

DeepGRN: Prediction of transcription factor binding site across cell-types using attention-based deep neural networks

Chen Chen¹, Jie Hou², Xiaowen Shi³, Hua Yang³, James A. Birchler³, and Jianlin Cheng^{1,*}

¹Electrical Engineering and Computer Science Department, University of Missouri, Columbia, MO 65211, USA

²Department of Computer Science, Saint Louis University, St. Louis, MO, 63103, USA.

³Division of Biological Sciences, University of Missouri, Columbia, MO 65211, USA.

* To whom correspondence should be addressed.

Email addresses:

CC: ccm3x@mail.missouri.edu

JH: jie.hou@slu.edu

XS: shix@missouri.edu

HY: yanghu@missouri.edu

JB: birchlerj@missouri.edu

JC: chengji@missouri.edu

1 **Abstract**

2 **Background**

3 Due to the complexity of the biological systems, the prediction of the potential DNA
4 binding sites for transcription factors remains a difficult problem in computational
5 biology. Genomic DNA sequences and experimental results from parallel sequencing
6 provide available information about the affinity and accessibility of genome and are
7 commonly used features in binding sites prediction. The attention mechanism in deep
8 learning has shown its capability to learn long-range dependencies from sequential

9 data, such as sentences and voices. Until now, no study has applied this approach in
10 binding site inference from massively parallel sequencing data. The successful
11 applications of attention mechanism in similar input contexts motivate us to build and
12 test new methods that can accurately determine the binding sites of transcription
13 factors.

14 **Results**

15 In this study, we propose a novel tool (named DeepGRN) for transcription factors
16 binding site prediction based on the combination of two components: single attention
17 module and pairwise attention module. The performance of our methods is evaluated
18 on the ENCODE-DREAM in vivo Transcription Factor Binding Site Prediction
19 Challenge datasets. The results show that DeepGRN achieves higher unified scores
20 in 6 of 13 targets than any of the top four methods in the DREAM challenge. We also
21 demonstrate that the attention weights learned by the model are correlated with
22 potential informative inputs, such as DNase-Seq coverage and motifs, which provide
23 possible explanations for the predictive improvements in DeepGRN.

24 **Conclusions**

25 DeepGRN can automatically and effectively predict transcription factor binding sites
26 from DNA sequences and DNase-Seq coverage. Furthermore, the visualization
27 techniques we developed for the attention modules help to interpret how critical
28 patterns from different types of input features are recognized by our model.

29 **Keywords**

30 Transcription factor, Attention mechanism, DNA binding site prediction.

31

32 **Background**

33 Transcription factors (TFs) are proteins that bind to specific genomic sequences and
34 affect numerous cellular processes. They regulate the rates of transcriptional
35 activities of downstream genes through such binding events, thus acting as activators
36 or repressors in the gene regulatory networks by controlling the expression level and
37 the protein abundance of their targeted genes [1]. Chromatin immunoprecipitation-
38 sequencing (ChIP-Seq) is the golden standard to determine the interactions of a TF
39 and all its potential binding regions on genomic sequences. However, ChIP-Seq
40 experiments usually require reagents and materials that are infeasible to acquire,
41 such as antibodies targeting specific TF of interest. Thus, predictions of potential
42 binding sites through computational methods are considered as alternative solutions.
43 Also, the prediction of binding sites of TFs would facilitate many biological studies by
44 providing resources as reference for experimental validation.

45 Many algorithms have been developed to infer the potential binding sites of different
46 TFs, including hidden Markov models [2, 3], hierarchical mixture models [4], support
47 vector machines [5, 6], discriminative maximum conditional likelihood [7] and random
48 forest [8, 9]. These methods usually rely on prior knowledge about sequence
49 preference, such as position weight matrix [10]. However, these features may be less
50 reliable if they are generated from inference based methods (such as de-novo motif
51 discovery) when no prior knowledge is available [7].

52 More recently, methods based on deep neural networks (DNNs), such as DeepBind,
53 TFImpute, and DeepSEA, have shown performances superior to traditional models
54 [11-13]. Compared with the conventional methods, deep learning models have their
55 advantages at learning high-level features from data with huge sizes. This property
56 makes them ideal for the binding site prediction task since a genome-wide binding
57 profile of a TF can be generated from each ChIP-Seq experiment. Unlike many
58 existing models that rely on the quality of the input data and labor-intensive feature

59 engineering, deep learning requires less domain knowledge or data pre-processing
60 and is more powerful when there is little or no prior knowledge of potential binding
61 regions. Current studies in the protein binding site prediction tasks usually involve the
62 combination of two deep learning architectures: convolutional neural networks (CNN)
63 and recurrent neural networks (RNN). The convolutional layer has the potential to
64 extract local features from different genomic signals and regions [14], while the
65 recurrent layer is better at utilizing useful information across the entire sequences of
66 data. Several popular methods for binding prediction, such as DanQ [15],
67 DeeperBind [16], and FactorNet [17], are built on such model architecture.

68 Recently, the concept of attention mechanism has achieved great success in neural
69 machine translation [18] and sentiment analysis [19]. It enhances the ability of DNNs
70 by focusing on the information that is highly valuable to successful prediction.
71 Combining with RNNs, it allows models to learn the high-level representations of
72 input sequences with long-range dependencies. For example, long short-term
73 memory (LSTM) models with attention mechanism have been proposed in relation
74 classification [20] and sentence compression [21]. Because of the input context
75 similarities between language processing (sentences) and the DNA binding site
76 prediction (sequences and results from massively parallel sequencing), similar
77 approaches can be applied improve the performance of existing methods [22-24].

78 Interrogating the input-output relationships for complex models is another important
79 task in machine learning. The weights of a deep neural network are usually difficult to
80 interpret directly due to their redundancy and nonlinear relationship with the output.
81 Saliency maps and feature importance scores are conventional approaches for
82 model interpretation in machine learning involving genomics data [25]. With the
83 application of attention mechanism, we are also interested in testing its ability to
84 enhance the interpretability of existing CNN-RNN architecture models.

85 In this paper, we develop a TF binding prediction tool (DeepGRN) that is based on
86 deep learning with attention mechanism. The experimental results demonstrate that

87 our approach is competitive among the current state-of-the-art methods. Also, our
88 work can be extended to explain the input-output relationships through the learning
89 process. We show that the utilization of informative patterns in both DNase-Seq and
90 DNA sequences is important for accurate prediction.

91 **Implementation**

92 **Datasets from ENCODE-DREAM Challenge**

93 The datasets used for model training and benchmarking are from the 2016
94 ENCODE-DREAM *in vivo* Transcription Factor Binding Site Prediction Challenge.
95 The detailed description of the pre-processing of the data can be found at
96 <https://www.synapse.org/#!Synapse:syn6131484/>.

97 For all TF and cell-types provided in the challenge datasets, the label of the binding
98 status of the TFs is generated from ChIP-Seq experiments and used as ground truth.
99 Chromatin accessibility information (DNase-Seq data), and RNA-Seq data are
100 provided as input features for model training.

101 For model training, we follow the rules and restrictions of the DREAM challenge: the
102 models are trained on all chromosomes except 1, 8, and 21, and chromosome 11 is
103 used as validation. The model with the best performance in validation data is used for
104 final prediction if no “leaderboard” dataset is provided by the challenge. The
105 leaderboard data are available for some TFs for benchmarking, and each participant
106 can test the performance on these TFs with up to ten submissions. Thus, if such data
107 are provided, we pick the top 10 best models from the first step as an optional model
108 selection step. The final performance of our models is reported based on the final
109 test data that are used to determine the rank of the submissions in the challenge
110 (Figure S1 and Table S1, see Additional file 1). We use the similar organization of
111 input features introduced by FactorNet [17]: DNA Primary sequence, Chromatin
112 accessibility information (DNase-Seq data) are transformed into sequential features

113 and become the input of the convolution layers at the first part of the models. Gene
114 expression and annotations are transformed into non-sequential features and feed
115 into the intermediate dense layers of the model (Details are described in the “Deep
116 neural network models with attention modules” section).

117 We also collected DNase and ChIP profiles for additional cell lines from the Encode
118 Project (<https://www.encodeproject.org>) and Roadmap Epigenomics databases
119 (<http://www.roadmapepigenomics.org/data/>) to improve the capability of
120 generalization of our model. The performance of models trained with and without
121 external datasets are evaluated separately.

122 **Transcription factor binding data**

123 Transcription factor binding data from ChIP-Seq experiments is the target for our
124 prediction. The whole genome is divided into bins of 200bp with a sliding step size of
125 50bp (i.e., 250-450bp, 300-500bp). Each bin falls into one of the three types: bound,
126 unbound, or ambiguous, which is determined from the ChIP-Seq results. Bins
127 overlapping with peaks and passing the Irreproducible Discovery Rate (IDR) check
128 with a threshold of 5% [26] are labeled as bound. Bins that overlap with peaks but fail
129 to pass the reproducibility threshold are labeled as ambiguous. All other bins are
130 labeled as unbound. We do not use any ambiguous bins during the training or
131 validation process according to the common practice. Therefore, each bin in the
132 genomic sequence will either be a positive site (bounded) or a negative site
133 (unbounded).

134 **DNA primary sequence**

135 Human genome release hg19/GRCh37 is used as the reference genome. In
136 concordance with the common practice of algorithms that perform feature extraction
137 from chromatin profile, such as FactorNet[17], DeepSea[12], and DanQ[27], we

138 expand each bin by 400bp in both upstream and downstream, resulting in a 1000bp
139 input region. In addition, we have evaluated the performance of different selections of
140 input ranges and showed that range above 600bp is sufficient to acquire stable
141 prediction performance (Figure S2). The sequence of this region is represented by a
142 1000×4 bit matrix by 1-hot encoding, with each row represented a nucleotide. Since
143 low mappability sequences may introduce bias in parallel sequencing experiments,
144 sequence uniqueness (also known as “mappability”) is closely related to the quality of
145 sequencing data [28]. Thus, we select Duke 35bp uniqueness score
146 (<https://genome.ucsc.edu/cgi-bin/hgFileUi?db=hg19&g=wgEncodeMapability>) as an
147 extra feature. Scores ranging from 0 to 1 are assigned to each position as the inverse
148 of occurrences of a sequence with the exceptions that the scores of unique
149 sequences are 1 and scores of sequences occurring more than four times are 0 [29].
150 As a result, the sequence uniqueness is represented by a 1000×1 vector for each
151 input bin. The ENCODE Project Consortium has provided a blacklist of genomic
152 regions that produce artifact signals in NGS experiments [30]. We exclude input
153 bins overlapping with these regions from training data and set their prediction scores
154 to 0 automatically if they are in target regions of prediction.

155 **DNase-Seq data**

156 Chromatin accessibility refers to the accessibility of regions on a chromosome and is
157 highly correlated with TF binding events [4]. DNase-Seq experiment can be used to
158 obtain genome-wide maps of chromatin accessibility information as chromatin
159 accessible regions are usually more sensitive to the endonuclease DNase-I than
160 non-accessible regions [31]. DNase-Seq results for all cell-types are provided in the
161 Challenge datasets in the BigWig format. Normalized 1x coverage score is generated
162 from the BAM files using deepTools [32] with bin size = 1 and is represented by a
163 1000×1 vector for each input bin.

164 **Gene expression and annotation**

165 The annotation feature for each bin is encoded as a binary vector of length 6, with
166 each value represent if there is an overlap between the input bin and each of the six
167 genomic features (coding regions, intron, promoter, 5'/3'-UTR, and CpG island). We
168 also include RNA-Seq data since they can be used to characterize the differences in
169 gene expression levels among different cell-types. Principal Component Analysis
170 (PCA) is performed on the Transcripts per Million (TPM) normalized counts from
171 RNA-Seq data of all cell-types provided by the Challenge. The first eight principal
172 components of a cell-type are used as expression scores for all inputs from that cell-
173 type, generating a vector of length 8. The processed data files for these features are
174 provided in the FactorNet Repository ([https://github.com/uci-
175 cbcl/FactorNet/tree/master/resources](https://github.com/uci-cbcl/FactorNet/tree/master/resources)). These non-sequential features are fused into
176 the first dense layer in the model.

177 **PhastCons Genome Conservation tracks**

178 We use the 100-way PhastCons conservation tracks [33] as a feature for additional
179 models. The PhastCons scores are represented as base-by-base conservation
180 scores generated from multiple alignments of 99 vertebrates to the human genome.
181 Conserved elements along the genome are recognized from phylogenetic models,
182 and the conservation score for each base is computed as the probability that it
183 locates in such conserved regions. For each input bin, the PhastCons scores are
184 represented as a vector of $L \times 1$ with a range from 0 to 1.

185 **CpG island feature profiling**

186 We use the CGI score derived from Mocap [34] to profile the epigenomic
187 environment for each input region. The CGI score can be calculated as:

$$188 \quad CGI(N_{CPG}, N_C, N_G, L) = \begin{cases} 1 & \text{if } \frac{N_{CPG}L}{((N_C + N_G)/2)^2} > 0.6 \text{ and } \frac{N_C + N_G}{L} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

189 For each input bin, the CGI scores are represented as a vector of $L \times 1$ with binary
190 values of 0 or 1.

191 **Deep neural network models with attention modules**

192 The shape of each sequential input is $L \times (4+1+1)$ for each region with length L after
193 combining all sequential features (DNA sequence, sequence uniqueness, and
194 Chromatin accessibility). Sequential inputs are generated for both the forward strand
195 and the reverse complement strand. The weights in all layers of the model are
196 shared between both inputs to form a “Siamese” architecture [17, 35, 36]. Vectors of
197 non-sequential features from gene expression data and genomic annotations are
198 fused into the model at the first dense layer. The overall architecture of our model is
199 shown in Figure 1. The model is built with two major modules: single attention and
200 pairwise attention. They use the same input and architecture except for their internal
201 attention mechanism. The final result of our model is the average of the output of two
202 modules.

203 The first part of our model is a 1D convolutional layer, which is a common practice for
204 feature extraction in deep learning models involving genomics data [13, 17]. We use
205 Bidirectional Long Short-term Memory (Bi-LSTM) nodes as recurrent units in our
206 model. The computation steps in an LSTM unit can be written as:

$$207 \quad f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (1)$$

$$208 \quad i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (2)$$

$$209 \quad \tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \quad (3)$$

$$210 \quad \tilde{C}_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (4)$$

$$211 \quad o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (5)$$

$$212 \quad h_t = o_t * \tanh(\tilde{C}_t) \quad (6)$$

213 Where f_t , i_t , and o_t are the forget gate, input gate, and output gate. h_{t-1} and h_t
214 are the hidden state vectors at position $t-1$ and t . x_t is the input vector at
215 position t . $[h_{t-1}, x_t]$ stands for vector concatenation operation. C_{t-1} , \tilde{C}_t and C_t are
216 output cell state at position $t-1$, new cell state at position t , and output cell state at
217 position t , respectively. W_f , W_i , W_C , and W_o are learned weight matrices. b_f , b_i ,
218 b_C , and b_o are learned bias vector parameters for each gate. σ and \tanh are
219 sigmoid function and hyperbolic tangent function, respectively.

220 In Bi-LSTM layers, two copies of the inputs of LSTM are rearranged into two
221 directions: one for the forward direction and one for the backward direction, and they
222 go into the LSTM unit separately. The outputs from two directions are concatenated
223 at the last dimension. Thus, the last dimension of the Bi-LSTM output is two times of
224 the last dimension of the input.

225 In the single attention module, suppose its input vector h has shape l by r , we first
226 computed the unnormalized attention score $e = M \times h$ where M is a weight
227 matrix with shape l by l , and e has shape l by r . A learned bias of shape l by r
228 is added to e after the multiplication. This can be summarized as a dense layer
229 operation $f_{att,r}$ on input h . Then, we apply the Softmax function along the first
230 dimension of e in order to get the normalized attention score α . Finally, the
231 weighted output Z will be computed based on the attention weight α . At dimension
232 r of input h , these steps can be written as follows:

$$233 \quad e_r = f_{att,r}(h_{1,r}, h_{2,r}, \dots, h_{N,r}) \quad (7)$$

$$234 \quad \alpha_{i,r} = \exp(e_{i,r}) / \sum_{k=1}^N \exp(e_{k,r}) \quad (8)$$

$$235 \quad \alpha_i = (\sum_{r=1}^R \alpha_{i,r}) / D \quad (9)$$

$$236 \quad z_{i,r} = h_{i,r} * \alpha_i \quad (10)$$

237 Here, e_r is the unnormalized attention score at dimension r . Vector $\alpha_{i,r}$ represents
238 attention weight at dimension r of position i and is normalized by Softmax function.
239 The attention dimension r in our model will stay unchanged during the

240 transformations. The dimension of the attention weights can be reduced from $N \times r$
 241 to $N \times 1$ by averaging at each position. The final output $z_{i,r}$ is computed based on
 242 the corresponding attention score. After the attention layers, the prediction scores are
 243 computed from dense layers with sigmoid activation function and merged from both
 244 forward and reverse complement inputs.

245 In the pairwise attention module, there are three components: Q(query), K(key) and
 246 V(value). Their values are computed from LSTM output from three different trainable
 247 weight matrices. The dimension of the trained weights for Q, K and V are l by d_k , l
 248 by d_k and l by d_v where d_k and d_v are set as 64 as the default setup described
 249 in [37]. The multiplication of Q and transpose of K are used to compute the attention
 250 weights for each position of V after Softmax conversion and dimension normalization.
 251 The multiplication of V and attention weights are the output of the pairwise attention
 252 module. The output of the pairwise attention module is computed as:

$$253 \quad Z = \text{Softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \times V \quad (11)$$

254 Since each position in the sequential features simultaneously flows through the
 255 pairwise attention module, the pairwise attention module itself is not able to sense
 256 the position and order from the sequential input. To address this, we add the
 257 positional encodings to the input of the pairwise attention. We expect this additional
 258 encoding will enhance the ability of the model to make use of the order of the
 259 sequence. The positional encodings have the same dimension d as the input of the
 260 pairwise attention module. In this work, we choose different frequencies sine and
 261 cosine functions [38] to encode the positional information:

$$262 \quad PE_{(pos,2i)} = \sin(pos/10000^{2i/d}) \quad (12)$$

$$263 \quad PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d}) \quad (13)$$

264 where pos is the position in the sequential input, and i is the index of the last
 265 dimension of the model. The resulting positional encodings vector is added to its
 266 input. Through such encoding technique, the relative position information can be

267 learned by the model since for any fixed offset k , $PE_{(pos+k)}$ can be represented as
268 $PE_{(pos,2i)}\cos(10000^{2k/d}) + PE_{(pos,2i+1)}\sin(10000^{2k/d})$, which is the linear
269 combination of $PE_{(pos)}$. Similarly, this also applies to dimensions of $2i + 1$ as well.
270 The single attention module is designed to represent the importance of different
271 regions along with the sequential input, while the pairwise attention module seeks to
272 attend the importance between each pair of positions across the sequential input. We
273 expect this difference in architecture will help to improve the learning ability of the
274 model in a complementary manner.
275 We tested different configurations for typical hyperparameters (learning rate, network
276 depth, dropout rates) and the hyperparameters specific to our model (the dimension
277 of attention weights, merging function the two output scores) during training. The
278 complete description of hyperparameters and their possible options are summarized
279 in Table S2 [see Additional file 1]. We train one model for each TF, resulting in 12
280 models in total. The single and pairwise attention module will always use the same
281 configuration rather than train separately.
282 There are 51676736 bins in total on training chromosomes in the labels, resulting in
283 $51676736 \times n$ potential training samples for each TF, where n is the number of
284 available cell-types for training. Due to limited computing capacity, we use the
285 iterative training process. During training, the training data is the mixture of all
286 positives (labeled as “B”) with downsampled negatives (labeled as “U”) [17]. In
287 the traditional model training in deep learning, all input data are used to update the
288 model weights exactly once for each epoch. However, this is not applicable in our
289 task since the negative samples (regions do not bind to TFs) are much more
290 abundant than the positive samples (regions bind to TFs), and use all negative
291 samples for training in one epoch is not practical since the number of them is
292 extremely huge (as they cover most of the human genome). Thus, in each epoch
293 during model training, we first sample negative samples with numbers proportional to

294 the number of all positive samples, and combine these negative samples with all
295 positive samples for training. We will re-sample the negative bins and start another
296 round of model training (next epoch). To make the training process more effective,
297 we use a different strategy to generate positive training samples for transcription
298 factors that have a large number of positive labels (CTCF, FOXA1, HNF4A, MAX,
299 REST and JUND). For these TFs, we randomly sample a 200-bp region from each
300 ChIP-Seq peak in the narrow peak data as positive instances for training instead of
301 using all positive samples in each epoch. We use Adam [39] the optimizer, and
302 binary cross-entropy as the loss function. The default number of epochs is set to 60,
303 but the training will be early stopped if there are no improvements in validation
304 auPRC for five consecutive epochs. For detailed instructions about data retrieving,
305 training, prediction, and visualization with our programs, please see Additional file 2.

306 **Results**

307 **Overall benchmarking on evaluation data**

308 We list the performance of our model as four metrics used in the DREAM Challenge
309 (Table 1) and compare them with the unified score from the top four teams in the final
310 leaderboard of the ENCODE-DREAM Challenge (Table 2). The unified score for
311 each TF and cell-type is based on the rank of each metric and is computed as:
312 $\sum \ln(r/(6))$ where r is the rank of the method for one specific performance measure
313 (auROC, auPRC, Recall at 50% FDR and Recall at 10% FDR). Thus, smaller scores
314 indicate better performance. The TFs, chromosomes, and cell-types for evaluation
315 are the same as those used for the final rankings. DeepGRN typically achieves
316 auROC scores above 98% for most of the TF/cell type pairs, reaching as low
317 as 97.1% for HNF4A/liver. The scores of auPRC have a more extensive range of
318 values, from 40.4% for E2F1/ K562 to 90.2% for CTCF/iPSC.

319

320 **Table 1. The performance of DeepGRN with four metrics used in the DREAM**
 321 **Challenge.**

TF Name	Cell-type	auROC	auPRC	Recall at 50% FDR	Recall at 10% FDR
CTCF	PC-3	0.987	0.767	0.766	0.603
CTCF	induced pluripotent stem cell	0.998	0.902	0.945	0.744
E2F1	K562	0.989	0.404	0.388	0.100
EGR1	liver	0.993	0.405	0.318	0.021
FOXA1	liver	0.985	0.546	0.584	0.164
FOXA2	liver	0.984	0.548	0.588	0.143
GABPA	liver	0.991	0.516	0.488	0.154
HNF4A	liver	0.971	0.636	0.700	0.263
JUND	liver	0.983	0.535	0.585	0.027
MAX	liver	0.990	0.425	0.349	0.004
NANOG	induced pluripotent stem cell	0.996	0.499	0.515	0.035
REST	liver	0.986	0.482	0.527	0.030
TAF1	liver	0.989	0.424	0.393	0.000

322

323 For each TF and cell-type combination, our attention model has better performance
 324 on 69% (9/13) of the prediction targets than Anchor [40], 85% (11/13) than FactorNet
 325 [17], 85% (11/13) than Cheburashka [7], and 77% (10/13) than Catchitt [41]. Among
 326 all methods benchmarked, our method has the highest ranking in 7 out of 13 targets
 327 (CTCF/iPSC, FOXA1/liver, FOXA2/liver, GABPA/liver, HNF4A/liver, NANOG/iPSC,
 328 and REST/liver), with the best average score (0.31) across all TF/ cell-types pairs
 329 (Table 2).

330 **Table 2. The unified scores of DeepGRN and the top four algorithms in the**
 331 **DREAM Challenge.**

TF	cell	Anchor	FactorNet	Cheburashka	Catchitt	DeepGRN
CTCF	PC-3	0.67	0.17	0.83	0.5	0.33
CTCF	induced pluripotent stem cell	0.83	0.33	0.67	0.5	0.17
E2F1	K562	0.5	0.83	0.67	0.17	0.33
EGR1	liver	0.17	0.83	0.67	0.33	0.5
FOXA1	liver	0.67	0.33	0.83	0.5	0.17
FOXA2	liver	0.33	0.83	0.67	0.5	0.17
GABPA	liver	0.33	0.83	0.67	0.5	0.17
HNF4A	liver	0.67	0.33	0.83	0.5	0.17
JUND	liver	0.17	0.83	0.67	0.5	0.33
MAX	liver	0.17	0.83	0.33	0.67	0.5
NANOG	induced pluripotent stem cell	0.33	0.5	0.83	0.67	0.17
REST	liver	0.67	0.33	0.83	0.5	0.17
TAF1	liver	0.17	0.5	0.67	0.33	0.83

332 To precisely evaluate the capability of deepGRN under the restrictions of the
333 ENCODE DREAM Challenge, we also compared the performance of deepGRN
334 trained using datasets provided by the challenge with four available features:
335 Genomic sequence features, DNase-Seq and RNA-Seq data. The results are
336 summarized in Table S3 and S4. DeepGRN still achieves the highest ranking in 6 out
337 of 13 targets, with the best average unified score (0.33) across all targets.

338 **Performance comparison between two attention modules**

339 In addition to the comparisons with the top 4 methods in the challenge, we also
340 benchmarked the individual performance of the single and pairwise attention module
341 (Table S5, see Additional file 1). In general, the results extracted from the single
342 attention module have similar performances. For all 13 TF and cell-type pairs, the
343 single attention module has higher auROC in 6 targets while the pairwise attention
344 module has higher auROC in 3 targets. The rest of the targets are tied. The final
345 output of the model is the ensemble of these two modules by averaging, and it
346 outperforms any of the individual attention modules in 10 of 13 targets (Table 1). The
347 largest improvements from ensemble (as auPRC) come from FOXA2 (0.34) , REST
348 (0.09) and FOXA1 (0.09). We also found that the performance of the two attention
349 modules have the same trend across all TF and cell-types in all four performance
350 measures (Figure 2), suggesting that the capability of learning from features are
351 coherent between the two modules.

352 We evaluated the importance of each feature between single and pairwise attention
353 mechanism. For the prediction of each target, we set the values of each sequential
354 feature (DNase-Seq, sequence, or uniqueness) to zero, or randomly switch the order
355 of the vector for a non-sequential feature (genomic elements or RNA-Seq). The
356 decrease of auPRC from these new predictions is used as the importance score of
357 each feature (Figure 3). We found that across all TF and cell-types, the sequential
358 features have the largest average importance scores: DNase-Seq (0.36), DNA

359 sequence (0.21), and 35bp uniqueness (0.21) while the scores for other features are
360 much smaller. Similar trends have also been found using individual single and pair
361 attention modules.

362 **Interpretation of attention scores with DNase-Seq and ChIP-Seq**

363 In the single attention module, the output is a weighted sum of the input from the
364 attention layer, and the attention scores are used as weights. These scores
365 characterize a unified mapping between the importance of input feature with its
366 relative position in the sequential input. To analyze the relationship between attention
367 weights and the position of TF binding events, we extract the attention scores from
368 the single attention module for both forward strand and reverse complement strand
369 and compare them with the corresponding normalized ChIP-Seq fold changes in the
370 same region that are predicted as positive (score>0.5). Similarly, we computed the
371 saliency scores for the same input regions (The implementation details are described
372 in Additional file 1). We found that the attention scores on the two DNA strands have
373 a higher correlation ($\rho=0.90$, $\sigma=0.79$) than the saliency scores ($\rho=0.78$, $\sigma=0.51$)
374 (Figure 4a, 4b). Across all TF and cell-type pairs, we found that there is a positive
375 correlation between the attention weights and normalized ChIP-Seq Fold (Figure 4c),
376 and such relationship is not detected globally in saliency scores (Figure 4d). For all
377 TF and cell-types in the benchmark datasets, we select at least four different
378 genomic regions that have a clear ChIP-Seq peak signal in each target for
379 demonstration. We show that the averaged attention weights put more focus on the
380 actual binding region for each cell-type and these focusing points shift along with the
381 shift of TF binding signals (see Additional file 3).

382 Since the accessibility of the genome plays an important role in TF binding, it is
383 expected to find high DNase coverage for those openly accessible areas that can
384 explain the binding event detected by the ChIP-Seq experiment. We run a genome-
385 wide analysis on regions with high DNase-Seq peaks in the single attention module

386 for transcription factor JUND, which is one of the most susceptible targets to DNase-
387 Seq. We illustrate the distribution of normalized DNase coverage values from both
388 the true positives that are false negatives without attention and true negatives that
389 are false positives without attention (Figure 5). The predictions without attention
390 models are generated from the deep learning model in FactorNet
391 (<https://github.com/uci-cbcl/FactorNet/tree/master/models>). The results show that
392 the true positives that only recognized by attention models generally have a smaller
393 DNase peak signal/noise ratio than those recognized by both models, while the
394 negative bins that only recognized by attention models have a larger scale of values.
395 Both two cases will increase the difficulty of correct classification. This observation
396 indicates that the predictive improvements of attention models may result from
397 focusing on more informative DNase-Seq coverage values while ignoring irrelevant
398 regions in negative samples.

399 **Motif detection over high attention scores regions**

400 For those positive samples without distinct DNase-Seq peaks, the patterns of
401 genomic sequences are critical information for successful prediction. To test the
402 ability of attention weights to recognize motifs that contribute to binding events from
403 the genomic sequences, we use an approach similar to DeepBind [13]. For the model
404 trained for each TF, we first acquire the coordinates on the relative positions of
405 maximum column sum of the attention weights from all positive bins in test datasets
406 and extract a subsequence with a length of 20bp around each coordinate. To exclude
407 samples that can be easily classified from patterns of DNase-Seq signal, we only
408 select positive bins that have no significant coverage peaks (ratio between the
409 highest score and average scores < 15). Then we run FIMO [42] to detect known
410 motifs relevant to the TF of the model in the JASPAR database [43]. From the
411 extracted subsequences, we discover motif MA0139.1 (CTCF) in the prediction for
412 CTCF/induced pluripotent cell and MA0148.4 (FOXA1) in the prediction for

413 FOXA1/liver cell. Figure 6a and 6b show the comparison between the sequence logo
414 of the motif rebuilt from the subsequences and the actual known motifs. We also plot
415 the attention scores of the samples that contain these subsequences (Figure 6c, 6f)
416 and the relative location of the regions with detected motifs in FIMO (Figure 6d, 6g).
417 Furthermore, we show that these maximum attention weights do not come from the
418 DNase-Seq peaks near the motif regions by coincidence since no similar pattern is
419 detected from the normalized DNase scores in the same regions (Figure 6e, 6h). We
420 illustrate the similar trends found in the single attention module in Figure S3 [see
421 Additional file 1].

422 **Discussion**

423 The attention mechanism is attractive in various machine learning studies and has
424 achieved superior performance in image caption generation and natural language
425 processing tasks [38, 44]. Recurrent neural network models with attention
426 mechanism are particularly good at tasks with long-range dependency in input data.
427 Inspired by these works, we introduce the attention mechanism to DNN models for
428 TF binding site prediction.

429 The benchmark result using ENCODE-DREAM Challenge datasets shows that the
430 performances of our model are competitive with the current state-of-the-art methods.
431 It is worth mentioning that the DNase-Seq scores are the most critical feature in the
432 attention mechanism from our experiments according to the feature importance
433 analysis. Many prediction tools for binding site prediction before the challenge, such
434 as DeepBind or TFImpute, are not able to utilize the DNase-Seq data and are not as
435 suitable as the four methods that we used for benchmarking in this study. However,
436 the methods we benchmarked in this study share the similar concepts with these
437 existing tools (For example, FactorNet is built with similar architecture as the

438 TFImpute with additional support for the DNase-Seq data) and may reflect the
439 potential of them using the same set of features.

440 The attention weights learned by the models provide an alternative approach to
441 exploring the dependencies between input and output other than saliency maps. By
442 comparing true ChIP-Seq fold change peaks with attention weights, we show how
443 attention weights shift when the fold change peaks move along the DNA sequence.
444 We also demonstrate that our attention model has the ability to learn from known
445 motifs related to specific TFs.

446 Due to the rules of the DREAM Challenge, we only use very limited types of features
447 in this work. However, if more types of features (such as sequence conservation or
448 epigenetic modifications) are available, they can possibly be transformed into
449 sequential formats and may further improve the prediction performance through our
450 attention architecture. The attention mechanism itself is also evolving rapidly. For
451 example, the multi-head attention introduced by Transformer [38] showed that high-
452 level features could be learned by attention without relying on any recurrent or
453 convolution layers. We expect that better prediction for the TF binding may also be
454 benefited from these novel deep learning architectures in both accuracy and efficacy.

455 **Conclusions**

456 In this study, we propose a new tool (DeepGRN) that incorporates the attention
457 mechanism with the CNNs-RNNs based architecture. The result shows that the
458 performances of our models are competitive with the top 4 methods in the Challenge
459 leaderboard. We demonstrate that the attention modules in our model help to
460 interpret how critical patterns from different types of input features are recognized.

461 **Abbreviations**

462 **TF:** Transcription Factor

463 **Bi-LSTM:** Bidirectional Long Short-Term Memory

464 **DNase-Seq:** DNase I hypersensitive sites Sequencing

465 **ChIP-Seq:** Chromatin Immunoprecipitation Sequencing

466 **Declarations**

467 **Ethics approval and consent to participate**

468 Not applicable.

469 **Consent for publication**

470 Not applicable.

471 **Availability of data and materials**

472 The datasets used in this study and the source code of DeepGRN are available at

473 <https://github.com/jianlin-cheng/DeepGRN>.

474 **Competing interests**

475 The authors declare they have no conflict of interest.

476 **Funding**

477 This work has been supported by NSF grants (IOS1545780 and DBI1149224) and

478 the U.S. Department of Energy (DOE) grant “Deep Green: Structural and Functional

479 Genomic Characterization of Conserved Unannotated Green Lineage Proteins” (DE-

480 SC0020400). The funders (NSF and DOE) does not play a role in conducting this

481 research.

482 **Authors Contributions**

483 JC conceived of the project. CC and JH designed the experiment. CC implemented
484 the method and gathered the results. CC, JH, XS, HY, and JB wrote the manuscript.
485 All authors edited and approved the manuscript.

486 **Acknowledgements**

487 We wish to thank the organizers of ENCODE-DREAM in vivo Transcription Factor
488 Binding Site Prediction Challenge.

489 **Availability and Requirements**

490 Project name: DeepGRN

491 Project home page: <https://github.com/jianlin-cheng/DeepGRN>

492 Operating system(s): Linux, Mac OS, Windows

493 Programming language: Python, R

494 Other requirements: Python version 3.6.0 or higher, R version 3.3.0 or higher

495 License: GNU GPL

496 Any restrictions to use by non-academics: None

497 **References**

- 498 1. Hobert O: **Gene regulation by transcription factors and microRNAs.**
499 *Science* 2008, **319**(5871):1785-1786.
- 500 2. Mehta P, Schwab D, Sengupta A: **Statistical Mechanics of Transcription-**
501 **Factor Binding Site Discovery Using Hidden Markov Models.** *J Stat Phys*
502 2011, **142**(6):1187-1205.
- 503 3. Mathelier A, Wasserman WW: **The next generation of transcription factor**
504 **binding site prediction.** *PLoS Comput Biol* 2013, **9**(9):e1003214.

- 505 4. Pique-Regi R, Degner JF, Pai AA, Gaffney DJ, Gilad Y, Pritchard JK:
506 **Accurate inference of transcription factor binding from DNA sequence**
507 **and chromatin accessibility data.** *Genome Res* 2011, **21**(3):447-455.
- 508 5. Zhou TY, Shen N, Yang L, Abe N, Horton J, Mann RS, Bussemaker HJ,
509 Gordan R, Rohs R: **Quantitative modeling of transcription factor binding**
510 **specificities using DNA shape.** *P Natl Acad Sci USA* 2015, **112**(15):4654-
511 4659.
- 512 6. Djordjevic M, Sengupta AM, Shraiman BI: **A biophysical approach to**
513 **transcription factor binding site discovery.** *Genome Res* 2003, **13**(11):2381-
514 2390.
- 515 7. Keilwagen J, Posch S, Grau J: **Accurate prediction of cell type-specific**
516 **transcription factor binding.** *Genome Biology* 2019, **20**(1):9.
- 517 8. Xiao Y, Segal MR: **Identification of yeast transcriptional regulation**
518 **networks using multivariate random forests.** *PLoS Comput Biol* 2009,
519 **5**(6):e1000414.
- 520 9. Hooghe B, Broos S, Van Roy F, De Bleser P: **A flexible integrative**
521 **approach based on random forest improves prediction of transcription**
522 **factor binding sites.** *Nucleic acids research* 2012, **40**(14):e106-e106.
- 523 10. Sherwood RI, Hashimoto T, O'Donnell CW, Lewis S, Barkal AA, van Hoff
524 JP, Karun V, Jaakkola T, Gifford DK: **Discovery of directional and**
525 **nondirectional pioneer transcription factors by modeling DNase profile**
526 **magnitude and shape.** *Nat Biotechnol* 2014, **32**(2):171-178.
- 527 11. Zeng H, Edwards MD, Liu G, Gifford DKJB: **Convolutional neural network**
528 **architectures for predicting DNA–protein binding.** *Bioinformatics* 2016,
529 **32**(12):i121-i127.
- 530 12. Zhou J, Troyanskaya OG: **Predicting effects of noncoding variants with**
531 **deep learning-based sequence model.** *Nat Methods* 2015, **12**(10):931-934.
- 532 13. Alipanahi B, DeLong A, Weirauch MT, Frey BJ: **Predicting the sequence**
533 **specificities of DNA- and RNA-binding proteins by deep learning.** *Nat*
534 *Biotechnol* 2015, **33**(8):831-838.
- 535 14. Kalkatawi M, Magana-Mora A, Jankovic B, Bajic VB: **DeepGSR: an**
536 **optimized deep-learning structure for the recognition of genomic signals**
537 **and regions.** *Bioinformatics* 2019, **35**(7):1125-1132.
- 538 15. Quang D, Xie X: **DanQ: a hybrid convolutional and recurrent deep neural**
539 **network for quantifying the function of DNA sequences.** *Nucleic acids*
540 *research* 2016, **44**(11):e107-e107.
- 541 16. Hassanzadeh HR, Wang M: **DeeperBind: Enhancing prediction of sequence**
542 **specificities of DNA binding proteins.** In: *IEEE International Conference on*
543 *Bioinformatics and Biomedicine (BIBM): 2016.* 178-183.
- 544 17. Quang D, Xie X: **FactorNet: A deep learning framework for predicting**
545 **cell type specific transcription factor binding from nucleotide-resolution**
546 **sequential data.** *Methods* 2019.
- 547 18. Luong M-T, Pham H, Manning CD: **Effective approaches to attention-based**
548 **neural machine translation.** In: *Proceedings of the 2015 Conference on*
549 *Empirical Methods in Natural Language Processing: 2015.*
- 550 19. Wang Y, Huang M, Zhao L: **Attention-based lstm for aspect-level**
551 **sentiment classification.** In: *Proceedings of the 2016 conference on empirical*
552 *methods in natural language processing: 2016.* 606-615.
- 553 20. Zhou P, Shi W, Tian J, Qi Z, Li B, Hao H, Xu B: **Attention-Based**
554 **Bidirectional Long Short-Term Memory Networks for Relation**

- 555 **Classification.** In: *aug 2016; Berlin, Germany*. Association for Computational
556 Linguistics: 207-212.
- 557 21. Tran N-T, Luong V-T, Nguyen NL-T, Nghiem M-Q: **Effective attention-**
558 **based neural architectures for sentence compression with bidirectional**
559 **long short-term memory.** In: *Proceedings of the Seventh Symposium on*
560 *Information and Communication Technology; Ho Chi Minh City, Vietnam.*
561 3011111: ACM 2016: 123-130.
- 562 22. Singh R, Lanchantin J, Sekhon A, Qi Y: **Attend and Predict: Understanding**
563 **Gene Regulation by Selective Attention on Chromatin.** *Adv Neural Inf*
564 *Process Syst* 2017, **30**:6785-6795.
- 565 23. Shen Z, Bao W, Huang D-S: **Recurrent Neural Network for Predicting**
566 **Transcription Factor Binding Sites.** *Scientific Reports* 2018, **8**(1):15270.
- 567 24. Park S, Koh Y, Jeon H, Kim H, Yeo Y, Kang J: **Enhancing the**
568 **interpretability of transcription factor binding site prediction using**
569 **attention mechanism.** *Scientific Reports* 2020, **10**(1):13413.
- 570 25. Eraslan G, Avsec Ž, Gagneur J, Theis FJ: **Deep learning: new computational**
571 **modelling techniques for genomics.** *Nature Reviews Genetics* 2019,
572 **20**(7):389-403.
- 573 26. Li QH, Brown JB, Huang HY, Bickel PJ: **Measuring Reproducibility of**
574 **High-Throughput Experiments.** *Ann Appl Stat* 2011, **5**(3):1752-1779.
- 575 27. Quang D, Xie X: **DanQ: a hybrid convolutional and recurrent deep neural**
576 **network for quantifying the function of DNA sequences.** *Nucleic Acids Res*
577 2016, **44**(11):e107.
- 578 28. Sholtis SJ, Noonan JP: **Gene regulation and the origins of human biological**
579 **uniqueness.** *Trends Genet* 2010, **26**(3):110-118.
- 580 29. Derrien T, Estelle J, Marco Sola S, Knowles DG, Raineri E, Guigo R, Ribeca
581 P: **Fast computation and applications of genome mappability.** *PLoS One*
582 2012, **7**(1):e30377.
- 583 30. Consortium EP: **An integrated encyclopedia of DNA elements in the**
584 **human genome.** *Nature* 2012, **489**(7414):57-74.
- 585 31. Madrigal P, Krajewski P: **Current bioinformatic approaches to identify**
586 **DNase I hypersensitive sites and genomic footprints from DNase-seq data.**
587 *Front Genet* 2012, **3**:230.
- 588 32. Ramirez F, Dundar F, Diehl S, Gruning BA, Manke T: **deepTools: a flexible**
589 **platform for exploring deep-sequencing data.** *Nucleic Acids Res* 2014,
590 **42**(Web Server issue):W187-191.
- 591 33. Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A: **Detection of nonneutral**
592 **substitution rates on mammalian phylogenies.** *Genome Res* 2010,
593 **20**(1):110-121.
- 594 34. Chen X, Yu B, Carriero N, Silva C, Bonneau R: **Mocap: large-scale**
595 **inference of transcription factor binding sites from chromatin**
596 **accessibility.** *Nucleic Acids Res* 2017, **45**(8):4315-4329.
- 597 35. Mueller J, Thyagarajan A: **Siamese recurrent architectures for learning**
598 **sentence similarity.** In: *Thirtieth AAAI Conference on Artificial Intelligence:*
599 *2016.*
- 600 36. Qin Q, Feng J: **Imputation for transcription factor binding predictions**
601 **based on deep learning.** *PLoS Comput Biol* 2017, **13**(2):e1005403.
- 602 37. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł,
603 Polosukhin I: **Attention is all you need.** In: *Proceedings of the 31st*

- 604 *International Conference on Neural Information Processing Systems; Long*
605 *Beach, California, USA.* Curran Associates Inc. 2017: 6000–6010.
- 606 38. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł,
607 Polosukhin I: **Attention is all you need.** In: *Advances in Neural Information*
608 *Processing Systems: 2017.* 5998-6008.
- 609 39. Kingma DP, Ba J: **Adam: A method for stochastic optimization.** *CoRR*,
610 **abs/1412.6980.**
- 611 40. Li H, Quang D, Guan Y: **Anchor: trans-cell type prediction of**
612 **transcription factor binding sites.** *Genome Res* 2019, **29(2):281-292.**
- 613 41. **Preselection of training cell types improves prediction of transcription**
614 **factor binding sites**
- 615 42. Bailey TL, Elkan C: **Fitting a mixture model by expectation maximization**
616 **to discover motifs in biopolymers.** *Proc Int Conf Intell Syst Mol Biol* 1994,
617 **2:28-36.**
- 618 43. Khan A, Fornes O, Stigliani A, Gheorghe M, Castro-Mondragon JA, van der
619 Lee R, Bessy A, Cheneby J, Kulkarni SR, Tan G *et al*: **JASPAR 2018:**
620 **update of the open-access database of transcription factor binding**
621 **profiles and its web framework.** *Nucleic Acids Res* 2018, **46(D1):D260-**
622 **D266.**
- 623 44. Yang Z, Yang D, Dyer C, He X, Smola A, Hovy E: **Hierarchical attention**
624 **networks for document classification.** In: *Proceedings of the 2016*
625 *Conference of the North American Chapter of the Association for*
626 *Computational Linguistics: Human Language Technologies: 2016.* 1480-
627 1489.

628 **Figures**

629 **Figure 1 – The general framework of the two attention modules of DeepGRN.**

630 The diagram of the deep neural network architecture. Convolutional and
631 bidirectional LSTM layers use both forward and reverse complement features
632 as inputs. In the single attention module, attention weights are computed from
633 hidden outputs of LSTM and are used to generate the weighted
634 representation through an element-wise multiplication. In the pairwise
635 attention module, three components: Q(query), K(key), and V(value) are
636 computed from LSTM output. The multiplication of Q and transpose of K are
637 used to calculate the attention weights for each position of V. The
638 multiplication of V and attention scores is the output of the pairwise attention
639 module. Outputs from attention layers are flattened and fused with non-
640 sequential features (genomic annotation and gene expression). The final

641 score is computed through dense layers with sigmoid activation and merging
642 of both forward and reverse complement inputs. The dimensions of each
643 layer are shown beside each component.

644 **Figure 2 – Performance comparison between single and pairwise attention**
645 **mechanism.**

646 The performance of each TF and cell-type pairs of the output of the individual
647 module are shown in four measures: (auROC, auPRC, recall at 50% FDR
648 and Recall at 10% FDR). ρ : Pearson Correlation Coefficient, σ : Spearman
649 Correlation Coefficient.

650 **Figure 3 – Importance score of features between single and pairwise attention**
651 **mechanism.**

652 The values represented as the decrease of auPRC without using the specific
653 feature for prediction. The negative value represents an increase of auPRC.

654 **Figure 4 – Analysis of attention weights and saliency scores.**

- 655 (a) Scatterplot of attention weights from positive strand and reverse strand.
656 (b) Scatterplot of saliency scores from positive strand and reverse strand.
657 (c) Scatterplot of ChIP-Seq fold change and mean attention weights from both
658 strands. Z-score transformation is applied to both axes.
659 (d) Distribution of the correlation between attention weights/saliency scores and
660 ChIP-Seq fold change. The dashed line represents the mean of each group.
661 The p-value is calculated using the Wilcoxon signed-rank test. The attention
662 weights and saliency scores on the reverse complement strand are reversed
663 before plotting.
664 ρ : Spearman Correlation Coefficient, σ : Pearson Correlation Coefficient.
665 The correlation between normalized ChIP-Seq Fold change and normalized
666 saliency scores is 0.40 (Spearman) and 0.49 (Pearson).

667 **Figure 5 – Distribution of average normalized DNase coverage values along**
668 **with the inputs of JUND.**

669 Due to a large number of candidates, we only randomly sample 1000
670 samples from each case for plotting. For true positives, only bins that have
671 the largest fold change values in the center are selected to ensure they are
672 aligned. For true negatives, we exclude those in the blacklist as their labels
673 will always be 0.

674 (a) True positives in both models.

675 (b) True positives in the single attention module only.

676 (c) True negatives in both models.

677 (d) True negatives in the single attention module only.

678 **Figure 6 – Comparisons of known motifs and matching motifs learned by**
679 **pairwise attention module in CTCF and FOXA1.**

680 (a) Sequence logo built from subsequences detected in CTCF/induced
681 pluripotent cell prediction (left) and motif MA0139.1/ CTCF (right).

682 (b) The attention scores of the samples selected from CTCF/induced pluripotent
683 cell prediction with hits of MA0139.1/ CTCF in FIMO.

684 (c) The relative positions of the detected motifs in the same region of (b).

685 (d) The normalized DNase-Seq scores in the same region of (b).

686 (e) Sequence logo built from subsequences detected in FOXA1/liver cell
687 prediction (left) and motif MA0148.4/ FOXA1 (right).

688 (f) The attention scores of the samples selected from FOXA1/liver cell prediction
689 with hits of MA0148.4/ FOXA1 in FIMO.

690 (g) The relative positions of the detected motifs in the same region of (f).

691 (h) The normalized DNase-Seq scores in the same region of (f).

692 **Tables**

693 **Table 1 – The performance of DeepGRN with four metrics used in the DREAM**
694 **Challenge.**

695 **Table 2 – The unified scores of DeepGRN and the top four algorithms in the**
696 **DREAM Challenge.**

697 Bold scores denote the TF and cell-types that DeepGRN rank as the highest.

698 **Additional files**

699 **Additional file 1.pdf – Supplementary figures and tables.**

700 Including all supplementary figures and tables referenced in the main text.

701 **Additional file 2.pdf – Instructions of training, prediction, and visualization data**
702 **with DeepGRN.**

703 Including data retrieving, training, prediction with DeepGRN and the
704 implementation details of the visualization scripts used in the main text.

705 **Additional file 3.pdf – Visualization of the relationship between ChIP-Seq peak**
706 **and attention weights.**

707 For each genomic region, the plot on the left represents the attention weights,
708 and the plot on the right represents the enrichment of ChIP-Seq signal fold
709 changes in the same region. Since the lengths of attention weights are
710 reduced by the convolution and pooling layers, their lengths are less than the
711 fold change values. Thus, the plots are aligned on the X-axis to represent the
712 relative position of fold change and averaged attention weights.