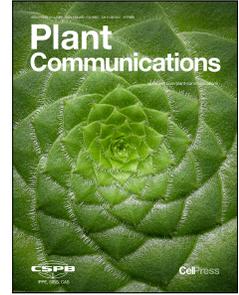


Journal Pre-proof

iCREPCP: a deep learning-based web server for identifying base-resolution *cis*-regulatory elements within plant core promoters

Kaixuan Deng, Qizhe Zhang, Yuxin Hong, Jianbing Yan, Xuehai Hu



PII: S2590-3462(22)00292-9

DOI: <https://doi.org/10.1016/j.xplc.2022.100455>

Reference: XPLC 100455

To appear in: *PLANT COMMUNICATIONS*

Received Date: 31 July 2022

Revised Date: 16 September 2022

Accepted Date: 23 September 2022

Please cite this article as: Deng, K., Zhang, Q., Hong, Y., Yan, J., Hu, X., iCREPCP: a deep learning-based web server for identifying base-resolution *cis*-regulatory elements within plant core promoters, *PLANT COMMUNICATIONS* (2022), doi: <https://doi.org/10.1016/j.xplc.2022.100455>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022

1 **iCREPCP: a deep learning-based web server for identifying base-**
2 **resolution *cis*-regulatory elements within plant core promoters**
3

4 Kaixuan Deng^{1, #}, Qizhe Zhang^{1, #}, Yuxin Hong¹, Jianbing Yan^{2, *} and Xuehai Hu^{1, *}
5

6 ¹College of Informatics, Hubei Engineering Technology Research Center of Agricultural Big Data, Huazhong
7 Agricultural University, Wuhan, Hubei, P.R. of China

8 ²National Key Laboratory of Crop Genetic Improvement, Huazhong Agricultural University, Wuhan 430070, China

9 *Correspondence: Xuehai Hu (huxuehai@mail.hzau.edu.cn) and Jianbing Yan (yjianbing@mail.hzau.edu.cn)

10 # These authors contributed equally to this article.

11 Tel.: +86-18171282783; Fax: +86-27-87288509.
12

13 Dear editor,

14 A central question of plant biology is to specify the temporal and spatial patterns as well as the quantitative
15 level of gene expressions, which are significantly associated with important agronomic traits. There is a
16 growing consensus in the past decade that two key factors determining gene expression level are *cis*-
17 regulatory modules (CRMs) and *trans*-acting factors (TAFs) (Schmitz et al., 2022). Common CRMs
18 include gene-proximal promoters and distal enhancers, which are all considered as the complex assemblies
19 of *cis*-regulatory elements (CREs). It is the binding or interaction between CREs and TAFs (often are
20 transcription factors, TFs) in a ubiquitous or cell-specific manner that determines in which cell, at what
21 time and at what level a gene is expressed. Therefore, the identification of plant CRMs or critical CREs
22 will not only help us understand transcriptional regulatory mechanisms in plants, but also is an essential
23 prerequisite for plant breeding 4.0—breeding by genome editing (Gao, 2021).

24 However, comparing with rich data resources on CREs in mammalian genomes (Fornes et al., 2020),
25 related works in plants have lagged far behind (Schmitz *et al.*, 2022). The bottleneck mainly lies in two
26 aspects: (1) the lack of a big project like ENCODE in plants makes epigenomic features absent or
27 fragmented, leading only a handful of putative plant CREs from genome-wide identification; (2) too few
28 transient transfection systems (only two of protoplasts and tobacco leaves (Jores et al., 2021)), together
29 with difficult validation assays like self-transcribing active regulatory region-sequencing (STARR-seq) in
30 plants, make fewer experimental-validated CREs.

31 **Plant Core Promoter (PCP)**, with the minimal sequence region of 50-100bp around transcription start
32 site (TSS), is a large group of CRMs that are rich in CREs, and can drive basal level of target gene
33 transcriptions (Schmitz *et al.*, 2022). The promoter strength of PCP is defined as the ability to drive
34 expression of a barcoded green fluorescent protein (GFP) reporter gene via transient transfection systems.
35 To our best knowledge, there is no existing computational tool for identifying CREs within PCPs. Here,
36 we developed a deep learning-based web server (<http://www.hzau-hulab.com/icrepcp/>) to identify which
37 **CREs** a given **Plant Core Promoter (iCREPCP)** contains, with a focus on base-resolution position of
38 each CRE and its contribution to the promoter strength.

39 We first downloaded a large-scale PCP dataset of 18,329 *Arabidopsis*, 34,415 maize and 27,094
40 sorghum core promoters, whose strengths were measured by STARR-seq assays in six transient
41 transfection systems (tobacco leaves with enhancer in dark, tobacco leaves without enhancer in dark,
42 tobacco leaves with enhancer in light, tobacco leaves without enhancer in light, maize protoplasts with
43 enhancer in dark and maize protoplasts without enhancer in dark) (Jores *et al.*, 2021). We will take

44 ‘sequence’ as input and ‘enrichment’ as output of a total of about 76,000 samples from all three species
45 for training and testing deep learning models.

46 We next trained a deep learning architecture of ‘DenseNet’ (Huang, 2017) to fit promoter strengths with
47 their DNA sequences. DenseNet has won the best paper award of CVPR-2017, and it can alleviate the
48 vanishing-gradient problem (Figure 1A, Supplementary Information). As expected, **iCREPCP** can
49 accurately fit the experimental results in all six transfection systems: the mean training R^2 ranges from
50 0.490 to 0.782, and all models have low variances, implying their feasibilities (Figure 1B). We next
51 investigate its generalizability by an independent testing dataset (Supplementary Information).
52 Remarkably, **iCREPCP** achieves good testing R^2 ranging from 0.420 to 0.752 and obviously improves
53 the previous work who employed a simple convolutional neural network (Jores *et al.*, 2021) (Figure 1B),
54 also implying its strong generalizability. Moreover, the small differences between training R^2 and testing
55 R^2 (ranging from 0.03 to 0.07) demonstrate that **iCREPCP** have little problems of overfitting, further
56 suggesting that they have potential transfer abilities for other plant species.

57 To investigate the biological interpretability and practicability of **iCREPCP**, we here are more
58 concerned on the contribution of each base during the promoter strength prediction of PCP rather than the
59 prediction accuracy. Because several successive bases having high contributions are potential critical
60 CREs, which are ideal targets of genome editing engineering (Gao, 2021). To this end, we employed a
61 powerful interpretability tool of DeepLIFT (Shrikumar, 2017) to assign a DeepLIFT contribution score to
62 each base of a given PCP. We employed two known PCP examples of maize YIGE1 gene and rice IPA1
63 gene for demonstrating the detecting power of **iCREPCP** together with DeepLIFT (DeepLIFT
64 contribution scores are visualized as high characters with colors that help readers easily find critical bases).
65 YIGE1 is a newly-reported maize gene contributing to ear length and grain yield, and a single-nucleotide
66 polymorphism (SNP) located in its regulatory region had a large effect on its promoter strength (Luo *et*
67 *al.*, 2022). Using the trained model of tobacco leaves without enhancer in light, **iCREPCP** successfully
68 located a large-contributed regulatory region flanking the important SNP (also repeatedly detected by two
69 additional interpretability tools of in-silico tilling deletion and in-silico mutagenesis, Figure 1C),
70 suggesting its detecting power. For trans-species circumstance, IPA1 is a rice star gene that is a master
71 regulator of rice plant architecture. Its’ function was known to increase grains per panicle but reduce tillers,
72 but a recent breakthrough reported that a 54-base pair *cis*-regulatory deletion can both increase grains per
73 panicle and tiller number (Song *et al.*, 2022). Surprisingly, **iCREPCP** successfully detected a 12bp region
74 (-128~-117) with large contributions that exactly covers the An-1 binding site within the deletion (Figure

75 1D, Supplementary Figure 1), implying that **iCREPCP** has great potentials for trans-species
76 identifications of base-resolution critical CREs.

77 For a rough estimation of precision and recall of iCREPCP, we constructed a benchmark of *Arabidopsis*
78 CREs, which was employed for an evaluation: precision and recall are 0.447 and 0.344 respectively
79 (Supplementary Figure 3 and Supplementary materials).

80 To investigate biological implications of several successive bases with high DeepLIFT contribution
81 scores, we next naturally ask whether they are TF motifs and then employed a new motif discovery
82 algorithm of TF-MoDISco (Shrikumar, 2018), that was specifically developed for deep learning, to
83 identify high-quality, non-redundant TF motifs within PCPs (Supplementary Information). For the trained
84 model of tobacco leaves without enhancer in light, TF-MoDISco totally identified 21 clustered seqlets, 14
85 out of which have perfect matching in JASPAR database (Figure 1E, Supplementary Figure 2 and Table
86 1). To further quantify the population-level effect size of the 14 enriched TF motifs, we performed a global
87 importance analysis (Koo et al., 2021) and found that 8 (including TATATA motif, TCP8 and AP1) out
88 of 14 have positive global importance, whereas 6 (including ERF3 and ABI3) out of 14 have negative
89 effects (Figure 1F). Finally, we scanned all 75,375 PCPs using the 14 PWMs of the enriched TF motifs
90 and gave a comprehensive statistic about their occurrence numbers in each PCP sample (Figure 1G,
91 Supplementary Table 2). Notably, the TATATA motif has the most occurrence numbers within PCPs
92 having large promoter strengths in all three species, whereas the ERF3 motif has more occurrences within
93 PCPs having small promoter strengths in both sorghum and maize, which is consistent with their global
94 importance analysis results.

95 In summary, iCREPCP (Figure 1H) provides a user-friendly platform to identify critical CREs that
96 importantly contribute to the promoter strength of any given PCPs with base resolution. These resources,
97 including the six trained prediction models and a powerful visualization tool, will greatly help plant
98 scientists in at least two respects: (i) easily obtain an accurate prediction value of the promoter strength
99 with the only need of the 170bp DNA sequence around its TSS; (ii) precisely detect the base-resolution
100 position of each CRE and its contribution to the promoter strength. The later function will provide
101 important candidate targets of genome editing and will be of general interests in the plant community. The
102 main limitation of iCREPCP is that it was trained with promoter strength measured in vitro via tobacco
103 leaves or maize protoplasts, implying that iCREPCP might work not well on some genes needing distinct
104 expression pattern in vivo. Another limitation is that the prediction accuracy is sensitive to the boundary
105 (Supplementary Figure 4), imply that our models only can be used on (-165, +5) of TSS. Further

106 improvements of iCREPCP will focus on the accurate identification of distal CREs: (i) take longer
 107 genomic sequences as the inputs, covering more distal CREs (such as enhancers) influencing gene
 108 expressions; (ii) develop more sophisticated models for capturing long-range dependency information.

109 **Data availability**

110 The datasets and codes used to build the DenseNet model, to compute the DeepLIFT contribution scores
 111 and to perform TF-MoDISco analysis are available at <https://github.com/kaixuanDeng95/iCREPCP>.

112 **Acknowledgments**

113 We acknowledge Prof. Weibo Xie for helpful discussions and we thank Dr. Yun Luo for providing the
 114 example of maize YIGE1 gene. We also thank four anonymous reviewers for their helpful suggestions
 115 that have greatly improved the original manuscript.

116 **Funding**

117 This work was supported by the National Natural Science Foundation of China (32070689 to Xuehai Hu)
 118 and the Joint Funds of the National Natural Science Foundation of China (U1901201 to Jianbing Yan).

119 **Conflicts of interest**

120 The authors declare no conflicts of interest.

121 **Authors contributions**

122 X.-H.H. and J.-B.Y. designed the research and wrote the manuscript. K.-X.D., Q.-Z.Z. and Y.-X.H.
 123 collected the data and built the model. K.-X.D. and Q.-Z.Z. performed the DeepLIFT analysis. K.-X.D.
 124 performed the motif analysis and developed the web server of iCREPCP. All authors read and approved
 125 the final manuscript.

126 **References**

- 127 **Fornes, O., Castro-Mondragon, J.A., Khan, A., van der Lee, R., Zhang, X., Richmond, P.A., Modi, B.P., Correard, S., Gheorghe, M., Baranašić, D.,**
 128 **et al.** (2020). JASPAR 2020: update of the open-access database of transcription factor binding profiles. *Nucleic acids research* **48**:D87-d92.
 129 10.1093/nar/gkz1001.
- 130 **Gao, C.** (2021). Genome engineering for crop improvement and future agriculture. *Cell* **184**:1621-1635. 10.1016/j.cell.2021.01.005.
- 131 **Huang, G.L., Z; Van Der Maaten, L; Weinberger, KQ.** (2017). Densely Connected Convolutional Networks. 2017 IEEE Conference on Computer Vision
 132 and Pattern Recognition (CVPR):17355312. DOI: 10.1109/CVPR.2017.243.
- 133 **Jores, T., Tonnies, J., Wrightsman, T., Buckler, E.S., Cuperus, J.T., Fields, S., and Queitsch, C.** (2021). Synthetic promoter designs enabled by a
 134 comprehensive analysis of plant core promoters. *Nature plants* **7**:842-855. 10.1038/s41477-021-00932-y.
- 135 **Koo, P.K., Majdandzic, A., Ploenzke, M., Anand, P., and Paul, S.B.** (2021). Global importance analysis: An interpretability method to quantify importance
 136 of genomic features in deep neural networks. *PLoS computational biology* **17**:e1008925. 10.1371/journal.pcbi.1008925.
- 137 **Luo, Y., Zhang, M., Liu, Y., Liu, J., Li, W., Chen, G., Peng, Y., Jin, M., Wei, W., Jian, L., et al.** (2022). Genetic variation in YIGE1 contributes to ear
 138 length and grain yield in maize. *The New phytologist* **234**:513-526. 10.1111/nph.17882.
- 139 **Schmitz, R.J., Grotewold, E., and Stam, M.** (2022). Cis-regulatory sequences in plants: Their importance, discovery, and future challenges. *The Plant cell*
 140 **34**:718-741. 10.1093/plcell/koab281.

141 **Shrikumar, A., Greenside, P., Kundaje, A.** (2017). Learning Important Features Through Propagating Activation Differences. Proceedings of the 34th
142 International Conference on Machine Learning **70**:3145-3153.
143 **Shrikumar, A., Tian, K., Avsec, Ž., Shcherbina, A., Banerjee, A., Sharmin, M., Nair, S., Kundaje, A.** (2018). TF-MoDISco v0.4.2.2-alpha: technical
144 note. Preprint at arXiv <https://arxiv.org/abs/1811.00416>.
145 **Song, X., Meng, X., Guo, H., Cheng, Q., Jing, Y., Chen, M., Liu, G., Wang, B., Wang, Y., Li, J., et al.** (2022). Targeting a gene regulatory element
146 enhances rice grain yield by decoupling panicle number and size. Nature biotechnology 10.1038/s41587-022-01281-7.
147

148

149 **Figure 1.** The workflow of iCREPCP.

150 **(A)** The deep learning architecture of DenseNet.

151 **(B)** The prediction performances via R^2 on training sets and on independent testing sets for six transient systems.

152 **(C)** The example of maize YIGE1 gene for demonstrating the detecting power of iCREPCP. Top panel, a snapshot of core
153 promoter region: chr1_51127917-51128086; The second panel is the FIMO scanning results; The third panel is the DeepLIFT
154 contribution score; The fourth and fifth panels are used to demonstrate the results of in-silico tilling deletion, which measure
155 the difference of predicted promoter strength with a sliding window of 5-bp deletion across the whole sequence; The bottom
156 panel is the heatmap for demonstrating in-silico mutagenesis results.

157 **(D)** A trans-species example of rice IPA1 gene with the same layout of (C).

158 **(E)** A total of 14 seqlets identified by TF-MoDISco of the model of tobacco leaves without enhancer in light and their similar TF
159 motifs in JASPAR.

160 **(F)** Motif occurrence frequencies and global importances of 14 enriched TF motifs of the model of tobacco leaves without enhancer
161 in light.

162 **(G)** The heatmap for demonstrating occurrence numbers of 14 enriched TF motifs within all 75,375 PCPs. Each row
163 represents a PCP and each column represents a specific TF motif. The row order (from top to bottom) is based on promoter
164 strength (from high to low) within each species and the column order (from left to right) is based on the total occurrence
165 number of TF motifs across three species (from more to less).

166 **(H)** The homepage of iCREPCP.

167

168

