

1 Supplemental Materials and Methods and Results

2

3 Supplemental materials and methods

4 Construction of databases of upstream and first-intron sequences from maize

5 Flat files of putative maize promoter and first-intron sequences were created using custom Perl script
6 programs. For the promoter database, the 5' end of each predicted immature mRNA from the maize,
7 rice or Arabidopsis genomes were used to define the transcriptional start site (+1 - TSS) of each gene.
8 The TSS are defined from the predicted cDNAs for each gene in the 3 plant genomes. For each gene, 1 kb
9 of upstream sequence was extracted and used to create a flat file of predicted maize, rice and
10 Arabidopsis promoters. Whenever a sequence gap was identified, only the relevant downstream
11 sequence was extracted. If an upstream sequences available from a genome was less than 40 bp, it was
12 discarded from the flat file, as the motif discovery algorithms need a minimal sequence size to
13 accurately discover motifs. Three upstream sequence databases were created, each representing
14 different sequence lengths (1000 bp, 500 bp, 200 bp). For the database of first-intron sequences (for
15 maize only), a Perl script was written to recognize and retrieve predicted intron 1 sequences (based on
16 lower case annotation).

17 The sequences used to generate the databases were extracted from: MaizeSequence.org (release
18 5b.60 Working Set - <http://ftp.maizesequence.org>), the rice sequence from MSU.6.14 (downloaded from
19 Gramene - www.gramene.org) and the Arabidopsis sequence from TAIR 10.14
20 (<ftp://ftp.arabidopsis.org/home/tair/Genes/>). The sequence data originated from the
21 Genome Sequencing Center at the Washington University School of Medicine, St. Louis, USA [1].

22 Comparisons of current motif discovery tools

23 In order to determine whether existing motif discovery programs give equivalent results, Weeder,
24 MEME and BioProspector programs were compared. Each program was used to discover 213 known

25 TRANSFAC[®] motifs [2] embedded in 125 promoter data sets known as the benchmark data set,
26 previously generated to help researchers make direct comparisons of the effectiveness of different
27 motif discovery tools [3]. The data sets are grouped into three types: synthetic (Algorithm Markov, AM),
28 semi-synthetic (Algorithm Real, AR), and real biological promoters (Model Real, MR). The success rate of
29 each motif discovery program for each benchmark data set was measured using the following statistical
30 outputs either generated by the benchmark web application or calculated:

- 31 • the nucleotide level sensitivity (nSn):

$$32 \quad nSn = \frac{nTP}{nTP + nFN}$$

- 33 • the nucleotide False Discovery Ratio ($nFDR$):

$$34 \quad nFDR = \frac{nFP}{nTP + nFP}$$

- 35 • the nucleotide level correlation coefficient (nCC):

$$36 \quad nCC = \frac{nTP \cdot nTN - nFN \cdot nFP}{\sqrt{(nTP + nFN)(nTN + nFP)(nTP + nFP)(nTN + nFN)}}$$

37

38 Where nTP is the number of true positive motif nucleotides found; nTP , the number of true negative
39 motif nucleotides found; nFP , the number of false positive motif nucleotides found; nFN , the number of
40 false negative motif nucleotides found. nCC is a measure of the correlation between the known
41 nucleotide positions and the predicted nucleotide positions [4].

42 For each benchmark data set, the nucleotide level correlation coefficient score (nCC) of motif
43 prediction was compared between pairs of motif discovery programs. The results of all data sets were
44 plotted and compared using the Spearman correlation coefficient (Prism 5, GraphPad Software, USA).

45 Filters for each standalone program

46 As each program generates different sets of false-positives, a custom filter was designed for each
47 motif discovery tool to reduce the nFDR while preserving nTPs. In order to optimize the filter
48 parameters, candidate filters were applied to the Sandve et al. (2007) benchmark data set described
49 above and the best filters were chosen based on comparisons of nFDR and nCC (see above) using the
50 Friedman test (non-parametric repeated measures ANOVA) (Prism 5, GraphPad Software, USA).

51 Each filter was based on limiting the probability (p) that the frequency of the candidate motif in the
52 user data set (with sample size N) could occur randomly if a genome was repeatedly sampled using
53 sample size N . Two sampling algorithms were tested. The first one was the motif “enrichment”, without
54 replacement of the subject (promoter sequences) that uses the hypergeometric distribution [5-7]. The
55 second was with replacement of the sample subject following the binomial distribution [7, 8].

56 For MEME, the significance level (P_H) based on the hypergeometric distribution was used for the
57 motif filtering and was calculated as follows:

$$P_H = \sum_{i=n}^{\min(N,g)} \frac{\binom{N}{i} \binom{G-N}{g-i}}{\binom{G}{g}}$$

58
59 Where n is the number of benchmark data set sequences containing the predicted motif out of the
60 total number of sequences (N) belonging to that data set; G is the number of random sequences from
61 the organism that is the most overrepresented in the data set (G was set at 300); g is the number of
62 random sequences containing the predicted motif; n is the size of the motif in base pairs; and i is the
63 position within the motif. To retrieve predicted motifs in both the benchmark and random data sets,
64 FIMO was used [9]; the FIMO e-value threshold was set at $1e^{-4}$. Predicted motifs were filtered out when
65 P_H was > 0.05 . These threshold levels (FIMO value and P_H) were chosen based on optimization runs
66 using the benchmark data sets that increased the nSn but decreased nFDR.

67 For Weeder, there were two filters applied. The first filter removed predicted motifs with low
68 complexity DNA stretches (> 75% of the same nucleotide). The second filter was based on the binomial
69 distribution (P_B) based on previous works [7, 8], which gives an estimate of the probability that a motif
70 is non-random, calculated as follows:

$$P_B = \binom{N}{n} p^n \cdot (1-p)^{N-n}$$

71
72 Where n is the number of benchmark data set sequences containing the predicted motif out of the
73 total number of sequences (N) belonging to that data set; and P is the ratio of the number of random
74 sequences containing the predicted motif compared to the total number of random promoter
75 sequences (set at 300). To retrieve predicted motifs in both the benchmark and random data sets, Pscan
76 was used [10]; the Pscan score threshold was set at 0.97. Predicted motifs were filtered out when P_B
77 was > 0.3. Pscan and P_B significance levels were also selected after optimization runs using the
78 benchmark data sets that increased the nSn and decreased nFDR.

79 For BioProspector, the same binomial probability (P_B) used for Weeder was applied, except that the
80 Pscan score threshold was set at 0.90, and predicted motifs were filtered out when P_B was >0.7. The
81 significance levels were selected using the same method as Weeder.

82 To test each potential filter threshold, the average of three run results was used; for each run, a new
83 set of random promoter sequences was generated. This multiple-run method also helped to buffer
84 against the fact that BioProspector uses a stochastic algorithm, each run generating a different
85 prediction.

86 Combining multiple programs

87 As each standalone program appeared to predict different but overlapping sets of motifs, the effect
88 of combining all three filtered motif discovery programs was tested using the Sandve et al. (2007)
89 benchmark data set. The performance of the filtered combination against each standalone program was

90 compared using the nSn, nCC and nFDR scores (see above) with the Friedman test (non-parametric
91 repeated measures ANOVA) (Prism 5, GraphPad Software, USA).

92 Motif ranking using the MNCP score

93 The occurrence of motif m_x is determined in each of the promoters/first introns of the regulated user
94 data set u belonging to the regulated promoter/first intron population N_u . Each promoter/first intron
95 within the regulated data set is ranked according to the occurrence of motif m_x : promoter(s)/first
96 intron(s) with the highest motif occurrence are given the 1st rank. In parallel, the occurrence of motif m_x
97 is also determined in the random promoter/first intron data set r (regulated and non-regulated)
98 belonging to the random promoter/first intron population N_r . If the motif m_x is a regulator of the user
99 data set N_u , its occurrences should be higher than in the random data set N_r . Each promoter/first intron
100 (p_x) in the regulated data set has a rank $R_u(p_x)$ and another rank in the random data set $R_r(p_x)$. A
101 normalized ratio of the two ranks (C) for each promoter/first intron p_x is hence:

$$C(p_x) = \frac{R_u(p_x)/N_u}{R_r(p_x)/N_r}$$

102
103 C is calculated for each promoter/first introns containing motif m_x in the user data set. MNCP is the
104 mean of all the C values. If MNCP for motif m_x is greater than 1, that motif is more represented in the
105 regulated data set compared to the random data set. In Promzea, each motif is ranked according to its
106 relative MNCP score. Clover software is used to retrieve the motif and estimate its occurrence in the
107 user and random data sets from the maize genome.

108 Using user input sequences to extract corresponding promoter (and first-intron) regions

109 The Perl script was written to allow each user to generate a list of promoters (and first-intron
110 sequences) corresponding to only those genes of interest (e.g. promoters of co-expressed genes or
111 genes in the same biochemical pathway). The program accepts the following inputs: cDNA FASTA
112 sequence files, microarray probe-set ID or Gramene maize ID list (Figure 1). The GeneChip Maize

113 Genome Array (Affymetrix) is currently supported by Promzea. In the case of the cDNA FASTA files, the
114 Perl script matches each input cDNA to its corresponding MaizeSequence.org cDNA using standalone
115 BLAST (NCBI, version 2.2.23). The BLAST parameters were chosen empirically based on training data
116 (data not shown). A cDNA sequence in the genome is considered similar to an input cDNA if the
117 percentage identity is > 85% and the e-value is < 1e-50. The selected cDNAs from MaizeSequence.org
118 are then used to retrieve the corresponding upstream (or first intron sequences). For
119 microarray/Gramene ID inputs, the Perl script generates a list of the corresponding upstream (or intron
120 1) sequences directly.

121 Motif discovery, filtering, ranking and graphical output in Promzea

122 The user lists of promoters and first intron sequences generated above were used by our Perl script
123 as inputs into three complimentary motif discovery programs shown to retrieve different types of
124 motifs: MEME [11], BioProspector [12] and Weeder [13]. Predicted motifs from each standalone
125 program were filtered using the parameters described above. All the filtered results are regrouped.

126 Understanding motif function in Promzea using functional gene annotation

127 Gene annotations (e.g. anthocyanin pathway) can be used to help users understand the biological
128 function of a motif. The annotation can also be used by the user as a second form of motif validation as
129 the annotation-defined trait should relate to the user experiment. A flat file of well-described gene
130 annotations was first created using the "Functional-Annotations" files from MaizesSequence.org. Clover
131 [14] is used to search the maize promoter/first-intron flat file for each predicted motif which is then
132 matched to its corresponding gene annotation (Figure 1). The iGA Perl program [15] is used to calculate
133 if an annotation is overrepresented for a given motif. Retrieved annotations are represented as a pie
134 chart for the user using Chart:Clicker where each slice is $-\log(\text{annotation } p\text{-value})$. For a p-value equal

135 to zero, $-\log(p\text{-value})$ is equal to *infinite* which cannot be represented on a pie-chart. To circumvent this
136 problem, the choice has been made to replace zero p-values with p-values equal to 10^{-8} .

137 Validation of Promzea predictions using experimentally defined motifs

138 Promzea was further validated by searching a data set of promoters regulated by transcription
139 factors *C1* and *P* which activate the maize anthocyanin and phlobaphene biosynthetic pathways,
140 respectively [16]. To generate the input for Promzea, all cDNAs from Genbank that were annotated as
141 corresponding to the co-regulated genes were gathered in a FASTA file (Additional File 2); a promoter
142 sequence list was generated as described above. The 200 base promoter option was used for Promzea
143 analysis, as the literature shows that motifs important for the expression of anthocyanin biosynthetic
144 enzymes are within the first 200 bases of the promoter [17].

145 Testing of Promzea with co-expression data from the Maize Development Atlas

146 This gene expression data are available from the GEO database (Gene Expression Omnibus -
147 <http://www.ncbi.nlm.nih.gov/geo/>). The normalized microarray data were extracted from the GSE27004
148 experiment of the GEO database. The data for 60 different tissues were normalized using the RMA
149 method. The authors had produced tissue specific clusters of gene expression using the unweighted
150 pair-group method with the arithmetic mean (UPGMA) approach and Pearson's correlation. The cluster
151 gene lists were used from the Additional Table S4 of the publication {Sekhon, 2011 #80}. Promzea was
152 fed with the different tissue-specific gene lists. The similarities between each predicted Promzea motif
153 with experimentally defined motifs were determined by using the default setting of the STAMP software
154 used with the plant databases PLACE, Athamap and Agris.

155

156

157

158 Supplemental Results

159 Comparisons of current motif discovery programs using benchmark data sets

160 The first objective was to evaluate whether each standalone motif discovery program predicted the
161 known motif nucleotides in each data set to a similar extent. For this analysis, a previously generated
162 benchmark data set was used consisting of 213 known motifs embedded into sets of promoter
163 sequences [3]. For each data set, the nCC score was calculated, a measure of the correlation between
164 the known motif nucleotide positions and the predicted motif nucleotide positions [4]. When software
165 predicted nucleotides that exactly matched with the known binding sites (true positives, nTP), the nCC
166 score was +1, whereas an nCC score of ≤ 0 indicated a random prediction. For every paired program
167 comparison, Spearman correlations (r) of the benchmark nCC scores ranged from 0.14 to 0.36
168 (Additional Figure S1). In a large subset of benchmark data sets, Weeder predicted known motifs
169 effectively ($nCC > 0$), whereas BioProspector and MEME did not ($nCC \leq 0$) (Additional Figure 1B, C). These
170 results suggest that each motif discovery program retrieves a distinct set of motifs.

171

172

173

173 Supplemental References

174

- 175 1. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, Liang C, Zhang J, Fulton L, Graves
176 TA *et al*: **The B73 maize genome: complexity, diversity, and dynamics**. *Science* 2009,
177 **326**(5956):1112-1115.
- 178 2. Matys V, Kel-Margoulis OV, Fricke E, Liebich I, Land S, Barre-Dirrie A, Reuter I, Chekmenev D,
179 Krull M, Hornischer K *et al*: **TRANSFAC® and its module TRANSCompel®: transcriptional gene**
180 **regulation in eukaryotes**. *Nucleic Acids Research* 2006, **34**(suppl 1):D108-D110.
- 181 3. Sandve G, Abul O, Walseng V, Drablos F: **Improved benchmarks for computational motif**
182 **discovery**. *BMC Bioinformatics* 2007, **8**(1):193.
- 183 4. Tompa M, Li N, Bailey TL, Church GM, De Moor B, Eskin E, Favorov AV, Frith MC, Fu Y, Kent WJ *et*
184 *al*: **Assessing computational tools for the discovery of transcription factor binding sites**. *Nature*
185 *Biotechnology* 2005, **23**(1):137-144.
- 186 5. Sinha S, Ling X, Whitfield CW, Zhai C, Robinson GE: **Genome scan for cis-regulatory DNA motifs**
187 **associated with social behavior in honey bees**. *Proceedings of the National Academy of Sciences*
188 *USA* 2006, **103**(44):16352-16357.
- 189 6. Harbison CT, Gordon DB, Lee TI, Rinaldi NJ, Macisaac KD, Danford TW, Hannett NM, Tagne J-B,
190 Reynolds DB, Yoo J *et al*: **Transcriptional regulatory code of a eukaryotic genome**. *Nature* 2004,
191 **431**(7004):99-104.
- 192 7. Linhart C, Halperin Y, Shamir R: **Transcription factor and microRNA motif discovery: The**
193 **Amadeus platform and a compendium of metazoan target sets**. *Genome Research* 2008,
194 **18**(7):1180-1189.
- 195 8. Van Helden J, André B, Collado-Vides J: **Extracting regulatory sites from the upstream region of**
196 **yeast genes by computational analysis of oligonucleotide frequencies**. *Journal of Molecular*
197 *Biology* 1998, **281**(5):827-842.
- 198 9. Grant CE, Bailey TL, Noble WS: **FIMO: scanning for occurrences of a given motif**. *Bioinformatics*
199 2011, **27**(7):1017-1018.
- 200 10. Zambelli F, Pesole G, Pavesi G: **Pscan: finding over-represented transcription factor binding site**
201 **motifs in sequences from co-regulated or co-expressed genes**. *Nucleic Acids Research* 2009,
202 **37**(suppl 2):W247-W252.
- 203 11. Bailey TL, Elkan C: **Fitting a mixture model by expectation maximization to discover motifs in**
204 **biopolymers**. In: *Proceedings of the Second International Conference on Intelligent Systems for*
205 *Molecular Biology; Menlo Park, California*. AAAI Press 1994: 28-36.
- 206 12. Liu X, Brutlag D, Liu J: **BioProspector: discovering conserved DNA motifs in upstream regulatory**
207 **regions of co-expressed genes**. In: *Pacific Symposium on Biocomputing 2001*: Edited by Altman
208 RB, Dunker AK, Hunter L, Klein TE. 2001: 127-138.
- 209 13. Pavesi G, Zambelli F, Pesole G: **WeederH: an algorithm for finding conserved regulatory motifs**
210 **and regions in homologous sequences**. *BMC Bioinformatics* 2007, **8**(1):46.
- 211 14. Frith MC, Fu Y, Yu L, Chen JF, Hansen U, Weng Z: **Detection of functional DNA motifs via**
212 **statistical over-representation**. *Nucleic Acids Research* 2004, **32**(4):1372-1381.
- 213 15. Breitling R, Amtmann A, Herzyk P: **Iterative Group Analysis (iGA): A simple tool to enhance**
214 **sensitivity and facilitate interpretation of microarray experiments**. *BMC Bioinformatics* 2004,
215 **5**(1):34.
- 216 16. Dooner HK, Robbins TP, Jorgensen RA: **Genetic and developmental control of anthocyanin**
217 **biosynthesis**. *Annual Review of Genetics* 1991, **25**(1):173-199.
- 218 17. Bodeau JP, Walbot V: **Structure and regulation of the maize Bronze2 promoter**. *Plant Molecular*
219 *Biology* 1996, **32**(4):599-609.