

PLOS ONE
Nucleotide patterns aiding in prediction of eukaryotic promoters
 --Manuscript Draft--

Manuscript Number:	PONE-D-17-24929
Article Type:	Research Article
Full Title:	Nucleotide patterns aiding in prediction of eukaryotic promoters
Short Title:	Nucleotide patterns aiding in prediction of eukaryotic promoters
Corresponding Author:	Tatiana V Tatarinova, PhD University of La Verne La Verne, CA UNITED STATES
Keywords:	deep learning; promoters; classification
Abstract:	Computational analysis of promoters is hindered by their complex architecture. In less studied genomes with complex organization, false positive promoter predictions are common. Accurate identification of TSS and core promoter regions remains an unsolved problem. In this paper, we present a comprehensive analysis of genomic features associated with the promoters and show that probabilistic integrative algorithm models built upon the maps of the distributions of SNPs, RNA sequencing reads on genomic DNA, methylated nucleotides, gene models, TFBS as well as nucleotides and their combinations, allows one achieve accurate DNA sequence classification into "promoter" and "non-promoter", even in absence of the full-length cDNA sequences. Positional clustering of TBFS clearly demonstrates that the cells of <i>Oryza sativa</i> utilize three distinct classes of transcription factors: those that bind preferentially to the [-500,0] region ("promoter-specific"), those that bind preferentially to the [0,500] region ("5' UTR-specific"), and promiscuous transcription factors, that have weak or no location preference with respect to TSS. For the most informative motifs, their positional preferences were conserved between dicots and monocots.
Order of Authors:	Martin Triska Victor Solov'yev Ancha Baranova Alexander Kel Tatiana V Tatarinova, PhD
Opposed Reviewers:	
Additional Information:	
Question	Response
Financial Disclosure Please describe all sources of funding that have supported your work. This information is required for submission and will be published with your article, should it be accepted. A complete funding statement should do the following: Include grant numbers and the URLs of any funder's website. Use the full name, not acronyms, of funding institutions, and use initials to identify authors who received the funding. Describe the role of any sponsors or funders in the study design, data	AK was supported by a grant of the Federal Targeted Program "Research and development on priority directions of science and technology in Russia, 2014-2010", Contract № 14.604.21.0101, unique identifier of the applied scientific project: RFMEFI60414X0101. AK work was also supported by the following grants of the EU FP7 program: "SYSCOL", "SysMedIBD", "RESOLVE" and "MIMOMICS". TT, AB and MT were supported by the NSF Division of Environmental Biology (1456634) and NSF STTR award 1622840. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

<p>collection and analysis, decision to publish, or preparation of the manuscript. If the funders had no role in any of the above, include this sentence at the end of your statement: "<i>The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.</i>"</p> <p>However, if the study was unfunded, please provide a statement that clearly indicates this, for example: "<i>The author(s) received no specific funding for this work.</i>"</p>	
<p>* typeset</p>	
<p>Competing Interests</p> <p>You are responsible for recognizing and disclosing on behalf of all authors any competing interest that could be perceived to bias their work, acknowledging all financial support and any other relevant financial or non-financial competing interests.</p>	<p>None</p>
<p>Do any authors of this manuscript have competing interests (as described in the PLOS Policy on Declaration and Evaluation of Competing Interests)?</p> <p>If yes, please provide details about any and all competing interests in the box below. Your response should begin with this statement: <i>I have read the journal's policy and the authors of this manuscript have the following competing interests:</i></p>	
<p>If no authors have any competing interests to declare, please enter this statement in the box: "<i>The authors have declared that no competing interests exist.</i>"</p>	
<p>* typeset</p>	
<p>Ethics Statement</p> <p>You must provide an ethics statement if your study involved human participants, specimens or tissue samples, or vertebrate animals, embryos or tissues.</p>	<p>NA</p>

All information entered here should **also be included in the Methods section** of your manuscript. Please write "N/A" if your study does not require an ethics statement.

Human Subject Research (involved human participants and/or tissue)

All research involving human participants must have been approved by the authors' Institutional Review Board (IRB) or an equivalent committee, and all clinical investigation must have been conducted according to the principles expressed in the [Declaration of Helsinki](#). Informed consent, written or oral, should also have been obtained from the participants. If no consent was given, the reason must be explained (e.g. the data were analyzed anonymously) and reported. The form of consent (written/oral), or reason for lack of consent, should be indicated in the Methods section of your manuscript.

Please enter the name of the IRB or Ethics Committee that approved this study in the space below. Include the approval number and/or a statement indicating approval of this research.

Animal Research (involved vertebrate animals, embryos or tissues)

All animal work must have been conducted according to relevant national and international guidelines. If your study involved non-human primates, you must provide details regarding animal welfare and steps taken to ameliorate suffering; this is in accordance with the recommendations of the Weatherall report, "[The use of non-human primates in research](#)." The relevant guidelines followed and the committee that approved the study should be identified in the ethics statement.

If anesthesia, euthanasia or any kind of animal sacrifice is part of the study, please include briefly in your statement which substances and/or methods were applied.

<p>Please enter the name of your Institutional Animal Care and Use Committee (IACUC) or other relevant ethics board, and indicate whether they approved this research or granted a formal waiver of ethical approval. Also include an approval number if one was obtained.</p> <p>Field Permit</p> <p>Please indicate the name of the institution or the relevant body that granted permission.</p>	
<p>Data Availability</p> <p>PLOS journals require authors to make all data underlying the findings described in their manuscript fully available, without restriction and from the time of publication, with only rare exceptions to address legal and ethical concerns (see the PLOS Data Policy and FAQ for further details). When submitting a manuscript, authors must provide a Data Availability Statement that describes where the data underlying their manuscript can be found.</p> <p>Your answers to the following constitute your statement about data availability and will be included with the article in the event of publication. Please note that simply stating 'data available on request from the author' is not acceptable. If, however, your data are only available upon request from the author(s), you must answer "No" to the first question below, and explain your exceptional situation in the text box provided.</p> <p>Do the authors confirm that all data underlying the findings described in their manuscript are fully available without restriction?</p>	Yes - all data are fully available without restriction
<p>Please describe where your data may be found, writing in full sentences. Your answers should be entered into the box below and will be published in the form you provide them, if your manuscript is accepted. If you are copying our sample text below, please ensure you replace any instances of XXX with the appropriate details.</p> <p>If your data are all contained within the paper and/or Supporting Information files, please state this in your answer below. For example, "All relevant data are within</p>	All relevant data are within the paper and its Supporting Information files.

the paper and its Supporting Information files."

If your data are held or will be held in a public repository, include URLs, accession numbers or DOIs. For example, "All XXX files are available from the XXX database (accession number(s) XXX, XXX)." If this information will only be available after acceptance, please indicate this by ticking the box below.

If neither of these applies but you are able to provide details of access elsewhere, with or without limitations, please do so in the box below. For example:

"Data are available from the XXX Institutional Data Access / Ethics Committee for researchers who meet the criteria for access to confidential data."

"Data are from the XXX study whose authors may be contacted at XXX."

* typeset

Additional data availability information: Tick here if the URLs/accession numbers/DOIs will be available only after acceptance of the manuscript for publication so that we can ensure their inclusion before publication.

Dear Editors

We have analyzed characteristic features of plant promoter regions, such as drop in DNA methylation and SNP density, peak in TFBS frequency and CG-skew. We analyzed evolutionary conservation of these features and developed a machine learning method for classification of DNA sequence into promoter and non-promoter regions. This careful analysis will be instrumental for development of high confidence TSS prediction method. Accuracy of TSS identification affects the quality of regulatory region analysis. Intelligent integration of multiple types of genomic information (DNA composition, regulatory elements, DNA methylation, RNA-seq coverage data, SNP distribution etc.) can improve annotation of tissue- and developmental stage-specific genes that are often mis-predicted due to atypical sequence composition of stress-related genes in grasses. Integration of multiple noisy features of promoter regions can result in 99% classification accuracy. Features, identified as important by our classification process using deep learning approach, can be used to build a scoring function for promoter prediction.

We hope that our method will be of interest for a wide range of researchers, since the presented approach is easily scalable to other species.

Best regards

Tatiana Tatarinova

**Nucleotide patterns aiding in prediction of eukaryotic promoters**1
2 Martin Triska

3 University of Southern California, Los Angeles, CA, USA

4 martin.triska@usc.edu

5 Victor Solovyev

6 Softberry, Inc. Novosibirsk, Russia,

7 solovictor@softberry.com

8 Ancha Baranova

9 George Mason University, USA, and

10 Research Centre for Medical Genetics, Moscow, Russia

11 abaranov@gmu.edu12 Alexander Kel[#]

13 geneXplain GmbH, Wolfenbuettel, Germany

14 Institute of Chemical Biology and Fundamental Medicine, Novosibirsk, Russia,

15 alexander.kel@geneexplain.com

16 Tatiana V. Tatarinova*#

17 University of La Verne, 1950 3rd St, La Verne, CA 91750,

18 AA Kharkevich Institute for Information Transmission Problems RAS, Moscow, Russia,

19 ttatarinova@laverne.edu

20 #joint senior authors; *corresponding author

21

22 **Abstract**

23 Computational analysis of promoters is hindered by their complex architecture. In less studied genomes
24 with complex organization, false positive promoter predictions are common. Accurate identification of TSS
25 and core promoter regions remains an unsolved problem. In this paper, we present a comprehensive analysis
26 of genomic features associated with the promoters and show that probabilistic integrative algorithm models
27 built upon the maps of the distributions of SNPs, RNA sequencing reads on genomic DNA, methylated
28 nucleotides, gene models, TFBS as well as nucleotides and their combinations, allows one achieve accurate
29 DNA sequence classification into “promoter” and “non-promoter”, even in absence of the full-length cDNA
30 sequences.

31 Positional clustering of TBFS clearly demonstrates that the cells of *Oryza sativa* utilize three distinct
32 classes of transcription factors: those that bind preferentially to the [-500,0] region (“promoter-specific”),
33 those that bind preferentially to the [0,500] region (“5' UTR-specific”), and promiscuous transcription factors,
34 that have weak or no location preference with respect to TSS. For the most informative motifs, their positional
35 preferences were conserved between dicots and monocots.

36 **Introduction**

37 Promoters are the 5' regions adjacent to the transcriptional start site (TSS) that contain transcription factor
38 binding sites (TFBS) regulating transcription. Each gene contains one or more TSS, and, respectively, one or
39 more promoters, which initiate transcription of a gene. Computational analysis of the eukaryotic promoters is
40 hindered by their complex architecture (Sandelin, Carninci et al. 2007, Solovyev, Shahmuradov et al. 2010,
41 Shahmuradov, Umarov et al. 2017). In 30-60% of eukaryotic genes, depending on species, the proximal part
42 of the promoter, also known as core promotes is contains the TATA motif located approximately 30
43 nucleotides upstream from the TSS. Most commonly, TATA-containing core promoters are associated with
44 stress-related, tissue-specific and/or highly expressed genes (Troukhan, Tatarinova et al. 2009). Other, TATA-
45 less promoters commonly found in association with broadly expressed genes may have a relatively broad
46 transcription start region (TSR) instead of a pronounced TSS (Tatarinova, Kryshchenko et al. 2014,
47 Shahmuradov, Umarov et al. 2017). To predict the position of the TSR computationally, characteristic

48 promoter initiation region (Inr) or the downstream promoter element (DPE) may be used (Sandelin, Carninci
49 et al. 2007, Solovyev, Shahmuradov et al. 2010).

50 Many computational approaches for promoter prediction have been developed. A majority of TSS
51 prediction software tools use sophisticated computational algorithms (e.g. oligonucleotide content-based
52 neural network and linear discriminant approaches), and focus on specific sequence features of promoter
53 region (e.g. TATA-box or CA-motif) (Fickett and Hatzigeorgiou 1997). Genome complexity affects quality
54 of promoter finding methods: for example, presence of several tissue-specific, alternative TSSs negatively
55 affects the prediction accuracy. For the model plant *Arabidopsis thaliana*, modern algorithms identify TATA-
56 containing promoters with sensitivities up to 95% and specificities up to 97% TATA+ promoters
57 (Shahmuradov, Solovyev et al. 2005, Anwar, Baker et al. 2008, Troukhan, Tatarinova et al. 2009, Solovyev,
58 Shahmuradov et al. 2010, Azad, Shahid et al. 2011, Tatarinova, Kryshchenko et al. 2014, Umarov and
59 Solovyev 2017); while for *Homo sapiens* and *Oryza sativa*, prediction accuracy is substantially lower (Umarov
60 and Solovyev 2017). In case of less studied genomes with complex organization (such as the oil palm), false
61 positive as well as false negative error rates can be large, with a false positive promoter prediction occurring
62 once per every 700 - 1000 nucleotides of the genome (Solovyev 2003).

63 Even the best modern methods of promoter prediction, including genomic sequencing coupled with
64 full-length cDNA capture and ascertainment (Alexandrov, Troukhan et al. 2006, Alexandrov, Brover et al.
65 2009, Troukhan, Tatarinova et al. 2009), CAGE (Kawaji, Kasukawa et al. 2006, Kawaji, Lizio et al. 2014),
66 3PEAT(Morton, Petricka et al. 2014), or RAMPAGE (Batut and Gingeras 2013), do not predict TSS positions
67 with the desired 100% accuracy (Tatarinova, Kryshchenko et al. 2014). For example, the mapping of CAGE
68 tags relative to existing human cDNA/mRNA sequences revealed that less than 10% of these tags are located
69 within 10 nucleotides from TSS (Dieterich 2003). When RNA-seq reads covering 12 experimentally validated
70 *O. sativa* promoters were randomly selected from Plant Prom DB (Shahmuradov, Abdulazimova et al. 2012,
71 Shahmuradov, Umarov et al. 2017) and mapped onto the regions TSS+-1000 nt, the resultant plots showed
72 that the peaks of RNA-Seq coverage and exact positions of known TSS were mismatched, supporting the
73 correct mapping of eukaryotic promoters must be supported by multiple sources of data (Figure 1).

74 Figure 1: Examples of RNA-seq coverage near 12 randomly selected promoters with experimentally validated
75 transcription start sites.

76

77 Two most commonly used techniques to predict eukaryotic promoter by distribution of transcription
78 factors binding sites were proposed in 1995 by Kondrakhin, Kel et al. (1995) and by Prestridge (1995). The
79 method of Kondrakhin, Kel et al. (1995) paired up the detection of TATA boxes with computed weight
80 matrices of transcription factor binding site (TFBS) localizations and achieved much better accuracy in
81 comparison to approaches based on detection of the TATA box alone. Prestridge (1995) combined individual
82 density ratios of all TFBSs into a scoring profile further augmented by a weighted TATA matrix for a TATA
83 box and reported relatively low false positive rate. For both tools, however, applicability remains limited due
84 to insufficiency of reliable for species-specific model training and inability to point at the location of TSS.

85 In the last decade, several improvements in the promoter prediction were made. Troukhan (2009)
86 combined positional frequency of 5' EST matches onto genomic DNA with the gene models. This approach,
87 known as TSSer, is, in a nutshell, a deterministic method that predicts one transcription start site per locus.
88 For *Arabidopsis thaliana* promoters, it achieves remarkable accuracy. However, even the most reliable
89 prediction of a single promoter per gene cannot adequately reflect biological complexity underlying its
90 regulation due to common occurrence of alternative promoters, which are often tissue-specific or responsive
91 to the changes in architecture of chromatin (Rye, Sandve et al. 2014). In 2013, the TSSer approach was
92 augmented by non-parametric maximum likelihood approaches to be reborn as NPEST (Tatarinova,
93 Kryshchenko et al. 2014), which allows one to predict the positions of alternative TSSs in the *A. thaliana*
94 genome with better accuracy than the sequences identified in the several “gold standard” databases, such as
95 TAIR (Huala, Dickerman et al. 2001, Berardini, Reiser et al. 2015), Plant Prom DB (Shahmuradov,

96 Abdulazimova et al. 2012) and Plant Promoter Database (Hieno, Naznin et al. 2014). For example, for the set
97 of 15,875 Arabidopsis promoters derived by both TAIR and NPEST, 11,304 (71%) were predicted within 50
98 nucleotides of each other and 7,192 (45%) within 10 nucleotides of each other. Thirty percent of TAIR-
99 predicted and 44% of NPEST-predicted promoters identified the “TATA” sequence within the interval [-40,
100 -20] nucleotides upstream from the respective TSS. Nucleotide consensus scores at the TSS (46% of T and
101 49% of C followed by 65% of A) were stronger for NPEST than for TAIR (43% of T and 35% of C followed
102 by 53% of A). When NPEST predictions were compared to experimentally conformed promoters from other
103 databases, similar data patterns were observed.

104 Recently, many more types of experimental and computational observations highlighting the TSS
105 positions became available. For example, forty million single nucleotide polymorphisms (SNPs) from the
106 3,000 Rice Genomes Project (<http://snp-seek.irri.org>), the largest and highest density SNP collection for any
107 higher plant (Alexandrov, Tai et al. 2015), were shared to facilitate an analysis of genetic variants across the
108 *Oryza sativa* cultivars (Tatarinova, Chekalin et al. 2016). Observed clusters of nucleotide variability were
109 shown to highlight functionally important genomic regions. Interestingly, sharp declines in SNP density were
110 noted about 250 nucleotides upstream of TSS elements; these declines reach their minima exactly at the TSS.

111 Precise mapping of plant TSSs in the genes with multiple promoters requires incorporation of diverse
112 data types including of tissue/stress specificity of each transcript. Unfortunately, most currently available
113 techniques cannot incorporate diverse data types and simply ignore alternative promoters. Therefore, accurate
114 identification of TSS and core promoter regions remains an unsolved problem. Since evidences for location
115 of TSS are imprecise, the best approach for promoter production should embed probabilistic integrative
116 algorithms. In this paper, we present a comprehensive analysis of genomic features associated with the
117 promoters and show that models built upon the maps of the distributions of SNPs, RNA sequencing reads on
118 genomic DNA, methylated nucleotides, gene models, TFBS as well as nucleotides and their combinations,
119 allows one achieve classification of DNA sequence into “promoter” and “non-promoter”, that is accurate even
120 in absence of the full-length cDNA sequences.

121 **RESULTS**

122 **Selection of the “gold standard” gene prediction models**

123 *Figure 2: Nucleotide consensus around TSS. A) Frequency of CA, B) Frequency of TATA, D) Frequencies of*
124 *nucleotides A, C, G, T around TSS for Fgenesh, E) Frequencies of A, C, G, T around TSS for MSU*

125 To aid a selection of best available promoter annotation in the rice genome, Fgenesh and MSU mRNA-

126 based gene prediction models were compared. Fgenesh gene prediction set contains 18,389 high quality (5'
127 full, with mRNA support) gene models, while the MSU gene prediction set contains 20,367 high quality gene
128 models (Kawahara, de la Bastide et al. 2013, Tatarinova, Chekalin et al. 2016). For each gene in each model,
129 we extracted a 1,000 nt long sequence centered at TSS and calculated distributions of genomic features
130 previously associated with the start of transcription: (1) frequency of dinucleotide CA (Shinya and Shimada
131 1994, Pullen and Friesen 1995, Sandelin, Carninci et al. 2007); (2) frequency of TATA (Kiran, Ansari et al.
132 2006, Sandelin, Carninci et al. 2007, Troukhan, Tatarinova et al. 2009); (3) nucleotide consensus around TSS
133 (Alexandrov, Troukhan et al. 2006, Alexandrov, Brover et al. 2009, van Heeringen, Akhtar et al. 2011). Figure
134 2 (C and D) shows that the promoters annotated by Fgenesh shared more pronounced consensus as compared
135 to the promoters annotated by MSU. Fgenesh promoters also had higher frequency of TATA boxes at -30 (B),
136 and more CA dinucleotides at TSS (A). Based on the assumption that these features are reliable quality
137 indicators, Fgenesh annotation was selected for further analysis.

138 **Distribution of transcription factor binding sites**

139 The distributions of transcription factor binding sites (TFBS) in promoters and UTRs of high-
140 confidence rice genes in the regions of -1000 +1000 around TSS around TSS were investigated by utilizing
141 the MATCH software (Kel, Gossling et al. 2003) incorporated in geneXplain platform (www.genexplain.com)
142 and TRANSFAC database (Matys, Kel-Margoulis et al. 2006) comprising 764 plant position weight matrices
143 (PWM) with a strict score threshold of 0.95. MATCH scans the targets promoter sequences with a sliding
144 window equal to the length of the PWM and calculates a score for each of the windows. The maximum value
145 of the score (1.0) corresponds to the sequence that fully fits to the consensus of the PWM. Score threshold of
146 0.95 allows very little mismatches to the consensus located in the positions with lower degree of conservation

147 between known binding sites for this TF. The MATCH scores are calculated using an algorithm which
148 considers the nucleotide position-specific entropy measures. In a recent study, MATCH performed with
149 accuracy superior to that of other algorithms (Kondrakhin, Valeev et al. 2016).

150 In the Fgenesh-predicted rice promoters, MATCH search against the TRASNFAC database revealed
151 3.2 million potential TF binding sites corresponding to 667 PWMs, while 97 plant PWMs were matchless,
152 possibly due to their exclusive regulatory role in the dicots locations or binding to distal promoters excluded
153 from present study. Interestingly, 487 out of 667 TBFS (73%) were found in proximal promoters of *Oryza*
154 *sativa* more than 1000 times; the most frequent sites were that for the transcription factors ASR1, DOF56 and
155 PBF. When the frequencies of TFBS found in the proximal promoters were compared with the frequencies
156 for the same PWMs found in randomly shuffled sequences, the most significant promoter-specific enrichments
157 (twice or more) of TFBS were for SPL12, SPL5, GBF1, ABI5, BZIP68, LEC2, and GT1 transcription factors.

158 Additionally, another set of randomly shuffled sequences was generated using the dinucleotide
159 statistics matching that of the Fgenesh rice promoter regions as described by Stepanova, Tiazhelova et al.
160 (2005). Briefly, the 2000 nt regions [TSS-1000, TSS+1000] were divided onto non-overlapping 100 nt
161 windows, then the dinucleotide statistics were calculated for each window. For each promoter, a 2000 nt long
162 sequence with matching dinucleotide composition was generated, where the positions of TFBS were predicted
163 using the MATCH tool (Kel, Gossling et al. 2003). Kolmogorov-Smirnov tests were applied to find
164 significantly over-represented sequence motifs after selecting only those motifs that occur at least in 100
165 different rice promoters. Figure 3 shows examples of TBFS that differ and not differ significantly between
166 real and simulated sequences. The most pronounced differences ($p\text{-value} < 0.002$) were detected for the
167 binding sites of TCP15, LIM1, HBP1A, and TCP23. On the other hand, occurrences of binding sites for
168 CMTA2, GATA1, SBF1, and WRKY48 in real and simulated sequences were not different ($p\text{-value} >$
169 0.99999).

170 *Figure 3: Examples of observed and expected occurrences of TFBS in rice promoters. Different: TCP15,*
171 *LIM1, HBP1A, TCP23, ARALY493022, AT1G26610, TFIIAL, BZIP910, CBF1, DREB1F, STY1. Observations*
172 *agree with expectations: CMTA2, GATA1, SBF1, WRKY48.*

173 **Positional Specificity of TFBS distribution**

174 A phenomenon of the positional preference in TF binding was previously described by Weirauch,
175 Yang et al. (2014), who showed that, across evolutionary kingdoms, positions of TBFS are not random with
176 respect to the start of transcription. Therefore, we have computed Pearson's correlation coefficients and then
177 clustered the profiles of positional distributions of TFBS across the set of Fgenesh promoters of *Oryza sativa*.

178 Positional clustering of TBFS clearly demonstrates that the cells of *Oryza sativa* utilize three distinct
179 classes of transcription factors: those that bind preferentially to the [-500,0] region ("promoter-specific"),
180 those that bind preferentially to the [0,500] region ("5' UTR-specific"), and predominantly promiscuous
181 transcription factors with weak or no location preference for respective TSS. Note that some Class 3 TFs
182 cannot be classified as promiscuous (Figure 5), with more regular pattern of positional distribution around the
183 translation start rather than around the transcription start. Examples of the position frequency preference are

184 *Figure 4: Positional specificity of TFBS distribution*

185 shown in the Figure 6.

185 Comparative gene ontology analysis of Class 1, 2 and 3 transcription factors (see

186 Table 1) showed that Class 1 TFs are enriched in the following GO terms: “sequence-specific DNA
187 binding”, “protein dimerization activity”, “systemic acquired resistance, salicylic acid mediated signaling
188 pathway”, “regulation of transcription from RNA polymerase II promoter”, “response to bacterium”,
189 “jasmonic acid mediated signaling pathway”, “carpel development”, “protein binding”, “negative regulation
190 of defense response”, “protein targeting to membrane”, “regulation of plant-type hypersensitive response”,
191 “plant ovule development”, “response to ozone”. Class 2 TFs were enriched in “DNA binding”, “ethylene-
192 activated signaling pathway”, “response to water deprivation”. Class 3 TFs were enriched in “cellular response
193 to nitrogen levels”.

194 *Figure 5: The distribution pattern for MADSB binding sites highlight ATG codon rather than the respective*
195 *TSS.*

196

197 *Table 1: GO categories that are significantly different between three TF classes. Total number of genes with*
 198 *GO categories for Class 1, 2 and 3 is 130,164, and 144 correspondingly.*

GO	Class 1	Class 2	Class 3	PVAL
<i>sequence-specific DNA binding</i>	43	22	61	7.78E-06
<i>protein dimerization activity</i>	21	4	18	0.000438
<i>systemic acquired resistance, salicylic acid mediated signaling pathway</i>	15	2	6	0.000504
<i>regulation of transcription from RNA polymerase II promoter</i>	5	2	16	0.000558
<i>response to bacterium</i>	14	3	4	0.000955
<i>jasmonic acid mediated signaling pathway</i>	15	4	5	0.001897
<i>carpel development</i>	6	1	13	0.002606
<i>protein binding</i>	51	31	35	0.002904
<i>negative regulation of defense response</i>	11	1	5	0.003011
<i>DNA binding</i>	90	146	85	0.007339
<i>protein targeting to membrane</i>	14	4	7	0.010643
<i>regulation of plant-type hypersensitive response</i>	14	4	7	0.010643
<i>ethylene-activated signaling pathway</i>	10	20	4	0.012518
<i>response to water deprivation</i>	48	76	38	0.016138
<i>plant ovule development</i>	12	4	15	0.017373
<i>nucleus</i>	76	126	75	0.017489
<i>response to ozone</i>	9	4	14	0.033868
<i>cellular response to nitrogen levels</i>	13	11	23	0.044478

199

200 These three classes of TFs also differ in their expression patterns (Table 2): genes encoding TFs of
 201 Class 1 are predominantly expressed in petals, sepals and plant embryos, while mRNAs encoding Class 2 TFs
 202 are overrepresented in roots. This observation may explain why TATA boxes were the only cis-element

203 statistically significantly associated with expression in plant roots (Troukhan, Tatarinova et al. 2009, Triska,
204 Grocott et al. 2013): possibly, most root-specific transcription factors bind to the 5' UTR region rather than
205 the region upstream of TSS.

206 *Table 2: Expression specificity of TF from Class 1 and 2*

Expression pattern	Class 1 (99)	Class 2 (134)	Z-score
<i>root</i>	66	114	-3.31344
<i>pollen</i>	57	60	1.93163
<i>carpel</i>	72	77	2.39883
<i>seed</i>	70	74	2.40454
<i>leaf lamina base</i>	61	61	2.43144
<i>cauline leaf</i>	64	65	2.4497
<i>collective leaf structure</i>	80	87	2.65976
<i>petal</i>	72	74	2.73045
<i>plant embryo</i>	78	81	2.97259
<i>sepal</i>	76	76	3.17706

207
208 Figure 6 shows frequency profiles for TFBSSs for ARALY493022_04 (Class 1, left panel), RAP26_03
209 (Class 2, middle panel), and MYB111_02 (Class 3, right panel). ARALY493022 is basic helix-loop-helix
210 factor, with GGGCCC consensus sequence, its presence in the upstream region is associated with the overall
211 strength of gene expression (Troukhan, Tatarinova et al. 2009, Viola, Uberti Manassero et al. 2011, Weirauch,
212 Yang et al. 2014, Mathelier, Fornes et al. 2016). RAP2.6 is a defense-related, ethylene response transcription
213 factor, recognizing the GCC-box, and showing high binding affinity DNA sequence GCGCCGCCG (Franco-
214 Zorrilla, Lopez-Vidriero et al. 2014). Ali, Abbas et al. (2013) experimentally showed that RAP2.6 works both
215 in tissue-specific and stress-specific manner. Under normal conditions, expression of RAP26 is elevated in
216 roots and stems, while significantly reduced when plant is infected with pathogenic nematodes, such as *H.*

217 *schachtii*. To suppress resistance responses of the host plant, the nematodes downregulate expression of
218 RAP2.6. One of the better studied transcription factors with not-position-specific motifs is MYB111, which
219 is involved in the regulation of several genes in flavonoid biosynthesis pathway in cotyledons and leaves
220 (Stracke, Ishihara et al. 2007, Stracke, Jahns et al. 2010) and confers tolerance to UV-B (Stracke, Favery et
221 al. 2010). Binding site MYB111_02 has consensus G[G/T]TAGGT[A/G] (Franco-Zorrilla, Lopez-Vidriero et
222 al. 2014). Such less position-specific motifs often have quite conserved consensus and they could be found in
223 the promoters with relatively low frequency. Such motifs usually provide condition-specific regulation of
224 genes and don't have specific location relative TSS. From this Figure, it is obvious that Class 1 and Class 2
225 TBFS are useful for delineating potential TSSs, while the mapping of the Class 3 TBFS does not convey
226 positional information that can be used for TSS prediction.

227 *Figure 6: Frequency distribution of TFBS may have different patterns around the start of transcription*
228 (*position 0 on the horizontal axis*). X-axis shows distance from TSS, Y-axis is frequency of motif in each
229 window. Frequencies of ARALY493022_04 TFBS (Class 1) are plotted on the left panel, RAP26_03 TFBS
230 (Class 2) on the middle panel, and MYB111_02 (Class 3) on the right panel.

231 According to the Kolmogorov-Smirnov test, three classes of TFs differ in the significance of over-
232 representation of their TFBS in promoters as compared to randomly shuffled sequences: in the Class 1 TFs
233 with motifs located predominantly upstream of TSS, 37% of members were significantly overrepresented as
234 pointed by p-values <0.05, while in the Class 2 TFs with TFBS located in 5' UTRs overrepresentation was
235 confirmed for 20%, and in the Class 3 TFs associated more or less evenly distributed motifs, he over-
236 representation was detected for 15% of class members.

237 In summary, three classes of transcription factor binding sites differ in their position specificity,
238 significance of detected over-representation, functional classification and patterns of gene expression.

239 **Evolutionary conservation of motif informativeness**

240 We have analyzed evolutionary conservation of the position informativeness (a measure of unevenness
241 of the motif distribution along promoter regions, see Method section) for TFBS motifs in monocots *Oryza*
242 *sativa* and *Zea mays* and the dicot *Arabidopsis thaliana* (Figure 7). We identified the following correlations

243 between these measures:

244

$$I_{corn} = -0.0006 + 0.4422 \times I_{rice}$$

245 Multiple R²=0.9786, Adjusted R²=0.9785, F-statistic: 7732 on 1 and 169 DF, p-value: < 2.2E-16

246

$$I_{arabidopsis} = 0.001 + 0.6150 \times I_{rice}$$

247 Multiple R²=0.6627, Adjusted R²= 0.6619, F-statistic: 893.8 on 1 and 455 DF, p-value: < 2.2E-16.

248

249 *Figure 7: Relationship between informativeness of TFBS in rice, corn and Arabidopsis. Each point*
250 *corresponds to one transcription factor; X axis shows informativeness in rice, Y axis – informativeness in corn*
251 *and Arabidopsis.*

252 The correlations of the position informativeness for TFBS motifs in two monocots (rice and corn) were
253 higher than that for rice and a dicot plant Arabidopsis. By extracting TFBS with more than 10,000 matches in
254 each of three plant genomes, a list of 171 “common” TFBS was compiled. Each of these TFBS was classified
255 into “promoter-specific” or “5’ UTR-specific”. Between rice and corn, 67% of 171 “common” TBFS are
256 consistent in position preference, and 94% of the 90 most informative TBFS in rice had the same positional
257 preference in corn. Between rice and Arabidopsis, 86% of 171 TBFS agreed on their positional preference,
258 while 99% of the 90 most informative TBFS in rice had the same positional preference in Arabidopsis (see
259 Supplemental Data). From these numbers, it appears that rice has better agreement with Arabidopsis than with
260 corn in terms of promoter organization. We hypothesize that this discrepancy is due to less reliable TSS
261 prediction in corn genome as compared to Arabidopsis and rice (Figure 8). Importantly, incorrect prediction
262 of TSS in corn does not affect the informativeness of a TBFS, as defined by deviation from a uniform
263 distribution, but may lead to a systematic “shifting” of the TBFS peaks from upstream of TSS to “5’ UTR
264 specific” and vice versa. In summary, positional preference of the most informative motifs remains conserved
265 between dicots and monocots.

266 *Figure 8: Assessment of promoter prediction quality in Arabidopsis (left) and corn (right). Arabidopsis shows*
267 *more pronounced consensus at TSS, with higher frequency of TATA at -30 and CA at TSS.*

268

269 **Identification of similar TBFS**270 *Figure 9: An example of five distinct TBFS entries in the TRANSFAC database with very similar position*
271 *weight matrices (PWMs).*272 Since TRANSFAC database tends to collect all published motifs, some of its motifs appear to be
273 redundant. For example, several independently published PWMs may be independently built and reported for
274 the same transcription factor (Figure 9). Also, many of the transcription factors that belong to the same protein
275 family may recognize highly similar motifs, which will be reflected by respective PWMs. Figure 9 shows a
276 group of transcription factors with highly similar motifs. Although these transcription factors may differ in
277 their regulatory functions, or participate in differing regulatory networks, for a practical use in promoter
278 prediction, a non-redundant set of these motifs could be made by the clustering based on similarity of their
279 PWMs. By clustering 764 plant PWMs from TRANSFAC, a non-redundant set of 376 sequence motifs was
280 obtained, among which, forty-six were found informative with the scores above 0.0138 (see Supp. Table 1).281 **SNPs resulting in the TFBS loss and gain**282 Both the core promoter and the 5' UTR regions located within 200 bp around the TSS are protected
283 against genomic variations (Figure 10). This protective effect is due to selection constraints on regulatory
284 elements, i.e. TBFS, located near TSS, which prevents their disruption by neutral or near-neutral genomic
285 variants. Analysis of comprehensive collection of plant TBFS (Kondrakhin, Valeev et al. 2016) together with
286 an extensive dataset of the genomic variants in various rice cultivars (Alexandrov, Tai et al. 2015) allowed us
287 to classify regulatory elements of the plants according to their tolerance to the mutations.

288

289 *Figure 10: Frequency of SNPs located near the TSS in rice.*290 To achieve that, we considered distribution of SNPs and their effect on TF binding site loss and gain.
291 For each nucleotide change, we have calculated $\Delta = |q - q^*|$ for the TBFS scores before (q) and after (q^*)
292 nucleotide change, and compared its value to empirically determined threshold. Calculation of the scores q

293 and q^* was done per the MATCH score formula (see Materials and Methods). If $|\Delta| \geq |\Delta_0|$, the site was
294 considered as “lost” or “gained” depending which score value was larger, q or q^* .

295 Frequencies of site losses and site gains for the promoters and for the random subset of 18,389
296 intergenic sequences, each 2,000 nt in length, were compared. We hypothesized that functionally important
297 motifs in promoters will have lower amount of variation causing site loss. For each TF, we calculated ratio of
298 the number of site losses in promoters to the number of site losses in intergenic sequences, then the entries in
299 the list were ranked. The binding sites for ABF (CACGTGGC) and CBF4 transcription factors were the most
300 “protected” from the site loss. ABF factors govern osmotic stress response through modulation of the gene
301 expression downstream of SnRK2 kinases in abscisic acid signaling, while CBF4 regulates adaptation to
302 drought. For several important transcription factors, such as MADS8 (involved in the control of flowering
303 time), GT-1 and GATA-1 (response to light), we observed that variation was avoided in positions where
304 nucleotide change can lead to the site gain. Additional data and analysis of SNPs in TFBS is contained in the
305 Table 3, Table 4, and Supplemental Table 5.

306 The binding sites for AT2G20350 and ARF1 transcription factors were the most “protected” from the site
307 loss. AT2G20350 factors regulate activity of ethylene-activated signaling pathway. The plant hormone
308 ethylene is involved in many aspects of the plant life cycle, including seed germination, root hair development,
309 root nodulation, flower senescence, abscission, and fruit ripening (Johnson and Ecker, 1998). ARF1 is a
310 member of the auxin response factor family, involved in hyperosmotic salinity response. For several important
311 transcription factors, such as WRKY23 (involved in hyperosmotic salinity response and response to auxin),
312 FUS3 (plays a role in embryonic development ending in seed dormancy and response to auxin stimulus), we
313 observed that variation was avoided in positions where nucleotide change can lead to the site gain.

314 *Table 3: Suppression of site loss caused by SNPs in promoters. Frequency intergenic/ Frequency promoters*
315 *– ratio between frequencies of site loss due to SNPs in intergenic regions and site loss due to SNPs in*
316 *promoters.*

ID	Frequency intergenic/ Frequency promoters	#Promoter Sites	#Intergenic Sites	P-value
P\$AT2G20350_01	1.856	1909	3660	6.15E-112
P\$ARF1_01	1.814	856	1604	3.30E-47

P\$DREBIII4_01	1.677	870	1507	2.78E-35
P\$AT2G41690_01	1.648	1072	1825	5.37E-40
P\$CBF1_03	1.540	2803	4459	8.89E-74
P\$DREB1F_01	1.534	4873	7720	4.52E-124
P\$ORA47_01	1.529	4129	6521	8.35E-104
P\$RAP210_01	1.527	4519	7128	1.11E-112
P\$DEAR3_01	1.526	4525	7134	1.51E-112
P\$RAP210_02	1.526	4525	7134	1.51E-112
P\$ERF019_01	1.526	4199	6620	1.36E-104
P\$RAP21_01	1.526	4236	6675	3.17E-105
P\$AT1G71520_01	1.525	4242	6682	3.57E-105
P\$DREB1B_01	1.449	2544	3808	1.16E-48
P\$HSF3_01	1.423	1262	1855	1.08E-22
P\$AT4G16610_01	1.374	8852	12563	7.08E-118
P\$AT4G16750_01	1.347	1175	1635	2.62E-15
P\$AT2G44940_01	1.338	1231	1701	2.99E-15
P\$MADS17_01	1.307	7134	9628	1.37E-66

317

318

319

Table 4: Suppression of site gain caused by SNPs in promoters. Frequency intergenic/ Frequency promoters – ratio between frequencies of site gain due to SNPs in intergenic regions and site gain due to SNPs in promoters.

320

321

ID	Frequency intergenic/ Frequency promoters	#Promoter Sites	#Intergenic Sites	P-value
P\$WRKY23_01	1.376	1691	2403	2.59E-24
P\$FUS3_Q2	1.331	2249	3091	2.07E-25
P\$BHLH112_01	1.319	1813	2471	1.14E-19
P\$MYB46_02	1.290	1015	1352	4.37E-10
P\$WRKY_Q2	1.286	10925	14507	1.56E-88
P\$TGA2_Q2	1.275	6140	8084	3.08E-47
P\$CDC5_01	1.269	3653	4787	8.44E-28
P\$MADS4_01	1.266	1948	2548	1.89E-15

322

323

324 **Distribution of RNA-Seq reads**

325 *Figure 11: RNA-Seq coverage near the transcription start site*

326 Predictably, an analysis of mapped RNA-seq reads near TSS [-1000; +1000] showed that, on average,
327 coverage peaks are observed immediately downstream of TSS (Figure 11). However, some genes lack a peak
328 of RNS-seq at their TSS. Notably, only 26% of genes display a maximum of the coverage in the range [-50,
329 +250], and only 60% of genes display this maximum in the range [-50, +550].

330

331 **R-loop forming sequences (RLFS)**

332 Three-stranded nucleic acid R-loop structure is formed between nascent RNA transcript and DNA template
333 (Wongsurawat, Jenjaroenpun et al. 2012). R-loops help prevent methylation of promoters (Ginno, Lott et al.
334 2012, Ginno, Lim et al. 2013, Sanz, Hartono et al. 2016). Length of the R-loop sequence can be between 150
335 to 650 nt. R-loop forming sequences are associated with initiation of transcription and other important genic
336 features (Wongsurawat, Jenjaroenpun et al. 2012). R-loops were found to accumulate at the G-rich 5'-UTR
337 regions immediately downstream of the CpG-non-methylated human promoters (Ginno, Lott et al. 2012). We
338 used the QmRRFS tool (Wongsurawat, Jenjaroenpun et al. 2012, Jenjaroenpun, Chew et al. 2015,
339 Jenjaroenpun, Wongsurawat et al. 2015) to map the R-loop forming structures in the area [TSS-1000,
340 TSS+1000]. QmRRFS partitions RLFS three segments, RIZ (DNA region of initiation of R-loops containing
341 at least three contiguous guanines), linker (a spacer up to 50 nt between RIZ and REZ), and REZ (G-rich
342 region supporting extension of R-loop, up to 2000 nt long). According to QmRRFS, 22% of rice genes have
343 at least one RLFS in the area [TSS-1000, TSS+1000], and their distribution is localized to 5'-UTR, in
344 agreement with Ginno, Lott et al. (2012). Distribution of RLFS is unimodal with peak density position around
345 200 nt downstream from the TSS; over a half of RLFS (52%) are found in the 5'-UTR [TSS, TSS+400]
346 (Supplemental table 6). This position is associated with the polymerase pause region typically occurring after
347 the initiation of transcription (Wongsurawat, Jenjaroenpun et al. 2012, Jenjaroenpun, Chew et al. 2015,
348 Jenjaroenpun, Wongsurawat et al. 2015).

349 **DNA methylation**

350 *Figure 12: Methylation around transcription start site in rice in different contexts: red – CG, green - CHG,*
351 *blue – CHH.*

352 In the intergenic regions and within functional classes of genes and their promoters, the patterns of DNA
353 methylation predictably differ (Tatarinova, Elhaik et al. 2013, Elhaik, Pellegrini et al. 2014). The most
354 pronounced effect was observed for the methylated CpGs (see Figure 12). Intergenic level of CpG methylation
355 was at 0.27, with sharp decline starting around 600 bp upstream of TSS to about 50% of that in intergenic
356 region level at the position of -170, then proceeds to its minimum (0.01) at 8 bp upstream from TSS.

357 **Combining the characteristic features of TSS into promoter classifier**

358 We used 18,389 “promoter” (positives) and 18,389 “non-promoter” (negatives) sequences. To train the
359 model, we used 14,711 positives and negatives each; and for testing 3,678 positives and negatives each. The
360 binary classifier interrogates the candidate sequence and reports whether the sequence is “promoter” or
361 “non-promoter”. The best combination of features was: composition of DNA sequence, GC-skew value and
362 presence/absence of the CA-motif in every position. It achieved the best accuracy (0.9995) and has the
363 Matthews correlation coefficient of 0.9989 (see

364 Table 5). Other features also improve the classification accuracy in comparison with the DNA sequence alone,
365 however, not performing as well as the combination of DNA sequence, GC-skew and CA-motif distribution.
366

367 Table 5: Promoter classification accuracy

Features	TP	TN	FP	FN	Accuracy	Sensitivity	Specificity	CC
DNA sequence	3424	3030	648	254	0.8774	0.9309	0.8238	0.7591
DNA sequence + CG skew	3635	3653	25	43	0.9907	0.9883	0.9932	0.9832
DNA sequence + CG skew + frequency of CA motif	3674	3678	0	4	0.9994	0.9989	1.0	0.9989
DNA sequence + CG skew + RNA-seq coverage	3658	3666	12	20	0.9956	0.9945	0.9967	0.9913
DNA sequence + CG skew + frequency of TATA motif	3653	3608	70	25	0.9870	0.9932	0.9810	0.9742
DNA sequence + CG skew + DNA methylation	3657	3563	115	21	0.9815	0.9942	0.9687	0.9633
DNA sequence + all TFBS	3241	3386	292	437	0.9009	0.8812	0.9206	0.8024
DNA sequence + all TFBS +CG skew	3619	3668	10	59	0.9906	0.9839	0.9973	0.9813
DNA sequence + selected TFBS+CG skew	3628	3674	4	50	0.9927	0.9864	0.9989	0.9854
DNA sequence + SNP	3430	3138	540	248	0.8929	0.9326	0.8532	0.7882
DNA sequence + SNP+CG skew	3348	3296	382	330	0.9032	0.9103	0.8962	0.8065

368
369

370 **DISCUSSION**

371 We have investigated several features of promoter area and identified characteristic patterns of their
372 distribution in this area and assessed their utility for identification of TSS location. Promoter-UTR boundary
373 is marked by the following pronounced trends: (1) drop in SNP density, (2) evolutionary conserved peaks and
374 valleys of positions of regulatory elements, (3) peak of RNA-Seq coverage immediately downstream from the
375 TSS, (4) peak of CG skew, (5) drop in DNA methylation density in CpG, CHH and CHG contexts.

376 Analysis of SNPs in the context of TBFS in promoter and non-promoter region indicated that TBFS
377 differ by tolerance to polymorphisms. The binding sites for AT2G20350 and ARF1 transcription factors were
378 the most “protected” from the site loss. AT2G20350 factors regulate activity of ethylene-activated signaling
379 pathway. The plant hormone ethylene is involved in many aspects of the plant life cycle, including seed
380 germination, root hair development, root nodulation, flower senescence, abscission, and fruit ripening
381 (Johnson and Ecker 1998). ARF1 is a member of the auxin response factor family, involved in hyperosmotic
382 salinity response. For several important transcription factors, such as WRKY23 (involved in hyperosmotic
383 salinity response and response to auxin), FUS3 (plays a role in embryonic development ending in seed
384 dormancy and response to auxin stimulus), we observed that variation was avoided in positions where
385 nucleotide change can lead to the site gain Binding sites for several important transcription factors, such as
386 AT2G20350 and ARF1 found to be significantly less frequently lost due to polymorphic mutations then it is
387 expected by random chance. We can conclude that sites for such transcription factors are “protected” in
388 evolution from being lost since their high general importance for regulation of plant genes. Indeed, these
389 transcription factors are of utmost importance. They are involved in regulation of many plant developmental
390 processes in response to the plant hormones ethylene and auxin. Both these hormones play a central role in
391 coordination of many growth and behavioral processes in plants and are essential for plant body development.
392 It was interesting to observe that for several transcription factors nucleotide variations were avoided in
393 positions where nucleotide change can lead to the site gain. Among such factors we detect for instance
394 WRKY23 and FUS3. We propose that sporadic creation of new sites for such transcription factors can
395 significantly alter cellular timing, therefore such mutations are avoided in these classes of TBFS. It is

396 interesting that these factors are also involved in gene regulation in response to the plant hormone auxin. We
397 conclude that TRANSFAC analysis results in interesting observations of the architecture of rice promoters
398 and provide clear avoidance of interplay between SNPs and TF binding sites in rice genome.

399 Accuracy of TSS identification affects the quality of regulatory region analysis. Intelligent integration of
400 multiple types of genomic information (DNA composition, regulatory elements, DNA methylation, RNA-seq
401 coverage data, SNP distribution etc.) can improve annotation of tissue- and developmental stage-specific
402 genes that are often mis-predicted due to atypical sequence composition of stress-related genes in grasses
403 (Tatarinova 2010, Tatarinova, Elhaik et al. 2013, Elhaik, Pellegrini et al. 2014). Integration of multiple noisy
404 features of promoter regions can result in 99% classification accuracy. Features, identified as important by
405 our classification process using deep learning approach, can be used to build a scoring function for promoter
406 prediction.

407 MATERIALS AND METHODS

408 Fgenesh++ Rice Gene Prediction

409 Fgenesh++ (Find genes using Hidden Markov Models) (Yao, Guo et al. 2005, Bajic, Brent et al. 2006,
410 Solovyev, Kosarev et al. 2006) is a HMM-based *ab initio* gene prediction program (Salamov and Solovyev
411 2000). We used the rice chromosomes (version MSU 7, (Kawahara, de la Bastide et al. 2013)) to make the
412 initial gene prediction set, applying the Fgenesh gene finder with generic parameters for monocot plants. From
413 this set, we selected a subset of predicted genes that encode highly homologous proteins (using BLAST with
414 E-value cut-off 1.0E-10) to known plant proteins from the NCBI non-redundant (NR) database. Based on this
415 subset, we computed gene-finding parameters, optimized for the rice genome, and executed the Fgenesh++
416 pipeline to annotate the genes in the genomic scaffolds. The Fgenesh++ pipeline used all available supporting
417 data, such as known transcripts and homologous protein sequences. NR plant and, specifically, rice transcripts
418 were mapped to the rice genomic sequences, therefore identifying a set of potential splice sites. Plant proteins
419 were mapped to the rice genomic contigs, and the high scoring matches were selected to generate protein-
420 supported gene predictions, so that only the highly homologous proteins were used in gene identification.

421 Amino acid sequences from predicted rice genes were then compared to the protein sequences from
422 plant NR database using the 'bl2seq' routine, and the similarity was significant if it had a blast percent identity
423 ≥ 50 , blast score ≥ 100 , coverage of predicted protein $\geq 80\%$ and coverage of homologous protein $\geq 80\%$.
424 BLAST analysis of the predicted sequences was also carried out against the *O. sativa* mRNA dataset, using
425 an identify cutoff of $>90\%$. Predictions that have both NR plant RefSeq and *O. sativa* mRNA support, as well
426 as the 5' UTR longer than 20 nucleotides and shorter than 1000 were selected for the analysis.

427 GFF file with Fgenesh++ gene prediction is available as a Supplemental Data file.

428 MSU Rice Gene Models

429 The current MSUv7 annotation (<http://rice.plantbiology.msu.edu>) of rice genome contains 55,986
430 predicted genes and 66,338 gene models (Kawahara, de la Bastide et al. 2013). Upon exclusion of
431 pseudogenes, transposable elements, and genes with atypical lengths of 5' UTR (below 20 nt or above 1000 nt
432 long), a high-confidence set contains 20,367 expressed protein-coding rice genes.

433 Arabidopsis Gene and Promoter Models

434 Genome annotation files for TAIR 10 version and sequences for 3000 nucleotides upstream from ATG
435 were obtained from The Arabidopsis Information Resource ([TAIR](#)) (Lamesch, Berardini et al. 2012, Berardini,
436 Reiser et al. 2015). The upstream sequences were truncated based on the position of the nearest upstream
437 locus. 290,085 EST sequences were obtained from NCBI and TAIR and mapped onto the 27,199 upstream
438 sequences using nucleotide BLAST + (minimum identity percent: 95%; maximum query start of alignment:
439 5; only plus strand alignments were used). Using the text search, we removed ESTs annotated as 3' or partial.
440 NPEST (Tatarinova, Kryshchenko et al. 2014) algorithm was used and resulted in prediction of 17,452
441 transcription start sites for 16,520 protein-coding loci.

442 Corn Gene and Promoter Models

443 Genome annotation of maize ([B73, 6a](#)) contains 40,602 predicted protein-coding genes (Law, Childs
444 et al. 2015). We excluded genes with atypical lengths of 5' UTR (below 20 nt or above 1000 nt long), genes

445 without full-length mRNA support, without valid start and stop codon, or no PFAM annotation. This filtering
446 resulted in 16,180 putative corn TSS.

447 **Positional Informativeness of Transcription Factor Binding Sites**

448 We selected TFBS that occur at least 10,000 times in promoters of a given species. In rice, it amounted
449 to 487, in Arabidopsis -559, and in corn - 171 TBFS. To calculate informativeness of each TFBS, we divided
450 the region around the start of transcription (TSS-1000, TSS+1000) into 100 nt long non-overlapping windows,
451 and calculated the observed frequency of TFBS matches in every window as a ratio of matches within the
452 window to the total number of matches $f_o = \frac{m}{T}$. The expected frequency is calculated as $f_e =$
453 $\frac{1}{\text{Number of windows}}$. The informativeness defined as $I = \sum_{\text{Windows}} (f_o - f_e)^2$. The binding sites were ranked
454 from highest to lowest informativeness.

455 **RNA-seq Data**

456 We used following publicly available RNA-seq datasets: SRR034580, SRR034581, SRR034582,
457 SRR034583, SRR034584, SRR034585, SRR034586, SRR034587, SRR034588, SRR034589, SRR034590,
458 SRR034591, SRR034592, SRR034593, SRR034594, SRR034595, SRR034596, SRR034597, SRR034598,
459 SRR034599, SRR042529, SRR074125, SRR074126, SRR074127, SRR074128, SRR074129, SRR074130,
460 SRR074131, SRR074132, SRR074133, SRR074134, SRR074135, SRR074136, SRR074137, SRR074139,
461 SRR074140, SRR074142, SRR074143, SRR074144, SRR074145, SRR074146, SRR074147, SRR074149,
462 SRR074150.

463 The datasets were processed using the following protocol:

- 464 1. Duplicates were removed using tool *clumpify* (<http://ggi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/clumpify-guide/>) allowing for up to two errors per read.
- 465 2. Quality trimmed using *trimmmomatic* (Bolger, Lohse et al. 2014) with minimum read length=16,
466 minimum quality 28 (sliding window of length 10)
- 467 3. Aligned to the MSU 7 rice genome using *Hisat2* (Pertea, Kim et al. 2016) aligner.

468 Summary statistics is shown in the Table 6.

470 *Table 6: RNA-seq dataset quality.*

Experiments	Reads	Read length	Quality	Aligned
SRR034580-SRR034599	~5.5 M	35	Poor	67-72%
SRR042529	8.5M	36	Good	84%
SRR074125-SRR074150	2-5M	26	Good	~1.5% (!)

471

472 **Identification of Transcription factor binding sites**

473 The prediction of TF binding sites is done using the MATCH tool, which is based on the usage of information
 474 vector-based PWM model. This model calculates the *matrix similarity score* (mSS) defined in (Kel, Gossling
 475 et al. 2003). This model is a common additive model, which uses a transformed matrix instead of an initial
 476 matrix, where each column of the transformed matrix is determined with the help of weighting the
 477 corresponding initial column by information content. More specifically, the j^{th} column of the weight matrix is

478 equal (up to the constant $\frac{-MIN}{MAX - MIN}$) to the product of the j^{th} column of the frequency matrix and the

479 value $\frac{I(j)}{MAX - MIN}$, $i = 1, \dots, l$ where $I(j)$, MIN , and MAX were defined in Kel, Gossling et al. (2003).

480

481 **Site loss and gain**

482 We analyzed distribution of SNPs and their effect on TF binding site loss and gain. The effect of a
 483 SNP on TF binding sites was computed as the follows. For each SNP and for each PWM model we computed
 484 two TFBS scores: q and q^* corresponding to two nucleotides in the SNP – the reference and alternative

485 nucleotides. The score was calculated as:
$$q = \frac{\sum_{i=1}^l I(i)f(b_i, i) - \sum_{i=1}^l I(i)f^{\min}(i)}{\sum_{i=1}^l I(i)f^{\max}(i)}, \quad \text{where}$$

486 $I(i) = \sum_{b \in \{A,T,G,C\}} f(b,i) \ln(4f(b,i)) .$

487 Next, we calculated $\Delta = |q - q^*|$, and compared its value to empirically determined threshold. If

488 $|\Delta| \geq |\Delta_0|$, the site was considered as “lost” or “gained” depending on sign of the Δ .

489 We then calculated frequencies of site loss and site gain for all considered SNPs to identify loss or gain of
 490 which transcription factor binding sites (TFBS) are significantly enriched by the effect of nucleotide changes
 491 in SNPs analyzed. As a background, we considered random nucleotide changes in random genomic positions.
 492 We denote study and background sets briefly as “Yes” and “No” sets (the “Yes” set is the set of TFBS
 493 sequences overlapping SNPs with either the reference nucleotide or alternative nucleotide; the “No” set is the
 494 set created by random nucleotide substitutions in random genomic positions). The algorithm for TFBS
 495 enrichment analysis, called F-Match, has been described in Kel, Konovalova et al. (2006) and Koschmann,
 496 Bhar et al. (2015). Briefly, the procedure finds a critical value (a threshold) for the differences between scores
 497 q and q^* (*the threshold* Δ_0) of each PWM in the library that maximizes the “Yes/No” ratio R_{YN} as defined
 498 in Equation (1) under the constraint of statistical significance:

$$R_{YN} = \frac{\#Sites_{Yes} / \#Sites_{No}}{\#Seq_{Yes} / \#Seq_{No}} \quad (1)$$

500 In Equation (1), $\#Sites$ and $\#Seq$ are the sites and sequences counted in “Yes” and “No” sets. A high
 501 “Yes/No” ratio indicates strong enrichment of binding sites for a given PWM in the “Yes” sequences. The
 502 statistical significance is computed as follows:

$$\begin{aligned} P(X \geq x) &= \sum_{n=x}^N \binom{N}{n} \cdot p^n \cdot (1-p)^{N-n} \\ p &= \#Seq_{Yes} / (\#Seq_{Yes} + \#Seq_{No}) \\ N &= \#Sites_{Yes} + \#Sites_{No} \\ n &= \#Sites_{Yes} \end{aligned} \quad (2)$$

503 The Yes/No ratio and P-value is computed separately for the site gain and for the site loss. If “Yes/No”

505 ratio >1 and a P-value < 0.01 for a given PWM we consider this as an indication of an enrichment of SNPs by
506 the sites for the given PWM. We can say that sites of this PWM are frequently effected by the SNPs and
507 therefore the gene regulation by the respective TFs is significantly altered by the considered SNPs.

508 **Matrix clustering**

509 Many matrices in the TRANSFAC database are highly similar, up to the point being undistinguishable. To
510 lower the complexity of the training data, we performed hierarchical clustering and used only one matrix from
511 each cluster for promoter classification. The distance between two motifs is calculated as sum of squared
512 differences between all matrix elements. If matrices were not the same size, we slide the shorter matrix over
513 the longer one and take minimal distance. The cut-off for merging clusters was determined empirically by
514 considering the sequence logos of matrices to be merged at each step and deciding which matrices we consider
515 duplicates.

516 **Classification of Promoter Regions**

517 There are many network architectures and the task is to choose a suitable one for a given research problem.
518 We used Convolutional Neural Networks (CNN) architecture for building promoter recognition models
519 developed by Umarov and Solovyev (2017). The software consists of several modules. In the *learnCNN.py*
520 modules the CNN model was implemented using *Keras* - a minimalist, highly modular neural networks
521 library, written in Python. It uses the *Theano* library as a backend and utilizes GPU for fast neural network
522 training. *Adam* optimizer was used for training with categorical cross-entropy as a loss function. Our CNN
523 architecture (Figure 13) consists of one convolutional layer with 200 filters of length 21. After the
524 convolutional layer, there is a standard Max-Pooling layer. The output from the Max-Pooling layer is fed into
525 a standard fully connected ReLU layer with 128 neurons. Pooling size was equal to 2. The ReLU layer is
526 connected to the output layer with sigmoid activation, where neurons correspond to promoter and non-
527 promoter classes. The batch size used for training was 16.

528

529 *Figure 13: Basic CNN architecture that was used in building promoter models implemented in the*

530 *learnCNN.py* program (Shahmuradov, Umarov et al. 2017, Umarov and Solovyev 2017).

531

532 Input of the network consisted of nucleotide sequences where each nucleotide is encoded by a four-
533 dimensional vector A (1,0,0,0), T (0,1,0,0), G (0,0,1,0) and C (0,0,0,1) and other dimensions filled by other
534 promoter features such as: GC-skew, DNA methylation, SNP, presence of CA motif, presence of TATA
535 motifs, TFBS. The output is a two-dimensional vector: “promoter” (1, 0) and “non-promoter” (0, 1) prediction.

536 *learnCNN.py* learns parameters of the CNN model and outputs the accuracy of promoter prediction for the
537 test set of sequences. It also writes the computed CNN Model into a file, which can be used later in programs
538 for promoter identification in each sequence. We used 70% of these examples for learning, 10% for validation
539 (to find an optimal number of learning epochs) and 20% for testing.

540 We have extracted 18,389 sequences around transcription start site determined by full-length mRNA.
541 Sequence [TSS-199, TSS+50], containing 200 nucleotides from promoter and 50 nucleotides from 5' UTR,
542 was designated as the “promoter” region, and sequence [TSS+751, TSS+1000], from the coding part of the
543 gene, as “non-promoter”.

544 Quality of prediction was assessed using the following measures: True Positives (TP), True Negatives
545 (TN), False Positive (FP), False Negative (FN), Accuracy, Sensitivity, Specificity, Matthews correlation
546 coefficient (CC):

$$547 \quad Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$548 \quad Sensitivity = \frac{TP}{TP + FN}$$

$$549 \quad Specificity = \frac{TN}{TN + FP}$$

$$550 \quad CC = \frac{TP \times TN + FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

551

552

553 **Acknowledgments**

554 TT and MT were supported by the NSF Division of Environmental Biology (1456634) and NSF STTR award
555 1622840. AB was supported by NSF STTR award 1622840. AK was supported by a grant of the Federal
556 Targeted Program “Research and development on priority directions of science and technology in Russia,
557 2014-2010”, Contract № 14.604.21.0101, unique identifier of the applied scientific project:
558 RFMEFI60414X0101. AK was also supported by the following grants of the EU FP7 program: “SYSCOL”,
559 “SysMedIBD”, “RESOLVE” and “MIMOMICS”.
560

- 561 **References**
- 562 Alexandrov, N., S. Tai, W. Wang, L. Mansueto, K. Palis, R. R. Fuentes, V. J. Ulat, D. Chebotarov, G.
563 Zhang, Z. Li, R. Mauleon, R. S. Hamilton and K. L. McNally (2015). "SNP-Seek database of SNPs derived
564 from 3000 rice genomes." *Nucleic Acids Res* **43**(Database issue): D1023-1027.
- 565 Alexandrov, N. N., V. V. Brover, S. Freidin, M. E. Troukhan, T. V. Tatarinova, H. Zhang, T. J. Swaller, Y.
566 P. Lu, J. Bouck, R. B. Flavell and K. A. Feldmann (2009). "Insights into corn genes derived from large-scale
567 cDNA sequencing." *Plant Mol Biol* **69**(1-2): 179-194.
- 568 Alexandrov, N. N., M. E. Troukhan, V. V. Brover, T. Tatarinova, R. B. Flavell and K. A. Feldmann (2006).
569 "Features of Arabidopsis genes and genome discovered using full-length cDNAs." *Plant Mol Biol* **60**(1): 69-
570 85.
- 571 Ali, M. A., A. Abbas, D. P. Kreil and H. Bohlmann (2013). "Overexpression of the transcription factor
572 RAP2.6 leads to enhanced callose deposition in syncytia and enhanced resistance against the beet cyst
573 nematode *Heterodera schachtii* in *Arabidopsis* roots." *BMC Plant Biol* **13**: 47.
- 574 Anwar, F., S. M. Baker, T. Jabid, M. Mehedi Hasan, M. Shoyaib, H. Khan and R. Walshe (2008). "Pol II
575 promoter prediction using characteristic 4-mer motifs: a machine learning approach." *BMC Bioinformatics*
576 **9**: 414.
- 577 Azad, A. K., S. Shahid, N. Noman and H. Lee (2011). "Prediction of plant promoters based on hexamers and
578 random triplet pair analysis." *Algorithms Mol Biol* **6**: 19.
- 579 Bajic, V. B., M. R. Brent, R. H. Brown, A. Frankish, J. Harrow, U. Ohler, V. V. Solovyev and S. L. Tan
580 (2006). "Performance assessment of promoter predictions on ENCODE regions in the EGASP experiment."
581 *Genome Biol* **7 Suppl 1**: S3 1-13.
- 582 Batut, P. and T. R. Gingeras (2013). "RAMPAGE: promoter activity profiling by paired-end sequencing of
583 5'-complete cDNAs." *Curr Protoc Mol Biol* **104**: Unit 25B 11.
- 584 Berardini, T. Z., L. Reiser, D. Li, Y. Mezheritsky, R. Muller, E. Strait and E. Huala (2015). "The
585 *Arabidopsis* information resource: Making and mining the "gold standard" annotated reference plant
586 genome." *Genesis* **53**(8): 474-485.
- 587 Bolger, A. M., M. Lohse and B. Usadel (2014). "Trimmomatic: a flexible trimmer for Illumina sequence
588 data." *Bioinformatics* **30**(15): 2114-2120.
- 589 Dieterich, C., Wang,H., Rateitschak,K., Luz,H. and Vingron,M. (2003). "CORG: a database for
590 COmparative Regulatory Genomics." *Nucleic Acids Res* **31**: 55-57.
- 591 Elhaik, E., M. Pellegrini and T. V. Tatarinova (2014). "Gene expression and nucleotide composition are
592 associated with genic methylation level in *Oryza sativa*." *BMC Bioinformatics* **15**: 23.
- 593 Fickett, J. W. and A. G. Hatzigeorgiou (1997). "Eukaryotic promoter recognition." *Genome Res* **7**(9): 861-
594 878.
- 595 Franco-Zorrilla, J. M., I. Lopez-Vidriero, J. L. Carrasco, M. Godoy, P. Vera and R. Solano (2014). "DNA-
596 binding specificities of plant transcription factors and their potential to define target genes." *Proc Natl Acad
597 Sci U S A* **111**(6): 2367-2372.
- 598 Ginno, P. A., Y. W. Lim, P. L. Lott, I. Korf and F. Chedin (2013). "GC skew at the 5' and 3' ends of human
599 genes links R-loop formation to epigenetic regulation and transcription termination." *Genome Res* **23**(10):
600 1590-1600.
- 601 Ginno, P. A., P. L. Lott, H. C. Christensen, I. Korf and F. Chedin (2012). "R-loop formation is a distinctive
602 characteristic of unmethylated human CpG island promoters." *Mol Cell* **45**(6): 814-825.
- 603 Hieno, A., H. A. Naznin, M. Hyakumachi, T. Sakurai, M. Tokizawa, H. Koyama, N. Sato, T. Nishiyama, M.
604 Hasebe, A. D. Zimmer, D. Lang, R. Reski, S. A. Rensing, J. Obokata and Y. Y. Yamamoto (2014). "ppdb:
605 plant promoter database version 3.0." *Nucleic Acids Res* **42**(Database issue): D1188-1192.
- 606 Huala, E., A. W. Dickerman, M. Garcia-Hernandez, D. Weems, L. Reiser, F. LaFond, D. Hanley, D.
607 Kiphart, M. Zhuang, W. Huang, L. A. Mueller, D. Bhattacharyya, D. Bhaya, B. W. Sobral, W. Beavis, D.
608 W. Meinke, C. D. Town, C. Somerville and S. Y. Rhee (2001). "The *Arabidopsis* Information Resource
609 (TAIR): a comprehensive database and web-based information retrieval, analysis, and visualization system
610 for a model plant." *Nucleic Acids Res* **29**(1): 102-105.
- 611 Jenjaroenpun, P., C. S. Chew, T. P. Yong, K. Choowongkomon, W. Thammasorn and V. A. Kuznetsov
612 (2015). "The TTSMI database: a catalog of triplex target DNA sites associated with genes and regulatory

- elements in the human genome." *Nucleic Acids Res* **43**(Database issue): D110-116.
- Jenjaroenpun, P., T. Wongsurawat, S. P. Yenamandra and V. A. Kuznetsov (2015). "QmRLFS-finder: a model, web server and stand-alone tool for prediction and analysis of R-loop forming sequences." *Nucleic Acids Res* **43**(20): 10081.
- Johnson, P. R. and J. R. Ecker (1998). "The ethylene gas signal transduction pathway: a molecular perspective." *Annu Rev Genet* **32**: 227-254.
- Kawahara, Y., M. de la Bastide, J. P. Hamilton, H. Kanamori, W. R. McCombie, S. Ouyang, D. C. Schwartz, T. Tanaka, J. Wu, S. Zhou, K. L. Childs, R. M. Davidson, H. Lin, L. Quesada-Ocampo, B. Vaillancourt, H. Sakai, S. S. Lee, J. Kim, H. Numa, T. Itoh, C. R. Buell and T. Matsumoto (2013). "Improvement of the *Oryza sativa* Nipponbare reference genome using next generation sequence and optical map data." *Rice (N Y)* **6**(1): 4.
- Kawaji, H., T. Kasukawa, S. Fukuda, S. Katayama, C. Kai, J. Kawai, P. Carninci and Y. Hayashizaki (2006). "CAGE Basic/Analysis Databases: the CAGE resource for comprehensive promoter analysis." *Nucleic Acids Res* **34**(Database issue): D632-636.
- Kawaji, H., M. Lizio, M. Itoh, M. Kanamori-Katayama, A. Kaiho, H. Nishiyori-Sueki, J. W. Shin, M. Kojima-Ishiyama, M. Kawano, M. Murata, N. Ninomiya-Fukuda, S. Ishikawa-Kato, S. Nagao-Sato, S. Noma, Y. Hayashizaki, A. R. Forrest, P. Carninci and F. Consortium (2014). "Comparison of CAGE and RNA-seq transcriptome profiling using clonally amplified and single-molecule next-generation sequencing." *Genome Res* **24**(4): 708-717.
- Kel, A., T. Konovalova, T. Valeev, E. Cheremushkin, O. Kel-Margoulis and E. Wingender (2006). "Composite Module Analyst: a fitness-based tool for identification of transcription factor binding site combinations." *Bioinformatics* **22**(10): 1190-1197.
- Kel, A. E., E. Gossling, I. Reuter, E. Cheremushkin, O. V. Kel-Margoulis and E. Wingender (2003). "MATCH: A tool for searching transcription factor binding sites in DNA sequences." *Nucleic Acids Res* **31**(13): 3576-3579.
- Kiran, K., S. A. Ansari, R. Srivastava, N. Lodhi, C. P. Chaturvedi, S. V. Sawant and R. Tuli (2006). "The TATA-box sequence in the basal promoter contributes to determining light-dependent gene expression in plants." *Plant Physiol* **142**(1): 364-376.
- Kondrakhin, Y., T. Valeev, R. Sharipov, I. Yevshin, F. Kolpakov and A. Kel (2016). "Prediction of protein-DNA interactions of transcription factors linking proteomics and transcriptomics data." *EuPA Open Proteomics* **13**: 14-23.
- Kondrakhin, Y. V., A. E. Kel, N. A. Kolchanov, A. G. Romashchenko and L. Milanesi (1995). "Eukaryotic promoter recognition by binding sites for transcription factors." *Comput Appl Biosci* **11**(5): 477-488.
- Koschmann, J., A. Bhar, P. Stegmaier, A. E. Kel and E. Wingender (2015). ""Upstream Analysis": An Integrated Promoter-Pathway Analysis Approach to Causal Interpretation of Microarray Data." *Microarrays (Basel)* **4**(2): 270-286.
- Lamesch, P., T. Z. Berardini, D. Li, D. Swarbreck, C. Wilks, R. Sasidharan, R. Muller, K. Dreher, D. L. Alexander, M. Garcia-Hernandez, A. S. Karthikeyan, C. H. Lee, W. D. Nelson, L. Ploetz, S. Singh, A. Wensel and E. Huala (2012). "The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools." *Nucleic Acids Res* **40**(Database issue): D1202-1210.
- Law, M., K. L. Childs, M. S. Campbell, J. C. Stein, A. J. Olson, C. Holt, N. Panchy, J. Lei, D. Jiao, C. M. Andorf, C. J. Lawrence, D. Ware, S. H. Shiu, Y. Sun, N. Jiang and M. Yandell (2015). "Automated update, revision, and quality control of the maize genome annotations using MAKER-P improves the B73 RefGen_v3 gene models and identifies new genes." *Plant Physiol* **167**(1): 25-39.
- Mathelier, A., O. Fornes, D. J. Arenillas, C. Y. Chen, G. Denay, J. Lee, W. Shi, C. Shyr, G. Tan, R. Worsley-Hunt, A. W. Zhang, F. Parcy, B. Lenhard, A. Sandelin and W. W. Wasserman (2016). "JASPAR 2016: a major expansion and update of the open-access database of transcription factor binding profiles." *Nucleic Acids Res* **44**(D1): D110-115.
- Matys, V., O. V. Kel-Margoulis, E. Fricke, I. Liebich, S. Land, A. Barre-Dirrie, I. Reuter, D. Chekmenev, M. Krull, K. Hornischer, N. Voss, P. Stegmaier, B. Lewicki-Potapov, H. Saxel, A. E. Kel and E. Wingender (2006). "TRANSFAC and its module TRANSCompel: transcriptional gene regulation in eukaryotes." *Nucleic Acids Res* **34**(Database issue): D108-110.
- Morton, T., J. Petricka, D. L. Corcoran, S. Li, C. M. Winter, A. Carda, P. N. Benfey, U. Ohler and M.

- 666 Megraw (2014). "Paired-end analysis of transcription start sites in Arabidopsis reveals plant-specific
667 promoter signatures." *Plant Cell* **26**(7): 2746-2760.
- 668 Pertea, M., D. Kim, G. M. Pertea, J. T. Leek and S. L. Salzberg (2016). "Transcript-level expression analysis
669 of RNA-seq experiments with HISAT, StringTie and Ballgown." *Nat Protoc* **11**(9): 1650-1667.
- 670 Prestridge, D. S. (1995). "Predicting Pol II promoter sequences using transcription factor binding sites." *J
671 Mol Biol* **249**(5): 923-932.
- 672 Pullen, S. S. and P. D. Friesen (1995). "The CAGT motif functions as an initiator element during early
673 transcription of the baculovirus transregulator ie-1." *J Virol* **69**(6): 3575-3583.
- 674 Rye, M., G. K. Sandve, C. O. Daub, H. Kawaji, P. Carninci, A. R. Forrest, F. Drablos and F. consortium
675 (2014). "Chromatin states reveal functional associations for globally defined transcription start sites in four
676 human cell lines." *BMC Genomics* **15**: 120.
- 677 Salamov, A. A. and V. V. Solovyev (2000). "Ab initio gene finding in Drosophila genomic DNA." *Genome
678 Res* **10**(4): 516-522.
- 679 Sandelin, A., P. Carninci, B. Lenhard, J. Ponjavic, Y. Hayashizaki and D. A. Hume (2007). "Mammalian
680 RNA polymerase II core promoters: insights from genome-wide studies." *Nat Rev Genet* **8**(6): 424-436.
- 681 Sanz, L. A., S. R. Hartono, Y. W. Lim, S. Steyaert, A. Rajpurkar, P. A. Ginno, X. Xu and F. Chedin (2016).
682 "Prevalent, Dynamic, and Conserved R-Loop Structures Associate with Specific Epigenomic Signatures in
683 Mammals." *Mol Cell* **63**(1): 167-178.
- 684 Shahmuradov, I. A., A. Abdulazimova, F. Z. Khan, V. Solovyev, N. Mustafaev, Y. Akbarova, R. Qamar and
685 J. Aliyev (2012). The PlantProm DB: Recent Updates. *2012 International Conference on Biomedical
686 Engineering and Biotechnology (iCBEB)*. IEEE. Macau, Macao.
- 687 Shahmuradov, I. A., V. V. Solovyev and A. J. Gammerman (2005). "Plant promoter prediction with
688 confidence estimation." *Nucleic Acids Res* **33**(3): 1069-1076.
- 689 Shahmuradov, I. A., R. K. Umarov and V. V. Solovyev (2017). "TSSPlant: a new tool for prediction of plant
690 Pol II promoters." *Nucleic Acids Res*.
- 691 Shinya, E. and T. Shimada (1994). "Identification of two initiator elements in the bidirectional promoter of
692 the human dihydrofolate reductase and mismatch repair protein 1 genes." *Nucleic Acids Res* **22**(11): 2143-
693 2149.
- 694 Solovyev, V., P. Kosarev, I. Seledsov and D. Vorobyev (2006). "Automatic annotation of eukaryotic genes,
695 pseudogenes and promoters." *Genome Biol* **7 Suppl 1**: S10 11-12.
- 696 Solovyev, V., Shahmuradov, I. (2003). "PromH: promoters identification using orthologous genomic
697 sequences." *Nucleic Acids Research* **31**(13).
- 698 Solovyev, V. V., I. A. Shahmuradov and A. A. Salamov (2010). "Identification of promoter regions and
699 regulatory sites." *Methods Mol Biol* **674**: 57-83.
- 700 Stepanova, M., T. Tiazhelova, M. Skoblov and A. Baranova (2005). "A comparative analysis of relative
701 occurrence of transcription factor binding sites in vertebrate genomes and gene promoter areas."
702 *Bioinformatics* **21**(9): 1789-1796.
- 703 Stracke, R., J. J. Favory, H. Gruber, L. Bartelniewoehner, S. Bartels, M. Binkert, M. Funk, B. Weisshaar and
704 R. Ulm (2010). "The Arabidopsis bZIP transcription factor HY5 regulates expression of the PFG1/MYB12
705 gene in response to light and ultraviolet-B radiation." *Plant Cell Environ* **33**(1): 88-103.
- 706 Stracke, R., H. Ishihara, G. Huep, A. Barsch, F. Mehrrens, K. Niehaus and B. Weisshaar (2007).
707 "Differential regulation of closely related R2R3-MYB transcription factors controls flavonol accumulation
708 in different parts of the Arabidopsis thaliana seedling." *Plant J* **50**(4): 660-677.
- 709 Stracke, R., O. Jahns, M. Keck, T. Tohge, K. Niehaus, A. R. Fernie and B. Weisshaar (2010). "Analysis of
710 PRODUCTION OF FLAVONOL GLYCOSIDES-dependent flavonol glycoside accumulation in
711 Arabidopsis thaliana plants reveals MYB11-, MYB12- and MYB111-independent flavonol glycoside
712 accumulation." *New Phytol* **188**(4): 985-1000.
- 713 Tatarinova, T., Alexandrov, N., Bouck, J., Feldmann, K. (2010). "GC3 Biology in Corn, Rice, Sorghum and
714 other grasses." *BMC Genomics* **11**(308).
- 715 Tatarinova, T., E. Elhaik and M. Pellegrini (2013). "Cross-species analysis of genic GC3 content and DNA
716 methylation patterns." *Genome Biol Evol* **5**(8): 1443-1456.
- 717 Tatarinova, T., A. Kryshchenko, M. Triska, M. Hassan, D. Murphy, M. Neely and A. Schumitzky (2014).
718 "NPEST: a nonparametric method and a database for transcription start site prediction." *Quant Biol* **1**(4):

- 719 261-271.
- 720 Tatarinova, T. V., E. Chekalin, Y. Nikolsky, S. Bruskin, D. Chebotarov, K. L. McNally and N. Alexandrov
721 (2016). "Nucleotide diversity analysis highlights functionally important genomic regions." *Sci Rep* **6**:
722 35730.
- 723 Triska, M., D. Grocott, J. Southern, D. J. Murphy and T. Tatarinova (2013). "cisExpress: motif detection in
724 DNA sequences." *Bioinformatics* **29**(17): 2203-2205.
- 725 Troukhan, M., T. Tatarinova, J. Bouck, R. B. Flavell and N. N. Alexandrov (2009). "Genome-wide
726 discovery of cis-elements in promoter sequences using gene expression." *OMICS* **13**(2): 139-151.
- 727 Troukhan, M., Tatarinova, T. Bouck, J., Flawell, R., Alexandrov, N. (2009). "Genome-wide discovery of
728 cis-elements in promoter sequences using gene expression data." *Oomics* **13**(1).
- 729 Umarov, R. K. and V. V. Solovyev (2017). "Recognition of prokaryotic and eukaryotic promoters using
730 convolutional deep learning neural networks." *PLoS One* **12**(2): e0171410.
- 731 van Heeringen, S. J., W. Akhtar, U. G. Jacobi, R. C. Akkers, Y. Suzuki and G. J. Veenstra (2011).
732 "Nucleotide composition-linked divergence of vertebrate core promoter architecture." *Genome Res* **21**(3):
733 410-421.
- 734 Viola, I. L., N. G. Uberti Manassero, R. Ripoll and D. H. Gonzalez (2011). "The Arabidopsis class I TCP
735 transcription factor AtTCP11 is a developmental regulator with distinct DNA-binding properties due to the
736 presence of a threonine residue at position 15 of the TCP domain." *Biochem J* **435**(1): 143-155.
- 737 Weirauch, M. T., A. Yang, M. Albu, A. G. Cote, A. Montenegro-Montero, P. Drewe, H. S. Najafabadi, S. A.
738 Lambert, I. Mann, K. Cook, H. Zheng, A. Goity, H. van Bakel, J. C. Lozano, M. Galli, M. G. Lewsey, E.
739 Huang, T. Mukherjee, X. Chen, J. S. Reece-Hoyes, S. Govindarajan, G. Shaulsky, A. J. Walhout, F. Y.
740 Bouget, G. Ratsch, L. F. Larrondo, J. R. Ecker and T. R. Hughes (2014). "Determination and inference of
741 eukaryotic transcription factor sequence specificity." *Cell* **158**(6): 1431-1443.
- 742 Wongsurawat, T., P. Jenjaroenpun, C. K. Kwoh and V. Kuznetsov (2012). "Quantitative model of R-loop
743 forming structures reveals a novel level of RNA-DNA interactome complexity." *Nucleic Acids Res* **40**(2):
744 e16.
- 745 Yao, H., L. Guo, Y. Fu, L. A. Borsuk, T. J. Wen, D. S. Skibbe, X. Cui, B. E. Scheffler, J. Cao, S. J. Emrich,
746 D. A. Ashlock and P. S. Schnable (2005). "Evaluation of five ab initio gene prediction programs for the
747 discovery of maize genes." *Plant Mol Biol* **57**(3): 445-460.
- 748

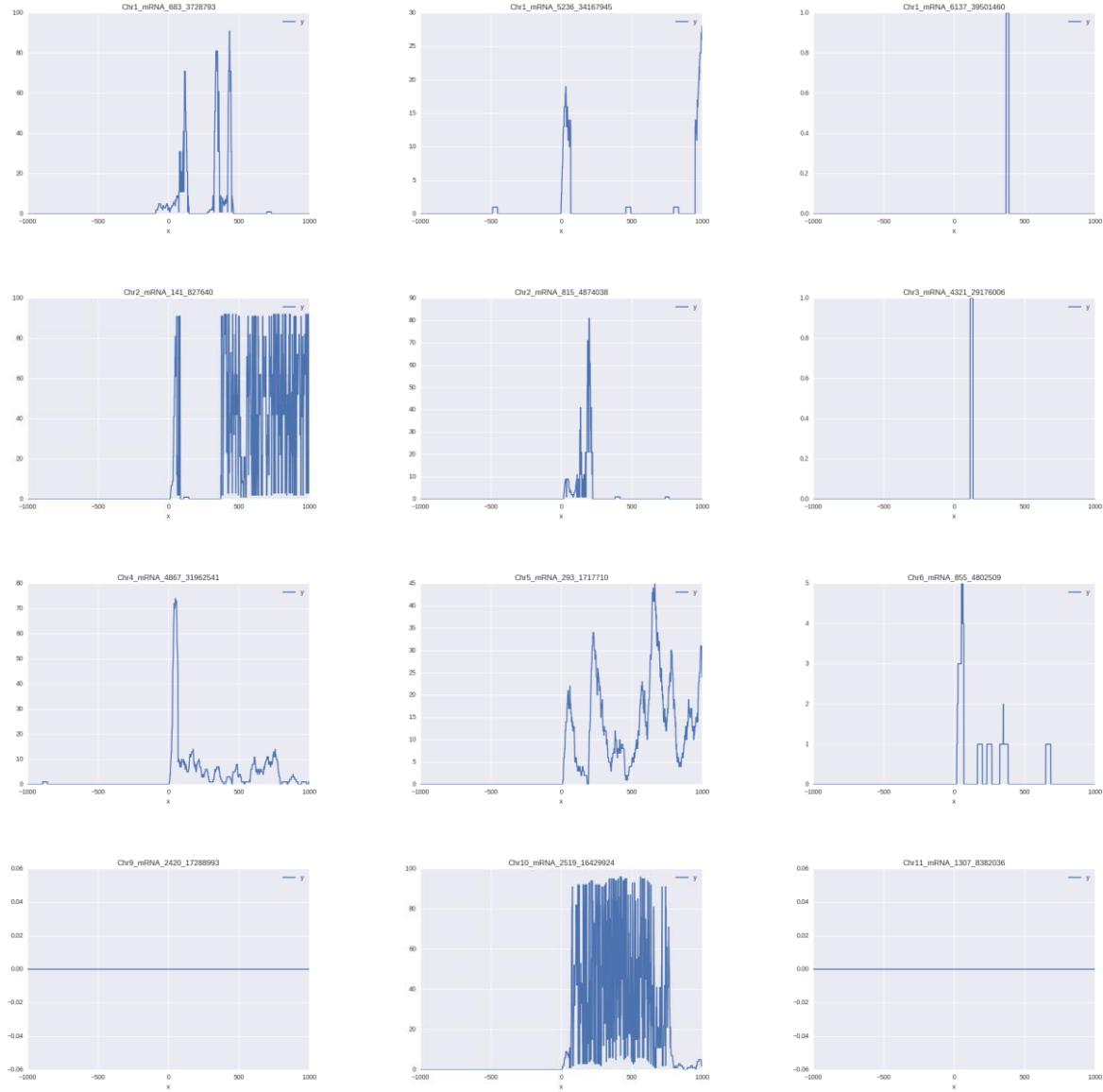


Figure 2

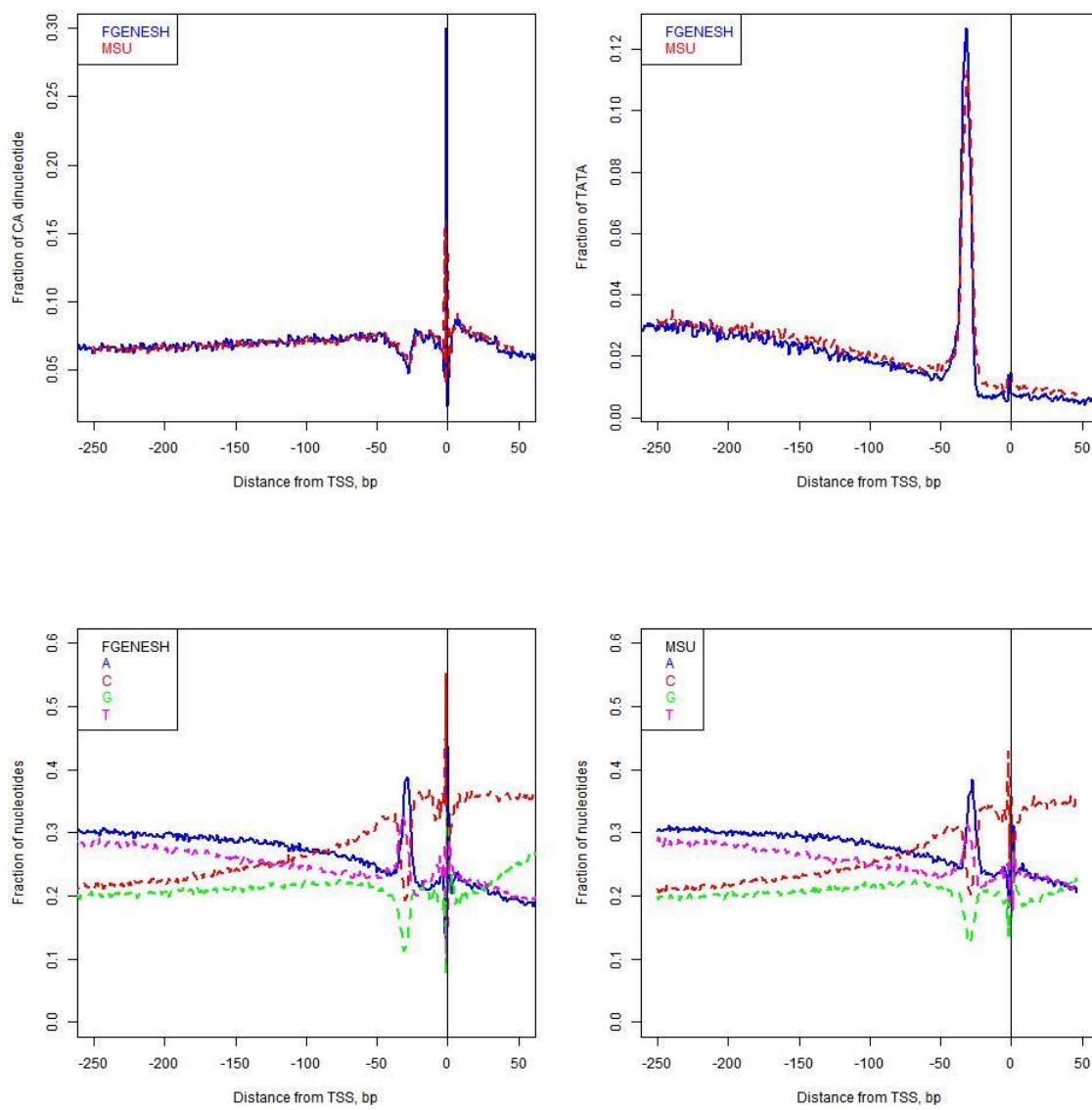
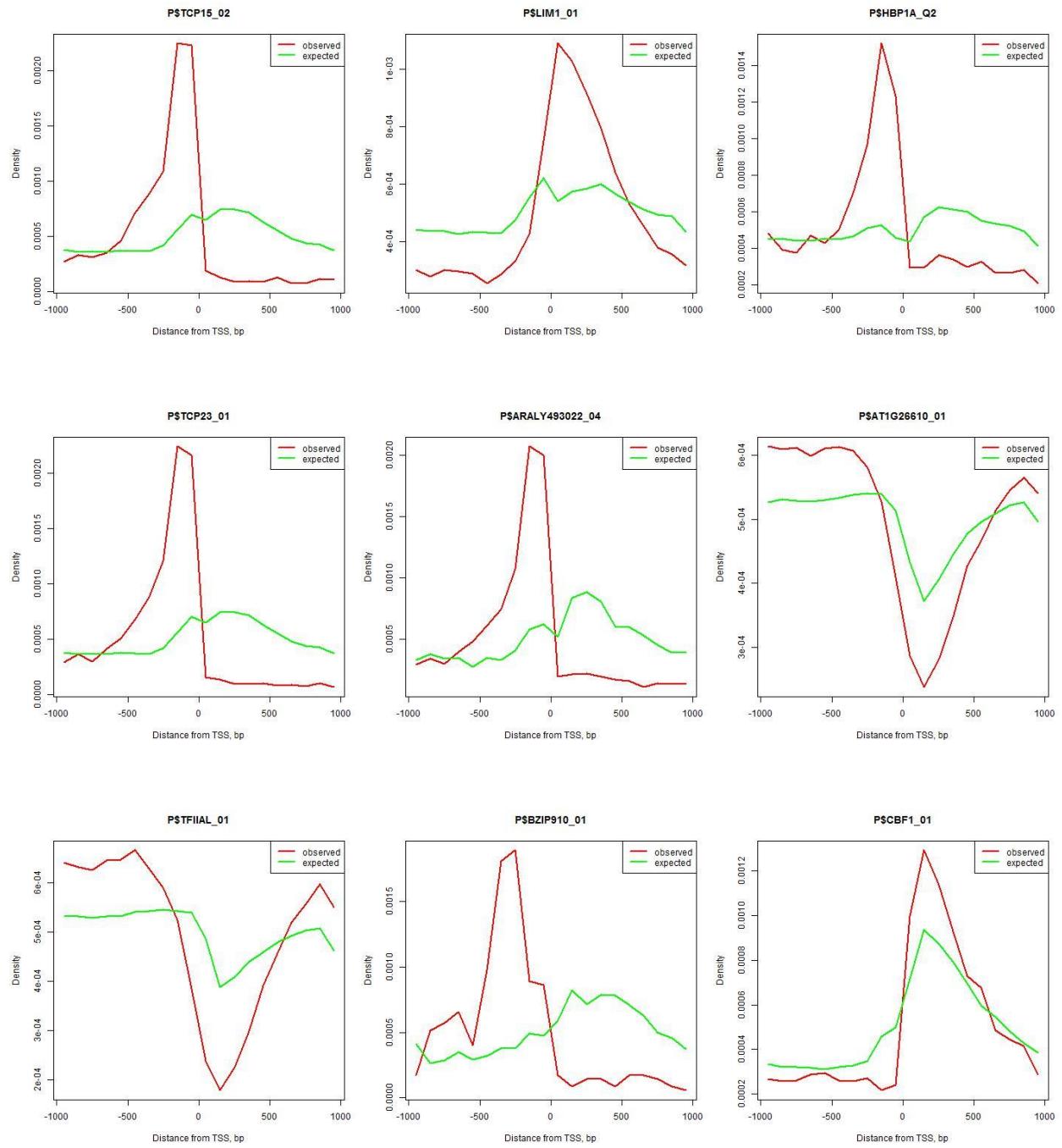
[Click here to download Figure Figure 2.docx](#)

Figure 3

[Click here to download Figure Figure 3.docx](#)

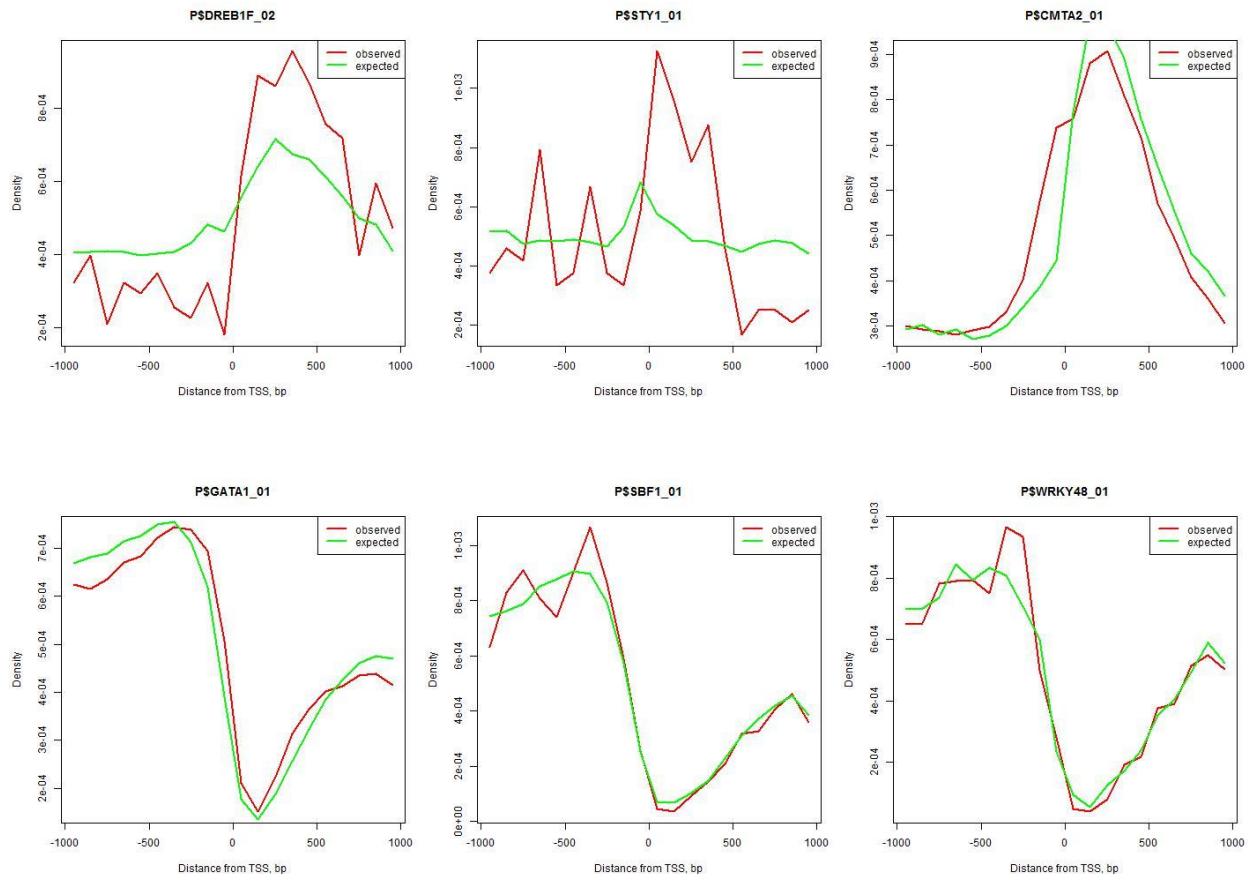


Figure 4

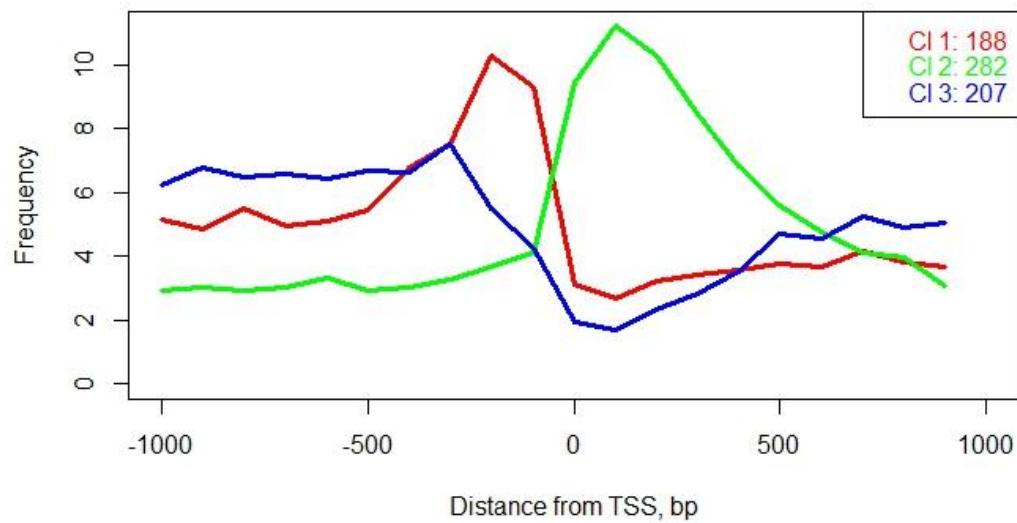
[Click here to download Figure Figure 4.docx](#)

Figure 5

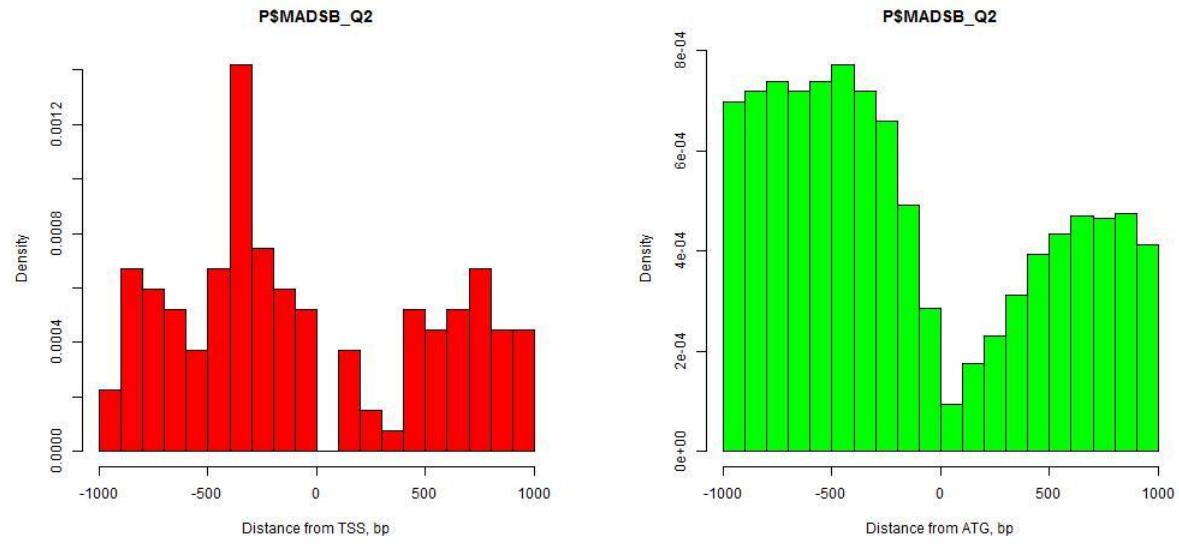
[Click here to download Figure Figure 5.docx](#)

Figure 6

[Click here to download Figure Figure 6.docx](#)

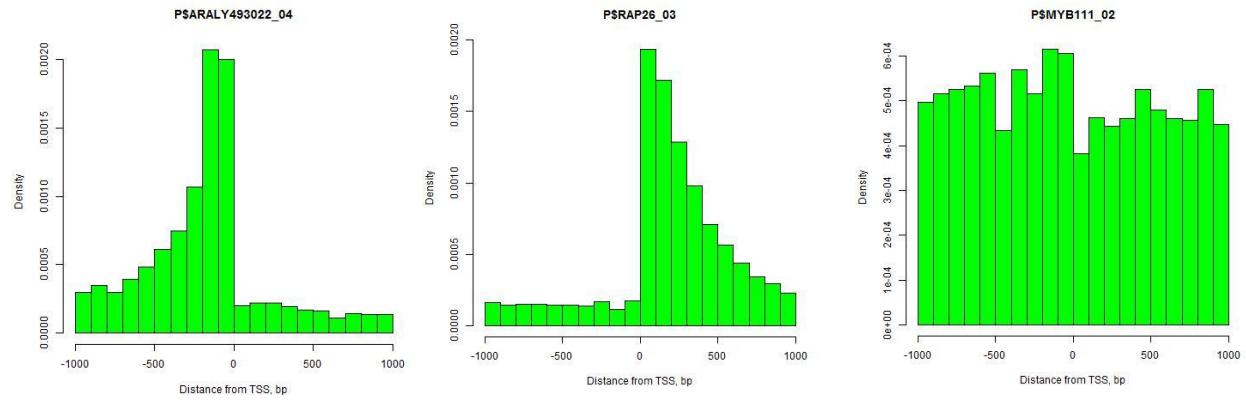


Figure 7

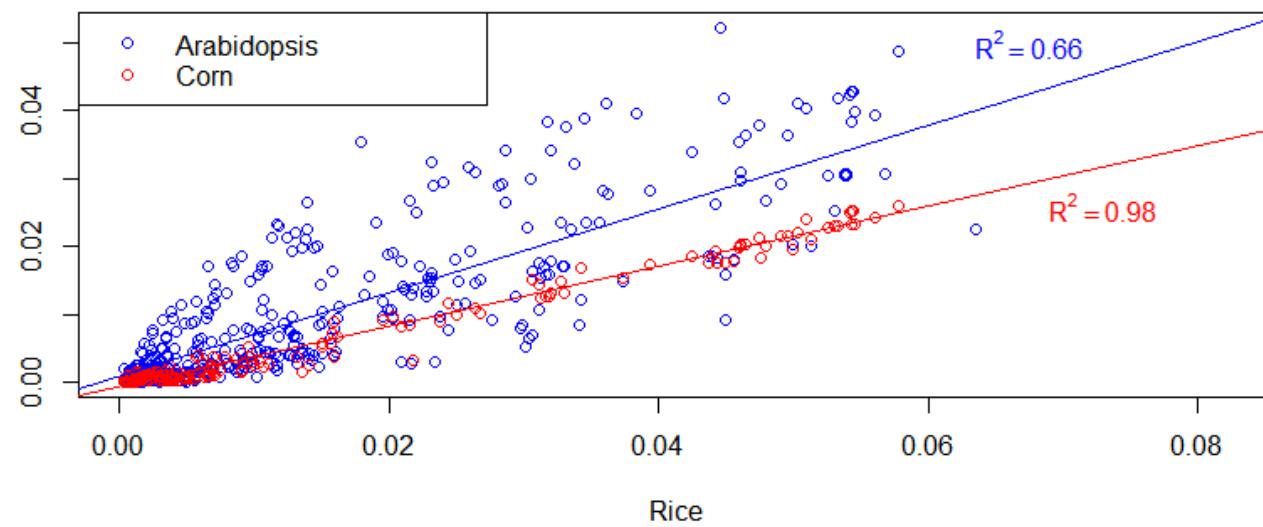
[Click here to download Figure Figure 7.docx](#)

Figure 8

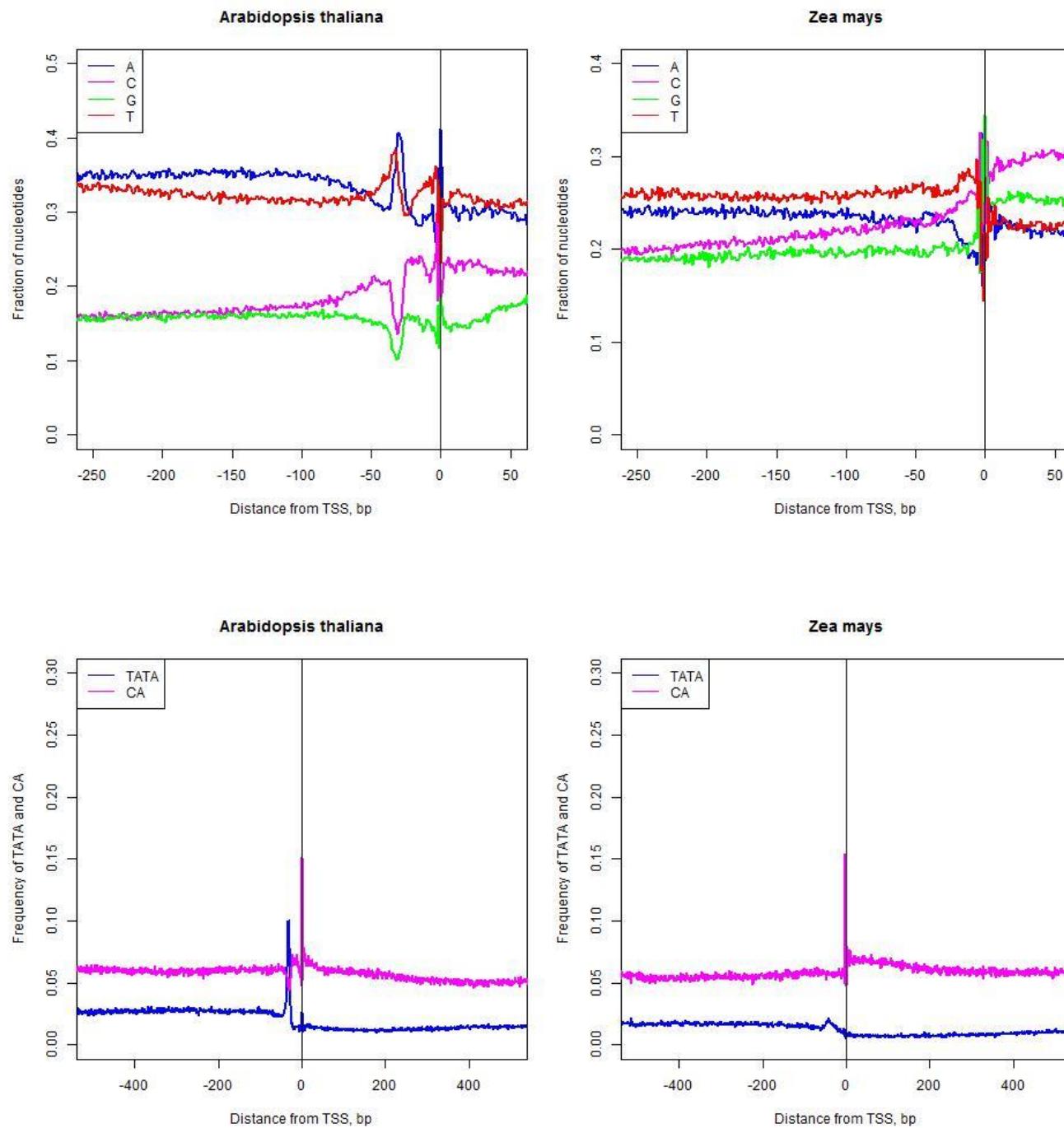
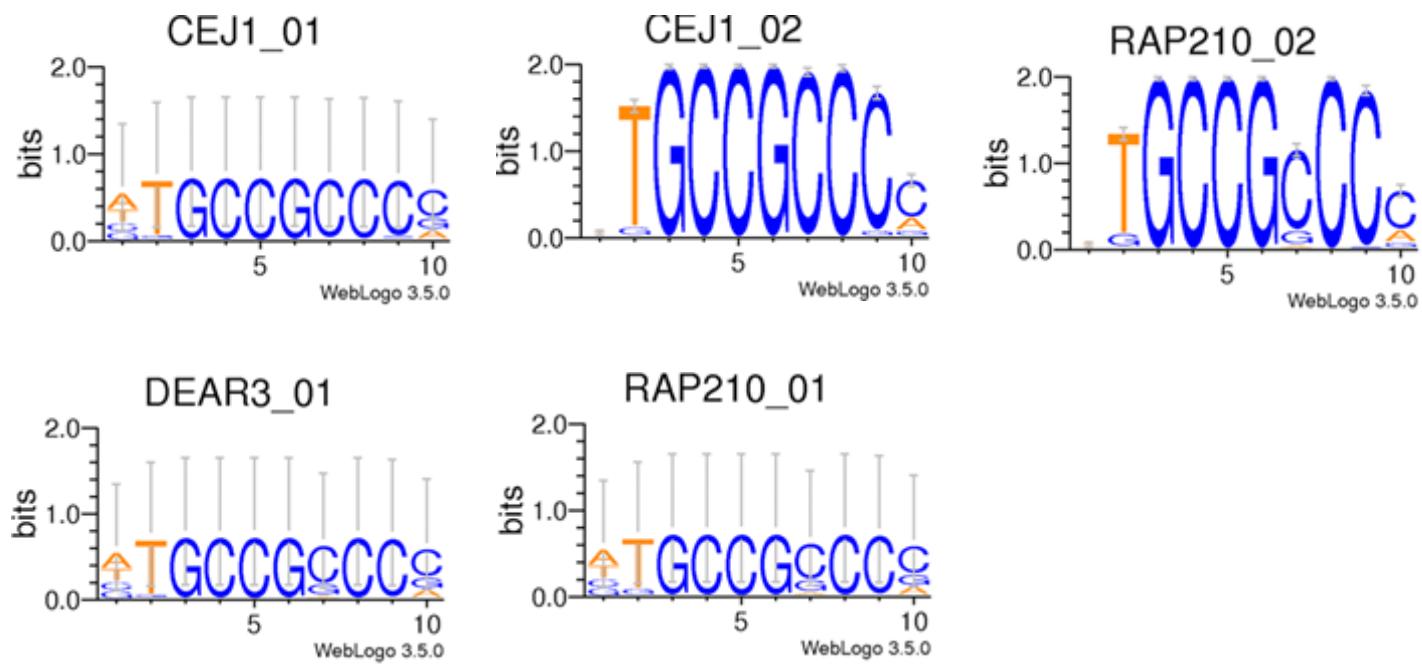
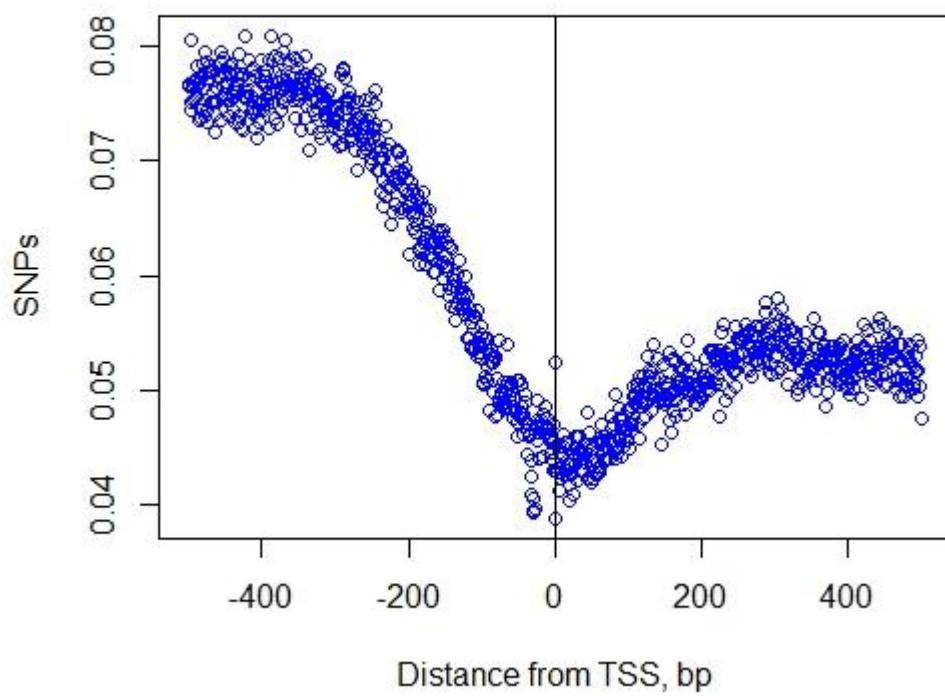
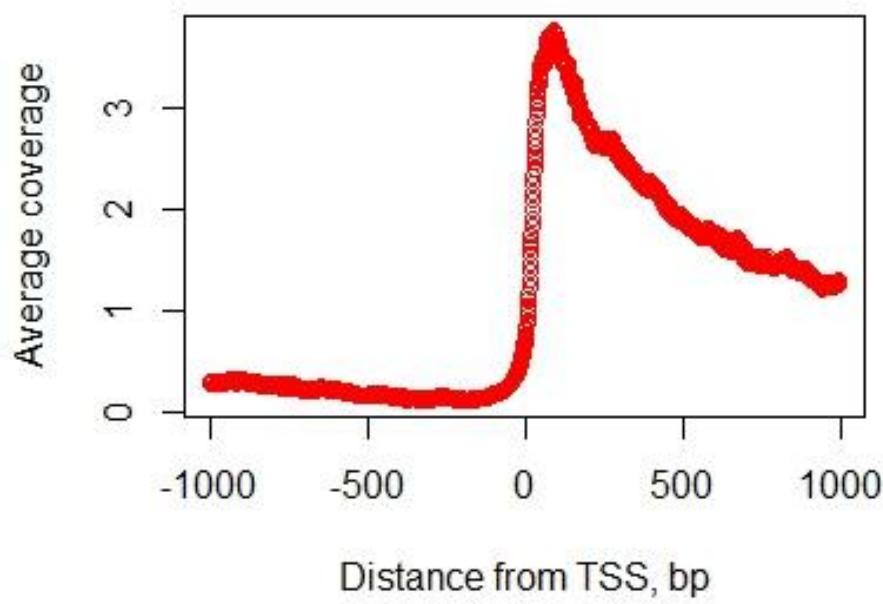
[Click here to download Figure Figure 8.docx](#)

Figure 9

[Click here to download Figure Figure 9.docx](#)





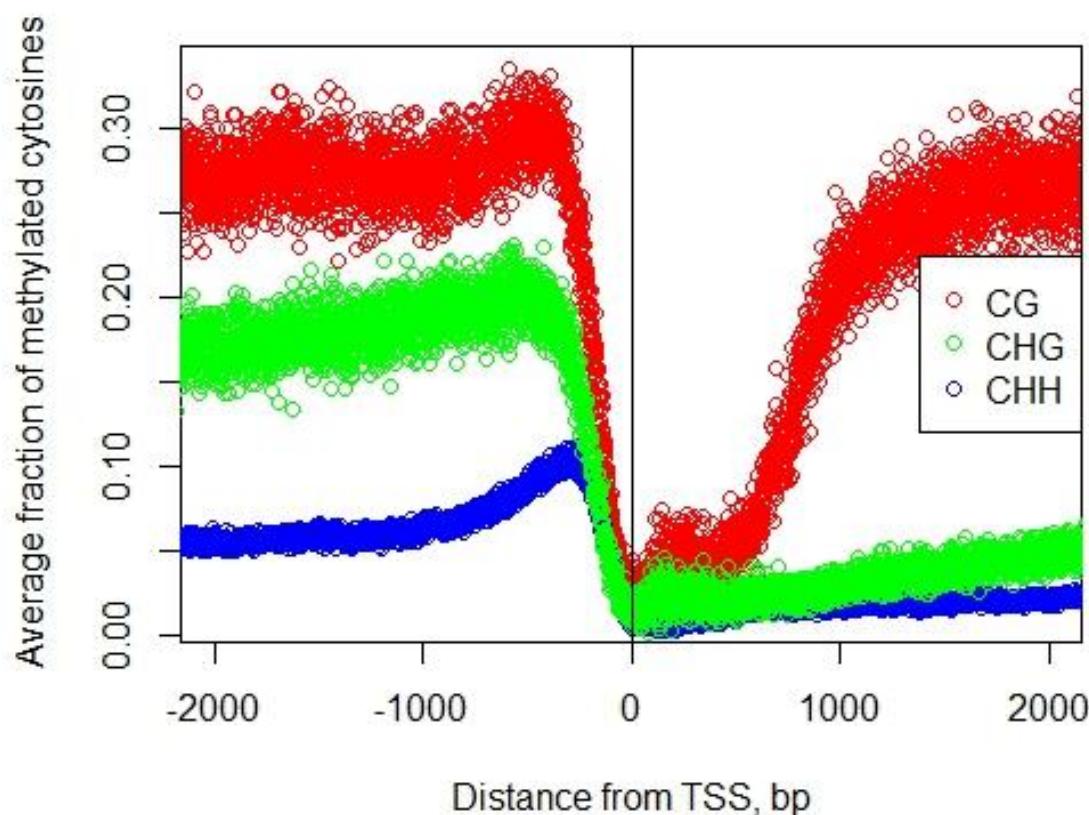
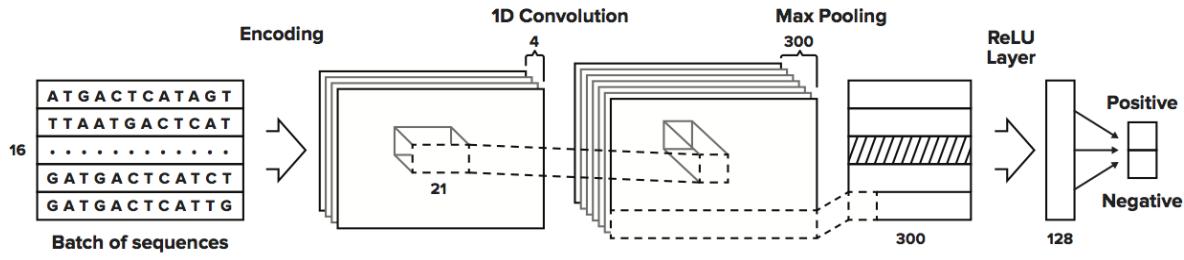


Figure 13

[Click here to download Figure 13.docx](#)



Click here to access/download
Supporting Information
FGEHESH_RICE.gff3





Click here to access/download
Supporting Information
Supplemental data.xlsx

