype, we defined an association score, which is the product of the relative frequency direction
282 and $-\log_{10}\text{FDR}$.

283
284
285

286 Results

287 Pan-cancer modeling of transcriptional programs

288 To systematically interpret somatic alterations within context-specific transcriptional programs
289 and to identify disrupted TFs that drive tumor-specific gene expression patterns across multiple
290 cancer types, we developed CITRUS (**Fig. 1**). CITRUS traces biological signaling from somatic
291 alterations to signaling pathways, to TFs, and finally to target gene expression (mRNA levels). To
292 enable this tracing, CITRUS employs an encoder-decoder architecture (**Fig. 1**). The encoder
293 module compresses input somatic alterations into a latent vector variable called a tumor
294 embedding. The decoder predicts TF activities from the tumor embedding and then predicts target
295 gene expression. We used sparse TF-target gene priors based on tumor ATAC-seq data. Briefly,
296 we started with an atlas of chromatin accessible genomic locations derived from the tumor types
297 to be analyzed using ATAC-seq profiling data (see Methods). We then represented every gene
298 by its feature vector of TF-binding scores, where motif information was summarized across all
299 promoter, intronic, and intergenic chromatin accessible sites assigned to the gene (see Methods).

300
301 We applied this approach to 17 tumors from TCGA and identified key TFs associated with somatic
302 alterations. Our dataset included samples from 17 different tumor types for which mRNA, somatic
303 mutation, copy number variation, and ATAC-seq data were available: bladder urothelial
304 carcinoma (BLCA, n=371), breast cancer (BRCA, n=719), cervical squamous cell carcinoma and
305 endocervical adenocarcinoma (CESC, n=267), colorectal adenocarcinoma (COAD, n=271),
306 esophageal carcinoma (ESCA, n=170), glioblastoma multiforme (GBM, n=143), head and neck
307 squamous carcinoma (HNSC, n=475), kidney renal cell-clear carcinoma (KIRC, n=357), kidney
308 renal papillary cell carcinoma (KIRP, n=272), liver hepatocellular carcinoma (LIHC, n=336), lung
309 adenocarcinoma (LUAD, n=459), lung squamous cell carcinoma (LUSC, n=430),
310 pheochromocytoma and paraganglioma (PCPG, n=109), prostate cancer (PRAD, n=449),
311 stomach adenocarcinoma (STAD, n=373), thyroid carcinoma (THCA, n=216), and uterine corpus
312 endometrial carcinoma (UCEC, n=361).

313
314 For statistical evaluation, we computed the mean Spearman correlation between predicted and
315 measured gene expression profiles on the testing set (see Methods). CITRUS achieved
316 significantly better performance than a regularized bilinear regression algorithm called affinity
317 regression (AR) (20-22) that was trained independently for each cancer type. and explain gene
318 expression across tumors in terms of somatic alteration status and presence of TF binding sites
319 based on a pan-cancer ATAC-seq atlas (**Fig. 2A**).

320
321 To identify somatic alterations that influenced gene expression programs, we compared the
322 relationship of overall attention weights (inferred by CITRUS) and the frequencies of somatic
323 alterations (used as the control group) across all cancer types and within each cancer type (**Fig.**
324 **2B and Supplementary Fig. 1**). In general, attention weights were positively correlated with the
325 frequency of somatic alteration. For example, the top altered genes *TP53* and *PIK3CA* had high
326 attention weights. However, our self-attention mechanism assigned low attention weights to many
327 frequently altered genes, indicating that these genes may be cancer passengers. Indeed, we
328 found genes with high attention weights were enriched for known cancer drivers using the
329 IntOGen⁹ database. We first grouped all the genes into two parts with the threshold of 2
330 ($\log(\text{attention}+1) \geq 2$ as the more attended group, and $\log(\text{attention}+1) < 2$ as the less attended
331 group). Using Fisher's exact test, we verified that known cancer driver genes were enriched in
332 the highly attended group ($P = 4.48 \times 10^{-41}$) in the pan-cancer analysis. We also observed a few
333 infrequently altered genes with high attention weights. For example, the H3K4 methyltransferase
334 *KMT2C* had a high attention weight in BRCA but was infrequently altered. Indeed, *KMT2C* is a
335 key regulator of ER α activity and anti-estrogen response in breast cancer (23,24).

336 We used CITRUS to infer patient-specific TF activities across tumor types. Clustering tumors by
337 these inferred TF activities largely recovered the distinction between major tumor types (**Fig. 2C**).
338 Interestingly, samples with squamous morphology components (BLCA, CESC, ESCA, HNSC,
339 and LUSC) clustered together. Tumors with tissue or organ similarities or proximity were also
340 clustered together. These included neuroendocrine and glioma tumors (GBM and PCPG), clear
341 cell and papillary renal carcinomas (KIRC and KIRP), a gastrointestinal group (COAD, and
342 STAD), and breast and endometrial cancer (BRCA and UCEC). We also observed similar
343 clustering of the tumor embeddings (**Supplementary Fig. 2**).
344

345 Next, we assessed TF-tumor type associations by t-test and compared inferred TF activities
346 between samples in each tumor type versus those in all other tumor types. We corrected for false
347 discovery rate (FDR) across TFs and identified significant shared and cancer-specific TFs, which
348 are listed in **Supplementary Data 1**. The average TF activity and significance of the four most
349 significant TFs in each cancer are shown in **Fig. 3**. Our results highlight both known and novel
350 cancer-specific TF regulators. For example, FUBP1, which regulates *c-Myc* gene transcription,
351 had significantly higher inferred activity in many cancer types, including LIHC, HNSC, BLCA,
352 ESCA, CESC, LUSC, PRAD, BRCA, and UCEC. Consistent with previous reports, IRF3 activity
353 was significantly higher in GBM(25). KLF8 had decreased activity in GBM, LIHC, and KIRC, which
354 is consistent with its role in suppressing cell apoptosis during tumor progression (26). Additionally,
355 YY1, which regulates various developmental processes (27), had increased activity in CESC and
356 COAD.
357

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

359 Next, we asked whether CITRUS could identify cancer subtypes based on the TF activity
360 associated with somatic alterations. We conducted *k*-means clustering of inferred TF activities for
361 each cancer type to define subtypes, and then we conducted hierarchical clustering of both the
362 cancer subtypes and TF activities. **Fig. 4** shows the clustering of subtypes by CITRUS-inferred
363 mean TF activities and corresponding somatic alteration associations (see Methods). We
364 observed major differences in mean TF activities across cancer types and minor but significant
365 differences within cancer types. Variations within a cancer type may arise from distinct mutation
366 or CNV profiles of subgroups. For example, clustering by TF activities revealed subclasses of
367 CESC enriched with *KRAS*; KIRC enriched with *VHL*, *BAP1*, *PBRM1*, and *TP53*; LIHC enriched
368 with *CTNNB1*, *BAP1*, and *TP53*; THCA enriched with *NRAS*, *HRAS*, and *BRAF*; and PCPG
369 enriched with *HRAS*.
370

371 As our goal was to decipher cancer-specific downstream effects of targeted therapies and to
372 discover secondary targets for combination drug strategies, we developed a systematic statistical
373 approach for modeling the impact of somatic alterations on TF activity. We implemented an *in*
374 *silico* knock out approach that removes a specific somatic mutation (or CNV) *g* from all carrier
375 tumor samples in each TCGA cancer study and then predicts altered TF activity (see Methods).
376 Using this approach, we were able to identify TFs whose inferred activity was significantly
377 dysregulated by somatic alterations in known cancer driver genes. **Fig. 5A** demonstrates TF
378 activities that were associated with somatic alterations in UCEC. CITRUS identified mutations in
379 *PIK3CA*, *PTEN*, *KRAS*, *TP53*, and *CTNNB1* that were significantly associated with various TF
380 activities across UCEC tumors (~66% of tumors have *PTEN* inactivating mutations, ~50% have
381 *PIK3CA* activating mutations, ~38% have *TP53* mutations, ~26% have *CTNNB1* mutations, and
382 ~20% have *KRAS* mutations). UCEC samples with *PTEN* mutations were mutually exclusive with
383 *TP53*, *CTNNB1*, and *KRAS* mutations and showed distinct TF activity patterns. Mutations in
384 *PTEN* that inactivate its phosphatase activity result in increased PI3K signaling. Consistent with
385 this effect, TFs associated with *PTEN* mutations were involved in cell cycle and differentiation,
386 including E2F5, TP63, ELF3, DBP, ZKSCAN3, LHX2, HOXB6, SOX9, DBP, MYLB1, and GLIS1.

387 TFs associated with *CTNNB1* mutant status were involved in WNT and TGF-beta signaling
388 including TCF7, TCF7L2, TCF7L1, FOXH1, EMX1, and MYBL1.

389
390 Similarly, CITRUS identified TF activities that were associated with somatic alterations in BRCA
391 (**Fig. 5B**). Mutations in *PIK3CA*, *PTEN*, *MAP2K4*, *GATA3*, *TP53*, and *CDH1* were significantly
392 associated with various TF activities. In BRCA, ~36% of tumors have *PIK3CA* activating
393 mutations, ~35% have *TP53* mutations, ~15% have *GATA3* mutations, ~15% have *CDH1*
394 mutations, ~10% have *PTEN* mutations, and ~7% have *MAP2K4* mutations. Activating mutations
395 in *PIK3CA* often occur in one of three hotspot locations (E545K, E542K, and H1047R) and
396 promote constitutive signaling through the pathway. TFs associated with *PIK3CA* mutations were
397 involved in WNT signaling, epithelial–mesenchymal transition, and cancer stem cell transition,
398 including ELF3, TFEC, STAT4, STAT5B, NFATC1, GLIS1, CDC5L, and AR. BRCA samples with
399 *PIK3CA* and *TP53* mutations were mutually exclusive, and our *in silico* knock out analysis
400 associated distinct TFs with these mutations. *TP53* mutant tumors were associated with increased
401 activity of TFs that have roles in tumor growth, such as ETS2 and FOSB, growth modulation, such
402 as THAP1, CREB3L1, and CEBPZ, and development, such as MEF2C/D, MEOX1, and MSX1.
403 We performed similar analyses for other cancer types (**Supplementary Fig. 3**).

404
405 Although the TFs affected by some somatic alterations differed between cancer types, mutation
406 of *TP53* was associated with similar TFs across cancer types (**Supplementary Fig. 4**). *TP53* is
407 one of the most frequently inactivated tumor suppressor genes that suffers from missense
408 mutations in human cancer. These missense mutations result in the expression of a mutant form
409 of p53 protein. Mutant p53 protein can disable other tumor suppressors (e.g., p63 and p73) or
410 enable oncogenes, such as ETS2 (28). Indeed, the inferred TF activity of ETS2 was higher in
411 mutant versus WT *TP53* tumors across cancers (**Fig. 5C**); however, these differences were not
412 as significant at the gene expression level (**Supplementary Fig. 5**).

413 414 **Discussion**

415 Analysis of the regulatory network in tumor datasets is challenging due to the complexity of the
416 cancer genome (e.g., aneuploidy, CNVs, structural variation, and mutations). CITRUS provides a
417 systematic framework for integrating regulatory genomics with tumor expression and somatic
418 alterations to better understand how expression programs are affected by somatic alterations in
419 cancers and to infer patient-specific TF activities. Our method uses a deep learning framework
420 called a self-attention mechanism to capture the complex contextual interactions between somatic
421 alterations. For a more accurate representation of TF-target gene relationships, we leveraged
422 ATAC-seq tumor data from TCGA patients. CITRUS is designed to capture the flow of information
423 from altered genes (e.g., signaling proteins) to TFs to target genes, and our *in silico* knock out
424 analysis predicts the causal impacts of somatic alterations. Joint modeling across different tumor
425 types also revealed patient subgroups associated with somatic alterations. In cases where a
426 somatic alteration is associated with the activity of a targetable TF or their upstream/downstream
427 component, it may be possible to identify combination therapies using CITRUS.

428
429 One limitation of the TF binding motif approach utilized by CITRUS is that TFs of the same family
430 often share a similar motif and thus are difficult to disambiguate. Therefore, TF motifs may
431 encompass the activities of multiple TFs. Moreover, co-binding TF binding patterns (e.g., AP-
432 1–IRF complexes) can be biologically meaningful for gene expression and are not currently
433 represented in our model. Future models will work to represent these composite elements as
434 features. Another limitation is that we do not represent directionality in the TF-target gene priors
435 (i.e., whether a gene is activated or repressed by a TF). Prior knowledge of whether the TF is
436 acting as an activator or as a repressor would add meaningful interpretation to inferred TF
437 activities. These limitations may confound the interpretation of the activity of TFs with context-

438 specific activator and repressor roles. Further, regulatory network analysis of tumor datasets is
439 also complicated by the presence of stromal/immune cells within the tumor and the heterogeneity
440 of the cancer cells themselves. However, our framework can be extended to model single-cell
441 RNA-seq or deconvoluted RNA-seq via computational methods.

442
443 Despite these limitations, modeling the impact of somatic alterations on transcriptional programs
444 may ultimately enable the development of individualized therapies, aid in understanding
445 mechanisms of drug resistance, and allow the identification of biomarkers of response. We
446 anticipate that computational modeling of transcriptional regulation across different tumor types
447 will emerge as an important tool in precision oncology, aiding in the eventual goal of selecting the
448 best therapeutic option for individual patients.

449 **Data availability**

451 ATAC-seq data are available in the public repository Genomic Data Commons
452 (<https://gdc.cancer.gov/about-data/publications/ATACseq-AWG>). RNA-seq gene expression,
453 somatic mutation, copy number variation, and clinical data are available in a public repository
454 from TCGA's Firehose data run ([https://confluence.broadinstitute.org/display/GDAC/Dashboard-](https://confluence.broadinstitute.org/display/GDAC/Dashboard-Stddata)
455 [Stddata](https://confluence.broadinstitute.org/display/GDAC/Dashboard-Stddata)). Only the samples 'whitelisted' by TCGA for the Pan-Cancer Analysis Working Group
456 were used in the study. For our analysis, we only used samples with parallel RNA-seq, somatic
457 mutation, and GISTIC copy number data. Processed input and output files have been made
458 available at the supplementary website for the paper: <http://www.pitt.edu/~xim33/CITRUS>.

459 **Code availability**

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

462 **Funding**

464 This study was funded by support through the National Institutes of Health (R00 CA207871 to
465 H.U.O.); the Fellowship in Digital Health from the Center for Machine Learning and Health at
466 Carnegie Mellon University (to Y.T.); the UPMC-ITTC fund (to H.S.); the National Institutes of
467 Health (R01HG010589 and R21CA216452 to R.S.); the Pennsylvania Department of Health
468 (FP00003273 to R.S.); the Mario Lemieux Foundation (to R.S.); the AWS Machine Learning
469 Research Award (to R.S.). The content of this manuscript is solely the responsibility of the authors
470 and does not necessarily represent the official views of the National Institutes of Health or other
471 funding agencies. The Pennsylvania Department of Health specifically disclaims responsibility for
472 any analyses, interpretations, or conclusions. Funding for open access charge: National Institutes
473 of Health.

474 **Acknowledgements**

476 The results published here are based on data generated by The Cancer Genome Atlas project
477 established by the NCI and NHGRI (accession number: phs000178.v7p6). Information about
478 TCGA and the investigators and institutions that constitute TCGA research network can be found
479 at <http://cancergenome.nih.gov/>. We thank Jacob Stewart-Ornstein for helpful discussions

480 **References**

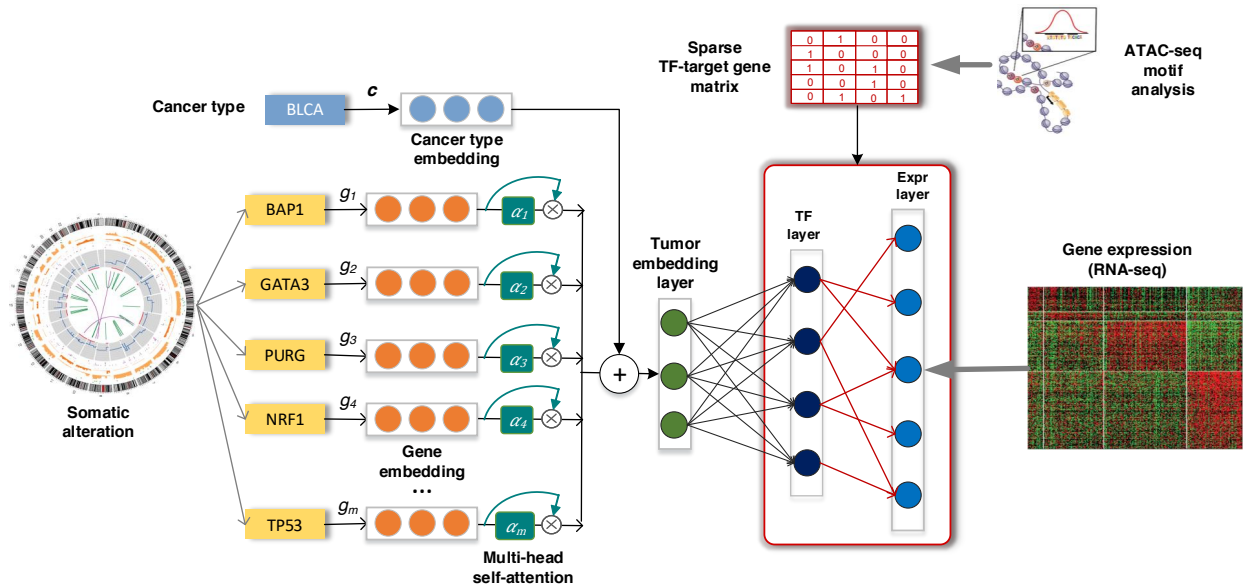
- 482 1. Consortium, I.T.P.-C.A.o.W.G. (2020) Pan-cancer analysis of whole genomes. *Nature*,
483 **578**, 82-93.
- 484 2. Cancer Genome Atlas Research, N., Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw,
485 K.R., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C. and Stuart, J.M. (2013)
486 The Cancer Genome Atlas Pan-Cancer analysis project. *Nat Genet*, **45**, 1113-1120.

- 487 3. Wang, Z., Ng, K.S., Chen, T., Kim, T.B., Wang, F., Shaw, K., Scott, K.L., Meric-
488 Bernstam, F., Mills, G.B. and Chen, K. (2018) Cancer driver mutation prediction through
489 Bayesian integration of multi-omic data. *PLoS One*, **13**, e0196939.
- 490 4. Cai, C., Cooper, G.F., Lu, K.N., Ma, X., Xu, S., Zhao, Z., Chen, X., Xue, Y., Lee, A.V.,
491 Clark, N. *et al.* (2019) Systematic discovery of the functional impact of somatic genome
492 alterations in individual tumors through tumor-specific causal inference. *PLoS Comput*
493 *Biol*, **15**, e1007088.
- 494 5. Hofree, M., Shen, J.P., Carter, H., Gross, A. and Ideker, T. (2013) Network-based
495 stratification of tumor mutations. *Nat Methods*, **10**, 1108-1115.
- 496 6. Paull, E.O., Carlin, D.E., Niepel, M., Sorger, P.K., Haussler, D. and Stuart, J.M. (2013)
497 Discovering causal pathways linking genomic events to transcriptional states using Tied
498 Diffusion Through Interacting Events (TieDIE). *Bioinformatics*, **29**, 2757-2764.
- 499 7. Basha, O., Mauer, O., Simonovsky, E., Shpringer, R. and Yeger-Lotem, E. (2019)
500 ResponseNet v.3: revealing signaling and regulatory pathways connecting your proteins
501 and genes across human tissues. *Nucleic Acids Res*, **47**, W242-W247.
- 502 8. Bashashati, A., Haffari, G., Ding, J., Ha, G., Lui, K., Rosner, J., Huntsman, D.G., Caldas,
503 C., Aparicio, S.A. and Shah, S.P. (2012) DriverNet: uncovering the impact of somatic
504 driver mutations on transcriptional networks in cancer. *Genome Biol*, **13**, R124.
- 505 9. Ng, S., Collisson, E.A., Sokolov, A., Goldstein, T., Gonzalez-Perez, A., Lopez-Bigas, N.,
506 Benz, C., Haussler, D. and Stuart, J.M. (2012) PARADIGM-SHIFT predicts the function
507 of mutations in multiple cancers using pathway impact analysis. *Bioinformatics*, **28**, i640-
508 i646.
- 509 10. Corces, M.R., Granja, J.M., Shams, S., Louie, B.H., Seoane, J.A., Zhou, W., Silva, T.C.,
510 Groeneveld, C., Wong, C.K., Cho, S.W. *et al.* (2018) The chromatin accessibility
511 landscape of primary human cancers. *Science*, **362**.
- 512 11. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.
513 and Polosukhin, I. (2017), *Advances in neural information processing systems*, pp. 5998-
514 6008.
- 515 12. Tao, Y., Cai, C., Cohen, W.W. and Lu, X. (2020) From genome to phenome: Predicting
516 multiple cancer phenotypes based on somatic genomic alterations via the genomic
517 impact transformer. *Pac Symp Biocomput*, **25**, 79-90.
- 518 13. Cadow, J., Born, J., Manica, M., Oskooei, A. and Rodriguez Martinez, M. (2020)
519 PaccMann: a web service for interpretable anticancer compound sensitivity prediction.
520 *Nucleic Acids Res*, **48**, W502-W508.
- 521 14. Marquart, K.F., Allam, A., Janjuha, S., Sintsova, A., Villiger, L., Frey, N., Krauthammer,
522 M. and Schwank, G. (2021) Predicting base editing outcomes with an attention-based
523 deep learning algorithm trained on high-throughput target library screens. *Nat Commun*,
524 **12**, 5114.
- 525 15. Bailey, T.L., Boden, M., Buske, F.A., Frith, M., Grant, C.E., Clementi, L., Ren, J., Li,
526 W.W. and Noble, W.S. (2009) MEME SUITE: tools for motif discovery and searching.
527 *Nucleic Acids Res*, **37**, W202-208.
- 528 16. Weirauch, M.T., Yang, A., Albu, M., Cote, A.G., Montenegro-Montero, A., Drewe, P.,
529 Najafabadi, H.S., Lambert, S.A., Mann, I., Cook, K. *et al.* (2014) Determination and
530 inference of eukaryotic transcription factor sequence specificity. *Cell*, **158**, 1431-1443.
- 531 17. Grant, C.E., Bailey, T.L. and Noble, W.S. (2011) FIMO: scanning for occurrences of a
532 given motif. *Bioinformatics*, **27**, 1017-1018.
- 533 18. Zhu, L.J., Gazin, C., Lawson, N.D., Pages, H., Lin, S.M., Lapointe, D.S. and Green, M.R.
534 (2010) ChIPpeakAnno: a Bioconductor package to annotate ChIP-seq and ChIP-chip
535 data. *BMC Bioinformatics*, **11**, 237.

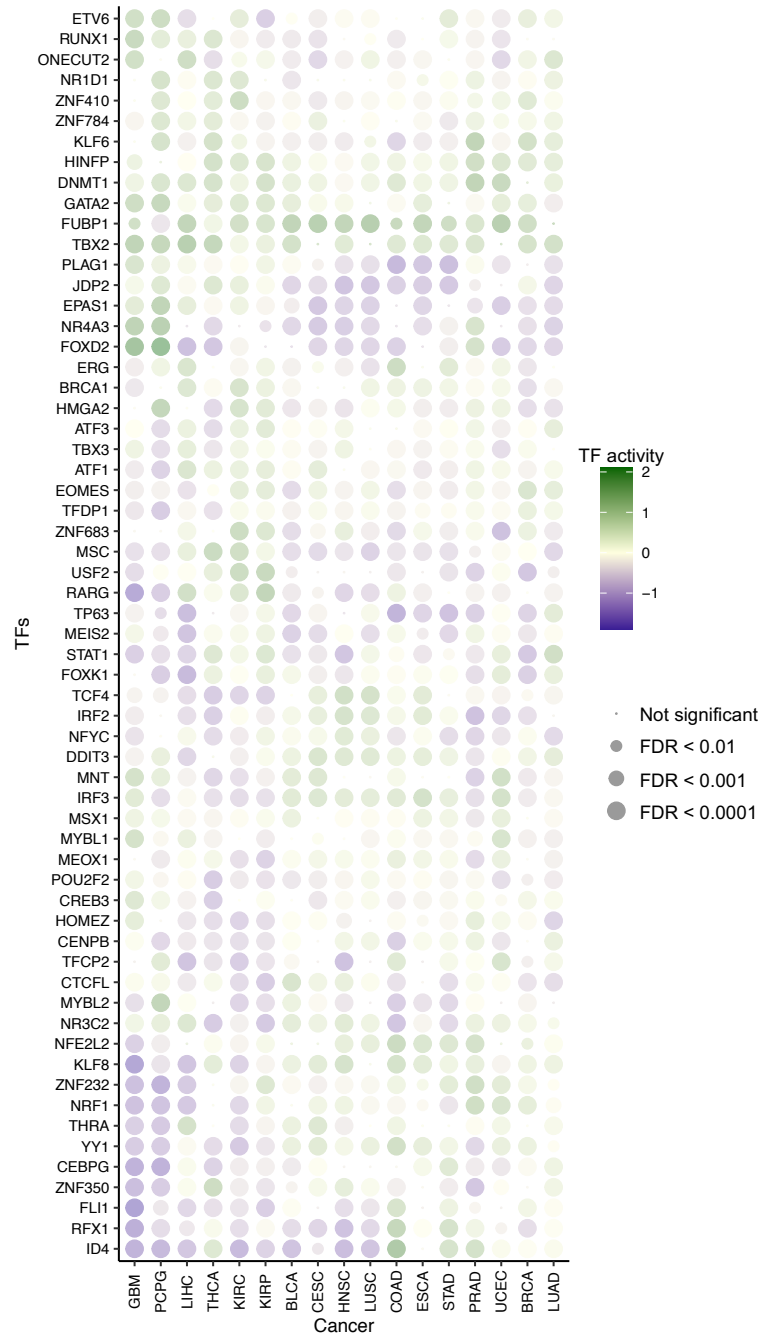
- 536 19. Osmanbeyoglu, H.U., Pelossof, R., Bromberg, J.F. and Leslie, C.S. (2014) Linking
537 signaling pathways to transcriptional programs in breast cancer. *Genome Res*, **24**, 1869-
538 1880.
- 539 20. Pelossof, R., Singh, I., Yang, J.L., Weirauch, M.T., Hughes, T.R. and Leslie, C.S. (2015)
540 Affinity regression predicts the recognition code of nucleic acid-binding proteins. *Nat*
541 *Biotechnol*, **33**, 1242-1249.
- 542 21. Osmanbeyoglu, H.U., Toska, E., Chan, C., Baselga, J. and Leslie, C.S. (2017)
543 Pancancer modelling predicts the context-specific impact of somatic mutations on
544 transcriptional programs. *Nat Commun*, **8**, 14249.
- 545 22. Ma, X., Somasundaram, A., Qi, Z., Hartman, D.J., Singh, H. and Osmanbeyoglu, H.U.
546 (2021) SPaRTAN, a computational framework for linking cell-surface receptors to
547 transcriptional regulators. *Nucleic Acids Res*, **49**, 9633-9647.
- 548 23. Gala, K., Li, Q., Sinha, A., Razavi, P., Dorso, M., Sanchez-Vega, F., Chung, Y.R.,
549 Hendrickson, R., Hsieh, J.J., Berger, M. *et al.* (2018) KMT2C mediates the estrogen
550 dependence of breast cancer through regulation of ERalpha enhancer function.
551 *Oncogene*, **37**, 4692-4710.
- 552 24. Jozwik, K.M., Chernukhin, I., Serandour, A.A., Nagarajan, S. and Carroll, J.S. (2016)
553 FOXA1 Directs H3K4 Monomethylation at Enhancers via Recruitment of the
554 Methyltransferase MLL3. *Cell Rep*, **17**, 2715-2723.
- 555 25. Tarassishin, L. and Lee, S.C. (2013) Interferon regulatory factor 3 alters glioma
556 inflammatory and invasive properties. *J Neurooncol*, **113**, 185-194.
- 557 26. Wang, M.D., Xing, H., Li, C., Liang, L., Wu, H., Xu, X.F., Sun, L.Y., Wu, M.C., Shen, F.
558 and Yang, T. (2020) A novel role of Kruppel-like factor 8 as an apoptosis repressor in
559 hepatocellular carcinoma. *Cancer Cell Int*, **20**, 422.
- 560 27. Zhang, Q., Stovall, D.B., Inoue, K. and Sui, G. (2011) The oncogenic role of Yin Yang 1.
561 *Crit Rev Oncog*, **16**, 163-197.
- 562 28. Martinez, L.A. (2016) Mutant p53 and ETS2, a Tale of Reciprocity. *Front Oncol*, **6**, 35.

563
564

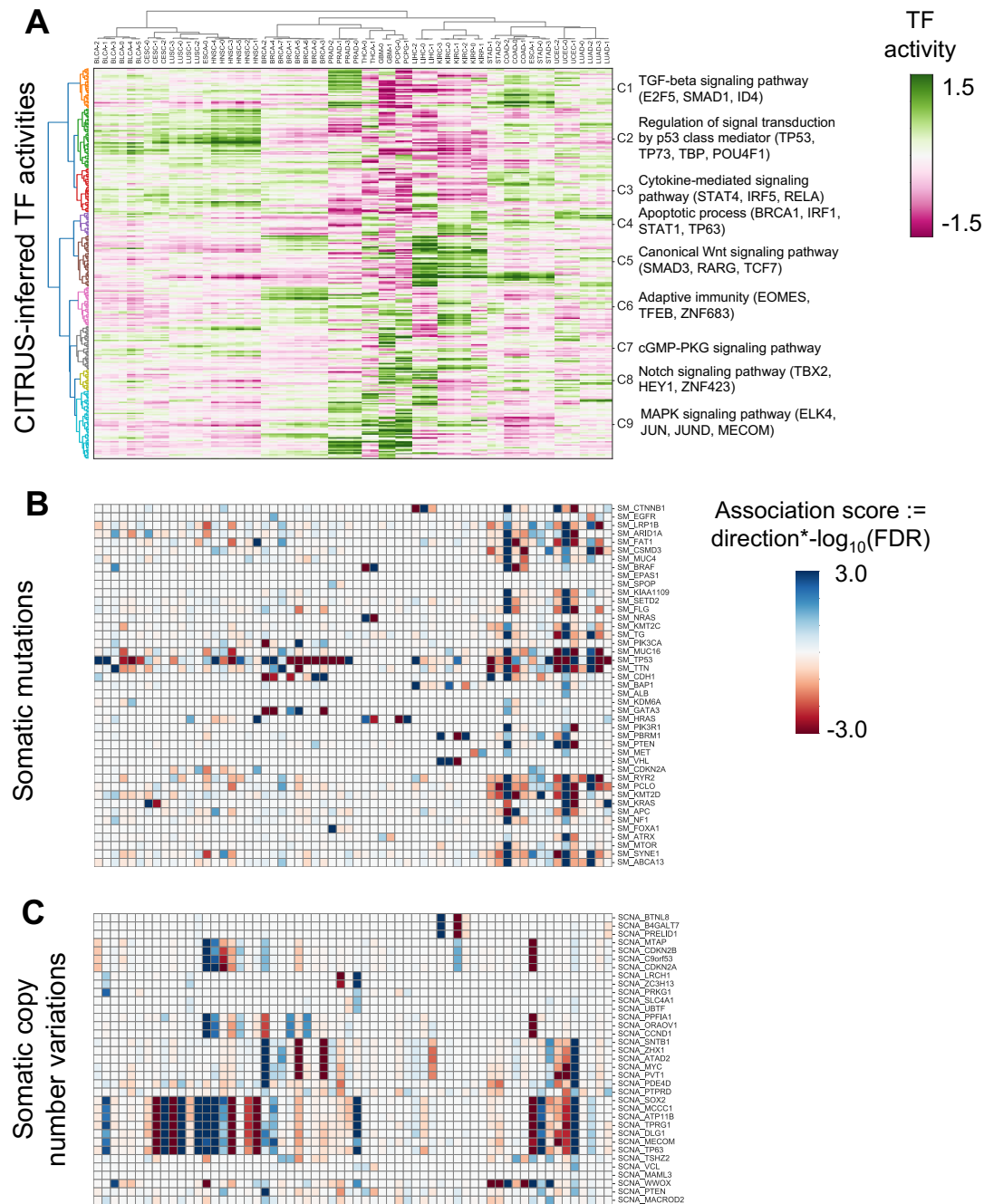
565 **Figures**
566



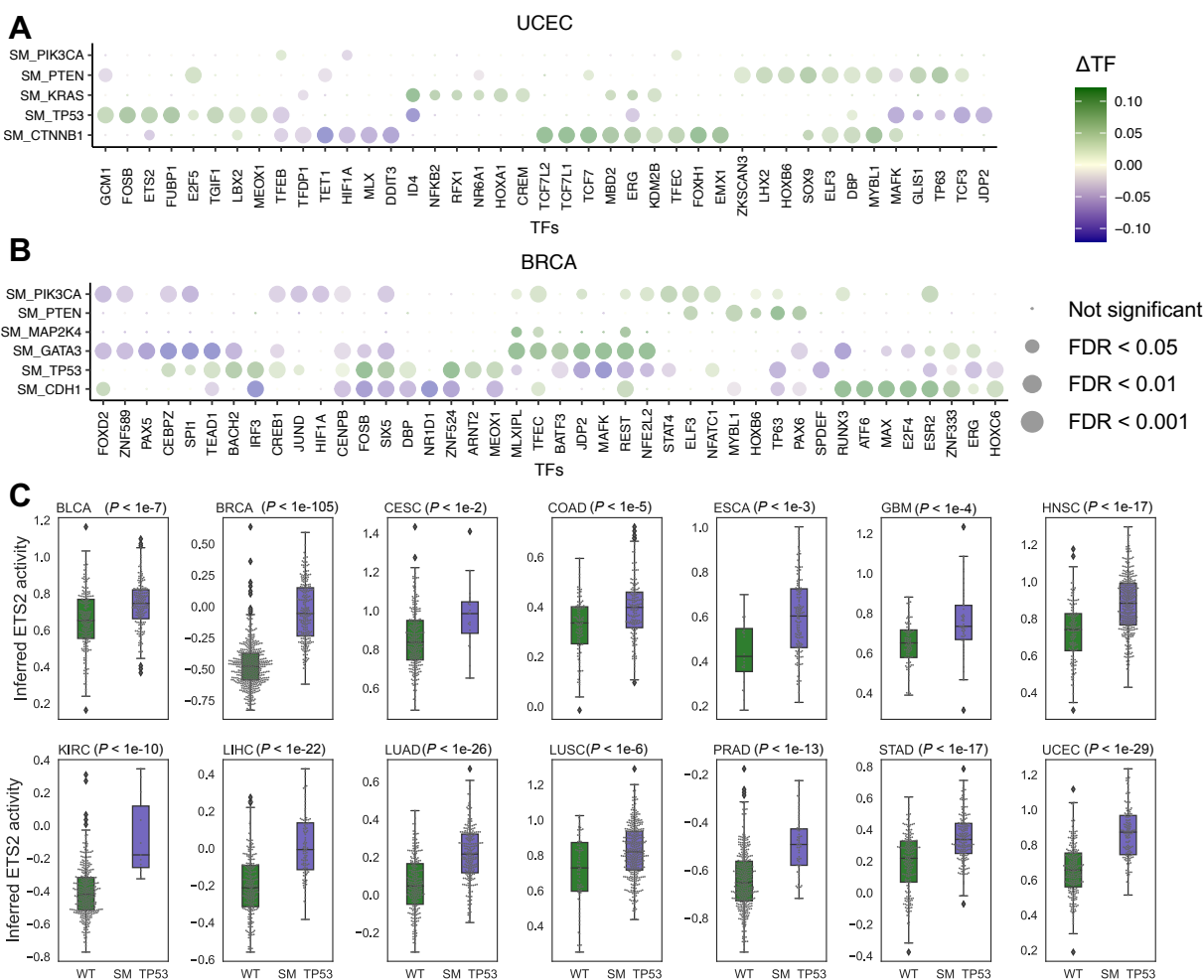
567
568 **Fig. 1: Overview of CITRUS: An attention-based model with TF-target gene priors.** The input
569 to our framework includes somatic alteration and copy number variation, assay for transposase-
570 accessible chromatin with high-throughput sequencing (ATAC-seq), tumor expression datasets
571 and TF recognition motifs. CITRUS takes somatic alteration and copy number variation data as
572 input and encodes them as a tumor embedding using a self-attention mechanism. Additional
573 cancer type information is used to stratify the confounding factor of tissue type. The middle layer
574 further transforms the tumor embeddings into a TF layer, which represents the inferred activities
575 of 320 TFs. Finally, gene expression levels are predicted from the TF activities through a TF-
576 target gene priors constrained sparse layer based on ATAC-seq.
577



590
 591 **Fig. 3: CITRUS identifies regulatory features of tumor types.** Dot plot shows the mean inferred
 592 TF activity differences between samples in a given tumor type versus those in all other tumor
 593 types by t-test. We corrected for FDR across TFs for each pairwise comparison and identified
 594 significant TFs. The complete results are included in **Supplementary Data 1**. The dot size
 595 indicates $-\log_{10}(\text{FDR})$. For clarity, the union of the top four significant TFs in each cancer type is
 596 shown.



597
598 **Fig. 4: Landscape of somatic alterations and inferred TF activities. (A)** Heatmap shows tumor
599 subtypes clustered by mean inferred TF activity. The color scale is proportional to TF activity. **(B–**
600 **C)** Heatmaps of association scores for **(B)** mutations and **(C)** copy number variations. Association
601 scores were calculated by multiplying the $-\log_{10}$ FDR by the direction derived from Fisher's exact
602 test.



603
 604 **Fig. 5: Somatic alterations are associated with dysregulated TF activity.** Impact of somatic
 605 alterations on individual TFs based on *in silico* knock out experiments in (A) UCEC and (B) BRCA
 606 datasets from TCGA. The dot plot shows mean TF activity, and dot size indicates $-\log_{10}(\text{FDR})$.
 607 See **Supplementary Fig. 3** for the full list of cancer types. (C) Inferred ETS2 activity in TCGA
 608 studies and impact of *TP53* mutations. Tumors with mutant *TP53* have significantly higher ETS2
 609 activity than WT tumors ($P < 0.01$, t-test). This association is not significant using mRNA levels of
 610 *ETS2* (**Supplementary Fig. 5**). Box edges represent the upper and lower quantile with median
 611 value shown as a bold line in the middle of the box. Whiskers extend to 1.5 times the quantile.
 612
 613