

1 **CNNC: Convolutional Neural Networks for Co-Expression Analysis**

2

3 Ye Yuan¹, Ziv Bar-Joseph^{1,2*}

4 ¹Machine Learning Department, School of Computer Science, Carnegie Mellon

5 University, Pittsburgh, PA 15213, USA. ²Computational Biology Department, School

6 of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

7 *e-mail: zivbj@cs.cmu.edu

8

9 **Abstract**

10 Co-expression analysis has been extensively used in genomics studies and tools for
11 over two decades. To date, most methods for such analysis are unsupervised and
12 symmetric. Such methods cannot infer causality and are prone to both overfitting and
13 false negatives resulting from differences between cells in bulk studies. Here we
14 present a new, supervised method based on convolutional neural networks (CNNs)
15 for co-expression analysis. We use a normalized histogram image of gene pair
16 co-expression as the input to the CNN. Testing our method on several co-expression
17 prediction tasks we show that it outperforms prior methods and that scRNA-Seq data
18 leads to more accurate results when compared to bulk data. The method can be
19 directly extended to integrate sequence and epigenetic data and to infer causal
20 relationships.

21

22 **Supporting website with software and data:** <https://github.com/xiaoyeye/CNNC>.

23

24

25

26

27

28

29

30 **Introduction**

31 Co-expression analysis, which seeks to identify genes that are correlated or
32 anti-correlated across a large number of samples or time points, has been a key
33 research area of computational genomics for almost two decades¹⁻⁵. In addition to the
34 identification of pairs of related genes, co-expression analysis serves as an initial step
35 in many of the most widely used computational methods for the analysis of genomics
36 data including various clustering methods⁶, network inference and reconstruction
37 approaches⁷⁻¹¹, methods for classification based on genes expression³ and many
38 more.

39 Given its centrality for several downstream applications, much work has focused on
40 improving the ability to infer correlated and anti-correlated genes. The most popular
41 method is based on Pearson correlation analysis¹². Such analysis focuses on shared
42 trends rather than exact values. Other popular and widely used methods involve
43 mutual information (MI)¹³⁻¹⁵, nonparametric methods, for example Spearman
44 correlation¹⁶ methods based on alignment¹⁷ and more¹⁸⁻²⁰.

45 While the above methods were shown to be useful in many applications, they also
46 suffer from serious drawbacks. The first major issue is overfitting. Given the large
47 number of genes that are profiled, and the often relatively small (at least in
48 comparison) number of samples, several genes that are determined to be
49 co-expressed may only reflect chance or noise in the data²¹. In addition, to date most
50 co-expression analysis utilized bulk gene expression data (either array or RNA-Seq).
51 In such data, correlations may be obscured by the different cell populations such that
52 even if two genes appear highly correlated, it may be that they are actually never
53 expressed in the same cell at the same time²². Finally, most of the widely used
54 co-expression analysis methods are symmetric which means that each pair has only
55 one co-expression value. While this is advantageous for some applications (for
56 example, clustering) it may be problematic for methods that aim at inferring causality
57 (for example, network reconstruction methods).

58 To address these issues we developed a new tool, CNNC which provides a

59 supervised way (that can be tailored to the condition / question of interest) to perform
60 co-expression analysis. The method utilizes both bulk and single cell RNA-Seq
61 (scRNA-Seq) data from tens of thousands of experiments, allowing us to overcome
62 cell population confounders. Our method utilizes CNNs which we tailor for the gene
63 expression analysis by representing input data for each pair of genes as an (image)
64 histogram. The network is trained with positive and negative examples for the specific
65 domain of interest (for example, known targets of a TF in a specific cell type, known
66 pathways in a specific biological process etc.). Depending on the input data the
67 training can be either symmetric or directed (for example, training the network to infer
68 that TF A regulate gene B but not vice versa). To reduce overfitting CNNC determines
69 specific thresholds based on the training for calling a pair correlated or anti-correlated
70 or for inferring causality.

71 We applied CNNC to data from tens of thousands of single cell and bulk experiments.
72 We noticed that scRNA-Seq data greatly improves performance when compared to
73 bulk RNA-Seq. Using the same expression data to learn different CNNs (by varying
74 the labels based on the specific domain the network was applied to) we show that
75 CNNC outperforms prior co-expression analysis methods both for directly inferring
76 interactions (including TF-gene and protein-protein interactions) and as a component
77 in algorithms for the reconstruction of known pathways and clustering. Finally, we
78 discuss the accuracy of the directionality predictions which are unique to CNNC and
79 shown that these predictions provide important information for determining missing
80 interactions in known pathways.

81

82

83

84

85

86

87

88 **Results**

89 We developed a general computational framework for supervised co-expression
90 analysis (**Fig. 1**). CNNC is based on CNN which is used to analyze summarized
91 co-expression histograms from pairs of genes from bulk and scRNA-Seq data. Given
92 a relatively small labeled set of positive pairs (with either negative or random pairs
93 serving as negative) the CNN learns to discriminate between interacting and / or
94 causal pairs and negative pairs. Once trained the CNN can be used to predict
95 co-expression scores for all gene pairs.

96

97 **Learning a CNNC model**

98 We used processed scRNA-Seq data that was collected from over 500 different
99 studies representing a wide range of cell types, conditions etc²³. All raw data was
100 uniformly processed and assigned to a pre-determined set of more than 20,000
101 mouse genes (Methods). We also used bulk RNA-Seq RPKM data from Encyclopedia
102 of DNA Elements (ENCODE) project²⁴, which contains 58 mouse tissues or cell types.
103 For both datasets we first generated a normalized empirical probability distribution
104 function (NEPDF) for each gene pair (genes *a* and *b*) based on their expression in the
105 scRNA-Seq or bulk RNA-Seq data (**Fig. 1**, left). We calculated their normalized 2-
106 dimension (2D) histogram and fixed its size at 32X32, where columns represent gene
107 *a* expression levels and rows represent gene *b* such that entries in the matrix
108 represent the (normalized) co-occurrences of these values. See Methods for details.
109 Bulk and sc NEPDF were then either used separately or concatenated to form a
110 combined NEPDF with dimension of 32X64. Next, the histogram matrix is used as
111 input to a CNN which is trained using a N-dimension (ND) output label vector, where
112 N depends on specific tasks. In our case N can either be 1 (interacting or not) or 3 in
113 which case label 0 indicates that genes *a* and *b* are not interacting and label 1 (2)
114 indicates that gene *a* (*b*) regulates gene *b* (*a*). Because of the final 'softmax' layer
115 classification utilized by CNNs, for a three-label task CNNC's output is a vector
116 consisting of three respective probabilities, [p_0 , p_1 , p_2], which sum to 1. In general,

117 our CNN model consists of one 32X32 or 32X64 input layer, ten intermediate layers
118 including six convolutional layers, three maxpooling layers, one flatten layer, and a
119 final ND 'softmax' layer or one scalar 'Sigmoid' layer (**Methods** and **Supplementary**
120 **Fig. 1**).

121 In addition to gene expression data, we can easily integrate other data types including
122 Dnase-seq, PWM, etc. For this, we concatenated the additional information as a
123 vector to the intermediate output of the gene expression data and continued with the
124 standard CNN architecture. See **Methods** and **Supplementary Fig. 1** for complete
125 details and **Supplementary Table 1** for information on training and run time.

126

127 **Using CNNC to predict TF-gene interactions**

128 Chromatin immunoprecipitation (ChIP)-seq has been widely used as a gold standard
129 for studying cell-specific protein-DNA interactions²⁵. Here we first evaluated CNNC's
130 performance on regulator-target prediction based on data from the GTRD ChIP-seq
131 database²⁶.

132 We extracted data from GTRD for 41 TFs for which ChIP-seq experiments were
133 performed in mouse embryonic stem cell (mESC). To determine targets for each TF
134 based on the ChIP-seq data from GTRD, we followed ref 27 and 28^{27, 28} and defined
135 the promotor region as 10KB upstream to 1KB downstream from the transcription
136 start site (TSS) for each gene. If a TF X has at least one detected peak signal with
137 p-value smaller than 10^{-300} in or overlapping the promotor region of gene Y, we say
138 that TF X regulates gene Y. We also used this data to compare CNNC with the two
139 most popular co-expression analysis methods: Pearson correlation (PC) and mutual
140 information (MI) and to compare the accuracy of predictions based on the sc and bulk
141 RNA-Seq data. Since the prior methods used for comparison are symmetric, we
142 focused here on the two labels setting (interacting or not). We later discuss causality
143 inference on this data. To compare the methods and data types we performed
144 leave-one-TF-out cross validation analysis. For each of the 41TFs, we trained CNNC
145 with all other TFs and used the left out TF for testing (**Methods**).

146 **Fig. 2** presents the results of this comparison analysis. First, we see that for all
147 methods scRNA-Seq data (left column) provides much more information when
148 compared to bulk (middle column). Note that while we had more scRNA-Seq profiles
149 in our training set when compared to bulk experiments, these actually represent much
150 fewer cells and conditions than those used in the bulk data. We have also tested the
151 performance when using the same number of bulk and scRNA profiles
152 (**Supplementary Fig. 2**). We found that even with this very small number of
153 scRNA-Seq profiles (with much fewer cells than the bulk) CNNC performs better when
154 using scRNA-Seq. These results support prior claims about convolution effects
155 resulting from population of cells that make target inference harder when using bulk
156 data^{22, 29}. Still, bulk data did include some useful information and for all methods since
157 the joint sc and bulk data performed best when compared to individual data type on its
158 own. As for the methods themselves, for all types of input data, CNNC outperforms
159 the other two methods. This is especially noticeable when using the scRNA-Seq (and
160 combined) data where CNNC is 15% more accurate than MI and close to 20% more
161 accurate than PC. The difference is even more pronounced for the top ranked
162 predictions. Here, for CNNC we see almost no false negatives for the first 15% of
163 ranked pairs (inset, **Fig. 2i**).

164

165 **Data Integration further improves TF target gene prediction**

166 The above analysis was only based on using expression values. However, as noted
167 above, gene co-expression is often used as a component in more extensive
168 procedures that often integrate different types of genomics data. To test how the use
169 of the NN based method can aid such procedures we combined the co-expression
170 values obtained by our method and the other methods with sequence and DNase
171 hypersensitivity data. For sequence, we used PWMs for the TFs we tested from the
172 Jaspar website³⁰. We have also used Dnase-seq data for the same cell line from the
173 mouse ENCODE project²⁴. While there are several different methods for integrating
174 expression, sequence and DNase data, since our main focus here is on the

175 co-expression analysis methods we used a simple strategy for processing the PWM
176 and DNase data (**Methods**) which resulted in an additional 2D vector as input for
177 each pair which we embedded to create a 512D vector (**Fig. 1 and Methods**). We
178 next extended the CNN to utilize the additional data by concatenating it with the
179 NEPDF's 512D vector in the flatten layer to form a 1024D vector as shown in **Fig. 1**
180 and **Supplementary Fig. 1**.

181 Results, presented in **Fig. 2j**, show that these additional data sources indeed improve
182 the ability to predict TF-gene interactions. However, as before a combined framework
183 utilizing our CNNC method for co-expression analysis outperformed a method that
184 used both MI and PC. Thus, the NN based approach can successfully replace other
185 methods as a component in a more elaborated systems biology framework for
186 inferring interactions.

187

188 **CNNC can predict pathway regulator-target gene pairs**

189 While TFs usually directly impact the resulting expression profile of their target genes
190 (and so co-expression analysis seems like a natural option to study such interactions)
191 several methods have also utilized RNA-Seq data to infer pathways that combine
192 protein-protein and protein-DNA interactions³¹⁻³⁵. To test whether CNNC can serve as
193 a component in pathway inference methods we selected two representative pathway
194 databases, KEGG³⁶ and Reactome³⁷ as gold standard and used these to train and
195 test our co-expression framework. Since we are interested in causal relationships we
196 only used directed edges with activation or inhabitation edge types and filtered out
197 cyclic gene pairs where genes regulate each other mutually (to allow for a unique
198 label for each pair). As for the negative data, here we limited the negative set to a
199 random set of pairs where both genes appear in pathways in the database but do not
200 interact. Here leave-one-gene-out cross validation strategy requires extremely large
201 computing resources due to the large number of genes with outgoing directed edges
202 (3,057 for KEGG and 2,519 for Reactome). Instead, we performed a three-fold cross
203 validation where we kept the set of genes for which we predicted interactions

204 completely separated (so a gene in the test set does not have any interaction in the
205 training set). The positive data was uniformly divided by the outgoing gene into three
206 equal sized outgoing subsets, CNNC was trained using any two subsets and
207 evaluated using the left subset, and then the test ROCs for each outgoing gene in the
208 three subsets were calculated (**Methods**). Results are presented in **Fig. 3**. As can be
209 seen, CNNC performs very well on the KEGG pathways (See **Supplementary Fig. 3**
210 for the different folds) and also performs quite well on Reactome pathways (see also
211 **Supplementary Fig. 3**). In contrast, the other co-expression analysis methods do not
212 perform as well on these datasets (**Supplementary Fig. 4**).

213

214 **Using CNNC for casualty prediction**

215 So far we focused on general interaction predictions, which is what most symmetric
216 co-expression analysis are aimed at. However, as discussed above CNNC can also
217 be used to infer directionality by changing the output of the NN.

218 We next used CNNC to infer causal edges for all three datasets we studied (for
219 TF-gene interactions causal relationships are clear, for the pathway database we only
220 analyzed directed edges and so had the ground truth for that data as well). As can be
221 seen in **Fig. 4**, when using the TF GTRD dataset, CNNC achieves a median AUROC
222 of 0.8227 (**Fig. 4a**), with 32 of the 41 TFs obtaining an AUROC of more than 0.5 on
223 this leave-one-TF-out classification task. Interestingly, as can be seen in
224 **Supplementary Fig. 5**, when only using bulk RNA-Seq data performance on the
225 GTRD data prediction is very weak. Thus, for the causality inference task scRNA-Seq
226 data is the only one that can provide enough information. For KEGG, CNNC is very
227 successful achieving a median AUROC of 0.9949 (**Fig. 4c**) (See **Supplementary Fig.**
228 **6** for the different folds). For Reactome (**Fig. 4e**) we see that the most confident
229 predictions are correct, but beyond the top predictions performance levels off (See
230 **Supplementary Fig. 6** for the different folds). To try to understand the process used
231 by the NN to distinguish causal directions we plotted two NEPDF inputs to the NN
232 (**Figs. 4g and 4h**) which were correctly predicted as two different labels (1 for **4g** and

233 2 for **4h**). As can be seen, in both inputs the two genes display partial correlations and
234 there are places where both are up or down concurrently. However, the main
235 difference between the histograms in **4g** and **4h** are cases where one gene is up and
236 the other is not. In **4g** gene 2 is up while gene 1 is not indicating that the causal
237 relationship is likely $g1 \rightarrow g2$. The opposite holds for **4h** and so the method infers that
238 $g2 \rightarrow g1$ for that input. Thus, unlike the prior symmetric methods, the NEPDF that
239 serves as input provides important clues that the NN can utilize to infer both
240 interactions and causality.

241

242 **Pathway application**

243 Given the results for KEGG we asked whether we can use the CNNC method to infer
244 missing edges in current pathways. There have been several attempts to utilize
245 expression and other data to further refine known pathways and many of these are
246 based on co-expression analysis^{9, 19-21, 32, 33, 38, 39}. Since our method provides both
247 direction and score we can extract all predicted directed edges above a certain score
248 and compare the resulting pathway to the database pathway to see if any additional
249 edges, that do not appear in the database, are predicted by our method. For this we
250 focused on the interleukin 17 (IL-17) pathway from KEGG database, which plays
251 crucial roles in inflammatory responses. We extracted 6 proteins and 4 directed edges
252 from this pathway by only using directed edges with activation or inhibition edge
253 types and filtering out cyclic gene pairs (**Fig. 5a**). We applied CNNC trained on all
254 database pathway edges that do not contain any of these 6 proteins. As can be seen
255 in **Supplementary Fig. 7**, CNNC predicted 8 of the possible 30 edges for this path
256 (15 pairs with two possible directions each). 4 of these 8 were the original 4 directed
257 edges annotated in the database itself. The other 4 edges were not present in the
258 KEGG as causal interactions for this pathway. However, all are either supported by
259 their presence in other KEGG pathways or by recent publications (we reiterate that
260 interactions for all six proteins in other KEGG pathways were excluded from the
261 training data as mentioned above so these predictions are not contaminated by their

262 presence in other pathways in KEGG). Specifically, (traf6, nfkb1) and (map3k7, nfkb1)
263 have been annotated as causal pairs in KEGG's 'Pertussis pathway' and 'RIG-I-like
264 receptor signaling pathway'³⁶, respectively. (rela, nfkb1) is the known p50/p65
265 heterodimer of NF- κ B⁴⁰. As for the (nfkb1a, traf6) pair, it was found that traf6's binding
266 to MAP3K7 activates ikkbk which in turn phosphorylates nfkb1a⁴¹.

267

268 **CNNC output as a similarity matrix for clustering**

269 To evaluate CNNC's performance on the downstream analysis, we used it to generate
270 a gene-gene similarity matrix. We next used this matrix as an input for a hierarchical
271 clustering algorithm.

272 We extracted the top 2,000 (top 1,000, see **Supplementary Fig. 8**) differentially
273 expressed genes based on the expression data used in this paper using fano factor
274 (FF) (**Supplementary Note**). Since we have trained CNNC using the KEGG database,
275 we removed KEGG genes from the test set. Next we performed hierarchical
276 clustering⁴² using CNNC and PC based on all sc and bulk data (**Fig. 5b**, and **5c**). For
277 comparison, we selected the top 8 clusters for the resulting hierarchical clustering tree
278 for all inputs (see also **Supplementary Fig. 8**). Next, for each input we calculated the
279 significant GO terms (q-values < 0.05)⁴³ and plotted the results in **Fig. 5d**. As can be
280 seen, using CNNC as the input led to the identification of more significant GO terms
281 for the same set of genes indicating that the clustering obtained using this input is
282 more aligned with current biological knowledge.

283

284

285

286

287

288

289

290

291 **Discussion and conclusion**

292 Gene co-expression analysis has been a widely successful method for the analysis of
293 gene expression data starting over two decades ago with the introduction of
294 microarrays. Several methods have been suggested for this task and several other
295 methods use co-expression analysis as a component in a more elaborate modeling
296 framework.

297 While co-expression analysis performs well in some cases, it suffers from a number of
298 drawbacks that often led to overfitting (false positives) or missing key relationships
299 (false negatives). The former can be attributed to the unsupervised nature of most
300 co-expression methods making it hard to ‘train’ them on a labeled dataset. The latter
301 often resulted from the nature of the data used for co-expression analysis (bulk or
302 population of cells data) which led to masking of relationships that existed in single
303 cells. Moreover, while certain more sophisticated methods attempted to utilize gene
304 expression to infer causality (for example, Bayesian network based methods⁴⁴) these
305 were only able to detect directed interactions, were based on very specific
306 probabilistic modeling assumptions, and did not directly provide a confidence score
307 for the resulting edges.

308 To address these issues we presented CNNC, a general framework for co-expression
309 analysis which is based on convolutional NN (CNN). The key idea here is to convert
310 the input data into a co-expression histogram which is very suitable for CNNs. Unlike
311 most prior methods our method is supervised which allows the CNN to zoom in on
312 subtle differences between positive and negative pairs. Supervision also helps fine
313 tune the scoring function based on the different application. For example, different
314 features may be important for analyzing TF-gene interactions when compared to
315 inferring proteins in the same pathway. In addition to the supervised approach the fact
316 that the network can utilize the large volumes of scRNA-Seq data allows it to better
317 overcome masking issues reducing false negative.

318 Analysis of several different interaction prediction tasks indicates that CNNC can
319 improve upon prior, unsupervised methods. It can also be naturally extended to

320 integrate complementary data including epigenetic and sequence information. Finally,
321 CNNC is easy to use either with general data or with condition specific data. For the
322 former, users can download the data and implementation from the supporting website
323 (**Supplementary Fig. 9**), provide a list of labels (positive and negative pairs for their
324 system of interest) and retrieve the scores for all possible gene pairs. These in turn
325 can be used for any downstream application including clustering, network analysis
326 etc.

327 In addition to comparing CNNC to prior methods we have also used it to evaluate the
328 advantages conferred by scRNA-Seq data. Models trained with scRNA-Seq data
329 outperformed those trained with bulk data for all systems we looked at. This supports
330 prior findings^{45, 46} and addressed a key criticism of co-expression analysis – that many
331 interactions are observed or missed due to aggregation effects from the collection of
332 cells rather than because they truly represent specific molecular events. While the
333 scRNA-Seq data we used contained two orders of magnitude more samples, the total
334 number of cells profiled is smaller (each bulk experiment often profiles at least three
335 orders of magnitude more cells than a single scRNA-Seq profile⁴⁷). In addition,
336 scRNA-Seq coverage is often two orders of magnitude less than bulk experiments so
337 that total number of reads in the two datasets is not very different. Even when
338 comparing with the same number of profiles for bulk and sc we find that CNNC
339 performs better when using scRNA-Seq data. This result seems to indicate that
340 despite the much greater noise associated with scRNA-Seq, such data can provide
341 more accurate models for the same overall costs, coverage and sample size.

342 CNNC is implemented in Python and both data and an open source version of the
343 software are available from the supporting website. We hope that this method would
344 become a useful component in future co-expression studies.

345

346

347

348

349 **Online methods**

350 **Dataset sources and pre-process pipelines**

351 We used mouse scRNA-Seq dataset collected by Alavi et al²³. The dataset consists of
352 uniformly processed 43,261 expression profiles from over 500 different scRNA-Seq
353 studies. For each profile, expression values are available for the same set of 20,463
354 genes. Mouse bulk RNA-Seq dataset were downloaded from Mouse Encode project
355 ²⁴. That data included 249 samples and we only utilized genes that are present in the
356 scRNA-Seq dataset leading to the same number of genes for both datasets. mESC
357 DNase data was also downloaded from Mouse Encode project²⁴
358 (ENCFF096WRW.bed). Mouse TF motif information is from TRANSFAC database⁴⁸.
359 PWM values were calculated by Python package 'Biopython'⁴⁹.

360 For the DNase and PWM analysis we followed prior papers and defined the
361 transcription start site (TSS) region as 10KB upstream to 1KB downstream from the
362 TSS for each gene^{27, 28}. For each TF and gene pair, using Biopython package we
363 calculated the score between the TF motif sequence and both the '+/-' sequences at
364 all possible positions along the TSS region of the gene, and then selected the
365 maximum one as the final PWM score. The maximum DNase peak signal in the TSS
366 region was calculated as the scalar DNase value for each gene.

367

368 **Labeled data:**

369 mESC ChIP-seq peak region data was downloaded from GTRD database, and we
370 used peaks with threshold p value $< 10^{-300}$. If one TF X has at least one ChIP-seq
371 peak signal in or partially in the TSS region of gene Y, as defined above, we say that X
372 regulates Y.

373 KEGG and Reactome pathway data were downloaded by the R package 'graphite'⁵⁰.

374 KEGG contains 290 pathways and Reactome contains 1581 pathways. For both, we
375 only select directed edges with either activation or inhibition edge types and filter
376 out cyclic gene pairs where genes regulate each other mutually (to allow for a unique
377 label for each pair). In total, we have 3,057 proteins with outgoing directed edges in

378 KEGG and the total number of directed edges is 33,127. For Reacotome the
379 corresponding numbers are 2,519 and 33,641.

380

381

382 **Constructing the input histogram**

383 For any gene pair a and b , we first log transformed their expression, and then
384 uniformly divided the expression range of each gene to 32 bins. Next we created the
385 32X32 histogram by assigning each sample to an entry in the matrix and counting the
386 number of samples for each entry. Due to the very low expression levels and even
387 more so to dropouts in scRNA data, the value in zero-zero position is always very
388 large and often dominates the entire matrix. To overcome this, we added
389 pseudocounts to all entries. We combined bulk and scRNA-Seq NEPDFs by
390 concatenating them as a 32X64 matrix to achieve better performance.

391

392 **CNN for RPKM data**

393 We followed VGGnet⁵¹ to build our convolutional neural networks (CNN) model
394 (**Supplementary Fig. 1**). The CNN consists of stacked layers of x 3X3 convolutional
395 filters (equation (1)) (x is a power of 2, ranging from 32 to 64 to 128) and interleaved
396 layers of 2X2 maxpooling (equation (2)). We used the constructed input data as input
397 to CNN. Each convolution layer computes the following function:

$$398 \text{ Convolution } (X)_{i,j}^k = \sum_{m=1}^3 \sum_{n=1}^3 W_{i,j}^k X_{i+m,j+n} \quad (1)$$

399 Where X is the input from the previous layer, (i,j) is output position, k is convolutional
400 filter index and W is the filter matrix of size 3X3. In other words, each convolutional
401 layer computes a weighted average of the prior layer values where the weights are
402 determined based on training. The maxpooling layer computes the following function:

$$403 \text{ maxpooling } (X)_{i,j}^k = \max(\{X_{i,j}^k, X_{i+1,j}^k, X_{i,j+1}^k, X_{i+1,j+1}^k\}) \quad (2)$$

404 Where X is input, (i,j) is output position and k is the convolutional filter index. In other
405 words, the layer selects one of the values of the previous layer to move forward.

406

407 **Overall structure**

408 The overall structure of the CNN is presented in **Supplementary Fig. 1**. The input
409 layer of the CNN is either 32X32 or 32X64 as discussed above. In addition, the CNN
410 contains 10 intermediate layers and a single one or three-dimension output layer. The
411 ten layers include both convolutional and maxpooling layers, and the exact
412 dimensions of each layer are shown in **Supplementary Fig. 1**. Following ref 52⁵² we
413 used rectified linear activation function (ReLU) as the activation function (equation (3))
414 across the whole network, except the final classification layers where 'sigmoid'
415 function (equation (4)) was used for two categories classification and 'softmax'
416 function (equation (5)) for multiple categories classification. These functions are
417 defined below.

418
$$\text{ReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \quad (3)$$

419
$$\text{Sigmoid}_{\theta}(x) = \frac{1}{1 + e^{\theta x}} \quad (4)$$

420
$$\text{Softmax}_{\theta}(x) = \frac{1}{\sum_{j=1}^k e^{\theta_j x}} \begin{bmatrix} e^{\theta_1 x} \\ e^{\theta_2 x} \\ \dots \\ e^{\theta_k x} \end{bmatrix} \quad (5)$$

421

422

423 **Training and testing strategy**

424 We evaluated the CNN using cross validation. In these, training and test datasets are
425 strictly separated to avoid information leakage. See **Supplementary Note**,
426 **Supplementary Fig. 10** and **Supplementary Table 1** for details. For the three labels
427 (causality analysis) we did the following: for each gene, we generated (*a*, *b*) (label1)
428 and (*b*, *a*)'s (label2) NEPDF matrices. For the 0 label we generated a (*a*, *N*) NEPDF
429 matrices for GTRD where *N* was a random gene and *a* was the TF. 0 labels for KEGG
430 and Reactome were generated from random (*M*, *N*) gene pairs. After training, we
431 used $p1(a, b) + p2(a, b)$ as the probability that *a* interacts *b*, $p2(a, b) - p2(b, a)$ as the
432 pseudo probability that *b* regulates *a*.

433

434 **Integrating expression, sequence and DNase data**

435 To integrate Dnase and PWM data with the processed RNA-Seq data, we first
436 computed the max value for a PWM scan and DNase accessibility for each promotor
437 region. We next generated a two-value vector from this data for each pair and
438 embedded it to a 512D vector using one fully connected layer containing 512 nodes.
439 Next these are concatenated with the expression processed data to form a 1024D
440 vector which serves as input to a fully connected 512-node plus 128-node layer neural
441 network classifier. See **Supplementary Fig. 1** for details. Early stopping strategy by
442 monitoring validation loss function is used to avoid overfitting.

443

444 **Selection of edges for the IL-17 pathway analysis**

445 We performed leave-one-pathway-out validation to evaluate CNNC' performance for
446 predicting edges for individual pathways. We selected a relatively small pathway
447 ('IL-17' from KEGG) to improve our ability to present it visually. We discuss more
448 general results for KEGG as well (**Fig. 4**). For this analysis we only selected directed
449 edges with either activation or inhabitation types and filtered out cyclic gene pairs
450 where genes regulate each other mutually to purify the edge types. In total, we had 6
451 nodes and 4 directed edges for the IL-17 pathway. Next, we trained CNNC with the
452 entire KEGG dataset excluding any interactions for the six 'IL-17' pathway proteins.

453

454 **Hierarchical clustering and GO term enrichment analysis**

455 We performed hierarchical clustering followed by GO term enrichment analysis to
456 evaluate CNNC' performance in downstream analysis. We selected the top 2,000 (or
457 1,000 (**Supplementary Fig. 8**)) genes with highest Fano factor (**Supplementary**
458 **Note**) We obtained the similarity matrices for the filtered gene list based on CNNC, sc
459 PC, bulk PC and sc&bulk PC respectively. We cut the tree at 8 clusters for all inputs.
460 Next, we performed GO term enrichment analysis using fisher's exact test and
461 counted the significant GO terms for each of the cluster result. Significance of
462 difference between different inputs was computed using one-side

463 Wilcoxon-Mann-Whitney test for the q-values of the four strategies (**Fig. 5d**).

464

465 **Acknowledgements**

466 Work partially supported by NIH grant 1R01GM122096, US National Science
467 Foundation (DBI-1356505) to ZBJ. and a James S. McDonnell Foundation Scholars
468 Award in Studying Complex Systems to Z.B.-J.

469

470 **Author contributions**

471 Y.Y. and Z.B conceived the method. Y.Y. implemented CNNC and the support website.
472 Y.Y. and Z.B designed the experiments. Y.Y. and Z.B wrote the manuscript.

473

474 **Competing interests**

475 None

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493 **References**

494

- 495 1. De Smet, R. & Marchal, K. Advantages and limitations of current network inference
496 methods. *Nat Rev Microbiol* **8**, 717-729 (2010).
- 497 2. Stuart, J.M., Segal, E., Koller, D. & Kim, S.K. A gene-coexpression network for global
498 discovery of conserved genetic modules. *Science* **302**, 249-255 (2003).
- 499 3. van Dam, S., Vosa, U., van der Graaf, A., Franke, L. & de Magalhaes, J.P. Gene
500 co-expression analysis for functional classification and gene-disease predictions. *Brief*
501 *Bioinform* (2017).
- 502 4. Butte, A.J. & Kohane, I.S. Unsupervised knowledge discovery in medical databases
503 using relevance networks. *Proc AMIA Symp*, 711-715 (1999).
- 504 5. Eisen, M.B., Spellman, P.T., Brown, P.O. & Botstein, D. Cluster analysis and display
505 of genome-wide expression patterns. *Proc Natl Acad Sci U S A* **95**, 14863-14868
506 (1998).
- 507 6. Newman, M.E. Spectral methods for community detection and graph partitioning.
508 *Phys Rev E Stat Nonlin Soft Matter Phys* **88**, 042822 (2013).
- 509 7. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network
510 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 511 8. Bartlett, T.E., Muller, S. & Diaz, A. Single-cell Co-expression Subnetwork Analysis. *Sci*
512 *Rep* **7**, 15066 (2017).
- 513 9. Chan, T.E., Stumpf, M.P.H. & Babbie, A.C. Gene Regulatory Network Inference from
514 Single-Cell Data Using Multivariate Information Measures. *Cell Syst* **5**, 251-267 e253

- 515 (2017).
- 516 10. Huynh-Thu, V.A., Irrthum, A., Wehenkel, L. & Geurts, P. Inferring regulatory networks
517 from expression data using tree-based methods. *PLoS One* **5** (2010).
- 518 11. Faith, J.J. et al. Large-scale mapping and validation of Escherichia coli transcriptional
519 regulation from a compendium of expression profiles. *PLoS Biol* **5**, e8 (2007).
- 520 12. Lee, H.K., Hsu, A.K., Sajdak, J., Qin, J. & Pavlidis, P. Coexpression analysis of human
521 genes across many microarray data sets. *Genome Res* **14**, 1085-1094 (2004).
- 522 13. Song, L., Langfelder, P. & Horvath, S. Comparison of co-expression measures:
523 mutual information, correlation, and model based indices. *BMC Bioinformatics* **13**, 328
524 (2012).
- 525 14. Daub, C.O., Steuer, R., Selbig, J. & Kloska, S. Estimating mutual information using
526 B-spline functions--an improved similarity measure for analysing gene expression
527 data. *BMC Bioinformatics* **5**, 118 (2004).
- 528 15. Basso, K. et al. Reverse engineering of regulatory networks in human B cells. *Nat*
529 *Genet* **37**, 382-390 (2005).
- 530 16. Kotlyar, M., Fuhrman, S., Ableson, A. & Somogyi, R. Spearman correlation identifies
531 statistically significant gene expression clusters in spinal cord development and injury.
532 *Neurochem Res* **27**, 1133-1140 (2002).
- 533 17. Qian, J., Dolled-Filhart, M., Lin, J., Yu, H. & Gerstein, M. Beyond synexpression
534 relationships: local clustering of time-shifted and inverted gene expression profiles
535 identifies new, biologically relevant interactions. *J Mol Biol* **314**, 1053-1066 (2001).
- 536 18. Reverter, A. & Chan, E.K. Combining partial correlation and an information theory

- 537 approach to the reversed engineering of gene co-expression networks. *Bioinformatics*
538 **24**, 2491-2497 (2008).
- 539 19. Kumari, S. et al. Evaluation of gene association methods for coexpression network
540 construction and biological knowledge discovery. *PLoS One* **7**, e50411 (2012).
- 541 20. Wang, Y.X., Waterman, M.S. & Huang, H. Gene coexpression measures in large
542 heterogeneous samples using count statistics. *Proc Natl Acad Sci U S A* **111**,
543 16371-16376 (2014).
- 544 21. Freytag, S., Gagnon-Bartsch, J., Speed, T.P. & Bahlo, M. Systematic noise degrades
545 gene co-expression signals but can be corrected. *BMC Bioinformatics* **16**, 309 (2015).
- 546 22. Crow, M., Paul, A., Ballouz, S., Huang, Z.J. & Gillis, J. Exploiting single-cell
547 expression to characterize co-expression replicability. *Genome Biol* **17**, 101 (2016).
- 548 23. Amir Alavi, M.R., Aiyappa Parvangada, Zhilin Huang, Ziv Bar-Joseph scQuery: a web
549 server for comparative analysis of single-cell RNA-seq data. (2018).
- 550 24. Yue, F. et al. A comparative encyclopedia of DNA elements in the mouse genome.
551 *Nature* **515**, 355-364 (2014).
- 552 25. Johnson, D.S., Mortazavi, A., Myers, R.M. & Wold, B. Genome-wide mapping of in
553 vivo protein-DNA interactions. *Science* **316**, 1497-1502 (2007).
- 554 26. Yevshin, I., Sharipov, R., Valeev, T., Kel, A. & Kolpakov, F. GTRD: a database of
555 transcription factor binding sites identified by ChIP-seq experiments. *Nucleic Acids*
556 *Res* **45**, D61-D67 (2017).
- 557 27. Schulz, M.H. et al. Reconstructing dynamic microRNA-regulated interaction networks.
558 *Proc Natl Acad Sci U S A* **110**, 15686-15691 (2013).

- 559 28. Schulz, M.H. et al. DREM 2.0: Improved reconstruction of dynamic regulatory
560 networks from time-series expression data. *BMC Syst Biol* **6**, 104 (2012).
- 561 29. Wang, J. et al. Single-Cell Co-expression Analysis Reveals Distinct Functional
562 Modules, Co-regulation Mechanisms and Clinical Outcomes. *PLoS Comput Biol* **12**,
563 e1004892 (2016).
- 564 30. Khan, A. et al. JASPAR 2018: update of the open-access database of transcription
565 factor binding profiles and its web framework. *Nucleic Acids Res* **46**, D260-D266
566 (2018).
- 567 31. Gitter, A., Carmi, M., Barkai, N. & Bar-Joseph, Z. Linking the signaling cascades and
568 dynamic regulatory networks controlling stress responses. *Genome Res* **23**, 365-376
569 (2013).
- 570 32. van Noort, V., Snel, B. & Huynen, M.A. Predicting gene function by conserved
571 co-expression. *Trends Genet* **19**, 238-242 (2003).
- 572 33. Parikshak, N.N. et al. Integrative functional genomic analyses implicate specific
573 molecular pathways and circuits in autism. *Cell* **155**, 1008-1021 (2013).
- 574 34. Willsey, A.J. et al. Coexpression networks implicate human midfetal deep cortical
575 projection neurons in the pathogenesis of autism. *Cell* **155**, 997-1007 (2013).
- 576 35. Young, J.H. & Marcotte, E.M. Predictability of Genetic Interactions from Functional
577 Gene Modules. *G3 (Bethesda)* **7**, 617-624 (2017).
- 578 36. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new
579 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res* **45**,
580 D353-D361 (2017).

- 581 37. Fabregat, A. et al. The Reactome Pathway Knowledgebase. *Nucleic Acids Res* **46**,
582 D649-D655 (2018).
- 583 38. Angelini, C. & Costa, V. Understanding gene regulatory mechanisms by integrating
584 ChIP-seq and RNA-seq data: statistical solutions to biological problems. *Front Cell*
585 *Dev Biol* **2**, 51 (2014).
- 586 39. D'Haeseleer, P., Liang, S. & Somogyi, R. Genetic network inference: from
587 co-expression clustering to reverse engineering. *Bioinformatics* **16**, 707-726 (2000).
- 588 40. Chen, F.E., Huang, D.B., Chen, Y.Q. & Ghosh, G. Crystal structure of p50/p65
589 heterodimer of transcription factor NF-kappaB bound to DNA. *Nature* **391**, 410-413
590 (1998).
- 591 41. Kramer, I. Signal Transduction. (Elsevier, 2015).
- 592 42. Bar-Joseph, Z., Gifford, D.K. & Jaakkola, T.S. Fast optimal leaf ordering for
593 hierarchical clustering. *Bioinformatics* **17 Suppl 1**, S22-29 (2001).
- 594 43. Ashburner, M. et al. Gene ontology: tool for the unification of biology. The Gene
595 Ontology Consortium. *Nat Genet* **25**, 25-29 (2000).
- 596 44. Spirtes, P. & Zhang, K. Causal discovery and inference: concepts and recent
597 methodological advances. *Appl Inform (Berl)* **3**, 3 (2016).
- 598 45. Wagner, A., Regev, A. & Yosef, N. Revealing the vectors of cellular identity with
599 single-cell genomics. *Nat Biotechnol* **34**, 1145-1160 (2016).
- 600 46. Wang, Y. & Navin, N.E. Advances and applications of single-cell sequencing
601 technologies. *Mol Cell* **58**, 598-609 (2015).
- 602 47. Ozsolak, F. & Milos, P.M. RNA sequencing: advances, challenges and opportunities.

603 *Nat Rev Genet* **12**, 87-98 (2011).

604 48. Wingender, E., Dietze, P., Karas, H. & Knuppel, R. TRANSFAC: a database on
605 transcription factors and their DNA binding sites. *Nucleic Acids Res* **24**, 238-241
606 (1996).

607 49. Cock, P.J. et al. Biopython: freely available Python tools for computational molecular
608 biology and bioinformatics. *Bioinformatics* **25**, 1422-1423 (2009).

609 50. Sales, G., Calura, E., Cavalieri, D. & Romualdi, C. graphite - a Bioconductor package
610 to convert pathway topology to gene network. *BMC Bioinformatics* **13**, 20 (2012).

611 51. Karen Simonyan, A.Z. Very Deep Convolutional Networks for Large-Scale Image
612 Recognition. (2014).

613 52. Xavier Glorot, A.B., Yoshua Bengio Deep Sparse Rectifier Neural Networks.
614 *Proceedings of the Fourteenth International Conference on Artificial Intelligence and*
615 *Statistics, PMLR*, 315-323 (2011).

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637 **Figure legends**

638

639 **Figure 1 The CNNC architecture**

640 For each gene pair, expression levels from bulk and sc RNA-Seq are transformed into
641 two 32X32 normalized empirical probability function (NEPDF) matrices, and the two
642 are concatenated to form a combined 32X64 NEPDF matrix (left). The combined
643 NEPDF serves as an input to a convolutional neural network (CNN). The intermediate
644 layer of the CNN can be further concatenated with input vectors representing
645 Dnase-seq and PWM data (top). The output layer contains three probability nodes
646 where p_0 is the probability that genes a and b are not interacting, p_1 encodes the
647 case that gene a regulates gene b , and p_2 is the probability that gene b regulates
648 gene a .

649

650 **Figure 2 GTRD TF-target prediction**

651 (a) ROCs for Pearson Correlation (PC) based on scRNA-Seq. Light gray lines
652 represent the performance for each TF. Black line represents the median ROC, and
653 light green region represents the 25~75 quantile. (b) PC for bulk RNA-Seq. (c) PC for
654 combined bulk and scRNA-Seq. (d) ROCs for Mutual Information (MI) when using
655 scRNA-Seq. (e) MI when using bulk RNA-Seq. (f) MI using the combined bulk and
656 scRNA-Seq. (g) ROCs for CNNC using scRNA-Seq. (h) CNNC using bulk RNA-Seq.
657 (i) CNNC for combined data. Inset – Top ranking CNNC pairs are much more likely to
658 be correct pairs when compared to other methods. (j) Comparison of TF-target
659 predictions with additional data. Columns 1-3 show median AUROC of PC, MI, and
660 CNNC using scRNA-Seq, bulk and the combine data, respectively. 4th column shows
661 prediction of TF-gene interactions using only PWM or Dnase. 5th column shows
662 performance when integrating expression, sequence and DNase data.

663

664

665

666 **Figure 3 Predicting undirected pathway edges**

667 (a) Overall ROCs for CNNC performance on KEGG pathway undirected edge
668 prediction with bulk and scRNA-Seq. (b) The Area Under the Receiver Operating
669 Characteristic curve (AUROC) histogram for (a). (c) The overall ROCs for
670 performance of CNNC on Reactome pathway undirected edge prediction with bulk
671 and scRNA-Seq. (d) The AUROC histogram for (c).

672

673 **Figure 4 Directed (causal) edge prediction**

674 (a) Overall ROCs for performance of CNNC on GTRD directed prediction with bulk
675 and scRNA-Seq. (b) The AUROC histogram for (a). (c) Overall ROCs for performance
676 of CNNC on KEGG pathway directed edge prediction with bulk and scRNA-Seq. (d)
677 The AUROC histogram for (c). (e) Overall ROCs for performance of CNNC on
678 Reactome pathway directed edge prediction with bulk and scRNA-Seq. (f) The

679 AUROC histogram for (e). (g) A typical NEPDF sample from a KEGG interaction that
680 is correctly predicted as label 1. (h) A typical NEPDF sample that is correctly predicted
681 as label 2.

682

683 **Figure 5 Downstream applications using CNNC**

684 CNNC can be used as a component in downstream analysis algorithms including for
685 pathway analysis and clustering. (a) Top: Directed edges annotated in KEGG for the
686 IL-17 pathway Bottom: predicted directed edges for the pathway. (b) Hierarchical
687 clustering of top 2,000 DE genes based on CNNC similarity matrix score. The number
688 under the horizontal line represents the distance between the two groups, and the
689 black horizontal line shows the resulting eight clusters groups. (c) Hierarchical
690 clustering based using PC as the input. (d) GO term analysis of the clustering results
691 from (Figs. 5b, 5c, Supplementary Figs. 8a and 8b).

692









